

# Fast Kernel-Based Independent Component Analysis

Hao Shen\*, *Member, IEEE*, Stefanie Jegelka and Arthur Gretton

## Abstract

Recent approaches to Independent Component Analysis (ICA) have used kernel independence measures to obtain highly accurate solutions, particularly where classical methods experience difficulty (for instance, sources with near-zero kurtosis). FastKICA (Fast HSIC-based Kernel ICA) is a new optimisation method for one such kernel independence measure, the Hilbert-Schmidt Independence Criterion (HSIC). The high computational efficiency of this approach is achieved by combining geometric optimisation techniques, specifically an approximate Newton-like method on the orthogonal group, with accurate estimates of the gradient and Hessian based on an incomplete Cholesky decomposition. In contrast to other efficient kernel-based ICA algorithms, FastKICA is applicable to any twice differentiable kernel function. Experimental results for problems with large numbers of sources and observations indicate that FastKICA provides more accurate solutions at a given cost than gradient descent on HSIC. Comparing with other recently published ICA methods, FastKICA is competitive in terms of accuracy, relatively insensitive to local minima when initialised far from independence, and more robust towards outliers. An analysis of the local convergence properties of FastKICA is provided.

## Index Terms

Independent component analysis, Hilbert-Schmidt independence criterion, kernel methods, approximate Newton-like methods, the orthogonal group.

Copyright (c) 2008 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

H. Shen is with the Institute for Data Processing, Technische Universität München, 80290 München, Germany. (email: [hao.shen@tum.de](mailto:hao.shen@tum.de)).

S. Jegelka and A. Gretton are with Department of Empirical Inference for Machine Learning and Perception, Max Planck Institute for Biological Cybernetics, 72076 Tübingen, Germany. (email: [Stefanie.Jegelka,Arthur.Gretton}@tuebingen.mpg.de](mailto:{Stefanie.Jegelka,Arthur.Gretton}@tuebingen.mpg.de})).

## I. INTRODUCTION

The problem of Independent Component Analysis (ICA) involves the recovery of linearly mixed, statistically independent sources, in the absence of information about the source distributions beyond their mutual independence [1], [2]. The performance of ICA algorithms thus depends on the choice of the contrast function measuring the degree of statistical independence of the recovered signals, and on the optimisation technique used to obtain the estimated mixing coefficients.

Classical approaches, also referred to as parametric ICA approaches, construct their independence criteria according to certain hypothetical properties of the probability distributions, either by an explicit parametric model of these distributions via maximum likelihood [3], or by maximising certain statistics of the unmixed sources (often measures of non-Gaussianity, such as the kurtosis) [2], [4]. These approaches can therefore be less powerful than methods which explicitly model the source distributions, and can even fail completely when the modelling assumptions are not satisfied (e.g. a kurtosis-based contrast will not work for sources with zero kurtosis).

More recently, several approaches to ICA have been proposed that directly optimise nonparametric independence criteria. One option is to minimise the mutual information between the sources, as in [5]–[8]. Another approach is to use a characteristic function-based measure of mutual independence due to Kankainen [9] based on the pairwise criterion of Feuerverger [10], which was applied to ICA in [11], [12], and to ICA with post-nonlinear mixing in [13], [14].

Finally, a variety of kernel independence criteria have been employed in ICA. These criteria measure dependence using the spectrum of a covariance operator between mappings of the variables to high dimensional feature spaces, specifically reproducing kernel Hilbert spaces (RKHSs) [15]. The various kernel independence criteria differ in the way they summarise the covariance operator spectrum, and in the normalisation they use. They include the kernel canonical correlation [16], the kernel generalised variance [16], the spectral norm of the covariance operator (COCO) [17], the kernel mutual information [17], and the Hilbert-Schmidt Independence Criterion (HSIC) [18]. A biased empirical estimate of the HSIC statistic is in fact identical (as a function of its kernel argument) to the characteristic function-based criterion of [10], which is in turn identical to the  $\ell_2$  distance between Parzen window estimates of the joint density and the product of the marginals: see Rosenblatt [19]. When a Gaussian kernel is used and the sample size is fixed, the three statistics correspond exactly. As pointed out elsewhere [9], [10], however, the characteristic function-based statistic is more general than Rosenblatt's, since it admits a wider range of kernels while remaining an independence measure (a further difference is that the kernel

bandwidth may remain fixed for increasing sample size). Likewise, there exist universal kernels (in the sense of [20]: that is, kernels for which HSIC is zero iff the variables are independent, for any probability distribution [17, Theorem 6]) which have no equivalence with the characteristic function-based criterion of [9], [10]: examples are given in [20, Section 3] and [21, Section 3].<sup>1</sup> Thus, the RKHS criterion is a more general dependence measure than the characteristic function criterion, which is in turn more general than the  $\ell_2$  distance between Parzen window density estimates. Since the HSIC-based algorithm performs as well as or better than the remaining kernel dependence criteria for large sample sizes [18] on the benchmark data of [16], we use it as the contrast function in our present algorithm.

While the above studies report excellent demixing accuracy, efficient optimisation of these dependence measures for ICA remains an ongoing problem,<sup>2</sup> and a barrier to using nonparametric methods when the number of sources,  $m$ , is large. The main focus of the present work is thus on more efficient optimisation of kernel dependence measures. ICA is generally decomposed into two sub-problems [1], [12]: signal decorrelation or whitening, which is straightforward and is not discussed further, and optimisation over the set of orthogonal matrices (the orthogonal group,  $O(m)$ ), which is a differentiable manifold, and for which the bulk of the computation is required. The approach of [12], [16]–[18] is to perform gradient descent on  $O(m)$  in accordance with [24], choosing the step width by a Golden search. This is inefficient on two counts: gradient descent can require a very large number of steps for convergence even on relatively benign cost functions, and the Golden search requires many costly evaluations of the dependence measure. Although [23] propose a cheaper local quadratic approximation to choose the step size, this does not address the question of better search direction choice. An alternative solution is to use a Jacobi-type method [5], [6], [11], where the original optimisation problem on  $O(m)$  is decomposed into a sequence of one-dimensional sub-problems over a set of pre-determined curves, parameterised by the Jacobi angles. While the theoretical convergence properties of a Jacobi approach as compared with direct optimisation on  $O(m)$  are beyond the scope of this work, we perform an empirical evaluation against algorithms employing optimisation over Jacobi angles in our experiments.

In the present study, we develop an approximate Newton-like method for optimising the HSIC-based ICA contrast over  $O(m)$ , namely Fast HSIC-based Kernel ICA (FastKICA). A key feature of our approach is its computational efficiency, due to both the Newton-like optimisation and accurate low

<sup>1</sup>The RKHS approach also allows dependence testing on more general structures such as strings and graphs [22].

<sup>2</sup>Most of the effort in increasing efficiency has gone into cheaply and accurately approximating the independence measures [8], [12], [16], [23].

rank approximations of the independence measure and its derivatives. Importantly, these techniques do not require particular mathematical properties of the kernel (e.g. compact support, or that it be Laplace), but can be applied directly for any twice differentiable kernel function. The optimisation strategy follows recent studies on Newton-like methods for numerical optimisation on smooth manifolds in [25]. Approximate Newton-like algorithms have previously been developed in the case of classical ICA contrast functions [26], [27], where the authors use the diagonal structure of the Hessian at independence to greatly reduce complexity and computational cost. These earlier methods share the significant property of local quadratic convergence to a solution with correct source separation. We show the HSIC-based ICA contrast likewise has a diagonal Hessian at independence (this analysis originally appeared in [28]; note also that the diagonal property does not hold for the multivariate characteristic function-based counterpart [9] to the HSIC-based contrast), and that FastKICA is locally quadratically convergent to a correct unmixing matrix. Moreover, our experiments suggest that in the absence of a good initialisation, FastKICA converges more often to a correct solution than gradient descent methods. Previous kernel algorithms require either a large number of restarts or a good initial guess [12], [16], [17]. The current work is built on an earlier presentation by the authors in [29]. Compared with [29], the present study contains proofs of the main theorems (which were omitted in [29] due to space constraints); a proof of local quadratic convergence in the neighbourhood of the global solution; additional experiments on ICA performance vs “smoothness” of the departure from independence; and experiments on outlier resistance, for which our method strongly outperforms the other tested approaches.

The paper is organised as follows. In Section II, we briefly introduce the instantaneous noise-free ICA model, the HSIC-based ICA contrast, and a Newton-like method on  $O(m)$ . In Section III, we analyse the critical point condition and the structure of the Hessian of this contrast. We describe our ICA method, FastKICA, in Section IV, and prove local quadratic convergence. We also present an efficient implementation of FastKICA, based on the incomplete Cholesky decomposition [30]. Finally, our experiments in Section V compare FastKICA with several competing nonparametric approaches: RADICAL [5], MILCA [6], mutual information-based ICA (MICA) [31], and KDICA [8]. Experiments address performance and runtimes on large-scale problems, performance for decreasing smoothness of the departure of the mixture from independence (which makes demixing more difficult for algorithms that assume smooth source densities), and outlier resistance. Matlab code for FastKICA may be downloaded at [www.kyb.mpg.de/bs/people/arthur/fastkica.htm](http://www.kyb.mpg.de/bs/people/arthur/fastkica.htm)

## II. PRELIMINARIES: ICA, HSIC AND NEWTON-LIKE METHOD ON $O(m)$

### A. Linear Independent Component Analysis

The instantaneous noise-free ICA model takes the form

$$Z = AS, \quad (1)$$

where  $S \in \mathbb{R}^{m \times n}$  is a matrix containing  $n$  observations of  $m$  sources,  $A \in \mathbb{R}^{m \times m}$  is the mixing matrix (assumed here to have full rank),<sup>3</sup> and  $Z \in \mathbb{R}^{m \times n}$  contains the observed mixtures. Denote as  $s$  and  $z$  single columns of the matrices  $S$  and  $Z$ , respectively, and let  $s_i$  be the  $i$ -th source in  $s$ . ICA is based on the assumption that the components  $s_i$  of  $s$ , for all  $i = 1 \dots m$ , are mutually statistically independent. This ICA model (1) is referred to as instantaneous as a way of describing the dual assumptions that the observation vector  $z$  depends only on the source vector  $s$  at that instant, and the source samples  $s$  are drawn independently and identically from  $\text{Pr}_s$ . As a consequence of the first assumption, the mixture samples  $z$  are likewise drawn independently and identically from  $\text{Pr}_z$ .

The task of ICA is to recover the independent sources via an estimate  $B$  of the inverse of the mixing matrix  $A$ , such that the recovered signals  $Y = BAS$  have mutually independent components. It is well known that if at most one of the sources  $s$  is Gaussian, the mixing matrix  $A$  can be identified up to an ordering and scaling of the recovered sources [1]. This means the unmixing matrix  $B$  is the inverse of  $A$  up to an  $m \times m$  permutation matrix  $P$  and an  $m \times m$  diagonal (scaling) matrix  $D$ , i.e.,  $B = PDA^{-1}$ . To reduce the computational complexity the mixtures  $Z$  are usually pre-whitened via principal component analysis (PCA) [1], [12]. Whitening corresponds to finding a matrix  $V \in \mathbb{R}^{m \times m}$  such that  $W = VZ = VAS \in \mathbb{R}^{m \times n}$  with  $\mathbb{E}[ww^\top] = I$ , where  $W$  are referred to as the whitened observations. While this pre-whitening step is less statistically efficient than solving directly for the unconstrained mixing matrix [4, Section VI.B], the optimisation problem in the pre-whitened case is easier. Assuming the sources  $s_i$  have zero mean and unit variance, we find  $VA \in \mathbb{R}^{m \times m}$  to be orthogonal. Therefore, the whitened noise-free ICA unmixing model becomes

$$Y = X^\top W, \quad (2)$$

where  $X \in \mathbb{R}^{m \times m}$  is an orthogonal unmixing matrix (i.e.,  $X^\top X = I$ ), and  $Y \in \mathbb{R}^{m \times n}$  contains our estimates of the sources. Let  $O(m)$  denote the orthogonal group:

$$O(m) := \{X \in \mathbb{R}^{m \times m} | X^\top X = I\}. \quad (3)$$

<sup>3</sup>In other words, we do not address the more difficult problems of undercomplete or overcomplete ICA (corresponding to more mixtures than sources, or fewer mixtures than sources, respectively).

We focus in the remainder of this work on the problem of finding  $X \in O(m)$  so as to recover the mutually statistically independent sources via the model (2). Thus, we next describe our measure of independence.

### B. The Hilbert-Schmidt Independence Criterion

The Hilbert-Schmidt Independence Criterion (HSIC) is a bivariate independence measure obtained as the squared Hilbert-Schmidt (HS) norm of the covariance operator between mappings to RKHSs [18], and generalises the characteristic function-based criterion originally proposed by Feuerverger [10]. The Hilbert space  $\mathcal{F}$  of functions from a compact subset  $\mathcal{U} \subset \mathbb{R}$  to  $\mathbb{R}$  is an RKHS if at each  $u \in \mathcal{U}$ , the point evaluation operator  $\delta_u : \mathcal{F} \rightarrow \mathbb{R}$ , which maps  $f \in \mathcal{F}$  to  $f(u) \in \mathbb{R}$ , is a continuous linear functional. To each point  $u \in \mathcal{U}$ , there corresponds an element  $\alpha_u \in \mathcal{F}$ , also called the feature map, such that  $\langle \alpha_u, \alpha_{u'} \rangle_{\mathcal{F}} = \psi(u, u')$ , where  $\psi : \mathcal{U} \times \mathcal{U} \rightarrow \mathbb{R}$  is a unique positive definite kernel. We also define a second RKHS  $\mathcal{G}$  with respect to  $\mathcal{U}$ , with feature map  $\beta_v \in \mathcal{G}$  and corresponding kernel  $\langle \beta_v, \beta_{v'} \rangle_{\mathcal{G}} = \widehat{\psi}(v, v')$ .

Let  $\Pr_{u,v}$  be a joint measure on  $(\mathcal{U} \times \mathcal{U}, \Gamma \times \Lambda)$  (here  $\Gamma$  and  $\Lambda$  are Borel  $\sigma$ -algebras on  $\mathcal{U}$ ), with associated marginal measures  $\Pr_u$  and  $\Pr_v$ . The covariance operator  $C_{uv} : \mathcal{G} \rightarrow \mathcal{F}$  is defined as

$$\langle f, C_{uv}(g) \rangle_{\mathcal{F}} = \mathbb{E}[f(u)g(v)] - \mathbb{E}[f(u)]\mathbb{E}[g(v)] \quad (4)$$

for all  $f \in \mathcal{F}$  and  $g \in \mathcal{G}$ . The squared HS norm of the covariance operator  $C_{uv}$ , denoted as HSIC, is then

$$\|C_{uv}\|_{\text{HS}}^2 = \mathbb{E}_{u,u',v,v'} \left[ \psi(u, u') \widehat{\psi}(v, v') \right] \quad (5a)$$

$$+ \mathbb{E}_{u,u'} \left[ \psi(u, u') \right] \mathbb{E}_{v,v'} \left[ \widehat{\psi}(v, v') \right] \quad (5b)$$

$$- 2\mathbb{E}_{u,v} \left[ \mathbb{E}_{u'} \left[ \psi(u, u') \right] \mathbb{E}_{v'} \left[ \widehat{\psi}(v, v') \right] \right] \quad (5c)$$

(see [18] for details), where  $(u, v) \sim \Pr_{u,v}$  and  $(u', v') \sim \Pr_{u,v}$  are independent random variables drawn from the same distribution, and  $\mathbb{E}[\cdot]$  denotes the expectation over the corresponding random variables.

As long as the kernels  $\psi(u, \cdot) \in \mathcal{F}$  and  $\widehat{\psi}(u, \cdot) \in \mathcal{G}$  are universal in the sense of [20], e.g., the Gaussian and Laplace kernels,  $\|C_{uv}\|_{\text{HS}}^2 = 0$  if and only if  $u$  and  $v$  are statistically independent [18, Theorem 4].

In this work, we confine ourselves to a Gaussian kernel, and use the same kernel for both  $\mathcal{F}$  and  $\mathcal{G}$ ,

$$\psi(a, b) = \widehat{\psi}(a, b) := \phi(a - b) = \exp\left(-\frac{(a-b)^2}{2\lambda^2}\right). \quad (6)$$

As discussed in the introduction, the empirical expression for HSIC in [18] is identical to Feuerverger's independence criterion [10] and Rosenblatt's  $\ell^2$  independence statistic [19] for a Gaussian kernel at a given sample size.

We now construct an HSIC-based ICA contrast for more than two random variables. In the ICA model (1), the components  $s_i$  of the sources  $s$  are mutually statistically independent if and only if their probability distribution factorises as  $\Pr_s = \prod_{i=1}^m \Pr_{s_i}$ . Although the random variables are pairwise independent if they are mutually independent, where pairwise independence is defined as  $\Pr_{s_i} \Pr_{s_j} = \Pr_{s_i, s_j}$  for all  $i \neq j$ , the reverse does not generally hold: pairwise independence does not imply mutual independence. Nevertheless, Theorem 11 of [1] shows that in the ICA setting, unmixed components can be uniquely identified using only the pairwise independence between components of the recovered sources  $Y$ , since pairwise independence between components of  $Y$  in this case implies their mutual independence (and thus recovery of the sources  $S$ ).<sup>4</sup> Hence, by summing all unique pairwise HSIC measures, an HSIC-based contrast function over the estimated signals  $Y \in \mathbb{R}^{m \times n}$  is defined as

$$H : O(m) \rightarrow \mathbb{R},$$

$$H(X) := \sum_{1 \leq i < j \leq m} \mathbb{E}_{k,l} \left[ \phi \left( x_i^\top \bar{w}_{kl} \right) \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \quad (7a)$$

$$+ \mathbb{E}_{k,l} \left[ \phi \left( x_i^\top \bar{w}_{kl} \right) \right] \mathbb{E}_{k,l} \left[ \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \quad (7b)$$

$$- 2 \mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi \left( x_i^\top \bar{w}_{kl} \right) \right] \mathbb{E}_l \left[ \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \right], \quad (7c)$$

where  $X := [x_1, \dots, x_m] \in O(m)$ ,  $\bar{w}_{kl} = w_k - w_l \in \mathbb{R}^m$  denotes the difference between  $k$ -th and  $l$ -th samples of the whitened observations, and  $\mathbb{E}_{k,l}[\cdot]$  represents the empirical expectation over all  $k$  and  $l$ .

### C. Newton-like Methods on $O(m)$

In this section, we briefly review some basic concepts regarding Newton-like methods on the orthogonal group  $O(m)$ . We refer to [32], [33] for an excellent introduction to differential geometry, and to [34] for an introduction to optimisation algorithms on differentiable manifolds. We will review both the classical approach to Newton-type optimization on smooth manifolds [24], and then describe a more recently developed Newton-like method [25], which we apply on  $O(m)$ .

We consider the orthogonal group  $O(m)$  as an  $m(m-1)/2$  dimensional embedded submanifold of  $\mathbb{R}^{m \times m}$ , and denote the set of all  $m \times m$  skew-symmetric matrices by  $\mathfrak{so}(m) := \{\Omega \in \mathbb{R}^{m \times m} | \Omega = -\Omega^\top\}$ . Note that  $\mathfrak{so}(m)$  is isomorphic to  $\mathbb{R}^{m(m-1)/2}$ , written  $\mathfrak{so}(m) \cong \mathbb{R}^{m(m-1)/2}$ . The tangent space  $T_X O(m)$

<sup>4</sup>That said, in the finite sample setting, the statistical performance obtained by optimizing over a pairwise independence criterion might differ from that of a mutual independence criterion.

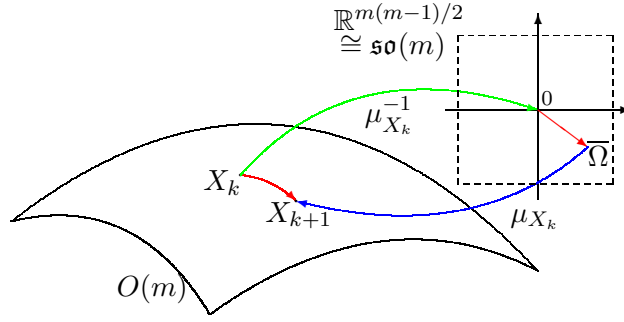


Fig. 1. Illustration of a Newton-like method on  $O(m)$ .

of  $O(m)$  at point  $X \in O(m)$  is given by

$$T_X O(m) := \{ \Xi \in \mathbb{R}^{m \times m} \mid \Xi = X\Omega, \Omega \in \mathfrak{so}(m) \}. \quad (8)$$

A typical approach in developing a Newton-type method for optimising a smooth function  $H: O(m) \rightarrow \mathbb{R}$  is to endow the manifold  $O(m)$  with a Riemannian structure: see [24]. Rather than moving along a straight line as in the Euclidean case, a Riemannian Newton iteration moves along a geodesic<sup>5</sup> in  $O(m)$ . For a given tangent space direction  $\Xi = X\Omega \in T_X O(m)$ , the geodesic  $\gamma_X$  through  $X \in O(m)$  with respect to the Riemannian metric  $\langle X\Omega_1, X\Omega_2 \rangle := -\text{tr} \Omega_1 \Omega_2$ , for  $X\Omega_1, X\Omega_2 \in T_X O(m)$ , is

$$\gamma_X: \mathbb{R} \rightarrow O(m), \quad \varepsilon \mapsto X \exp(\varepsilon X^\top \Xi), \quad (9)$$

with  $\gamma_X(0) = X$  and  $\dot{\gamma}_X(0) = \Xi$ . Here,  $\exp(\cdot)$  denotes matrix exponentiation. It is well known that this method enjoys the significant property of local quadratic convergence.

More recently, a novel Newton-like method on smooth manifolds was proposed [25]. This method has lower complexity than the classical approach, but retains the property of local quadratic convergence. We adapt the general formulation from [25] to the present setting, the orthogonal group  $O(m)$ .

For every point  $X \in O(m)$ , there exists a smooth map

$$\mu_X: \mathbb{R}^{m(m-1)/2} \rightarrow O(m), \quad \mu_X(0) = X, \quad (10)$$

which is a local diffeomorphism around  $0 \in \mathbb{R}^{m(m-1)/2}$ , i.e., both the map  $\mu_X$  and its inverse are locally smooth around 0. Let  $X^* \in O(m)$  be a nondegenerate critical point of a smooth contrast function

<sup>5</sup>The geodesic is a concept on manifolds analogous to the straight line in a Euclidean space. It allows parallel transportation of tangent vectors on manifolds.



$H: O(m) \rightarrow \mathbb{R}$ . If there exists an open neighborhood  $\mathcal{U}(X^*) \subset O(m)$  of  $X^* \in O(m)$  and a smooth map

$$\tilde{\mu}: \mathcal{U}(X^*) \times \mathbb{R}^{m(m-1)/2} \rightarrow O(m), \quad (11)$$

such that  $\tilde{\mu}(X, \bar{\Omega}) = \mu_X(\bar{\Omega})$  for all  $X \in \mathcal{U}(X^*)$  and  $\bar{\Omega} \in \mathbb{R}^{m(m-1)/2}$ , we call  $\{\mu_X\}_{X \in \mathcal{U}(X^*)}$  a locally smooth family of parametrisations around  $X^*$ .

Let  $X_k \in O(m)$  be the  $k$ -th iteration point of a Newton-like method for minimising the contrast function  $H$ . A local cost function can then be constructed by composing the original function  $H$  with the local parametrisation  $\mu_{X_k}$  around  $X_k$ , i.e.,  $H \circ \mu_{X_k}: \mathbb{R}^{m(m-1)/2} \rightarrow \mathbb{R}$ , which is a smooth function locally defined on the Euclidean space  $\mathbb{R}^{m(m-1)/2}$ . Thus, one Euclidean Newton step  $\bar{\Omega} \in \mathbb{R}^{m(m-1)/2}$  for  $H$ , expressed in local coordinates, is the solution of the linear equation

$$\mathcal{H}(H \circ \mu_{X_k})(0)\bar{\Omega} = -\nabla(H \circ \mu_{X_k})(0), \quad (12)$$

where  $\nabla(H \circ \mu_{X_k})(0)$  and  $\mathcal{H}(H \circ \mu_{X_k})(0)$  are respectively the gradient and Hessian of  $H \circ \mu_{X_k}$  at  $0 \in \mathbb{R}^{m(m-1)/2}$  with respect to the standard Euclidean inner product on the parameter space  $\mathbb{R}^{m(m-1)/2}$ . Finally, projecting the Newton step  $\bar{\Omega}$  in the local coordinates back to  $O(m)$  using the local parametrisation<sup>6</sup>  $\mu_{X_k}$  completes a basic iteration of a Newton-like method on  $O(m)$  (see Fig. 1 for an illustration of the method).

We now describe our choice of local parametrisation  $\mu_X$  of  $O(m)$ , which follows directly from the geodesic expression (9). We define  $\Omega = (\omega_{ij})_{i,j=1}^m \in \mathfrak{so}(m)$  as before and let  $\bar{\Omega} = (\omega_{ij})_{1 \leq i < j \leq m} \in \mathbb{R}^{m(m-1)/2}$  in a lexicographical order. A local parametrisation of  $O(m)$  around a point  $X \in O(m)$  is given by

$$\mu_X: \mathbb{R}^{m(m-1)/2} \cong \mathfrak{so}(m) \rightarrow O(m), \quad \bar{\Omega} \mapsto X \exp(\Omega), \quad (13)$$

which is a local diffeomorphism around  $0 \in \mathbb{R}^{m(m-1)/2}$ , i.e.  $\mu_X(0) = X$ .

To summarise, a Newton-like method for minimising the contrast function  $H: O(m) \rightarrow \mathbb{R}$  can be stated as follows:

<sup>6</sup>In the general setting, the second parametrisation can be from a different family to the first: see [25] for details. This is not the case for our algorithm, however.

Newton-like method on $O(m)$
------------------------------

Step 1: Given an initial guess  $X_0 \in O(m)$  and set  $k = 0$ .

Step 2: Calculate  $H \circ \mu_{X_k} : \mathbb{R}^{m(m-1)/2} \cong \mathfrak{so}(m) \rightarrow \mathbb{R}$ .

Step 3: Compute the Euclidean Newton step, i.e., solve the linear system for  $\bar{\Omega} \in \mathbb{R}^{m(m-1)/2}$ ,

$$\mathcal{H}(H \circ \mu_{X_k})(0)\bar{\Omega} = -\nabla(H \circ \mu_{X_k})(0).$$

Step 4: Set  $X_{k+1} = \nu_{X_k}(\bar{\Omega})$ .

Step 5: If  $\|X_{k+1} - X_k\|_F$  is small enough, stop.

Otherwise, set  $k = k + 1$  and go to Step 2.

Here,  $\|\cdot\|_F$  is the Frobenius norm of matrices. According to Theorem 1 in [25], this Newton-like method is locally quadratically convergent to  $X^* \in O(m)$ . For an approximate Newton-like method on  $O(m)$ , we replace the true Hessian  $\mathcal{H}(H \circ \mu_X)(0)$  by an approximate Hessian  $\tilde{\mathcal{H}}(H \circ \mu_X)(0)$  which is more efficient to compute. We will prove that local quadratic convergence is still obtained with this approximation.

### III. THE HSIC-BASED ICA CONTRAST AT INDEPENDENCE

In this section we examine the critical points of the HSIC-based ICA contrast at independence. It turns out that any correct unmixing matrix which is a global minimum of the HSIC-based contrast is a nondegenerate critical point.

Let  $X = [x_1, \dots, x_m] \in O(m)$ . By the chain rule, the first derivative of  $H$  in the direction  $\Xi = [\xi_1, \dots, \xi_m] \in T_X O(m)$  is

$$\begin{aligned} \mathrm{D}H(X)\Xi &= \frac{\mathrm{d}}{\mathrm{d}\varepsilon}(H \circ \gamma_X)(\varepsilon)\Big|_{\varepsilon=0} \\ &= \sum_{i,j=1;i \neq j}^m \mathbb{E}_{k,l} \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) \xi_i^\top \bar{w}_{kl} \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \end{aligned} \quad (14a)$$

$$+ \mathbb{E}_{k,l} \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) \xi_i^\top \bar{w}_{kl} \right] \mathbb{E}_{k,l} \left[ \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \quad (14b)$$

$$- 2\mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) \xi_i^\top \bar{w}_{kl} \right] \mathbb{E}_l \left[ \phi \left( x_j^\top \bar{w}_{kl} \right) \right] \right]. \quad (14c)$$

Setting the above derivative to zero, we can characterise the critical points of the HSIC-based ICA contrast function defined in (7). Obviously, the critical point condition depends not only on the statistical characteristics of the sources, but also on the properties of the kernel function. It is hard to characterise

all critical points of HSIC in full generality: thus, we deal only with those critical points that occur at the ICA solution.

*Lemma 1:* Let  $X^* \in O(m)$  be a correct unmixing matrix of the model (2). Then  $X^*$  is a nondegenerate critical point of the HSIC-based ICA contrast (7), i.e.,  $DH(X^*)\Xi = 0$  for arbitrary  $\Xi \in T_{X^*}O(m)$ .

*Proof:* We recall that  $\|C_{uv}\|_{\text{HS}}^2 \geq 0$  and  $\|C_{uv}\|_{\text{HS}}^2 = 0$  if and only if  $\text{Pr}_{u,v} = \text{Pr}_u \text{Pr}_v$ . In other words,  $X^* \in O(m)$  is a global minimum of the HSIC-based contrast function. Theorem 4 in [18] shows that any small displacement of  $X^* \in O(m)$  along the geodesics will result in an increase in the score of the contrast function  $H$  between the recovered sources (only larger perturbations, e.g. a swap of the source orderings, would yield independent sources). According to Theorem 4.2 in [35], the Hessian of  $H$  at  $X^*$  is positive definite. Hence, any correct separation point  $X^* \in O(m)$  is a nondegenerate critical point of the HSIC-based contrast function  $H$ .  $\blacksquare$

In what follows, we investigate the structure of the Hessian of the HSIC-based contrast  $H$  at independence, i.e. at a correct unmixing matrix  $X^* \in O(m)$ . We first compute the second derivative of  $H$  at  $X \in O(m)$  in direction  $\Xi = [\xi_1, \dots, \xi_m] \in T_X O(m)$ ,

$$\begin{aligned} D^2H(X)(\Xi, \Xi) &= \left. \frac{d^2}{d\varepsilon^2} (H \circ \gamma_X)(\varepsilon) \right|_{\varepsilon=0} \\ &= \sum_{i,j=1; i \neq j}^m \mathbb{E}_{k,l} \left[ \phi''(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \bar{w}_{kl}^\top \xi_j \phi(x_j^\top \bar{w}_{kl}) \right] \end{aligned} \quad (15a)$$

$$- \mathbb{E}_{k,l} \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \Xi X^\top \bar{w}_{kl} \phi(x_j^\top \bar{w}_{kl}) \right] \quad (15b)$$

$$+ \mathbb{E}_{k,l} \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \bar{w}_{kl}^\top \xi_j \phi'(x_j^\top \bar{w}_{kl}) \right] \quad (15c)$$

$$+ \mathbb{E}_{k,l} \left[ \phi''(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \bar{w}_{kl}^\top \xi_i \right] \mathbb{E}_{k,l} \left[ \phi(x_j^\top \bar{w}_{kl}) \right] \quad (15d)$$

$$- \mathbb{E}_{k,l} \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \Xi X^\top \bar{w}_{kl} \right] \mathbb{E}_{k,l} \left[ \phi(x_j^\top \bar{w}_{kl}) \right] \quad (15e)$$

$$+ \mathbb{E}_{k,l} \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \right] \mathbb{E}_{k,l} \left[ \phi'(x_j^\top \bar{w}_{kl}) \xi_j^\top \bar{w}_{kl} \right] \quad (15f)$$

$$- 2\mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi''(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \bar{w}_{kl}^\top \xi_i \right] \mathbb{E}_l \left[ \phi(x_j^\top \bar{w}_{kl}) \right] \right] \quad (15g)$$

$$+ 2\mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \Xi X^\top \bar{w}_{kl} \right] \mathbb{E}_l \left[ \phi(x_j^\top \bar{w}_{kl}) \right] \right] \quad (15h)$$

$$- 2\mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi'(x_i^\top \bar{w}_{kl}) \xi_i^\top \bar{w}_{kl} \right] \mathbb{E}_l \left[ \phi'(x_j^\top \bar{w}_{kl}) \xi_j^\top \bar{w}_{kl} \right] \right]. \quad (15i)$$

Let  $\Omega = [\omega_1, \dots, \omega_m] = (\omega_{ij})_{i,j=1}^m \in \mathfrak{so}(m)$ . For  $X = X^*$ , a tedious but straightforward computation (see Appendix A) gives

$$D^2H(X)(X\Omega, X\Omega)|_{X=X^*} = \sum_{1 \leq i < j \leq m} \omega_{ij}^2 (\kappa_{ij} + \kappa_{ji}), \quad (16)$$

where

$$\kappa_{ij} = 2\mathbb{E}_k [\mathbb{E}_l[\bar{s}_{kli}]\mathbb{E}_l[\phi'(\bar{s}_{kli})]]\mathbb{E}_k [\mathbb{E}_l[\bar{s}_{klj}]\mathbb{E}_l[\phi'(\bar{s}_{klj})]] \quad (17a)$$

$$+ \mathbb{E}_{k,l}[\phi''(\bar{s}_{kli})]\mathbb{E}_{k,l}[(\bar{s}_{klj})^2\phi(\bar{s}_{klj})] \quad (17b)$$

$$+ 2\mathbb{E}_{k,l}[\phi''(\bar{s}_{kli})]\mathbb{E}_{k,l}[\phi(\bar{s}_{klj})] \quad (17c)$$

$$- 2\mathbb{E}_{k,l}[\phi''(\bar{s}_{kli})]\mathbb{E}_k [\mathbb{E}_l[(\bar{s}_{klj})^2]\mathbb{E}_l[\phi(\bar{s}_{klj})]] \quad (17d)$$

$$- \mathbb{E}_{k,l}[\phi'(\bar{s}_{kli})\bar{s}_{kli}]\mathbb{E}_{k,l}[\phi'(\bar{s}_{klj})\bar{s}_{klj}]. \quad (17e)$$

*Remark 1:* Without loss of generality, let  $\bar{\Omega} = (\omega_{ij})_{1 \leq i < j \leq m} \in \mathbb{R}^{m(m-1)/2}$  in a lexicographical order. The quadratic form (16) is a sum of pure squares, which indicates that the Hessian of the contrast function  $H$  in (7) at the desired critical point  $X^* \in O(m)$ , i.e., the symmetric bilinear form  $\mathcal{H}H(X^*): T_{X^*}O(m) \times T_{X^*}O(m) \rightarrow \mathbb{R}$ , is diagonal with respect to the standard basis of  $\mathbb{R}^{m(m-1)/2}$ . Furthermore, following the arguments in Lemma 1, the expressions  $\kappa_{ij} + \kappa_{ji}$  are positive.  $\square$

#### IV. FAST HSIC BASED ICA AND ITS IMPLEMENTATION

Having defined our independence criterion and its behaviour at independence, we now describe an efficient Newton-like method for minimising  $H(X)$ . We begin in Section IV-A with an overview of the method, including the approximate Hessian used to speed up computation. The subsequent two sections describe how the Hessian (Section IV-B), as well as  $H(X)$  (Section IV-C) and its gradient (Section IV-D), can be computed much faster using the incomplete Cholesky decomposition.

##### A. Fast HSIC Based ICA Method

We compute the gradient and Hessian of  $H \circ \mu_X$  in the parameter space  $\mathbb{R}^{m(m-1)/2}$ . By analogy with equation (14), the first derivative of  $H \circ \mu_X$  at  $0 \in \mathbb{R}^{m(m-1)/2}$  is

$$\begin{aligned} D(H \circ \mu_X)(0)\bar{\Omega} &= \frac{d}{d\varepsilon}(H \circ \mu_X)(\varepsilon\bar{\Omega})\Big|_{\varepsilon=0} \\ &= \sum_{1 \leq i < j \leq m}^m \omega_{ij}(\tau_{ij} - \tau_{ji}), \end{aligned} \quad (18)$$

where

$$\begin{aligned} \tau_{ij} &= \sum_{r=1; r \neq i}^m \mathbb{E}_{k,l} \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) x_j^\top \bar{w}_{kl} \phi \left( x_r^\top \bar{w}_{kl} \right) \right] \\ &+ \mathbb{E}_{k,l} \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) x_j^\top \bar{w}_{kl} \right] \mathbb{E}_{k,l} \left[ \phi \left( x_r^\top \bar{w}_{kl} \right) \right] \\ &- 2\mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi' \left( x_i^\top \bar{w}_{kl} \right) x_j^\top \bar{w}_{kl} \right] \mathbb{E}_l \left[ \phi \left( x_r^\top \bar{w}_{kl} \right) \right] \right]. \end{aligned} \quad (19)$$

Let  $Q := (q_{ij})_{i,j=1}^m \in \mathfrak{so}(m)$  and  $\nabla(H \circ \mu_X)(0) = (q_{ij})_{1 \leq i < j \leq m} \in \mathbb{R}^{m(m-1)/2}$ . Then the Euclidean gradient  $\nabla(H \circ \mu_X)(0)$  is given entry-wise as  $q_{ij} = \tau_{ij} - \tau_{ji}$  for all  $1 \leq i < j \leq m$ .

Similarly, the Hessian of the contrast function  $H$  can be computed directly from (15). The diagonal property of the Hessian does not hold true for arbitrary  $X \in O(m)$ , however, and it is clearly too expensive to compute the true Hessian at each step of our optimization. Nevertheless, since the Hessian is diagonal at  $X^*$  (Eq. (16)), a diagonal approximation of the Hessian makes sense in a neighborhood of  $X^*$ . Thus, we propose a diagonal matrix for the Hessian of  $H \circ \mu_X$  for an arbitrary  $X \in O(m)$  at  $0 \in \mathbb{R}^{m(m-1)/2}$ ; the result is a symmetric bilinear form  $\mathcal{H}(H \circ \mu_X)(0): \mathbb{R}^{m(m-1)/2} \times \mathbb{R}^{m(m-1)/2} \rightarrow \mathbb{R}$ . Due to the smoothness of both the contrast function  $H$  in (7) and the map  $\tilde{\mu}$  on  $O(m)$  in (11), the Hessian  $\mathcal{H}(H \circ \mu_X)(0)$  is smooth in an open neighborhood  $\mathcal{U}(X^*) \subset O(m)$  of  $X^* \in O(m)$ . Replacing the correct unmixing components  $s$  in (17) by the current estimates  $y = X^\top w$ , i.e.  $\bar{y}_{kli} = x_i^\top \bar{w}_{kl}$ , gives

$$\mathcal{H}(H \circ \mu_X)(0)(\bar{\Omega}, \bar{\Omega}) \approx \sum_{1 \leq i < j \leq m} \omega_{ij}^2 (\tilde{\kappa}_{ij} + \tilde{\kappa}_{ji}), \quad (20)$$

where

$$\tilde{\kappa}_{ij} = 2\mathbb{E}_k[\mathbb{E}_l[\bar{y}_{kli}]\mathbb{E}_l[\phi'(\bar{y}_{kli})]]\mathbb{E}_k[\mathbb{E}_l[\bar{y}_{klj}]\mathbb{E}_l[\phi'(\bar{y}_{klj})]] \quad (21a)$$

$$+ \mathbb{E}_{k,l}[\phi''(\bar{y}_{kli})]\mathbb{E}_{k,l}[(\bar{y}_{klj})^2\phi(\bar{y}_{klj})] \quad (21b)$$

$$+ 2\mathbb{E}_{k,l}[\phi''(\bar{y}_{kli})]\mathbb{E}_{k,l}[\phi(\bar{y}_{klj})] \quad (21c)$$

$$- 2\mathbb{E}_{k,l}[\phi''(\bar{y}_{kli})]\mathbb{E}_k[\mathbb{E}_l[(\bar{y}_{klj})^2]\mathbb{E}_l[\phi(\bar{y}_{klj})]] \quad (21d)$$

$$- \mathbb{E}_{k,l}[\phi'(\bar{y}_{kli})\bar{y}_{kli}]\mathbb{E}_{k,l}[\phi'(\bar{y}_{klj})\bar{y}_{klj}]. \quad (21e)$$

We emphasise that (by Remark 1) the Hessian at a correct unmixing matrix  $X^*$  is positive definite (i.e. the terms  $\tilde{\kappa}_{ij} + \tilde{\kappa}_{ji}$  are positive). Thus, the approximate Newton-like direction  $\tilde{\Omega} \in \mathbb{R}^{m(m-1)/2}$  with

$$\tilde{\omega}_{ij} = \frac{\tau_{ij} - \tau_{ji}}{\tilde{\kappa}_{ij} + \tilde{\kappa}_{ji}} \quad (22)$$

for  $1 \leq i < j \leq m$  is smooth, and is well defined within  $\mathcal{U}(X^*)$ . We call the Newton-like method arising from this approximation Fast HSIC-based Kernel ICA (FastKICA).

Although the approximation (20) can differ substantially from the true Hessian at an arbitrary  $X \in O(m)$ , they coincide at  $X^*$ . Thus, while the diagonal approximation is exact at the correct solution, it becomes less accurate as we move away from this solution. To ensure FastKICA is nonetheless well behaved as the global solution is approached, we provide the following local convergence result (we investigate the performance of our algorithm given arbitrary initialisation in our numerical experiments: see Section V).

*Corollary 1:* Let  $X^* \in O(m)$  be a correct unmixing matrix. Then FastKICA is locally quadratically convergent to  $X^*$ .

*Proof:* By considering the approximate Newton-like direction  $\tilde{\Omega}$  in (22) as  $\tilde{\Omega}: \mathcal{U}(X^*) \rightarrow \mathfrak{so}(m)$ , each iteration of FastKICA can be written as the map

$$\mathcal{A}: \mathcal{U}(X^*) \subset O(m) \rightarrow O(m), \quad X \mapsto X \exp(\tilde{\Omega}(X)). \quad (23)$$

A tedious but direct computation shows that  $\tilde{\Omega}(X^*) = 0$ , and  $X^*$  is a fixed point of  $\mathcal{A}$ . To prove the local quadratic convergence of FastKICA, we employ Lemma 2.9 in [36], bearing in mind that the required smoothness conditions (that the function be at least twice differentiable) are fulfilled by  $H(X)$ . According to this lemma, we only need to show that the first derivative of  $\mathcal{A}$ ,

$$D\mathcal{A}: T_X O(m) \rightarrow T_{\mathcal{A}(X)} O(m), \quad (24)$$

vanishes at a fixed point  $X^*$ . Thus we compute directly

$$D\mathcal{A}(X^*)(X^*\Omega) = X^*\Omega + X^* D\tilde{\Omega}(X^*)(X^*\Omega). \quad (25)$$

For the expression in (25) to vanish is equivalent to

$$\Omega = -D\tilde{\Omega}(X^*)(X^*\Omega), \quad (26)$$

i.e., for all  $1 \leq i < j \leq m$ ,

$$D\tilde{\omega}_{ij}(X^*)(X^*\Omega) = -\omega_{ij}. \quad (27)$$

Denoting the numerator  $\tau_{ij} - \tau_{ji}$  in (22) by  $\tilde{\omega}_{ij}^{(n)}(X)$  and the denominator  $\tilde{\kappa}_{ij} + \tilde{\kappa}_{ji}$  by  $\tilde{\omega}_{ij}^{(d)}(X)$ , we compute

$$\begin{aligned} D\tilde{\omega}_{ij}(X^*)(X^*\Omega) &= -\frac{\tilde{\omega}_{ij}^{(n)}(X^*)}{(\tilde{\omega}_{ij}^{(d)}(X^*))^2} D\tilde{\omega}_{ij}^{(d)}(X^*)(X^*\Omega) \\ &\quad + \frac{D\tilde{\omega}_{ij}^{(n)}(X^*)(X^*\Omega)}{\tilde{\omega}_{ij}^{(d)}(X^*)}. \end{aligned} \quad (28)$$

It can be shown that the first summand in (28) is equal to zero. Finally, following an argument almost identical to that in Appendix A,

$$D\tilde{\omega}_{ij}^{(n)}(X^*)(X^*\Omega) = -\tilde{\omega}_{ij}^{(d)}(X^*)\omega_{ij}. \quad (29)$$

Thus, we conclude condition (26) holds true at  $X^*$ , i.e.,  $D\mathcal{A}(X^*)$  vanishes as required. The result follows. ■

We emphasise that the local convergence result is proved in the population setting. In the finite sample case, this theoretical convergence rate might not be achievable [37]. For this reason, we experimentally

compare the convergence behavior of FastKICA with a gradient-based algorithm (see Section V-A). We will see that FastKICA converges significantly faster than gradient descent in practice.

As an additional caveat to Corollary 1, there is in general no practical strategy to guarantee that our ICA algorithm will be initialised within the neighbourhood  $\mathcal{U}(X^*)$ . Nevertheless, our experiments in Section V demonstrate that FastKICA often converges to the global solution even from arbitrary initialisation points, and is much more reliable in this respect than simple gradient descent.

### B. Incomplete Cholesky Estimate of the Hessian

In the following, we present an explicit formulation of the Hessian, and derive its approximation using an incomplete Cholesky decomposition of the kernel matrices. Sections IV-C and IV-D provide implementations of the contrast and gradient, respectively, that re-use the factors of the Cholesky decomposition.

First, we rewrite the pairwise HSIC terms (7) in a more convenient matrix form. Let  $K_i$  denote the kernel matrix for the  $i$ -th estimated source, i.e., its  $(k, l)$ -th entry is  $\phi(\bar{y}_{kli}) = \phi(x_i^\top \bar{w}_{kl})$ . We further denote by  $M \in \mathbb{R}^{n \times n}$  the centring operation,  $M = I - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$ , with  $\mathbf{1}_n$  being an  $n \times 1$  vector of ones.

*Lemma 2 (HSIC in terms of kernels [18]):* An empirical estimate of HSIC for two estimated sources  $y_i, y_j$  is

$$H_{uv}(X) := \frac{1}{(n-1)^2} \text{tr}(MK_i MK_j). \quad (30)$$

This empirical estimate is biased, however the bias in this estimate decreases as  $1/n$ , and thus drops faster than the variance (which decreases as  $1/\sqrt{n}$ : see [18]).

The first and second derivatives of the Gaussian kernel function  $\phi$  in (6) are

$$\begin{cases} \phi'(a, b) = -\frac{a-b}{\lambda^2} \phi(a, b), \\ \phi''(a, b) = \frac{(a-b)^2}{\lambda^4} \phi(a, b) - \frac{1}{\lambda^2} \phi(a, b). \end{cases} \quad (31)$$

Substituting the above terms into the approximate Hessian of  $H \circ \mu_X$  at  $0 \in \mathfrak{so}(m)$ , as computed in (20), yields

$$\mathcal{H}(H \circ \mu_X)(0)(\Omega, \Omega) \approx \sum_{1 \leq i < j \leq m} \omega_{ij}^2 \Delta_{ij}, \quad (32)$$

where

$$\Delta_{ij} = \frac{2}{\lambda^2} (\beta_i \zeta_j + \zeta_i \beta_j) + \frac{4}{\lambda^4} (\zeta_i \zeta_j - \eta_i \eta_j), \quad (33)$$

and

$$\begin{cases} \beta_i = \mathbb{E}_{k,l}[\phi(\bar{y}_{kli})], \\ \zeta_i = \mathbb{E}_{k,l}[\phi(\bar{y}_{kli}) y_{ki} y_{li}], \\ \eta_i = \mathbb{E}_{k,l}[\phi(\bar{y}_{kli}) y_{ki}^2]. \end{cases} \quad (34)$$

Here,  $y_{ki} = e_i^\top y_k$  is the  $i$ -th entry of the  $k$ -th sample of the estimates  $Y$ , selected by the  $i$ -th standard basis vector  $e_i$  of  $\mathbb{R}^m$ .

We now outline how the incomplete Cholesky decomposition [30] helps to estimate the approximate Hessian efficiently. An incomplete Cholesky decomposition of the Gram matrix  $K_i$  yields a low-rank approximation  $K_i \approx G_i G_i^\top$  that greedily minimises  $\text{tr}(K_i - G_i G_i^\top)$ . The cost of computing the  $n \times d$  matrix  $G_i$  is  $O(nd^2)$ , with  $d \ll n$ . Greater values of  $d$  result in a more accurate reconstruction of  $K_i$ . As pointed out in [16, Appendix C], however, the spectrum of a Gram matrix based on the Gaussian kernel generally decays rapidly, and a small  $d$  yields a very good approximation (the experiments in [23] provide further empirical evidence). With approximate Gram matrices, the empirical estimates of the three terms in Equation (34) become

$$\begin{aligned}\hat{\beta}_i &= \frac{1}{n^2} (\mathbf{1}_n^\top G_i) (\mathbf{1}_n^\top G_i)^\top, \\ \hat{\zeta}_i &= \frac{1}{n^2} (y_i^\top G_i) (y_i^\top G_i)^\top, \\ \hat{\eta}_i &= \frac{1}{n^2} ((y_i \odot y_i)^\top G_i) (G_i^\top \mathbf{1}_n),\end{aligned}$$

where  $y_i$  is the sample vector for the  $i$ th estimated source, and  $\odot$  the entry-wise product of vectors.

### C. Incomplete Cholesky Estimate of HSIC

Lemma 2 states an estimate of the pairwise HSIC as the trace of a product of centred kernel matrices  $MK_iM$ ,  $MK_jM$ . Reusing the Cholesky decomposition from above, i.e.  $MK_iM = MG_iG_i^\top M =: \tilde{G}_i\tilde{G}_i^\top$ , where  $\tilde{G}_i$  is  $n \times d_i$ , we arrive at an equivalent trace of a much smaller  $d_j \times d_j$  matrix for each pair of estimated sources  $(y_i, y_j)$ , and avoid any product of  $n \times n$  matrices:

$$\begin{aligned}\hat{H}(X) &= \frac{1}{(n-1)^2} \sum_{1 \leq i < j \leq m} \text{tr} \left( \tilde{G}_i \tilde{G}_i^\top \tilde{G}_j \tilde{G}_j^\top \right) \\ &= \frac{1}{(n-1)^2} \sum_{1 \leq i < j \leq m} \text{tr} \left( \left( \tilde{G}_j^\top \tilde{G}_i \right) \left( \tilde{G}_i^\top \tilde{G}_j \right) \right).\end{aligned}$$

### D. Incomplete Cholesky Estimate of the Gradient

We next provide an approximation of the gradient in terms of the same Cholesky approximation  $K = GG^\top$ . Williams and Seeger [38] describe an approximation of  $K$  based on an index set  $I$  of length  $d$  with unique entries from  $\{1, \dots, n\}$ :

$$K \approx K' := K_{:,I} K_{I,I}^{-1} K_{I,:}, \quad (35)$$



where  $K_{:,I}$  is the Gram matrix with the rows unchanged, and the columns chosen from the set  $I$ ; and  $K_{I,I}$  is the  $d \times d$  submatrix with both rows and columns restricted to  $I$ . We adapt the approach of Fine and Scheinberg [30] and choose the indices  $I$  in accordance with an incomplete Cholesky decomposition to minimise  $\text{tr}(K - K')$ .

For simplicity of notation, we will restrict ourselves to the gradient of HSIC for a pair  $(y_i, y_j)$  and ignore the normalisation by  $(n-1)^{-2}$ . The gradient of  $\widehat{H}(X)$  is merely the sum of all pairwise gradients. Let  $K$  be the Gram matrix for  $y_i$  and  $L$  for  $y_j$ . With centring matrices  $\widetilde{K} = MKM$  and  $\widetilde{L} = MLM$  the (unnormalised) differential of the pairwise HSIC is [23]

$$\begin{aligned}
d\text{tr}(\widetilde{K}\widetilde{L}') &= \text{tr}(\widetilde{K}'d(L')) + \text{tr}(\widetilde{L}'d(K')) \\
&= \text{tr}(\widetilde{K}'d(L')) + \text{tr}(\widetilde{L}'d(K_{:,I}K_{I,I}^{-1}K_{I,:})) \\
&= \text{tr}(\widetilde{K}'d(L')) \\
&\quad + 2\text{vec}(\widetilde{L}'K_{:,I}K_{I,I}^{-1})^\top \text{vec}(dK_{:,I}) \\
&\quad - \text{vec}(K_{I,I}^{-1}K_{I,:}\widetilde{L}'K_{:,I}K_{I,I}^{-1})^\top \text{vec}(dK_{I,I}).
\end{aligned} \tag{36}$$

Our expression for  $d(K')$  is derived in Appendix B-A. The expansion of  $dL'$  is analogous. The matrix decompositions shrink the size of the factors in (36). An appropriate ordering of the matrix products allows us to avoid ever generating or multiplying an  $n \times n$  matrix. In addition, note that for a column  $x$  of  $X$ ,

$$\text{vec}(dK_{:,I}) = d\text{vec}(K_{:,I}) = \left( \partial \text{vec}(K_{:,I}) / \partial x^\top \right) d\text{vec}(x).$$

The partial matrix derivative  $\partial \text{vec}(K_{:,I}) / \partial x^\top$  is defined in Appendix B-B and has size  $nd \times m$ , whereas the derivative of the full  $K$  has  $n^2 \times m$  entries. Likewise,  $\partial \text{vec}(K_{I,I}) / \partial x^\top$  is only  $d^2 \times m$ . We must also account for the rapid decay of the spectrum of Gram matrices with Gaussian kernels [16, discussion in Appendix C], since this can cause the inverse of  $K_{I,I}$  to be ill-conditioned. We therefore add a small ridge of  $10^{-6}$  to  $K_{I,I}$ , although we emphasise that our algorithm is insensitive to this value.

We end this section with a brief note on the overall computational cost of FastKICA. As discussed in [23, Section 1.5], the gradient and Hessian are computable in  $\mathcal{O}(nm^3d^2)$  operations. A more detailed breakdown of how we arrive at this cost may be found in [23], bearing in mind that the Hessian has the same cost as the gradient thanks to our diagonal approximation.

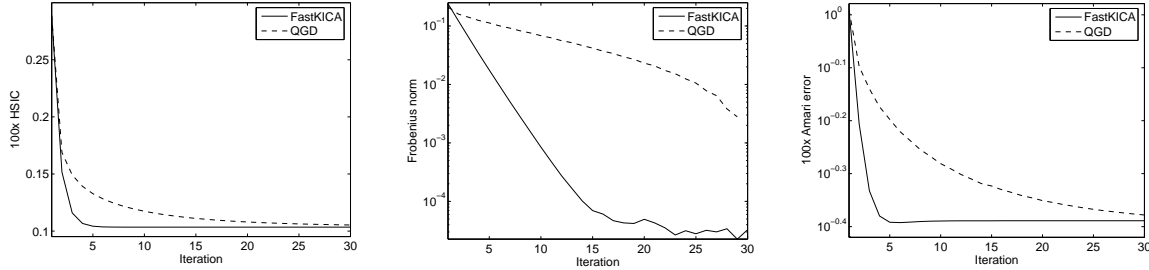


Fig. 2. Convergence measured by HSIC, by the Frobenius norm of the difference between the  $i$ -th iterate  $X^{(i)}$  and  $X^{(30)}$ , and by the Amari error. FastKICA converges faster. The plots show averages over 25 runs with 40,000 samples from 16 sources.

## V. NUMERICAL EXPERIMENTS

In our experiments, we demonstrate four main points: First, if no alternative algorithm is used to provide an initial estimate of  $X$ , FastKICA is resistant to local minima, and often converges to the correct solution. This is by contrast with gradient descent, which is more often sidetracked to local minima. In particular, if we choose sources incompatible with the initialising algorithm (so that it fails completely), our method can nonetheless find a good solution.<sup>7</sup> Second, when a good initial point is given, the Newton-like algorithm converges faster than gradient descent. Third, our approach runs sufficiently quickly on large-scale problems to be used either as a standalone method (when a good initialisation is impossible or unlikely), or to fine tune the solution obtained by another method. While not the fastest method tested, demixing performance of FastKICA achieves a reasonable compromise between speed and solution quality, as demonstrated by its performance as the departure of the mixture from independence becomes non-smooth. Finally, FastKICA shows better resistance to outliers than alternative approaches.

Our artificial data for Sections V-A (comparison of FastKICA with gradient descent) and V-B (computational cost benchmarks) were generated in accordance with [17, Table 3], which is similar to the artificial benchmark data of [16]. Each source was chosen randomly with replacement from 18 different distributions having a wide variety of statistical properties and kurtoses. Sources were mixed using a random matrix with condition number between one and two. Section V-C describes our experiments on source smoothness vs performance, and Section V-D contains our experiments on outlier resistance. We used the Amari divergence, defined by [39], as an index of ICA algorithm performance (we multiplied

<sup>7</sup>Note the criterion optimised by FastKICA is also the statistic of an independence test [10], [22]. This test can be applied directly to the values of HSIC between pairs of unmixed sources, to verify the recovered signals are truly independent; no separate hypothesis test is required.

this quantity by 100 to make the performance figures more readable). In all experiments, the precision of the incomplete Cholesky decomposition was  $10^{-6}n$ . Convergence was measured by the difference in HSIC values over consecutive iterations.

### A. Comparison with Gradient Descent

We first compare the convergence of FastKICA with a simple gradient descent method [23]. In order to find a suitable step width along the gradient mapped to  $O(m)$ , the latter uses a quadratic interpolation of HSIC along the geodesic. This requires HSIC to be evaluated at two additional points. Both FastKICA and quadratic gradient descent (QGD) use the same gradient and independence measure. Figure 2 compares the convergence for both methods (on the same data) of HSIC, the Amari error, and the Frobenius norm of the difference between the  $i$ -th iterate  $X^{(i)}$  and the solution  $X^{(30)}$  reached after 30 iterations. The results are averaged over 25 runs. In each run, 40,000 observations from 16 artificial, randomly drawn sources were generated and mixed. We initialised both methods with FastICA [2], and used a kernel width of  $\lambda = 0.5$ . As illustrated by the plots, FastKICA approaches the solution much faster than QGD. We also observe that the number of iterations to convergence decreases when the sample size grows. The plot of the Frobenius norm suggests that convergence of FastKICA is only linear for the artificial data set. While local quadratic convergence is guaranteed in the population setting, the required properties for Corollary 1 do not hold exactly in the finite sample setting, which can reduce the convergence rate [37].

For arbitrary initialisations, FastKICA is still applicable with multiple restarts, although a larger kernel width is more appropriate for the initial stages of the search (local fluctuations in FastKICA far from independence are then smoothed out, although the bias in the location of the global minimum increases). We set  $\lambda = 1.0$  and a convergence threshold of  $10^{-8}$  for both FastKICA and QGD. For 40,000 samples from 8 artificial sources, FastKICA converged on average for 37% of the random restarts with an average error ( $\times 100$ ) of  $0.54 \pm 0.01$ , whereas the QGD did not yield any useful results at all (mean error  $\times 100$ :  $74.14 \pm 1.39$ ). Here, averages are over 10 runs with 20 random initialisations each. The solution obtained with FastKICA can be refined further by shrinking the kernel width after initial convergence, to reduce the bias.

### B. Performance and Cost vs Other Approaches

Our next results compare the performance and computational cost of FastKICA, Jade [4], KDICA [8], MICA [31], MILCA [6], RADICAL [5], and quadratic gradient descent (QGD) [23]. The timing experiments for all methods except MILCA were performed on the same dual AMD Opteron (2x AMD

Opteron(tm) Processor 250, 64KiB L1 cache, 1MiB L2 cache, GiB System Memory). Since MILCA was much slower, its tests were run in parallel on 64 bit cluster nodes with 2–16 processors and 7.8–94.6 GB RAM, running Ubuntu 7.04: these nodes were generally faster than the one used for the other algorithms, and the runtime of MILCA is consequently an underestimate, relative to the remaining methods.

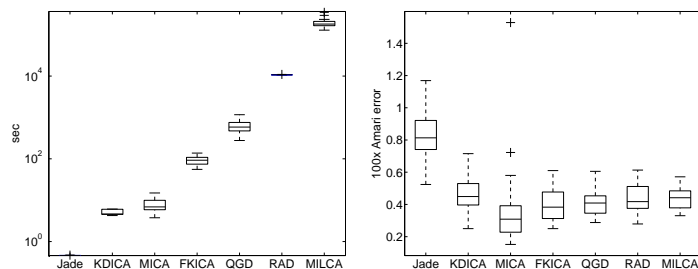
We demixed 8 sources and 40,000 observations of the artificial data. The run times include the initialisation by Jade for FastKICA, QGD, and KDICA, and are averaged over 10 repetitions for each data set (except for the slower methods RADICAL and MILCA). QGD was run for 10 iterations, and the convergence threshold for FastKICA was  $10^{-5}$  ( $\lambda = 0.5$ ). MICA was allowed a maximum number of 50 iterations, since with the default 10 iterations, the Amari error was more than twice as high ( $0.84 \pm 0.74$ ), and worse than that all other algorithms besides Jade.

Figure 3 displays the error and time for 24 data sets. We see that demixing performance is very similar across all nonparametric approaches, which perform well on this benchmark. MICA has the best median performance, albeit with two more severely misconverged solutions. While small, the median performance difference is statistically significant according to a level 0.05 sign test. FastKICA and QGD provide the next best result, and exhibit a small but statistically significant performance difference compared with KDICA, RADICAL and Jade, but not MILCA. The time differences between the various algorithms are much larger than their performance differences. In this case, the ordering is Jade, KDICA, MICA, FastKICA, QGD, RADICAL, and MILCA. The additional evaluations of HSIC for the quadratic approximation make QGD slower per iteration than FastKICA. As shown above, FastKICA also converges in fewer iterations than QGD, requiring 4.32 iterations on average.

We also compared KDICA and FastKICA when random initialisations were used. We see in Figure 5(a) that FastKICA solutions have a clear bivariate distribution, with a large number of initialisations reaching an identical global minimum: indeed, the correct solution is clearly distinguishable from local optima on the basis of its HSIC value. By contrast, KDICA appears to halt at a much wider variety of local minima even for these relatively simple data, as evidenced by the broad range of Amari errors in the estimated unmixing matrices. Thus, in the absence of a good initialising estimate (where classical methods fail), FastKICA is to be preferred. We will further investigate misconvergence behaviour of the different ICA algorithms in the next section, for a more difficult (non-smooth) demixing problem.

### *C. ICA performance as a function of problem smoothness*

While the foregoing experiments provide a good idea of computational cost, they do not address the performance of the various ICA methods as a function of the statistical properties of the sources (rather,



	Error $\times 100$	Time [s]
Jade	$0.83 \pm 0.17$	$0.46 \pm 0.002$
KDICA	$0.47 \pm 0.12$	$5.09 \pm 0.8$
MICA	$0.37 \pm 0.28$	$8.10 \pm 3.4$
FKICA	$0.40 \pm 0.10$	$91.27 \pm 22.9$
QGD	$0.41 \pm 0.08$	$618.17 \pm 208.7$
RAD	$0.44 \pm 0.09$	$10.72 \cdot 10^3 \pm 50.38$
MILCA	$0.44 \pm 0.06$	$18.97 \cdot 10^4 \pm 5.1 \cdot 10^4$

Fig. 3. Comparison of run times (left) and performance (middle) for various ICA algorithms. FastKICA is faster than MILCA, RADICAL, and gradient descent with quadratic approximation, and its results compare favorably to the other methods. KDICA is even faster, but performs less well than FastKICA. Both KDICA and MICA have higher variance than FastKICA. The values are averages over 24 data sets.

the performance is an average over sources with a broad variety of behaviours, and is very similar across the various nonparametric approaches). In the present section, we focus specifically on how well the ICA algorithms perform as a function of the smoothness of the ICA problem. Our source distributions take the form of Gaussians with sinusoidal perturbations, and are proportional to  $g(x)(1 + \sin(2\pi\nu\beta x))$ , where  $g(x)$  is a Gaussian probability density function with unit variance,  $\beta$  is the maximum perturbing frequency considered and  $\nu$  is a scaling factor ranging from 0.05 to 1.05 with spacing 0.05 (the choice of  $\beta$  will be addressed later). The sinusoidal perturbation is multiplied by  $g(x)$  to ensure the resulting density expression is everywhere non-negative. Plots of the source probability density function and its characteristic function are given in Figure 4. Bearing in mind that purely Gaussian sources can only be resolved up to rotation [1], the resulting ICA problem becomes more difficult (for algorithms making smoothness assumptions on the sources) as the departure from Gaussianity is encoded at increasing frequencies, which are harder to distinguish from random noise for a given sample size.<sup>8</sup> By contrast, the sources used in the previous section (taken from [17, Table 3]) yield very similar demixing performance when comparing across the nonparametric algorithms used in our benchmarks. One reason for this similarity is that the departure from Gaussianity of these sources has substantial amplitude at low

<sup>8</sup>This perturbation to the Gaussian distribution differs from that used for designing contrast functions in classical ICA studies, which employ Edgeworth [1] or Gram-Charlier [39] expansions. We shall see, however, that our perturbed sources are better suited to characterizing the interaction between ICA performance and the choice of kernel, for methods using kernel density estimates (MICA, and KDICA) and for FastKICA.

frequencies, resulting in ICA problems of similar difficulty for MICA, KDICA, and FastKICA.<sup>9</sup> We remark that linear mixtures where the departure from independence occurs only at high frequencies are not typical of real-life ICA problems. That said, such mixtures represent an important failure mode of ICA algorithms that make smoothness assumptions on the source densities (as for MICA, KDICA, and FastKICA). Thus, our purpose in this section is to compare the decay in unmixing performance across the various ICA algorithms as this failure mode is approached.

We decide on the base frequency  $\beta$  of the perturbation with reference to the parameters of MICA, to simplify the discussion of performance. The MICA algorithm optimizes a sum of entropies of each mixture component, where the entropies are computed using discretised empirical estimates of the mixture probability distributions. We can express the distribution estimates by first convolving the mixture sample by a B-spline kernel (of order 3, although other orders are also possible) and then downsampling to get probability estimates at the gridpoints. If we consider the population setting and perform these operations in the frequency domain, this corresponds to multiplying the Fourier transform of the source density by that of the B-spline, and then aliasing the frequency components that exceed the grid Nyquist frequency.

Given a baseline bandwidth  $b_0$ , the grid spacing is computed as a function of the sample size  $n$  according to  $b_n = b_0 \times 2.107683/n^{0.2}$ ; without loss of generality, we set  $b_0 = 1$ . The spline kernel bandwidth is also scaled by this factor, such that the zeros in the kernel spectrum occur at integer multiples of  $f_m := n^{0.2}/(2.107683)$ . To use these factors in setting  $\beta$ , consider two sources consisting of perturbed Gaussians with identical  $\beta$ . The characteristic function of the two mixtures resulting from a rotation with angle  $\pi/4$  has a distinctive peak at  $\beta\sqrt{2}$  when  $\nu = 1$ . Thus, by setting  $\beta = n^{0.2}/(2.107683\sqrt{2})$ , this peak will fall at the first minimum of the spline kernel spectrum. An illustration is provided in Figure 4. Note in particular the decay of the spline spectrum towards its first zero at  $f_m$ : any component of the mixture characteristic function at this frequency will be severely attenuated, and thus we expect that demixing sources with  $\nu$  approaching 1 will be difficult (the sources will appear Gaussian).

We sampled 25 data sets consisting of two sources with  $n = 1,000$  for each value of  $\nu$ , and mixed the sources with orthonormal matrices. To ensure comparable results, we used the same set of 25 mixing matrices across all  $\nu$ . The algorithms were run for a maximum of 50 iterations. The convergence threshold for FastKICA was a 0.5% change of HSIC, and the bandwidth  $\lambda = 0.5$ . MICA used a bandwidth of  $b_0 = 1$ . For these bandwidth choices, we emphasise that both FKICA and MICA reached chance level

<sup>9</sup>As we shall see, the spacings-based entropy estimates of RADICAL and the graph-based mutual information estimates of MILCA behave quite differently.

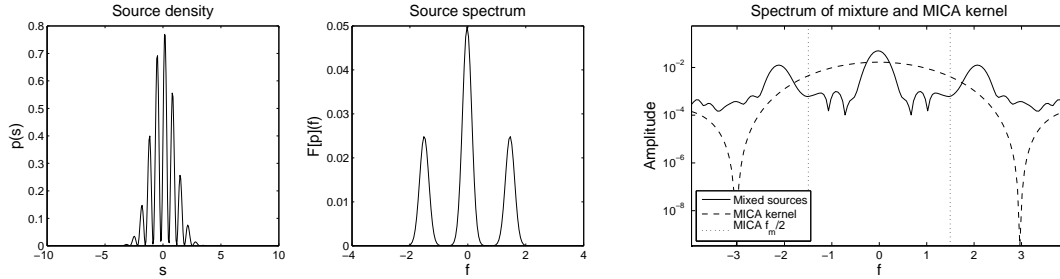


Fig. 4. **Left:** Source probability density function for a perturbation at frequency  $0.7f_m/\sqrt{2}$ , where  $f_m$  is the frequency of the first zero in the spectrum of the spline kernel used by MICA. **Middle:** Characteristic function of the source, showing peaks at the perturbation frequency. **Right:** Empirical (smoothed) characteristic function of the mixture of two sources with angle  $\pi/4$ . Two peaks are seen at locations  $0.7f_m$ . The spectrum of the MICA kernel (a 3rd order B-spline) is superposed. The dashed vertical lines are at  $f_m/2$ , which is the Nyquist frequency for the grid used by MICA. Thus, perturbations exceeding this frequency will be aliased.

performance (i.e. complete failure) at the same source perturbing frequency, corresponding to  $\nu = 1$ , making the behaviour of these two methods across the  $\nu$  range directly comparable. In other words, we report the relative performance of the two algorithms over the frequency range for which they operate at better than chance level. The Amari errors in Figure 5(b) are averages over 10 random initialisations and the 25 data sets for each frequency. Each algorithm was initialised with the same 10 random orthonormal matrices.

We note first of all that FastKICA has a longer  $\nu$  interval in which the average Amari error is very low, compared with MICA and KDICA. In addition, as  $\nu$  rises above 0.5, the average error of FastKICA is consistently below that of MICA and KDICA. On the other hand, for the lowest perturbing frequencies, KDICA and MICA perform better than FastKICA. The two most computationally costly methods, RADICAL and MILCA, perform best, with a low Amari divergence over all the high  $\nu$  values tested. This is as expected, for two reasons: first, both methods perform an exhaustive search over all pairs of Jacobi angles, and are not susceptible to local minima. Second, RADICAL is based on a spacings estimate of entropy, and MILCA on a  $k$ -nn estimate of the mutual information: in other words, both methods adapt automatically to the scale of the variations in the mixture densities. That said, efficient optimization techniques have yet to be developed for RADICAL and MILCA.

We next examine in more detail the convergence behaviour leading to the drop in average performance of FastKICA, MICA, and KDICA as  $\nu$  rises. First, as noted in the previous section, KDICA can be sensitive to local minima: thus its average performance degrades even for low values of  $\nu$  as a large number

of initialisations result in misconvergence. The behaviour of MICA is more complex, with performance dropping at  $\nu \approx 0.4$  but recovering for  $\nu \approx 0.55$  (two more such oscillations occur at higher  $\nu$ ). At  $\nu \approx 0.4$ , the MICA entropy score develops a local minimum at a rotation of  $\pi/4$  from the true unmixing matrix, resulting in a substantial number of initializations converging to this incorrect solution, as well as a group of correct solutions (this local minimum is also seen for other  $b_0$  values, but at different onset values of  $\nu$ ). The local minimum becomes less pronounced at  $\nu \approx 0.55$ , but then strengthens again at  $\nu \approx 0.65$ . By contrast, the results for FastKICA at moderate values of  $\nu$  more closely follow the histogram of Figure 5(a), with a large fraction of solutions at the global optimum, and the remaining misconverged solutions having a range of Amari errors. Taking the best solution over all 10 initializations (as measured by the HSIC or entropy score), rather than the average solution, the results of MICA and FastKICA at larger  $\nu$  both remain indistinguishable from RADICAL and MILCA until  $\nu \approx 0.8$ . For  $\nu > 0.8$ , performance worsens towards chance level as  $\nu$  rises to 1, at which point the global optimum of both contrast functions occurs at a random angle. The onset of this performance drop can be increased for both FastKICA and MICA by decreasing  $b_0$  or  $\lambda$ , respectively; but at a cost of worse mean performance due to more pronounced local minima. On the other hand, the best KDICA result continues to perform as well as RADICAL and MILCA, since the slow decaying Fourier transform of its Laplace kernel makes it sensitive to higher frequencies.

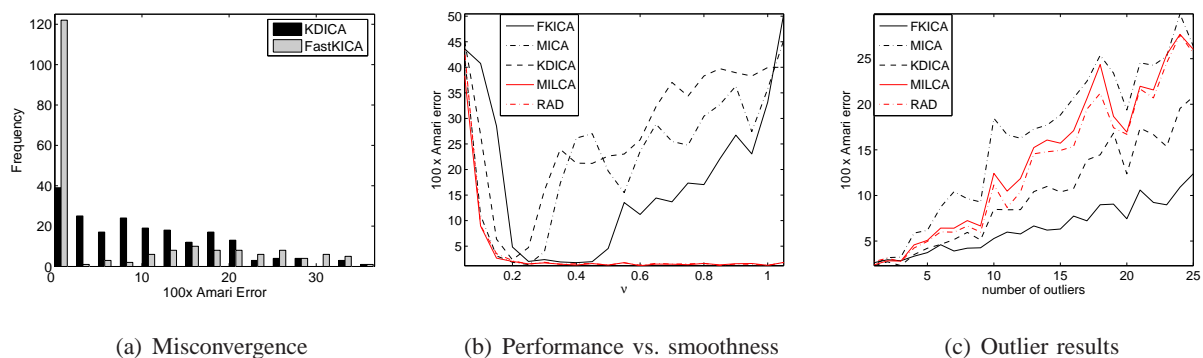


Fig. 5. (a) Comparison of performance for arbitrary initialisations ( $n = 40,000$ ,  $m = 8$ ). Amari error histograms are shown for FastKICA vs. KDICA with mixed artificial sources (10 data sets, 20 initialisations each). FastKICA reaches a global minimum far more often than KDICA. (b) Amari error depending on the sine frequencies for  $n = 1000$  samples and two sources. (c) Effect of outliers on the performance of the ICA algorithms, for two sources of length  $n = 1000$ , drawn independently with replacement from [17, Table 3], and corrupted at random observations with outliers at  $\pm 5$  (where each sign has probability 0.5). Each point represents an average over 100 independent experiments. The number of corrupted observations in both signals is given on the horizontal axis.



#### D. Resistance to Outliers

In our final experiment, we investigate the effect of outlier noise added to the observations. We selected two generating distributions from the benchmark data in [17, Table 3], randomly and with replacement. After combining these signals with a randomly generated matrix with condition number between 1 and 2, we generated a varying number of outliers by adding  $\pm 5$  (with equal probability) to both signals at random locations. We evaluated HSIC using a Gaussian kernel of size  $\lambda = 0.5$ . For FastKICA and MICA, we chose the best result out of 3 random initialisations, according to HSIC or the estimated entropy, respectively. The initialisation for KDICA was the FastKICA result, since the KDICA is sensitive to poor initialisation (as seen in Section V-B). Results are shown in Figure 5(c). It is clear that FastKICA substantially outperforms the alternatives in outlier resistance.

## VI. CONCLUSION

We demonstrate that an approximate Newton-like method, FastKICA, can improve the speed and performance of kernel/characteristic function-based ICA methods. We emphasise that FastKICA is applicable even if no good initialisation is at hand. With a modest number of restarts and a kernel width that shrinks near independence (on our data, from  $\lambda = 1.0$  to  $\lambda = 0.5$ ), the correct global optimum is consistently found. A good initialisation results in more rapid convergence, and we do not need to adapt the kernel size. Our method demonstrates much better outlier resistance than recently published competing approaches. Moreover, our optimization method can be applied to any twice differentiable RKHS kernel, rather than relying on the specific properties of particular kernels (to be Laplace in the case of [8], or to be a spline kernel with compact support in [31]).

Several directions for future work are suggested by the present study. First, the kernel bandwidth used is currently chosen heuristically. It would be of interest to develop more principled methods for choosing this bandwidth based on properties of the data. Second, it is notable that ICA methods based on spacings estimates of entropy, or nearest-neighbour estimates of mutual information, perform very well for ICA problems where the departure from independence is encoded at high frequencies. Unfortunately, efficient optimization techniques have yet to be developed for ICA using these dependence measures.

## APPENDIX A

### EVALUATION OF THE SECOND DERIVATIVE OF

## THE HSIC BASED ICA CONTRAST

We assume the unmixing matrix is correct, and thus  $X = X^*$ . The term (15a) for a fixed pair  $(i, j)$  can be computed as

$$\begin{aligned} & \mathbb{E}_{k,l} \left[ \phi''(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \bar{s}_{kl}^\top \omega_i \phi(\bar{s}_{klj}) \right] \\ &= \sum_{r,t=1;r,t \neq i}^m \omega_{ir} \omega_{it} \mathbb{E}_{k,l} \left[ \phi''(\bar{s}_{kli}) \bar{s}_{klr} \bar{s}_{klt} \phi(\bar{s}_{klj}) \right]. \end{aligned} \quad (37)$$

Under the assumptions of independence and whitened mixtures, the corresponding  $(r, t)$  expression can be written

$$\begin{cases} 0, & r \neq t; \\ 2\mathbb{E}_{k,l} [\phi''(\bar{s}_{kli}) \phi(\bar{s}_{klj})], & r = t \neq i, j; \\ \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_{k,l} [\bar{s}_{klj}^2 \phi(\bar{s}_{klj})], & r = t = j. \end{cases} \quad (38)$$

Thus the term (15a) can be further simplified as

$$\begin{aligned} (15a): & \mathbb{E}_{k,l} \left[ \phi''(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \bar{s}_{kl}^\top \omega_i \phi(\bar{s}_{klj}) \right] \\ &= \sum_{r=1;r \neq i,j}^m 2\omega_{ir}^2 \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})] \\ & \quad + \omega_{ij}^2 \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_{k,l} [\bar{s}_{klj}^2 \phi(\bar{s}_{klj})]. \end{aligned} \quad (39)$$

By applying the same techniques, the remaining terms (15b)–(15i) become

$$\begin{aligned} (15b): & \mathbb{E}_{k,l} \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \Omega \bar{s}_{kl} \phi(\bar{s}_{klj}) \right] \\ &= \sum_{r=1;r \neq i}^m -\omega_{ir}^2 \mathbb{E}_{k,l} [\phi'(\bar{s}_{kli}) \bar{s}_{kli}] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})], \end{aligned} \quad (40)$$

$$\begin{aligned} (15c): & \mathbb{E}_{k,l} \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \bar{s}_{kl}^\top \omega_j \phi'(\bar{s}_{klj}) \right] \\ &= -\omega_{ij}^2 \mathbb{E}_{k,l} [\phi'(\bar{s}_{kli}) \bar{s}_{kli}] \mathbb{E}_{k,l} [\phi'(\bar{s}_{klj}) \bar{s}_{klj}], \end{aligned} \quad (41)$$

$$\begin{aligned} (15d): & \mathbb{E}_{k,l} \left[ \phi''(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \bar{s}_{kl}^\top \omega_i \right] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})] \\ &= \sum_{r=1;r \neq i}^m 2\omega_{ir}^2 \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})], \end{aligned} \quad (42)$$

$$\begin{aligned} (15e): & \mathbb{E}_{k,l} \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \Omega \bar{s}_{kl} \right] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})] \\ &= \sum_{r=1;r \neq i}^m -\omega_{ir}^2 \mathbb{E}_{k,l} [\phi'(\bar{s}_{kli}) \bar{s}_{kli}] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})], \end{aligned} \quad (43)$$

$$(15f): \mathbb{E}_{k,l} \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \right] \mathbb{E}_{k,l} \left[ \phi'(\bar{s}_{klj}) \omega_j^\top \bar{s}_{kl} \right] = 0, \quad (44)$$

$$\begin{aligned}
(15g): & \mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi''(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \bar{s}_{kl}^\top \omega_i \right] \mathbb{E}_l [\phi(\bar{s}_{klj})] \right] \\
&= \sum_{r=1; r \neq i, j}^m 2\omega_{ir}^2 \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})] \\
&\quad + \omega_{ij}^2 \mathbb{E}_{k,l} [\phi''(\bar{s}_{kli})] \mathbb{E}_k [\mathbb{E}_l [\bar{s}_{klj}^2] \mathbb{E}_l [\phi(\bar{s}_{klj})]],
\end{aligned} \tag{45}$$

$$\begin{aligned}
(15h): & \mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \Omega \bar{s}_{kl} \right] \mathbb{E}_l [\phi(\bar{s}_{klj})] \right] \\
&= \sum_{r=1; r \neq i}^m -\omega_{ir}^2 \mathbb{E}_{k,l} [\phi'(\bar{s}_{kli}) \bar{s}_{kli}] \mathbb{E}_{k,l} [\phi(\bar{s}_{klj})],
\end{aligned} \tag{46}$$

$$\begin{aligned}
(15i): & \mathbb{E}_k \left[ \mathbb{E}_l \left[ \phi'(\bar{s}_{kli}) \omega_i^\top \bar{s}_{kl} \right] \mathbb{E}_l \left[ \phi'(\bar{s}_{klj}) \omega_j^\top \bar{s}_{kl} \right] \right] \\
&= -\omega_{ij}^2 \mathbb{E}_k [\mathbb{E}_l [\bar{s}_{kli}] \mathbb{E}_l [\phi'(\bar{s}_{kli})]] \mathbb{E}_k [\mathbb{E}_l [\bar{s}_{klj}] \mathbb{E}_l [\phi'(\bar{s}_{klj})]].
\end{aligned} \tag{47}$$

Substituting (39)–(47) into equation (15), the result in equation (16) follows directly from  $\omega_{ij} = -\omega_{ji}$ .

## APPENDIX B

### DERIVATION OF THE APPROXIMATE GRADIENT AND THE MATRIX PARTIAL DERIVATIVES

In this appendix, we first derive the differential  $dK'$  of the Cholesky approximation to the Gram matrix  $K$ , and use it to obtain the differential of the approximate HSIC in (36). We then give an expression for the differential of the factors of  $K'$ , which involves the entry-wise derivative of the Gram matrix with respect to a column  $x$  of  $X$ . Details have been published in [23].

#### A. Differential of the incomplete Cholesky approximation to HSIC

Recall that the differential of the low-rank approximation to HSIC is

$$d\text{tr}(\tilde{K}'\tilde{L}') = \text{tr}(\tilde{K}'d(L')) + \text{tr}(\tilde{L}'d(K')). \tag{48}$$

We expand the differential of  $K' = K_{:,I}K_{I,I}^{-1}K_{I,:}$  by the product rule and by rewriting  $d(K_{I,I}^{-1}) = K_{I,I}^{-1}(dK_{I,:})K_{I,I}^{-1}$ . Plugging the result into the second term  $\text{tr}(\tilde{L}'d(K'))$  from (48) yields

$$\begin{aligned}
& \text{tr}(\tilde{L}'d(K')) \\
&= \text{tr}(\tilde{L}'dK_{:,I}K_{I,I}^{-1}K_{I,:}) + \text{tr}(\tilde{L}'K_{:,I}K_{I,I}^{-1}dK_{I,:}) \\
&\quad - \text{tr}(\tilde{L}'K_{:,I}K_{I,I}^{-1}(dK_{I,I})K_{I,I}^{-1}K_{I,:}).
\end{aligned} \tag{49}$$

Using the symmetry of the Gram matrices, the second term on the right hand side can be transformed as

$$\begin{aligned}
\text{tr}(\tilde{L}'K_{:,I}K_{I,I}^{-1}dK_{I,:}) &= \text{tr}((dK_{:,I}K_{I,I}^{-1}K_{I,:})^\top \tilde{L}'^\top) \\
&= \text{tr}(\tilde{L}'dK_{:,I}K_{I,I}^{-1}K_{I,:}),
\end{aligned}$$

and (49) becomes

$$\begin{aligned} \text{tr} \left( \tilde{L}' d(K') \right) &= 2 \text{tr} \left( K_{I,I}^{-1} K_{I,:} \tilde{L}' dK_{:,I} \right) - \text{tr} \left( K_{I,I}^{-1} K_{I,:} \tilde{L}' K_{:,I} K_{I,I}^{-1} dK_{I,I} \right) \\ &= 2 \text{vec} \left( \tilde{L}' K_{:,I} K_{I,I}^{-1} \right)^\top \text{vec} \left( dK_{:,I} \right) \\ &\quad - \text{vec} \left( K_{I,I}^{-1} K_{I,:} \tilde{L}' K_{:,I} K_{I,I}^{-1} \right)^\top \text{vec} \left( dK_{I,I} \right). \end{aligned}$$

The other term  $\text{tr} \left( \tilde{K}' d(L') \right)$  in (48) is equivalent.

### B. Derivative of the Gram matrix with respect to $X$

The derivative of the Gram matrix entries with respect to a particular column  $x$  of the unmixing matrix  $X$  depends on the kernel. We employ a Gaussian kernel here, but one could easily obtain the derivatives of additional kernels: these can then be plugged straightforwardly into the equations in the previous section.

*Lemma 3 (Derivative of  $K$  with respect to  $X$ ):* Let  $K$  be the Gram matrix computed with a Gaussian kernel, and let  $x$  be an  $m \times 1$  column of the unmixing matrix, such that the  $(i, j)$ th entry of  $K$  is

$$k_{ij} = \phi(y_i, y_j) = \exp \left[ \frac{-1}{2\lambda^2} x^\top W_{ij} x \right],$$

where  $W_{ij} = (w_i - w_j)(w_i - w_j)^\top$ , and  $w_i$  is the  $i$ th sample of observations. Then the derivative of any  $k_{ij}$  with respect to  $x$  is

$$\frac{\partial k_{ij}}{\partial x^\top} = -\frac{k_{ij}}{\lambda^2} x^\top (w_i - w_j)(w_i - w_j)^\top.$$

Since the above derivative is a vector, we require appropriate notation to express the derivative of the entire Gram matrix in a tractable form. This is done using the  $\text{vec}(A)$  operation, which stacks the columns of the matrix  $A$  on top of each other. Thus, the resulting differential is

$$d(\text{vec}K) = \underbrace{\left[ \frac{\partial k_{11}}{\partial x}, \dots, \frac{\partial k_{n1}}{\partial x}, \frac{\partial k_{12}}{\partial x}, \dots, \frac{\partial k_{nn}}{\partial x} \right]^\top}_{\partial \text{vec}(K) / \partial x^\top} d(\text{vec}x), \quad (50)$$

where  $d(\text{vec}x) = dx$  and  $\partial k_{11} / \partial x = (\partial k_{11} / \partial x^\top)^\top$ . The derivatives of the submatrices  $\partial \text{vec}(K_{:,I}) / \partial x^\top$  and  $\partial \text{vec}(K_{I,I}) / \partial x^\top$  are submatrices of  $\partial \text{vec}(K) / \partial x^\top$ , restricted to the rows  $\partial k_{ij} / \partial x^\top$ , with  $1 \leq i \leq n$  and  $j \in I$ , or both  $i, j \in I$ , respectively.

### ACKNOWLEDGMENT

The authors would like to thank Aiyou Chen for providing us with his KDICA code, and Alex Smola and Knut Hüper for valuable discussions. This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778.

The present work started when H. Shen was with the Department of Information Engineering, The Australian National University, Australia and NICTA, Australia.

## REFERENCES

- [1] P. Comon, "Independent component analysis, a new concept?" Signal Processing, vol. 36, no. 3, pp. 287–314, 1994.
- [2] A. Hyvärinen, J. Karhunen, and E. Oja, Independent Component Analysis. New York: Wiley, 2001.
- [3] J.-F. Cardoso, "Infomax and maximum likelihood for blind source separation," IEEE Letters on Signal Processing, vol. 4, pp. 112–114, 1997.
- [4] —, "Blind signal separation: Statistical principles," Proceedings of the IEEE, Special Issue on Blind Identification and Estimation, vol. 86, no. 10, pp. 2009–2025, 1998.
- [5] E. G. Learned-Miller and J. W. Fisher III, "ICA using spacings estimates of entropy," Journal of Machine Learning Research, vol. 4, pp. 1271–1295, 2003.
- [6] H. Stögbauer, A. Kraskov, S. Astakhov, and P. Grassberger, "Least dependent component analysis based on mutual information," Phys. Rev. E, vol. 70, no. 6, p. 066123, 2004.
- [7] R. Boscolo, H. Pan, and V. P. Roychowdhury, "Independent component analysis based on nonparametric density estimation," IEEE Transactions on Neural Networks, vol. 15, no. 1, pp. 55–65, 2004.
- [8] A. Chen, "Fast kernel density independent component analysis," in Sixth International Conference on ICA and BSS, vol. 3889. Berlin/Heidelberg: Springer-Verlag, 2006, pp. 24–31.
- [9] A. Kankainen, "Consistent testing of total independence based on empirical characteristic function," Ph.D. dissertation, University of Jyväskylä, 1995.
- [10] A. Feuerverger, "A consistent test for bivariate dependence," International Statistical Review, vol. 61, no. 3, pp. 419–433, 1993.
- [11] J. Eriksson and K. Koivunen, "Characteristic-function based independent component analysis," Signal Processing, vol. 83, no. 10, pp. 2195–2208, 2003.
- [12] A. Chen and P. J. Bickel, "Consistent independent component analysis and prewhitening," IEEE Transactions on Signal Processing, vol. 53, no. 10, pp. 3625–3632, 2005.
- [13] S. Achard, D.-T. Pham, and C. Jutten, "Quadratic dependence measure for non-linear blind sources separation," in Fourth International Symposium on ICA and BSS, 2003, pp. 263–268.
- [14] S. Achard, C. Jutten, and D.-T. Pham, "Mutual information and quadratic dependence for post nonlinear blind source separation," in EUSIPCO 13, 2005.
- [15] B. Schölkopf and A. Smola, Learning with Kernels. Cambridge, MA: MIT Press, 2002.
- [16] F. R. Bach and M. I. Jordan, "Kernel independent component analysis," Journal of Machine Learning Research, vol. 3, pp. 1–48, 2002.
- [17] A. Gretton, R. Herbrich, A. J. Smola, O. Bousquet, and B. Schölkopf, "Kernel methods for measuring independence," Journal of Machine Learning Research, vol. 6, pp. 2075–2129, 2005.
- [18] A. Gretton, O. Bousquet, A. J. Smola, and B. Schölkopf, "Measuring statistical dependence with Hilbert-Schmidt norms," in ALT. Heidelberg: Springer-Verlag, 2005, pp. 63–78.
- [19] M. Rosenblatt, "A quadratic measure of deviation of two-dimensional density estimates and a test of independence," The Annals of Statistics, vol. 3, no. 1, pp. 1–14, 1975.

- [20] I. Steinwart, "On the influence of the kernel on the consistency of support vector machines," Journal of Machine Learning Research, vol. 2, pp. 67–93, 2002.
- [21] C. Micchelli, Y. Xu, and H. Zhang, "Universal kernels," Journal of Machine Learning Research, vol. 7, pp. 2651–2667, 2006.
- [22] A. Gretton, K. Fukumizu, C. H. Teo, L. Song, B. Schölkopf, and A. Smola, "A kernel statistical test of independence," in NIPS 20. Cambridge, MA: MIT Press, 2008, pp. 585–592.
- [23] S. Jegelka and A. Gretton, Large-Scale Kernel Machines. Cambridge, Massachusetts: MIT Press, 2007, ch. Brisk kernel ICA.
- [24] A. Edelman, T. A. Arias, and S. T. Smith, "The geometry of algorithms with orthogonality constraints," SIAM Journal on Matrix Analysis and Applications, vol. 20, no. 2, pp. 303–353, 1998.
- [25] K. Hüper and J. Trunpf, "Newton-like methods for numerical optimisation on manifolds," in Proceedings of Thirty-eighth Asilomar Conference on Signals, Systems and Computers, 2004, pp. 136–139.
- [26] H. Shen, K. Hüper, and A.-K. Seghouane, "Geometric optimisation and FastICA algorithms," in MTNS 17, Kyoto, Japan, 2006, pp. 1412–1418.
- [27] H. Shen and K. Hüper, "Newton-like methods for parallel independent component analysis," in MLSP 16, Maynooth, Ireland, 2006, pp. 283–288.
- [28] H. Shen, S. Jegelka, and A. Gretton, "Geometric analysis of Hilbert-Schmidt independence criterion based ICA contrast function," National ICT Australia Technical Report (PA006080), ISSN 1833-9646, October 2006.
- [29] —, "Fast kernel ICA using an approximate Newton method," in AISTATS 11. Microtome, 2007, pp. 476–483.
- [30] S. Fine and K. Scheinberg, "Efficient SVM training using low-rank kernel representations," Journal of Machine Learning Research, vol. 2, no. 22, pp. 243–264, 2001.
- [31] D.-T. Pham, "Fast algorithms for mutual information based independent component analysis," IEEE Transactions on Signal Processing, vol. 52, no. 10, pp. 2690–2700, 2004.
- [32] W. M. Boothby, An Introduction to Differentiable Manifolds and Riemannian Geometry, Revised, 2nd ed. Academic Press, 2002.
- [33] M. Spivak, A Comprehensive Introduction to Differential Geometry, Vol. 1 – 5, 3rd ed. Publish or Perish, Inc, 1999.
- [34] P.-A. Absil, R. Mahony, and R. Sepulchre, Optimization Algorithms on Matrix Manifolds. Princeton, NJ: Princeton University Press, 2008.
- [35] D. Gabay, "Minimizing a differentiable function over a differential manifold," Journal of Optimization Theory and Applications, vol. 37, no. 2, pp. 177–219, 1982.
- [36] M. Kleinsteuber, "Jacobi-type methods on semisimple Lie algebras – a Lie algebraic approach to the symmetric eigenvalue problem," Ph.D. dissertation, Bayerische Julius-Maximilians-Universität Würzburg, 2006.
- [37] E. Oja and Z. Yuan, "The FastICA algorithm revisited: Convergence analysis," IEEE Transactions on Neural Networks, vol. 17, no. 6, pp. 1370–1381, 2006.
- [38] C. K. I. Williams and M. Seeger, "Using the Nystrom method to speed up kernel machines," in NIPS 14. The MIT Press, 2001, pp. 682–688.
- [39] S. Amari, A. Cichocki, and H. H. Yang, "A new learning algorithm for blind signal separation," in NIPS 8. The MIT Press, 1996, pp. 757–763.



**Hao Shen** received his Bachelor degree in Mechanical Engineering and Applied Mathematics from Xi'an Jiaotong University, China, in 2000, and his Masters degrees in Computer Studies and Computer Science from the University of Wollongong, Australia, in 2002 and 2004, respectively; and received his PhD from the Australian National University, Australia, in 2008. Currently, he is a Post-doc at the Institute for Data Processing, Technische Universität München, Germany. His research interests focus on applying geometric optimisation techniques to linear independent component analysis and related topics in signal processing.

It includes the analysis of existing methods and the development of new efficient numerical algorithms to solve these kind of problems. He is a member of IEEE.



**Stefanie Jegelka** received her diploma in Computer Science with specialisation in Bioinformatics (with distinction) from the University of Tuebingen, Germany, in 2007. She also spent two semesters at the University of Texas at Austin as an exchange student. Since 2006, she has been a research assistant at the Max Planck Institute for Biological Cybernetics in Tuebingen, where she is now working towards a PhD. Her interests include machine learning and kernel methods, theoretical aspects of clustering and graph cuts with respect to statistical learning theory and approximation, and geometric and combinatorial optimization.

She has received a Google Anita Borg Scholarship and has been a scholar of the German National Academic Foundation during the studies for her diploma.



**Arthur Gretton** is a project scientist with the Machine Learning Department at Carnegie Mellon University since February 2009, and is affiliated as a research scientist with the Max Planck Institute for Biological Cybernetics, where he has worked since September 2002. He received degrees in physics and systems engineering from the Australian National University in 1996 and 1998, respectively; and completed his PhD with the Signal Processing and Communications Laboratory and Microsoft Research at the University of Cambridge in 2003. His research interests include machine learning, kernel methods, statistical learning

theory, nonparametric hypothesis testing, blind source separation, Gaussian processes, and non-parametric techniques for neural data analysis.