

# THE KERNEL MUTUAL INFORMATION

Arthur Gretton

Ralf Herbrich

Alexander J. Smola

MPI for Biological Cybernetics  
Spemannstr 38  
D-72076 Tübingen - Germany  
arthur@tuebingen.mpg.de

Microsoft Research  
7 J J Thomson Avenue  
Cambridge CB3 0FB, UK  
rherb@microsoft.com

RSISE, MLG  
Australian National University  
Canberra, Australia  
Alex.Smola@anu.edu.au

## ABSTRACT

We introduce a new contrast function, the kernel mutual information (KMI), to measure the degree of independence of continuous random variables. This contrast function provides an approximate upper bound on the mutual information, as measured near independence, and is based on a kernel density estimate of the mutual information between a discretised approximation of the continuous random variables. We show that Bach and Jordan's kernel generalised variance (KGV) is also an upper bound on the same kernel density estimate, but is looser. Finally, we suggest that the addition of a regularising term in the KGV causes it to approach the KMI, which motivates the introduction of this regularisation.

## 1. INTRODUCTION

The problem of separating mixtures of signals, so as to recover the original signals prior to mixing, is a much studied challenge in signal processing. Methods of solution generally depend on the nature of the signals, and the manner in which they are mixed; in particular, a criterion known as the *contrast function* is required to determine when the demixing is successful. We assume here that the original signals are generated i.i.d. according to some unknown probability distributions, and are combined in a scalar mixing process: demixing is then achieved by ensuring that the recovered signals are statistically independent. This is the framework for *instantaneous ICA*<sup>1</sup>, and has been used successfully in a wide variety of problems: for instance, the separation of linearly mixed audio signals, and the recovery of evoked potentials from EEG signals (see [10, 4], and references therein).

A measure of statistical independence between two random variables is the *mutual information* [5], which for random vectors  $\vec{x}, \vec{y}$  is zero if and only if the random vectors are independent. This may also be interpreted as the KL divergence  $D_{\text{KL}}(\mathbf{f}_{\vec{x}, \vec{y}} || \mathbf{f}_{\vec{x}} \mathbf{f}_{\vec{y}})$  between the joint density  $\mathbf{f}_{\vec{x}, \vec{y}}$  and the product of the marginal densities  $\mathbf{f}_{\vec{x}} \mathbf{f}_{\vec{y}}$ ; the latter quantity generalises readily to distributions of more than two random variables. We therefore propose two quantities, based on the mutual information, that may be used as contrast functions in ICA. The first, which we call the kernel covariance (KC), can be shown to be zero if and only if the random variables are independent. The second function, the kernel mutual information (KMI), is an upper bound on the Parzen window estimate of the mutual information, and is also zero if and only if the random variables are independent. Both functions bear a strong resemblance to the kernel canonical correlation (KCC) and kernel generalised variance (KGV) introduced by Bach and Jordan [3]:

<sup>1</sup>We shall in future refer to this problem simply as ICA.

indeed, we demonstrate that the KGV can also be thought of as a (looser) upper bound on the same Parzen window estimate. An important advantage of the derivation described herein, however, is that it addresses the behaviour of the contrast functions for finite kernel sizes, rather than relying on a limiting argument in which the kernel size approaches zero, as in [3]. Our approach thus allows us to apply well established methods for selecting kernel size as a function of the number of observations; see for instance [14].

In Section 2, we introduce the ICA problem, and describe our terminology. We then introduce the KC and KCC in Section 3, and derive the KMI and KGV in Section 4. Finally, we show in Section 5 that the performance of the KMI, when used in ICA, is competitive with that of the KGV, and that both the KMI and KGV outperform many traditional ICA algorithms.

## 2. ICA: PROBLEM STATEMENT

We begin by introducing the ICA problem. The discussion draws on the numerous existing surveys of ICA and related methods; see for instance [10, 4]. Suppose we have a random vector  $\vec{s}$  of dimension  $N$ , with independent identically distributed (i.i.d.) components (we use the sans serif to indicate random variables);

$$\mathbf{f}_{\vec{s}}(\vec{s}) = \prod_{i=1}^N \mathbf{f}_{s_i}(s_i),$$

where  $s_i \in \mathbb{R}$ . We do not observe  $\vec{s}$ , however: instead, we observe the random vector  $\vec{t}$ , such that

$$\vec{t} = \mathbf{A}\vec{s}, \quad (2.1)$$

where  $\mathbf{A}$  is an  $N \times N$  matrix<sup>2</sup>. Clearly, the components of  $\vec{t}$  will not be independent unless  $\mathbf{A} = \mathbf{P}\mathbf{S}$ , where  $\mathbf{P}$  is a permutation matrix and  $\mathbf{S}$  is a diagonal scaling matrix. Our goal is to find an approximation  $\mathbf{V}$  to the *inverse* of the matrix<sup>3</sup>  $\mathbf{A}$ , given  $m$  i.i.d. samples from  $\mathbf{f}_{\vec{t}}$ , and using *only* the model (2.1) and the fact that the unmixed components are independent. The determination of  $\mathbf{A}$  can only be made within certain identifiability constraints, however; in particular, no more than one source can be Gaussian.

Assume we have  $m$  observations  $\mathbf{t} := (\vec{t}_1, \dots, \vec{t}_m)$ . Our first step in computing  $\mathbf{V}$  is to subtract the mean of  $\mathbf{t}$  from each  $\vec{t}_i$ , and to whiten it,  $\vec{t} = \mathbf{Q}\vec{t}$ , such that the new observations  $\vec{t}$  have

<sup>2</sup>This corresponds to the number of sources being equal to the number of sensors. In fact, it is possible to recover (2.1) when the number of sources is less than the number of sensors by a change of basis, although the presence of noise makes this more difficult.

<sup>3</sup>Up to permutation and scaling.

a unit covariance matrix. Our estimate of the demixing matrix then becomes  $\mathbf{V} := \mathbf{W}\mathbf{Q}$ , where  $\mathbf{W}$  is an orthogonal matrix; our estimate of  $\tilde{\mathbf{s}}$  is  $\tilde{\mathbf{x}} := \mathbf{W}\tilde{\mathbf{l}}$ . Although the determination of  $\mathbf{W}$  remains difficult, there are only  $N(N-1)$  degrees of freedom involved in this problem, as opposed to the  $N^2$  degrees of freedom present in the estimation of  $\mathbf{V}$ .

### 3. THE KERNEL COVARIANCE AND CORRELATION

We now describe the kernel covariance, which is proposed as a measure of statistical independence of the random vectors  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$ , defined on  $\mathcal{X} := \mathbb{R}^{n_x}$  and  $\mathcal{Y} := \mathbb{R}^{n_y}$ . The generalisation to more than two vectors is addressed in [8]. We define the vectors  $\mathbf{x}$  and  $\mathbf{y}$  and the random vectors  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$  in the feature spaces  $\mathcal{F}_X$  and  $\mathcal{F}_Y$ , and the mappings  $\phi_x : \mathcal{X} \rightarrow \mathcal{F}_X$  and  $\phi_y : \mathcal{Y} \rightarrow \mathcal{F}_Y$  such that

$$\mathbf{x} := \phi_x(\tilde{\mathbf{x}}) \quad \text{and} \quad \mathbf{y} := \phi_y(\tilde{\mathbf{y}}).$$

The feature spaces may be the reproducing kernel Hilbert spaces (and subspaces of  $\mathcal{L}^2$ ) associated with particular kernels, which represent the inner products<sup>4</sup> on  $\mathcal{F}_X$  and  $\mathcal{F}_Y$ . We define

$$\mathbf{C}_{xy} := \mathbf{E}_{\mathbf{x},\mathbf{y}} \left( (\mathbf{x} - \mathbf{E}_x(\mathbf{x})) (\mathbf{y} - \mathbf{E}_y(\mathbf{y}))^\top \right), \quad (3.1)$$

$$\mathbf{C} := \begin{bmatrix} \mathbf{C}_{xx} & \mathbf{C}_{xy} \\ \mathbf{C}_{xy}^\top & \mathbf{C}_{yy} \end{bmatrix}, \quad (3.2)$$

where  $\mathbf{C}_{xx}$  and  $\mathbf{C}_{yy}$  are given by analogy. We observe  $m$  i.i.d. samples of data;  $\mathbf{z} = ((\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_m, \mathbf{y}_m))$ , where  $\mathbf{x}_i \in \mathcal{F}_X$  and  $\mathbf{y}_i \in \mathcal{F}_Y$ .

We may also define the *Gram matrices*  $\mathbf{K}_{mm}^{(x)}$ ,  $\mathbf{K}_{mm}^{(y)}$  of inner products between the mapped observations above, in the case where  $\mathcal{F}_X$  and  $\mathcal{F}_Y$  are reproducing kernel Hilbert spaces (RKHSs) with associated kernels  $k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j) := \mathbf{x}_i^\top \mathbf{x}_j = (\mathbf{K}_{mm}^{(x)})_{i,j}$  and  $k(\tilde{\mathbf{y}}_i, \tilde{\mathbf{y}}_j) := \mathbf{y}_i^\top \mathbf{y}_j = (\mathbf{K}_{mm}^{(y)})_{i,j}$ . According to [13], Gram matrices for the variables centred in feature space are  $\tilde{\mathbf{K}}_{mm}^{(x)} := \mathbf{H}\mathbf{K}_{mm}^{(x)}\mathbf{H}$ ,  $\tilde{\mathbf{K}}_{mm}^{(y)} := \mathbf{H}\mathbf{K}_{mm}^{(y)}\mathbf{H}$ , where  $\mathbf{H} = \mathbf{I}_m - m^{-1}\mathbf{1}_m\mathbf{1}_m^\top$ , and  $\mathbf{1}_m$  is an  $m \times 1$  vector of ones.

We can now introduce the kernel covariance (KC). In the population case, the KC is

$$\mathcal{J} = \sup_{f \in \tilde{\mathcal{F}}_X, g \in \tilde{\mathcal{F}}_Y} |\mathbf{E}_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}} [f(\tilde{\mathbf{x}})g(\tilde{\mathbf{y}})] - \mathbf{E}_{\tilde{\mathbf{x}}} [f(\tilde{\mathbf{x}})] \mathbf{E}_{\tilde{\mathbf{y}}} [g(\tilde{\mathbf{y}})]|,$$

where  $\tilde{\mathcal{F}}_X := \{f \in \mathcal{F}_X : \|f\|_{\mathcal{F}_X} \leq 1\}$ , and  $\tilde{\mathcal{F}}_Y$  is analogous. An empirical estimate  $\mathcal{J}(\mathbf{z})$  may be obtained from the finite sample  $\mathbf{z}$ , using the representer theorem (Schölkopf *et al.* [12]) to replace

$$f(\tilde{\mathbf{x}}) = \sum_{l=1}^m c_l k(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}_l) = \sum_{l=1}^m c_l \mathbf{x}^\top \mathbf{x}_l, \quad (3.3)$$

with a similar replacement for  $g(\tilde{\mathbf{y}})$ ; it follows that  $\mathcal{J}(\mathbf{z}) := \max_i \gamma_i$ , where  $\gamma_i$  are the eigenvalues of

$$\begin{bmatrix} \tilde{\mathbf{K}}_{mm}^{(x)} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{K}}_{mm}^{(y)} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} & \tilde{\mathbf{K}}_{mm}^{(x)} \tilde{\mathbf{K}}_{mm}^{(y)} \\ \tilde{\mathbf{K}}_{mm}^{(y)} \tilde{\mathbf{K}}_{mm}^{(x)} & \mathbf{0} \end{bmatrix}$$

We now describe the link between the kernel covariance and independence; details are given in [8].

<sup>4</sup>To be a kernel associated with a RKHS,  $k(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$  must satisfy the Mercer conditions [1]; these hold for Gaussian and Laplace kernels, among (many) others. Note also that the argument of the kernel specifies whether the kernel pertains to  $\mathcal{F}_X$  or  $\mathcal{F}_Y$ , although these kernels are identical in the present study.

**Theorem 1 (Kernel covariance and independence).**  $\mathcal{J} = 0$  if and only if  $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$  are independent.

Finally, we introduce the canonical correlation, as described in [7, 3, 11]; the final reference is a particularly insightful investigation of the canonical correlation in high dimensional spaces (such as RKHSs). We first define the canonical correlation in the general case, without reference to its interpretation when  $\mathcal{F}_X, \mathcal{F}_Y$  are RKHSs. We would like to find vectors  $\alpha_i, \beta_i$  onto which  $\mathbf{x}$  and  $\mathbf{y}$  respectively project, such that the correlation  $\rho_i$  between these projections is a stationary point with respect to  $\alpha_i, \beta_i$ . The canonical correlations,  $\rho_i$ , are thus given by

$$\rho_i := \frac{\alpha_i^\top \mathbf{C}_{xy} \beta_i}{\sqrt{(\alpha_i^\top \mathbf{C}_{xx} \alpha_i) (\beta_i^\top \mathbf{C}_{yy} \beta_i)}}. \quad (3.4)$$

When  $\mathcal{F}_X$  and  $\mathcal{F}_Y$  are RKHSs, then care must be taken when finding empirical estimates of the canonical correlates, to ensure that these estimates are data dependent. This may be done by confining  $\alpha_i, \beta_i$  to subspaces of the space spanned by the sample in  $\mathcal{F}_X, \mathcal{F}_Y$ , as in Kuss [11], or by regularising, as in [3]; in the latter case, the largest kernel canonical correlation may be used as a contrast function for ICA.

### 4. UPPER BOUNDS ON MUTUAL INFORMATION

We now apply both these definitions to derive an approximation of the mutual information between random variables  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$ , defined on the respective bounded intervals  $\mathcal{X}$  and  $\mathcal{Y}$  on  $\mathbb{R}$ . Full details of the proofs, and a generalisation to more than two random variables, may be found in [8]. We begin by introducing the Gaussian mutual information, and its relation with the canonical correlation. If  $\tilde{\mathbf{x}}_G, \tilde{\mathbf{y}}_G$  are Gaussian random variables in  $\mathbb{R}^{p_x}, \mathbb{R}^{p_y}$  respectively, then according to [3] the mutual information between them can be written

$$I(\tilde{\mathbf{x}}_G; \tilde{\mathbf{y}}_G) = -\frac{1}{2} \log \left( \prod_{i=1}^{\min(p_x, p_y)} (1 - \rho_i^2) \right), \quad (4.1)$$

where the  $\rho_i$  are given by the canonical correlations in (3.4).

Next, we consider a grid of size  $p_x \times p_y$  over  $\mathcal{X}$  and  $\mathcal{Y}$  respectively. Let the indices  $i, j$  denote the point  $(q_i, r_j) \in \mathcal{X} \times \mathcal{Y}$  on this grid, and let  $\mathbf{q} := (q_1, \dots, q_{p_x})$ ,  $\mathbf{r} := (r_1, \dots, r_{p_y})$  be the grid coordinates. The spacing between points along the  $x$  and  $y$  axes is respectively  $\Delta_x$  and  $\Delta_y$ . We define two multinomial random variables  $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$  with a distribution  $\mathbf{P}_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}}(i, j)$  over the grid (we write the complete  $p_x \times p_y$  matrix of such probabilities as  $\mathbf{P}_{xy}$ ), where

$$\mathbf{P}_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}}(i, j) := \int_{q_i}^{q_i + \Delta_x} \int_{r_j}^{r_j + \Delta_y} \mathbf{f}_{\mathbf{x},\mathbf{y}}(x, y) dx dy.$$

Thus  $\mathbf{P}_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}}(i, j)$  is a discretisation of  $\mathbf{P}_{xy}$ . We denote as  $\mathbf{p}_x$  the vector for which  $(\mathbf{p}_x)_i = \mathbf{P}_{\tilde{\mathbf{x}}}(i)$ , with a similar  $\mathbf{p}_y$  definition. We may always write  $\mathbf{P}_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}}(i, j) = \mathbf{P}_{\tilde{\mathbf{x}}}(i) \mathbf{P}_{\tilde{\mathbf{y}}}(j) (1 + \epsilon_{i,j})$  for an appropriate choice of  $\epsilon_{i,j}$ . If  $\epsilon_{i,j}$  is small, we approximate

$$I(\tilde{\mathbf{x}}; \tilde{\mathbf{y}}) \approx \frac{1}{2} \sum_{i=1}^{p_x} \sum_{j=1}^{p_y} \mathbf{P}_{\tilde{\mathbf{x}}}(i) \mathbf{P}_{\tilde{\mathbf{y}}}(j) \epsilon_{i,j}^2, \quad (4.2)$$

It is well known (see [5]) that  $I(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$  represents the upper bound on  $I(\tilde{\mathbf{x}}; \tilde{\mathbf{y}})$  as the discretisation becomes infinitely fine.

We next define an equivalent multidimensional representation  $\bar{x}, \bar{y}$  of  $\hat{x}, \hat{y}$  in the previous section, where  $\bar{x} \in \mathbb{R}^{p_x}$  and  $\bar{y} \in \mathbb{R}^{p_y}$ , such that  $\hat{x} = i$  is equivalent to  $(\bar{x})_i = 1$  and  $(\bar{x})_{j:j \neq i} = 0$ . Using

$$\mathbf{E}_{x,y}(\bar{x}\bar{y}^\top) = \mathbf{P}_{xy}, \quad \mathbf{E}_x(\bar{x}) = \mathbf{p}_x, \quad \mathbf{E}_x(\bar{x}\bar{x}^\top) = \mathbf{D}_x,$$

where  $\mathbf{D}_x = \text{diag}(\mathbf{p}_x)$ , it is possible to define covariances

$$\mathbf{C}_{xy} = \mathbf{P}_{xy} - \mathbf{p}_x \mathbf{p}_y^\top, \quad \mathbf{C}_{xx} = \mathbf{D}_x - \mathbf{p}_x \mathbf{p}_x^\top. \quad (4.3)$$

We define the Gaussian random variables  $\bar{x}_G, \bar{y}_G$  to have the same covariance structure as  $\bar{x}, \bar{y}$ , and with mutual information given by (4.1). The mutual information for this Gaussian case may then be approximated by (4.2) near independence; see [3, 8].

Given that we are not provided with the distribution  $\mathbf{P}_{\bar{x}, \bar{y}}(i, j)$ , but rather a finite sample  $\mathbf{z}$  of size  $m$ , we make use of a kernel density estimate of the mutual information for the discretised random variables. A detailed discussion of the properties and behaviour of such estimates may be found in [14], and previous work on their application to the computation of entropies in [9]. The kernel density (Parzen window) estimates of  $\mathbf{f}_x$  and  $\mathbf{f}_{x,y}$  are

$$\begin{aligned} \hat{\mathbf{f}}_x(x) &:= \frac{1}{m} \sum_{l=1}^m k(x_l, x), \\ \hat{\mathbf{f}}_{x,y}(x, y) &:= \frac{1}{m} \sum_{l=1}^m k(x_l, x) k(y_l, y). \end{aligned}$$

The kernels must be non-negative and continuous, with unit integral w.r.t. to its two arguments. We require approximations to the covariance matrices in the Gaussian mutual information, as described in (4.3). We therefore define the vectors  $\hat{\mathbf{p}}_x, \hat{\mathbf{p}}_y$ , and the matrix  $\hat{\mathbf{P}}_{xy}$ , using the expectations computed with these kernel expressions;

$$\hat{\mathbf{E}}_{xy}(\bar{x}\bar{y}^\top) = \hat{\mathbf{P}}_{xy}, \quad \hat{\mathbf{E}}_x(\bar{x}) = \hat{\mathbf{p}}_x, \quad \hat{\mathbf{E}}_x(\bar{x}\bar{x}^\top) = \hat{\mathbf{D}}_x.$$

Let  $\mathbf{K}_{pm}^{(x)}$  be the matrix of inner products in  $\mathcal{F}_X$  between the grid points and sample, with  $\mathbf{K}_{pm}^{(y)}$  defined by analogy; we assume  $p_x \gg m$  and  $p_y \gg m$ . The first subscript specifies whether the grid ( $q$  or  $r$ ) or the sample ( $x$  or  $y$ ) is used in the rows of this matrix, and the second subscript whether the grid or sample is used in the columns. By analogy, we may also define the matrices  $\mathbf{K}_{pp}^{(x)}, \mathbf{K}_{mm}^{(x)}, \mathbf{K}_{pp}^{(y)}, \mathbf{K}_{mm}^{(y)}$ . In the limit where  $\Delta_x, \Delta_y$  are small (and thus, by implication,  $p_x \gg m, p_y \gg m, \sigma \gg \Delta_x$ , and  $\sigma \gg \Delta_y$ , where  $\sigma$  defines the kernel size), we make the approximations

$$\begin{aligned} \hat{\mathbf{P}}_{xy} - \hat{\mathbf{p}}_x \hat{\mathbf{p}}_y^\top &= \frac{\Delta_x \Delta_y}{m} \left( \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(y)} \right)^\top - \frac{1}{m} \mathbf{K}_{pm}^{(x)} \mathbf{1}_m \left( \mathbf{K}_{pm}^{(y)} \mathbf{1}_m \right)^\top \right), \\ \hat{\mathbf{D}}_x &= \frac{\Delta_x}{m} \text{diag} \left( \mathbf{K}_{pm}^{(x)} \mathbf{1}_m \right) = \frac{\Delta_x^2}{m} \text{diag} \left( \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(x)} \right)^\top \mathbf{1}_{p_x} \right). \end{aligned}$$

The kernel density approximation to the discretised mutual information is then found by replacing the  $\rho_i$  in (4.1) with

$$\hat{\rho}_i := \frac{\hat{\mathbf{c}}_i^\top \left( \hat{\mathbf{P}}_{xy} - \hat{\mathbf{p}}_x \hat{\mathbf{p}}_y^\top \right) \hat{\mathbf{d}}_i}{\sqrt{\hat{\mathbf{c}}_i^\top \hat{\mathbf{D}}_x \hat{\mathbf{c}}_i \hat{\mathbf{d}}_i^\top \hat{\mathbf{D}}_y \hat{\mathbf{d}}_i}}. \quad (4.4)$$

This cannot easily be computed, however, since it is computationally prohibitive to evaluate the Gram matrices on a sufficiently fine grid. Noting that

$$\check{\nu}_m \hat{\mathbf{c}}_i^\top \mathbf{K}_{pm}^{(x)} \hat{\mathbf{c}}_i \leq \hat{\mathbf{c}}_i^\top \text{diag} \left( \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(x)} \right)^\top \mathbf{1}_p \right) \hat{\mathbf{c}}_i, \quad (4.5)$$

where  $\check{\nu}_m := \min_{j \in \{1 \dots p_x\}} \sum_{l=1}^m k(x_l, q_j)$ , we replace the matrix term in  $\hat{\mathbf{D}}_x$  by the left hand expression in (4.5), yielding a new quantity  $|\tilde{\gamma}_i| \geq |\hat{\rho}_i|$ ; it follows that replacing  $\hat{\rho}_i$  with  $\tilde{\gamma}_i$  yields an upper bound on (4.1). In fact,  $\tilde{\gamma}_i$  is simply the kernel covariance, but with the additional requirement that the functions  $f, g$  be projected in their respective feature spaces onto the basis spanned by the columns of the grid  $\mathbf{q}, \mathbf{r}$ , as well as an added scaling factor  $\check{\nu}_m$ . We use this insight to replace  $\tilde{\gamma}_i$  in (4.1) with an appropriately scaled  $\gamma_i$ , and  $\check{\nu}_m$  with  $\nu_m := \min_{j \in \{1 \dots m\}} \sum_{l=1}^m k(x_l, x_j)$ , to obtain the empirical *kernel mutual information* (KMI),

$$\mathcal{M}(\mathbf{z}) := -\frac{1}{2} \log \left( \left[ \mathbf{I} - (\nu_m \nu_y) \tilde{\mathbf{K}}_{mm}^{(x)} \tilde{\mathbf{K}}_{mm}^{(y)} \right] \right), \quad (4.6)$$

which is also an upper bound on (4.1). It follows from Theorem 1 that the random variables  $x, y$  are independent if and only if the *population* KMI satisfies  $\mathcal{M} = 0$ .

Bach and Jordan [3] propose a related quantity as a contrast function for ICA: the kernel generalised variance (KGV). In fact, the latter quantity may also be derived by finding an upper bound on (4.1): this is a different approach to the proof in [3], which uses a limit as the kernel becomes infinitely small. Using

$$\hat{\mathbf{c}}_i^\top \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(x)} \right)^\top \hat{\mathbf{c}}_i \leq \hat{\mathbf{c}}_i^\top \text{diag} \left( \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(x)} \right)^\top \mathbf{1}_{p_x} \right) \hat{\mathbf{c}}_i,$$

we replace the right hand term in the above with the left hand term in the denominator of (4.4) to get a set of kernel canonical correlations  $\tilde{\rho}_i \geq \hat{\rho}_i$ , restricted to the basis spanned by the grid. The mutual information computed using the unrestricted kernel canonical correlations  $\rho_i$  is therefore an upper bound on (4.1). The contrast function thus derived is never used in practice, since it is infinite; in other words, the approximation we made above is too loose. If we instead make the replacement

$$\frac{m}{\Delta_x^2} \hat{\mathbf{c}}_i^\top \hat{\mathbf{D}}_x \hat{\mathbf{c}}_i \Rightarrow \hat{\mathbf{c}}_i^\top \left( \theta_1 \mathbf{K}_{pm}^{(x)} \left( \mathbf{K}_{pm}^{(x)} \right)^\top + \theta_2 \nu_m \mathbf{K}_{pp}^{(x)} \right) \hat{\mathbf{c}}_i,$$

where  $\theta_1 \geq 0, \theta_2 \geq 0$ , and  $\theta_1 + \theta_2 \leq 1$ , we recover an expression which, for correct choice of  $\theta_1, \theta_2$ , yields the *regularised* KGV proposed in [3]<sup>5</sup>. We therefore expect the performance of both the KGV and KMI to be very similar when used for ICA: this is indeed the case in our experimental results.

We now briefly address the generalisation of the kernel covariance  $\mathcal{J}$  to the case of  $N$  random variables  $x_j$  on bounded subsets  $\mathcal{X}_j \subset \mathbb{R}$ , by analogy with derivation of [3]; this can be used to measure the pairwise independence of our estimate  $\bar{x}$  of the independent components  $\bar{s}$ . The KC is the largest eigenvalue  $\lambda_i$  of

$$\left( \tilde{\mathbf{K}} \tilde{\mathbf{K}}^\top - \begin{bmatrix} \tilde{\mathbf{K}}_1 \tilde{\mathbf{K}}_1 & & \\ & \ddots & \\ & & \tilde{\mathbf{K}}_N \tilde{\mathbf{K}}_N \end{bmatrix} \right) \mathbf{c}_i = \lambda_i \left( \text{diag}(\tilde{\mathbf{K}}) \mathbf{c}_i \right), \quad (4.7)$$

where  $\mathbf{c}_i = (c_{i,1}, \dots, c_{i,N})^\top$  and  $\tilde{\mathbf{K}} = [\tilde{\mathbf{K}}_1, \tilde{\mathbf{K}}_2, \dots, \tilde{\mathbf{K}}_N]^\top$ . To reduce computational cost, we use a reduced rank approximation of  $\tilde{\mathbf{K}}_j$ , via an incomplete Cholesky factorization with appropriate pivoting [6] (that is,  $\tilde{\mathbf{K}}_j \approx \mathbf{Z}_j \mathbf{Z}_j^\top$  with  $\mathbf{Z}_j \in \mathbb{R}^{m \times d}$  and  $d \ll m$ ). We set  $\mathbf{d}_i = [c_{i,1}^\top \mathbf{Z}_1, \dots, c_{i,N}^\top \mathbf{Z}_N]^\top$ , and rewrite (4.7) as

$$\left( \mathbf{Z} \mathbf{Z}^\top - \begin{bmatrix} \mathbf{Z}_1^\top \mathbf{Z}_1 & & \\ & \ddots & \\ & & \mathbf{Z}_N^\top \mathbf{Z}_N \end{bmatrix} \right) \mathbf{d}_i = \lambda_i \mathbf{d}_i, \quad (4.8)$$

<sup>5</sup>Specifically, the parameter denoting the amount of regularisation in [3] can be written  $\kappa = \theta_2 \nu_m / \theta_1$ , although we must be careful in our choice of  $\theta_1, \theta_2$  to ensure we still have an upper bound; see [8] for details.

where  $\mathbf{Z} = [\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_N]^\top$ . This transformation takes care of the nullspace inherent in  $\mathbf{Z}_j \mathbf{Z}_j^\top$  and reduces the eigenproblem to  $dN$  dimensions<sup>6</sup>. Finally, the KMI for more than two variables is

$$\mathcal{M}(z) := -\frac{1}{2} \log \left( \prod_{i=1}^N \nu_z \lambda_i \right) \quad (4.9)$$

where  $\nu_z = \min_j \nu_{z_j}$  (in our experiments, we simply set  $\nu_z = 1/m$ ; the performance remained satisfactory).

## 5. EXPERIMENTAL RESULTS

We now apply the KGV and KMI to the problem of ICA. Since the main purpose is to compare the performance with that reported in [3], we use identical settings and data. The mixing matrix  $\mathbf{A}$  was chosen randomly, with condition number between 1 and 2. We used the Gaussian RBF kernel,  $k(x, x') = \exp(-\frac{1}{2\sigma^2} \|x - x'\|^2)$ , with  $\sigma^2 = \frac{1}{4}$  and  $\kappa = 2 \times 10^{-3}$  for the KGV, except in the case of the 250 point sample, where  $\sigma^2 = 1$  and  $\kappa = 2 \times 10^{-2}$ . We only used  $\sigma^2 = 1$  for the KMI. The orthogonal component  $\mathbf{W}$  of the demixing matrix was found using gradient descent on the manifold of orthogonal matrices; see [3]. In order to measure the distance between the true ( $\mathbf{A}^{-1}$ ) and approximate ( $\mathbf{WQ}$ ) demixing matrices, we used the *Amari divergence* [4]. This metric is in the interval  $[0, 100]$ , is equal to zero if and only if  $\mathbf{A}^{-1}, \mathbf{WQ}$  are piecewise identical, and is invariant to permutation and scaling of  $\mathbf{A}^{-1}, \mathbf{WQ}$ .

Our experiment consisted in de-mixing data drawn independently from 2 – 16 distributions, chosen at random with replacement from 18 possible options; these include signals with both positive and negative kurtosis, and are described in detail in [3, 8]. Table 1 summarises our results; the KMI seems somewhat better in the case of larger  $m$  and  $N$ , although further refinement of the parameter choices in both methods might be possible. Further experiments are described in [8], most notably addressing the problem of recovering signals in the presence of noise, and in the case of low kurtosis. In these cases, the KMI and KGV again yield the best observed performance.

## 6. CONCLUSIONS

We have presented a novel derivation of several kernel based contrast functions for ICA (the KMI, KGV, and related), which yields useful insight both into the problem of model selection, and the function of the regularising term in these contrasts. The KMI and KGV are comparable in performance, and substantially outperform several alternative ICA approaches. Further work will focus on the application of kernel based contrasts to convolutive mixing, and to the recovery of random processes that are not i.i.d.; an application to graphical model estimation is given in [2].

## 7. REFERENCES

- [1] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68:337 – 404, 1950.
- [2] F. Bach and M. Jordan. Learning graphical models with Mercer kernels. In *NIPS 2002 (to appear)*.
- [3] F. Bach and M. Jordan. Kernel independent component analysis. *JMLR*, (3):1–48, 2002.

<sup>6</sup>Note that there is no reason why all  $\mathbf{Z}_j$  should have the same dimensionality: it may in fact be more computationally efficient in some circumstances to use different decompositions for different  $x_j$ .

**Table 1.** Illustration of the demixing of  $N$  randomly chosen signals; comparison with fast ICA, Jade, and extended Infomax. The best result is in boldface.

$N$	$m$	Rep.	fICA	Jade	lmax	KGV	KMI
2	250	1000	11.6± 0.4	10.6± 0.4	46.7± 0.9	5.4± <b>0.2</b>	6.2± <b>0.2</b>
2	1000	1000	6.2± 0.2	4.8± 0.2	10.9± 0.6	2.5± <b>0.1</b>	2.8± 0.1
4	1000	100	6.0± 0.4	5.5± 0.4	10.7± 0.9	3.5± <b>0.4</b>	3.7± 0.7
4	4000	100	3.3± 0.2	2.7± 0.1	6.2± 0.7	1.4± <b>0.1</b>	1.4± <b>0.05</b>
8	2000	50	4.0± 0.2	3.9± 0.3	8.2± 0.8	3.7± 0.9	2.9± <b>0.4</b>
8	4000	50	3.0± 0.3	2.5± 0.1	5.6± 0.7	1.5± 0.1	1.3± <b>0.04</b>
16	4000	28	3.1± 0.2	3.3± 0.2	11.1± 1.1	3.1± 0.9	2.2± <b>0.3</b>

- [4] A. Cichocki and S.-I. Amari. *Adaptive Blind Signal and Image Processing*. John Wiley and Sons, New York, 2002.
- [5] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley and Sons, New York, 1991.
- [6] S. Fine and K. Scheinberg. Efficient SVM training using low-rank kernel representation. Technical report, IBM Watson Research Center, New York, 2000.
- [7] M. Greenacre. *Theory and Applications of Correspondence Analysis*. Academic Press, London, 1984.
- [8] A. Gretton, R. Herbrich, and A. Smola. *The Kernel Mutual Information*. Max Planck Institute for Biological Cybernetics, 2002. Forthcoming (draft available on request).
- [9] L. Györfi and E. van der Meulen. An entropy estimate based on a kernel density estimation. *Colloquia Mathematica Societatis János Bolyai, 57: Limit Theorems in Probability and Statistics*:229–240.
- [10] A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. John Wiley and Sons, 2001.
- [11] M. Kuss. Kernel multivariate analysis. Master’s thesis, Technical University of Berlin, 2001.
- [12] B. Schölkopf, R. Herbrich, and A. Smola. A generalised representer theorem. In *Proceedings of the Annual Conference on Computational Learning Theory*, 2001.
- [13] B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319.
- [14] B. Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall, New York, 1986.