

# Kernel Independent Component Analysis

Jake Spiteri

University of Bristol

2020

# Table of Contents

## 1 Summary of ICA

## 2 Kernel ICA

- $\mathcal{F}$ -correlation
- Estimating the  $\mathcal{F}$ -correlation
  - Kernelization of CCA
- Outline of the algorithm

# Summary of Independent Component Analysis

Problem:

- We want to recover a latent random vector  $\mathbf{x} = (x_1, \dots, x_m)^\top$  from observations  $\mathbf{y} = (y_1, \dots, y_m)$  which are unknown linear functions of  $\mathbf{x}$ .
- The components of  $\mathbf{x}$  are modeled as mutually independent.
- An observation  $\mathbf{y}$  is modeled as

$$\mathbf{y} = \mathbf{A}\mathbf{x},$$

where  $\mathbf{A}$  is an  $m \times m$  matrix of parameters.

Given  $N$  observations of  $\mathbf{y}$ , we want to estimate  $\mathbf{A}$  and thus recover the latent vector  $\mathbf{x}$ .

# Seeking Independence

Our problem can be reduced to finding  $\mathbf{W} := \mathbf{A}^{-1}$  such that the components of  $\hat{\mathbf{x}} = \hat{\mathbf{W}}\mathbf{y}$  are *independent*.

We have previously performed ICA by maximizing the negentropy, which is a measure of non-Gaussianity.

To achieve independence we can estimate parameters by minimizing a *contrast function*, where a contrast function is defined to always be nonnegative and equal to zero if and only if variables  $x_1$  and  $x_2$  are independent.

# Kernel ICA

Kernel ICA uses kernel-based measures of statistical dependence.

## Definition ( $\mathcal{F}$ -correlation)

For a reproducing-kernel Hilbert space (RKHS)  $\mathcal{F}$ , the  $\mathcal{F}$ -correlation between the random variables  $f_1(x_1)$  and  $f_2(x_2)$ , where  $f_1, f_2 \in \mathcal{F}$  is:

$$\begin{aligned}\rho_{\mathcal{F}} &= \max_{f_1, f_2 \in \mathcal{F}} \text{corr}(f_1(x_1), f_2(x_2)), \\ &= \max_{f_1, f_2 \in \mathcal{F}} \frac{\text{cov}(f_1(x_1), f_2(x_2))}{(\text{var } f_1(x_1))^{1/2} (\text{var } f_2(x_2))^{1/2}}.\end{aligned}$$

Clearly if  $x_1$  and  $x_2$  are independent, then the  $\mathcal{F}$ -correlation is zero.

# Contrast function

We will use the following contrast function based on the  $\mathcal{F}$ -correlation

$$I_{\rho_{\mathcal{F}}} = -\frac{1}{2} \log(1 - \rho_{\mathcal{F}})$$

## Reproducing Property

Restricting the maximization in the  $\mathcal{F}$ -correlation to the RKHS allows us to exploit the *reproducing property*:

$$f(x) = \langle \Phi(x), f \rangle, \quad \forall f \in \mathcal{F},$$

where  $\Phi : \mathcal{X} \rightarrow \mathcal{F}$  is a map from our input space into the RKHS. This allows us to write

$$\begin{aligned} \rho_{\mathcal{F}} &= \max_{f_1, f_2 \in \mathcal{F}} \text{corr}(f_1(x_1), f_2(x_2)) \\ &= \max_{f_1, f_2 \in \mathcal{F}} \text{corr}(\langle \Phi(x_1), f_1 \rangle, \langle \Phi(x_2), f_2 \rangle) \end{aligned}$$

That is, the  $\mathcal{F}$ -correlation is the maximum correlation between one-dimensional linear projections of  $\Phi(x_1), \Phi(x_2)$ . This is the definition of the first *canonical correlation* between  $\Phi(x_1)$ , and  $\Phi(x_2)$ .

# Problem Setup

- To use the  $\mathcal{F}$ -correlation as a contrast function for ICA, we need to compute canonical correlations in our feature space.
- We need a kernelization of the canonical correlation. This will allow us to work with an empirical sample and work in the feature space.



# Kernelization of CCA

Let  $\{x_1^1, \dots, x_1^N\}$  and  $\{x_2^1, \dots, x_2^N\}$  denote sets of  $N$  empirical observations of  $x_1$  and  $x_2$ . The observations generate Gram matrices  $\mathbf{L}_1, \mathbf{L}_2$ , where  $\{\mathbf{L}_i\}_{r,k} := K(x_i^r, x_j^k)$ . We then compute the centered Gram matrices  $\mathbf{K}_1, \mathbf{K}_2$ . Our kernelized CCA problem becomes

$$\begin{aligned}\hat{\rho}_{\mathcal{F}}(\mathbf{K}_1, \mathbf{K}_2) &= \max_{\alpha_1, \alpha_2 \in \mathbb{R}^N} \text{corr}(\alpha_1^\top \mathbf{x}_1, \alpha_2^\top \mathbf{x}_2) \\ &= \max_{\alpha_1, \alpha_2 \in \mathbb{R}^N} \frac{\alpha_1^\top \mathbf{K}_1 \mathbf{K}_2 \alpha_2}{(\alpha_1^\top \mathbf{K}_1^2 \alpha_1)^{1/2} (\alpha_2^\top \mathbf{K}_2^2 \alpha_2)^{1/2}}\end{aligned}$$

## Kernelization of CCA

Based on the previous slide, we can perform a kernelized version of CCA by solving the generalized eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{K}_1\mathbf{K}_2 \\ \mathbf{K}_2\mathbf{K}_1 & \mathbf{0} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \rho \begin{pmatrix} \mathbf{K}_1^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix}$$

The  $\mathcal{F}$ -correlation is defined as the first (largest) eigenvalue of the kernelized CCA problem.

## Kernelization of CCA

Based on the previous slide, we can perform a kernelized version of CCA by solving the generalized eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{K}_1\mathbf{K}_2 \\ \mathbf{K}_2\mathbf{K}_1 & \mathbf{0} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \rho \begin{pmatrix} \mathbf{K}_1^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix}$$

The  $\mathcal{F}$ -correlation is defined as the first (largest) eigenvalue of the kernelized CCA problem. We can rewrite this as

$$\begin{pmatrix} \mathbf{K}_1^2 & \mathbf{K}_1\mathbf{K}_2 \\ \mathbf{K}_2\mathbf{K}_1 & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{K}_1^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix},$$

where  $\lambda = 1 + \rho$ . We can easily generalize this result to more than two variables.

## Kernelization of CCA

Based on the previous slide, we can perform a kernelized version of CCA by solving the generalized eigenvalue problem:

$$\begin{pmatrix} \mathbf{0} & \mathbf{K}_1\mathbf{K}_2 \\ \mathbf{K}_2\mathbf{K}_1 & \mathbf{0} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \rho \begin{pmatrix} \mathbf{K}_1^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix}$$

The  $\mathcal{F}$ -correlation is defined as the first (largest) eigenvalue of the kernelized CCA problem. We can rewrite this as

$$\begin{pmatrix} \mathbf{K}_1^2 & \mathbf{K}_1\mathbf{K}_2 \\ \mathbf{K}_2\mathbf{K}_1 & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{K}_1^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{K}_2^2 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix},$$

where  $\lambda = 1 + \rho$ . We can easily generalize this result to more than two variables. We will write this as

$$\mathcal{K}\alpha = \lambda\mathcal{D}\alpha$$

# Outline of the kernel ICA algorithm

---

## Algorithm KERNELICA-KCCA

**Input:** Data vectors  $y^1, y^2, \dots, y^N$   
Kernel  $K(x, y)$

1. Whiten the data
2. Minimize (with respect to  $W$ ) the contrast function  $C(W)$  defined as:
  - a. Compute the centered Gram matrices  $K_1, K_2, \dots, K_m$  of the estimated sources  $\{x^1, x^2, \dots, x^N\}$ , where  $x^i = Wy^i$
  - b. Define  $\hat{\lambda}_{\mathcal{F}}^{\kappa}(K_1, \dots, K_m)$  as the minimal eigenvalue of the generalized eigenvector equation  $\mathcal{K}_{\kappa}\alpha = \lambda\mathcal{D}_{\kappa}\alpha$
  - c. Define  $C(W) = \hat{I}_{\lambda_{\mathcal{F}}}(K_1, \dots, K_m) = -\frac{1}{2} \log \hat{\lambda}_{\mathcal{F}}^{\kappa}(K_1, \dots, K_m)$

**Output:**  $W$

---