

# Independent Component Analysis

## PGM Project

NATHAN DE LARA  
École polytechnique

FLORIAN TILQUIN  
ENS Cachan

VINCENT VIDAL  
ENS Ulm

January 11, 2016

## Contents

<b>1</b>	<b>Problem statement</b>	<b>1</b>
1.1	Introduction . . . . .	1
1.2	Information Theory . . . . .	2
1.3	ICA and Maximum Likelihood . . . . .	3
<b>2</b>	<b>Algorithms for ICA</b>	<b>3</b>
2.1	Hérault and Jutten (HJ) algorithm . . . . .	3
2.2	Jade algorithm . . . . .	3
2.3	FastICA algorithm . . . . .	4
2.4	Kernel ICA algorithm . . . . .	4
<b>3</b>	<b>Results</b>	<b>4</b>

## ABSTRACT

This paper is dedicated to the study of Independent Component Analysis. We intent to implement, apply and compare several algorithms while being presenting some theoretical aspects such as the link between the likelihood maximisation and the mutual information.

## 1 Problem statement

### 1.1 Introduction

The general Independent Component Analysis problem can be formalised this way: Suppose we have some random variables  $x \in \mathbb{R}^p$  which correspond to a mix of some primitive sources  $s \in \mathbb{R}^n$ . The aim is to extract from  $x$  every source  $s_i$ . To do so, we will suppose here that:

- the sources are independents.
- the mix is linear and instantaneous
- at most one source has a Gaussian distribution.

We define:

$$x = As \text{ and } y = Wx, \quad (1)$$

where  $A$  is the mixing matrix,  $W$  the separation matrix and  $y$  the estimation of the sources. The goal is then to find a matrix  $W$  that maximise a certain measure of independence of  $y$ .

As a measure of independence, we consider, for theoretical purpose, the mutual information:

$$I(Y) = \int_{\mathbb{R}^p} P(Y) \log \frac{P(Y)}{\prod_i P_i(Y_i)} dY. \quad (2)$$

However, as it is too hard to compute, we consider other contrast functions, invariant by permutation, scaling on coordinates and maximal for independent ones.

## 1.2 Information Theory

Let  $X \in \mathbb{R}^n$  be a random variable, we note  $P(X)$  his density and  $\Sigma_X$  his covariance matrix.

In the space of measures, let  $\mathcal{G}$  be the manifold of Gaussian distributions,  $\mathcal{P}$  the manifold of “product” distributions and  $\mathcal{P} \wedge \mathcal{G}$  the manifold of Gaussian “product” distributions. Note that these manifolds are exponential families.

The main advantage of this geometric point of view is that the Kullback-Leibler divergence allows the notion of projection on exponential families. The projection of  $P$  on the family  $\mathcal{E}$  is defined as the vector of  $\mathcal{E}$  that minimise the divergence to  $P$ . We write this projection  $P^\mathcal{E}$ .

Then, we define:

The **Kullback–Leibler divergence** distribution from  $Q$  to  $P$ :

$$K(P \parallel Q) = \int_{\mathbb{R}^n} P(x) \log \frac{P(x)}{Q(x)} dx. \quad (3)$$

The **entropy**:

$$H(P) = - \int_{\mathbb{R}^n} P(x) \log P(x) dx. \quad (4)$$

The **mutual information**:

$$I(Y) = K(P(Y) \parallel \prod_i P_i(Y_i)) = K(P(Y) \parallel P(Y)^\mathcal{P}). \quad (5)$$

The **negentropy**:

$$G_n(Y) = H(\mathcal{N}(\mathbb{E}[Y], \Sigma_Y)) - H(Y) = H(P(Y)^\mathcal{G}) - H(P(Y)). \quad (6)$$

The **non-gaussianity**:

$$G(Y) = K(Y \parallel \mathcal{N}(\mathbb{E}[Y], \Sigma_Y)) = K(P(Y) \parallel P(Y)^\mathcal{G}). \quad (7)$$

The **correlation**:

$$\begin{aligned} C(Y) &= K(\mathcal{N}(\mathbb{E}[Y], \Sigma_Y) \parallel \mathcal{N}(\mathbb{E}[Y], \text{Diag } \Sigma_Y)) \\ &= K(P(Y)^\mathcal{G} \parallel P(Y)^{\mathcal{P} \wedge \mathcal{G}}) \\ &= \frac{1}{2} \log \frac{\det(\text{Diag}(\Sigma_Y))}{\det(\Sigma_Y)}. \end{aligned} \quad (8)$$

Using the Pythagorean theorem and the two decompositions of  $K(P \parallel P^{\mathcal{P} \wedge \mathcal{G}})$ , through  $P^\mathcal{P}$  or  $P^\mathcal{G}$ , shown in the Figure 1, we can prove that:

$$I(Y) + \sum_i G(Y_i) = G(Y) + C(Y). \quad (9)$$

Because the non gaussianity is invariant under invertible affine transforms, minimising the mutual independence according to  $W$  is equivalent to minimise  $C(Y) - \sum_i G(Y_i)$ . We can then define a set of contrast function, for  $\alpha \geq 0$ :

$$\phi_\alpha(Y) = \alpha C(Y) - \sum_i G(Y_i). \quad (10)$$

For more information see [Car03].

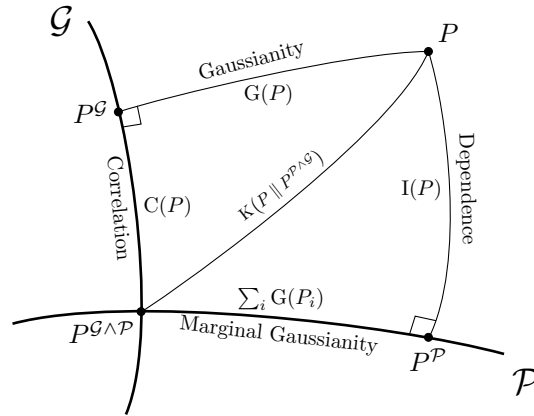


Figure 1: Representation of a distribution  $P$  and the different projections on the exponential families  $\mathcal{P}$  and  $\mathcal{G}$ . On the paths between the distributions are the quantities associated to the Kullback-Leibler divergence between those distributions.

### 1.3 ICA and Maximum Likelihood

As presented in [HO00a], it is possible to consider ICA as a maximum likelihood problem linked to the infomax principle. With the previously introduced notations, the log-likelihood is defined as:

$$L = \sum_t \sum_i \log f_i(w_i^T x(t)) + T \log(|\det(W)|) \quad (11)$$

Where  $f_i$  is the density function of  $s_i$ . The expectation of this likelihood is:

$$\mathbb{E}[L] = \sum_i \mathbb{E}[\log f_i(w_i^T x(t))] + \log(|\det(W)|) \quad (12)$$

In the case where  $f_i$  is the actual distribution of  $w_i^T x(t)$ , the first term becomes  $-\sum_i H(w_i^T x(t))$  which is one of the independence measures listed in 1.2.

## 2 Algorithms for ICA

### 2.1 Hérault and Jutten (HJ) algorithm

This method is based on the neural network principle. We write  $W = (I_n + \widetilde{W})^{-1}$  and for a pair of given functions  $(f, g)$ , we adapt  $\widetilde{W}$  as follows:

$$\widetilde{W}_{ij} = f(y_i)g(y_j). \quad (13)$$

### 2.2 Jade algorithm

Several methods are based on the cumulants. The goal here is to annul all the cross cumulants of order 4. Thus, the idea is to diagonalize the cumulant tensor which is equivalent to minimise the following contrast function:

$$c(x) = \sum_{i,k,l} |\text{Cum}(x_i, x_i^*, x_k, x_l)|^2. \quad (14)$$

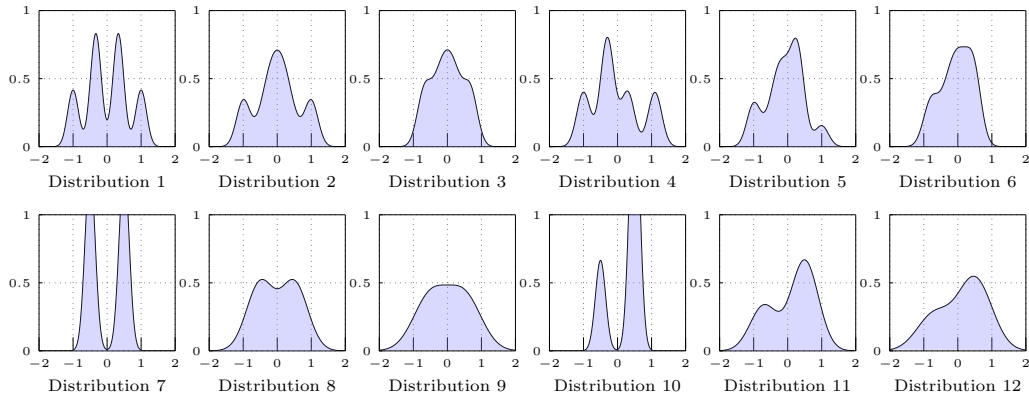


Figure 2: Distributions used to test the algorithms.

### 2.3 FastICA algorithm

The FastICA algorithm is based on the information theory. The goal here is to maximise the marginal non-gaussianity on the whitened data, relying on a non linear quadratic function  $f$  with the following rule:

$$\widetilde{W}_{t+1} = \mathbb{E}[X.f(W_t^T X)^T] - \mathbb{E}[f''(W_t^T X)] W_t, \quad (15)$$

with  $W_t$  the normalise vector of  $\widetilde{W}_t$ . In our experiments, we used  $f(x) = \frac{x^4}{4}$ . But it is possible to use  $f(x) = \log \cosh x$  or  $f(x) = \exp\left(-\frac{x^2}{2}\right)$  as well.

### 2.4 Kernel ICA algorithm

Given a reproducing kernel Hilbert space  $\mathcal{F}$ , this algorithm seeks to minimize the Kernel Generalized Variance defined as:

$$\widehat{\delta}_{\mathcal{F}} = -\frac{1}{2} \log \prod_i (1 - \rho_i^2) \quad (16)$$

where the  $\rho_i$  are the kernel canonical correlations between the observations components, obtained with computations over the observations Gram matrices.

## 3 Results

**Performance measure** The “Amari divergence”, equation 17, gives a criterion of proximity between two matrices, to evaluate the performance of an algorithm. If  $U$  and  $V$  are two  $n$ -by- $n$  matrices, the Amari error is defined by:

$$d(U, V) = \frac{1}{2n} \sum_{i=1}^n \left( \frac{\sum_{j=1}^n |a_{ij}|}{\max_j |a_{ij}|} - 1 \right) + \frac{1}{2n} \sum_{j=1}^n \left( \frac{\sum_{i=1}^n |a_{ij}|}{\max_i |a_{ij}|} - 1 \right) \quad (17)$$

with  $a_{ij} = (UV^{-1})_{ij}$ . This function, which is not an actual distance, has the advantage to be invariant by scaling factors and permutations of the matrices components.

	JADE	HJ	FastICA	Kernel ICA
1	5.67	6.37	6.51	<b>3.35</b>
2	9.47	8.15	11.54	<b>5.11</b>
3	6.76	<b>6.74</b>	8.39	10.55
4	6.21	6.82	7.66	<b>3.32</b>
5	5.99	6.42	7.30	<b>2.93</b>
6	9.32	<b>9.25</b>	12.18	9.37
7	2.90	<b>2.28</b>	3.68	2.74
8	8.52	<b>8.38</b>	10.31	12.53
9	17.52	<b>11.65</b>	21.95	28.21
10	13.94	10.30	17.71	<b>2.95</b>
11	9.57	8.98	12.31	<b>6.57</b>
12	19.19	<b>12.33</b>	22.93	13.76

m	N	JADE	HJ	FastICA	Kernel ICA
2	250	8.35	7.56	10.45	<b>6.06</b>
	1000	3.66	3.39	4.42	<b>2.38</b>
4	1000	11.88	58.52	12.20	<b>9.72</b>
	4000	5.33	83.07	5.80	<b>3.86</b>
8	2000	19.75	X	19.40	<b>19.15</b>
	4000	13.06	X	13.29	<b>9.71</b>
16	4000	31.81	X	<b>28.92</b>	X
	8000	<b>20.56</b>	X	27.89	X

Figure 3: **Left:** Average Amari divergence re-scaled by 100 obtained with the listed algorithms for random mix  $m = 2$  sources of size  $N = 250$  sampled with twelve different distributions. **Right:** Same measure for  $m$  sources of size  $N$  whose distributions are randomly selected among the twelve. The best results are in bold font. An X is put when a standard desktop computer could not compute the result.



Figure 4: Application of JADE algorithm to images separation. The first line presents the original sources, the second one the mix and the last one the estimations.

## References

- [ACY<sup>+</sup>96] Shun-ichi Amari, Andrzej Cichocki, Howard Hua Yang, et al. A new learning algorithm for blind signal separation. *Advances in neural information processing systems*, pages 757–763, 1996.
- [BJ03] Francis R Bach and Michael I Jordan. Kernel independent component analysis. *The Journal of Machine Learning Research*, 3:1–48, 2003.
- [BS95] Anthony J Bell and Terrence J Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural computation*, 7(6):1129–1159, 1995.
- [Car89] Jean-Francois Cardoso. Source separation using higher order moments. In *Acoustics, Speech, and Signal Processing, 1989. ICASSP-89., 1989 International Conference on*, pages 2109–2112. IEEE, 1989.
- [Car97] Jean-Francois Cardoso. Infomax and maximum likelihood for blind source separation. 1997.
- [Car03] Jean-François Cardoso. Dependence, correlation and gaussianity in independent component analysis. *The Journal of Machine Learning Research*, 4:1177–1203, 2003.
- [Com94] Pierre Comon. Independent component analysis, a new concept? *Signal processing*, 36(3):287–314, 1994.
- [HO00a] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4):411–430, 2000.
- [HO00b] Aapo Hyvärinen and Erkki Oja. Independent component analysis: algorithms and applications. *Neural networks*, 13(4):411–430, 2000.
- [JH91] Christian Jutten and Jeanny Herault. Blind separation of sources, part i: An adaptive algorithm based on neuromimetic architecture. *Signal processing*, 24(1):1–10, 1991.
- [LeB] Hervé LeBorgne. Analyse en composantes indépendantes. chapter 3.