

Mutual Information Approach to Blind Separation of Stationary Sources

Dinh Tuan Pham, *Member, IEEE*

Abstract—This paper presents a unified approach to the problem of blind separation of sources, based on the concept of mutual information. This concept is applied to the whole source sequences as stationary processes and thus provides a universal contrast applicable to both the instantaneous and convolutive mixture cases. For practical implementation, we introduce several degraded forms of this contrast, computable from a finite-dimensional distribution of the reconstructed source processes only. From them, we derive several sets of estimating equations, generalizing those considered earlier.

Index Terms—Contrast, convolution, entropy, independent component analysis, Kullback–Leibler divergence, mutual information, separation of sources, stationary process.

I. INTRODUCTION

BLIND separation of sources is a topic which has received much attention recently, as it has many important applications (see, e.g., [2] for a review). Basically, one observes several linear instantaneous or convolutive mixtures of independent signals,¹ called sources, and the problem is to recover them from the observations, *without relying on any specific knowledge of the sources*. In this blind context, a sensible method is to adopt the approach of an independent component analysis (ICA) in which a measure of dependence between the reconstructed sources is minimized. A natural such measure is the mutual information which has been introduced in [3] and implemented in [8]. However, this measure and others which have been proposed rely only on the marginal distribution of the sources, thereby ignoring their temporal properties. The exploitation of such properties can yield better separation in the case of instantaneous mixtures and is crucial in the case of convolutive mixtures. The last case has been less well investigated because of lack of a well-understood methodology: most works in this area adopt an *ad hoc* approach based on canceling cross cumulants. Recently, this author [10] has proposed the use of the mutual information between stationary processes as a contrast function for blind sources separation and discusses some implementation issues. In this work, we further develop this proposal and provide full proofs of results, as the cited paper was only a short version presented at the ICA'99 Workshop. We also provide several sets of

estimating equations, obtained by differentiating the above contrast and its variants, which can be related to earlier works [7], [9]. Our emphasis will be on the general ideas and concepts and therefore details of implementations of the methods will not be discussed (these implementations would depend on area of applications and can be the topics of subsequent works).

To proceed, let us describe the problem in mathematical terms and introduce some notations. We assume that K sequences of observations $\{X_k(t), t \in \mathbb{Z}\}$, $k = 1, \dots, K$, are available, each being a mixture of K independent sources, $\{S_j(t), t \in \mathbb{Z}\}$, $k = 1, \dots, K$, either instantaneously or through a convolution. More precisely, in the last case

$$\mathbf{X}(t) = \sum_{l=-\infty}^{\infty} \mathbf{A}(l)\mathbf{S}(t-l) = (\mathbf{A} \star \mathbf{S})(t) \quad (1)$$

where $\mathbf{X}(t)$ and $\mathbf{S}(t)$ denote the vectors $[X_1(t) \cdots X_K(t)]^T$ and $[S_1(t) \cdots S_K(t)]^T$, respectively, $\{\mathbf{A}(l), l \in \mathbb{Z}\}$ is a sequence of matrices, and \star denotes convolution. The instantaneous mixture case can be viewed as a particular case of the above where the sequence $\{\mathbf{A}(l), l \in \mathbb{Z}\}$ reduces to a single matrix of index 0, which we denote again by \mathbf{A} , so that

$$\mathbf{X}(t) = \mathbf{A}\mathbf{S}(t). \quad (2)$$

To separate the sources one naturally performs an “inverse” transformation on the sequence of observed vectors, namely

$$\mathbf{Y}(t) = \mathbf{B}\mathbf{X}(t) \quad (3)$$

in the instantaneous mixture case, \mathbf{B} denoting the separation matrix, and

$$\mathbf{Y}(t) = \sum_{l=-\infty}^{\infty} \mathbf{B}(l)\mathbf{X}(t-l) = (\mathbf{B} \star \mathbf{X})(t) \quad (4)$$

in the convolutive mixture case, $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ denoting the sequence of separation matrices.

The idea is to determine the separating matrix \mathbf{B} or sequence of matrices $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ such that the output sequence $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ in (3) or (4) has components as independent as possible. This is precisely the goal of the ICA, except that ICA thus far has been restricted to instantaneous transformations and the observed sequence $\{\mathbf{X}(t), t \in \mathbb{Z}\}$ is not necessarily a mixture of independent sources. Clearly, our approach requires a good measure of dependence between random stationary processes (as we assume that the sources are stationary), which we take to be the mutual information. This is a theoretical measure as it involves the density functions of the source processes, which in practice must be estimated from the data.

Manuscript received July 1, 2000; revised October 15, 2001.

The author is with the Laboratory of Modeling and Computation, IMAG–C.N.R.S., 38041 Grenoble Cedex 9, France (e-mail: Dinh-Tuan.Pham@imag.fr).

Communicated by U. Madhow, Associate Editor for Detection and Estimation.

Publisher Item Identifier S 0018-9448(02)05149-0.

¹In a more general setup, the output channels may be corrupted with noises, but in this paper we shall restrict ourselves to the pure mixture case.

The estimation problem will not be considered (for a simple case, see [8]) although we are fully aware of the difficulty of estimating the density in a high-dimensional space and we will try to avoid it as much as possible.

For ease of reading, proofs of results will be relegated to the Appendix.

II. MUTUAL INFORMATION BETWEEN STATIONARY PROCESSES

A. Some Definitions and Notations

Recall that the mutual information between a set of K random vectors Y_1, \dots, Y_K , with joint and marginal density functions p_{Y_1, \dots, Y_K} and p_{Y_1}, \dots, p_{Y_K} , is given by [4]

$$I(Y_1, \dots, Y_K) = -\mathbb{E} \log \frac{\prod_{i=1}^K p_{Y_i}(Y_i)}{p_{Y_1, \dots, Y_K}(Y_1, \dots, Y_K)}$$

$$= \sum_{i=1}^K h(Y_i) - h(Y_1, \dots, Y_K)$$

where

$$h(Y_1, \dots, Y_K) = -\mathbb{E} \log p_{Y_1, \dots, Y_K}(Y_1, \dots, Y_K)$$

and

$$h(Y_k) = -\mathbb{E} \log p_{Y_k}(Y_k)$$

are the (Shannon differential) joint and marginal entropies of Y_1, \dots, Y_K , respectively. Note that the notations

$$h([Y_1^T \dots Y_K^T]^T) \text{ and } h(Y_1, \dots, Y_K)$$

are the same and will be used interchangeably.

The mutual information is actually the Kullback–Leibler divergence between the joint density of Y_1, \dots, Y_K and their product densities. From the inequality $\log x \leq x - 1$ for all $x \geq 0$, with equality attained only at $x = 1$, one can see that $I(Y_1, \dots, Y_K) \geq 0$ with equality if only if Y_1, \dots, Y_K are independent. Thus, the mutual information is a measure of dependence between a set of random vectors.

1) *Entropy of Stationary Processes:* The entropy concept can be generalized to the case of stationary (vector) processes. For any process $\{Y(t), t \in \mathbb{Z}\}$ one can write [4]

$$h[Y(1), \dots, Y(m)] = \sum_{t=1}^m h[Y(t)|Y(t-1), \dots, Y(1)] \quad (5)$$

where

$$h[Y(t)|Y(t-1), \dots, Y(1)] = h[Y(t), \dots, Y(1)] - h[Y(t-1), \dots, Y(1)] \quad (6)$$

is the conditional entropy (that is the expected entropy of the conditional distribution) of $Y(t)$ given $Y(t-1), \dots, Y(1)$. But the conditioning decreases the entropy (see [4]), hence if the process is stationary

$$h[Y(t+1)|Y(t), \dots, Y(1)] \leq h[Y(t+1)|Y(t), \dots, Y(2)]$$

$$= h[Y(t)|Y(t-1), \dots, Y(1)] \quad (7)$$

One then deduces the following result, which is somewhat more precise than a result in [4, pp. 64–65 and 273].²

Lemma 1: For any stationary process $\{Y(t), t \in \mathbb{Z}\}$

$$h[Y(1), \dots, Y(m)]/m \geq h[Y(m)|Y(m-1), \dots, Y(1)]$$

and both sides of this inequality converge nonincreasingly to the same limit (possibly $-\infty$) as $m \rightarrow \infty$.

Following [4], we call the common limit in Lemma 1 the entropy (rate) of the process $\{Y(t), t \in \mathbb{Z}\}$ and denote it by $h[Y(\cdot)]$. If $Y(t)$ is a vector with components $Y_1(t), \dots, Y_K(t)$, it is also called joint entropy (rate) of the processes

$$\{Y_k(t), t \in \mathbb{Z}\}, \quad k = 1, \dots, K$$

and denoted by $h[Y_1(\cdot), \dots, Y_K(\cdot)]$.

2) *Mutual Information Between Processes:* The mutual information between the K jointly stationary processes $\{Y_k(t), t \in \mathbb{Z}\}, k = 1, \dots, K$ can now be defined as [4]

$$I[Y_1(\cdot), \dots, Y_K(\cdot)] = \sum_{i=1}^K h[Y_k(\cdot)] - h[Y_1(\cdot), \dots, Y_K(\cdot)]. \quad (8)$$

Clearly, $I[Y_1(\cdot), \dots, Y_K(\cdot)] \geq 0$ and vanishes if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are independent. For the converse, we are able to prove it only in the Markovian case (of arbitrary order, however) but we believe it holds much more generally. We can write

$$I[Y_1(\cdot), \dots, Y_K(\cdot)] = \sum_{i=1}^K \{h[Y_k(1)|Y_k(t), t < 1] - h[Y_k(1)|\mathbf{Y}(t), t < 1]\}$$

$$+ \left\{ \sum_{i=1}^K h[Y_k(1)|\mathbf{Y}(t), t < 1] - h[\mathbf{Y}(1)|\mathbf{Y}(t), t < 1] \right\} \quad (9)$$

where $h[Y_k(1)|Y_k(t), t < 1]$ stands for

$$\lim_{m \rightarrow \infty} h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)] = \lim_{m \rightarrow \infty} h[Y_k(1)|Y_k(0), \dots, Y_k(2-m)]$$

(the last equality coming from stationarity) and similarly for $h[Y_k(1)|\mathbf{Y}(t), t < 1]$ and $h[\mathbf{Y}(1)|\mathbf{Y}(t), t < 1]$, $\mathbf{Y}(t)$ denoting the vector $[Y_1(t) \dots Y_K(t)]^T$. Since the conditioning decreases the entropy, each term in the sum on the right-hand side of (9) is nonnegative. The last term on this right-hand side represents the mutual information between the components of the conditional distribution of $\mathbf{Y}(1)$ given $\mathbf{Y}(t), t < 1$ and hence is nonnegative as well. Thus, $I[Y_1(\cdot), \dots, Y_K(\cdot)] = 0$ implies that all these terms vanish, which entails that i) $Y_k(1)$ is independent of $\{Y_j(t), t < 1, j \neq k\}$, conditionally on $Y_k(t), t < 1$ and ii) $Y_1(1), \dots, Y_K(1)$ are independent conditionally on $\mathbf{Y}(t), t < 1$. For Markovian processes, ii) means that the transition probability factors into K factors and i) implies that each factor depends only on $\{Y_k(t), t < 1\}$ for an index k . Since the transition probability of a stationary ergodic Markov process determines its distribution entirely, the processes $\{Y_k(t), t \in \mathbb{Z}\}$ must be independent.

²Our result provides further an inequality and the monotonicity of the convergence of its left-hand side.

B. Calculation of Entropy

The computation of the entropy of a process through its definition is not practical as it involves a limiting operation. There are some special cases where this can be avoided, which we now consider.

1) *Gaussian Processes*: For Gaussian processes, a closed-form formula for the entropy is available. Indeed, let $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ be a Gaussian stationary vector process $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$, the conditional distribution of $\mathbf{Y}(m)$ given $\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)$ is a Gaussian distribution with covariance matrix \mathbf{G}_m , the error covariance matrix of the best linear predictor of $\mathbf{Y}(m)$ based on $\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)$, hence by a direct calculation (see also [4])

$$h[\mathbf{Y}(m)|\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)] = \frac{1}{2} \log \det(2\pi e \mathbf{G}_m)$$

where $e = \exp(1)$. Then letting m go to infinity and using the extension of Szegö's theorem to the multivariate case (see, e.g., [6, p. 162]), one gets

$$\begin{aligned} h[\mathbf{Y}(\cdot)] &= \frac{1}{2} \log \det(2\pi e \mathbf{G}) \\ &= \frac{1}{4\pi} \int_{-\pi}^{\pi} \log \det[4\pi^2 e \mathbf{f}(\lambda)] d\lambda \end{aligned} \quad (10)$$

where \mathbf{G} is the error covariance matrix of the best linear predictor of $\mathbf{Y}(t)$ based on $\mathbf{Y}(s), s < t$, and \mathbf{f} is the spectral density matrix of the process.

2) *Temporally Independent and Markovian Processes*: For such processes, the following result, which is an easy consequence of the definition of entropy and Lemma 1, is quite useful.

Corollary 1: For a stationary process $\{Y(t), t \in \mathbb{Z}\}$

$$h[Y(\cdot)] \leq h[Y(m)|Y(m-1), \dots, Y(1)]$$

with equality if and only if it is Markovian of order $m-1$, that is, the conditional distribution of $Y(t)$ given $Y(\tau), \tau < t$ depends only on $Y(t-1), \dots, Y(t+1-m)$. In particular, $h[Y(\cdot)] \leq h[Y(1)]$ with equality if and only if the process is temporally independent.

Thus, the entropy of a temporally independent process is simply its marginal entropy. Of greater interest is the fact that the entropy of a stationary Markovian process $\{Y(t), t \in \mathbb{Z}\}$ of order m equals

$$\begin{aligned} h[Y(m)|Y(m-1), \dots, Y(1)] \\ = h[Y(1), \dots, Y(m)] - h[Y(1), \dots, Y(m-1)] \end{aligned}$$

Because of stationarity, the last right-hand side also equals

$$\begin{aligned} h[Y(1), \dots, Y(m)] - h[Y(2), \dots, Y(m)] \\ = h[Y(1)|Y(2), \dots, Y(m)]. \end{aligned}$$

3) *Filtered Processes*: A general class of processes which includes the widely used autoregressive moving average (ARMA) processes is the class of *linear processes*, defined as follows.

Definition 1: A process $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ is called linear if it can be represented in the form

$$\mathbf{Y}(t) = \sum_{l=-\infty}^{\infty} \mathbf{A}(l) \mathbf{e}(t-l) \quad (11)$$

where $\{\mathbf{A}(l), l \in \mathbb{Z}\}$ is some sequence of matrices and $\{\mathbf{e}(t), t \in \mathbb{Z}\}$ is a sequence of independent and identically distributed random vectors.

In other words, a linear process is the output of some filter applied to a temporally independent process. Note that we allow the filter to be noncausal ($\mathbf{A}(l) \neq \mathbf{0}$ for $l < 0$), the causal case will receive some special attention later. Since we already know the entropy of a temporally independent process, it would be a simple matter to compute that of a linear process if we knew how to relate the entropy of a filtered process to that of the original process. To derive such a result, we will need to restrict ourselves to a class of "well-behaved" filters. We call a *sequence of square matrices* $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of class \mathcal{A} if

$$\sum_{l=-\infty}^{\infty} \|\mathbf{B}(l)\| < \infty$$

and

$$\det \left[\sum_{l=-\infty}^{\infty} \mathbf{B}(l) e^{il\lambda} \right] \neq 0$$

for all λ . It can then be seen that for such a sequence, the process $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ defined in (4) is well defined for any stationary process $\{\mathbf{X}(t), t \in \mathbb{Z}\}$ with finite α th absolute moment ($\alpha \geq 1$) and is itself stationary with finite α th absolute moment. Further, the class \mathcal{A} is closed with respect to the convolution in the sense that if it contains the sequences $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ and $\{\mathbf{C}(l), l \in \mathbb{Z}\}$, then it also contains their convolution.³ Another interesting property of the class \mathcal{A} is that any sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ in this class admits an inverse, with respect to the convolution, which is also of this class, the inverse being precisely the sequence of the Fourier coefficients of the function

$$\lambda \mapsto \left[\sum_{l=-\infty}^{\infty} \mathbf{B}(l) e^{il\lambda} \right]^{-1}.$$

This result follows from a result of Wiener which says that if f is a 2π -periodic function, nonzero everywhere, and has absolutely summable Fourier coefficients, then the same is true for $1/f$ (see, e.g., [14, p. 245]).

We shall further need a lower semi-continuity condition of the entropy functional.

Definition 2: The entropy functional is said to be lower semi-continuous, with respect to the convolution, at the process $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ if for any integer $m > 0$, real $\epsilon > 0$, there exists $\delta > 0$ such that

$$h[(\mathbf{C} \star \mathbf{Y})(1), \dots, (\mathbf{C} \star \mathbf{Y})(m)] \leq h[\mathbf{Y}(1), \dots, \mathbf{Y}(m)] + \epsilon$$

for all sequences $\{\mathbf{C}(l), l \in \mathbb{Z}\}$ satisfying

$$\|\mathbf{C}(0) - \mathbf{I}\| + \sum_{l \neq 0} \|\mathbf{C}(l)\| \leq \delta.$$

³This can be seen from

$$\begin{aligned} \left\| \sum_l \left[\sum_m \mathbf{B}(l-m) \mathbf{C}(m) \right] \right\| &\leq \sum_l \sum_m \|\mathbf{B}(l-m)\| \|\mathbf{C}(m)\| \\ &= \left[\sum_l \|\mathbf{B}(l)\| \right] \left[\sum_m \|\mathbf{C}(m)\| \right] \end{aligned}$$

and the fact that the Fourier transform transforms a convolution into a multiplication.

Proposition 1: Let $\{\mathbf{X}(t), t \in \mathbb{Z}\}$ be a vector random stationary process admitting α th absolute moment ($\alpha \geq 1$) and $\mathbf{Y}(t) = (\mathbf{B} \star \mathbf{X})(t)$, where $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ is a sequence of matrices of class \mathcal{A} . Assume that the entropy functional is lower semi-continuous at both $\{\mathbf{X}(t), t \in \mathbb{Z}\}$ and $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$, then

$$h[\mathbf{Y}(\cdot)] = h[\mathbf{X}(\cdot)] + \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}(l) e^{i l \lambda} \right| \frac{d\lambda}{2\pi}.$$

Proposition 1 can be viewed as an extension of the following result, which can be easily obtained along the same lines as in [8] based on [8, Lemma A1]: *Let \mathbf{X} be a random vector and \mathbf{B} be an invertible matrix, then the entropy of $\mathbf{Y} = \mathbf{B}\mathbf{X}$ equals*

$$h(\mathbf{Y}) = h(\mathbf{X}) + \log |\det \mathbf{B}|. \quad (12)$$

The lower semi-continuity condition is admittedly hard to verify, but it is very mild. It holds under the following condition, which we believe to be far from necessary.

Lemma 2: Let $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ be a vector random stationary process admitting the α th absolute moment ($\alpha \geq 1$) such that the joint density p_m of $\mathbf{Y}(1), \dots, \mathbf{Y}(m)$ exists (for all m) and is differentiable with $\|\nabla \log p_m(\mathbf{y})\| \leq C(1 + \|\mathbf{y}\|^{\alpha-1})$, ∇ denoting the gradient operator and C being a constant.⁴ Then the entropy functional is lower semi-continuous at $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$.

Proposition 1 is fundamental in that it describes how the entropy changes when the process is filtered. In particular, it provides the entropy of a linear process or more generally a filtered Markov process. Further, it provides a method for the *deconvolution* of a linear (or a filtered Markov) process, through the use of the following corollary.

Corollary 2: Let $\{\mathbf{Y}(t), t \in \mathbb{Z}\}$ be a vector stationary process such that the entropy functional is lower semi-continuous at $\{(\mathbf{B} \star \mathbf{Y})(t), t \in \mathbb{Z}\}$, for all sequences $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of class \mathcal{A} . Then

$$h[\mathbf{Y}(\cdot)] \leq h[(\mathbf{B} \star \mathbf{Y})(1)] - \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}(l) e^{i l \lambda} \right| \frac{d\lambda}{2\pi} \quad (13)$$

with equality if and only if the process $\{(\mathbf{B} \star \mathbf{Y})(t), t \in \mathbb{Z}\}$ is temporally independent. The same inequality holds with $h[(\mathbf{B} \star \mathbf{Y})(1)]$ replaced by

$$h[(\mathbf{B} \star \mathbf{Y})(m) | (\mathbf{B} \star \mathbf{Y})(m-1), \dots, (\mathbf{B} \star \mathbf{Y})(1)];$$

in this case, equality is attained if and only if the process $\{(\mathbf{B} \star \mathbf{Y})(t), t \in \mathbb{Z}\}$ is Markovian of order $m-1$.

Clearly, Corollary 2 still holds if the class \mathcal{A} is replaced by a smaller subclass. A subclass of interest is the subclass \mathcal{A}^+ of sequences $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ which are causal and have causal inverses, in the sense that it and its inverse sequence vanish at negative indexes. It is well known that the last condition is equivalent to the minimum-phase condition: $\det[\sum_{l=0}^{\infty} \mathbf{B}(l) z^l] \neq 0$ for all complex number z of modulus not exceeding 1. Under this condition

$$\oint_C \left\{ \log \det \left[\mathbf{I} + \sum_{l=1}^{\infty} \mathbf{B}(0)^{-1} \mathbf{B}(l) z^l \right] \right\} dz / (2\pi i z) = 0$$

⁴This constant can depend on m as the δ in Definition 2 can depend on m .

the integration being made along the unit circle C of the complex plane. Thus,

$$\int_{-\pi}^{\pi} \log \left| \det \sum_{l=0}^{\infty} \mathbf{B}(l) e^{i l \lambda} \right| \frac{d\lambda}{2\pi} = \log \det |\mathbf{B}(0)|. \quad (14)$$

Therefore, by restricting to the class \mathcal{A}^+ , one gets the same result as in Corollary 2 with the last integral in (13) replaced by $\log \det |\mathbf{B}(0)|$.

III. CONTRASTS

We shall assume throughout that the sequence $\{\mathbf{A}(l), l \in \mathbb{Z}\}$ in (1) is of class \mathcal{A} , hence in the reconstruction formula (4) we will restrict ourselves to sequences $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of this class. As mentioned earlier, to separate the source, one may minimize the mutual information $I[Y_1(\cdot), \dots, Y_K(\cdot)]$, where $Y_k(t)$ are the components of $\mathbf{Y}(t)$ defined by (3) or (4) according to model (2) or (1) is considered. But by Proposition 1, this criterion equals, up to a constant term

$$C_{\infty} = \sum_{k=1}^K h[Y_k(\cdot)] - \int_{-\pi}^{\pi} \log \left| \det \sum_{l=-\infty}^{\infty} \mathbf{B}(l) e^{i l \lambda} \right| \frac{d\lambda}{2\pi}. \quad (15)$$

in the case of model (1), or the same expression but with the integral replaced by $\log |\det \mathbf{B}|$, in the case of model (3).

The above criterion, by construction, is a *contrast* [3] in the sense that it is minimized if the reconstructed sources $\{Y_k(t), t \in \mathbb{Z}\}$, $k = 1, \dots, K$, coincide with the true sources *up to a permutation and a filtering* (or a scaling in the instantaneous mixture case). This ambiguity is *inherent* to the blind source separation problem (since it relies only on the independence assumption of the sources) and is manifested in the invariance property of C_{∞} : it is unchanged when one preconvolves the sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ with a sequence of diagonal matrices of class \mathcal{A} and premultiplies the result with a permutation matrix, as can be easily seen by applying Proposition 1 to scalar filtered processes. To be useful, however, C_{∞} should be discriminating, in the sense of [3], that is, it should attain its minimum *only* when the reconstructed sources coincide with the true sources *modulo the above ambiguity*. But minimizing C_{∞} can only ensure the independence between the reconstructed sources and thus for this contrast to be discriminating one may need some further conditions (such as non-Gaussianity), but since we will not actually use it, we do not pursue this question.

A. Contrasts for Instantaneous Mixtures

The contrast C_{∞} is of theoretical interest only, since its computation requires the complete knowledge of the distribution of each process $\{Y_k(t), t \in \mathbb{Z}\}$. Although one can always approximate $h[Y_k(\cdot)]$ by $h[Y_k(1), \dots, Y_k(m)]/m$, the number m might be very large for the approximation to be accurate, leading to the problem of estimation of a density in a *high-dimensional* space, which we would like to avoid.⁵ Therefore, it is of interest to obtain a simplified version of C_{∞} .

⁵The amount of data needed for a "good" density estimation in a high-dimensional space grows *exponentially* with the dimension.

1) *Contrasts Based on Finite Joint Distribution:* Instead of considering the mutual information between processes, we consider the mutual information between segments of processes. Explicitly, we consider the criterion

$$\frac{1}{m} I\{[Y_1(1) \cdots Y_1(m)]^T, \dots, [Y_K(1) \cdots Y_K(m)]^T\} \quad (16)$$

where m is a (small) integer and $Y_k(t)$ are the components of $\mathbf{Y}(t)$, defined in (3). But from (12)

$$h[\mathbf{Y}(1), \dots, \mathbf{Y}(m)] = m \log \det(\mathbf{B}) + h[\mathbf{X}(1), \dots, \mathbf{X}(m)] \quad (17)$$

hence this criterion can be seen to be equal, up to a constant term, to

$$C_m(\mathbf{B}) = \frac{1}{m} \sum_{k=1}^K h[Y_k(1), \dots, Y_k(m)] - \log |\det \mathbf{B}|. \quad (18)$$

By construction, this is a contrast, which, in the case where $m = 1$, has been shown to be discriminating if no more than one source can be Gaussian [3]. For $m > 1$, one can allow the sources to be Gaussian if the covariance matrices of m consecutive observations of the Gaussian sources are not proportional (see [12], [13]).

One can view C_m as an approximation to C_∞ , in which the entropy $h[Y_k(\cdot)]$ is replaced by $h[Y_k(1), \dots, Y_k(m)]/m$. An alternative approach is to replace it by the conditional entropy $h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)]$, which is a better approximation by Lemma 1. This leads to the criterion

$$C_m^*(\mathbf{B}) = \sum_{k=1}^K h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)] - \log |\det \mathbf{B}|. \quad (19)$$

The following result shows that $C_m^*(\mathbf{B})$ is a contrast. It can be shown to be discriminating under the same conditions as for the contrast C_m (see [12], [13]).

Lemma 3: Under the model (2), the criterion (19) equals the expected Kullback–Leibler divergence between the conditional distribution of $\mathbf{Y}(m)$ given $\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)$ and the product of the conditional distribution of $Y_k(m)$ given $Y_k(m-1), \dots, Y_k(1)$, plus the constant term $h[\mathbf{X}(m)|\mathbf{X}(m-1), \dots, \mathbf{X}(1)]$.

2) *The Gaussian Mutual Information:* To avoid the difficulty in calculating the entropy and mutual information, we introduce the Gaussian entropy mutual information, defined as before but with the random vectors or processes involved replaced by the Gaussian random vectors or processes having the same covariance structure. From (8) and (10), the Gaussian mutual information between the stationary processes $\{Y_k(t), t \in \mathbb{Z}\}$, $k = 1, \dots, K$ is

$$I_g[Y_1(\cdot), \dots, Y_K(\cdot)] = \frac{1}{4\pi} \int_{-\pi}^{\pi} [\log \det \text{diag } \mathbf{f}_Y(\lambda) - \log \det \mathbf{f}_Y(\lambda)] d\lambda \quad (20)$$

where \mathbf{f}_Y is the spectral density matrix of the vector process $\{\mathbf{Y}(t) = [Y_1(t) \cdots Y_K(t)]^T, t \in \mathbb{Z}\}$ and diag denotes the diagonal matrix with the same diagonal as its argument.

The criterion (20) is a joint diagonalization criterion, since it is nonnegative and can be zero if and only if \mathbf{f}_Y is diagonal almost everywhere. This is easily seen from the Hadamard inequality which says that for a positive-definite matrix \mathbf{f} , $\det \mathbf{f} < \det \text{diag } \mathbf{f}$ unless \mathbf{f} is diagonal in which case one has equality (see, e.g., [4, p. 502]). Since this criterion involves only the correlations between the sources, it would *not permit the separation of a convolutive mixture*: It is easy to see that it vanishes as soon as $[\sum_{l=-\infty}^{\infty} \mathbf{B}(l)e^{il\lambda}] \mathbf{f}_Y^{1/2}(\lambda)$ is a unitary matrix for almost all λ , $\mathbf{f}_Y^{1/2}$ being the Cholesky factor in the decomposition $\mathbf{f}_Y = \mathbf{f}_Y^{1/2}(\mathbf{f}_Y^{1/2})^T$. However, in the instantaneous mixture case, this contrast is discriminating, provided that there exists no pair of sources which have proportional spectral densities [11]. Since $\log \det \mathbf{f}_Y = 2 \log \det \mathbf{B} + \log \det \mathbf{f}_X$ in this case, this contrast is equivalent to

$$C_g(\mathbf{B}) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \log[\det \text{diag } \mathbf{f}_Y(\lambda)] d\lambda - \log \det \mathbf{B}.$$

One can further degrade the above contrast by considering the Gaussian analogs of the criteria C_m and C_m^* . This yields, after dropping a constant term

$$C_{g,m}(\mathbf{B}) = \frac{1}{2m} \sum_{k=1}^K \log |\det \text{cov}\{[Y_k(1) \cdots Y_k(m)]^T\}| - \log |\det \mathbf{B}| \quad (21)$$

$$C_{g,m}^*(\mathbf{B}) = \frac{1}{2} \sum_{k=1}^K \log \text{var}\{Y_k(m) - Y_k(m|1:m-1)\} - \log |\det \mathbf{B}| \quad (22)$$

where $\text{cov}\{\cdot\}$ refers to covariance matrix, $\text{var}\{\cdot\}$ refers to variance, and $Y_k(m|1:m-1)$ denotes the best linear predictor of $Y_k(m)$ based on $Y_k(1), \dots, Y_k(m-1)$. It can be shown that (21) and (22) are discriminating contrasts provided that there exists no pair j, k such that the covariance matrices of $[S_j(1) \cdots S_j(m)]^T$ and $[S_k(1) \cdots S_k(m)]^T$ are proportional [11].

B. Convolutive Mixtures and/or Linear Source Processes

Unfortunately, the above approach cannot be generalized to the case of convolutive mixtures. The reason is that the convolution is a transformation on the whole process, not a finite segment of it. The criterion (16) is still a contrast (although we are not sure if it is discriminating), but, unlike the instantaneous mixture case, would involve the joint entropy of $\mathbf{Y}(1), \dots, \mathbf{Y}(m)$; Proposition 1 is not applicable since it concerns the entropy of a whole process, not of a finite segment of it. Thus, the use of (16) would require the estimation of the entropy of an mK -dimensional distribution, which we would like to avoid. By the same reason, there is no analog of the contrast (19) for the convolutive mixture case; Lemma 3 applies only to the instantaneous mixture case.

However, if one restricts oneself to the class of linear or Markovian source processes, then simple contrasts can be constructed.

1) Convolutional Mixtures of Linear Sources:

Proposition 2: Assume that the sources are linear processes, specifically

$$S_k(t) = \sum_{l=-\infty}^{\infty} a_k(l)e_k(t-l) \quad (23)$$

where $\{e_k(t), t \in \mathbb{Z}\}$ are temporally independent processes and $\{a_k(l), l \in \mathbb{Z}\}$ are sequences of class \mathcal{A} , then the criterion

$$C_1[\mathbf{B}(\cdot)] = \sum_{k=1}^K h[Y_k(1)] - \int_{-\pi}^{\pi} \log \left| \det \sum_{l=-\infty}^{\infty} \mathbf{B}(l)e^{i\lambda l} \right| \frac{d\lambda}{2\pi}. \quad (24)$$

is minimized if and only if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are independent among themselves and are temporally independent.

The preceding result shows that the criterion (24) is a contrast since it is minimized when the reconstructed sources coincide with the true sources up to a permutation and a convolution. However, it is not necessarily discriminating since there is still the possibility that the processes $\{Y_k(t), t \in \mathbb{Z}\}$, despite being independent, do not coincide with the sources up to a permutation and a filtering (an example is the case where the sources are Gaussian). Nevertheless, it can be proved by a different method in [12] (and mentioned in [13]) that the contrast (24) is discriminating if no more than one source can be Gaussian.

As made clear by Proposition 2, minimizing (24) not only separates the sources but deconvolves them as well. If there can be no more than one Gaussian source, one would recover the sequences $\{e_k(t), t \in \mathbb{Z}\}$ in (23) up to a scaling, a permutation, and a time shift. More precisely, minimizing (24) yields $Y_k(t) = \alpha_k e_{\pi_k}(t - \tau_k)$ for some permutation $\{\pi_1, \dots, \pi_K\}$, some nonzero constants $\alpha_1, \dots, \alpha_k$, and some integers τ_1, \dots, τ_k . Note that, unlike the contrast (15) which is invariant with respect to filtering, the contrast (24) is not.

A more restrictive assumption on the distribution of the sources is that they are *linear causal processes with minimum phase*. By this we mean that the sequences $\{a_k(l), l \in \mathbb{Z}\}$ in the representation (23) are of class \mathcal{A}^+ . Assume further that the sequence $\{\mathbf{A}(l), l \in \mathbb{Z}\}$ in (1) is also of this class; then it makes sense to restrict the sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ in (4) to this class as well. Therefore, the contrast (24), by (14), reduces to

$$C_1^+[\mathbf{B}(\cdot)] = \sum_{k=1}^K h[Y_k(1)] - \log |\det \mathbf{B}(0)|. \quad (25)$$

As before, minimizing it *among all sequences of class \mathcal{A}^+* not only separates the sources but deconvolves them as well. The contrast (25), however, relies on somewhat artificial assumptions on the sources and the mixing matrix sequence. But it has the advantage of being simple.

2) Convolutional Mixtures of Markovian Sources: A weaker assumption on the sources is that they are filtered Markov processes. More precisely, it is assumed that the sources can be represented by (23) but with $\{e_k(t), t \in \mathbb{Z}\}$ now being an

$(m-1)$ th-order Markov process. Then, similarly to Proposition 2, one can show that the criterion

$$C_m^*[\mathbf{B}(\cdot)] = \sum_{k=1}^K h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)] - \int_{-\pi}^{\pi} \log \left| \det \sum_{l=-\infty}^{\infty} \mathbf{B}(l)e^{i\lambda l} \right| \frac{d\lambda}{2\pi} \quad (26)$$

is minimized if and only if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are independent among themselves and are Markovian of order $m-1$.

The proof of this result is very similar to that of Proposition 2, substituting $h[Y_k(1)]$ by $h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)]$.

As before, minimizing $C_m^*[\mathbf{B}(\cdot)]$ not only separates the sources but actually extracts the underlying Markov processes which generate them. In practice, it is likely that the sources are themselves Markovian and not filtered Markov processes; in this case, they are recovered exactly up to a permutation, a scaling, and a time shift. The ambiguity with respect to filtering is lifted because one has the *a priori* information that the sources are Markovian.

3) Instantaneous Mixture of Linear Sources: In the previous subsection, we have focused on the convolutional mixture case, but the approach there can be also applied to the instantaneous mixture case. By Corollary 2, C_∞ in this case is bounded above by

$$\sum_{k=1}^K \inf \left\{ h[(b_k \star Y_k)(1)] - \int_{-\pi}^{\pi} \log \left| \sum_{l=-\infty}^{\infty} b_k(l)e^{i\lambda l} \right| \frac{d\lambda}{2\pi} \right\} - \log |\det \mathbf{B}| \quad (27)$$

where the infimum is taken over all sequences $\{b_k(l), l \in \mathbb{Z}\}$ of class \mathcal{A} . Further, equality can be achieved if and only if the process $\{Y_k(t), t \in \mathbb{Z}\}$ is linear. Since C_∞ is a contrast and the sources are linear processes, this shows that (27) is indeed a contrast.

Clearly, by (14), C_∞ is also bounded above by

$$\sum_{k=1}^K \inf \{ h[(b_k \star Y_k)(1)] - \log |b_k(0)| \} - \log |\det \mathbf{B}| \quad (28)$$

where the infimum is taken over all sequences $\{b_k(l), l \in \mathbb{Z}\}$ of class \mathcal{A}^+ . Further, by Corollary 2 again, equality can be achieved if and only if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are linear causal with minimum phase. Thus, in the case where the sources are linear causal with minimum phase processes, (28) is indeed a contrast.

It is worthwhile to note that the coefficients $b_k(0)$ in (27) and (28) can be taken equal to 1, that is, the infima there are taken over all sequences $\{b_k(l), l \in \mathbb{Z}\}$ in \mathcal{A} or \mathcal{A}^+ with $b_k(0) = 1$. This is because multiplying the sequence $\{b_k(l), l \in \mathbb{Z}\}$ by a constant does not change the expression inside the curly bracket $\{ \}$ in (27) and (28).

C. Discussion

The contrast C_1 is well known (see, e.g., [3], [8]). But it exploits only the marginal distribution of the sources at a given time point. Our contrasts C_m and C_m^* , $m > 1$, involve their temporal dependence as well and thus could have better per-

formance especially in the case where the sources are strongly temporally dependent (note that if they are white, $C_m = C_m^* = C_1$). However, m should be small due to the difficulty of estimating the entropy of a high-dimensional random vector. Thus, it might be of interest to consider the contrasts (27) and (28) which requires only the entropy of random variables and yet taking into account the temporal dependence of the sources. The drawback is that they rely on the linearity assumption of the sources and require an extra minimization. Another possibility is to focus only on the second-order dependence of the sources, as implied by the use of the contrasts $C_{m,g}$, $C_{m,g}^*$, and C_g . The use of correlations only, but including lagged correlations, for blind sources separation, has been proposed, for example, in [1] and [9]. The use of C_1 and C_m^* in the convolutive mixtures case is new.

IV. ESTIMATING EQUATIONS

By differentiating the above contrasts, one obtains a system of equations to be satisfied, called estimating equations (see [5]). For this purpose, the following result plays a central role.

Lemma 4: Let \mathbf{Y} and \mathbf{Z} be two random vectors admitting absolute α th moment for some $\alpha \geq 1$. Assume that \mathbf{Y} and $\mathbf{Y} + \mathcal{E}\mathbf{Z}$, \mathcal{E} being a matrix for which the product $\mathcal{E}\mathbf{Z}$ makes sense and has the same dimension as \mathbf{Y} , admit densities $p_{\mathbf{Y}}$ and $p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}$ satisfying the following conditions.

C1) As $\mathcal{E} \rightarrow \mathbf{0}$

$$\int \log[p_{\mathbf{Y}}(u)/p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u)] p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u) du \rightarrow 0$$

faster than $\|\mathcal{E}\|$.

C2) The function $-\log p_{\mathbf{Y}}$ admits almost everywhere a gradient (column) vector $\psi_{\mathbf{Y}}$ such that

$$\|\psi_{\mathbf{Y}}(u)\| \leq C(1 + \|u\|^{\alpha-1}), \quad \text{for all } u$$

for some constant C .

Then as $\mathcal{E} \rightarrow \mathbf{0}$

$$h(\mathbf{Y} + \mathcal{E}\mathbf{Z}) - h(\mathbf{Y}) = \mathbf{E}(\psi_{\mathbf{Y}}^T \mathcal{E}\mathbf{Z}) + o(\mathcal{E})$$

where $o(\mathcal{E})$ denotes a term tending to $\mathbf{0}$ faster than \mathcal{E} .

Note: Condition C1) could be hard to verify, but it is quite reasonable. Indeed

$$\begin{aligned} & \int \log \frac{p_{\mathbf{Y}}(u)}{p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u)} p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u) du \\ &= \int \left[\log \frac{p_{\mathbf{Y}}(u)}{p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u)} - \frac{p_{\mathbf{Y}}(u)}{p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u)} + 1 \right] p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}(u) du. \end{aligned}$$

For small \mathcal{E} , one would expect that the expression inside the bracket $[\]$ is of the order $\|\mathcal{E}\|^2$ and thus the whole integral would be of this order. The difficulty is that $p_{\mathbf{Y}}$ and $p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}$ converge to zero at infinity and hence the behavior of the ratio $p_{\mathbf{Y}}/p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}$ near infinity is difficult to predict. The expression inside the bracket is in general of the order $\|\mathcal{E}\|^2$ for fixed u , but not uniformly in u . This uniformity is, however, not at all necessary since we will integrate with respect to $p_{\mathbf{Y}+\mathcal{E}\mathbf{Z}}$, which can be expected to converge to zero with a fast rate. But we have been unable to find simple conditions to ensure that C1) is satisfied.

The function $\psi_{\mathbf{Y}}$ will play a fundamental role in the sequel. In the case of a real random variable, it is usually referred to (in the statistical literature) as the score function. For a random vector, we therefore call $\psi_{\mathbf{Y}}$ the multivariate score function of the density of \mathbf{Y} .

A. Instantaneous Mixtures

We now apply the above result to obtain necessary conditions for the contrasts (18) and (19) to be minimized.

Proposition 3: A necessary condition for C_m to be minimized at \mathbf{B} is

$$\mathbf{E}\{[Y_j(1) \cdots Y_j(m)] \psi_{k,m}[Y_k(1), \dots, Y_k(m)]\} = 0, \quad 1 \leq j \neq k \leq K \quad (29)$$

and for C_m^* to be minimized at \mathbf{B} is

$$\mathbf{E}\{[Y_j(1) \cdots Y_j(m)] \psi_{k,m}^*[Y_k(1), \dots, Y_k(m)]\} = 0, \quad 1 \leq j \neq k \leq K \quad (30)$$

where $\psi_{k,m}$ and $\psi_{k,m}^*$ are the multivariate score functions of the joint density of $Y_k(1), \dots, Y_k(m)$ and of the conditional density of $Y_k(m)$ given $Y_k(m-1), \dots, Y_k(1)$. (Here the conditional density is considered as a function of both the dependent and the conditioning variables.)

Note: It can be seen from the proof of the preceding result that the conditions

$$\mathbf{E}\{[Y_k(1) \cdots Y_k(m)] \psi_{k,m}[Y_k(1), \dots, Y_k(m)]\} = m$$

$$\mathbf{E}\{[Y_k(1) \cdots Y_k(m)] \psi_{k,m}^*[Y_k(1), \dots, Y_k(m)]\} = 1$$

are also necessary. But these conditions are actually always satisfied because of the definitions of $\psi_{k,m}$ and $\psi_{k,m}^*$. This is an easy consequence of the following result which can be obtained through an integration by parts.

Lemma 5: Let Y be a random vector having a density f_Y such that $f_Y(y)y \rightarrow 0$ as $y \rightarrow \pm\infty$ and $-\log f_Y$ admits a gradient ψ_Y . Then $\mathbf{E}[Y_k \psi_{k,Y}(Y)] = 1$, Y_k , and $\psi_{k,Y}$ denotes the k th component of Y and ψ_Y .

Consider now the Gaussian contrasts (21), (22), and (20).

Proposition 4: A necessary condition for $C_{g,m}$ to be minimized at \mathbf{B} is

$$\sum_{t=1}^m \text{cov} \left\{ Y_j(t), \sum_{s=1}^m b_{k,m}(t,s) Y_k(s) \right\} = 0, \quad 1 \leq j \neq k \leq K \quad (31)$$

for $C_{g,m}^*$ to be minimized at \mathbf{B} is

$$\sum_{t=1}^m \text{cov} \left\{ Y_j(t), a_{k,m-1}(m-t) \sum_{l=0}^{m-1} a_{k,m-1}(l) Y_k(m-l) \right\} = 0, \quad 1 \leq j \neq k \leq K \quad (32)$$

and for C_g^* to be minimized at \mathbf{B} is

$$\int_{-\pi}^{\pi} [f_{Y_k Y_j}(\lambda) / f_{Y_k}(\lambda)] d\lambda = 0, \quad 1 \leq j \neq k \leq K \quad (33)$$

where $b_{k,m}(t, s)$ are the general elements of the inverse of $\text{cov}\{[Y_k(1) \cdots Y_k(m)]^T\}$, $a_{k,m-1}(l)$ are the coefficients in the representation $\sum_{l=0}^{m-1} a_{k,m-1}(l)Y_k(m-l)$ of $Y_k(m) - Y_k(m|1 : m-1)$ and $f_{Y_k Y_j}$ is the cross-spectral density between the processes $\{Y_k(t), t \in \mathbb{Z}\}$ and $\{Y_j(t), t \in \mathbb{Z}\}$.

One can see that (33) is a limiting form of (32) as $m \rightarrow \infty$. Indeed, $a_{k,m-1}(l)$ converges to $a_k(l)$ such that $\sum_{l=0}^{\infty} a_k(l)Y_k(t-l)$ is the error of the best linear predictor of $Y_k(t)$ based on $Y_k(s)$, $s < t$. Hence, the right-hand side of (32) converges to a constant times that of (33), since $f_{Y_k Y_k}(\lambda)$ is proportional to $1/|\sum_{l=0}^{\infty} a_k(l)e^{i\lambda l}|^2$. It is also possible to prove that (33) is a limiting form of (31) as well.

B. Linear and Markovian Sources

We first consider the convolutive mixture case.

Proposition 5: Assume that the sources are linear processes. A necessary condition for the contrast (24) to be minimized at the sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of class \mathcal{A} is

$$\begin{aligned} E\{Y_j(1-\tau)\psi_k[Y_k(1)]\} &= 0, \\ j, k &= 1, \dots, K, \tau \in \mathbb{Z}, j \neq k \text{ or } \tau \neq 0 \end{aligned} \quad (34)$$

and if the sources are also causal with minimum phase, a necessary condition for the contrast (25) to be minimized at a sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of class \mathcal{A}^+ is

$$\begin{aligned} E\{Y_j(1-\tau)\psi_k[Y_k(1)]\} &= 0, \\ j, k &= 1, \dots, K, \tau \geq 0, j \neq k \text{ or } \tau > 0 \end{aligned} \quad (35)$$

where, in both cases, ψ_k denotes the score function of the density of $Y_k(t)$.

A similar result, concerning the contrast (26), can be obtained by a combination of the proofs of Propositions 5 and 3.

Proposition 6: Assume that the sources are filtered $(m-1)$ th-order Markov processes. A necessary condition for the contrast (26) to be minimized at the sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ of class \mathcal{A} is that

$$\begin{aligned} E\{[Y_j(1-\tau) \cdots Y_j(m-\tau)]\psi_{k,m}^*[Y_k(1), \dots, Y_k(m)]\} &= 0, \\ j, k &= 1, \dots, K, \tau \in \mathbb{Z}, j \neq k \text{ or } \tau \neq 0 \end{aligned} \quad (36)$$

where $\psi_{k,m}^*$ is as in Proposition 3.

For the instantaneous mixtures case, the estimating equations associated with (27) and (28) are somewhat more complex.

Proposition 7: Assume that the source processes are linear. If \mathbf{B} minimizes (27) and the infimum of

$$h[(b_k \star Y_k)(1)] - \int_{-\pi}^{\pi} \log \left| \sum_{l=-\infty}^{\infty} b_k(l)e^{i\lambda l} \right| d\lambda / (2\pi)$$

is attained at some sequence $\{b_k^*(l), l \in \mathbb{Z}\}$ of class \mathcal{A} , then

$$E\{(b_k^* \star Y_j)(1)\psi_k[(b_k^* \star Y_k)(1)]\} = 0, \quad 1 \leq j \neq k \leq K \quad (37)$$

and if the sources are also causal with minimum phase and \mathbf{B} minimizes (25) and the infimum of $h[\sum_{l=0}^{\infty} b_k(l)Y_k(1-l)]$,

among all sequences $\{b_k(l), l \in \mathbb{Z}\}$ of class \mathcal{A}^+ with $b_k(0)=1$, is attained at some sequence $\{b_k^*(l), l \in \mathbb{Z}\}$, then (37) is again satisfied. In both cases, ψ_k denotes the score function of the density of $(b_k^* \star Y_k)(1)$.

C. Discussion

It can be seen from the preceding results that, in the instantaneous mixtures case, the estimating equations (29)–(32) are of the form

$$\begin{aligned} E\{[Y_j(1) \cdots Y_j(m)]\varphi_{k,m}[Y_k(1), \dots, Y_k(m)]\} &= 0, \\ 1 \leq j \neq k \leq K \end{aligned} \quad (38)$$

where $\varphi_{k,m}$ is a function from \mathbb{R}^m to \mathbb{R}^m . The estimating equations (33) associated with the Gaussian contrast C_g can be viewed as a limiting form of the above, as shown before. This is also true for the estimating equations (37) associated with the contrasts (27) and (28), as they can be put into the form (38) if one truncates the sequence $\{b_k^*(l), l \in \mathbb{Z}\}$ to a finite sequence (which one must in practice). In this case, $\varphi_{k,m}$ takes the form

$$\varphi_{k,m}(y(1), \dots, y(m)) = \begin{bmatrix} \beta_{k,1} \\ \vdots \\ \beta_{k,m} \end{bmatrix} \varphi_k \left[\sum_{l=1}^m b_{k,l} y(l) \right] \quad (39)$$

for some real function φ_k of a real variable and some real numbers $\beta_{k,1}, \dots, \beta_{k,m}$. In all cases, the functions $\varphi_{k,m}$, referred to as separating functions, are related to the densities of the sources in a specific way (they are linear in the Gaussian case).

Turning to the convolutive mixture case, we see that the estimating equations (34)–(36) are of the form

$$\begin{aligned} E\{[Y_j(1-\tau) \cdots Y_j(m-\tau)]\varphi_{k,m}[Y_k(1), \dots, Y_k(m)]\} &= 0, \\ j, k &= 1, \dots, K, \tau \in \mathbb{Z}, j \neq k \text{ or } \tau \neq 0 \end{aligned} \quad (40)$$

where $\varphi_{k,m}$ is a function from \mathbb{R}^m to \mathbb{R}^m , with $m=1$ in the case of linear sources and τ constrained to be nonnegative in the case where the sources are further causal with minimum phase and the reconstruction sequence of matrices $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ is restricted to the class \mathcal{A}^+ . Note that the system (40) contains an infinite number of equations with an infinite number of unknowns! In practice, one may restrict $\mathbf{B}(l)$ to be zero for l outside some given range $[L_1, L_2]$ and restrict τ to the same range, so as to have just K equations less than the number of unknowns (which accounts for the indeterminacy of scale). Note, however, that taking $L_1 = 0$ is not enough to ensure that the sequence $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ is of class \mathcal{A}^+ . This constraint is actually not easy to enforce.

The use of a system of estimating equations of the form (38) or (40) is much more flexible than that of contrasts, since such a system *needs not arise from the differentiation of a contrast*. In the context of blind source separation, it is simply a system of equations, which is satisfied when the reconstructed sources are independent [5]. (Note that the system generally includes the expectation operator, which should be replaced by appropriate sample average before being solved to obtain the estimates of the parameters.) It can be easily seen that any system of the form (38) or (40) is a system of estimating equations, as soon as the sources (or the $\varphi_{k,m}[Y_k(1), \dots, Y_k(m)]$) have

zero mean. Note that taking $m = 1$ in (38) yields the set of estimating functions introduced in [9], which can be traced back to the method in [7], while taking $\varphi_{k,m}$ of the form (39) with φ_k being the identity function yields the method for separating correlated sources in [9]. Also, many *ad hoc* methods for blind source separation consist in equating to zero the cross cumulants of higher order, possibly with lag, between the reconstructed sources. This amounts roughly to solving a system of the form (38) or (40).

But there is a price to pay for the above flexibility. First, the system of estimating equations constitutes only a necessary condition, it can (and often does) lead to spurious reconstructed sources, as such equations often have multiple solutions. We believe that by deriving them from a contrast one has a better chance of avoiding this problem. The fact that they come from a contrast makes it possible to monitor the calculation algorithm so as to ensure that the contrast is decreased at each step of the algorithm (and thus the reconstructed sources are closer to independence after each step in some sense). Second, the choice of the separating functions $\varphi_{k,m}$ can have great impact on the performance of the method: a bad choice could severely degrade the performance. Our results provide a set of good candidates for the separating functions, as they are derived from the mutual information contrast, which is related to the maximum-likelihood principle (see [2]). These functions need not be exactly the ones given in our propositions though. They can be simply some rough estimates of them. Note that the general form (38) or (40) requires the specification of Km real functions of m real variables and thus allows many degrees of freedom. If one believes that the sources are linear processes, or may be well approximated by such processes, one may settle for separating functions of the form (39) or (40) with $m = 1$, which requires the specification of only K real functions, but K linear filters need also to be estimated or specified.

APPENDIX PROOFS OF RESULTS

A. Proof of Lemma 1

By (7)

$$h[Y(t)|Y(t-1), \dots, Y(1)] \geq h[Y(m)|Y(m-1), \dots, Y(1)],$$

$$1 \leq t < m. \quad (41)$$

Therefore, from (5) one gets the inequality of the lemma. On the other hand, by (5) again

$$\begin{aligned} & h[Y(1), \dots, Y(m)]/m - h[Y(1), \dots, Y(m-1)]/(m-1) \\ &= \frac{1}{m} h[Y(m)|Y(m-1), \dots, Y(1)] \\ &\quad - \frac{1}{m(m-1)} \sum_{t=1}^{m-1} h[Y(t)|Y(t-1), \dots, Y(1)]. \end{aligned}$$

But from (41) the above right-hand side is nonpositive. Thus, $h[Y(1), \dots, Y(m)]/m$ is nonincreasing in m , as well as $h[Y(m)|Y(m-1), \dots, Y(1)]$, as implied by (7). This implies their convergence (a result already proved in [4]). \square

B. Proof of Proposition 1

The proof relies on the following result.

Lemma 6: Let $\{\mathbf{X}(t), t \in \mathbb{Z}\}$ be a vector random stationary process admitting α th absolute moment ($\alpha \geq 1$) and $\mathbf{Y}(t) = (\mathbf{B} \star \mathbf{X})(t)$ where $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ is a sequence of matrices with only a finite number of nonzero terms. Then

$$h[\mathbf{Y}(\cdot)] \geq h[\mathbf{X}(\cdot)] + \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}(l) e^{i l \lambda} \right| \frac{d\lambda}{2\pi}.$$

Proof: Consider the random vectors

$$\begin{aligned} \mathbf{Y}^{(m)}(t) &= \sum_{l=-\infty}^{\infty} \mathbf{B}(l) \mathbf{X}[(t-l) \pmod{m}] \\ &= \sum_{\tau=0}^{m-1} \left[\sum_{l=-\infty}^{\infty} \mathbf{B}(t-\tau+lm) \right] \mathbf{X}(\tau). \end{aligned}$$

This defines a linear transformation from $\mathbf{X}(0), \dots, \mathbf{X}(m-1)$ to $\mathbf{Y}^{(m)}(0), \dots, \mathbf{Y}^{(m)}(m-1)$ with the transformation matrix $\mathbf{B}^{(m)}$ being block circular Toeplitz with $\sum_{l=-\infty}^{\infty} \mathbf{B}(t-\tau+lm)$ at the (τ, t) place. Thus, by (12)

$$\begin{aligned} h[\mathbf{Y}^{(m)}(0), \dots, \mathbf{Y}^{(m)}(m-1)] \\ &= h[\mathbf{X}(0), \dots, \mathbf{X}(m-1)] + \log |\det \mathbf{B}^{(m)}|. \end{aligned}$$

To compute the determinant of $\mathbf{B}^{(m)}$, note that if \mathbf{U} is the block matrix with $e^{-i2\pi t\tau/m} \mathbf{I} / \sqrt{m}$ at the (τ, t) place, then \mathbf{U} is unitary and $\mathbf{U}^{-1} \mathbf{B}^{(m)} \mathbf{U}$ is block diagonal with diagonal blocks $\sum_{l=-\infty}^{\infty} e^{i2\pi l t/m} \mathbf{B}(l)$, $t = 0, \dots, m-1$. Therefore,

$$\log |\det \mathbf{B}^{(m)}| = \sum_{t=0}^{m-1} \log \left| \det \sum_{l=-\infty}^{\infty} e^{i2\pi l t/m} \mathbf{B}(l) \right|.$$

On the other hand, by assumption, $\mathbf{B}(l) = 0$ as soon as $|l|$ is greater than some integer, say q . Then it is easily seen that for $m > 2q$, $\mathbf{Y}^{(m)}(t) = \mathbf{Y}(t)$ if $t \in \{q, \dots, m-1-q\}$. Thus, since the mutual information is nonnegative

$$\begin{aligned} h[\mathbf{Y}^{(m)}(0), \dots, \mathbf{Y}^{(m)}(m-1)] \\ &\geq h[\mathbf{Y}(q), \dots, \mathbf{Y}(m-1-q)] + \sum_{0 \leq t < q, m-q \leq t < m} h[\mathbf{Y}^{(m)}(t)]. \end{aligned}$$

By assumption, the random vectors $\mathbf{Y}^{(m)}(t)$ admit α th absolute moment bounded by a constant not depending on m and t . We shall show in what follows that for a random vector Z with bounded α th absolute moment, $h(Z)$ is bounded above, regardless of the density of Z . From this and the preceding results

$$\begin{aligned} h[\mathbf{Y}(q), \dots, \mathbf{Y}(m-1-q)] &\geq h[\mathbf{X}(0), \dots, \mathbf{X}(m-1)] \\ &\quad + \sum_{t=0}^{m-1} \log \left| \det \sum_{l=-\infty}^{\infty} e^{i2\pi l t/m} \mathbf{B}(l) \right| - C \end{aligned}$$

C being a constant not depending on m . Dividing both sides of the preceding inequality by m , then letting $m \rightarrow \infty$, one gets the result of the lemma.

To complete the proof, we need to show the assertion mentioned earlier. Let p_Z be the density of Z and put

$q(z) = C_\alpha e^{-\|z\|^\alpha}$ where $\|\cdot\|$ denotes a vector norm and C_α is the normalizing constant so that q is a density. Then

$$\begin{aligned} h(Z) &= \mathbb{E} \log[q(Z)/p_Z(Z)] - \mathbb{E} \log q(Z) \\ &\leq \mathbb{E}[q(Z)/p_Z(Z) - 1] + \mathbb{E}\|Z\|^\alpha - \log C_\alpha. \end{aligned}$$

But the first term on the last right-hand side vanishes because q is a density, yielding the announced result. \square

Proof of the Proposition: Let m, ϵ , and δ be as in Definition 2 and $\{\mathbf{B}^\dagger(l), l \in \mathbb{Z}\}$ be the inverse sequence (with respect to the convolution) of $\{\mathbf{B}(l), l \in \mathbb{Z}\}$. Write $\mathbf{B}(l) = \mathbf{B}_n(l) + \bar{\mathbf{B}}_n(l)$ where $\mathbf{B}_n(l) = \bar{\mathbf{B}}(l)$ if $|l| \leq n$, = $\mathbf{0}$ otherwise. Then

$$\mathbf{Y}_n(t) = (\mathbf{B}_n * \mathbf{X})(t) = \mathbf{Y}(t) - (\bar{\mathbf{B}}_n * \mathbf{B}^\dagger * \mathbf{Y})(t).$$

Therefore, since

$$\sum_{l=-\infty}^{\infty} \|(\bar{\mathbf{B}}_n * \mathbf{B}^\dagger)(l)\| \leq \left[\sum_{|l|>n} \|\bar{\mathbf{B}}(l)\| \right] \left[\sum_{l=-\infty}^{\infty} \|\mathbf{B}^\dagger(l)\| \right]$$

one can choose n sufficiently large such that it is bounded by δ , and the continuity condition of Proposition 1 entails that $h[\mathbf{Y}_n(1), \dots, \mathbf{Y}_n(m)] \leq h[\mathbf{Y}(1), \dots, \mathbf{Y}(m)] + \epsilon$. Therefore, by Lemma 1, $h[\mathbf{Y}(1), \dots, \mathbf{Y}(m)] + \epsilon \geq mh[\mathbf{Y}_n(\cdot)]$. Applying now Lemma 6 to the process $\{\mathbf{Y}_n(t), t \in \mathbb{Z}\}$, one gets

$$\begin{aligned} h[\mathbf{Y}(1), \dots, \mathbf{Y}(m)]/m + \epsilon/m &\geq h[\mathbf{X}(\cdot)] \\ &+ \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}_n(l) e^{i\lambda l} \right| \frac{d\lambda}{2\pi}. \end{aligned}$$

Letting $n \rightarrow \infty$ and then $m \rightarrow \infty$, one gets the same inequality as in Lemma 6. But since $\mathbf{X}(t) = (\mathbf{B}^\dagger * \mathbf{Y})(t)$, one may apply the result just proved and obtains the reverse inequality

$$\begin{aligned} h[\mathbf{X}(\cdot)] &\geq h[\mathbf{Y}(\cdot)] + \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}^\dagger(l) e^{i\lambda l} \right| \frac{d\lambda}{2\pi} \\ &= h[\mathbf{Y}(\cdot)] - \int_{-\pi}^{\pi} \log \left| \det \sum_{j=-\infty}^{\infty} \mathbf{B}(l) e^{i\lambda l} \right| \frac{d\lambda}{2\pi}. \end{aligned}$$

It follows that the last inequality is an equality. \square

C. Proof of Lemma 2

For convenience, put

$$\begin{aligned} \mathbf{Y} &= [\mathbf{Y}^T(1) \dots \mathbf{Y}^T(m)]^T \\ \mathbf{Y}' &= [(\mathbf{C} * \mathbf{Y})^T(1) \dots (\mathbf{C} * \mathbf{Y})^T(m)]^T \end{aligned}$$

and denote by $p_{\mathbf{Y}}, p_{\mathbf{Y}'}$ their densities. Then

$$\begin{aligned} h(\mathbf{Y}') - h(\mathbf{Y}) &= \mathbb{E}[\log p_{\mathbf{Y}}(\mathbf{Y}) - \log p_{\mathbf{Y}'}(\mathbf{Y}')] \\ &\quad - \mathbb{E} \log [p_{\mathbf{Y}'}(\mathbf{Y}')/p_{\mathbf{Y}}(\mathbf{Y})]. \end{aligned}$$

The last term is nonnegative since it is a Kullback–Leibler divergence. As for the first term, by the mean value theorem and our assumption and the fact that

$$(a+b)^\lambda \leq 2^{\max(\lambda-1, 0)}(a^\lambda + b^\lambda)$$

for positive a, b, λ , one has

$$\begin{aligned} |\log p_{\mathbf{Y}}(\mathbf{Y}) - \log p_{\mathbf{Y}'}(\mathbf{Y}')| \\ \leq C[1 + 2^{\max(\alpha-2, 0)}(\|\mathbf{Y}\|^{\alpha-1} + \|\mathbf{Y}' - \mathbf{Y}\|^{\alpha-1})\|\mathbf{Y}' - \mathbf{Y}\|. \end{aligned}$$

But by the Hölder inequality, the last expression has expectation bounded by a constant times $\|\mathbf{C}(0) - \mathbf{I}\| + \sum_{l \neq 0} \|\mathbf{C}(l)\|$, which yields the result. \square

D. Proof of Lemma 3

Let $p_{Y_k(m)|Y_k(m-1), \dots, Y_k(1)}$ and $p_{\mathbf{Y}(m)|\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)}$ denote the conditional densities of $Y_k(m)$ given $Y_k(m-1), \dots, Y_k(m-1)$ and of $\mathbf{Y}(m)$ given $\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)$. Then the expected divergence mentioned in Lemma 3 is

$$\begin{aligned} &\mathbb{E} \log \frac{\prod_{k=1}^K p_{Y_k(m)|Y_k(m-1), \dots, Y_k(1)}[Y_k(1), \dots, Y_k(m)]}{p_{\mathbf{Y}(m)|\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)}[\mathbf{Y}_k(1), \dots, \mathbf{Y}_k(m)]}. \end{aligned}$$

But this expression can be easily seen to be

$$\begin{aligned} &\sum_{k=1}^K h[Y_k(m)|Y_k(m-1), \dots, Y_k(1)] \\ &\quad - h[\mathbf{Y}(m)|\mathbf{Y}(m-1), \dots, \mathbf{Y}(1)] \end{aligned}$$

and the last term (without the minus sign), by (6) and (12), equals $h[\mathbf{X}(m)|\mathbf{X}(m-1), \dots, \mathbf{Y}(1)] - \log \det \mathbf{B}$, yielding the result. \square

E. Proof of Proposition 2

Let $\{\mathbf{D}(l), l \in \mathbb{Z}\}$ be any sequence of diagonal matrices of class \mathcal{A} and put $\tilde{\mathbf{B}}(l) = (\mathbf{D} * \mathbf{B})(l)$, $\tilde{Y}_k(t) = (d_k * Y_k)(t)$ where $d_k(l)$ are the diagonal elements of $\mathbf{D}(l)$. Then one can write

$$\begin{aligned} C_1[\mathbf{B}(\cdot)] &= C_\infty[\tilde{\mathbf{B}}(\cdot)] + \sum_{k=1}^K \left\{ h[Y_k(1)] - h[\tilde{Y}_k(\cdot)] \right. \\ &\quad \left. + \int_{-\pi}^{\pi} \log \left| \det \sum_{l=-\infty}^{\infty} d_k(l) e^{i\lambda l} \right| \frac{d\lambda}{2\pi} \right\}. \end{aligned}$$

Each term in the sum in the above right-hand side is nonnegative since, by Corollary 2

$$h[\tilde{Y}_k(\cdot)] \leq h[Y_k(1)] + \int_{-\pi}^{\pi} \log \left| \sum_{l=-\infty}^{\infty} d_k(l) e^{i\lambda l} \right| \frac{d\lambda}{2\pi}. \quad (42)$$

Further, $C_\infty(\tilde{\mathbf{B}})$ is no other than $I[\tilde{Y}_1(\cdot), \dots, \tilde{Y}_K(\cdot)] + h[\mathbf{X}(\cdot)]$, hence, $C_1[\mathbf{B}(\cdot)]$ is bounded below by (the constant) $h[\mathbf{X}(\cdot)]$. This bound will be attained if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are temporally independent and independent among themselves and it is possible to choose $\{\mathbf{B}(l), l \in \mathbb{Z}\}$ so that these processes are so, since the sources are independent linear processes and the sequences $\{a_k(l), l \in \mathbb{Z}\}$ in their representation (23) are of class \mathcal{A} .

It remains to show that $C_1[\mathbf{B}(\cdot)]$ can attain its minimum only if the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are temporally independent and independent among themselves. We observe that, by Corollary 2, (42) can be an equality only if the process $\{Y_k(t), t \in \mathbb{Z}\}$ is temporally independent. Thus, $C_1[\mathbf{B}(\cdot)]$ can attain its minimum only if this happens for all k and $I[\tilde{Y}_1(\cdot), \dots, \tilde{Y}_K(\cdot)] = 0$. But

by Proposition 1, $I[\tilde{Y}_1(\cdot), \dots, \tilde{Y}_k(\cdot)] = I[Y_1(\cdot), \dots, Y_k(\cdot)]$ and because the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are temporally independent

$$I[Y_1(\cdot), \dots, Y_k(\cdot)] = \frac{1}{m} I\{[Y_1(1) \cdots Y_1(m)]^T, \dots, [Y_K(1) \cdots Y_K(m)]^T\}$$

for all m . Thus, the last right-hand side vanishes for all m and hence the processes $\{Y_k(t), t \in \mathbb{Z}\}$ are independent. \square

F. Proof of Lemma 4

Put $\mathbf{Y}' = \mathbf{Y} + \mathcal{E}\mathbf{Z}$, then by the same calculation as in the proof of Lemma 2 and using C1)

$$h(\mathbf{Y}') - h(\mathbf{Y}) = \mathbb{E}[\log p_{\mathbf{Y}}(\mathbf{Y}) - \log p_{\mathbf{Y}}(\mathbf{Y}')] + o(\mathcal{E}).$$

Therefore, one gets the result of the proposition if one has proved that

$$\mathbb{E}\{\log p_{\mathbf{Y}}(\mathbf{Y}) - \log p_{\mathbf{Y}}(\mathbf{Y}') - \psi_{\mathbf{Y}}^T(\mathbf{Y})\mathcal{E}\mathbf{Z}/\|\mathcal{E}\|\} \rightarrow 0, \quad \text{as } \|\mathcal{E}\| \rightarrow 0.$$

Since the random variable inside the curly bracket converges to 0 almost surely as $\|\mathcal{E}\| \rightarrow 0$, by the Lebesgue dominated convergence theorem, one needs only to show that it is bounded for all \mathcal{E} small enough by a fixed integrable random variable. But, repeating again the argument in the proof of Lemma 2, this random variable is bounded by

$$C[2 + 2^{\max(\alpha-2, 0)}(\|\mathbf{Y}\|^{\alpha-1} + c\|\mathbf{Z}\|^{\alpha-1}) + \|\mathbf{Y}\|^{\alpha-1}\|\mathbf{Z}\|].$$

for all \mathcal{E} such that $\|\mathcal{E}\| \leq c$. The last random variable is integrable by the Hölder inequality, yielding the result. \square

G. Proof of Proposition 3

For C_m to be minimized at \mathbf{B} one must have $C_m(\mathbf{B} + \mathcal{E}\mathbf{B}) \geq C_m(\mathbf{B})$ for all matrices \mathcal{E} . Take \mathcal{E} having a single nonzero element, say \mathcal{E}_{kj} , and put $\mathbf{Y}_k = [Y_k(1) \cdots Y_k(m)]^T$. One gets from (18)

$$[h(\mathbf{Y}_k + \mathcal{E}_{kj}\mathbf{Y}_j) - h(\mathbf{Y}_k)]/m - \log|\det(\mathbf{I} + \mathcal{E})| \geq 0.$$

But by Proposition 4, the first term in the preceding left-hand side equals $\mathcal{E}_{kj}\mathbb{E}[\mathbf{Y}_j\psi_{k,m}(\mathbf{Y}_k)]/m + o(\mathcal{E}_{kj})$ and if $k \neq j$, $\det(\mathbf{I} + \mathcal{E}) = 1$. Therefore, for \mathcal{E}_{kj} small enough

$$\mathcal{E}_{kj}\mathbb{E}[\mathbf{Y}_j\psi_{k,m}(\mathbf{Y}_k)] \geq 0$$

and since \mathcal{E}_{kj} is arbitrary, one must have $\mathbb{E}[\mathbf{Y}_j\psi_{k,m}(\mathbf{Y}_k)] = 0$, which yields (29).

Since $C_m^* = mC_m - (m-1)C_{m-1}$, a completely similar argument yields that a necessary condition for it to be minimized at \mathbf{B} is that (30) holds with

$$\psi_{k,m}^*[y(1), \dots, y(m)] = \psi_{k,m}[y(1), \dots, y(m)] - \begin{bmatrix} \psi_{k,m-1}[y(1), \dots, y(m-1)] \\ 0 \end{bmatrix}.$$

But one can easily see that $\psi_{k,m}^*$ as defined here is the same as the one in the proposition. \square

H. Proof of Proposition 4

Observe that

$$\log|\det(\mathbf{M} + \delta)| = \text{tr}(\mathbf{M}^{-1}\delta) + o(\|\delta\|)$$

as $\delta \rightarrow \mathbf{0}$, tr denoting the trace. Hence, when \mathbf{B} is changed to $\mathbf{B} + \mathcal{E}\mathbf{B}$ where \mathcal{E} is a matrix with only a nonzero term \mathcal{E}_{kj} with $j \neq k$, one gets from (21) that the change of $C_{g,m}$ is

$$\mathcal{E}_{kj}\text{tr}[\text{cov}(\mathbf{Y}_k)^{-1}\text{cov}(\mathbf{Y}_j, \mathbf{Y}_k)]/m + o(\mathcal{E}_{kj})$$

where we have put $\mathbf{Y}_k = [Y_k(1) \cdots Y_k(m)]^T$ while $\text{cov}(\cdot)$ and $\text{cov}(\cdot, \cdot)$ refer to the covariance and cross covariance matrices, respectively. For \mathbf{B} to maximize $C_{g,m}$, it is necessary that the last expression be nonpositive for all \mathcal{E}_{kj} . Therefore,

$$\text{tr}[\text{cov}(\mathbf{Y}_k)^{-1}\text{cov}(\mathbf{Y}_j, \mathbf{Y}_k)] = 0$$

which is no other than the condition (31).

On the other hand, by (22), the change of $C_{g,m}^*$ corresponding to the same above change of \mathbf{B} , is, putting $e_k(m) = Y_k(m) - Y_k(m|1 : m-1) = \sum_{l=0}^{m-1} a_{k,m-1}(l)Y_k(m-l)$

$$\text{cov}\left\{e_k(m), \sum_{l=1}^{m-1} [a_{k,m-1}(l)\mathcal{E}_{kj}Y_j(m-l) + \delta_{k,m-1}(l)Y_k(m-l)]\right\} / \text{var}\{e_k(m)\} + o(\mathcal{E}_{kj})$$

where $\delta_{k,m-1}(l)$ represents the change of $a_{k,m-1}(l)$ when $Y_k(t)$ is changed to $Y_k(t) + \mathcal{E}_{kj}Y_j(t)$. But the terms involving the $\delta_{k,m-1}(l)$ disappear since $\text{cov}\{e_k(m), Y_k(m-l)\} = 0$ by the definition of $Y_k(m|1 : m-1)$. Therefore, by the same argument as before, one obtains (32).

Finally, with the same change of \mathbf{B} as before, the process $\{Y_k(t), t \in \mathbb{Z}\}$ is changed to $\{Y_k(t) + \mathcal{E}_{kj}Y_j(t), t \in \mathbb{Z}\}$ while the processes $\{Y_i(t), t \in \mathbb{Z}\}$, $i \neq k$, are unchanged. It then follows from (20) that the corresponding change of C_g^* is

$$\mathcal{E}_{kj} \int_{-\pi}^{\pi} [f_{Y_k Y_j}(\lambda) / f_{Y_k}(\lambda)] d\lambda / (2\pi) + o(\mathcal{E}_{kj}).$$

This yields (33). \square

I. Proof of Proposition 5

Suppose that the contrast (34) is minimized at $\{\mathbf{B}(l), l \in \mathbb{Z}\}$, then adding to this sequence the sequence $\{(\mathcal{E} \star \mathbf{B})(l), l \in \mathbb{Z}\}$, where $\{\mathcal{E}(l), l \in \mathbb{Z}\}$ is any sequence of matrices, must not decrease this contrast. Choose this sequence to have only one nonzero term $\mathcal{E}(\tau)$ which has only one nonzero element $\mathcal{E}_{kj}(\tau)$, then the change of this contrast is

$$h[Y_k(1) + \mathcal{E}_{kj}(\tau)Y_j(1-\tau)] - h[Y_k(t)] - \int_{-\pi}^{\pi} \log|\det[\mathbf{I} + \mathcal{E}(\tau)e^{i\tau\lambda}]| \frac{d\lambda}{2\pi}.$$

If either $\tau \neq 0$ or $j \neq k$, the last term of the above expression is $o[\mathcal{E}_{kj}(\tau)]$ as $\mathcal{E}_{kj}(\tau) \rightarrow 0$. As for the contribution of the other terms, using Proposition 4, it can be seen to be $\mathcal{E}_{kj}(m)\mathbb{E}\{Y_j(1-\tau)\psi_k[Y_k(1)]\} + o[\mathcal{E}_{kj}(\tau)]$. This yields the first result of the proposition. The proof for the other result is similar. \square

J. Proof of Proposition 7

For the contrast (27), when \mathbf{B} is changed to $\mathbf{B} + \mathcal{E}\mathbf{B}$ with \mathcal{E} having a single nonzero element \mathcal{E}_{kj} , the corresponding change of this contrast is the infimum of

$$h[(b_k \star Y_k)(1) + \mathcal{E}_{kj}(b_k \star Y_j)(1)] - h[(b_k^* \star Y_k)(1)] \\ + \int_{-\pi}^{\pi} \log \left| \frac{\sum_{l=-\infty}^{\infty} b_k^*(l)e^{il\lambda}}{\sum_{l=-\infty}^{\infty} b_k(l)e^{il\lambda}} \right| \frac{d\lambda}{2\pi} - \log |\det(\mathbf{I} + \mathcal{E})|$$

over all sequences $\{b_k(l), l \in \mathbb{Z}\}$ of class \mathcal{A} . As \mathbf{B} minimizes the contrast, the above expression must be nonnegative; hence, taking $j \neq k$ (so that $\det(\mathbf{I} + \mathcal{E}) = 1$) and $b_k(l) = b_k^*(l)$

$$h[(b_k^* \star Y_k)(t) + \mathcal{E}_{kj}(b_k^* \star Y_j)(t)] - h[(b_k^* \star Y_k)(t)] \geq 0.$$

The above right-hand side, by Proposition 4, is

$$E\{(b_k^* \star Y_j)(1)\psi_k[(b_k^* \star Y_k)(1)]\} + o(\mathcal{E}_{kj})$$

yielding the first result of the proposition. The proof for the other result is similar. \square

REFERENCES

- [1] A. Belouchrani, K. Abed-Meraim, J. F. Cardoso, and E. Moulines, "A blind source separation technique using second-order statistics," *IEEE Trans. Signal Processing*, vol. 45, pp. 434–444, Feb. 1997.
- [2] J.-F. Cardoso, "Blind signal separation: Statistical principles," *Proc. IEEE*, vol. 86, pp. 2009–2026, Aug. 1998.
- [3] P. Comon, "Independent components analysis, a new concept," *Signal Processing*, vol. 36, no. 3, pp. 287–314, 1994.

- [4] T. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [5] V. P. Godambe, "Conditional likelihood and unconditional optimum estimating equations," *Biometrika*, vol. 63, pp. 277–284, 1963.
- [6] E. J. Hannan, *Multiple Time Series*. New York: Wiley, 1970.
- [7] C. Jutten and J. Héroult, "Blind separation of sources, Part I: An adaptive algorithm based on neuromimetic structure," *Signal Processing*, vol. 24, pp. 1–10, 1991.
- [8] D. T. Pham, "Blind separation of instantaneous mixture of sources via an independent component analysis," *IEEE Trans. Signal Processing*, vol. 44, pp. 2768–2779, Nov. 1996.
- [9] D. T. Pham and P. Garat, "Blind separation of mixtures of independent sources through a quasi maximum likelihood approach," *IEEE Trans. Signal Processing*, vol. 45, pp. 1712–1725, July 1997.
- [10] D. T. Pham, "Mutual information approach to blind separation of sources," in *Proc. ICA'99 Workshop*, Aussois, France, Jan. 1999, pp. 215–220.
- [11] —, "Blind separation of instantaneous mixture of sources via the Gaussian mutual information criterion," *Signal Processing*, vol. 81, pp. 855–870, 2001.
- [12] —, Contrast functions for ICA and sources separation. Tech. Rep. [Online]. Available: <http://www-lmc.imag.fr/lmc-sms/Dinh-Tuan.Pham/BSS/contrast.ps.gz>
- [13] —, "Contrast functions for blind separation and deconvolution of sources," in *Proc. ICA 2001 Conf.*, San Diego, CA, Dec. 2001.
- [14] A. Zygmund, *Trigonometric Series*. Cambridge, U.K.: Cambridge Univ. Press, 1968, vol. I.

Dinh-Tuan Pham (M'88) was born in Hanoi, VietNam, on February 10, 1945. He graduated from the School of Applied Mathematics and Computer Science (ENSIMAG) of the Polytechnic Institute of Grenoble, France, in 1968. He received the Ph.D. degree in statistics in 1975 from the University of Grenoble.

He was a Postdoctoral Fellow in the Department of Statistics, University of California, Berkeley during 1977–1978 and a Visiting Professor in the Department of Mathematics, University of Indiana at Bloomington during 1979–1980. He is currently Director of Recherche at the French Centre National de Recherche Scientifique, Grenoble. His researches includes time-series analysis, signal modeling, blind source separation, array processing, and biomedical signal processing.