

## What Is the Relation Between Slow Feature Analysis and Independent Component Analysis?

**Tobias Blaschke**

*t.blaschke@biologie.hu-berlin.de*

**Pietro Berkes**

*berkes@gatsby.ucl.ac.uk*

**Laurenz Wiskott**

*l.wiskott@biologie.hu-berlin.de*

*Institute for Theoretical Biology, Humboldt University Berlin, D-10115 Berlin, Germany*

**We present an analytical comparison between linear slow feature analysis and second-order independent component analysis, and show that in the case of one time delay, the two approaches are equivalent. We also consider the case of several time delays and discuss two possible extensions of slow feature analysis.**

### 1 Introduction ---

In data analysis, it is often desirable to transform the input signals into a new representation that recovers as much information as possible about the underlying processes. In the classical example of two people speaking simultaneously while being recorded with two microphones, for instance, the observed signal is a mixture of their voices. A more useful representation here would be one where each signal component contains only the information about a single speaker. In the visual domain, one might be interested in a representation that is invariant to typical transformations, such as translation or zoom. A variety of linear and nonlinear methods have been developed to extract the interesting features from an observed signal.

In this letter, we focus on two methods that consider different properties of the observed signal: independent component analysis (ICA) (see Hyvärinen, Karhunen, & Oja, 2001, for an overview) and slow feature analysis (SFA) (Wiskott & Sejnowski, 2002). ICA finds a representation of the data such that signal components are mutually statistically independent, which can be used to separate the two speakers in the example above. SFA extracts slowly-varying features, which can be used in the second example to learn visual invariances. At first glance, these two methods are very different and even seem to be conflicting, since two slowly varying signals of finite length are intuitively more likely to have statistical dependencies

than quickly varying ones. However, we will see that ICA and SFA have common properties, which we are going to point out by comparing the two algorithms mathematically.

To carry out the comparison, we have to apply some restrictions. SFA is constrained to nonwhite signals with a temporal structure (e.g., speech signals), and it is based on second-order statistics. We therefore compare it to ICA algorithms that use only second-order information and need a temporally structured signal as well (Molgedey & Schuster, 1994; Belouchrani, Abed Meraim, Cardoso, & Moulines, 1997; Ziehe & Müller, 1998; Zibulevsky & Pearlmutter, 2000; Nuzillard & Nuzillard, 2003). SFA is usually applied as a nonlinear method: it uses a nonlinear expansion to map the input signal into a feature space and then solves a linear problem there. ICA, on the other hand, is typically a linear method, since in the nonlinear case, the problem is in general underdetermined (because the solution is not unique), and there is thus no guarantee of recovering the original sources (Hyvärinen & Pajunen, 1999; Jutten & Karhunen, 2003). (However, there do exist some nonlinear approaches that make additional assumptions about the nonlinear mapping or the input data.) To make a comparison between the two methods possible, we will restrict SFA to the linear case. Nevertheless, all calculations in this letter are essentially the same for linear or nonlinear SFA.

## 2 Linear Mixing and Unmixing

---

Let  $\mathbf{x}(t) = [x_1(t), \dots, x_N(t)]^T$  be a linear mixture of a multidimensional source signal  $\mathbf{s}(t) = [s_1(t), \dots, s_N(t)]^T$ ,

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

where  $\mathbf{A}$  is a square mixing matrix and different components  $s_i$  come from statistically independent sources. In the following, we will assume that  $\mathbf{s}(t)$  and  $\mathbf{x}(t)$  have zero mean, without loss of generality. A common linear preprocessing step in many ICA algorithms as well as in linear SFA is the whitening of the input signal  $\mathbf{x}(t)$ . Whitening results in a signal  $\mathbf{y}(t) = \mathbf{W}\mathbf{x}(t)$  with mutually uncorrelated components,  $\langle y_i(t)y_j(t) \rangle = 0 \quad \forall i \neq j$ , unit variance,  $\langle y_i(t)^2 \rangle = 1$ , and zero mean,  $\langle y_i(t) \rangle = 0$ , where  $\langle \cdot \rangle$  denotes averaging over time. It can be shown that after the whitening step, an orthogonal transformation  $\mathbf{Q}$  on  $\mathbf{y}$  is sufficient to yield independent components (Comon, 1994) or slowly-varying features (Wiskott & Sejnowski, 2002). Therefore, the output signal  $\mathbf{u}(t)$  can be obtained by combining the whitening matrix  $\mathbf{W}$  and a rotation matrix  $\mathbf{Q}$ :

$$\mathbf{u}(t) = \mathbf{Q}\mathbf{y}(t) = \mathbf{Q}\mathbf{W}\mathbf{x}(t). \quad (2.2)$$

In the following, we will always assume whitened data  $\mathbf{y}(t)$  and focus on finding  $\mathbf{Q}$ . Since zero mean and whitening are preserved under any orthogonal transformation, the components of  $\mathbf{u}(t)$  also satisfy these conditions:

$$\langle u_i(t) \rangle = 0 \quad (\text{zero mean}) \quad (2.3)$$

$$\langle u_i(t)^2 \rangle = 1 \quad (\text{unit variance}) \quad (2.4)$$

$$\forall i \neq j : \langle u_i(t)u_j(t) \rangle = 0 \quad (\text{decorrelation}). \quad (2.5)$$

These properties fulfill the constraints imposed by SFA (cf. section 4) and are a good prerequisite for ICA because they constrain the output signals  $u_i(t)$  to be statistically independent in the first and second order.

### 3 Second-Order Independent Component Analysis

---

Given the linear mixture 2.1, ICA tries to retrieve the source signal components  $\mathbf{s}(t)$  from the input signal  $\mathbf{x}(t)$ . The mixing matrix  $\mathbf{A}$  is unknown, and the source signal components are assumed to be mutually independent. The typical approach is to define an objective function that is a measure of independence of the estimated source signal components  $u_i$ . The problem is then solved by optimizing this function with respect to  $\mathbf{Q}$ .

There exist different measures of independence. Most algorithms are based on the assumption that two signals are independent if their joint distribution is equal to the product of their marginals (e.g., Cardoso & Souloumiac, 1993; Hyvärinen, 1999; Lee, Girolami, & Sejnowski, 1999). A corresponding measure in this case is the Kullback-Leibler divergence. We will refer to this approach as *higher-order ICA*.

This definition, however, does not capture all aspects of independence. Consider a signal without temporal autocorrelation (e.g., white noise) and a second signal that is equal to the first one but shifted in time. Applying the measure of independence, the two signals appear to be independent, although they are actually a time-shifted copy of each other and thereby intuitively strongly dependent. This dependence across time can be taken into account using a different measure where two signals are considered statistically independent if all time-delayed correlations are zero (*second-order ICA*) (Molgedey & Schuster, 1994; Belouchrani et al., 1997; Ziehe & Müller, 1998). In order to successfully apply this measure, the source signals need to have a time structure (must be nonwhite), which is also a necessary condition for SFA. An alternative formulation of this idea is to use a model of the sources that includes a dynamics in time and assume that the time series are independent as a whole (Pearlmutter & Parra, 1996). In this letter, we are going to study algorithms based on this latter definition of independence, following the formulation by Molgedey and Schuster (1994).

To derive an objective function for second-order ICA, we first introduce time-delayed correlation matrices of the estimated source signal  $\mathbf{u}(t)$ ,

$$\mathbf{C}^{(\mathbf{u})}(\tau) := \langle \mathbf{u}(t)\mathbf{u}(t + \tau)^T \rangle, \quad (3.1)$$

where  $\tau$  is the time delay between two signals. We denote the entries of  $\mathbf{C}^{(\mathbf{u})}(\tau)$  as  $C_{ij}^{(\mathbf{u})}(\tau)$ . For a signal  $\mathbf{u}(t)$  with independent components,  $\mathbf{C}^{(\mathbf{u})}(\tau)$  should be diagonal for all  $\tau$ . We are therefore looking for an objective function that, when optimized, jointly diagonalizes those matrices.

It is common in practice to use a symmetrized version of the correlation matrices:<sup>1</sup>

$$\mathbf{C}^{(\mathbf{u})}(\tau) := \frac{1}{2} [\langle \mathbf{u}(t)\mathbf{u}(t + \tau)^T \rangle + \langle \mathbf{u}(t + \tau)\mathbf{u}(t)^T \rangle]. \quad (3.2)$$

Computing the symmetrized matrices is equivalent to applying the algorithm to the original input data and to the data reversed in time (because  $\langle \mathbf{u}(t + \tau)\mathbf{u}(t)^T \rangle = \langle \mathbf{u}(t)\mathbf{u}(t - \tau)^T \rangle$ ). This reflects the fact that with respect to the unmixing problem, the time direction is not important. Moreover, the symmetric form can always be diagonalized with a rotation matrix (while the nonsymmetric matrices can have complex eigenvalues and eigenvectors) and has better numerical properties. Note, however, that in some pathological cases, the cross-correlation terms can cancel out each other: For example, if  $\mathbf{u}(t) = [\sin(t), \cos(t)]^T$ , there clearly are cross-correlations but in the symmetrized version, the off-diagonal terms in equation 3.2 are zero for all  $\tau$ . The two signals are thus considered independent by the algorithm.

We first focus on the case of a single time delay  $\tau$  (Molgedey & Schuster, 1994). The extension to more than one time-delayed correlation matrix is straightforward and will be described in section 5. Because of the whitening step, equation 2.5, the correlation matrix with time delay zero is already diagonal. With one time delay, the ICA algorithm thus reduces to diagonalizing a single time-delayed correlation matrix  $\mathbf{C}^{(\mathbf{u})}(\tau)$ . This can be achieved by using the method of Jacobi (Cardoso & Souloumiac, 1996) to minimize the sum of the squared off-diagonal entries, a technique used in several second-order ICA algorithms (Belouchrani et al., 1997; Ziehe & Müller, 1998) as well as in methods based on higher-order statistics (Cardoso & Souloumiac, 1993). Using this method, we can define a simple objective

---

<sup>1</sup> In Ziehe and Müller (1998) the correlation matrices are not explicitly defined in the article, but the Matlab implementation made available by the authors uses the symmetric form.

function subject to minimization,

$$\Psi_{\text{ICA}}(\tau) := \sum_{\substack{i,j=1 \\ i \neq j}}^N (C_{ij}^{(\mathbf{u})}(\tau))^2 \tag{3.3}$$

$$= \sum_{i \neq j} (\mathbf{q}_i^T \mathbf{C}^{(\mathbf{y})}(\tau) \mathbf{q}_j)^2, \tag{3.4}$$

where  $\mathbf{q}_i$  is the  $i$ th row of  $\mathbf{Q}$ .  $\Psi_{\text{ICA}}$  is a function of the vectors  $\mathbf{q}_i$ , which are subject to learning, and of the whitened signal  $\mathbf{y}(t)$ , which is given. This objective function is optimized by a sequence of elementary rotations within the plane spanned by two axes. A possible optimization procedure has been described by Cardoso and Souloumiac (1996); a more efficient optimization schedule has been derived by Blaschke and Wiskott (2004a).

#### 4 Linear Slow Feature Analysis

---

Given a whitened input signal  $\mathbf{y}(t) = [y_1(t), \dots, y_N(t)]^T$ , linear SFA finds a rotation matrix  $\mathbf{Q}$  such that the components  $u_i$  of the output signal  $\mathbf{u}(t) = \mathbf{Q}\mathbf{y}(t)$  vary as slowly as possible in time and are ordered by decreasing slowness (the first one being the slowest possible, the second one the next slowest uncorrelated to the first, and so on). As a measure of slowness, we define (small values indicating slowly varying signals)

$$\Delta(u_i) := \langle \dot{u}_i(t)^2 \rangle, \tag{4.1}$$

which has to be minimized (Wiskott & Sejnowski, 2002). Due to the earlier whitening step, each output signal  $u_i(t)$  has zero mean and unit variance (see equations 2.3 and 2.4). This ensures that the solution will not be the trivial solution  $u_i(t) = \text{const}$ . The decorrelation of the output signals, equation 2.5, guarantees that different components carry different information.

We first show how to solve the optimization problem of SFA in a way similar to that described by Wiskott and Sejnowski (2002) and then establish a link between SFA and second-order ICA. For discrete time series, the first derivative of  $\mathbf{u}(t)$  can be approximated in the first order by

$$\dot{\mathbf{u}}(t) \approx \mathbf{u}(t + 1) - \mathbf{u}(t). \tag{4.2}$$

Using this approximation, we can rewrite the SFA objective function, equation 4.1, as

$$\begin{aligned} \Delta(u_i) &\approx \langle (u_i(t + 1) - u_i(t))^2 \rangle \\ &= \langle u_i(t + 1)u_i(t + 1) \rangle + \langle u_i(t)u_i(t) \rangle \end{aligned} \tag{4.3}$$

$$- \langle u_i(t)u_i(t + 1) \rangle - \langle u_i(t + 1)u_i(t) \rangle \tag{4.4}$$

$$= 2\langle u_i(t)^2 \rangle - 2\langle u_i(t)u_i(t + 1) \rangle \tag{4.5}$$

(since  $\langle u_i(t + 1)^2 \rangle = \langle u_i(t)^2 \rangle$  because we average over all  $t$ )

$$= 2 - 2\langle u_i(t)u_i(t + 1) \rangle \tag{4.6}$$

(since  $\langle u_i(t)^2 \rangle = 1$  because  $\mathbf{u}(t)$  is white (see equation 2.4))

Since the constant factor does not matter during optimization, instead of minimizing  $\Delta(u_i)$ , we can maximize

$$\tilde{\Delta}(u_i) := 1 - \frac{1}{2} \Delta(u_i) \tag{4.7}$$

$$= \langle u_i(t)u_i(t + 1) \rangle \tag{4.8}$$

$$= C_{ii}^{(u)}(1) \tag{4.9}$$

$$= \mathbf{q}_i^T \mathbf{C}^{(y)}(1) \mathbf{q}_i. \tag{4.10}$$

The objective function  $\tilde{\Delta}(u_i)$  is a function of the rotation matrix  $\mathbf{Q}$ , and we are thus searching for the orthogonal weight vectors  $\mathbf{q}_i$  in equation 4.10 that maximize  $\tilde{\Delta}(u_i)$ . The solution for  $i = 1$  is obviously the eigenvector of the largest eigenvalue of  $\mathbf{C}^{(y)}(1)$ , which yields the slowest component  $u_1(t) = \mathbf{q}_1^T \mathbf{y}(t)$ . The following eigenvectors in order of decreasing eigenvalue yield the next-slowest components,  $u_2(t)$ ,  $u_3(t)$ , and so forth.

Therefore, to extract all slow components, the maximization problem, equation 4.10, can be formulated as an eigenvalue problem,

$$\mathbf{C}^{(y)}(1) \mathbf{Q}^T = \mathbf{Q}^T \mathbf{\Lambda}, \tag{4.11}$$

where  $\mathbf{\Lambda}$  denotes a diagonal matrix with  $\Lambda_{ii}$  being the  $i$ th largest eigenvalue and  $\mathbf{q}_i$  the corresponding eigenvectors.

In order to allow a better comparison with second-order ICA, we now want to deduce an alternative formulation of SFA; that is, we want to construct an objective function similar to that of second-order ICA. First, we show the equivalence of solving the eigenvalue problem, equation 4.11, and the diagonalization of  $\mathbf{C}^{(u)}(1)$ . If we multiply both sides of equation 4.11 with  $\mathbf{Q}$ , we obtain

$$\mathbf{C}^{(u)}(1) = \mathbf{Q} \mathbf{C}^{(y)}(1) \mathbf{Q}^T = \mathbf{\Lambda}. \tag{4.12}$$

Since  $\mathbf{\Lambda}$  is diagonal  $\mathbf{C}^{(u)}(1)$  is diagonal too. Therefore, solving the eigenvalue problem for  $\mathbf{C}^{(y)}(1)$  is equivalent to finding a rotation matrix  $\mathbf{Q}$  such that the time-delayed correlation matrix  $\mathbf{C}^{(u)}(1)$  is diagonal. Second, to

perform the diagonalization, we minimize all off-diagonal entries of  $\mathbf{C}^{(u)}(1)$  using the same Jacobi scheme as for second-order ICA (see section 3) and define the following objective function for SFA:

$$\tilde{\Psi}_{\text{SFA}} := \sum_{i \neq j} (C_{ij}^{(u)}(1))^2 \quad (4.13)$$

$$= \sum_{i \neq j} (\mathbf{q}_i^T \mathbf{C}^{(y)}(1) \mathbf{q}_j)^2. \quad (4.14)$$

Minimizing this expression produces the same slow components  $u_1(t), \dots, u_N(t)$  as obtained by the eigenvalue problem, equation 4.11, again assuming an additional sorting step. Note also that this is equivalent to a decorrelation of the time derivatives of the output signal components  $u_i(t)$  (cf. Wiskott, 2003) since  $\langle \dot{u}_i \dot{u}_j \rangle = -2 C_{ij}^{(u)}(1)$  for  $i \neq j$ .

Interestingly, the objective function, equation 4.14, is identical to the one for ICA, equation 3.4. With this observation, we arrive at the important result that linear SFA is formally equivalent to second-order ICA with time delay one.

To bring equation 4.13 into a form that can be understood more intuitively in the sense of SFA, we can use the fact that the sum of all squared entries of correlation matrices with a given time delay  $\tau$  is invariant under orthogonal transformations,

$$\sum_{i,j} (C_{ij}^{(u)}(\tau))^2 = \sum_{i,j} (C_{ij}^{(y)}(\tau))^2 = \text{const.} \quad (4.15)$$

We can split this sum in two terms,

$$\sum_{i,j} (C_{ij}^{(u)}(\tau))^2 = \sum_i (C_{ii}^{(u)}(\tau))^2 + \sum_{i \neq j} (C_{ij}^{(u)}(\tau))^2 = \text{const.}, \quad (4.16)$$

so that it is easy to see that the minimization of  $\tilde{\Psi}_{\text{SFA}}$  is equivalent to the maximization of

$$\Psi_{\text{SFA}} := \sum_i (C_{ii}^{(u)}(1))^2 \quad (4.17)$$

$$= \sum_i (\mathbf{q}_i^T \mathbf{C}^{(y)}(1) \mathbf{q}_i)^2. \quad (4.18)$$

Having started from minimizing temporal variations, equation 4.1, as an objective for SFA, we now arrive at an objective for maximizing squared autocorrelations, equation 4.7, at time delay one. This relation can be

interpreted intuitively. A signal component with a large squared autocorrelation has a high temporal predictability. If the autocorrelation is positive (i.e.,  $C_{ii}^{(u)}(1) > 0$ ), predictability implies that the signal component has to vary slowly.

What if the autocorrelation is negative? This could happen if, for example,  $u_i(t)$  has alternating signs for successive data points. Consider the signal

$$u_i(t) := \begin{cases} -1 & \text{for } t \text{ odd} \\ 1 & \text{for } t \text{ even} \end{cases}, \quad (4.19)$$

with  $1 \leq t \leq T$ . This signal has zero mean and unit variance and thus fulfills constraints 2.3 and 2.4. Furthermore, it is favorable in terms of the objective 4.17, since  $C_{ii}^{(u)}(1)$  has a large absolute value. On the other hand, this is a very quickly varying component, which might seem paradoxical since maximizing equation 4.17 should result in slowly varying components. This apparent contradiction can be resolved by studying the constraints imposed on the optimization of equation 4.17. Since  $\mathbf{Q}$  is an orthogonal matrix, the trace of  $\mathbf{C}^{(u)}(1)$  is invariant under the transformation  $\mathbf{u}(t) = \mathbf{Q}\mathbf{y}(t)$  (e.g., Zurmühl & Falk, 1997). If we consider all  $N$  possible components in the optimization procedure, the decrease of one correlation  $C_{ii}^{(u)}(1)$  implies the increase of at least one other correlation  $C_{jj}^{(u)}(1)$ . Therefore, extracting the most slowly varying signals implies that other extracted components correspond to the most quickly varying signals. Hence, it is reasonable to further minimize negative correlations since this implies that other correlations will be maximized. As above, a successive sorting step is required to bring the components in order of increasing temporal variation.

## 5 More Than One Time Delay

---

**5.1 Second-Order ICA.** We know that second-order ICA can always be solved with a single time delay (Tong, Liu, Soon, & Huang, 1991). However, the delay  $\tau$  has to be chosen properly so that all eigenvalues of  $\mathbf{C}^{(y)}(\tau)$  are distinct. To obtain a more robust method, one can consider a certain number  $T$  of time-delayed correlation matrices with respective time delays  $\tau = 1, 2, \dots, T$  and diagonalize them jointly (Belouchrani et al., 1997; Ziehe & Müller, 1998). This leads to a straightforward extension of objective 3.3, subject to minimization,

$$\Psi_{\text{ICA}_j} := \sum_{\tau=1}^T \kappa_{\tau} \Psi_{\text{ICA}}(\tau) \quad (5.1)$$

$$= \sum_{\tau} \kappa_{\tau} \sum_{i \neq j} (C_{ij}^{(u)}(\tau))^2 \tag{5.2}$$

$$= \sum_{\tau} \kappa_{\tau} \sum_{i \neq j} (\mathbf{q}_i^T \mathbf{C}^{(y)}(\tau) \mathbf{q}_j)^2, \tag{5.3}$$

where we introduced positive factors  $\kappa_{\tau}$  that allow us to weight correlation matrices with different time delays differently. In equation 5.1, we write ICAj for *joint-diagonalization ICA*. Pham and Garat (1997) have derived a formula closely related to equation 5.3 with a maximum likelihood approach.

Extending the objective function of ICA in this way leads to the joint diagonalization of several correlation matrices with different time delays. Decorrelation is thus achieved over a time window of length  $T$ . It is intuitively clear that by enlarging the window length, the unmixing performance should improve until the width of the autocorrelation function is reached. Exceeding this limit would introduce matrices consisting entirely of zero mean noise, which would degrade the unmixing performance.

## 5.2 Linear SFA

*5.2.1 Joint Diagonalization.* We can use an argument similar to the one used for second-order ICA in order to extend SFA to more than a single time delay. Adding more time-delayed autocorrelations increases the temporal predictability of the signal. Knowing the amplitude of a signal at a given time can give a good prediction for the next  $T$  time points since they are strongly correlated. Signals with large temporal predictability are in turn likely to be slowly varying (cf. the end of section 4). Thus, an intuitive extension of the normal SFA objective, equation 4.17, subject to maximization, is

$$\Psi_{\text{SFAj}} := \sum_{\tau} \kappa_{\tau} \Psi_{\text{SFA}}(\tau) \tag{5.4}$$

$$= \sum_{\tau} \kappa_{\tau} \sum_i (C_{ii}^{(u)}(\tau))^2 \tag{5.5}$$

$$= \sum_{\tau} \kappa_{\tau} \sum_i (\mathbf{q}_i^T \mathbf{C}^{(y)}(\tau) \mathbf{q}_i)^2. \tag{5.6}$$

As in equations 5.1 to 5.3, we have introduced weighting factors  $\kappa_{\tau}$  for the delayed correlation matrices. Note that this new objective, equations 5.5 and 5.6, is again equivalent to the ICA objective, equations 5.2 and 5.3, due to the constancy of the sum of all squared entries of each time-delayed correlation matrix, equation 4.16.

We must be careful with this definition for two reasons. Firstly, while the definition of slowness based on  $C_{ii}^{(u)}(1)$  corresponds to our intuition of what a slow signal is,  $C_{ii}^{(u)}(2)$  can have a large, positive value for signal components that we would not consider to be slow at all. In fact, the alternating signal, equation 4.19, would yield a maximal value for  $C_{ii}^{(u)}(2)$ . Secondly, consider the case where two time-delayed autocorrelations have opposite signs, for example,  $C_{ii}^{(u)}(1) < 0$  and  $C_{ii}^{(u)}(2) > 0$ . Maximizing objective function 5.5 would favor a decreasing value of  $C_{ii}^{(u)}(1)$  (since it is negative) and an increasing value of  $C_{ii}^{(u)}(2)$ . The former would intuitively tend to make the signal faster, while the latter would make it slower. Thus, if the autocorrelations of a component have different signs for different time delays, the objective function appears to be inconsistent, at least for that component. This conflict cannot be solved as easily as the one discussed at the end of section 4. However, one can at least monitor the signs of the autocorrelations and diagnose the inconsistent cases. It is not clear to us how often these two problems arise in practice. We believe that by weighting the first autocorrelation stronger than the others, for example, with an exponential decay of the weights, the inconsistencies can be largely avoided.

5.2.2 *Linear Filtering.* An alternative to the joint diagonalization of several correlation matrices with different time delays in analogy to second-order ICA is to average over a range of time delays within one correlation matrix and diagonalize just this one matrix. To do so, we introduce the following new measure of slowness (cf. equations 4.7 to 4.10):

$$\tilde{\Sigma}(u_i) := \left\langle u_i(t) \left( \sum_{\tau} \kappa_{\tau} u_i(t + \tau) \right) \right\rangle \tag{5.7}$$

$$= \sum_{\tau} \kappa_{\tau} \langle u_i(t) u_i(t + \tau) \rangle \tag{5.8}$$

$$= \sum_{\tau} \kappa_{\tau} C_{ii}^{(u)}(\tau) \tag{5.9}$$

$$= \mathbf{q}_i^T \left( \sum_{\tau} \kappa_{\tau} \mathbf{C}^{(y)}(\tau) \right) \mathbf{q}_i \tag{5.10}$$

$$=: \mathbf{q}_i^T \tilde{\mathbf{C}}^{(y)} \mathbf{q}_i, \tag{5.11}$$

with constants  $\kappa_{\tau}$  that weight different time delays differently. This definition differs from that of equations 4.7 to 4.10 in that  $u_i(t)$  should not only be well correlated to the next data point but to a weighted average over the next  $T$  data points. This is a straightforward way of taking several timescales into account. Note that the weighted averaging is a linear filter operation.

As in the joint diagonalization extension, exponentially decaying weights  $\kappa_\tau := \exp(-\gamma\tau)$  for different time delays seem to be a suitable choice. With such weights, this measure of slowness is similar to the objective of temporal smoothness used by Stone (1995) and somewhat related to the trace learning rules introduced by Földiák (1991).

Because of the similarity of equation 5.11 with equation 4.10, we can apply the steps that led from equation 4.10 to equation 4.18 and derive the following objective function to be maximized,

$$\Psi_{\text{SFAI}} := \sum_i (\tilde{C}_{ii}^{(u)})^2 \tag{5.12}$$

$$= \sum_i (\mathbf{q}_i^T \tilde{\mathbf{C}}^{(y)} \mathbf{q}_i)^2, \tag{5.13}$$

where  $\tilde{\mathbf{C}}^{(u)}$  is defined analogously to  $\tilde{\mathbf{C}}^{(y)}$  and SFAI stands for *linear-filtering SFA*. Since this objective function is based on just one correlation matrix, it does not have the problems mentioned above for the joint diagonalization extension (see section 5.2.1).

Blaschke (2005, sec. 8.2.2) also considered extending SFA by simultaneously minimizing the variance not only of the first but also of higher-order derivatives, which could result in even more stable signals. This would also lead to equation 5.13, because discrete approximations of higher-order derivatives involve multiple time delays. In this case, with positive weights for all derivatives, the constants  $\kappa_\tau$  in equation 5.10 would have values with alternating signs (positive for odd  $\tau$  and negative for even  $\tau$ ), which is somewhat counterintuitive. We do not fully understand the implications of this effect but believe that higher-order derivatives do not offer a good way of extending SFA to longer timescales, even though unmixing performance was good in some simple examples.

## 6 Conclusion

---

The main result of this work is that linear SFA and second-order ICA with time delay one are formally equivalent (see equations 3.4 and 4.14). This is surprising, because SFA and ICA are based on two very different principles: slowness versus statistical independence. These principles might seem to contradict each other, because two analog signals of finite length would typically become more statistically dependent if they are more slowly varying.

The formal equivalence of linear SFA and second-order ICA with time delay allows us to apply the intuition we have gained for one algorithm to deepen our understanding of the other. For example, it is known that higher-order ICA applied to natural images learns linear filters similar to Gabor wavelets (e.g., Bell & Sejnowski, 1997; van Hateren & van der Schaaf,

1998), which in turn resemble receptive fields of simple cells in V1. On the other hand, linear SFA (and therefore also second-order ICA with time delay one) applied to natural image sequences learns filters similar to the principal components of natural images, the first of which are effectively spatial low-pass filters and therefore also generate slowly varying output signals. This suggests that the solutions found by second-order ICA and higher-order ICA can be very different in practice even though both methods try to maximize statistical independence.

Despite the formal equivalence in the linear case and for time delay one, SFA and ICA have different objectives and differ in the more general case.

Firstly, while in standard SFA the time delay is fixed to 1 due to the approximation of the time derivative, in ICA it can be chosen freely, or one can use several correlation matrices with different time delays simultaneously for optimal unmixing (see section 5.1). We have seen (see section 5.2.1) that the same extension to several time delays can also be used for SFA, but that the algorithm then becomes inconsistent with respect to the slowness objective if the entries of the time-delayed correlation matrices have different signs for different delays. An extension more consistent with the slowness objective is based on linear filtering before computing the time derivative (see section 5.2.2). This also introduces several time delays, but in a different way than used for ICA. Thus, when taking several time delays into account, the conceptual differences between ICA and SFA become relevant.

Secondly, in the nonlinear case, many output signal components can be extracted from a lower-dimensional input signal. With SFA, they would all be uncorrelated and ordered by slowness, in agreement with the definition in equations 2.3 to 2.5 and 4.1. With second-order ICA, they would not be ordered in any way and would not be statistically independent for dimensionality reasons. The results would therefore be inconsistent with the ICA objective. Thus, in the nonlinear case, the conceptual differences between ICA and SFA also matter.

We believe that the close relation between linear SFA and second-order ICA will lead to a way to combine the two algorithms into a nonlinear method for extracting slowly varying and statistically independent components and thereby perform nonlinear blind source separation. This is the subject of current research (Blaschke & Wiskott, 2004b, 2005).

## Acknowledgments

---

This work has been supported by a grant to L.W. from the Volkswagen Foundation.

## References

---

- Bell, A. J., & Sejnowski, T. J. (1997). The “independent components” of natural scenes are edge filters, *Vision Research*, 37(23), 3327–3338.

- Belouchrani, A., Abed Meraim, K., Cardoso, J.-F., & Moulines, E. (1997). A blind source separation technique based on second order statistics. *IEEE Transactions on Signal Processing*, 45(2), 434–44.
- Blaschke, T. (2005). *Independent component analysis and slow feature analysis: Relations and combination*. Doctoral dissertation, Humboldt-University Berlin. Available online at <http://edoc.huberlin.de/docviews/abstract.php?lang=ger&id=25458>.
- Blaschke, T., & Wiskott, L. (2004a). CuBICA: Independent component analysis by simultaneous third- and fourth-order cumulant diagonalization. *IEEE Transactions on Signal Processing*, 52(5), 1250–1256.
- Blaschke, T., & Wiskott, L. (2004b). Independent slow feature analysis and nonlinear blind source separation. In *Proc. of the 5th Int. Conf. on Independent Component Analysis and Blind Signal Separation*. Berlin: Springer-Verlag.
- Blaschke, T., & Wiskott, L. (2005). Nonlinear blind source separation by integrating independent component analysis and slow feature analysis. In L. K. Saul, Y. Weiss, & L. Bottou (Eds.), *Advances in neural information processing systems*, 17, (pp. 177–184). Cambridge, MA: MIT Press.
- Cardoso, J.-F., & Souloumiac, A. (1993). Blind beamforming for non gaussian signals, *IEEE Proceedings-F*, 140, 362–370.
- Cardoso, J.-F., & Souloumiac, A. (1996). Jacobi angles for simultaneous diagonalization. *SIAM J. Mat. Anal. Appl.*, 17(1), 161–164.
- Comon, P. (1994). Independent component analysis, a new concept? *Signal Processing*, 36(3), 287–314.
- Földiák, P. (1991). Learning invariance from transformation sequences. *Neural Computation*, 3(2), 194–200.
- Hyvärinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10(3), 626–634.
- Hyvärinen, A., Karhunen, J., & Oja, E. (2001). *Independent component analysis*. New York: Wiley.
- Hyvärinen, A., & Pajunen, P. (1999). Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, 12(3), 429–439.
- Jutten, C., & Karhunen, J. (2003). Advances in nonlinear blind source separation. In *Proc. of the 4th Int. Symposium on Independent Component Analysis and Blind Signal Separation* (pp. 245–256). Available online at <http://www.kecl.ntt.co.jp/icl/signal2003/index.html>.
- Lee, T.-W., Girolami, M. & Sejnowski, T. (1999). Independent component analysis using an extended Infomax algorithm for mixed sub-gaussian and super-gaussian sources. *Neural Computation*, 11(2), 409–433.
- Molgedey, L., & Schuster, G. (1994). Separation of a mixture of independent signals using time-delayed correlations. *Physical Review Letters*, 72(23), 3634–3637.
- Nuzillard, D., & Nuzillard, J.-M. (2003). Second-order blind source separation in the Fourier space of data. *Signal Processing*, 83(3), 627–631.
- Pearlmutter, B., & Parra, L. (1996). A context-sensitive generalization of ICA. In *Proc. of the International Conference on Neural Information Processing*. Berlin: Springer-Verlag.
- Pham, D., & Garat, P. (1997). Blind separation of mixtures of independent sources through a maximum likelihood approach. *IEEE Transactions on Signal Processing*, 45(7), 1712–1725.

- Stone, J. (1995). A learning rule for extracting spatio-temporal invariances. *Network*, 6(3), 1–8.
- Tong, L., Liu, R., Soon, V. C., & Huang, Y.-F. (1991). Indeterminacy and identifiability of blind identification. *IEEE Transactions on Circuits and Systems*, 38(5), 499–509.
- van Hateren, J., & van der Schaaf, A. (1998). Independent component filters of natural images compared with simple cells in primary visual cortex. *Proc. R. Soc. Lond. B*, 265, 359–366.
- Wiskott, L. (2003). Slow feature analysis: A theoretical analysis of optimal free responses. *Neural Computation*, 15(9), 2147–2177.
- Wiskott, L., & Sejnowski, T. (2002). Slow feature analysis: Unsupervised learning of invariances. *Neural Computation*, 14(4), 715–770.
- Zibulevsky, M., & Pearlmutter, B. (2000). Second order blind source separation by recursive splitting of signal subspaces. In *Proc. of the 2nd Int. Workshop on Independent Component Analysis and Blind Signal Separation* (pp. 489–491). Available online at <http://www.cis.hut.fi/ica2000>.
- Ziehe, A., & Müller, K.-R. (1998). TDSEP—an efficient algorithm for blind separation using time structure. In *Proc. of the 8th Int. Conference on Artificial Neural Networks* (pp. 675–680). Berlin: Springer-Verlag.
- Zurmühl, A., & Falk, S. (1997). *Matrizen und ihre Anwendungen* (Vol. 1, 6th ed.). Berlin: Springer-Verlag.