

Blind Separation of Nonlinear Mixtures by Variational Bayesian Learning

Antti Honkela^{a,*} Harri Valpola^b Alexander Ilin^a
Juha Karhunen^a

^a*Adaptive Informatics Research Centre, Helsinki University of Technology
P.O. Box 5400, FI-02015 TKK, Finland*

^b*Laboratory of Computational Engineering, Helsinki University of Technology
P.O. Box 9203, FI-02015 TKK, Finland*

Abstract

Blind separation of sources from nonlinear mixtures is a challenging and often ill-posed problem. We present three methods for solving this problem: an improved nonlinear factor analysis (NFA) method using multilayer perceptron (MLP) network to model the nonlinearity, a hierarchical NFA (HNFA) method suitable for larger problems and a post-nonlinear NFA (PNFA) method for more restricted post-nonlinear mixtures. The methods are based on variational Bayesian learning, which provides the needed regularisation and allows for easy handling of missing data. While the basic methods are incapable of recovering the correct rotation of the source space, they can discover the underlying nonlinear manifold and allow reconstruction of the original sources using standard linear independent component analysis (ICA) techniques.

Key words: Bayesian learning, blind source separation, nonlinear mixtures, post-nonlinear mixtures, variational Bayes

PACS:

1 Introduction

Independent component analysis (ICA) [1], that is blind separation of independent sources from their linear mixtures, forms the basis for most source sep-

* Corresponding author. Tel. +358-9-451 5283, Fax +358-9-451 3277

Email address: antti.honkela@tkk.fi (Antti Honkela).

URL: <http://www.cis.hut.fi/projects/bayes/> (Antti Honkela).

aration work today. Introduction of Bayesian methods based on the variational approximation for models capable of solving the ICA problem in 1999 independently by Attias [2] and Valpola¹ [3] was an important step for Bayesian source separation. Later, slightly different methods based on the same principles have been proposed by a number of authors [4–7].

The Bayesian modelling approach is very useful in source separation, as in addition to providing tools for model order selection, it is open to all kinds of extensions simply by extending the model [8,6,7]. This was illustrated, for instance, by the nonlinear factor analysis (NFA) and nonlinear independent factor analysis (NIFA) models for blind separation of sources from nonlinear mixtures by Valpola and Honkela [9]. While the original NFA/NIFA model could be used to solve the difficult nonlinear separation problem in many cases, we have developed new variants of the model to better deal with certain specific instances.

The hierarchical nonlinear factor analysis (HNFA) method [10] was introduced to overcome the quadratic computational complexity with respect to the number of the sources in original NFA/NIFA, thus allowing solving larger problems. The special case of post-nonlinear mixtures was addressed by post-nonlinear factor analysis (PNFA) [11]. Other important general extensions include ability to handle missing data [12,13], improved initialisation strategies [14] as well as more accurate approximation of the nonlinearity in NFA/NIFA leading to increased stability [15]. In this paper we summarise these new results presented in scattered conference papers, and present the improved methods in full detail together with new experimental comparisons.

The paper is organised as follows. In Sec. 2, the nonlinear blind source separation problem and some basic theory behind it are introduced, followed by a review of other existing methods in Sec. 3. This is followed by presentation of our models to solve the problem in Sec. 4. Learning the models is discussed in Sec. 5. Some experimental results on the methods are presented in Sec. 6 with discussion and conclusion in Secs. 7 and 8.

2 Nonlinear Blind Source Separation

Blind separation of sources from their nonlinear mixtures—known as nonlinear blind source separation (BSS)—is generally a difficult problem, both from a theoretical and a practical point of view [1,16,17]. The task is to extract the sources $\mathbf{s}(t)$ that have generated the observations $\mathbf{x}(t)$ through a nonlinear

¹ At the time his surname was Lappalainen.

mapping $\mathbf{f}(\cdot)$:

$$\mathbf{x}(t) = \mathbf{f}[\mathbf{s}(t)] + \mathbf{n}(t). \quad (1)$$

Here $\mathbf{x}(t)$ is the observed N -dimensional data (mixture) vector at time instant or index value t , $\mathbf{f} : \mathbb{R}^M \rightarrow \mathbb{R}^N$ is an unknown mixing function, $\mathbf{s}(t)$ is an M -vector whose elements are the M unknown source signals $s_i(t)$, $i = 1, 2, \dots, M$, and $\mathbf{n}(t)$ is an N -dimensional vector containing the additive noise terms.

In the basic linear case

$$\mathbf{x}(t) = \mathbf{A}\mathbf{s}(t) \quad (2)$$

where \mathbf{A} is an $N \times M$ mixing matrix, the BSS problem can be solved using ICA [1] under certain conditions. Linear ICA methods utilise the strong but often realistic assumption that the source signals are statistically independent. Furthermore, all but possibly one of the source signals are assumed to be non-Gaussian, there is no or only a little noise and the number of sources M equals the number N of mixtures [1]. Under these conditions, the sources in (2) can be estimated using ICA up to permutation, scaling, and sign.

Linear ICA and other BSS methods have been extended for the cases where the number of the sources is different from the number of mixtures, and for the noisy case. But the extension to the nonlinear mixture model (1) is difficult since both the nonlinear mapping and the underlying sources must be learned from the data in a blind manner, and the problem is highly ill-posed without suitable regularisation [18,1,16]. The independence assumption used in ICA is no longer sufficient for achieving separation. It can be shown that in the nonlinear case x and y can be mixed and still be statistically independent [18]. Additionally if x and y are two independent random variables, any of their functions $f(x)$ and $g(y)$ are also independent. These uniqueness issues have been considered in more detail in [18,16].

A related problem is that it is often quite difficult to infer the number of sources and the structure of the mapping $\mathbf{f}(\cdot)$. From a practical point of view, efficiency and reliability of nonlinear BSS algorithms are critical issues. They have restricted the number of sources that can be separated in practice to be quite small in many instances. An important reason for this is that the computational load of many nonlinear BSS methods increases quite rapidly with the dimensionality of the problem, preventing in practice their application to high-dimensional data sets [1].

An important special case of the general nonlinear mixing model (1) consists of so-called post-nonlinear mixtures. There each mixture has the form

$$x_i(t) = f_i \left(\sum_{j=1}^M A_{ij} s_j(t) \right), \quad i = 1, \dots, N. \quad (3)$$

Thus the sources $s_j(t)$, $j = 1, \dots, M$, are first mixed linearly according to

the basic linear ICA/BSS model (2), but after that nonlinear functions f_i are applied to them to get the final observations x_i . It can be shown [19,16] that for the noiseless post-nonlinear mixtures, the indeterminacies are usually the same as for the basic linear instantaneous ICA model (2). That is, the sources can be separated or the independent components estimated up to the scaling, permutation, and sign indeterminacies assuming the post-nonlinearities f_i are invertible and under weak conditions on the mixing matrix \mathbf{A} and source distributions. The post-nonlinearity assumption is useful and reasonable in many signal processing applications, because it can be thought of as a model for a nonlinear sensor distortion. In more general situations, it is a restrictive and somewhat arbitrary constraint.

3 Existing methods

Existing nonlinear BSS methods have been reviewed in [16,17], and earlier in Chapter 17 of the book [1]. The paper [16] reviews especially uniqueness results on nonlinear ICA and BSS, separation methods for post-nonlinear mixtures, and our variational Bayesian estimation methods, referring to papers in which the detailed results have been presented. In this section, we will briefly review methods for BSS in nonlinear mixtures with i.i.d. samples.

The nonlinear blind source separation problem has been approached using auto-associative multilayer perceptron (MLP) networks, that is MLP networks trained with input-output pairs (\mathbf{x}, \mathbf{x}) [20]. The number of nodes in a hidden layer is restricted to be smaller than the number of inputs and outputs, thus creating a bottleneck. The extracted nonlinear features can be retrieved from the values of the hidden nodes. With standard back-propagation this approach is very prone to overfitting and local minima, but more advanced learning methods such as flat minimum search can provide a method for nonlinear BSS [21]. The computational complexity from explicitly learning both mixing and demixing models can still be high.

The MISEP method [22] is a generalisation of the infomax method of linear ICA for nonlinear mixtures using an MLP network to model the nonlinearity. The method provides good separation results in several low-dimensional artificial examples as well as in a real nonlinear image mixture problem [23], but it has not been demonstrated for more than 4 mixtures of 4 sources. Nonlinear denoising source separation [24] can also be applied to separate low-dimensional nonlinear mixtures.

Kernel methods have recently become popular in producing nonlinear counterparts for many linear statistical methods based on second-order statistics. Kernel principal component analysis (KPCA) [25] is an extension of linear

principal component analysis (PCA) that operates on a high-dimensional feature space formed by a nonlinear transformation of the data space. The inner product features required for PCA can be evaluated efficiently in data space using the kernel trick. The method is computationally very efficient, but the results are highly dependent on the choice of the kernel.

The Bayesian kernel method of Gaussian Process Latent Variable Model (GPLVM) [26] offers a probabilistic alternative to KPCA. The presentation of the method seems to focus on use in visualisation and it has not been demonstrated for more than 3 dimensional latent spaces.

Post-nonlinear ICA has been approached, for instance, by direct minimisation of mutual information [19] as well as by heuristic methods. These are based on compensating for the invertible post-nonlinearities by first transforming the marginal densities of the data to be as Gaussian as possible and then applying standard ICA [27]. The downside of all of these methods is that they only work if the post-nonlinearities are invertible.

4 Models

In this section, the different models used in this work are presented. First, the multilayer perceptron (MLP) network based NFA model for general nonlinear mixtures from [9] is presented. Then, the hierarchical HNFA model based on the Bayes Blocks framework [33] and the PNFA model for post-nonlinear mixtures are introduced.

4.1 Multilayer Perceptron (MLP) Nonlinearity

Let us assume the observed data \mathbf{X} follows the nonlinear mixing model of Eq. (1). The nonlinearity \mathbf{f} is parametrised with a multilayer perceptron (MLP) network with single hidden layer with H hidden nodes [20]. This allows for writing a generative model for a data vector $\mathbf{x}(t)$ as [9]

$$\mathbf{x}(t) = \mathbf{f}(\mathbf{s}(t), \boldsymbol{\theta}_{\mathbf{f}}) + \mathbf{n}(t) = \mathbf{B}\phi(\mathbf{A}\mathbf{s}(t) + \mathbf{a}) + \mathbf{b} + \mathbf{n}(t), \quad (4)$$

where $\boldsymbol{\theta}_{\mathbf{f}} = (\mathbf{A}, \mathbf{B}, \mathbf{a}, \mathbf{b})$, $\mathbf{A} \in \mathbb{R}^{H \times M}$ and $\mathbf{B} \in \mathbb{R}^{N \times H}$ are the weight matrices and $\mathbf{a} \in \mathbb{R}^H$ and $\mathbf{b} \in \mathbb{R}^N$ are the bias vectors of the first and second layer of the MLP network, respectively. The weight matrices and bias vectors will hence be collectively called the weights of the MLP. The activation function ϕ is the standard hyperbolic tangent. It is applied component-wise to its argument vector. This approach is perfectly general because MLP networks are universal

function approximators that can model any nonlinear mapping to arbitrary accuracy [28].

The noise $\mathbf{n}(t)$ and all the weight matrices and bias vectors of the MLP are *a priori* assumed to be Gaussian and independent of each other. The noise is assumed to have a general diagonal covariance and a hierarchical log-normal variance prior of the form

$$p(\mathbf{n}(t)|\mathbf{v}_n) = N(\mathbf{n}(t); \mathbf{0}, \text{diag}(\exp(2\mathbf{v}_n))) \quad (5)$$

$$p(v_{n_i}|m_{v_n}, v_{v_n}) = N(v_{n_i}; m_{v_n}, \exp(2v_{v_n})), \quad (6)$$

where $\mathbf{v}_n = (v_{n_1}, \dots, v_{n_N})$ and $\text{diag}(\mathbf{v})$ is a diagonal matrix with the elements of vector \mathbf{v} on the main diagonal. The log-normal parametrisation is used because it is easy to define a hierarchical prior. As the variance model is not that important in this application, a conjugate inverse Gamma variance prior could be used here as well. Restriction to an isotropic noise model with the noise covariance of the form $\lambda\mathbf{I}$ is also possible. The noise model and the generative model (4) imply the likelihood

$$p(\mathbf{x}(t)|\boldsymbol{\theta}_f, \mathbf{s}(t), \mathbf{v}_n) = N(\mathbf{x}(t); \mathbf{f}(\mathbf{s}(t), \boldsymbol{\theta}_f), \text{diag}(\exp(2\mathbf{v}_n))). \quad (7)$$

The hierarchical priors of the MLP weights are similar to those of the noise except that the prior of \mathbf{A} is fixed to unit variance to resolve the scaling ambiguity between \mathbf{A} and \mathbf{s} , and the different columns of \mathbf{B} have their own priors:

$$p(A_{ij}) = N(A_{ij}; 0, 1) \quad (8)$$

$$p(B_{ij}|v_{B_j}) = N(B_{ij}; 0, \exp(2v_{B_j})) \quad (9)$$

$$p(a_i|m_a, v_a) = N(a_i; m_a, \exp(2v_a)) \quad (10)$$

$$p(b_i|m_b, v_b) = N(b_i; m_b, \exp(2v_b)) \quad (11)$$

$$p(v_{B_j}|m_{v_B}, v_{v_B}) = N(v_{B_j}; m_{v_B}, \exp(2v_{v_B})). \quad (12)$$

The highest-level hyperparameters m_{v_n} , v_{v_n} , m_a , v_a , m_b , v_b , m_{v_B} and v_{v_B} have vague Gaussian priors $N(0, 100^2)$ ².

To complete the definition of the model, the prior of the sources \mathbf{s} must still be determined. In this case there are several possibilities, leading to models with different benefits and drawbacks. The differences between the alternatives are discussed in more detail in Section 4.4 below.

The simplest source model is the Gaussian model used in nonlinear factor

² The data is usually preprocessed by scaling to approximately unit variance to ensure proper scaling of the weights so that the priors really are vague.

analysis (NFA). This is achieved by

$$p(\mathbf{s}(t)|\mathbf{v}_s) = N(\mathbf{s}(t); \mathbf{0}, \text{diag}(\exp(2\mathbf{v}_s))) \quad (13)$$

$$p(v_{s_i}|m_{v_s}, v_{v_s}) = N(v_{s_i}; m_{v_s}, \exp(2v_{v_s})), \quad (14)$$

where the hyperparameters m_{v_s} and v_{v_s} again have noninformative prior $N(0, 100^2)$. With these definitions, the sets of observations $\mathbf{X} = \{\mathbf{x}(t)|t\}$ and sources $\mathbf{S} = \{\mathbf{s}(t)|t\}$ are defined as usual. The parameter vector $\boldsymbol{\theta}$ contains everything described above, that is $\boldsymbol{\theta} = (\boldsymbol{\theta}_f, \mathbf{v}_n, m_{v_n}, v_{v_n}, (v_{B_j}), m_a, v_a, m_b, v_b, m_{v_B}, v_{v_B}, \mathbf{v}_s, m_{v_s}, v_{v_s})$ including all j in v_{B_j} .

The Gaussian model is computationally simple, but it suffers from the same rotation indeterminacy as the linear factor analysis (FA) model. This can be corrected in the same way linear independent factor analysis [2] extends linear FA: by using a mixture of Gaussians as the source prior. Introducing a new latent variable $M_i(t)$ to denote the active mixture component for source i at sample t leads to a nonlinear independent factor analysis (NIFA) model with a prior of the form

$$p(s_i(t)|M_i(t) = l, m_{sil}, v_{sil}) = N(s_i(t); m_{sil}, \exp(2v_{sil})) \quad (15)$$

$$p(m_{sil}|v_{m_s}) = N(m_{sil}; 0, \exp(2v_{m_s})) \quad (16)$$

$$p(v_{sil}|m_{v_s}, v_{v_s}) = N(v_{sil}; m_{v_s}, \exp(2v_{v_s})). \quad (17)$$

The mixing proportions have a logistic normal prior given by the softmax function

$$p(M_i(t) = l|c_i) = \exp(c_{il}) / \sum_{l'} \exp(c_{il'}) \quad (18)$$

$$p(c_{il}|v_c) = N(c_{il}; 0, \exp(2v_c)). \quad (19)$$

All the highest level parameters v_{m_s} , m_{v_s} , v_{v_s} and v_c again have a noninformative prior $N(0, 100^2)$. The parameters $\boldsymbol{\theta}$ can be defined similarly as above.

In the following, the simpler NFA model is used instead of the more complex NIFA model. More detailed justification for this is presented in Sec. 4.4. Most of the discussion generalises fairly easily to NIFA as well, as is shown in [9,29].

4.2 Hierarchical Nonlinearity

The computational complexity of learning the NFA model is quadratic with respect to the number of the sources. This arises from the need to evaluate the Jacobian matrix of the outputs of the MLP network with respect to the inputs for each sample. The Jacobian is needed to reliably approximate the nonlinearity [15].

One way of defining a more efficiently learnable new model is to introduce additional latent variables to the hidden nodes of the MLP-like network as is done in the hierarchical nonlinear factor analysis (HNFA) model [10]. The HNFA model is defined by the equations

$$\mathbf{s}(t) \sim N(\mathbf{0}, \mathbf{I}) \quad (20)$$

$$\mathbf{h}(t) \sim N(\mathbf{A}\mathbf{s}(t) + \mathbf{a}, \text{diag}(\exp(-\mathbf{v}_h))) \quad (21)$$

$$\mathbf{x}(t) \sim N(\mathbf{C}\mathbf{s}(t) + \mathbf{B}\phi(\mathbf{h}(t)) + \mathbf{b}, \text{diag}(\exp(-\mathbf{v}_n))), \quad (22)$$

where $\mathbf{h}(t)$ are the latent variables modelling the values of the hidden nodes, the vectors \mathbf{v}_h and \mathbf{v}_n parametrise the variances of the noise or innovation of the hidden nodes and the observations, respectively, ϕ is a vector of activation functions and \mathbf{A} , \mathbf{B} , \mathbf{C} , \mathbf{a} , \mathbf{b} are the weights of the mapping. The activation function is chosen to be $\phi(y) = \exp(-y^2)$, for computational simplicity. The structure of the model is illustrated in Fig. 1.

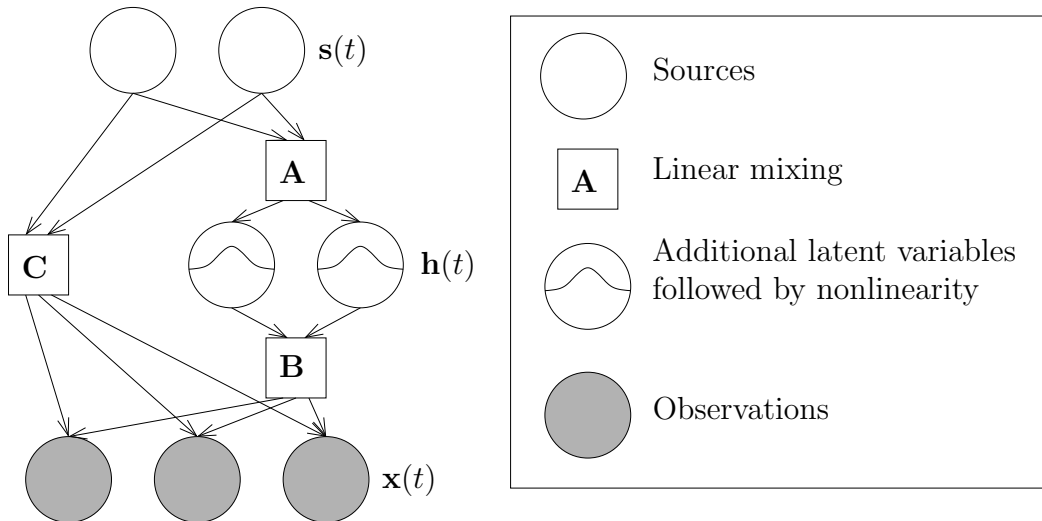


Fig. 1. The HNFA model is illustrated. Square nodes correspond to weight matrices and round nodes to the variables with shaded nodes being observed and unshaded latent.

When assigning hierarchical priors to the parameters of the mapping, we have taken into account the role of different parameters in the model. The weights between different layers, that is \mathbf{A} , \mathbf{B} and \mathbf{C} , can be expected to scale according to the scales of inputs and outputs. Since all $s_i(t)$ are assumed normally distributed, there should be no significant differences in the scales of the weights leaving different nodes but the hidden nodes and observations can

have different scales. This is reflected in the following parametrisation:

$$p(A_{ij}|v_{I,h_i}) = N(A_{ij}; 0, \exp(-v_{I,h_i})) \quad (23)$$

$$p(B_{ij}|v_{I,x_i}, v_{O,h_j}) = N(B_{ij}; 0, \exp(-v_{I,x_i} - v_{O,h_j})) \quad (24)$$

$$p(C_{ij}|v_{I,x_i}) = N(C_{ij}; 0, \exp(-v_{I,x_i})) \quad (25)$$

$$p(a_i) = N(a_i; 0, 1) \quad (26)$$

$$p(b_i) = N(b_i; 0, 1). \quad (27)$$

The parameters v_{O,h_j} reflect the different scales of weights leaving the hidden nodes and v_{I,h_i} and v_{I,x_i} correspond to the different scales of weights entering the hidden nodes and observations, respectively.

Finally, the parameters which model variances on logarithmic scale, that is \mathbf{v}_h , \mathbf{v}_n , v_{I,h_i} , v_{I,x_i} and v_{O,h_j} , each have hyperparameters for mean and log-variance. Their priors are $N(0, \exp(7))$ and $N(0, \exp(4))$, respectively.

MLP networks with sigmoid activation function are known to be universal approximators [28]. Similarly, the set of deterministic HNFA mappings without the linear shortcut

$$f(\mathbf{s}) = \mathbf{B}\phi(\mathbf{A}\mathbf{s}(t) + \mathbf{a}) + \mathbf{b}, \quad (28)$$

where $\phi(y) = \exp(-y^2)$, has the universal approximation property exactly analogous to that of MLP networks [30].

The HNFA model can be implemented efficiently using the building-blocks framework [31–33]. As the additional latent variables $\mathbf{h}(t)$ are independent by definition of the factorial posterior approximation, there is no need to evaluate any nontrivial Jacobians. The disadvantage of HNFA is that the approximation of the posterior density is farther away from the true posterior density. This may occasionally lead to inferior performance [34]. The linear shortcut is included in the model to partially help this, because representing nearly linear mappings would otherwise be significantly more difficult than with the MLP model.

4.3 Post-nonlinear Model

As discussed in Section 2, post-nonlinear (PNL) mixtures are restricted non-linear mixtures of the form

$$x_i(t) = f_i \left(\sum_{j=1}^M A_{ij} s_j(t) \right) + n_i(t), \quad i = 1, \dots, N \quad (29)$$

where the functions $f_i : \mathbb{R} \rightarrow \mathbb{R}$ are called the post-nonlinearities and n_i is additive Gaussian noise. PNL mixtures are a theoretically important special

case of nonlinear ICA, because they can be proven to be separable in the noiseless case if the post-nonlinearities are invertible and the mixing matrix satisfies certain regularity conditions [19].

A variational Bayesian approach to post-nonlinear mixtures was first presented in [11]. The post-nonlinear factor analysis (PNFA) model is based on the generative model

$$x_i(t) = f_i \left(\sum_{j=1}^M A_{ij} s_j(t); \boldsymbol{\theta}_{f_i} \right) + n_i(t), \quad (30)$$

where the post-nonlinearities f_i are modelled with MLP networks with weights $\boldsymbol{\theta}_{f_i}$ as

$$f_i(y, \boldsymbol{\theta}_{f_i}) = \mathbf{D}_i \phi(\mathbf{C}_i y + \mathbf{c}_i) + d_i, \quad (31)$$

where $\boldsymbol{\theta}_{f_i} = (\mathbf{C}_i, \mathbf{c}_i, \mathbf{D}_i, d_i)$ includes the column weight and bias vectors \mathbf{C}_i , \mathbf{c}_i , row weight vector \mathbf{D}_i and bias scalar d_i . This approach also allows for non-invertible post-nonlinearities.

Most of the weights have priors governed by higher level hyperparameters

$$p(A_{ij}) = N(A_{ij}; 0, 1) \quad (32)$$

$$p(C_{ij}|v_{C_i}) = N(C_{ij}; 0, \exp(2v_{C_i})) \quad (33)$$

$$p(c_{ij}|m_{c_i}, v_{c_i}) = N(c_{ij}; m_{c_i}, \exp(2v_{c_i})) \quad (34)$$

$$p(D_{ij}|v_{D_i}) = N(D_{ij}; 0, \exp(2v_{D_i})) \quad (35)$$

$$p(d_i|m_d, v_d) = N(d_i; m_d, \exp(2v_d)), \quad (36)$$

with hierarchical priors for $v_{C_i}, m_{c_i}, v_{c_i}, v_{D_i}$ similar to those used in the NFA model in Eq. (12). The highest level parameters again have vague priors $N(0, 100^2)$.

The source model used in PNFA is Gaussian. It could in principle be rather easily replaced with mixtures-of-Gaussians, but simple post-processing of the extracted sources with linear ICA is often sufficient, as discussed below in Sec. 4.4.

4.4 On different source models

Two different source models were presented above: a Gaussian nonlinear factor analysis (NFA) model in Eq. (13) and a mixture-of-Gaussians based nonlinear *independent* factor analysis (NIFA) model in Eq. (15). From purely theoretical perspective, the NIFA model is preferable as its non-Gaussian model is able to resolve the rotation indeterminacy inherent to the Gaussian model.

In practice, things are not quite that simple. The NIFA model with its fully factorial posterior approximation seems to suffer from similar problems of not being able to extract the true independent sources as linear algorithms based on similar principles [34].

Using the more complicated NIFA model slows down the computation and requires careful initialisation of the mixture components. As this seems to cause more trouble than offer benefits, we have usually chosen not to do it and only use plain NFA. In order to achieve BSS, standard linear ICA can be applied as post-processing to the extracted sources. This approach resembles the one used in [35] for linear separation, where linear factor analysis is used as pre-processing for ICA instead of PCA. Use of factor analysis pre-processing has been shown to lead to good results in case of noisy linear mixtures. The nonlinear case is, however, more difficult as the model may try to nonlinearly transform the sources to be more Gaussian than they should be. Some evidence of this behaviour can be seen in the experiments in Sec. 6.2. Even with a Gaussian source model the sources are, however, not fully Gaussian and using linear ICA to determine the rotation is thus possible.

5 Learning

In this section, methods for learning the models presented in Sec. 4 are presented. The learning method for NFA is again heavily based on earlier work [9] but uses a new method of linearising the MLP [15] that leads to significantly more stable and reliable algorithm.

All the models presented in Sec. 4 are learned using variational Bayesian (VB) learning [36]. Given the data \mathbf{X} and a model with parameters $\boldsymbol{\theta}$ and latent variables or sources \mathbf{S} , the method is based on approximating the posterior distribution $p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{X}, \mathcal{H})$ with a factorial approximation $q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi}) = q(\mathbf{S} | \boldsymbol{\xi}_{\mathbf{S}})q(\boldsymbol{\theta} | \boldsymbol{\xi}_{\boldsymbol{\theta}})$, where $\boldsymbol{\xi}$ are the variational parameters and \mathcal{H} represents the background assumptions of a particular model. The approximation is fitted by minimising the Kullback–Leibler divergence between q and p , or equivalently maximising a lower bound on marginal log-likelihood of the data. By changing the sign, this can be formulated as minimisation of the cost

$$\begin{aligned} \mathcal{C} &= \left\langle \log \frac{q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi})}{p(\mathbf{S}, \boldsymbol{\theta}, \mathbf{X} | \mathcal{H})} \right\rangle \\ &= D_{KL}(q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi}) || p(\mathbf{S}, \boldsymbol{\theta} | \mathbf{X}, \mathcal{H})) - \log p(\mathbf{X} | \mathcal{H}) \geq -\log p(\mathbf{X} | \mathcal{H}), \end{aligned} \quad (37)$$

where $\langle \cdot \rangle$ denotes the expectation over q and $D_{KL}(q || p)$ is the Kullback–Leibler divergence between q and p .

Most applications of variational Bayesian learning [37] involve only models in the conjugate exponential family. The models described in Sec. 4 are clearly not in the exponential family, so the learning methods are somewhat different and require numerical minimisation of the cost or even only its approximation.

In almost all the models, we use the naïve mean field approximation, that is a fully factorial Gaussian distribution

$$q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi}) = q(\mathbf{S} | \boldsymbol{\xi}_{\mathbf{S}}) q(\boldsymbol{\theta} | \boldsymbol{\xi}_{\boldsymbol{\theta}}) = \prod_{i,t} q(s_i(t) | \boldsymbol{\xi}_{s_i(t)}) \prod_j q(\theta_j | \boldsymbol{\xi}_{\theta_j}). \quad (38)$$

The individual factors are parametrised with variational parameters corresponding to the posterior mean and variance of the variable as

$$q(s_i(t) | \boldsymbol{\xi}_{s_i(t)}) = N(s_i(t); \bar{s}_i(t), \tilde{s}_i(t)) \quad (39)$$

$$q(\theta_j | \boldsymbol{\xi}_{\theta_j}) = N(\theta_j; \bar{\theta}_j, \tilde{\theta}_j). \quad (40)$$

For parameters modelling the mean of a Gaussian, the Gaussian distribution is a conjugate prior and the optimal free-form approximation is of this form, assuming the factorisation. For parameters modelling the log-variance of a Gaussian, the Gaussian prior is not conjugate and the posterior approximation is only an approximation.

The definition of the model and the approximating distribution allow evaluating the cost (37) as a function of the variational parameters $\boldsymbol{\xi}$. The cost can be split into two parts as the negative entropy of the approximation and an expected energy

$$\mathcal{C} = \mathcal{C}_q + \mathcal{C}_p = \langle \log q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi}) \rangle + \langle -\log p(\mathbf{X}, \mathbf{S}, \boldsymbol{\theta} | \mathcal{H}) \rangle. \quad (41)$$

Using the factorisation of q , the negative entropy term \mathcal{C}_q splits into

$$\begin{aligned} \mathcal{C}_q &= \langle \log q(\mathbf{S}, \boldsymbol{\theta} | \boldsymbol{\xi}) \rangle = \left\langle \log \prod_{i,t} q(s_i(t) | \boldsymbol{\xi}_{s_i(t)}) \prod_j q(\theta_j | \boldsymbol{\xi}_{\theta_j}) \right\rangle \\ &= \sum_{i,t} \langle \log q(s_i(t) | \boldsymbol{\xi}_{s_i(t)}) \rangle + \sum_j \langle \log q(\theta_j | \boldsymbol{\xi}_{\theta_j}) \rangle. \end{aligned} \quad (42)$$

The remaining terms are negative entropies of univariate Gaussians having values depending only on the variance

$$\langle \log q(\theta_j | \boldsymbol{\xi}_{\theta_j}) \rangle = -\frac{1}{2} - \frac{1}{2} \log(2\pi\tilde{\theta}_j). \quad (43)$$

The expected energy term \mathcal{C}_p is more model specific and is presented separately for the different models.

5.1 MLP nonlinearity

While basic NFA uses a fully factorial posterior approximation, this is not the case for NIFA. Instead, the dependences between $M_i(t)$ and $s_i(t)$ are modelled so that the approximation for them is of the form

$$q(M_i(t), s_i(t)) = q(s_i(t)|M_i(t))q(M_i(t)) \quad (44)$$

which yields a Gaussian mixture approximation for $s_i(t)$

$$q(s_i(t)) = \sum_l q(s_i(t)|M_i(t) = l)q(M_i(t) = l). \quad (45)$$

Otherwise the approximation is similar to NFA.

5.1.1 Evaluating the cost

Evaluating the expected energy term \mathcal{C}_p is slightly more difficult than term \mathcal{C}_q . Using the definition of the model, it can for NFA and PNFA models be factored into terms

$$\begin{aligned} \mathcal{C}_p &= \langle -\log p(\mathbf{X}|\mathbf{S}, \boldsymbol{\theta}, \mathcal{H}) \rangle + \langle -\log p(\mathbf{S}|\boldsymbol{\theta}, \mathcal{H}) \rangle + \langle -\log p(\boldsymbol{\theta}|\mathcal{H}) \rangle \\ &= \sum_{i,t} \langle -\log p(x_i(t)|\mathbf{s}(t), \boldsymbol{\theta}, \mathcal{H}) \rangle + \sum_{i,t} \langle -\log p(s_i(t)|\boldsymbol{\theta}, \mathcal{H}) \rangle + \langle -\log p(\boldsymbol{\theta}|\mathcal{H}) \rangle. \end{aligned} \quad (46)$$

The terms in the second and third summand are expectations of negative logarithm of Gaussian pdf over Gaussian mean and log-variance parameters. These can be evaluated for example parameter $\theta \sim N(m, \exp(2v))$ through integrals of the form [38]

$$\begin{aligned} \langle -\log p(\theta|m, v, \mathcal{H}) \rangle &= \iiint -\log N(\theta; m, \exp(2v)) q(\theta) q(m) q(v) d\theta dm dv \\ &= \frac{1}{2} \log(2\pi) + \bar{v} + \frac{1}{2} [(\bar{\theta} - \bar{m})^2 + \tilde{\theta} + \tilde{m}] \exp(2\bar{v} - 2\bar{v}). \end{aligned} \quad (47)$$

The terms of the first summand of Eq. (46),

$$\begin{aligned} \langle -\log p(x_i(t)|\mathbf{s}(t), \boldsymbol{\theta}, \mathcal{H}) \rangle &= \langle -\log N(x_i(t); f_i(\mathbf{s}(t), \boldsymbol{\theta}_f), \exp(v_{n_i})) \rangle \\ &= \frac{1}{2} \log(2\pi) + \bar{v}_{n_i} + \frac{1}{2} [(x_i(t) - \bar{f}_i(t))^2 + \tilde{f}_i(t)] \exp(2\bar{v}_{n_i} - 2\bar{v}_{n_i}), \end{aligned} \quad (48)$$

are more difficult as they depend on the mean $\bar{f}_i(t)$ and variance $\tilde{f}_i(t)$ of the outputs of the MLP network. Techniques for approximating these are presented in Appendix A.

5.1.2 Update algorithm

The posterior approximations of the parameters of the hierarchical model, that is those of the type m_θ and v_θ , can be updated using a standard VB EM algorithm [38,9]. The update rules for the sources \mathbf{S} as well as the weights $\boldsymbol{\theta}_f$ of the MLP network are more difficult and they are therefore presented in more detail³.

The variational parameters corresponding to the means of \mathbf{S} and $\boldsymbol{\theta}_f$ are updated with a conjugate gradient algorithm. The required derivatives can be computed from the expression of the cost function. In case of inputs and weights of the MLP, this leads to similar computation of partial derivatives as in back-propagation. The parameters corresponding to the variances are updated using a fixed point algorithm that can be derived by differentiating the split cost (41) with respect to $\tilde{\theta}_j$ and using the evaluated result (43). This yields

$$\frac{\partial \mathcal{C}}{\partial \tilde{\theta}_j} = \frac{\partial \mathcal{C}_q}{\partial \tilde{\theta}_j} + \frac{\partial \mathcal{C}_p}{\partial \tilde{\theta}_j} = -\frac{1}{2\tilde{\theta}_j} + \frac{\partial \mathcal{C}_p}{\partial \tilde{\theta}_j}. \quad (49)$$

Setting this to zero leads to a fixed point update rule for the variances of the sources and MLP network weights

$$\tilde{\theta}_j = \frac{1}{2} \left(\frac{\partial \mathcal{C}_p}{\partial \tilde{\theta}_j} \right)^{-1}. \quad (50)$$

Blindly applying this rule may in some cases lead to instability. This can be corrected by some form of dampening, such as halving the step length on logarithmic scale until the update does not increase the value of the cost function. The variance must also not be set to a negative value, even if the derivative is negative.

5.1.3 Initialisation

The MLP network and the gradient-based learning algorithms are notoriously prone to local minima [40,20]. In order to achieve good results, the methods require a reasonable initialisation.

The NFA method is typically initialised by setting the means of the sources to values given by suitable number of principal components of the data. The means of the MLP weights are initialised to random values while all the variances are initialised to small constant values. After this, only the MLP weights are updated during the first 20 iterations of the update algorithm so that the

³ A Matlab toolbox [39] implementing the NFA method is available at <http://www.cis.hut.fi/projects/bayes/software/>.

model can learn a mapping from the PCA sources to the observations. The hyperparameters of the model are only updated after the first 100 iterations.

The PCA initialisation is easy to compute and sufficient for many purposes, but its linearity may sometimes lead to suboptimal results. To resolve this, kernel PCA (KPCA) was used in the initialisation in [14]. The KPCA initialisations are sensitive to the choice of the kernel and its parameters, but with suitable choices they may help NFA yield significantly better results while using less time. The kernel can also be selected with the variational Bayesian criterion by running the learning algorithm for a few iterations with different initialisations and comparing the cost function values.

5.2 Hierarchical nonlinearity

The HNFA model is defined by a graphical model of Gaussian variables combined using addition, multiplication and nonlinearities directly following the variables. Such models can be implemented using the building-blocks framework introduced in [33]. The learning scheme is designed to minimise the cost function

$$\mathcal{C} = \left\langle \log \frac{q(\mathbf{S}, \mathbf{H}, \boldsymbol{\theta} | \boldsymbol{\xi})}{p(\mathbf{S}, \mathbf{H}, \boldsymbol{\theta}, \mathbf{X} | \mathcal{H})} \right\rangle, \quad (51)$$

where $\mathbf{H} = \{\mathbf{h}(t) | t\}$ is the set of all hidden nodes. The cost function is otherwise the same as Eq. (37), except that the hidden nodes \mathbf{H} have been included as unknown variables.

The basic operation during learning is an iteration where all the terms $q(\theta_i | \boldsymbol{\xi}_{\theta_i})$, $q(s_i(t) | \boldsymbol{\xi}_{s_i(t)})$ and $q(h_i(t) | \boldsymbol{\xi}_{h_i(t)})$ of the approximation

$$\begin{aligned} q(\mathbf{S}, \boldsymbol{\theta}, \mathbf{H} | \boldsymbol{\xi}) &= q(\mathbf{S} | \boldsymbol{\xi}_{\mathbf{S}}) q(\boldsymbol{\theta} | \boldsymbol{\xi}_{\boldsymbol{\theta}}) q(\mathbf{H} | \boldsymbol{\xi}_{\mathbf{H}}) \\ &= \prod_{i,t} q(s_i(t) | \boldsymbol{\xi}_{s_i(t)}) \prod_j q(\theta_j | \boldsymbol{\xi}_{\theta_j}) \prod_{i,t} q(h_i(t) | \boldsymbol{\xi}_{h_i(t)}) \end{aligned} \quad (52)$$

are updated one at a time by minimising (51) with respect to the corresponding variational parameters $\boldsymbol{\xi}$. The approximation is again otherwise the same as Eq. (38), except for the inclusion of the hidden nodes \mathbf{H} . Details of the minimisation of the q -terms are presented in [33, Appendices A and C]⁴.

In addition to the basic updates, several other operations are performed in order to help avoid local optima and speed up learning:

- addition of hidden nodes;

⁴ The Bayes Blocks toolbox [41] and the code for the HNFA method are available at <http://www.cis.hut.fi/projects/bayes/software/>.

- addition of weights;
- pruning of weights; and
- line search.

Line search has been explained in [42]. The idea is to monitor the individual updates during one iteration and then perform a line search simultaneously for all terms of the approximation. We applied the line search after every tenth iteration.

The addition and pruning operations aim at optimising the model structure. Incremental addition of hidden nodes to the HNFA model is necessary because the hidden nodes are “expensive” and thus easily pruned out by making their outgoing weights approach to zero. The cost function (37) relates to the marginal likelihood $p(\mathbf{X}|\mathcal{H})$ which can be used to find the most likely model structure.

In general, addition takes place randomly and pruning is based on estimating whether the cost function can be decreased by removing a weight. The motivation for this is that variational Bayesian learning can effectively prune out parts of the model which are not needed. The weights in the matrix \mathbf{B} corresponding to one hidden node can for instance approach zero. The cost function can usually be decreased by removing such weights. If all outgoing weights of a hidden node have been removed, the hidden node becomes useless and can be removed from the model. Variational Bayesian learning cannot, however, actively make room for a part of the model which may be added in the future. It usually takes some time for the rest of the model to accommodate to additions.

5.2.1 Evidence node

In order to avoid local minima in variational Bayesian learning, it is necessary to initialise some variables and keep them fixed for a while until other parts of the model have accommodated appropriately. In the building-blocks framework, we have adopted the virtual-evidence approach [43]. This is implemented by attaching evidence nodes [44], as we call them, to a variable θ_i , whose value we want to set. The node provides a term for the likelihood $p(\text{children}|\theta_i)$. When $q(\theta_i|\boldsymbol{\xi}_{\theta_i})$ is updated, the distribution of θ_i will be close to the value set by evidence node if the likelihood term has a narrow peak but θ_i can accommodate to other parts of the model if the likelihood term is wide. Since all the prior distributions in the present model are Gaussian, it is convenient to use the same functional form for the evidence nodes:

$$-\log p(\text{virtual child}|\theta_i; \alpha_i(\tau_i), \mu_i, \sigma_i^2) = \alpha_i(\tau_i) \frac{(\theta_i - \mu_i)^2}{2\sigma_i^2}, \quad (53)$$

where μ_i and σ_i^2 are the mean and variance induced by the evidence node. The term $\alpha(\tau_i)$ controls the decay of the virtual evidence. It decreases linearly from $\alpha(0) = 1$ to $\alpha(T) = 0$ after which the evidence node is removed. The persistence of the initialisation can thus be controlled by the life-span T (iterations) of the evidence node.

5.2.2 Phases of learning

The model is built in stages. First, only the linear mapping $\mathbf{C}\mathbf{s}(t)$ is used, so that there are no hidden nodes. The sources $\mathbf{s}(t)$ are initialised by principal component analysis (PCA) using evidence nodes with a life-span of 40 iterations and variance $\sigma^2 = 10^{-3}$ in order to estimate a reasonable \mathbf{C} . The linear model is learned for 100 iterations.

After that, 50 randomly initialised hidden nodes are added to the model and estimation of the model structure begins. That is, weights are added and pruned and hidden nodes are added every now and then. Every time new hidden nodes are added, five of them are selected from a pool of 1,000 random candidates. After each addition of hidden nodes, there is a period of 30 iterations during which no pruning is applied. This gives the new hidden nodes enough time to fit themselves into the model.

Hidden nodes are added a limited number of times. After that, learning continues with pruning and random additions of weights. The number of weights to be added decreases with time. Finally, only line searches are applied for the last 1,000 iterations. The total number of iterations in the HNFA simulations is 10,000 unless otherwise stated.

5.2.3 Addition of hidden nodes

The hidden nodes are latent variables which can independently represent some aspects of the observations. Due to our model structure, this usually corresponds to a local minimum of the cost function. It is better that the sources $\mathbf{s}(t)$ represent the data since they can share their information with all hidden nodes. Local minima can be avoided by evidence nodes which keep the variance of the Gaussian noise, innovation, associated to each newly added hidden node $h_i(t)$ through the corresponding parameter v_{h_i} low. Once the sources take the responsibility for the representation, the variances of hidden nodes no longer grow significantly. In our experiments, the life-span of these evidence nodes was 500 iterations and their mean and variance were set to 9 and unity, respectively. Since the parameter controls the logarithm of the inverse variance, value of 9 corresponds to standard deviation of approximately 1/90.

Fig. 2 shows a clear correlation of the quality of the results for the experimental

setting described in Sec. 6.1, and a measure of the amount of independent innovation in the hidden nodes, the latter of which can be influenced by introducing the evidence nodes. More detailed discussion is presented in [45].

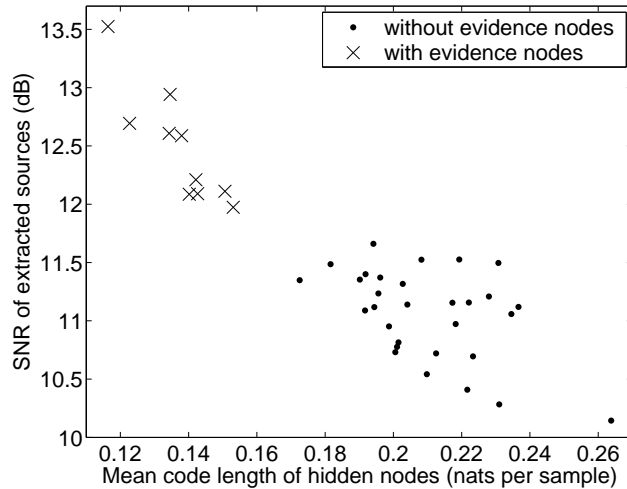


Fig. 2. The attained signal-to-noise ratio of the sources separated using HNFA with and without evidence nodes. (From [45].)

In addition to restricting the innovations, the incoming weights \mathbf{A} of the hidden nodes are initialised to random values by evidence nodes with variance $\sigma^2 = 10^{-2}$ and life time of 40 iterations, when new nodes are added. However, after the first addition, the added hidden nodes are selected from a large pool of candidate initialisations. The hidden nodes which correlate best with the remaining modelling error of the observations are selected. The first hidden nodes are able to model many of the large scale nonlinearities in the data. It is much more difficult to find useful hidden nodes by random initialisations later on since the new nodes tend to be quickly pruned away.

5.3 Post-nonlinear model

The evaluation of the cost function and the learning process of the PNFA model are similar to those of the general NFA. Approximating the nonlinearity is a little easier, as there is only one “source” input with larger posterior variance. The posterior variances of the MLP weights are typically smaller as a single weight affects and thus gains evidence from several observations. This allows for using a hybrid of a Gaussian quadrature with respect to the MLP inputs $y_i(t) = \sum_{j=1}^M A_{ij}s_j$ and a Taylor approximation with respect to the weights θ_{f_i} , as presented in Appendix A. The more general approach used with NFA could of course be used as well.

Using the approximation of the nonlinearity presented in Appendix A, the PNFA cost function can be evaluated using Eqs. (41)–(43) and (46)–(48). The

sources \mathbf{s} and all the weights \mathbf{A} and $\boldsymbol{\theta}_{f_i}$ are updated by conjugate gradient and all the other parameters by VB EM, as described for the NFA model in Sec. 5.1.2.

The initialisation of the PNFA method is similar to NFA: the sources are initialised with PCA. It would be possible to use other simple post-nonlinear methods such as Gaussianisation [27]. After the initialisation, the sources and the hyperparameters are kept fixed for the first 100 and 150 iterations, respectively, while the linear mapping \mathbf{A} and the post-nonlinearities f_i are updated.

5.4 *Missing Values*

The Bayesian framework facilitates principled handling of missing or partially missing elements in the data matrix [44]. Assuming the elements are missing at random, they provide no evidence and can hence be simply ignored in learning. After learning, the missing values can be reconstructed by using the posterior predictive distribution.

Missing values in nonlinear FA have been studied using NFA and HNFA models in [12,13].

6 Experiments

In this section, we present the results of two comparisons of the methods using artificial data sets. The artificial data sets are necessary for comparing different methods as the true sources are not known in most real applications and comparison of different methods is thus practically impossible.

6.1 *General nonlinear mixtures*

The NFA and HNFA methods were tested with an artificial example. The data set consisted of 1,000 samples from nonlinear mixtures of eight sources. Four of the sources were super-Gaussian and the remaining four were sub-Gaussian. The nonlinear mixing was a randomly initialised MLP network with \sinh^{-1} as the activation function for the 30 hidden nodes. The standard deviation of the additive Gaussian noise was one tenth of that of the signal. This corresponds to a signal-to-noise ratio (SNR) of 20 dB. The same data set was also used in [10] and a noiseless version of it in [9].

Linear models (implemented as HNFA models without any hidden nodes) as

well as NFA and HNFA models were tested with varying number of sources. The values of the cost function attained after learning were compared. The best linear model had 14 sources while the best NFA and HNFA models had eight sources. Hence both NFA and HNFA methods are able to infer the correct subspace dimension, while the linear model tries to model nonlinear effects using extra dimensions.

In addition to the proposed variational methods, the experiment was repeated using MAP estimation for the NFA model (MAP-NFA). In this case all the priors and the model structure were fixed to the their generative values.

Unfortunately most other nonlinear BSS methods do not scale to problems of this size. For Kernel PCA [25] the size is not a problem, but out of several Gaussian and polynomial kernels tested, the best results were attained with a linear kernel which is equivalent to linear PCA. Finding a better kernel may certainly be possible, but at least very difficult. We also tried the MISEP method [22], but were unable to attain better results than using linear methods. Again, this may be caused by the authors' lack of experience with MISEP.

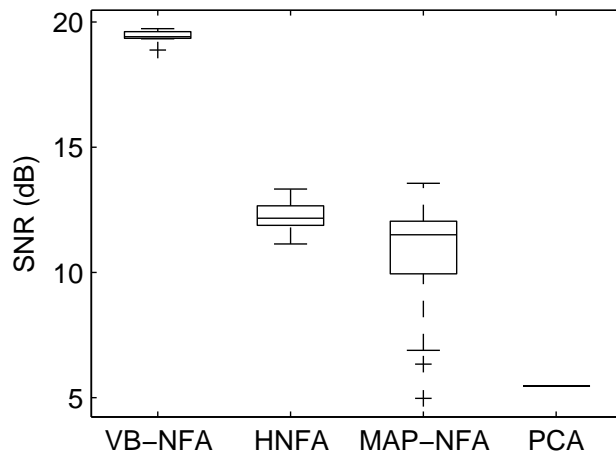


Fig. 3. Boxplot showing variation in the results of different algorithms in the general artificial mixture example.

The experiment was repeated a number of times with different random initialisations. Because of the flexibility of the nonlinear models, each of these ended in a different local optimum. The variation among the results is illustrated in Fig. 3 in terms of signal-to-noise ratio (SNR) with respect to the true sources after rotation of the extracted sources by symmetric FastICA [1]. The results shown are from 10 experiments with VB-NFA, 42 experiments with HNFA and 20 experiments with MAP-NFA. The results show that VB-NFA consistently finds very good solutions while there is lot more variation in the results of HNFA and especially MAP-NFA.

The figure also shows VB-NFA to be clearly superior to MAP-NFA, which

suffers from serious overfitting. Additionally, there is relatively strong correlation between the VB cost function value and SNR, making it easy to find the best solutions among the alternatives. For MAP-NFA, the correlation between posterior density value and SNR is significantly weaker.

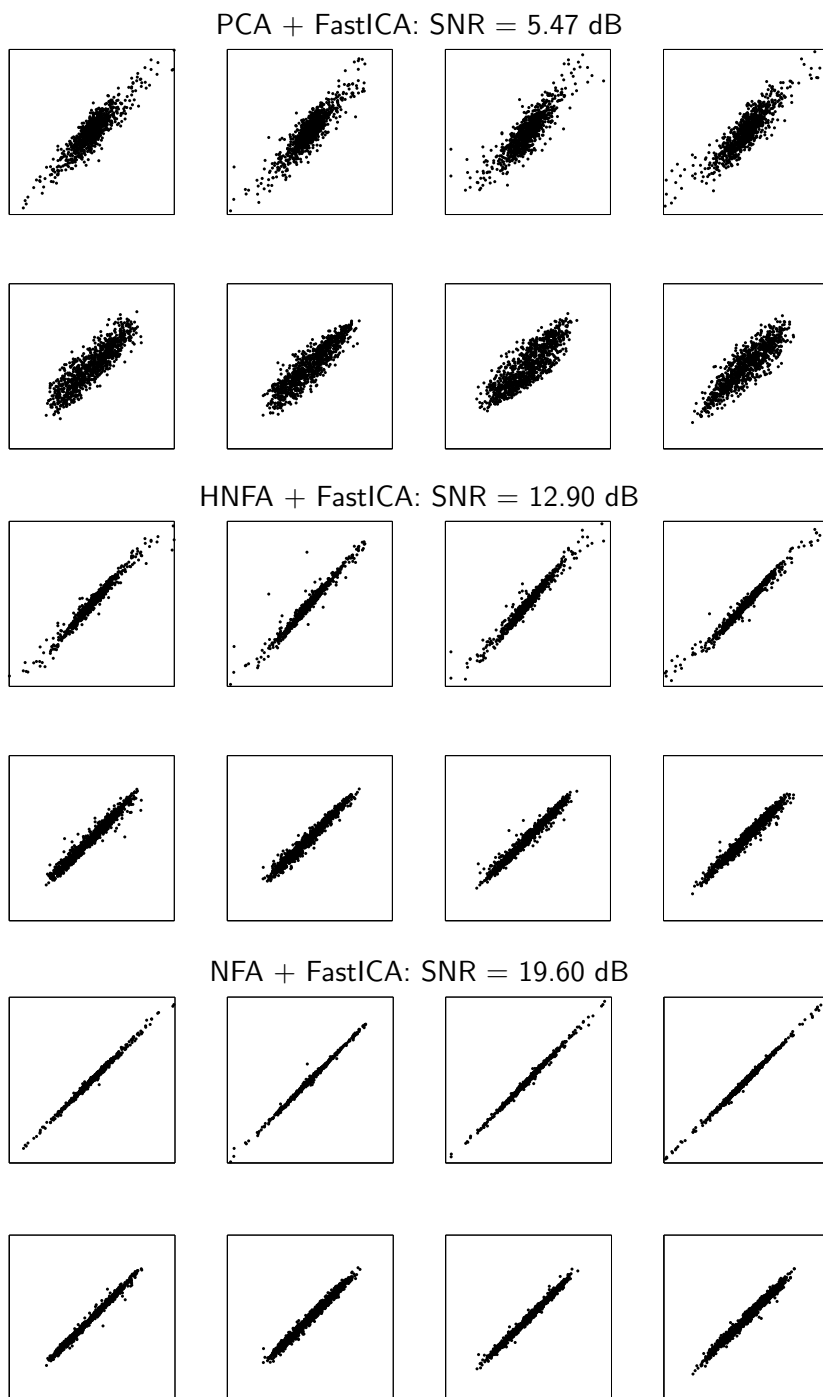


Fig. 4. Each scatter plot shows the values of one original source signal plotted against the best corresponding estimated source signal after a rotation with FastICA. The perfect result would be a straight line.

After a nonlinear subspace has been estimated by NFA or HNFA, standard linear ICA algorithms [1] can be used for rotating the subspace to obtain independent source signals. Fig. 4 shows the scatter plots of the original sources and the sources obtained after a rotation by symmetric FastICA. The plots are from the best results in Fig. 3 according to the VB cost function, hence they are not necessarily the ones with the best overall SNR. On the first two rows, standard linear PCA has been used to estimate the subspace. The middle two rows show the results when the subspace was estimated using the HNFA method and the two bottom rows when the NFA method was used.

The results show that linear methods alone are not sufficient to separate the sources. The results attained by HNFA are reasonably good while those of NFA are already very good. The difference in performance between the methods may be partly due to the fact that the generative nonlinearity more closely resembles that used in NFA, even though it is not the same. The computation times required by NFA and HNFA in this experiment are comparable.

6.2 Post-nonlinear (PNL) mixtures

The proposed NFA and PNFA methods were also tested with a three-dimensional PNL mixture of two independent sources. The sources were a sine wave and uniformly distributed white noise. The PNL transformation used for generating the data contained two non-invertible post-nonlinear distortions:

$$\mathbf{y} = \begin{bmatrix} 1.2 & 0.2 \\ 1 & 0.7 \\ 0.2 & 0.8 \end{bmatrix} \mathbf{s}; \quad \mathbf{x} = \begin{bmatrix} (y_1 - 0.5)^2 \\ (y_2 + 0.4)^2 \\ \tanh(2y_3) \end{bmatrix}. \quad (54)$$

The observations were centred and normalised to unit variance and observation noise with variance 0.01 was added. The number of samples was 400. The data set is illustrated in Fig. 5.

The NFA and PNFA models were trained by trying different model structures, i.e. different numbers of hidden nodes in the MLPs, and several random initialisations of weights of the linear and nonlinear mappings. The PNFA model with the lowest cost had 5 nodes in the hidden layers of all MLPs. The corresponding NFA model with the lowest cost had 10 hidden nodes in the single MLP. The PNFA model was trained for 10000 iterations and NFA for 5000 iterations.

The sources found by NFA and PNFA were further rotated by the FastICA algorithm to obtain independent signals (see Fig. 6). The scatter plots in Fig. 6a show how well the original sources were reconstructed. Each point

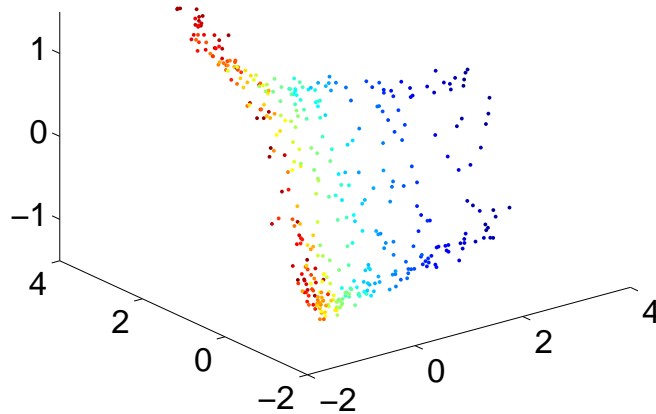


Fig. 5. The data set from the PNL mixture experiment.

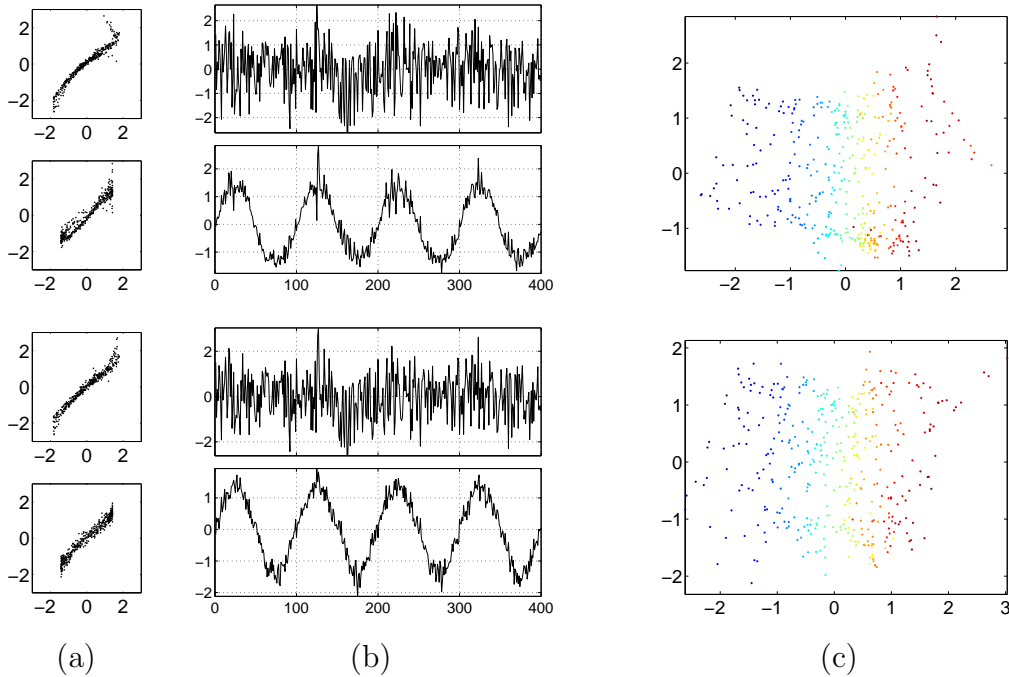


Fig. 6. The sources found by the NFA (top) as well as PNFA (bottom) and further rotated with the FastICA algorithm. (a) The scatter plots; (b) the estimated time series; (c) the distribution of the sources. The signal-to-noise ratio is 10.92 dB for NFA and 12.95 dB for PNFA.

corresponds to one source $s_i(t)$. The abscissa of a point is the original source which was used for generating the data and the ordinate is the estimated source. The optimal result would be a straight line which would mean that the estimated values of the sources coincide with the true values. For PNFA, the sources were estimated quite well except for some points at the edges,

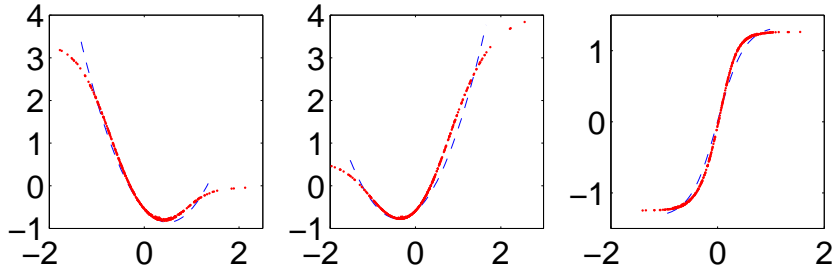


Fig. 7. The estimated post-nonlinear distortions f_i against the functions used for generating the data (the dashed line). Each point in the figure corresponds to a single observation.

while the sources recovered by NFA are not as good.

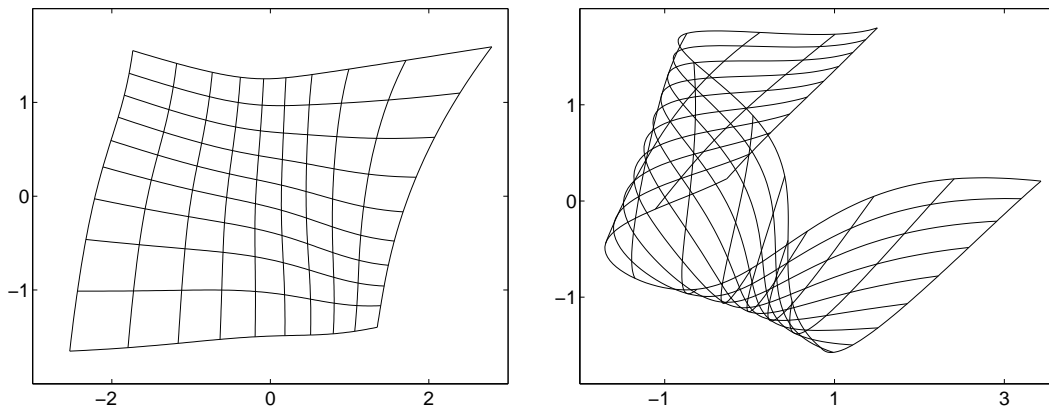


Fig. 8. A grid of original sources \mathbf{s} as mapped through the mixing to \mathbf{x} and demixing to their reconstructions $\hat{\mathbf{s}}$: the perfect solution would be a regular square. Because of the noise, the mappings are only approximate. The results are shown for the PNFA method (left) and Taleb–Jutten PNL ICA algorithm [19] based on non-invertible post-nonlinearities f_i (right).

The PNL distortions learned by the best PNFA model are presented in Fig. 7. The post-nonlinearities f_i are estimated quite well except for some points at the edges. The difficulties mostly affect the two quadratic functions which are difficult to model with such small MLP networks and relatively few observations, especially at those edges.

As expected, regular PNL ICA methods using non-invertible post-nonlinearities cannot recover the true sources at all. This is illustrated in Fig. 8, which shows a relatively clean separation by PNFA compared to severely tangled solution by Taleb–Jutten algorithm [19].

The mediocre overall quality of the results is understandable due to the great difficulty of the test problem: There are only two bounded sub-Gaussian sources in the mixture and their linear combinations are quite far from Gaussianity assumed by NFA and PNFA. Another difficulty is the complex PNL mapping with a small number of observations and several non-invertible post-

nonlinear distortions. Removing any of the observations from the mixture would make the mixing process non-injective and the separation problem unsolvable. The results of PNFA are slightly better than those of regular NFA, which is also natural as the data follows the more restricted model.

7 Discussion

Nonlinear BSS is a very difficult problem, as illustrated by the non-uniqueness results for nonlinear FA and ICA. From a Bayesian perspective, having a multitude of potential solutions that all should somehow be taken into account does not sound so unreasonable. The Bayesian averaging procedure together with reasonable priors on the mapping will help avoid most overly complicated solutions. This was also illustrated experimentally in Sec. 6.1. Using a slightly simpler EM-like approach with point estimates for the parameters of the nonlinear mapping would probably be enough to avoid most overfitting problems. However, as integration over the sources is the most complicated operation, the computational savings would be limited and the ability to infer the structure of the nonlinearity would be lost.

The variational approach used in this work has also other properties that further help regularise the problem. The flexible nonlinear model has many internal symmetries that are reflected in the true posterior. For purposes of the source separation problem, these symmetries are, however, not interesting. The variational approximation breaks the symmetry by ignoring the dependencies between different groups of parameters and the sources. The unimodal posterior approximation will find a broad region of potential solutions, thus returning essentially a single most plausible solution.

The regularisation power of the variational Bayesian approach is nicely demonstrated in [49], where HNFA is shown to prefer a linear model for a magnetoencephalographic (MEG) brain imaging data set where there are theoretical grounds to assume a linear mixing process. The example with approximately 50 sources extracted from 122 channels also demonstrates the scalability of HNFA, using NFA in a problem of this size would be infeasible.

Instead of a nonlinear model with a non-Gaussian prior capable of separating independent sources, simpler Gaussian source prior and post-processing by linear ICA has been used in this work to separate the independent sources. The main reason for using this suboptimal process is practical: even with a non-Gaussian source prior the factorial source posterior approximation would not allow determining the proper rotation of the independent sources [34]. This could be corrected by more advanced approximation techniques [5]. Application of such techniques for nonlinear models is an interesting topic of

future research.

In this work layered mappings have been used to model the nonlinearities. This gives a quite strong prior that allows us to solve the difficult example cases demonstrated in Sec. 6. The broad literature on applications of MLP networks on a wide variety of real world problems provides empirical evidence that such a prior may also be useful in real world applications. Still, having a more easily interpretable prior for the nonlinearities through a Gaussian process [26], for instance, would be interesting.

The presented HNFA and NFA methods have also been compared in reconstructing missing values in speech spectrograms in [33]. According to those experiments, the performance of HNFA was between linear FA and NFA when there were major differences. This seems to indicate that HNFA may emphasise less nonlinear solutions than NFA through stronger regularisation. Thus HNFA is likely better suited for mildly nonlinear problems.

8 Conclusion

We have presented a number of methods for BSS of nonlinear mixtures based on variational Bayesian learning. The basic NFA method uses a MLP network to model the nonlinearity. With a suitable linearisation technique, the variational Bayesian methodology can be easily applied to perform learning and inference on the model much more reliably than before [9]. The computational complexity of the resulting algorithm is quadratic in the number of the sources. To avoid this, an alternative HNFA method based on a hierarchical nonlinearity was proposed. The attained speedup comes at the expense of less accurate results as there are more latent variables in the model and the variational approximate posterior provides a poorer fit to the true posterior. The speeds of NFA and HNFA are comparable in the examples with approximately 10 sources presented in Sec. 6, but better scalability allows applying HNFA to larger problems, where no other alternatives are available. Additionally, a variant of NFA specifically designed for more restricted post-nonlinear mixtures was proposed.

Although HNFA seems less accurate in comparison to NFA in the experiment in Sec. 6.1, the same is not true in general. The comparison of Sec. 6.1 most likely favours NFA somewhat, because the structure of the generative nonlinearity is closer to that used in NFA, even though it is not the same. Other comparisons, such as some of the missing data experiments reported in [33] show HNFA outperform NFA.

Most of the proposed methods employ a Gaussian prior on the sources and are

thus unable to perform source separation by themselves. Still, post-processing the recovered sources by standard linear ICA was enough to recover the underlying independent sources in the presented examples. Using a non-Gaussian source model does not necessarily solve the problem, as ignoring the posterior correlations between the sources in the variational approximation may cause the model to prefer more orthogonal PCA-like solutions. Extending the model by using a Gaussian mixture source model and modelling the posterior correlations of the sources would be an interesting line of future work. The computational complexity of such a model is exponential in the number of the sources, but it would probably help in low dimensional cases.

In order to allow others use the presented methods more easily, free software implementations of the NFA and HNFA methods [39,41] are available on the web at <http://www.cis.hut.fi/projects/bayes/software/>.

Acknowledgements

The authors wish to thank Tapani Raiko for useful comments and discussions and Tomas Östman for his help with the HNFA experiments. This work was supported in part by the IST Programme of the European Community, under the project BLISS, IST-1999-14190, and under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors' views.

A Approximating the MLP nonlinearity

The mean \bar{f}_i and variance \tilde{f}_i of the outputs of the MLP used in Eq. (48) can be evaluated by multidimensional Gaussian integrals

$$\bar{f}_i(t) = \iint f_i(\mathbf{s}(t), \boldsymbol{\theta}_f) q(\mathbf{s}(t), \boldsymbol{\theta}_f | \boldsymbol{\xi}) d\mathbf{s}(t) d\boldsymbol{\theta}_f \quad (\text{A.1})$$

$$\tilde{f}_i(t) = \iint (f_i(\mathbf{s}(t), \boldsymbol{\theta}_f) - \bar{f}_i(t))^2 q(\mathbf{s}(t), \boldsymbol{\theta}_f | \boldsymbol{\xi}) d\mathbf{s}(t) d\boldsymbol{\theta}_f. \quad (\text{A.2})$$

These integrals depend on all the inputs and weights of the MLP network, thus leading to cases of the order of thousands of dimensions.

A.1 Multivariate nonlinearities

In [9], the moments were evaluated by using a first order Taylor approximation of \mathbf{f} for (A.2) and second order approximation for (A.1). This method was

later found to produce unreliable estimates in cases of large source posterior variance, which are common when trying to extract a large number of sources. This caused instability of the algorithm in such cases. This can be corrected by a more accurate approximation based on global linearisation of \mathbf{f} . This can be implemented by replacing the derivatives of the activation function appearing in the Taylor scheme by global linearisations evaluated using Gauss–Hermite quadratures [15].

The Gauss–Hermite quadrature is a general method for approximating integrals of the form

$$I(\phi) = \int_{-\infty}^{\infty} \phi(y) N(y; 0, 1) dy \approx \sum_{i=1}^n w_i \phi(t_i), \quad (\text{A.3})$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar function. For an approximation using n points, the weights w_i and abscissas t_i can be selected so that the result is exact for all polynomials up to order $2n$. We have used a 3-point approximation as it provides a good compromise between accuracy and efficiency. The evaluation points can be easily transformed to handle general mean and variance of the input distribution to get general expectation of $\phi(y)$

$$\bar{\phi}(y)_{\text{GH}} := \sum_{i=1}^n w_i \phi(\bar{y} + t_i \sqrt{\tilde{y}}) \approx \langle \phi(y) \rangle = \int_{-\infty}^{\infty} \phi(y) N(y; \bar{y}, \tilde{y}) dy. \quad (\text{A.4})$$

Variance of $\phi(y)$ can be evaluated through

$$\tilde{\phi}(y)_{\text{GH}} := \sum_{i=1}^n w_i \left[\phi(\bar{y} + t_i \sqrt{\tilde{y}}) - \bar{\phi}(y)_{\text{GH}} \right]^2 \approx \left\langle \left[\phi(y) - \langle \phi(y) \rangle \right]^2 \right\rangle. \quad (\text{A.5})$$

Both the evaluated mean and variance of $\phi(y)$ depend on the mean and variance of y in a nonlinear manner capable of taking into account the specific form of function ϕ .

The quadrature depends on the mean and variance of $\mathbf{y}(t) = \mathbf{A}\mathbf{s}(t) + \mathbf{a}$. These can be evaluated exactly because of the linearity of the mapping as

$$\bar{y}_i(t) = \sum_j \left(\bar{A}_{ij} \bar{s}_j(t) \right) + \bar{a}_i \quad (\text{A.6})$$

$$\tilde{y}_{i,\text{tot}}(t) = \sum_j \left(\tilde{A}_{ij} (\bar{s}_j(t))^2 + \tilde{s}_j(t) \right) + \bar{A}_{ij}^2 \tilde{s}_j(t) + \tilde{a}_i, \quad (\text{A.7})$$

where $\bar{\theta}$ denotes the mean and $\tilde{\theta}$ the variance of θ . Here it is assumed that the posterior approximations $q(\mathbf{S}|\boldsymbol{\xi}_{\mathbf{S}})$ and $q(\boldsymbol{\theta}_{\mathbf{f}}|\boldsymbol{\xi}_{TT})$ have diagonal covariances. Full covariances can be used instead without too much difficulty, if necessary. In order to get more accurate approximation, another quadrature is evaluated

using the variance of $y_i(t)$ originating mainly from $\boldsymbol{\theta}_f$,

$$\tilde{y}_{i,\text{weight}}(t) = \sum_j \tilde{A}_{ij}(\bar{s}_j(t)^2 + \tilde{s}_j(t)) + \tilde{a}_i, \quad (\text{A.8})$$

and using the implied $\tilde{\phi}(y_j(t))_{\text{GH, weight}}$ in the evaluation of the effects of these variances. The total variance (A.7) is still used in evaluation of the means and the evaluation of the effects of the variance of $\mathbf{s}(t)$.

The evaluated mean and variance can be used to define an effective linearisation of the hidden nodes by finding a corresponding linear function that would yield the same mean and variance. This yields the effective slope

$$\langle \phi'(y_i(t)) \rangle := \sqrt{\frac{\tilde{\phi}(y_i(t))_{\text{GH}}}{\tilde{y}_i(t)}}. \quad (\text{A.9})$$

The effective linearisation is able to take into account the variance of the input, thus following the form of the function more globally when the variance is large.

Thus the final approximation for the mean is

$$\bar{f}_i(t) = \sum_j \bar{B}_{ij} \bar{\phi}(y_i(t))_{\text{GH}} + \bar{b}_i \quad (\text{A.10})$$

and for the variance

$$\begin{aligned} \langle (f_i(t) - \bar{f}_i(t))^2 \rangle = \sum_j & \left(\tilde{B}_{ij} \left(\bar{\phi}(y)_{\text{GH}}^2 + \tilde{\phi}(y_j(t))_{\text{GH}} \right) + \bar{B}_{ij}^2 \tilde{\phi}(y_j(t))_{\text{GH, weight}} \right) \\ & + \tilde{b}_i + \tilde{\nabla}_{\mathbf{s}} f_i \text{diag}(\tilde{\mathbf{s}}) \tilde{\nabla}_{\mathbf{s}} f_i^T \end{aligned} \quad (\text{A.11})$$

where the effective gradient is

$$\tilde{\nabla}_{s_j} f_i(t) = \sum_l \bar{B}_{il} \langle \phi'(y_l(t)) \rangle \bar{A}_{lj}. \quad (\text{A.12})$$

A.2 Scalar post-nonlinearities

In case of the MLP networks modelling the scalar post-nonlinearities in PNFA, the approximation is slightly easier as it is possible to directly use the Gauss–Hermite quadrature with respect to the inputs of the MLP. As the variances of the weights are usually small, their effects are represented sufficiently well by using first-order Taylor approximation of the network with respect to them. Thus the mean of the output is approximated as

$$\bar{f}_i(t) = \sum_j w_j f_i(\hat{y}_j(t), \boldsymbol{\theta}_{f_i}) \quad (\text{A.13})$$

where $\mathbf{y} = \mathbf{A}\mathbf{s}$, w_j are the weights and $\hat{y}_j(t) = \bar{y}(t) + t_j\sqrt{\bar{y}(t)}$ are the basis points of the Gauss-Hermite quadrature corresponding to the abscissas t_j , and $\bar{\boldsymbol{\theta}}_{\mathbf{f}_i}$ denotes the mean of the weights $\boldsymbol{\theta}_{\mathbf{f}_i}$.

Correspondingly, the variance is approximated by a combined Gauss-Hermite and Taylor approximation

$$\begin{aligned} \langle (f_i(t) - \bar{f}_i(t))^2 \rangle = \sum_j w_j \left[\left(f_i(\hat{y}_j(t), \bar{\boldsymbol{\theta}}_{\mathbf{f}_i}) - \bar{f}_i(t) \right)^2 \right. \\ \left. + \nabla_{\boldsymbol{\theta}_{\mathbf{f}_i}} f_i(\hat{y}_j(t), \bar{\boldsymbol{\theta}}_{\mathbf{f}_i}) \text{diag}(\tilde{\boldsymbol{\theta}}_{\mathbf{f}_i}) \nabla_{\boldsymbol{\theta}_{\mathbf{f}_i}} f_i(\hat{y}_j(t), \bar{\boldsymbol{\theta}}_{\mathbf{f}_i})^T \right]. \quad (\text{A.14}) \end{aligned}$$

References

- [1] A. Hyvärinen, J. Karhunen, E. Oja, Independent Component Analysis, J. Wiley, 2001.
- [2] H. Attias, Independent factor analysis, *Neural Computation* 11 (4) (1999) 803–851.
- [3] H. Lappalainen, Ensemble learning for independent component analysis, in: Proc. Int. Workshop on Independent Component Analysis and Signal Separation (ICA'99), Aussois, France, 1999, pp. 7–12.
- [4] J. Miskin, Ensemble learning for independent component analysis, Ph.D. thesis, University of Cambridge, UK (2000).
- [5] P. Højen-Sørensen, O. Winther, L. K. Hansen, Mean-field approaches to independent component analysis, *Neural Computation* 14 (4) (2002) 889–918.
- [6] R. A. Choudrey, S. J. Roberts, Variational mixture of Bayesian independent component analyzers, *Neural Computation* 15 (1) (2003) 213–252.
- [7] K. Chan, T.-W. Lee, T. J. Sejnowski, Variational Bayesian learning of ICA with missing data, *Neural Computation* 15 (8) (2003) 1991–2011.
- [8] H. Attias, Independent factor analysis with temporally structured sources, in: S. Solla, T. Leen, K.-R. Müller (Eds.), *Advances in Neural Information Processing Systems 12*, MIT Press, Cambridge, MA, USA, 2000, pp. 386–392.
- [9] H. Lappalainen, A. Honkela, Bayesian nonlinear independent component analysis by multi-layer perceptrons, in: M. Girolami (Ed.), *Advances in Independent Component Analysis*, Springer-Verlag, Berlin, 2000, pp. 93–121.
- [10] H. Valpola, T. Östman, J. Karhunen, Nonlinear independent factor analysis by hierarchical models, in: Proc. 4th Int. Symp. on Independent Component Analysis and Blind Signal Separation (ICA2003), Nara, Japan, 2003, pp. 257–262.

- [11] A. Ilin, A. Honkela, Postnonlinear independent component analysis by variational Bayesian learning, in: C. G. Puntonet, A. Prieto (Eds.), Proc. of the Fifth Int. Conf. on Independent Component Analysis and Blind Signal Separation (ICA 2004), Vol. 3195 of Lecture Notes in Computer Science, Springer-Verlag, Berlin, Granada, Spain, 2004, pp. 766–773.
- [12] T. Raiko, H. Valpola, Missing values in nonlinear factor analysis, in: Proc. of the 8th Int. Conf. on Neural Information Processing (ICONIP'01), Shanghai, 2001, pp. 822–827.
- [13] T. Raiko, H. Valpola, T. Östman, J. Karhunen, Missing values in hierarchical nonlinear factor analysis, in: Proc. of the Int. Conf. on Artificial Neural Networks and Neural Information Processing (ICANN/ICONIP 2003), Istanbul, Turkey, 2003, pp. 185–189.
- [14] A. Honkela, S. Harmeling, L. Lundqvist, H. Valpola, Using kernel PCA for initialisation of variational Bayesian nonlinear blind source separation method, in: C. G. Puntonet, A. Prieto (Eds.), Proc. of the 5th Int. Conf. on Independent Component Analysis and Blind Signal Separation (ICA 2004), Vol. 3195 of Lecture Notes in Computer Science, Springer-Verlag, Berlin, Granada, Spain, 2004, pp. 790–797.
- [15] A. Honkela, H. Valpola, Unsupervised variational Bayesian learning of nonlinear models, in: L. Saul, Y. Weiss, L. Bottou (Eds.), Advances in Neural Information Processing Systems 17, MIT Press, Cambridge, MA, USA, 2005, pp. 593–600.
- [16] C. Jutten, J. Karhunen, Advances in blind source separation (BSS) and independent component analysis (ICA) for nonlinear mixtures, International Journal of Neural Systems 14 (5) (2004) 267–292.
- [17] L. B. Almeida, Nonlinear Source Separation, Synthesis Lectures on Signal Processing, Morgan&Claypool Publishers, 2006.
- [18] A. Hyvärinen, P. Pajunen, Nonlinear independent component analysis: Existence and uniqueness results, Neural Networks 12 (3) (1999) 429–439.
- [19] A. Taleb, C. Jutten, Source separation in post-nonlinear mixtures, IEEE Trans. on Signal Processing 47 (10) (1999) 2807–2820.
- [20] S. Haykin, Neural Networks – A Comprehensive Foundation, 2nd ed., Prentice-Hall, 1999.
- [21] S. Hochreiter, J. Schmidhuber, LOCOCODE performs nonlinear ICA without knowing the number of sources, in: Proc. Int. Workshop on Independent Component Analysis and Blind Signal Separation (ICA'99), Aussois, France, 1999, pp. 149–154.
- [22] L. B. Almeida, MISEP – linear and nonlinear ICA based on mutual information, Journal of Machine Learning Research 4 (2003) 1297–1318.
- [23] L. B. Almeida, Separating a real-life nonlinear image mixture, Journal of Machine Learning Research 6 (2005) 1199–1229.

- [24] M. S. C. Almeida, H. Valpola, J. Särelä, Separation of nonlinear image mixtures by denoising source separation, in: Proc. 6th Int. Conf. on Independent Component Analysis and Blind Source Separation (ICA 2006), Charleston, South Carolina, USA, 2006, pp. 8–15.
- [25] B. Schölkopf, A. Smola, K.-R. Müller, Nonlinear component analysis as a kernel eigenvalue problem, *Neural Computation* 10 (5) (1998) 1299–1319.
- [26] N. Lawrence, Probabilistic non-linear principal component analysis with Gaussian process latent variable models, *Journal of Machine Learning Research* 6 (2005) 1783–1816.
- [27] A. Ziehe, M. Kawanabe, S. Harmeling, K.-R. Müller, Blind separation of post-nonlinear mixtures using linearizing transformations and temporal decorrelation, *Journal of Machine Learning Research* 4 (2003) 1319–1338.
- [28] K. Hornik, M. Stinchcombe, H. White, Multilayer feedforward networks are universal approximators, *Neural Networks* 2 (5) (1989) 359–366.
- [29] H. Valpola, Nonlinear independent component analysis using ensemble learning: theory, in: Proc. Int. Workshop on Independent Component Analysis and Blind Signal Separation (ICA2000), Helsinki, Finland, 2000, pp. 251–256.
- [30] M. Stinchcombe, H. White, Universal approximation using feedforward networks with non-sigmoid hidden layer activation functions, in: Proc. International Joint Conference on Neural Networks (IJCNN '89), 1989, pp. I–613–617.
- [31] H. Valpola, T. Raiko, J. Karhunen, Building blocks for hierarchical latent variable models, in: Proc. 3rd Int. Conf. on Independent Component Analysis and Signal Separation (ICA2001), San Diego, USA, 2001, pp. 710–715.
- [32] M. Harva, T. Raiko, A. Honkela, H. Valpola, J. Karhunen, Bayes Blocks: An implementation of the variational Bayesian building blocks framework, in: Proc. of the 21st Conf. on Uncertainty in Artificial Intelligence (UAI 2005), Edinburgh, Scotland, 2005, pp. 259–266.
- [33] T. Raiko, H. Valpola, M. Harva, J. Karhunen, Building blocks for variational Bayesian learning of latent variable models, Tech. Rep. E4, Helsinki University of Technology, Publications in Computer and Information Science, Espoo, Finland, available at <http://www.cis.hut.fi/Publications/>. Accepted for publication in *Journal of Machine Learning Research* conditioned on minor revisions. (2006).
- [34] A. Ilin, H. Valpola, On the effect of the form of the posterior approximation in variational learning of ICA models, *Neural Processing Letters* 22 (2) (2005) 183–204.
- [35] S. Ikeda, ICA on noisy data: A factor analysis approach, in: M. Girolami (Ed.), *Advances in Independent Component Analysis*, Springer-Verlag, Berlin, 2000, pp. 201–215.

- [36] T. S. Jaakkola, Tutorial on variational approximation methods, in: M. Opper, D. Saad (Eds.), *Advanced Mean Field Methods: Theory and Practice*, The MIT Press, Cambridge, MA, USA, 2001, pp. 129–159.
- [37] J. Winn, C. M. Bishop, Variational message passing, *Journal of Machine Learning Research* 6 (2005) 661–694.
- [38] H. Lappalainen, J. Miskin, Ensemble learning, in: M. Girolami (Ed.), *Advances in Independent Component Analysis*, Springer-Verlag, Berlin, 2000, pp. 75–92.
- [39] A. Honkela, H. Valpola, X. Giannakopoulos, Nonlinear factor analysis Matlab package, <http://www.cis.hut.fi/projects/bayes/software/> (2004).
- [40] K. Fukumizu, S. Amari, Local minima and plateaus in hierarchical structures of multilayer perceptrons, *Neural Networks* 13 (3) (2000) 317–327.
- [41] H. Valpola, A. Honkela, M. Harva, A. Ilin, T. Raiko, T. Östman, Bayes Blocks software library, <http://www.cis.hut.fi/projects/bayes/software/> (2003).
- [42] A. Honkela, H. Valpola, J. Karhunen, Accelerating cyclic update algorithms for parameter estimation by pattern searches, *Neural Processing Letters* 17 (2) (2003) 191–203.
- [43] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann, 1988.
- [44] T. Raiko, Partially observed values, in: *Proc. Int. Joint Conf. on Neural Networks (IJCNN'04)*, Budapest, Hungary, 2004, pp. 2825–2830.
- [45] A. Honkela, H. Valpola, Variational learning and bits-back coding: an information-theoretic view to Bayesian learning, *IEEE Transactions on Neural Networks* 15 (4) (2004) 800–810.
- [46] H. Valpola, J. Karhunen, An unsupervised ensemble learning method for nonlinear dynamic state-space models, *Neural Computation* 14 (11) (2002) 2647–2692.
- [47] T. Kohonen, *Self-Organizing Maps*, 3rd Edition, Springer, 2001.
- [48] J. Vesanto, J. Himberg, E. Alhoniemi, J. Parhankangas, SOM toolbox for Matlab 5, Tech. Rep. A57, Helsinki University of Technology, Publications in Computer and Information Science, Espoo, Finland, URL: <http://www.cis.hut.fi/projects/somtoolbox/> (2000).
- [49] A. Honkela, T. Östman, R. Vigário, Empirical evidence of the linear nature of magnetoencephalograms, in: *Proc. 13th European Symposium on Artificial Neural Networks (ESANN 2005)*, Bruges, Belgium, 2005, pp. 285–290.