

Elsevier Editorial System(tm) for Neurocomputing

Manuscript Draft

Manuscript Number:

Title: Flexible and Efficient Implementations of Bayesian Independent Component Analysis

Article Type: Full Length Article (FLA)

Section/Category:

Keywords: Independent Component Analysis, Empirical Bayes, Mean Field Methods, Variational methods

Corresponding Author: Mr Kaare Brandt Petersen, PhD

Corresponding Author's Institution: Technical University of Denmark

First Author: Ole Winther, PhD

Order of Authors: Ole Winther, PhD; Kaare Brandt Petersen, PhD

Manuscript Region of Origin:

Abstract: In this paper we present an empirical Bayes method for flexible and efficient Independent Component Analysis (ICA). The method is flexible with respect to choice of source prior, dimensionality and positivity of the mixing matrix, and structure of the noise covariance matrix. The efficiency is ensured using parameter optimizers which are more advanced than the expectation maximization (EM) algorithm, but still easy to implement. These optimizers are the overrelaxed adaptive EM algorithm and the easy gradient recipe. The required expectations over the source posterior are estimated with accurate mean field

methods: variational and the expectation consistent framework.

We demonstrate the usefulness of the approach with the publicly available Matlab toolbox `icaMF`.

Flexible and Efficient Implementations of Bayesian Independent Component Analysis

Ole Winther, Kaare Brandt Petersen

December 21, 2005

Abstract

In this paper we present an empirical Bayes method for flexible and efficient Independent Component Analysis (ICA). The method is flexible with respect to choice of source prior, dimensionality and positivity of the mixing matrix, and structure of the noise covariance matrix. The efficiency is ensured using parameter optimizers which are more advanced than the expectation maximization (EM) algorithm, but still easy to implement. These optimizers are the overrelaxed adaptive EM algorithm and the easy gradient recipe. The required expectations over the source posterior are estimated with accurate mean field methods: variational and the expectation consistent framework. We demonstrate the usefulness of the approach with the publicly available Matlab toolbox `icamf`.

Contents

1	Introduction	3
2	Instantaneous ICA	4
3	Optimization of Parameters	5
3.1	The EM Algorithm	5
3.2	Easy Gradient Recipe	6
3.3	Overrelaxed Adaptive EM	7
3.4	Dealing with Constrained Parameters	7
4	Estimating Source Statistics	8
4.1	Expectation Consistent	8
4.2	Variational	12
4.3	EC and Variational Comparison	13
5	Software	13
6	Testing the Efficiency of the Framework	15

7 Conclusion	17
A Solving the EC Equations	18

1 Introduction

Since Independent Component Analysis (ICA) in the early nineties caught the attention of the machine learning community, the interest and activities within this area have all but exploded. Although initially regarded as an example of a blind source separation problem for independent data, the focus has in recent years gradually shifted from different aspects of this instantaneous problem to the challenge of the convolutive case (mixing over time). A process fuelled by algorithms such as InfoMax [2] and FastICA [7] which, although not very flexible, are robust and fast.

But the completely general instantaneous case is far from solved: So far there exists no algorithm which can do noisy, non-square mixing with an arbitrary non-gaussian prior with the same robustness and speed as the above mentioned algorithms. Variational (so-called mean field and also known as ensemble learning) methods [8, 1, 23, 11, 4, 6] are attractive because they are very flexible general modelling tools. The mean field ICA method as described in [6], however, had two major difficulties: First, the flexibility with respect to the prior makes the inside of the black-box rather complicated and unattractive to wide range of the application-driven part of the research community. Second, it was slow to converge and no universal stopping criteria could be given. These two difficulties, however, can now be handled, as this paper demonstrates: The complexity by the availability of a Matlab toolbox with plug-and-play demos and examples and the convergence by efficient optimization schemes beyond the traditional expectation maximization (EM) algorithm .

In a broader context reaching well beyond ICA, the Mean field methods such as variational Bayes, loopy belief propagation, expectation propagation (EP) and expectation consistent (EC) have recently gathered much interest, see e.g. [8, 1, 15, 25, 16, 10, 17] because of their potential as approximate Bayesian inference engines. An undesirable property of the mean field methods in this context is that the approximation error is unattainable. One cannot in any quantitative manner say much about the deviation of marginal moments or likelihoods from their true values. But in many tests, however, the accuracy and the polynomial complexity of the mean field methods to intractable sums and integrals, is positioning the mean field methods as a high-end approximation.

For the last few years, it has been observed that an important cause for the non-robustness of the mean field methods rests not only upon the approximation error but rather on the slow convergence of the expectation maximization (EM) style algorithms, typically used for the joint parameter-latent variable inference problem [20]. Several alternatives to EM learning, also applicable to non-mean field based inference, have been proposed recently [22, 14] and analysis have been given to explain convergence failure both in general and specific settings [21]. Furthermore, these methods in some cases make the difference between convergence in finite time or not [18], but in most cases they at least give a huge speed-up. This new insight has thus taken the mean field methods one step closer to realizing their full potential for Bayesian inference.

In this paper we revisit mean field ICA to demonstrate that the newly found

insights can give us a much more efficient system while still retaining the flexibility of the Bayesian ICA approach. Compared to other methods, the flexibility is rather extensive and enables the user to handle both over-/underdetermined and square mixing, positive constraints on the mixing matrix, noise estimation and general source priors, e.g. positive or discrete. The paper is organized as follows: Section 2 formulates the instantaneous noisy ICA problem which is the basic model of interest. Section 3 explains how the given log likelihood (bound) is optimized giving three different approaches. Section 4 deals with estimation of the source statistics. Two mean field methods – expectation consistent and variational – are reviewed and in Section 5 we present the Matlab toolbox. Finally in Section 6 and 7 we wrap up with demonstrations of the developed methods and conclusions.

2 Instantaneous ICA

In this section, we give a quick recap of the empirical Bayes approach to instantaneous ICA with additive Gaussian noise – for a more detailed account the reader is referred to e.g. [6]. The observation model is given by

$$\mathbf{x}_t = \mathbf{A}\mathbf{s}_t + \mathbf{n}_t, \quad t = 1, \dots, N$$

with N being the number of samples. The noise is assumed zero-mean gaussian with covariance Σ , i.e. the Likelihood is $p(\mathbf{x}_t|\mathbf{A}, \mathbf{s}_t, \Sigma) = \mathcal{N}(\mathbf{x}_t; \mathbf{A}\mathbf{s}_t, \Sigma)$. The source prior factorizes in both sources and time steps. Denoting the stacked sources by the matrix \mathbf{S} , we can write the prior as $p(\mathbf{S}|\boldsymbol{\nu}) = \prod_{it} p_i(S_{it}|\nu_i)$, where $\boldsymbol{\nu}$ is shorthand for the parameters of the prior. The observation vectors \mathbf{x}_t are stacked as columns into one matrix \mathbf{X} : $p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \Sigma) = \prod_t p(\mathbf{x}_t|\mathbf{A}, \mathbf{s}_t, \Sigma)$ and the posterior is given by $p(\mathbf{S}|\mathbf{X}, \boldsymbol{\theta}) = \frac{p(\mathbf{X}|\mathbf{A}, \mathbf{S}, \Sigma)p(\mathbf{S}|\boldsymbol{\nu})}{p(\mathbf{X}|\boldsymbol{\theta})}$, where we have used the shorthand $\boldsymbol{\theta} = \{\mathbf{A}, \Sigma, \boldsymbol{\nu}\}$ for the parameters. In the empirical Bayes (or Maximum Likelihood II) approach applied to ICA, the noise realization and the unobserved sources are integrated out, leaving the parameters $\boldsymbol{\theta}$ to be determined by maximizing the (marginal) Likelihood: $p(\mathbf{X}|\boldsymbol{\theta})$. Alternatively, one may use a hierarchical Bayesian approach [1] marginalizing also over $\boldsymbol{\theta}$, see [11] for an application to ICA.

The log likelihood $\mathcal{L}(\boldsymbol{\theta}) = \ln p(\mathbf{X}|\boldsymbol{\theta})$ is for most priors too complicated for practical approaches and instead a lower bound is used as objective function. The lower bound \mathcal{B} is defined by

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &\equiv \ln p(\mathbf{X}|\boldsymbol{\theta}) = \ln \int q(\mathbf{S}|\boldsymbol{\phi}) \frac{p(\mathbf{X}, \mathbf{S}|\boldsymbol{\theta})}{q(\mathbf{S}|\boldsymbol{\phi})} d\mathbf{S} \\ &\geq \int q(\mathbf{S}|\boldsymbol{\phi}) \ln \frac{p(\mathbf{X}, \mathbf{S}|\boldsymbol{\theta})}{q(\mathbf{S}|\boldsymbol{\phi})} d\mathbf{S} \equiv \mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi}) . \end{aligned} \quad (1)$$

The bounding property is a simple consequence of Jensen's inequality and holds for *any* choice of variational distribution $q(\mathbf{S}|\boldsymbol{\phi})$. In fact it is easy to show that $\mathcal{L}(\boldsymbol{\theta}) = \mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi}) - KL(q, p)$, where $KL(q, p) \geq 0$ denotes the Kullback-Leibler

divergence between the variational distribution and the source posterior. Thus, if the variational distribution becomes equal to the source posterior, $KL(p, p) = 0$ and the bound is equal to the log likelihood.

The derivatives of the bound are easily derived for the ICA model

$$\frac{\partial \mathcal{B}(\boldsymbol{\theta}, \phi)}{\partial \mathbf{A}} = \boldsymbol{\Sigma}^{-1} (\mathbf{X} \langle \mathbf{S} \rangle_q^T - \mathbf{A} \langle \mathbf{S} \mathbf{S}^T \rangle_q) \quad (2)$$

$$\frac{\partial \mathcal{B}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\Sigma}} = \frac{N}{2} \boldsymbol{\Sigma}^{-1} - \frac{1}{2} \boldsymbol{\Sigma}^{-1} \langle (\mathbf{X} - \mathbf{A} \mathbf{S})(\mathbf{X} - \mathbf{A} \mathbf{S})^T \rangle_q \boldsymbol{\Sigma}^{-1} \quad (3)$$

$$\frac{\partial \mathcal{B}(\boldsymbol{\theta}, \phi)}{\partial \boldsymbol{\nu}} = \left\langle \frac{\partial \ln p(\mathbf{s} | \boldsymbol{\mu})}{\partial \boldsymbol{\nu}} \right\rangle_q. \quad (4)$$

Constrained variables are handled by reparametrization and considered in detail in Section 3.4. Since many different source priors are relevant depending on the application, the derivative involving the prior parameter is not specified in details but left open for easier modular fit to various problem-specific priors.

3 Optimization of Parameters

In this section we discuss a number of approaches for optimization of the lower bound of the log likelihood. The starting point is the EM algorithm for its simplicity but also variants thereof are of interest in order to speed up convergence.

3.1 The EM Algorithm

The traditional optimization applied in [6] is the Expectation-Maximization (EM) algorithm as presented in [12]. In their formulation, the EM algorithm is a coordinate-wise descend in the so-called free energy, which in this context is minus the lower bound function. In short, the EM algorithm for $\mathcal{B}(\boldsymbol{\theta}, \phi)$ is

E: Maximize $\mathcal{B}(\boldsymbol{\theta}, \phi)$ with respect to ϕ keeping $\boldsymbol{\theta}$ fixed.

M: Maximize $\mathcal{B}(\boldsymbol{\theta}, \phi)$ with respect to $\boldsymbol{\theta}$ keeping ϕ fixed.

This results holds for any variational distribution. However, if the choice of q is constrained to a family of distributions such that the E-step does *not* give $q(\mathbf{S} | \phi) = p(\mathbf{S} | \mathbf{X}, \boldsymbol{\theta})$ then we are only optimizing the bound and not the likelihood itself. This variational approximation works well in many cases, see Section 6.

M-step: In the M-step the lower bound $\mathcal{B}(\boldsymbol{\theta}, \phi)$ is maximized with respect to the model (hyper) parameters $\boldsymbol{\theta} = \{\mathbf{A}, \boldsymbol{\Sigma}, \boldsymbol{\nu}\}$. Setting the derivatives in eqs. (2) and (3) equal to zero, one obtains the following EM updates for the mixing matrix and the noise covariance

$$\mathbf{A} = \mathbf{X} \langle \mathbf{S} \rangle_q^T \langle \mathbf{S} \mathbf{S}^T \rangle_q^{-1} \quad (5)$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \langle (\mathbf{X} - \mathbf{A} \mathbf{S})(\mathbf{X} - \mathbf{A} \mathbf{S})^T \rangle_q. \quad (6)$$

The corresponding result for the prior parameters cannot be expressed explicitly without choosing what prior to use, and is therefore left for the user of whatever

specific prior. With these estimates we are ready to make another E-step with the new values for the parameters, then another M-step, and so on.

The EM algorithm is simple and have important theoretical convergence properties, but also sometimes unreasonably slow – a postulate reported and documented by a number of papers regarding the use of EM algorithm for ICA. The results of [20] is based on a small scale statistical investigation, and [3] and [19] provide analytically insight to the experience of slow convergence in the low noise limit.

The analysis in the two latter papers is based on a Taylor expansion of the true source posterior moments in the noise variance. That is, when for simplicity the noise is assumed isotropic with variance σ^2 , the source posterior moments $\langle \mathbf{s}_t \rangle$ and $\langle \mathbf{s}_t \mathbf{s}_t^T \rangle$ are expanded in σ^2 and inserted to give the update of the the mixing matrix. The result is

$$\mathbf{A}_{n+1} = \mathbf{A}_n + \mathcal{O}(\sigma^2) ,$$

where \mathbf{A}_n denotes the nth estimate of the mixing matrix. This shows that if the noise variance is very small, then so is the change in the mixing matrix and in the limit of zero noise, the EM algorithm becomes infinitely slow. Further analysis demonstrates that the first order correction term for the mixing matrix is proportional to the noiseless InfoMax update [2]. This renders the use of EM for the noise model even less attractive, since only in cases where the noise is large enough to make the $\mathcal{O}(\sigma^4)$ contribution significant, is the solution of the noise model different from the noiseless InfoMax model. A recent result in [21] shows that going to in to hierarchical variational Bayesian framework [1] does not solve the problem. The resulting Variational Bayes EM algorithm, is suffering from the exact same defect.

With this defect of the EM algorithm in mind, we regard two optimization alternatives which are both closely related to the original EM algorithm.

3.2 Easy Gradient Recipe

The very appealing property of the EM algorithm is the combination of the easily implementable scheme and a guaranteed increase of the likelihood. Often, more advanced optimization methods is more demanding either analytically, computationally or both, and thus not appealing for large class of complicated problems. But using the Easy Gradient Recipe [14], one can obtain the efficiency of state-of-the-art gradient based optimization methods for the cost of the EM algorithm. This is done by recycling the E and M-steps: Consider in pseudo-code some function which is given the model parameters $\boldsymbol{\theta}$ and returns the log likelihood bound and gradient of the log likelihood bound computed in three steps

- $$[\mathcal{B}, \frac{d\mathcal{B}}{d\boldsymbol{\theta}}] = \text{fct}(\boldsymbol{\theta})$$
- 1) Find $\boldsymbol{\phi}^*$ such that $\frac{\partial \mathcal{B}}{\partial \boldsymbol{\phi}} \Big|_{\boldsymbol{\phi}^*} = 0$ (E-step)
 - 2) Calculate $\mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi}^*)$
 - 3) Calculate $\frac{\partial \mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi}^*)}{\partial \boldsymbol{\theta}}$ (M-step)

Step 1) is optimizing the bound with respect to the variational parameters and is therefore equivalent to the E-step. Step 2) is to compute the bound – a task which in many hidden variable problems is easy given the E-step, and step 3) is essentially the same computation as in the M-step, since the M-step means solving the stationarity condition for the model parameters. Note that in step 3) we have exploited that we in step 1) set ϕ to its value at stationarity such that implicit θ dependence through ϕ in the gradient will vanish. The returned function value and gradient can be fed to any gradient based optimizer. In this paper we have chosen the so-called UCMINF described in [13], which is a quasi-Newton method using BFGS update, line search and trust-regions.

3.3 Overrelaxed Adaptive EM

The Easy Gradient Recipe is a very efficient approach for a modest number of parameters to be estimated, but when in the case of for example very overcomplete systems, $\text{length}(\mathbf{x}) \gg \text{length}(\mathbf{s})$, such as images, the number of parameters becomes too large for practical optimization. To deal with these situations, we need to introduce a third optimization approach which in some sense is a compromise between the two already presented. Among the variants of the EM algorithm that modifies the step length we find the so-called Overrelaxed Adaptive EM algorithm, which, in the M-step, takes a larger step in the direction proposed by EM,

$$\theta^{n+1} = \theta^n + \eta(\theta_{EM} - \theta^n)$$

where $\eta \geq 1$. For $\eta = 1$ we retain the ordinary EM algorithm, but for each time we take a successful step forward, the parameter η is increased by a factor above 1 e.g. 2. If the bound is decreasing in a certain step, we undo the step and reset $\eta = 1$. This speeds up the process significantly and a nice feature about the Adaptive Overrelaxed EM is that the computational time spent in each step is reasonable for also a large number of parameters.

3.4 Dealing with Constrained Parameters

Some of the parameters $\theta = \{\mathbf{A}, \Sigma, \nu\}$ are either by definition or application constrained to be positive, positive definite, etc. Within the hierarchical Bayesian framework this is dealt with by imposing priors on these and restricting these priors to be zero in the forbidden domains as is the case for the hidden sources \mathbf{S} . In the empirical Bayes setting, however, constraints can be implemented using a reparametrization and thereby avoid the trouble of the integrals involved in the Bayesian approach.

The mixing matrix \mathbf{A} is considered to be either unconstrained or with positive elements. The positivity constraint is constructed using the exponential function, $(\mathbf{A})_{ij} = e^{(\alpha)_{ij}}$. In this case, the parameter \mathbf{A} is essentially substituted by the underlying parameter α in the parameter set θ . Note that with this

parametrization it holds for any function \mathcal{B} that

$$\frac{d\mathcal{B}}{d[\boldsymbol{\alpha}]_{ij}} = \frac{\partial \mathbf{A}}{\partial [\boldsymbol{\alpha}]_{ij}} \frac{\partial \mathcal{B}}{\partial \mathbf{A}} = [\mathbf{A}]_{ij} \left[\frac{\partial \mathcal{B}}{\partial \mathbf{A}} \right]_{ij} .$$

Setting this derivative to zero can be solved by a simple iterative scheme [6, 9, 22] as long as both the data and the sources are positive:

$$[\mathbf{A}]_{ij} := [\mathbf{A}]_{ij} \frac{[\boldsymbol{\Sigma}^{-1} \mathbf{X} \langle \mathbf{S} \rangle_q^T]_{ij}}{[\boldsymbol{\Sigma}^{-1} \mathbf{A} \langle \mathbf{S} \mathbf{S}^T \rangle_q]_{ij}} .$$

When negative data/sources are encountered the problem has to be solved via somewhat slower quadratic programming techniques.

The noise covariance must be positive definite in order to serve as a covariance, but can be further constrained to be for example diagonal. We consider in this setup noise covariances of the simple isotropic form $\boldsymbol{\Sigma} = e^\beta \mathbf{I}$, the diagonal form $\boldsymbol{\Sigma} = \text{diag}(e^{\beta_1}, \dots, e^{\beta_m})$, and the full parametrization $\boldsymbol{\Sigma} = \boldsymbol{\beta} \boldsymbol{\beta}^T$. The parameters in the source prior may also be constrained, and similar reparametrizations may be implemented.

4 Estimating Source Statistics

Calculating the required statistics and the normalization of the source posterior will in most cases be intractable because the non-Gaussian source prior and multivariate Gaussian Likelihood makes the posterior non-Gaussian multivariate. In this context tractable means we can calculate the normalization constant of the posterior, the marginal Likelihood, $p(\mathbf{X}|\boldsymbol{\theta})$ and posterior moments exactly in polynomial time. In the intractable case, we therefore have to resort to approximate inference techniques. In this section we discuss two deterministic mean field approaches, expectation consistent (EC) and variational (Bayes). In both these approaches variational distributions are used to make the approximations to the marginal Likelihood tractable, but there is an important distinction to be made between the two. In the variational approach a restricted form of the variational distribution q is used to made the calculation of the bound $\mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi})$ tractable. In EC, on the other hand, the aim is only for an approximation $\mathcal{A}(\boldsymbol{\theta}, \boldsymbol{\phi})$ to the log of the marginal likelihood, which will typically be more precise than the bound. But as will be shown below in optimization of the parameters $\mathcal{A}(\boldsymbol{\theta}, \boldsymbol{\phi})$ is used in exactly the same way as $\mathcal{B}(\boldsymbol{\theta}, \boldsymbol{\phi})$.

4.1 Expectation Consistent

The basic idea behind the expectation consistent (EC) framework [17, 15, 16] is to use more than one variational distribution approximation to the posterior. These encode complementary aspects such as prior constraints and the Likelihood term. We will show below that the requirement of consistency between the distributions on the sufficient statistics, i.e. expectation consistency, follows

very naturally when we derive the EC approximation to the marginal Likelihood. We will also give a recipe for attaining the consistency by a sequential iterative approach that alternates between updating each of the distributions.

In instantaneous ICA we can get tractability by choosing a decomposition into two distributions (here $\mathbf{s} = \mathbf{s}_t$ denotes the sources of one time instance t):

$$q(\mathbf{s}) = \frac{1}{Z_q(\boldsymbol{\lambda}_q)} p(\mathbf{s}) \exp(\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \quad (7)$$

$$r(\mathbf{s}) = \frac{1}{Z_r(\boldsymbol{\lambda}_r)} p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) \exp(\boldsymbol{\lambda}_r^T \mathbf{g}(\mathbf{s})) , \quad (8)$$

where the exponential factors are chosen to contain the first and diagonal second moment $\mathbf{g}(\mathbf{s}) = (s_1, \dots, s_M, -\frac{s_1^2}{2}, \dots, -\frac{s_M^2}{2})$, the parameters are denoted by $\boldsymbol{\lambda} = (\gamma_1, \dots, \gamma_M, \Lambda_1, \dots, \Lambda_M)$. Both q and r have distinct vectors $\boldsymbol{\lambda}_q$ and $\boldsymbol{\lambda}_r$ containing these terms. The normalizers are

$$Z_q(\boldsymbol{\lambda}_q) = \int ds p(\mathbf{s}) \exp(\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \quad (9)$$

$$Z_r(\boldsymbol{\lambda}_r) = \int ds p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) \exp(\boldsymbol{\lambda}_r^T \mathbf{g}(\mathbf{s})) . \quad (10)$$

The purpose of the exponential factors in the approximate distributions is to compensate for the factor we have omitted. How to choose the parameters will be explained below. We see that we get tractability with this choice since $q(\mathbf{s})$ is a product of univariate distributions and $r(\mathbf{s})$ is a multivariate Gaussian. Of course the choice of decomposition should be guided not only by tractability but also by quality of the approximation. We expect from central limit theorem arguments that the EC approximation with this decomposition will become better the higher the number of sources with a ‘‘homogeneous’’ connectivity of the mixing matrix [15, 16, 5]. We expect the EC approximation to almost be precise than the variational approximation since it contains the variational since it is a more flexible approximation. To proceed we rewrite the exact marginal Likelihood as

$$\begin{aligned} p(\mathbf{x}|\mathbf{A}, \boldsymbol{\Sigma}) &= \int ds p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) p(\mathbf{s}) = \frac{Z_q(\boldsymbol{\lambda}_q)}{Z_q(\boldsymbol{\lambda}_q)} \int ds p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) p(\mathbf{s}) \\ &= Z_q(\boldsymbol{\lambda}_q) \left\langle p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) \exp(-\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \right\rangle_q , \end{aligned} \quad (11)$$

where

$$\langle \dots \rangle_q = \frac{1}{Z_q(\boldsymbol{\lambda}_q)} \int ds \dots p(\mathbf{s}) \exp(\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \quad (12)$$

denotes an average over $q(\mathbf{s})$. The first step in the EC approximation is to introduce a simpler distribution containing only the exponential factor (a product of univariate Gaussians in this case)

$$u(\mathbf{s}) = \frac{1}{Z_u(\boldsymbol{\lambda}_u)} \exp(\boldsymbol{\lambda}_u^T \mathbf{g}(\mathbf{s})) \quad (13)$$

and exchange the average over q with an average over u . If u shares some key properties with q , e.g. the two first moments, then in many cases the finer details of the distribution will not matter very much:

$$\left\langle p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) \exp(-\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \right\rangle_q \approx \left\langle p(\mathbf{x}|\mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) \exp(-\boldsymbol{\lambda}_q^T \mathbf{g}(\mathbf{s})) \right\rangle_u = \frac{Z_r(\boldsymbol{\lambda}_u - \boldsymbol{\lambda}_q)}{Z_u(\boldsymbol{\lambda}_u)}.$$

Inserting the approximation we arrive at the EC approximation

$$\mathcal{A}(\boldsymbol{\theta}, \boldsymbol{\phi}) \equiv \ln Z_q(\boldsymbol{\lambda}_q) + \ln Z_r(\boldsymbol{\lambda}_u - \boldsymbol{\lambda}_q) - \ln Z_u(\boldsymbol{\lambda}_u) \quad (14)$$

with $\boldsymbol{\phi} = \{\boldsymbol{\lambda}_q, \boldsymbol{\lambda}_u\}$. With a change of variables $\boldsymbol{\lambda}_r \equiv \boldsymbol{\lambda}_u - \boldsymbol{\lambda}_q$ we can also write

$$\ln Z_{\text{EC}}(\boldsymbol{\lambda}_q, \boldsymbol{\lambda}_r) = \ln Z_q(\boldsymbol{\lambda}_q) + \ln Z_r(\boldsymbol{\lambda}_r) - \ln Z_u(\boldsymbol{\lambda}_q + \boldsymbol{\lambda}_r). \quad (15)$$

The second step in the EC approximation is to determine the parameters from the stationarity condition [17] which gives the expectation consistent condition of the three distribution

$$\frac{\partial \ln Z_{\text{EC}}}{\partial \boldsymbol{\lambda}_q} = 0 \quad : \quad \langle \mathbf{g}(\mathbf{s}) \rangle_q = \langle \mathbf{g}(\mathbf{s}) \rangle_u \quad (16)$$

$$\frac{\partial \ln Z_{\text{EC}}}{\partial \boldsymbol{\lambda}_r} = 0 \quad : \quad \langle \mathbf{g}(\mathbf{s}) \rangle_r = \langle \mathbf{g}(\mathbf{s}) \rangle_u \quad (17)$$

with $\boldsymbol{\lambda}_u = \boldsymbol{\lambda}_q + \boldsymbol{\lambda}_r$. At this stationarity point we have the EC approximation

$$\ln p(\mathbf{x}|\mathbf{A}, \boldsymbol{\Sigma}) \approx \mathcal{A}(\boldsymbol{\theta}, \boldsymbol{\phi}).$$

Below we will test this empirically by looking at predictions for moments.

EC Message Passing

Before giving explicit expressions for the marginal Likelihood expression and parameter derivatives for the ICA model, we give a general recipe for attaining the expectation consistent fixed-point which is identical to Minka's expectation propagation (EP) for two approximating factors [10]. This algorithm very often has very good convergence properties, but is not guaranteed to converge. Alternative guaranteed convergent so-called double loop algorithms exist [17]. The details for the ICA-model are given in the Appendix. Iteration k of the algorithm can be sketched as follows:

1. Send messages from r to q
 - Calculate parameters of $u(\mathbf{s})$: Solve for $\boldsymbol{\lambda}_u$: $\langle \mathbf{g}(\mathbf{s}) \rangle_u = \boldsymbol{\mu}_r(k-1) \equiv \langle \mathbf{g}(\mathbf{s}) \rangle_{r(k-1)}$
 - Update $q(\mathbf{x})$: $\boldsymbol{\lambda}_q(k) := \boldsymbol{\lambda}_u - \boldsymbol{\lambda}_r(k-1)$
2. Send messages from q to r
 - Calculate parameters $u(\mathbf{s})$: Solve for $\boldsymbol{\lambda}_u$: $\langle \mathbf{g}(\mathbf{s}) \rangle_u = \boldsymbol{\mu}_q(k) \equiv \langle \mathbf{g}(\mathbf{s}) \rangle_{q(k)}$

- Update $r(\mathbf{s})$: $\boldsymbol{\lambda}_r(k) := \boldsymbol{\lambda}_u - \boldsymbol{\lambda}_q(k)$

$r(k)$ and $q(k)$ denote the distributions q and r computed with the parameters $\boldsymbol{\lambda}_r(k)$ and $\boldsymbol{\lambda}_q(k)$. Convergence is reached when $\boldsymbol{\mu}_r = \boldsymbol{\mu}_q$ since each parameter update ensures $\boldsymbol{\lambda}_r = \boldsymbol{\lambda}_u - \boldsymbol{\lambda}_q$.

EC for the ICA Model

In the following we give the explicit expressions for the EC marginal likelihood expression, eq. (14), moments and the derivatives of the marginal Likelihood approximation with respect to the parameters.

The moments and normalizer of the $q(\mathbf{s}) = \prod_i q_i(s_i)$, $i = 1, \dots, M$, will depend upon the choice of prior. We denote the mean by

$$m_{q,i}(\gamma, \Lambda) = \frac{1}{Z_q(\gamma, \Lambda)} \int ds s p_i(s) \exp(\gamma s - \frac{1}{2} \Lambda s^2) \quad (18)$$

and likewise for the variance $v_{q,i}(\gamma, \Lambda)$. The multivariate Gaussian r -distribution has covariance and mean

$$\boldsymbol{\chi}_r = (\boldsymbol{\Lambda}_r + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A})^{-1} \quad (19)$$

$$\mathbf{m}_r = \boldsymbol{\chi}_r (\boldsymbol{\gamma}_r + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}) \quad (20)$$

and normalizer

$$\ln Z_r = \frac{d-M}{2} \ln 2\pi - \frac{1}{2} \ln \det \boldsymbol{\Sigma} + \frac{1}{2} \ln \det \boldsymbol{\chi}_r + \frac{1}{2} \mathbf{m}_r^T \boldsymbol{\chi}_r^{-1} \mathbf{m}_r - \frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} . \quad (21)$$

The u distribution is a the product of the univariate normals with moments $m_{u,i} = \gamma_{u,i}/\lambda_{u,i}$ and $v_{u,i} = 1/\lambda_{u,i}$. In the propagation algorithm, above and in the Appendix, we need to solve for the parameters in terms of the moments: $\gamma_{u,i} = v_{u,i} m_{u,i}$ and $\lambda_{u,i} = 1/v_{u,i}$. Finally the contribution to the marginal Likelihood from u is given by:

$$\ln Z_u = -\frac{M}{2} \ln 2\pi + \frac{1}{2} \sum_i \ln v_{u,i} + \frac{1}{2} \sum_i \frac{m_{u,i}^2}{v_{u,i}} . \quad (22)$$

Next we consider the derivatives of the marginal Likelihood with respect to \mathbf{A} , $\boldsymbol{\Sigma}$ and $\boldsymbol{\nu}$. When expectation consistency holds then $\frac{\partial \ln Z_{\text{EC}}}{\partial \boldsymbol{\lambda}_q} = \frac{\partial \ln Z_{\text{EC}}}{\partial \boldsymbol{\lambda}_r} = 0$ and we only need to consider the explicit parameter dependence. All \mathbf{A} and $\boldsymbol{\Sigma}$ dependence is contained in $\ln Z_r$, eq. (10), and all $\boldsymbol{\nu}$ dependence in $\ln Z_q$, eq. (9). Stacking the variables, the result—which is most easily derived by considering $\ln Z_r$ as a moment generating function—is very close to eqs. (2) and (3). Compared to these expressions the only difference is that we should now take the average with respect to the r -distribution. Note that although q and r have the same diagonal second moments, they differ on the off-diagonal terms: q has zero covariance since it factorized and r has the Gaussian covariance, eq. (19). It thus makes an important difference what variational distribution we use when calculate the derivatives. Finally, the derivatives with respect to parameters of the prior will be given by eq. (4) with the average being over the q -distribution.

4.2 Variational

The variational approximation can be motivated by the need to find a tractable expression for the bound function eq. (1). This can be achieved choosing the variational distribution in a tractable family. We have basically two possibilities for the ICA model either fully factorized $q(\mathbf{S}) = \prod_{it} q_{it}(s_{it})$ or as a multivariate Gaussian. The latter choice only give tractable expression for the variational bound for some choices of the prior. Here we will consider the fully factorized.

Choosing the variational distribution to be completely factorized it is not likely that a perfect fit to the true source posterior is possible, but in many cases, the approximation will suffice for a successful estimation. We obtain the optimal q (in the factorized family) by setting the functional derivative $\delta\mathcal{B}/\delta q_{it}$ equal to zero (the so-called freeform derivation) [8]. The general and specific solutions are:

$$\begin{aligned} q_{it}(s_{it}|\phi_{it}) &= \frac{1}{c} \exp [\langle \ln p(\mathbf{X}, \mathbf{S}|\boldsymbol{\theta}) \rangle_{q \setminus q_{it}}] \\ &= \frac{1}{Z_q} p_i(s_{it}|\nu_i) \exp [-\frac{1}{2}\Lambda_i s_{it}^2 + \gamma_{it} s_{it}] \end{aligned} \quad (23)$$

where $\langle \dots \rangle_{q \setminus q_{it}} = \int \prod_{i't' \neq it} ds_{i't'} q_{i't'}(s_{i't'}) \dots$ denotes an average over the variational distribution excluding $q_{it}(s_{it})$ and $\mathbf{\Lambda}$ (a vector of length M) and $\boldsymbol{\gamma}$ (a $M \times N$ dimensional matrix) are defined by

$$\mathbf{\Lambda} = \text{diag}(\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A}) \quad (24)$$

$$\boldsymbol{\gamma} = \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} - (\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A} - \text{diag}(\mathbf{\Lambda})) \langle \mathbf{S} \rangle_q . \quad (25)$$

Note how this elegantly both provide us with the structural form of q by eq. (23) and the optimal values of the parametrization by equations for $\mathbf{\Lambda}$ and $\boldsymbol{\gamma}$. Note also, however, that the expression for $\boldsymbol{\gamma}$ depends on the variational mean value and the equations therefore not are closed. Using eq. (23) as a sequential update for $q(\mathbf{S})$ is the coordinate ascent algorithm for the factorized variational distribution and thus guaranteed to converge to a (local) optimum. The sufficient statistics for the variational distribution is the means because it is the only statistic which is necessary to determine the parameter $\boldsymbol{\gamma}$ and $\mathbf{\Lambda}$. We thus write the update equations for the variational distribution in terms of the mean function $\langle s_{it} \rangle_q = m_{q,i}(\gamma_{it}, \Lambda_i)$, eq. (18). Any convergent integral is formally a function and it may seem that we have gained little by reformulating the problem. But for a large and relevant set of source priors, the mean function $m_{q,i}$ have a nice closed form expression. In those cases we have substituted an intractable integral with a non-linear equation which evaluates much faster and more efficiently. The function $m_{q,i}$ is described for a variety of priors in [6] including binary, uniform, exponential (positive), Laplace (bi-exponential) and Gaussian.

A consequence of using the factorized variational distribution is that we will make trivial predictions for the non-diagonal second moments: $\langle s_{it} s_{i't} \rangle_q = \langle s_{it} \rangle_q \langle s_{i't} \rangle_q$ for $i \neq i'$. These second moments are used in in the derivatives of

the bound function with respect to parameters, eqs. (2) and (3). This should be contrasted to the EC and the linear response correction to the variational approximation [6]. The linear response expression for the covariance is given by eq. (19) with $\Lambda_{r,it}$ being dependent upon the solution to the variational equations: $\Lambda_{r,it} = 1/v_{q,it} - \Lambda_i$. This result can thus be seen as an intermediate step between the completely factorized variational approach and EC.

4.3 EC and Variational Comparison

In figure 1 we compare the mean squared approximation error on the first $\langle s_{it} \rangle$ and second moments $\chi_{ii't} = \langle s_{it}s_{i't} \rangle - \langle s_{it} \rangle \langle s_{i't} \rangle$ for a range of signal to noise ratios. The two RMS error measures are defined as

$$\text{Error}_1 = \left[\frac{1}{NM} \sum_{it} (\langle S_{it} \rangle_{\text{exact}} - \langle S_{it} \rangle_{\text{app}})^2 \right]^{1/2} \quad (26)$$

$$\text{Error}_2 = \left[\frac{1}{NM^2} \left(\sum_{ii't} \chi_{ii't,\text{exact}} - \chi_{ii't,\text{app}} \right)^2 \right]^{1/2}, \quad (27)$$

where M is the number of sources. The example is using artificial data from a mixture of Gaussians source prior with two components (with equal weight), with zero mean and variances $\sigma_1^2 = 1, \sigma_2^2 = 0.01$. The number of samples is $N = 2000$ and the sources are mixed with a 2×2 matrix with column vectors $[1 \ 0]^T$ and $[\sqrt{2}/2 \ \sqrt{2}/2]^T$. For each SNR level, defined as $SNR = \text{Tr}(\mathbf{A} \langle \mathbf{s} \mathbf{s}^T \rangle \mathbf{A}^T) / \sigma^2$, a data set with the appropriate noise level is generated and thereafter solved by the various mean field approximations as indicated. In short, Figure 1 shows that expectation consistent method (EC) is much (typically orders of magnitude) more accurate than the variational methods, and that among the variational methods linear response (VarLR) is more accurate than the simple factorized model (VarFct).

5 Software

In this section we will briefly describe the `icaMF` Matlab toolbox¹ that implements the algorithms described in this paper. The basic function call is

$$[\mathbf{S}, \mathbf{A}, \text{loglikelihood}, \text{Sigma}] = \text{icaMF}(\mathbf{X}, \text{par}),$$

where \mathbf{X} is the data matrix, `par` is a list of parameters to algorithm (some of which are described below), \mathbf{S} is the estimated sources, \mathbf{A} is the estimated mixing matrix, `loglikelihood` is the estimated log Likelihood per sample and `Sigma` is the estimated noise covariance (default is scalar valued).

The `par` argument defines a number of settings for the algorithm. Some of the most important are summarized in table 1.

¹The toolbox with demos are available from <http://mole.imm.dtu.dk/>.

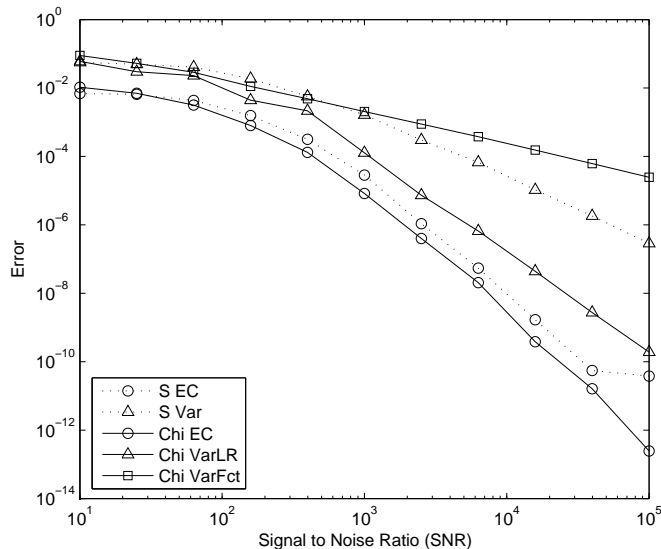


Figure 1: Accuracy of the different posterior approximations. For a Mixture of Gaussian prior, the exact source posterior moments can be calculated and compared to the approximations. The plot show that both for the first moment (S) and the covariance (Chi), the EC is at least one order of magnitude more precise than the variational approach. Furthermore, among the variational approaches is the linear response (VarLR) more accurate than the simple factorized model (VarFct). In [21] we also compare with the saddle-point method. It tends to give a worse approximation than variational.

An additional feature for model selection is a Bayesian information criterion (BIC) function which calculates

$$BIC = \mathcal{L}(\boldsymbol{\theta}) - \frac{|\boldsymbol{\theta}|}{2} \log N ,$$

where $|\boldsymbol{\theta}|$ is the number of parameters we estimate by maximum Likelihood, e.g. the number of free parameters in \mathbf{A} , $\boldsymbol{\Sigma}$ and $\boldsymbol{\nu}$. BIC is an asymptotic expansion for the log of the Likelihood marginalized over all parameter. Figure 2 shows 1) the result of positive ICA for three sources on fMRI brain image time series (described in figure caption) and 2) the output of the function call `icaMFbic(X,par,1:5)`. For comparison, figures 3 and 4 shows the result for positive ICA and standard ICA both for four sources. Although the data has been preprocessed such that it contains both negative and positive values the result coming out of the positive ICA is more clear-cut than standard ICA.

par.	Usage	Default	Examples options
<code>sources</code>	number of sources	<code>size(X,1)</code>	under-/overdetermined, square
<code>optimizer</code>	parameter optimizer	<code>'aem'</code>	<code>'em'</code> , <code>'conjgrad'</code> , <code>'bfgs'</code>
<code>solver</code>	source statistics solver	<code>'ec'</code>	<code>'ec'</code> , <code>'variational'</code>
<code>Sprior</code>	source prior	<code>'mog'</code>	<code>'exponential'</code> , <code>'binary'</code>
<code>method</code>	ICA method	<code>'free'</code>	<code>'constant'</code> , <code>'positive'</code> , <code>'fa'</code> , <code>'ppca'</code>

Table 1: Examples of the `par` settings in `icaMF(X,par)`. More detail are given in `help icaMF`. The ICA methods `par.method` explained: `'free'`, standard ICA, unconstrained mixing matrix and isotropic noise covariance $\Sigma = \sigma^2\mathbf{I}$ optimization and heavy-tailed source prior `'mog'` (fixed mixture of Gaussians); `'constant'`, for test sets, constant mixing matrix and noise covariance; `'positive'`, positive source prior `'exponential'` and `par.Aprior='positive'`; `'fa'`, factor analysis; and `'ppca'`, probabilistic PCA.

Interestingly, BIC for standard ICA tends prefer a much larger number ~ 20 of sources. This is probably because the model is more flexible than positive ICA. Note also that BIC is based upon the assumption of independent samples. This is not true in many case. In this example the samples are pixels in the image which tend to be quite correlated. So the effective number of samples are lower than the actual. Using the effective number of samples will lower the preferred number of sources but it is beyond our aim to estimate the effective number of samples.

6 Testing the Efficiency of the Framework

In this section we compare convergence properties between the optimization schemes proposed and demonstrate the strong improvement of the more advanced methods over the EM algorithm. We also shortly discuss the relation to the framework of non-negative matrix factorization, NMF.

The plots in Figure 5 are projections of the bound function contours and steps in the space of the mixing matrix. The mixing matrix is 2×2 and the space therefore 4 dimensional, but a plane is fitted to the path and bound and steps projected into this two dimensional plane. The noise is isotropic but not estimated in this example, a choice which that makes no difference to the points made. Note that the step size of the EM optimization (EM) is so small, that the dots form a solid line. In total, the EM algorithm use 729 steps to reach the optimal point. For the adaptive overrelaxed EM (AEM), the step size is far greater, and we can also see a failed test step as a detour from the region close to the optimal point. But the AEM cancels this step, resets the step size to 1 and investigate the region close to the optimal point more carefully to reach the optimal point in only 16 steps. The easy gradient approach with a quasi-Newton update (BFGS) is also doing well, reaching the optimal point in 25 steps.

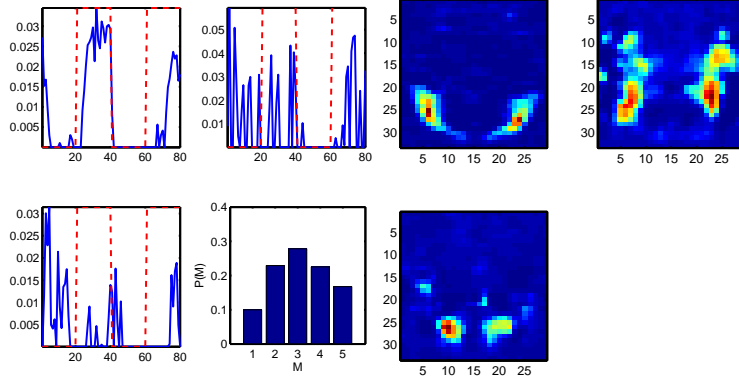


Figure 2: “Spatial” positive ICA on fMRI time-series $\mathbf{X} = \text{Time} \times \text{Image}$ for three sources. The data is described in more detail in [24]. The bar sub-plot shows the posterior probabilities for models computed by BIC for varying number of sources. The left plot is columns in mixing matrix: time-series with visual activation paradigm rest-activation-rest-activation superimposed (dashed line). Right plot shows associated sources images. Note that the decomposition is sorted according to “energy” $E_i = \sum_d A_{di}^2 \sum_t \langle s_{it} \rangle^2$.

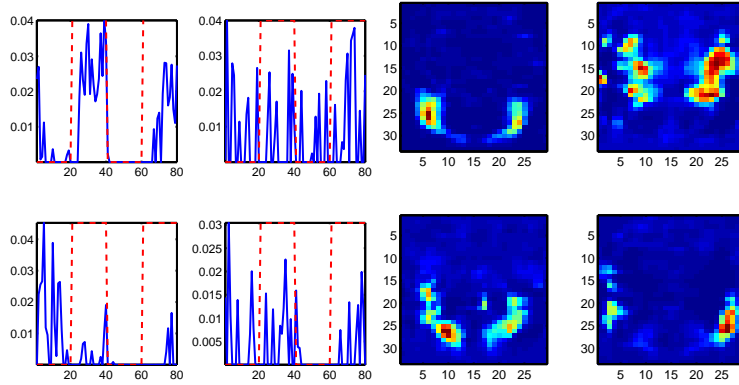


Figure 3: “Spatial” positive ICA on fMRI time-series $\mathbf{X} = \text{Time} \times \text{Image}$ for *four* sources.

Another example which links the slow convergence of the parameters with slow increase of the log likelihood can be seen in Figure 6. The data is an artificial mix of the two speech signals available as a demo in the toolbox and the source priors are chosen to be a mixtures of Gaussians with equal weights, zero mean and variances $\sigma_1^2 = 1, \sigma_2^2 = 0.01$. As the figure shows, the development of the log likelihood approximation is rather different for the three methods. While the easy gradient approach with quasi-Newton update (BFGS) is somewhat slower in the beginning, it quickly maximizes the log likelihood to a stable higher level. As an indicator of the development of the parameters in this process, the inserts show the data and the estimated mixing matrix for the BFGS at the stable level and the EM at iteration 15 and 45. Close inspection of the inserts reveals that the EM algorithm is not at all a perfect estimation at the latter insert and this final fine-tuning takes an excessive amount of iterations.

We have run simulation comparing non-negative matrix factorization (NMF) [9] with positive ICA. The results for the two algorithms are very similar although in some instances where the SNR is not very high, positive ICA gives the most reasonable results. Overall NMF, like InfoMax in the unconstrained case, seems quite robust to the addition of noise. The main advantages of the ICA model are precise estimates of noise level and the marginal likelihood.

7 Conclusion

In this paper we have combined the improved optimization methods with the more advanced source posterior statistics and presented it in a easy-to-use toolbox. Several examples demonstrates the drawback of the traditional EM algorithm and the improved optimization obtained with adaptive overrelaxed EM and the easy gradient recipe equipped with a suitably advanced optimizer. We have also demonstrated the improved accuracy on the estimated source posterior statistics obtained by the use of the expectation consistent (EC) method as compared to the more traditional variational methods.

The upshot is an efficient ICA method which is still sufficiently flexible to encompass constraints on the source priors, mixing matrix or noise covariance. In that sense, the potential of Bayesian ICA is being realized and we believe that the toolbox implementation is presenting an easy interface to an advanced method. With a user-friendly interface, we sincerely hope that practitioners in different research disciplines will also look to the more advanced and flexible ICA methods when analyzing data in the future.

Acknowledgements

We would like to thank Lars Kai Hansen for good discussions and inform the reader that this work is funded (in part) by the Danish Technical Research Council project No. 26-04-0092 Intelligent Sound. (www.intelligentsound.org).

A Solving the EC Equations

In this appendix we give the explicit expressions for iterative scheme for solve the EC equations for the sources at time t , $\mathbf{s} = \mathbf{s}_t$:

- Initialize covariance and mean of r -distribution:

$$\boldsymbol{\chi}_r := (\boldsymbol{\Lambda}_r + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A})^{-1} \quad (28)$$

$$\mathbf{m}_r := \boldsymbol{\chi}_r (\boldsymbol{\gamma}_r + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_t) \quad (29)$$

with $\boldsymbol{\gamma}_r = \mathbf{0}$ and $\boldsymbol{\Lambda}_r$ set such that the covariance is positive definite. It is sufficient to take $\boldsymbol{\Lambda}_r$ to be small positive since $\mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{A}$ is an outer-product form with only non-negative eigenvalues.

Run sequentially over the sources:

1. Send message from r to q_i
 - Calculate parameter of u_i : $\gamma_{u,i} := m_{r,i}/\chi_{r,ii}$ and $\Lambda_{u,i} := 1/\chi_{r,ii}$.
 - Update q_i : $\gamma_{q,i} := \gamma_{u,i} - \gamma_{r,i}$ and $\Lambda_{q,i} := \Lambda_{u,i} - \Lambda_{r,i}$.
 - Update moments of q_i : $m_{q,i} := m_{q,i}(\gamma_{q,i}, \Lambda_{q,i})$ and $\chi_{q,ii} = v_{q,i}(\gamma_{q,i}, \Lambda_{q,i})$.
2. Send message from q_i to r
 - Calculate parameters of u_i : $\gamma_{u,i} := m_{q,i}/\chi_{q,ii}$ and $\Lambda_{u,i} := 1/\chi_{q,ii}$.
 - Update r : $\gamma_{r,i} := \gamma_{u,i} - \gamma_{q,i}$, $\Delta\Lambda_{r,i} := \Lambda_{u,i} - \Lambda_{q,i} - \Lambda_{r,i}$ and $\Lambda_{r,i} := \Lambda_{u,i} - \Lambda_{q,i}$.
 - Update moments of r using Sherman-Morrison identity:

$$\boldsymbol{\chi}_r := \boldsymbol{\chi}_r - \frac{\Delta\Lambda_{r,i}}{1 + \Delta\Lambda_{r,i} [\boldsymbol{\chi}_r]_{ii}} [\boldsymbol{\chi}_r]_i [\boldsymbol{\chi}_r]_i^T \quad (30)$$

$$\mathbf{m}_r := \boldsymbol{\chi}_r (\boldsymbol{\gamma}_r + \mathbf{A}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}_t) . \quad (31)$$

Convergence is reached when and if $\mathbf{m}_r = \mathbf{m}_q$ and $\chi_{r,ii} = \chi_{q,ii}$, $i = 1, \dots, M$. The computational complexity of the algorithm is $\mathcal{O}(M^3 N_{\text{ite}})$, where M is the number of sources, because each Sherman-Morrison update is $\mathcal{O}(M^2)$ and we make M of those in each sweep over the nodes.

References

- [1] H. Attias. A variational Bayesian framework for graphical models. In T. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems 12*, pages 209–215. MIT Press, 2000.
- [2] Anthony J. Bell and Terrence J. Sejnowski. An Information-Maximization Approach to Blind Separation and Blind Deconvolution. *Neural Computation*, 7(6):1129–1159, 1995.

- [3] O. Bermond and Jean Francois Cardoso. Approximate Likelihood for Noisy Mixtures. In *Proceedings of the ICA Conference*, 1999.
- [4] Mark Girolami. A variational Method for Learning Sparse and Overcomplete Representations. *Neural Computation*, 13:2517–2532, 2001.
- [5] T. Heskes, M. Opper, W. Wiegierinck, O. Winther, and O. Zoeter. Approximate inference techniques with expectation constraints. *Journal of Statistical Mechanics: Theory and Experiment*, page P11015, 2005.
- [6] P. Hojen-Sorensen, O. Winther, and L. K. Hansen. Mean-field approaches to independent component analysis. *Neural Computation*, 14:889–918, 2002.
- [7] A. Hyvärinen. Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10:626–634, 1999.
- [8] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction to variational methods for graphical models. *Mach. Learn.*, 37:183–233, 1999.
- [9] Daniel D. Lee and H. Sebastian Seung. Algorithms for non-negative matrix factorization. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 13, pages 556–562, 2001.
- [10] T. Minka. *Expectation Propagation for Approximate Bayesian Inference*. Doctoral dissertation, MIT Media Lab (2001), 2001.
- [11] J. W. Miskin and D. MacKay. Ensemble learning for blind source separation. In S. Roberts and R. Everson, editors, *Independent Component Analysis: Principles and Practice*. Cambridge University Press, 2001.
- [12] Radford M. Neal and Geoffrey Hinton. A new view of the EM algorithm that justifies incremental and other variants. Technical report, Department of Computer Science, University of Toronto, 1993.
- [13] Hans Bruun Nielsen. UCMINF - An Algorithm for Unconstrained Nonlinear Optimization. Technical Report IMM-Rep-2000-19, Technical University of Denmark, 2000.
- [14] Rasmus Kongsgaard Olsson, Tue Lehn-Schiler, and Kaare Brandt Petersen. State-space models - from the EM algorithm to a gradient approach. *Submitted to Neural Computation*, 2005.
- [15] M. Opper and O. Winther. Gaussian processes for classification: Mean field algorithms. *Neural Computation*, 12:2655–2684, 2000.
- [16] M. Opper and O. Winther. Adaptive and self-averaging Thouless-Anderson-Palmer mean field theory for probabilistic modeling. *Phys. Rev. E*, 64:056131, 2001.

- [17] M. Opper and O. Winther. Expectation consistent approximate inference. *Journal of Machine Learning Research*, 2005.
- [18] Kaare Brandt Petersen. *Mean Field ICA*. PhD thesis, Technical University of Denmark, 2005.
- [19] Kaare Brandt Petersen and Ole Winther. Explaining slow convergence of EM in low noise linear mixtures. Technical Report 2005-2, Informatics and Mathematical Modelling, Technical University of Denmark, 2005.
- [20] Kaare Brandt Petersen and Ole Winther. The EM Algorithm in Independent Component Analysis. In *International Conference on Acoustics, Speech, and Signal Processing*, 2005.
- [21] Kaare Brandt Petersen, Ole Winther, and Lars Kai Hansen. On the convergence of EM and VBEM. *Neural Computation*, 17(9), 2005.
- [22] R. Salakhutdinov and S. Roweis. Adaptive overrelaxed bound optimization methods. In *Proceedings of International Conference on Machine Learning, ICML*. International Conference on Machine Learning, ICML, 2003.
- [23] H. Valpola. *Bayesian Ensemble Learning for Nonlinear Factor Analysis*. PhD thesis, Helsinki University of Technology, Espoo, 2000.
- [24] W. Vanduffel, D. Fize, J. B. Mandeville, K. Nelissen, P. Van Hecke, B. R. Rosen, R. B. Tootell, and G. A. Orban. Visual motion processing investigated using contrast agent-enhanced fmri in awake behaving monkeys. *Neuron*, 32:565–577, 2001.
- [25] J. S. Yedidia, W. T. Freeman, and Y. Weiss. Generalized belief propagation. In T. K. Leen, T. G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems (NIPS)*, volume 13, pages 689–695, 2000.

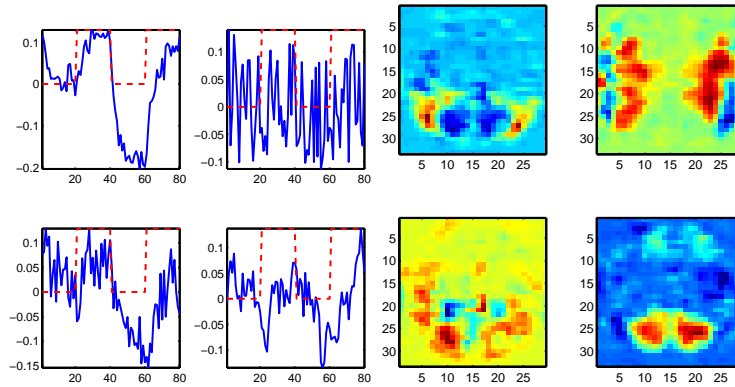


Figure 4: “Spatial” standard ICA on fMRI time-series $\mathbf{X} = \text{Time} \times \text{Image}$ for four sources. Standard ICA corresponds to unconstrained optimization of mixing matrix and heavy-tailed source prior.

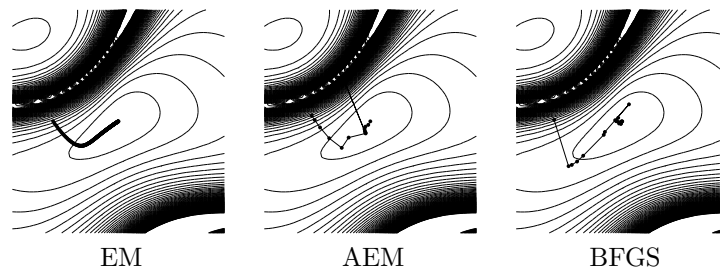


Figure 5: Visualization of convergence. In this plot, the paths of the three different optimization methods are plotted: The EM algorithm (EM), the adaptive overrelaxed EM (AEM), and the Easy gradient with a quasi-Newton optimizer (BFGS). The contours are the bound function projected into the plane of the path.

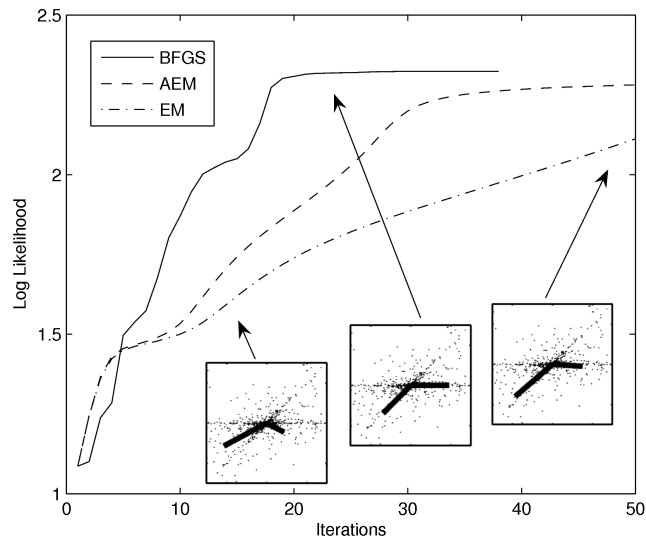


Figure 6: Log likelihood versus iterations for the three different optimization alternatives: The EM algorithm (EM), the over-relaxed adaptive EM (AEM), and the easy gradient recipe with a quasi-Newton method (BFGS). The data is a 2×2 mixing of the speech data. The inserts shows how the estimated mixing matrix fits to the data and demonstrates that the EM algorithm is converging only slowly to the right solution compared to the BFGS method.