

# Variational Bayes for Generalized Autoregressive Models

Stephen J. Roberts and Will D. Penny

**Abstract**—We describe a variational Bayes (VB) learning algorithm for generalized autoregressive (GAR) models. The noise is modeled as a mixture of Gaussians rather than the usual single Gaussian. This allows different data points to be associated with different noise levels and effectively provides robust estimation of AR coefficients. The VB framework is used to prevent overfitting and provides model-order selection criteria both for AR order and noise model order. We show that for the special case of Gaussian noise and uninformative priors on the noise and weight precisions, the VB framework reduces to the Bayesian evidence framework. The algorithm is applied to synthetic and real data with encouraging results.

**Index Terms**—Bayesian inference, generalized autoregressive models, model order selection, robust estimation.

## I. INTRODUCTION

THE STANDARD autoregressive (AR) model assumes that the noise is Gaussian and, therefore, that the AR coefficients can be set by minimizing a least squares cost function. Least squares, however, is known to be sensitive to outliers. Therefore, if the time series is even marginally contaminated by artifacts, the resulting AR coefficient estimates will be seriously degraded (see Bishop [1, p. 209] and Press *et al.* [2, p. 700] for a general discussion of this issue and a number of proposed solutions).

This paper tackles this problem by modeling the noise with a mixture of Gaussians (MoG) in which different data points can be associated with different noise levels. This provides a robust estimation of AR coefficients via a weighted least squares approach; data points associated with high noise levels are down-weighted in the AR estimation step. This approach is thus well suited to situations in which the signal is stationary (and may be modeled via an AR process), whereas the noise is non-Gaussian. The development of AR models with non-Gaussian excitation has a relatively long history. Work presented in [3]–[5] utilizes essentially the same model as we discuss here. In these papers, the issue of parameter inference is rightly tackled with a fully Bayesian approach using Markov chain Monte Carlo (MCMC) sampling methods. Such schemes may be made efficient by exploiting tractable integration for part of the AR model

[5], [6]. Recent use of reversible jump methods [7] applied to AR models [6] allows efficient exploration of model order in the Markov chain. Although these techniques are shown to be well suited to AR model analysis, this paper considers an alternative formalism, which is known as variational Bayes (VB) learning, which offers a tractable (nonsampling) approach. Although the Bayesian methodology has a long history, the use of VB is relatively new; the key idea of VB is to find a tractable approximation to the true posterior density that minimizes the Kullback–Leibler (KL) divergence [8]. Notable recent applications are to principal component analysis [9] and independent component analysis [10]. We have also published short conference papers summarizing some key results for standard autoregressive models [11] and non-Gaussian AR models [12]. To our knowledge, VB has not previously been applied to such models.

Section II describes the autoregressive process as a probabilistic Bayesian model. We then describe the VB method (Section III) and show how it can be applied to generalized AR models (Section IV). Section V describes the VB approach to the standard AR model and compares it with the evidence framework. Section VI presents results on synthetic and real data and compares the difference between model evidence for VB and from a sampling step. Appendices A and B, detailing some important results required elsewhere in the text, are included for completeness.

## II. GENERALIZED AUTOREGRESSIVE MODELS

We define an autoregressive model as

$$y_n = \mathbf{x}_n \mathbf{w} + e_n \quad (1)$$

where  $y_n$  is the  $n$ th value of a time series,  $\mathbf{w}$  is a column vector of AR coefficients (weights), and  $\mathbf{x}_n = [y_{n-1}, y_{n-2}, \dots, y_{n-p}]$  are the  $p$  previous time series values. The AR model has additive noise  $e_n$ , which is usually modeled as a Gaussian. In this paper, however, we model the noise as a one-dimensional (1-D) Gaussian mixture having  $m$  components. Component  $s$  has mixing coefficient  $\pi_s$ , mean  $\mu_s$ , and precision (inverse variance)  $\beta_s$ . We can write the parameters collectively as the vectors  $\mathbf{w}$ ,  $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_m]$ ,  $\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_m]$ , and  $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_m]$ . The weights are drawn from a Gaussian prior (see later) with precision  $\alpha$ . This choice of a Gaussian prior over the parameters is tantamount to a belief that the modeled data sequence has a smooth spectrum. For a more detailed discussion of this issue and the choice of priors in AR models, see [13] and [14].

We concatenate all the parameters into the overall parameter vector  $\boldsymbol{\theta} \stackrel{\text{def}}{=} \{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\beta}, \mathbf{w}, \alpha\}$ . We are given a set of data points

Manuscript received October 6, 2000; revised May 31, 2002. This work was supported, in part, by Grant EESD PRO264 from the U.K. Engineering and Physical Science Research Council (EPSRC). The associate editor coordinating the review of this paper and approving it for publication was Dr. Athina Petropulu.

S. J. Roberts is with the Robotics Research Group, Oxford University, Oxford, U.K. (e-mail: sjrob@robots.ox.ac.uk).

W. D. Penny is with the Wellcome Department of Imaging Neuroscience, University College London, London, U.K. (e-mail: wpenny@fil.ion.ucl.ac.uk).

Publisher Item Identifier 10.1109/TSP.2002.801921.

$D = \{y_n, \mathbf{x}_n\}$  with  $n = 1 \cdots N$ . The likelihood of a data point is given by the mixture model

$$p(y_n | \mathbf{x}_n, \boldsymbol{\theta}) = \sum_{s=1}^m p(s_n = s | \boldsymbol{\pi}) p(y_n | \mathbf{x}_n, s_n, \beta_s, \mu_s, \mathbf{w}) \quad (2)$$

where  $s_n$  is an indicator variable indicating which component of the noise mixture model is selected at presentation of the  $n$ th datum. These are chosen probabilistically according to

$$p(s_n = s | \boldsymbol{\pi}) = \pi_s. \quad (3)$$

The joint likelihood of a data point and indicator variable is

$$p(y_n, s_n | \mathbf{x}_n, \boldsymbol{\theta}) = p(s_n = s | \boldsymbol{\pi}) p(y_n | \mathbf{x}_n, s_n, \beta_s, \mu_s, \mathbf{w}) \quad (4)$$

which, given that the noise samples are assumed independent and identically distributed, gives

$$p(\mathbf{Y}, \mathbf{S} | \boldsymbol{\theta}) = \prod_{n=1}^N p(s_n = s | \boldsymbol{\pi}) p(y_n | \mathbf{x}_n, s_n, \beta_s, \mu_s, \mathbf{w}) \quad (5)$$

over the whole data set, where  $\mathbf{Y} = [y_1, y_2, \dots, y_N]^T$  and  $\mathbf{S} = [s_1, s_2, \dots, s_N]^T$ .

Each component in the noise model is a Gaussian, and hence

$$p(y_n | \mathbf{x}_n, s_n, \beta_{s_n}, \mu_{s_n}, \mathbf{w}) = (2\pi)^{-1/2} \beta_{s_n}^{1/2} \exp\left(-\frac{\beta_{s_n}}{2} (y_n - [\mu_{s_n} + \mathbf{x}_n \mathbf{w}])^2\right). \quad (6)$$

#### A. Model Priors

We choose the standard conjugate priors for each part of the model. These are either normal, gamma, or Dirichlet densities, which we define in Appendices A and B for completeness.

The prior on the model parameters is taken to factorize as

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\pi}) \prod_s p(\beta_s) \prod_s p(\mu_s) p(\mathbf{w} | \alpha) p(\alpha). \quad (7)$$

The prior over the mixing parameters of the MoG is a Dirichlet

$$p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}, \boldsymbol{\lambda}) \quad (8)$$

where the hyperparameters are  $\boldsymbol{\lambda} = [\lambda_0, \lambda_0, \dots, \lambda_0]$  (i.e., a symmetric Dirichlet). The prior over the precisions is a Gamma distribution

$$p(\beta_s) = \text{Ga}(\beta_s; b_0, c_0) \quad (9)$$

the prior over the means is a univariate normal

$$p(\mu_s) = \text{N}_1(\mu_s; m_0, v_0) \quad (10)$$

and the prior over the weights is a zero-mean Gaussian with an isotropic covariance having precision  $\alpha$

$$p(\mathbf{w} | \alpha) = \left(\frac{\alpha}{2\pi}\right)^{p/2} \exp(-\alpha E_W) \quad (11)$$

where

$$E_W = \frac{1}{2} \mathbf{w}^T \mathbf{w}. \quad (12)$$

Finally, the weight precision itself has a Gamma prior

$$p(\alpha) = \text{Ga}(\alpha; b_\alpha, c_\alpha). \quad (13)$$

#### B. Gaussian Noise Models

The standard autoregressive model is recovered when the noise mixture consists of a single Gaussian component with zero mean and precision (inverse variance)  $\beta$ . The parameters now consist of  $\boldsymbol{\theta} = \{\boldsymbol{\beta}, \mathbf{w}, \alpha\}$ . The likelihood of the data set [from (2)] now simplifies to

$$p(\mathbf{Y} | \boldsymbol{\theta}, \mathbf{X}) = \left(\frac{\beta}{2\pi}\right)^{N/2} \exp(-\beta E_D(\mathbf{w})) \quad (14)$$

where

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{x}_n \mathbf{w})^2 = \frac{1}{2} (\mathbf{Y} - \mathbf{X} \mathbf{w})^T (\mathbf{Y} - \mathbf{X} \mathbf{w}) \quad (15)$$

in which, as before,  $\mathbf{Y}$  is a column vector with entries  $(y_1, \dots, y_N)$ , and the  $n$ th row of the matrix  $\mathbf{X}$  contains  $\mathbf{x}_n$ .

#### C. Maximum Likelihood

The standard autoregressive model with Gaussian noise can be implemented using a maximum likelihood (ML) approach. The optimal AR coefficients, given by the maximum of (15) [15], are

$$\mathbf{w}_{ML} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}. \quad (16)$$

From the AR predictions  $\hat{y}_n = \mathbf{x}_n \mathbf{w}_{ML}$ , the ML noise variance  $\sigma_{ML}^2$  can be estimated as

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (\hat{y}_n - y_n)^2. \quad (17)$$

The covariance is then estimated as

$$\mathbf{C}_{ML} = \sigma_{ML}^2 (\mathbf{X}^T \mathbf{X})^{-1}. \quad (18)$$

This ML solution provides a method for initializing the Bayesian analysis.

### III. VARIATIONAL BAYES LEARNING

The central quantity of interest in Bayesian learning is the posterior distribution  $p(\boldsymbol{\theta}, \mathbf{S} | \mathbf{Y})$ , which fully describes our knowledge regarding the parameters of the model. In nonlinear or non-Gaussian models, however, the posterior is often difficult to estimate because, although one may be able to provide values for the posterior for a particular  $\boldsymbol{\theta}, \mathbf{S}$ , the partition function, or normalization term, may involve an intractable integral. To circumvent this problem, two approaches have been developed: the sampling framework and the parametric framework. In the sampling framework, integration is performed via a stochastic sampling procedure such as Markov chain Monte Carlo (MCMC). The latter, however, can be computationally intensive, and assessment of convergence is often problematic. Alternatively, the posterior can be assumed to be of a particular

parametric form. In the Laplace approximation, which is employed, for example, in the “evidence framework,” the posterior is assumed to be Gaussian [16]. This procedure is quick but is often inaccurate. Recently, an alternative parametric method has been proposed: variational Bayes (VB) or “ensemble learning.” A full tutorial on VB is given in [8]. In what follows, we briefly describe the key features.

Given a probabilistic model of the data with AR order  $p$  and  $m$  MoG components in the noise model, the “evidence” or “marginal likelihood” is given by

$$p(\mathbf{Y}|p, m) = \iint p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta}|p, m) d\boldsymbol{\theta} d\mathbf{S}. \quad (19)$$

The log evidence can be written as (dropping  $p, m$  in the conditioning for brevity)

$$\log p(\mathbf{Y}) = \log \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S} \quad (20)$$

where  $q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})$  is, as we will see, a hypothesized, or approximate, posterior density. This has been introduced in both denominator and numerator in (20). Noting that  $q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})$  is a density function (i.e., it integrates to unity), we may apply Jensen’s inequality to obtain a strict bound on the true log posterior as

$$\begin{aligned} \log p(\mathbf{Y}) &= \log \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S} \\ &\geq \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S}. \end{aligned} \quad (21)$$

To see this bound from another perspective, we may write

$$\begin{aligned} \log p(\mathbf{Y}) &= \iint d\boldsymbol{\theta} d\mathbf{S} q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log p(\mathbf{Y}) \\ &= \iint d\boldsymbol{\theta} d\mathbf{S} q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \left( p(\mathbf{Y}) \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})} \right) \\ &= \iint d\boldsymbol{\theta} d\mathbf{S} q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{p(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} \\ &= \iint d\boldsymbol{\theta} d\mathbf{S} q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) (\log p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta}) \\ &\quad - \log p(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) + \log q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) - \log q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})) \\ &= \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S} \\ &\quad + \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})}{p(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S}. \end{aligned} \quad (22)$$

We may write the latter equation as

$$\log p(\mathbf{Y}|p, m) = F(p, m) + KL_{post}(p, m) \quad (23)$$

where (reintroducing the dependence on  $p, m$ )

$$F(p, m) \stackrel{\text{def}}{=} \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta}|p, m)}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S} \quad (24)$$

is known as the negative variational free energy, and

$$KL_{post}(p, m) \stackrel{\text{def}}{=} \iint q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) \log \frac{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})}{p(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}, p, m)} d\boldsymbol{\theta} d\mathbf{S} \quad (25)$$

is the KL divergence [17] between the approximate posterior and the true posterior.

Equation (23) is the fundamental equation of the VB framework. Importantly, because the KL-divergence is always positive [17],  $F(p, m)$  provides a strict *lower bound* on the model evidence. Moreover, because the KL divergence is zero when the two densities are the same,  $F(p, m)$  will become equal to the model evidence when the approximating posterior is equal to the true posterior, i.e., if  $q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) = p(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})$ .

The aim of VB-learning is therefore to maximize  $F(p, m)$  and make the approximate posterior as close as possible to the true posterior. To obtain a practical learning algorithm, we must also ensure that the integrals in  $F(p, m)$  are tractable. One generic procedure for attaining this goal is to assume that the approximating density factorizes over groups of parameters (in physics, this is known as the mean field approximation). Thus, following [18], we consider

$$q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) = q(\mathbf{S}|\mathbf{Y}) \prod_i q(\boldsymbol{\theta}_i|\mathbf{Y}) \quad (26)$$

where  $\boldsymbol{\theta}_i$  is the  $i$ th group of parameters. The distributions that maximize the negative free energy can then be shown to be of the following form (see Appendix A), which, here, is shown for parameter group  $\boldsymbol{\theta}_i$

$$q(\boldsymbol{\theta}_i|\mathbf{Y}) = \frac{\exp[I(\boldsymbol{\theta}_i)]}{\int \exp[I(\boldsymbol{\theta}_i)] d\boldsymbol{\theta}_i} \quad (27)$$

where

$$I(\boldsymbol{\theta}_i) \stackrel{\text{def}}{=} \int q(\boldsymbol{\theta}^{\setminus i} | \mathbf{Y}) \log p(\mathbf{Y}, \mathbf{S}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}^{\setminus i} \quad (28)$$

and  $\boldsymbol{\theta}^{\setminus i}$  denotes all the parameters *not* in the  $i$ th group. For models having suitable priors, the above equations are available in closed analytic form. This leads to a set of coupled update rules. Iterated application of these leads to the desired maximization.

#### A. Model Order Selection

By computing the evidence for models of different order (i.e.,  $p$  and  $m$ ), we can estimate the posterior distribution over  $p$  and  $m$  using the negative free energy. We may thus use this to perform model order selection. The posterior over  $p, m$  is given via Bayes’ theorem as

$$P(p, m|\mathbf{Y}) = \frac{P(\mathbf{Y}|p, m)P(m, p)}{\sum_{p', m'} P(\mathbf{Y}|p', m')P(m', p')} \quad (29)$$

where, for example, we may have uniform priors over model orders  $p$  and  $m$ .

The evidence is estimated using  $F(p, m)$ . To see that  $F(p, m)$  is an intuitively suitable model order criterion, we decompose it as follows. Using  $p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta}) = p(\mathbf{Y}, \mathbf{S}|\boldsymbol{\theta})p(\boldsymbol{\theta})$  and assuming that the approximate posterior factorizes as in (26), i.e.,  $q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) = q(\boldsymbol{\theta}|\mathbf{Y})q(\mathbf{S}|\mathbf{Y})$ , we can write

$$F(p, m) = L_{av}(p, m) - KL_{prior} \quad (30)$$

where

$$L_{av}(p, m) = \iint q(\boldsymbol{\theta}|\mathbf{Y})q(\mathbf{S}|\mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{S}|\boldsymbol{\theta})}{q(\mathbf{S}|\mathbf{Y})} d\boldsymbol{\theta} d\mathbf{S} \quad (31)$$

corresponds to the average likelihood, and

$$KL_{prior} = \int q(\boldsymbol{\theta}|\mathbf{Y}) \log \frac{q(\boldsymbol{\theta}|\mathbf{Y})}{p(\boldsymbol{\theta})} d\boldsymbol{\theta} \quad (32)$$

is the KL divergence between the approximate posterior  $q(\boldsymbol{\theta}|\mathbf{Y})$  and the prior  $p(\boldsymbol{\theta})$ . Now, because the KL term increases with the number of model parameters, it acts as a penalty term in (30), which penalizes more complex models.

As the number of samples increases, the parameter posterior becomes sharply peaked about the most probable values (which are also the ML values)  $\hat{\boldsymbol{\theta}}$ . It can then be shown that in the large sample limit  $N \rightarrow \infty$ ,  $F(p, m)$  becomes equivalent to the Bayesian information criterion (BIC) [18], [19]

$$\text{BIC}(p, m) = \log p(\mathbf{Y}|\hat{\boldsymbol{\theta}}) - \frac{N_{par}}{2} \log N \quad (33)$$

where  $N_{par}$  is the total number of adjustable parameters in the model. The BIC is itself equal to the negative of the minimum description length (MDL) measure, i.e.,  $\text{BIC}(p, m) = -\text{MDL}(p, m)$ . These popular model order selection criteria can therefore be seen as limiting cases of the VB framework.

#### IV. VB FOR GENERALIZED AR MODELS

To apply VB to non-Gaussian autoregressive models, we approximate the posterior distribution over parameters with the factorized density

$$q(\boldsymbol{\theta}|\mathbf{Y}) = q(\boldsymbol{\pi}|\mathbf{Y})q(\boldsymbol{\beta}|\mathbf{Y})q(\boldsymbol{\mu}|\mathbf{Y})q(\mathbf{w}|\mathbf{Y})q(\boldsymbol{\alpha}|\mathbf{Y}) \quad (34)$$

and the posterior distribution over hidden variables by  $q(\mathbf{S}|\mathbf{Y})$ . We then set each distribution to maximize  $F(p, m)$  by applying (27) and substituting in the joint log-likelihood of the generalized AR model

$$\log p(\mathbf{Y}, \mathbf{S}|\boldsymbol{\theta}) = \sum_{n=1}^N \log \pi_s + \frac{1}{2} \log \beta_s - \frac{1}{2} \beta_s (y_n - [\mu_s + \mathbf{x}_n \mathbf{w}])^2. \quad (35)$$

Given our choice of prior and model likelihood, the optimal approximating densities take the following forms. For the precisions, we have  $q(\boldsymbol{\beta}|\mathbf{Y}) = \prod_s q(\beta_s|\mathbf{Y})$  and Gamma densities

$$q(\beta_s|\mathbf{Y}) = \text{Ga}(\beta_s; b_s, c_s) \quad (36)$$

for the means, we have  $q(\boldsymbol{\mu}) = \prod_s q(\mu_s)$  and univariate normal densities

$$q(\mu_s|\mathbf{Y}) = \text{N}_1(\mu_s; m_s, v_s) \quad (37)$$

for the mixing coefficients, we have a Dirichlet

$$q(\boldsymbol{\pi}|\mathbf{Y}) = \text{Dir}(\boldsymbol{\pi}; \lambda_1, \dots, \lambda_m) \quad (38)$$

for the weights, we have a multivariate normal

$$q(\mathbf{w}|\mathbf{Y}) = \text{N}(\mathbf{w}; \hat{\mathbf{w}}, \hat{\mathbf{C}}) \quad (39)$$

and for the weight precision, we have a Gamma density

$$q(\boldsymbol{\alpha}|\mathbf{Y}) = \text{Ga}(\boldsymbol{\alpha}; b'_\alpha, c'_\alpha). \quad (40)$$

We note finally that for the indicator posteriors, we have  $q(\mathbf{S}|\mathbf{Y}) = \prod_n q(s_n|\mathbf{Y})$ .

Each distribution is updated as described in what follows. First, the responsibilities (i.e., the probabilities over the indicator variable set) of the Gaussian components in the noise model are updated, along with the hyperparameters that govern their distribution. Second, the component means and precisions are updated as are the AR coefficients. All hyperparameters associated with the distributions over these variables are also updated during this second step. We detail the specific update equations for the variables in the following subsections.

##### A. Step 1

1) *Component Responsibilities:* We first define some intermediate variables

$$\begin{aligned} \log \tilde{\pi}_s &\stackrel{\text{def}}{=} \Psi(\lambda_s) - \Psi\left(\sum_{s'} \lambda_{s'}\right) \\ \log \tilde{\beta}_s &\stackrel{\text{def}}{=} \Psi(c_s) + \log b_s \\ \bar{\beta}_s &\stackrel{\text{def}}{=} b_s c_s \end{aligned} \quad (41)$$

where  $\Psi(\cdot)$  is the digamma function [2], and

$$\begin{aligned} \tilde{\sigma}_s^2(n) &= \int q(\mathbf{w}|\mathbf{Y})q(\mu_s|\mathbf{Y})(y_n - [\mu_s + \mathbf{x}_n \mathbf{w}])^2 d\mathbf{w} d\mu_s \\ &= y_n^2 - 2m_s y_n - 2\hat{y}_n y_n + m_s^2 + v_s \\ &\quad + 2m_s \hat{y}_n + \mathbf{x}_n \hat{\mathbf{C}} \mathbf{x}_n^T + \hat{y}_n^2 \end{aligned} \quad (42)$$

in which  $\hat{y}_n = \mathbf{x}_n \hat{\mathbf{w}}$  is the mean AR prediction. If we define  $\gamma_s^n \stackrel{\text{def}}{=} q(s_n|\mathbf{Y})$  to be the posterior of the indicator variable for the  $s$ th component of the MoG (i.e., the probability that component  $s$  is responsible for data point  $y_n$ ), then this step consists of updating the indicator posterior according to

$$\tilde{\gamma}_s^n = \tilde{\pi}_s \tilde{\beta}_s^{1/2} \exp[-\frac{1}{2} \bar{\beta}_s \tilde{\sigma}_s^2(n)]. \quad (43)$$

These variables are then normalized to give

$$\gamma_s^n = \frac{\tilde{\gamma}_s^n}{\sum_{s'} \tilde{\gamma}_{s'}^n}. \quad (44)$$

##### B. Step 2

We now define some intermediate variables

$$\begin{aligned} \bar{\pi}_s &\stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \gamma_s^n \\ \bar{N}_s &\stackrel{\text{def}}{=} N \bar{\pi}_s \\ \bar{y}_s &\stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \gamma_s^n y_n \\ \bar{\mathbf{x}}_s &\stackrel{\text{def}}{=} \frac{1}{N} \sum_{n=1}^N \gamma_s^n \mathbf{x}_n \end{aligned} \quad (45)$$

TABLE I  
PSEUDO-CODE FOR VB-AR UPDATES

---

```

initialize;
WHILE ( $\Delta F(p, m) > \text{tolerance}$ )
  update indicator posteriors,  $\tilde{\gamma}_s^n$  using Equations 43,44;
  update component means, precisions & associated hyperparameters using Equations 46-50;
  update AR coefficients and precisions using Equations 51, 52;
  calculate negative free energy,  $F(p, m)$  using Equation 55;
  calculate  $\Delta F(p, m) \stackrel{\text{def}}{=} |F(p, m)^{\text{new}} - F(p, m)^{\text{old}}|$ ;
END WHILE;

```

---

i.e.,

$\bar{\pi}_s$  proportion of data associated with component  $s$ ;  
 $\bar{N}_s$  number of data points associated with component  $s$  and the quantities  $\bar{\mathbf{x}}_s$ ;  
 $\bar{y}_s$  weighted data values.

1) *Component Mixing Fractions*: The hyperparameters,  $\{\lambda_s\}$  of the Dirichlet [Equation (38)] are updated in the standard manner by adding the data counts to the prior counts, i.e.,

$$\lambda_s = \bar{N}_s + \lambda_0. \quad (46)$$

2) *Component Precisions*: If we define the expected variance of component  $s$  as

$$\tilde{\sigma}_s^2 = \frac{1}{N} \sum_{n=1}^N \gamma_s^n \tilde{\sigma}_s^2(n) \quad (47)$$

then the hyperparameters for the precisions [see (36)] are updated as

$$\frac{1}{b_s} = \frac{N}{2} \tilde{\sigma}_s^2 + \frac{1}{b_0}$$

$$c_s = \frac{\bar{N}_s}{2} + c_0. \quad (48)$$

We can understand these equations by considering the corresponding mean variance (the inverse of the mean precision), which is given by  $1/(b_s c_s)$ . If we ignore terms involving the prior, this comes out to be  $\tilde{\sigma}_s^2/\bar{\pi}_s$ , which is the expected variance of that component reweighted according to the number of examples for which the component is responsible.

3) *Component Means*: In this step, the hyperparameters governing the posterior distribution over the component means are updated [see (37)]. If we first define the means and precisions (i.e., the hyperparameters) of the component means as *estimated from the data* as

$$m_{\text{data}}(s) = (\bar{\mathbf{x}}_s \hat{\mathbf{w}} - \bar{y}_s)/\bar{\pi}_s$$

$$\tau_{\text{data}}(s) = \bar{N}_s \bar{\beta}_s \quad (49)$$

then the posteriors of these hyperparameters are given by

$$\tau_s = \tau_0 + \tau_{\text{data}}(s)$$

$$m_s = \frac{\tau_0}{\tau_s} m_0 + \frac{\tau_{\text{data}}(s)}{\tau_s} m_{\text{data}}(s) \quad (50)$$

where  $\tau_0 = 1/v_0$  is the prior precision, and  $\tau_s = 1/v_s$  is the posterior precision.

4) *AR Coefficients*: The distribution of AR coefficients is multivariate normal with hyperparameters governing the mean

and covariance of the distribution [see (39)]. These hyperparameters are updated via

$$\hat{\mathbf{C}} = \left( \sum_s \bar{\beta}_s \mathbf{X}^T \Gamma_s \mathbf{X} + \bar{\alpha} \mathbf{I} \right)^{-1}$$

$$\hat{\mathbf{w}} = \hat{\mathbf{C}} \sum_s \bar{\beta}_s (\mathbf{X}^T \Gamma_s \mathbf{Y} + m_s \mathbf{X}^T \Gamma_s \mathbf{1}) \quad (51)$$

where  $\Gamma_s = \text{diag}([\gamma_s^1, \gamma_s^2, \dots, \gamma_s^N])$ , the  $n$ th row of  $\mathbf{X}$  contains  $\mathbf{x}_n$ , the  $n$ th entry in the column vector  $\mathbf{Y}$  is  $y_n$ , and  $\mathbf{1}$  is a column vector of 1s of length  $N$ .

5) *AR Coefficient Precisions*: The update equations for the coefficient precision and associated hyperparameters [see (40)] are

$$1/b'_\alpha = \frac{1}{2} \hat{\mathbf{w}}^T \hat{\mathbf{w}} + \frac{1}{2} \text{Tr}(\hat{\mathbf{C}}) + \frac{1}{b_\alpha}$$

$$c'_\alpha = \frac{p}{2} + c_\alpha$$

$$\bar{\alpha} = b'_\alpha c'_\alpha \quad (52)$$

in which  $\text{Tr}(\cdot)$  denotes the trace of a matrix.

### C. Co-Mean Noise Model

While it is possible to use a full Gaussian mixture model for the noise (as described previously), the numerical examples in this paper use a simplified model in which all the means are zero. The resulting model then resembles weighted least squares, where the weights are given by the precisions and responsibilities of the noise components. This simplifies a number of update equations. For step 1, as before, we now have [cf. (42)]

$$\tilde{\sigma}_s^2(n) = y_n^2 - 2\hat{y}_n y_n + \mathbf{x}_n \hat{\mathbf{C}} \mathbf{x}_n^T + \hat{y}_n^2. \quad (53)$$

In this case, (49) and (50) are dispensed with, and the second line in (51) simplifies to

$$\hat{\mathbf{w}} = \hat{\mathbf{C}} \sum_s \bar{\beta}_s \mathbf{X}^T \Gamma_s \mathbf{Y}. \quad (54)$$

### D. Pseudo-Code

The previous steps may be conveniently implemented in the algorithm shown in pseudo-code form in Table I.

### E. Practicalities

The prior distribution  $p(\boldsymbol{\pi})$  is chosen to be uniform, and we set  $\lambda_0 = 5$ . For the Gamma distributions, we use vague priors

(see, e.g., [16]); for  $p(\beta_s)$ , we set  $b_0 = 10^3$ ,  $c_0 = 10^{-3}$ , and for  $p(\alpha)$ , we set  $b_\alpha = 10^3$  and  $c_\alpha = 10^{-3}$ . It is noted that we do not find particular sensitivity to these values.

The posterior distributions  $q(\boldsymbol{\pi}|\mathbf{Y})$ ,  $q(\beta_s|\mathbf{Y})$ ,  $q(\alpha|\mathbf{Y})$ , and  $q(\mathbf{w}|\mathbf{Y})$  have the parameters  $\lambda_s$ ,  $b_s$ ,  $c_s$ ,  $b'_\alpha$ ,  $c'_\alpha$ ,  $\hat{\mathbf{w}}$ , and  $\hat{\mathbf{C}}$ , which are initialized as follows. The posterior for the AR coefficients is initialized using the ML solutions of (16)–(18). The posterior for the weight precision is then initialized using (52). The remaining posteriors are set as follows. We first calculate the errors  $e(n)$  from the ML model and then define a new variable  $z(n) = |e(n) - \bar{e}|$ , which is the absolute deviation of each error from the mean error  $\bar{e}$ . We then apply  $k$ -means clustering [1] to  $z(n)$ , which results in mixing coefficients  $\lambda_z(s)$  and means  $m_z(s)$ . We then set  $\lambda_s = 100\lambda_z(s)$ . The parameters  $b_s$  and  $c_s$  are then set to achieve means of  $1/m_z(s)^2$  and variances of  $\text{var}(1/m_z)$  (the mean and variance of a Gamma density are  $bc$  and  $b^2c$ , respectively).

The VB equations are then applied iteratively until a consistent solution is reached. Convergence is measured by evaluating the negative free energy

$$F(p, m) = L_{av} - KL(\mathbf{w}) - KL(\boldsymbol{\alpha}) - KL(\boldsymbol{\pi}) - KL(\boldsymbol{\beta}) \quad (55)$$

where we use the shorthand notation  $KL(a) \stackrel{\text{def}}{=} KL(q(a|\mathbf{Y})||p(a))$ . The first term of (55) is given as

$$L_{av} = H(q(\mathbf{S}|\mathbf{Y})) + \sum_{s=1}^m \bar{N}_s \left( \log \tilde{\pi}_s + \frac{1}{2} \log \tilde{\beta}_s \right) - \frac{1}{2} N \sum_{s=1}^m \bar{\beta}_s \tilde{\sigma}_s^2 - \frac{1}{2} N \log 2\pi. \quad (56)$$

The entropy over the hidden variables is

$$H(q(\mathbf{S}|\mathbf{Y})) = - \sum_{n=1}^N \sum_{s=1}^m \gamma_s^n \log \gamma_s^n. \quad (57)$$

The KL terms for normal, gamma, and Dirichlet densities are given in Appendices A and B. As the update for the AR coefficients is computationally more intensive than the updates for the other parameters (as it involves a matrix inversion), we perform this step only once every  $W_t$  iterations. In our experiments, we used  $W_t = 5$ . We evaluate  $F(p, m)$  every  $W_t$  iterations (after the AR updates) and terminate optimization if the proportionate increase from one evaluation to the next is less than 0.01%.

We note that the formalism adopted in this paper does not, unlike many AR modeling approaches, *guarantee* a stable model (i.e., the poles of the characteristic function do not lie outside the unit circle in the  $z$ -domain). Although it is possible to enforce stability by reflecting poles across the unit circle, we do not regard this as an elegant solution. Placing stability priors on the *reflection coefficients* of the model is another possibility (as described in [20]), but this leaves analytically troublesome descriptions of the AR coefficients themselves for which (slow) sample-based approaches must therefore be used. It is noted, however, that in all examples in this paper, and indeed all data analyzed thus far, no unstable models have been found.

## V. GAUSSIAN AUTOREGRESSIVE MODEL

The standard autoregressive model is recovered when the mixture consists of a single zero-mean Gaussian. Because of the factored form of the prior distribution, the form of the approximate posterior that maximizes the negative free energy is also of factorized form (we note that unlike the general case in variational Bayes methods, we do not have to *assume* this), i.e.,

$$q(\boldsymbol{\theta}|\mathbf{Y}) = q(\mathbf{w}|\mathbf{Y})q(\alpha|\mathbf{Y})q(\beta|\mathbf{Y}). \quad (58)$$

The weight posterior is a normal density  $q(\mathbf{w}|\mathbf{Y}) = \mathcal{N}(\mathbf{w}; \hat{\mathbf{w}}, \hat{\mathbf{C}})$ . By inspection of (51), noting that there is now only a single value of  $\beta$ , that  $m_1 = 0$ , and that the weighting matrix  $\Gamma_s$  is an identity matrix, we get

$$\hat{\mathbf{C}} = \left( \hat{\beta} \mathbf{X}^T \mathbf{X} + \hat{\alpha} \mathbf{I} \right)^{-1} \\ \hat{\mathbf{w}} = \hat{\mathbf{C}} \hat{\beta} \mathbf{X}^T \mathbf{Y}. \quad (59)$$

The weight precision posterior is a Gamma density with parameters as before [see (52)]. The noise precision posterior is  $q(\beta|\mathbf{Y}) = \text{Ga}(\beta; b'_\beta, c'_\beta)$ , where

$$1/b'_\beta = \bar{E}_D(\hat{\mathbf{w}}) + \frac{1}{b_\beta} \\ c'_\beta = \frac{N}{2} + c_\beta \\ \hat{\beta} = b'_\beta c'_\beta \quad (60)$$

and

$$\bar{E}_D(\hat{\mathbf{w}}) = E_D(\hat{\mathbf{w}}) + \frac{1}{2} \text{Tr} \left( \hat{\mathbf{C}} \mathbf{X}^T \mathbf{X} \right). \quad (61)$$

We note that in previous work, MacKay has derived VB updates for the weight and weight precisions in a linear regression model [21] and that, reassuringly, (52) and (59) are identical to his. We note that these update rules can also be derived directly, i.e., without considering the model as a special case of GAR.

For Gaussian noise models, the negative free energy simplifies to

$$F(p) = L_{av} - KL(\mathbf{w}) - KL(\alpha) - KL(\beta) \quad (62)$$

where

$$L_{av} = \frac{N}{2} \langle \log \beta \rangle - \hat{\beta} \bar{E}_D - \frac{N}{2} \log 2\pi \quad (63)$$

and

$$\langle \log \beta \rangle = \Psi(c'_\beta) + \log b'_\beta \quad (64)$$

where  $\Psi()$  is the digamma function.

An alternative Bayesian approach is that offered by the evidence framework [16]. Given that the observation noise is zero-mean Gaussian with precision  $\beta$  and that the weights are drawn from a prior distribution with zero mean and isotropic covariance having precision  $\alpha$ , the posterior distribution is Gaussian with mean and covariance, which is also given by (59). In the evidence framework, there are no priors over  $\alpha$  or

$\beta$ . The evidence for the model is therefore given by (see, e.g., [1, p. 409])

$$p(\mathbf{Y}) = \int p(\mathbf{Y}|\mathbf{w}, \beta)p(\mathbf{w}|\alpha) d\mathbf{w}. \quad (65)$$

The log of the evidence can then be written as

$$EV(p) = -\alpha E_W - \beta E_D(\hat{\mathbf{w}}) + \frac{1}{2} \log |\mathbf{C}| + \frac{p}{2} \log \alpha + \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi. \quad (66)$$

The precision parameters are then set to maximize the evidence as

$$\alpha = \frac{\gamma}{2E_W} \quad (67)$$

$$\beta = \frac{N - \gamma}{2E_D(\hat{\mathbf{w}})} \quad (68)$$

where  $\gamma$ , which is the number of ‘‘well-determined’’ coefficients, is given by

$$\gamma = p - \alpha \text{Tr}(\mathbf{C}) \quad (69)$$

where  $\gamma$  is calculated using the ‘‘old’’ value of  $\alpha$ . The update for  $\alpha$  is therefore an implicit equation. We can also write it as the explicit update

$$\alpha = \frac{p}{2E_W + \text{Tr}(\mathbf{C})}. \quad (70)$$

This is equivalent to the update for  $\hat{\alpha}$  from the VB framework [see (52)] if we choose uninformative priors on  $\alpha$ . Similarly, if we choose uninformative priors for  $\beta$ , then the VB update is [from (60)]

$$\hat{\beta} = \frac{N}{2E_D(\hat{\mathbf{w}}) + \text{Tr}(\hat{\mathbf{C}}\mathbf{X}^T\mathbf{X})} \quad (71)$$

which is equivalent to (68) as, from an eigendecomposition of  $\hat{\mathbf{C}}$ , it can be shown that  $\text{Tr}(\hat{\mathbf{C}}\mathbf{X}^T\mathbf{X}) = \gamma/\beta$ .

#### A. Model Order Selection

As mentioned Section III-A, the BIC model-selection criterion is a limiting case of the VB framework. For the AR model, we have

$$\text{BIC}(p) = \frac{N}{2} \log \beta_{ML} - \frac{p}{2} \log N \quad (72)$$

where  $\beta_{ML}$  is the precision of the ML model.

The model order selection criterion for the evidence framework [see (66)] can be rewritten in terms of the KL-divergence between the weight posterior and the weight prior  $KL(\mathbf{w})$ . For the prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{o}, (1/\alpha)\mathbf{I})$  and the posterior  $q(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \hat{\mathbf{w}}, \hat{\mathbf{C}})$ , it can be shown that

$$KL(\mathbf{w}) = -\frac{1}{2} \log |\mathbf{C}| - \frac{p}{2} \log \alpha + \alpha E_w - \frac{1}{2} \gamma. \quad (73)$$

Because  $\beta E_d = (1/2)(N - \gamma)$  [see (68)], we can combine (73) with (66) to give

$$EV(p) = \frac{N}{2} \log \beta - KL(\mathbf{w}) \quad (74)$$

where we have dropped the constant term  $-(1/2)N(1 + \log 2\pi)$ . The VB criterion is given by

$$F(p) = \frac{N}{2} \langle \log \beta \rangle - KL(\mathbf{w}) - KL(\alpha) - KL(\beta) \quad (75)$$

where, again, we have dropped the constant term  $-(1/2)N(1 + \log 2\pi)$ . The evidence and VB criteria are therefore identical, except for the divergences of  $\alpha$  and  $\beta$ . However, as these do not scale with  $p$ , we can ignore them. We also note that for large data sets, the distribution over  $\beta$  is sufficiently peaked for the first expectation term not to make a significant difference, i.e.,  $\langle \log \beta \rangle \approx \log \hat{\beta}$ . We may also thus see the evidence approach as a limiting case of the VB framework.

If we assume, *a priori*, that different model orders are equally likely, then the posterior probability distribution over model order is given by (29). This can also be applied to BIC/MDL.

## VI. RESULTS

### A. Synthetic Data: Gaussian Noise Model

We generated multiple three second-blocks of 128 Hz data from AR(6) models with varying signal-to-noise ratios (SNRs). We also generated data from AR(25) models, again with varying SNRs. We generated ten data sets of each type (each with a different realization of the noise process) and applied the VB and BIC/MDL model order selection methods. Typical results are shown in Figs. 1 and 2, which are for  $p = 6$ , SNR = 2 and  $p = 25$ , SNR = 20, respectively.<sup>1</sup> We have plotted the probability of each model order averaged over the ten runs (the averaging was done on the probabilities rather than the log probabilities in order to show a greater spread). For both the AR(6) and AR(25) data, model order selection is more difficult at lower SNR. For the AR(6) data, both methods essentially pick out the correct model order, with VB overestimating it on a small proportion of runs. For the AR(25) data, the BIC/MDL criterion significantly underestimates the model order. On longer blocks of data (e.g., 10 s), both methods were seen to converge to the same estimate of model order (as predicted by theory).

### B. Synthetic Data: Non-Gaussian Noise Model

We generated 3 s of synthetic data from an AR(5) model with coefficients

$$\mathbf{w}_{true} = [-1.8517, 1.3741, 0.1421, -0.6852, 0.3506]^T \quad (76)$$

and noise drawn from a two-component Gaussian mixture model with mixing coefficients  $\pi_1 = 0.9$ ,  $\pi_2 = 0.1$  and variances  $\sigma_1^2 = 1$ ,  $\sigma_2^2 = 100$ . We repeated the data generation process ten times and used the negative free energy  $F(p, m)$  as a model order selection criterion. The results in Figs. 3–5 show that the VB model order selection criterion selects both the correct AR order and the correct noise model order. For more on VB model order selection, including a comparison with the minimum description length (MDL) criterion, see [11].

We also calculated the error with which the AR coefficients were estimated using the measure  $\|\hat{\mathbf{w}} - \mathbf{w}_{true}\|$ . Over the ten

<sup>1</sup>The SNR was altered by adjusting the noise variance.

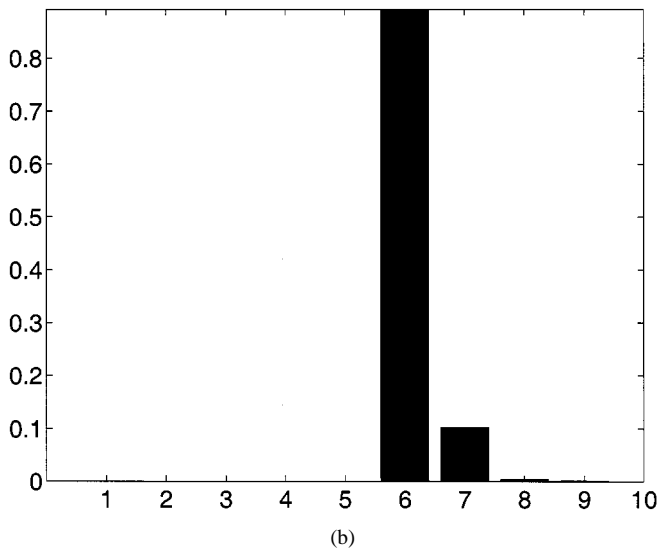
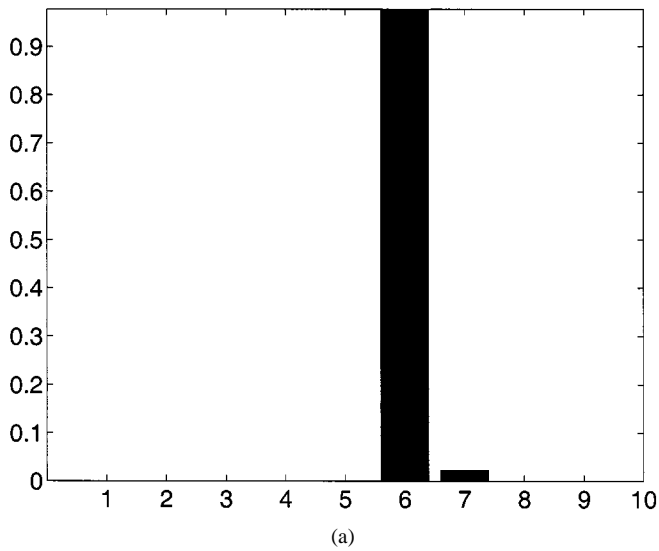


Fig. 1. Model order selection for AR(6) data (SNR = 2) with (a) BIC/MDL and (b) VB with 3-s blocks of data.

runs, the errors from the non-Gaussian AR model with two components were less than those from the Gaussian AR model by an average factor of six.

### C. EEG Data: Gaussian Noise Model

We applied AR models to short blocks of EEG data recorded while a subject was awake or in an anesthetized state. Fifty 1-s blocks were selected randomly from each 30-min recording, and BIC/MDL and VB model order selection criteria were compared.

The results from VB in Fig. 6 show that in the waking state, there are modes at  $p = 8$  and  $p = 13$ , and in the anesthetized state, the modes are at  $p = 3$  and  $p = 8/9$ . The results clearly show a decrease in model order from the waking state to the anesthetized state; according to a t-test, this decrease is significant at the 0.0004 level. The results from BIC/MDL, however, showed no such discrimination. The difference of means was not significant; from a t-test, we get a significance level of 0.062.

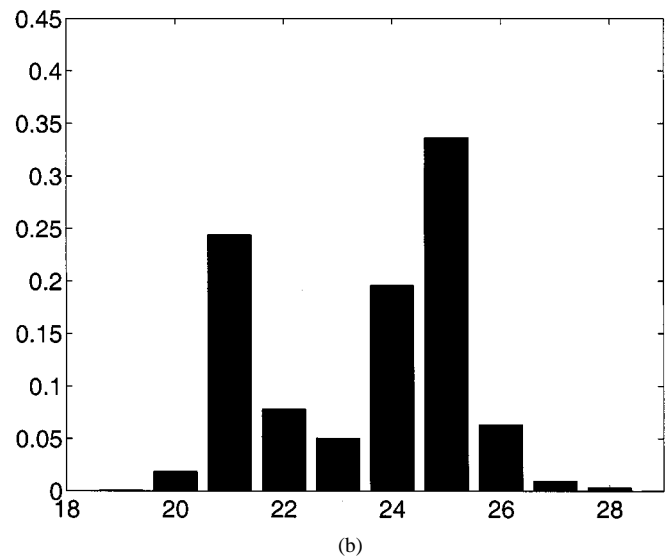
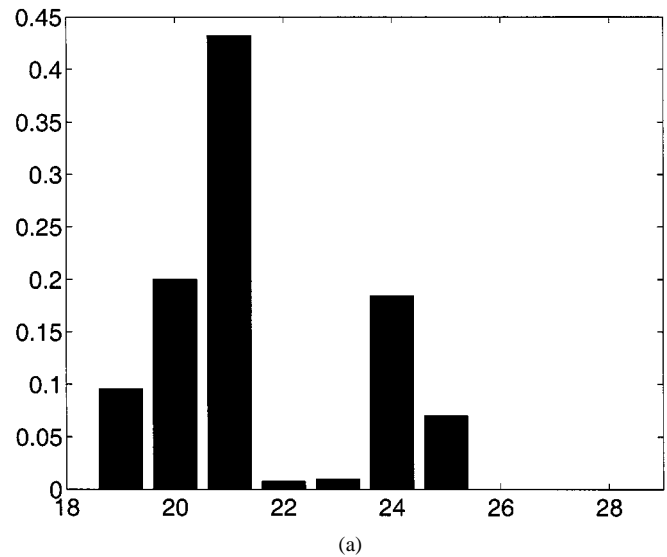
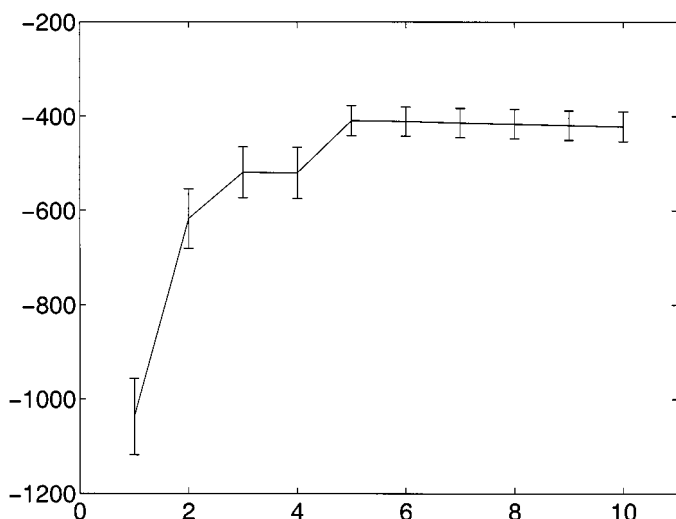


Fig. 2. Model order selection for AR(25) data (SNR = 20) with (a) BIC/MDL and (b) VB with 3-s blocks of data.

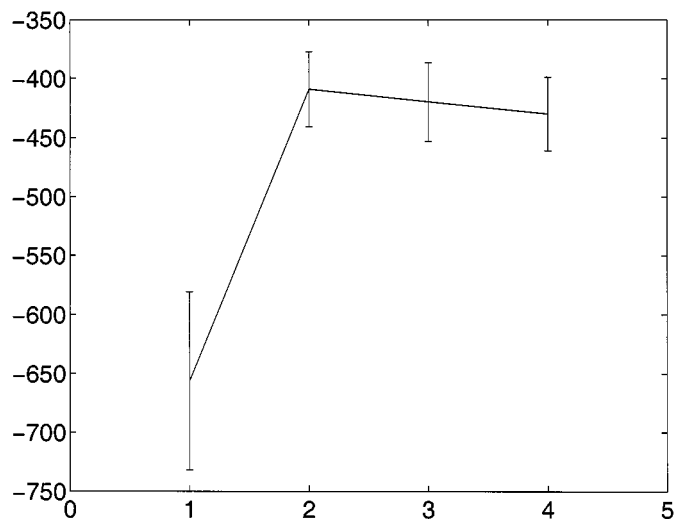
### D. EEG Data: Non-Gaussian Noise Model

We applied the model to a short sample of EEG data shown in Fig. 7, the middle section of which is contaminated by an eye-movement artifact. For a Gaussian noise model ( $m = 1$ ), the optimum AR model order was at  $p = 27$ , having  $F_{27,1} = -2681$ . For  $m = 2$  and  $m = 3$ , the best models were also at  $p = 27$  (this is not always the case) with  $F_{27,2} = -2661$  and  $F_{27,3} = -2671$ . This shows that the non-Gaussian AR process with two noise components is the preferred model. This model took 25 iterations to converge. The final mixing coefficients in the noise model were  $\pi_1 = 0.78$  and  $\pi_2 = 0.22$ , and the variances (inverse precisions) were  $\sigma_1^2 = 0.014$  and  $\sigma_2^2 = 0.066$ , i.e., a low and high variance component. Fig. 7 shows which data points “belong” to which noise component. The first component corresponds to the majority of the data (78% of it), and the second component picks out the outliers that are mostly in the middle section. These points are then effectively downweighted in the corresponding weighted least squares regression. The outliers therefore have a minimal influence on the estimation of the

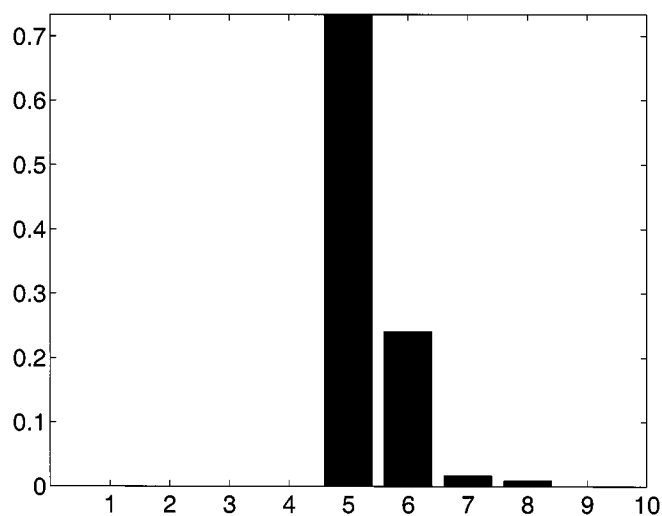




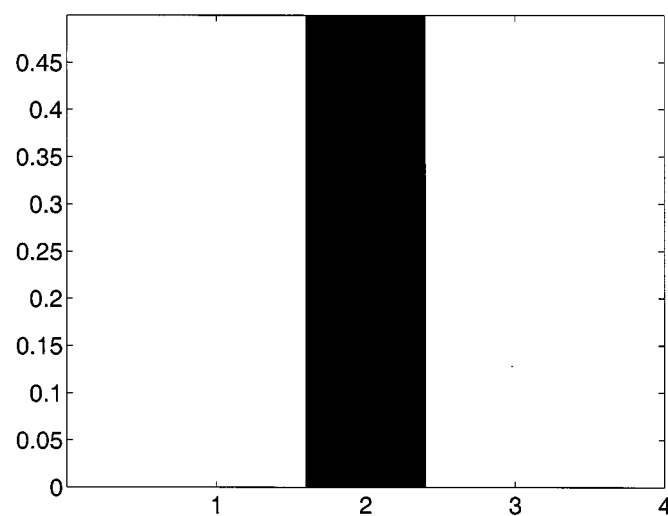
(a)



(a)



(b)



(b)

Fig. 3. AR model order selection (a)  $F(p, 2)$  versus  $p$  and (b) the corresponding probabilities  $P(p, 2)$  versus  $p$ .

Fig. 4. Noise model order selection (a)  $F(5, m)$  versus  $m$  and (b) the corresponding probabilities  $P(5, m)$  versus  $m$ .

AR coefficients; this can be seen from (54), where  $\bar{\beta}_s$  and  $\Gamma_s$  weight the samples accordingly.

#### E. How Good is the Variational Bound?

From (23) we saw that the negative free-energy term  $F(p, m)$  formed a strict lower bound to the true posterior (the KL term is strictly non-negative). How tight this bound is will clearly affect the confidence we have in the variational approach. We consider here a simple experimental validation of the bound and present results on the non-Gaussian noise example in Section VI.

Consider the marginal of (19), where once more, we drop, for notational convenience, dependence on  $p, m$

$$p(\mathbf{Y}) = \iint p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta}) d\mathbf{S} d\boldsymbol{\theta}. \quad (77)$$

We may write (77) as

$$\begin{aligned} p(\mathbf{Y}) &= \iint \left( \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})} \right) q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) d\mathbf{S} d\boldsymbol{\theta} \\ &= \iint w(\boldsymbol{\theta}, \mathbf{S}) q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y}) d\mathbf{S} d\boldsymbol{\theta} \end{aligned} \quad (78)$$

in which we define the weight function

$$w(\boldsymbol{\theta}, \mathbf{S}) \stackrel{\text{def}}{=} \frac{p(\mathbf{Y}, \mathbf{S}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})}. \quad (79)$$

As noticed by others [22, ch. 4], this enables the variational approximation  $q(\boldsymbol{\theta}, \mathbf{S}|\mathbf{Y})$  to be used as the proposal in an *importance sampling* step. We do not cover the details of importance sampling as they may be found in standard texts (e.g., [23, ch. 5]). Suffice to say, the true posterior may be estimated via importance weights, which are evaluated as an expectation under the proposal, i.e., (78) may be written as

$$p(\mathbf{Y}) \approx \mathbb{E}[w(\boldsymbol{\theta}, \mathbf{S})]_q. \quad (80)$$

We applied this procedure to the non-Gaussian AR( $p = 5, m = 2$ ) data analyzed in Section VI. A total of 3s of data (384 samples) was generated from the true model, as before. The VB-AR algorithm was applied to this data, and then, using the variational inference densities as proposals, 1000 samples were evaluated using importance sampling to

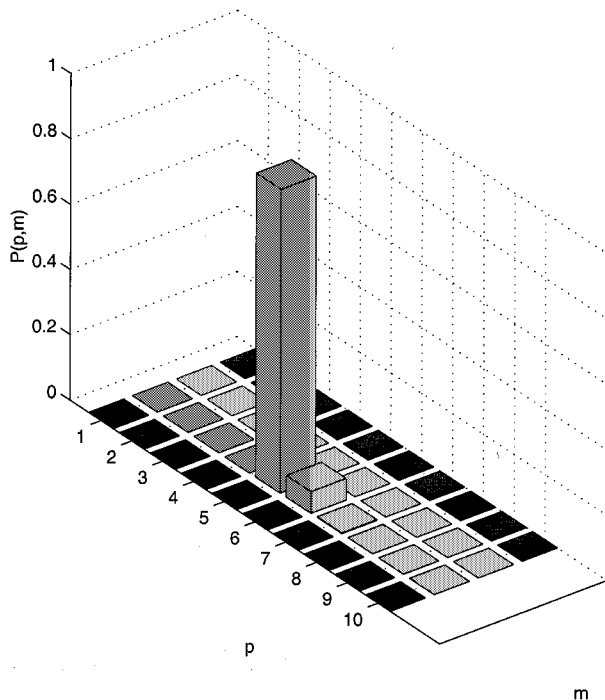


Fig. 5. Noise and AR model order selection. The bar plot shows  $P(p, m)$  against  $p, m$ .

estimate the true model posterior.<sup>2</sup> This was repeated over all model orders from  $p = 1$  to 10 and with the number of Gaussians in the noise model from  $m = 1$  to 5. Fig. 8 shows the relative increase in model evidence after the sampling procedure over the  $p, m$  grid. We note that the relative difference between the “true” and approximate posteriors is of order one part in  $10^5$  at most (for these models). We note that the maximum model evidence appeared at the true model order in both cases. We note that, as expected from theory, the variational approximation is always a strict lower bound. It appears that the variational bound is very tight for low model orders but rises with increased numbers of model parameters, both in terms of the AR model order ( $p$ ) and, importantly, in terms of the number of mixture components in the noise model ( $m$ ). This is to be expected due to the factored nature of the variational approximation, which gets worse with increasing numbers of parameters. It is encouraging that the evidence “peaks” at the correct order in both cases, indicating that the variational free energy may offer a realistic model-selection criterion. More importantly, the use of the variational approximation as a proposal for importance sampling makes the latter very fast, and this approach is to be commended in cases where evaluation of the bound is crucial.

## VII. DISCUSSION

We have proposed a non-Gaussian AR model for the modeling of stationary stochastic processes where the noise is modeled with a mixture of Gaussians. The main appli-

<sup>2</sup>As the dimensionality of the model increases, so does the tendency of the importance weights to be dominated by a small number of large weights. In the experiments reported here, this did not occur, although we note this tendency for high  $p$  models.

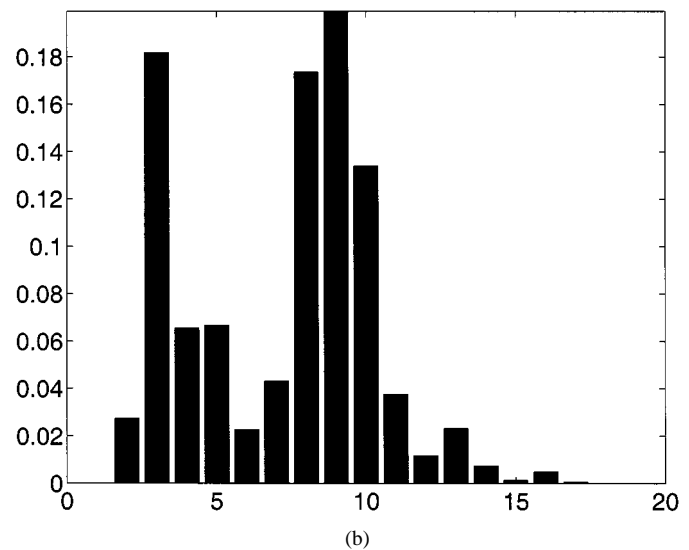
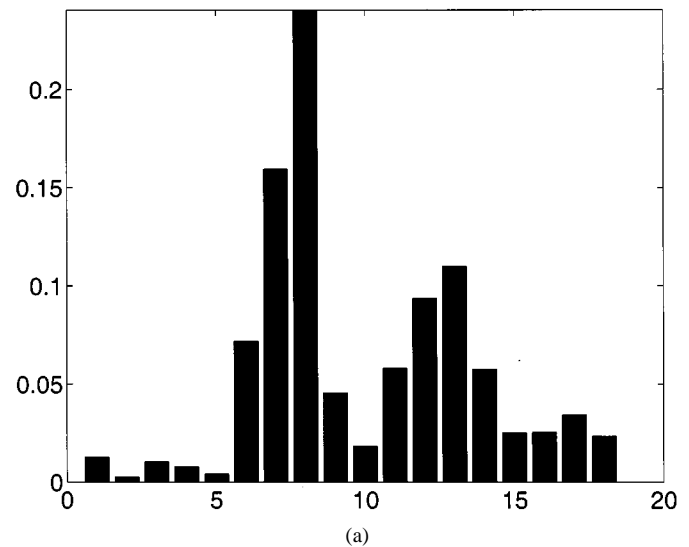


Fig. 6. VB model order selection on EEG data from (a) awake subject and (b) anesthetized subject.

cation is as a model for stationary signals contaminated by non-Gaussian and/or nonstationary noise processes. The VB framework prevents overfitting and provides a model order selection criterion both for the AR order and the noise model order. In the experiments presented, we restricted the means in the noise mixture to be identically zero. The resulting model resembles weighted least squares and, hence, provides robust estimation of AR coefficients.

Our model provides an alternative to cumulant methods for non-Gaussian AR modeling (see, e.g., [24] and [25]) for which model order estimation is problematic.

We envisage that a main application of the model is to EEG data; in an EEG recording, up to 30% [26] of data blocks are corrupted by muscle, eye-movement, or other artifacts, and traditionally, these data blocks are discarded. As we have shown, however, the non-Gaussian AR model is able to downweight the contaminated section of each block, thus allowing more reliable AR (and consequently spectral) estimation. We also envisage using the non-Gaussian AR algorithm to model the sources in an independent component analysis model; this would extend the work in [10] by allowing for dynamic sources at the same

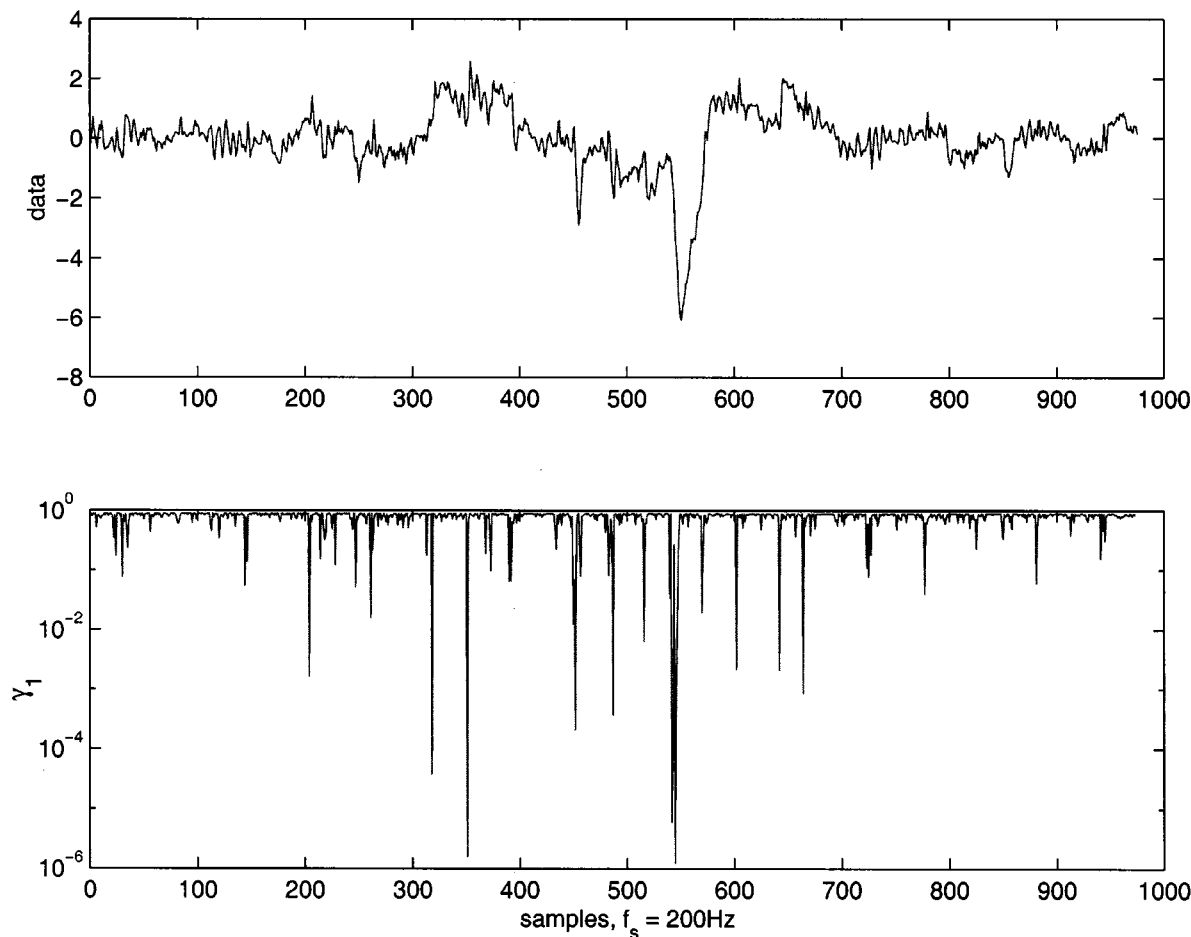


Fig. 7. Top plot: Original EEG trace. Bottom plot: Log-probability that the data belongs to the first (low-variance) noise component  $\log \gamma_1$  [calculated from (44)]. A high proportion of data in the middle section therefore clearly belongs to the second (high-variance) noise component. These points are downweighted in the estimation of the AR coefficients.

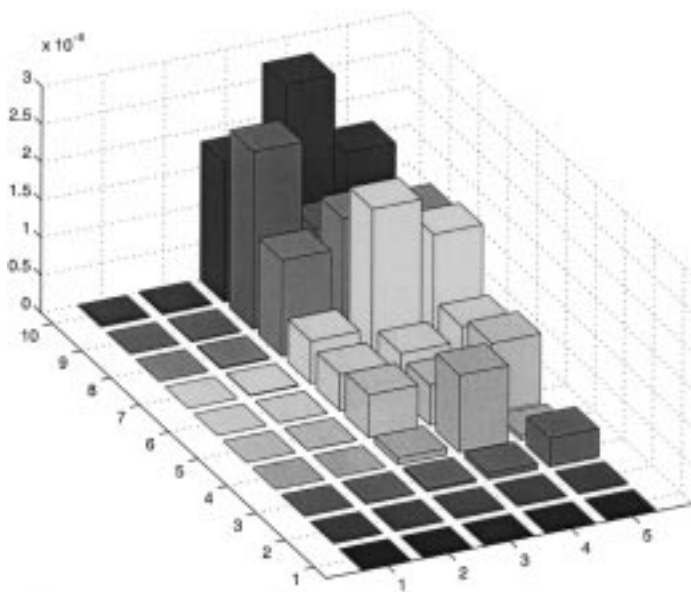


Fig. 8. Relative increase in evidence versus model order after importance sampling using the variational posterior as proposal.

time as using VB for model order estimation. It is noted that the model we use, in which the noise model has no dynamic struc-

ture, is equivalent to a 0th-order hidden Markov model (HMM) for the noise process with a number of hidden states equal to the number of components in the MoG. In the eye-movement artifact example, the high variance activity occurred sporadically, and the HMM(0) model is appropriate. It is arguable, however, that in some cases, where the high-variance noise occurs in temporally correlated bursts, that a higher order HMM is appropriate [a HMM(1), for example]. We have not considered this here, although it is an important area of future research.

For Gaussian AR models, our experiments show that VB model order selection is superior to the BIC/MDL criterion. We have shown that for noninformative priors on the precision of the coefficients  $\alpha$  and the precision of the noise  $\beta$ , the VB update rules are identical to those of the evidence framework. This is an extension of work by Mackay [21], who showed that for a linear regression model with known  $\beta$ , the updates for  $\alpha$  are the same. This link therefore provides further justification for the updating rules used in the evidence framework (as the VB update rules probably increase a lower bound on the log-evidence, whereas no such convergence proof previously existed in the evidence framework). We have also shown that apart from minor differences that have no practical impact, the VB and evidence model order selection criteria are identical (again, under the assumption of noninformative priors on  $\alpha$  and  $\beta$ ).

APPENDIX A  
DERIVATION OF VARIATIONAL LEARNING

We start by restating (27), which conjectures that the optimal form of the proposal ( $q$ ) distribution that maximizes the negative free energy (model evidence measure), with respect to parameter group  $\theta_i$ , is

$$q(\theta_i|\mathbf{Y}) = \frac{1}{Z_i} \exp[I(\theta_i)] \quad (81)$$

where  $Z_i$  is the partition function (normalizing constant), and where, as before, we define

$$I(\theta_i) \stackrel{\text{def}}{=} \int q(\theta^{\setminus i}|\mathbf{Y}) \log p(\mathbf{Y}, \mathbf{S}|\theta)p(\theta) d\theta^{\setminus i}. \quad (82)$$

This maximization is equivalent to stating that the negative KL-divergence between  $q(\theta_i|\mathbf{Y})$  and  $\exp[I(\theta_i)]/Z_i$  is maximized [this is strictly nonpositive and reaches its maximum when (81) holds]. The negative KL-divergence is given as

$$-KL(q(\theta_i|\mathbf{Y}), \exp[I(\theta_i)]/Z_i) = \int d\theta_i q(\theta_i|\mathbf{Y}) \cdot (\log \exp[I(\theta_i)] - \log q(\theta_i|\mathbf{Y}) - \log Z_i). \quad (83)$$

Substituting from (82), we obtain

$$\begin{aligned} & -KL(q(\theta_i|\mathbf{Y}), \exp[I(\theta_i)]/Z_i) \\ &= \int d\theta_i q(\theta_i|\mathbf{Y}) \log q(\theta_i|\mathbf{Y})^{-1} \int d\theta^{\setminus i} q(\theta^{\setminus i}|\mathbf{Y}) \\ & \quad \cdot \log[p(\mathbf{Y}, \mathbf{S}|\theta)p(\theta)] - \log Z_i \\ &= \iint d\theta_i d\theta^{\setminus i} q(\theta_i|\mathbf{Y})q(\theta^{\setminus i}|\mathbf{Y}) \\ & \quad \cdot \log \left[ \frac{p(\mathbf{Y}, \mathbf{S}|\theta)p(\theta)}{q(\theta_i|\mathbf{Y})} \right] - \log Z_i \\ &= \int d\theta q(\theta|\mathbf{Y}) \left[ \frac{p(\mathbf{Y}, \mathbf{S}|\theta)p(\theta)}{q(\theta_i|\mathbf{Y})} \right] - \log Z_i \end{aligned} \quad (84)$$

where the last step follows from the factored representation of the proposal, i.e., (26). Note that this is (to within an additive constant) the contribution to the free energy of (24) of the parameter group  $\theta_i$ . As the free energy is strictly non-negative, it is maximized by maximization of each contribution, e.g., maximization of (84). This is achieved when (81) holds, and hence, we obtain a simple methodology to obtain densities for each component of  $q$ .

APPENDIX B  
DENSITIES AND DIVERGENCES

For completeness, we include definitions of the normal, gamma, and Dirichlet densities and their entropies and KL-divergencies.

A. Univariate Normal Density

For univariate normal densities  $q(x) = N_1(x; \mu_q, \sigma_q^2)$  and  $p(x) = N_1(x; \mu_p, \sigma_p^2)$ , the KL-divergence is

$$KL(q||p) = \frac{1}{2} \log \frac{\sigma_p^2}{\sigma_q^2} + \frac{\mu_q^2 + \mu_p^2 + \sigma_q^2 - 2\mu_q\mu_p}{2\sigma_p^2} - \frac{1}{2}. \quad (85)$$

B. Multivariate Normal Density

The multivariate normal density is given by

$$N(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \cdot \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (86)$$

The KL-divergence for normal densities  $q(\mathbf{x}) = N(\mathbf{x}; \boldsymbol{\mu}_q, \mathbf{C}_q)$  and  $p(\mathbf{x}) = N(\mathbf{x}; \boldsymbol{\mu}_p, \mathbf{C}_p)$  is

$$\begin{aligned} KL(q||p) &= \frac{1}{2} \log \frac{|\mathbf{C}_p|}{|\mathbf{C}_q|} + \frac{1}{2} \text{Tr}(\mathbf{C}_p^{-1} \mathbf{C}_q) \\ & \quad + \frac{1}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^T \mathbf{C}_p^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) - \frac{d}{2} \end{aligned} \quad (87)$$

where  $|\mathbf{C}_p|$  denotes the determinant of the matrix  $\mathbf{C}_p$ .

C. Gamma Density

The Gamma density is defined as

$$\text{Ga}(x; b, c) = \frac{1}{\Gamma(c)} \frac{x^{c-1}}{b^c} \exp\left(-\frac{x}{b}\right). \quad (88)$$

For Gamma densities  $q(x) = \text{Ga}(x; b_q, c_q)$  and  $p(x) = \text{Ga}(x; b_p, c_p)$ , the KL-divergence is

$$\begin{aligned} KL(q||p) &= (c_q - 1)\Psi(c_q) - \log b_q - c_q - \log \Gamma(c_q) + \log \Gamma(c_p) \\ & \quad + c_p \log b_p - (c_p - 1)(\Psi(c_q) + \log b_q) + \frac{b_q c_q}{b_p} \end{aligned} \quad (89)$$

where  $\Psi(\cdot)$  is, as before, the digamma function [2].

D. Dirichlet Density

The Dirichlet density is given by

$$\text{Dir}(\boldsymbol{\pi}; \boldsymbol{\lambda}) = \frac{\Gamma\left(\sum_{s=1}^m \lambda_s\right)}{\prod_{s=1}^m \Gamma(\lambda_s)} \prod_{s=1}^m \pi_s^{\lambda_s - 1} \quad (90)$$

where  $\lambda_s$  is the  $s$ th element of  $\boldsymbol{\lambda}$ , and  $\Gamma(x)$  is the Gamma function [2].

For the  $m$ -state Dirichlet density with  $q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\lambda}_q)$ ,  $p(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi}; \boldsymbol{\lambda}_p)$  and  $\lambda_{qt} = \sum_{s=1}^m \lambda_q(s)$ ,  $\lambda_{pt} = \sum_{s=1}^m \lambda_p(s)$ , the KL-divergence is

$$\begin{aligned} KL(q||p) &= \log \frac{\Gamma(\lambda_{qt})}{\Gamma(\lambda_{pt})} + \sum_{s=1}^m (\lambda_q(s) - \lambda_p(s)) \\ & \quad \cdot (\Psi(\lambda_q(s)) - \Psi(\lambda_{qt})) + \sum_{s=1}^m \log \frac{\Gamma(\lambda_p(s))}{\Gamma(\lambda_q(s))}. \end{aligned} \quad (91)$$

ACKNOWLEDGMENT

The authors would like to thank R. Daniel and P. Sykacek for interesting discussions and helpful comments and criticism of this work. The comments of the anonymous referees were most valuable in the final revision of this paper.

## REFERENCES

- [1] C. M. Bishop, *Neural Networks for Pattern Recognition*. Oxford, U.K.: Oxford Univ. Press, 1995.
- [2] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. V. P. Flannery, *Numerical Recipes in C*. Cambridge, U.K.: Cambridge Univ. Press, 1992.
- [3] G. Barnett, R. Kohn, and S. Sheather, "Bayesian estimation of an autoregressive model using Markov chain Monte Carlo," *J. Econometr.*, vol. 74, no. 2, pp. 237–254, 1996.
- [4] S. Godshill, "Bayesian enhancement of speech and audio signals which can be modeled as ARMA processes," *Statist. Rev.*, vol. 65, no. 1, pp. 1–21, 1996.
- [5] S. Godshill and P. Rayner, "Statistical reconstruction and analysis of autoregressive signals in impulsive noise," *IEEE Trans. Speech Audio Processing*, vol. 6, pp. 352–372, July 1998.
- [6] P. Troughton and S. Godshill, "A reversible jump sampler for autoregressive time series, employing full conditionals to achieve efficient model space moves," Dept. Eng., Univ. Cambridge, Cambridge, U.K., Tech. Rep. CUED/F-INFENG/TR.304, 1997.
- [7] S. Richardson and P. J. Green, "On Bayesian analysis of mixtures with an unknown number of components," *J. R. Statist. Soc. B*, vol. 59, no. 4, pp. 731–758, 1997.
- [8] H. Lappalainen and J. W. Miskin, "Ensemble learning," in *Advances in Independent Components Analysis*, M. Girolami, Ed. Boston, MA: Kluwer, 2000.
- [9] C. M. Bishop, "Variational principal components," in *Proc. Int. Conf. Artif. Neural Networks*, 1999, pp. 509–514.
- [10] R. Choudrey, W. D. Penny, and S. J. Roberts, "An ensemble learning approach to independent component analysis," in *Proc. IEEE Int. Workshop Neural Networks Signal Process.*, Sydney, Australia, 2000.
- [11] W. D. Penny and S. J. Roberts, "Bayesian methods for autoregressive models," in *Proc. IEEE Int. Workshop Neural Networks Signal Process.*, Sydney, Australia, 2000.
- [12] —, "Variational Bayes for non-Gaussian autoregressive models," in *Proc. IEEE Int. Workshop Neural Networks Signal Process.*, Sydney, Australia, 2000.
- [13] G. Kitagawa and W. Gersch, "A smoothness priors time-varying ar-coefficient modeling of nonstationary covariance time series," *IEEE Trans. Automat. Contr.*, vol. AC-30, pp. 45–56, 1985.
- [14] —, "A smoothness priors long ar model method for spectral estimation," *IEEE Trans. Automat. Contr.*, vol. AC-30, pp. 57–65, 1985.
- [15] S. Weisberg, *Applied Linear Regression*. New York: Wiley, 1980.
- [16] J. J. K. O'Ruanaidh and W. J. Fitzgerald, *Numerical Bayesian Methods Applied to Signal Processing*. New York: Springer, 1996.
- [17] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [18] H. Attias *et al.*, "A variational Bayesian framework for graphical models," in *NIPS 12*, T. Leen, *et al.*, Eds. Cambridge, MA: MIT Press, 2000.
- [19] D. M. Chickering and D. Heckerman, "Efficient approximations for the marginal likelihood of Bayesian networks with hidden variables," Microsoft Research, Tech. Rep. MSR-TR-96-08, 1996.
- [20] P. Sykacek, I. Rezek, and S. Roberts. (2000) Markov chain Monte Carlo methods for Bayesian sensor fusion. Dept. Eng. Sci., Univ. Oxford, Oxford, U.K., Tech. Rep. PARG-00-10. [Online]. Available: <http://www.robots.ox.ac.uk/~parg/>
- [21] D. J. C. Mackay, "Ensemble learning and evidence maximization," Tech. Rep., Cavendish Lab., Univ. Cambridge, Cambridge, U.K., 1995.
- [22] J. Miskin, "Ensemble learning for independent component analysis," Ph.D. dissertation, Cavendish Lab., Univ. Cambridge, Cambridge, U.K., 2000.
- [23] J. M. Bernardo and A. F. M. Smith, *Bayesian Theory*. New York: Wiley, 1994.
- [24] S. Minfen, S. Lisha, and P. J. Beadle, "Parametric bispectral estimation of EEG signals in different functional states of brain," in *Advances in Medical Signal and Information Processing*. London, U.K.: IEE, 2000, pp. 66–72.
- [25] C. L. Nikias, *Higher-Order Spectra Analysis: A Nonlinear Signal Processing Framework*. Englewood Cliffs, NJ: Prentice-Hall, 1993.
- [26] A. Gevins, M. E. Smith, and D. Yu, "High-resolution EEG mapping of cortical activation related to working memory: Effects of task difficulty, type of processing, and practice," *Cerebral Cortex*, vol. 7, no. 4, pp. 374–385, 1997.



**Stephen J. Roberts** was born in London, U.K., in 1965. He received the Ph.D. degree in physics in 1991 from Oxford University, Oxford, U.K.

He studied physics at Oxford from 1983 to 1986 and then worked for a few years in the Research Department of Oxford Instruments. From 1989 to 1991 he researched his Ph.D. and worked as a postdoctoral researcher until 1994, when he was appointed to the Faculty at Imperial College, University of London. In 1999, he joined Oxford. His main area of research lies in machine learning approaches to data analysis. He

has particular interests in the application of machine learning methods to problems in mathematical biology. Current research focuses on Bayesian statistics, graphical models, independent component analysis, and information theory. He runs the Pattern Analysis and Machine Learning Research Group at Oxford.

Dr. Roberts is a Fellow of Somerville College, Oxford.



**Will D. Penny** received the Ph.D. degree in electrical engineering from Imperial College, University of London, London, U.K., in 1993.

He has since worked as a post-doctoral researcher at University College, London; Imperial College; and Oxford University, Oxford, U.K., and is now working as a brain-imaging statistician at University College. His research interests are in Bayesian statistics and data-driven machine learning and their application to biomedical data.