



**HAL**  
open science

# Variational Bayesian Approximation methods for inverse problems

Ali Mohammad-Djafari

► **To cite this version:**

Ali Mohammad-Djafari. Variational Bayesian Approximation methods for inverse problems. MaxEnt 2012, Jul 2012, Garching, Germany. pp.230. hal-00833298

**HAL Id: hal-00833298**

**<https://hal.science/hal-00833298>**

Submitted on 12 Jun 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Bayesian Inference Tools for Inverse Problems

Ali Mohammad-Djafari

*Laboratoire des Signaux et Systèmes,  
UMR 8506 CNRS-SUPELEC-UNIV PARIS SUD  
SUPELEC, Plateau de Moulon, 3 rue Joliot-Curie, 91192 Gif-sur-Yvette, France*

**Abstract.** In this paper, first the basics of the Bayesian inference for linear inverse problems are presented. The inverse problems we consider are, for example, signal deconvolution, image restoration or image reconstruction in Computed Tomography (CT). The main point to discuss then is the prior modeling of signals and images. We consider two classes of priors: *simple* or *hierarchical with hidden variables*. For practical applications, we need also to consider the estimation of the hyper parameters. Finally, we see that we have to infer simultaneously the unknowns, the hidden variables and the hyper parameters.

Very often, the expression of the joint posterior law of all the unknowns is too complex to be handled directly. Indeed, rarely we can obtain analytical solutions to any point estimators such the Maximum A posteriori (MAP) or Posterior Mean (PM). Three main tools can then be used: Laplace approximation (LAP), Markov Chain Monte Carlo (MCMC) and Bayesian Variational Approximations (BVA).

To illustrate all these aspects, we will consider a deconvolution problem where we know that the input signal is *sparse* and propose to use a Student-t distribution for that. Then, to handle the Bayesian computations with this model, we use the property of Student-t which is modelling it via an infinite mixture of Gaussians, introducing thus hidden variables which are the variances. Then, the expression of the joint posterior of the input signal samples, the hidden variables (which are here the inverse variances of those samples) and the hyper-parameters of the problem (for example the variance of the noise) is given. From this point, we will present the joint maximization by alternate optimization and the three possible approximation methods. Finally, the proposed methodology is applied in different applications such as mass spectrometry, spectrum estimation of quasi periodic biological signals and X ray computed tomography.

## INTRODUCTION

In many generic inverse problems in signal and image processing, the problem is to infer on an unknown signal  $f(t)$  or an unknown image  $f(r)$  with  $r = (x, y)$  through an observed signal  $g(t')$  or an observed image  $g(r')$  related between them through an operator  $\mathcal{H}$  such as convolution  $g = h * f$  or any other linear or nonlinear transformation  $g = \mathcal{H}f$ . When this relation is linear and we have discretized the problem, we arrive to the relation:  $\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\varepsilon}$  where  $\mathbf{f} = [f_1, \dots, f_n]'$  represents the unknowns,  $\mathbf{g} = [g_1, \dots, g_m]'$  the observed data,  $\boldsymbol{\varepsilon} = [\varepsilon_1, \dots, \varepsilon_m]'$  the errors of modelling and measurement and  $\mathbf{H}$  the matrix of the system response.

The Bayesian inference approach is based on the posterior law:

$$p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}) = \frac{p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f}|\boldsymbol{\theta})}{p(\mathbf{g}|\boldsymbol{\theta})} \propto p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta})p(\mathbf{f}|\boldsymbol{\theta}) \quad (1)$$

where the sign  $\propto$  stands for “proportional to”,  $p(g|f, \theta)$  is the likelihood,  $p(f|\theta)$  the prior probability law model and  $p(g|\theta)$  is called the evidence of the model.  $\theta$  is called the vector of the hyper-parameters of the problem. It can be divided in two independent parts  $\theta = (\theta_1, \theta_2)$  where  $\theta_1$  is only appearing in the likelihood term and  $\theta_2$  in the prior term. Then, we can write  $p(f|g, \theta) \propto p(g|f, \theta_1) p(f|\theta_2)$ .

When the parameters  $\theta$  have to be estimated too, a prior probability law  $p(\theta|\theta_0)$  with fixed values for  $\theta_0$  is assigned to them in such a way to obtain the joint posterior law:

$$p(f, \theta|g, \theta_0) = \frac{p(g|f, \theta) p(f|\theta) p(\theta|\theta_0)}{p(g|\theta_0)} \quad (2)$$

which can be used to infer them jointly. From this point (omitting  $\theta_0$  for simplicity of notations), at least three main approaches are commonly used to infer both unknowns  $f$  and  $\theta$ :

- Joint MAP estimation:

$$(\hat{f}, \hat{\theta}) = \arg \max_{(f, \theta)} \{p(f, \theta|g)\} \quad (3)$$

which is done in general by an alternate optimization algorithm such as:

$$\begin{cases} \hat{f} = \arg \max_f \{p(f, \hat{\theta}|g)\} \\ \hat{\theta} = \arg \max_{\theta} \{p(\hat{f}, \theta|g)\} \end{cases} \quad (4)$$

- Marginalization over  $f$  to obtain:

$$p(\theta|g) = \int p(f, \theta|g) \, df \quad (5)$$

which can be used to infer  $\theta$  which is then used for the estimation of  $f$ . Unfortunately, in general, the marginalization step can not be done analytically. The Expectation-Maximization (EM) algorithm is used to obtain  $\hat{\theta}$ . The EM algorithm can be summarized as follows:

$$\begin{cases} Q(\theta, \hat{\theta}) = \int p(f|\hat{\theta}, g) \ln p(f, \theta|g) \, df = \langle \ln p(f, \theta|g) \rangle_{p(f|\hat{\theta}, g)} \\ \hat{\theta} = \arg \max_{\theta} \{Q(\theta, \hat{\theta})\} \end{cases} \quad (6)$$

- Variational Bayesian Approximation (VBA):

The main idea here is to approximate  $p(f, \theta|g)$  by a separable one  $q(f, \theta) = q_1(f) q_2(\theta)$  which can then be used for inferring on  $f$  or  $\theta$  [3, 1, 13, 18, 2, 15, 17, 12, 9].

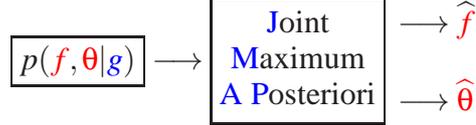
As we will see more in detail in the next section,  $q_1(f)$  and  $q_2(\theta)$  are obtained via the following iterative algorithms:

$$\begin{cases} q_1(f) \propto \exp \left\{ - \langle \ln p(f, \theta, g) \rangle_{q_2(\theta)} \right\} \\ q_2(\theta) \propto \exp \left\{ - \langle \ln p(f, \theta, g) \rangle_{q_1(f)} \right\} \end{cases} \quad (7)$$

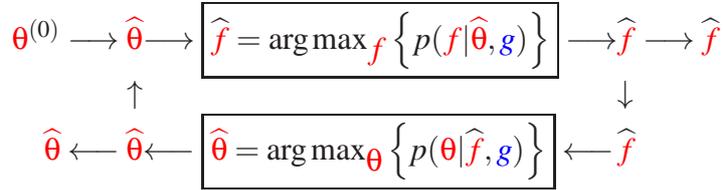
The two first approaches are very well known. The last one is less and is explained in more details in the next section. These approaches are compared in Figure 1.

**Joint MAP:**

Main idea:

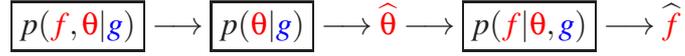


Algorithm:

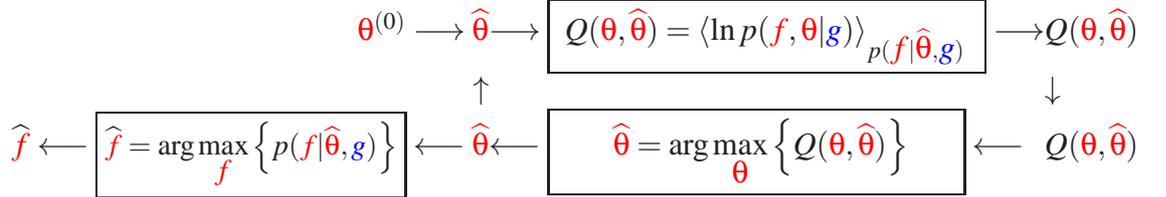


**Marginalization:**

Main idea:

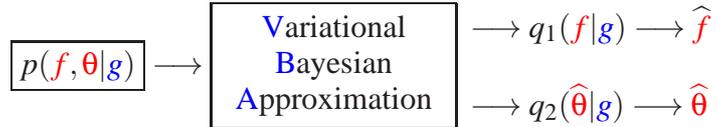


Algorithm:



**VBA:**

Main idea:



Algorithm:

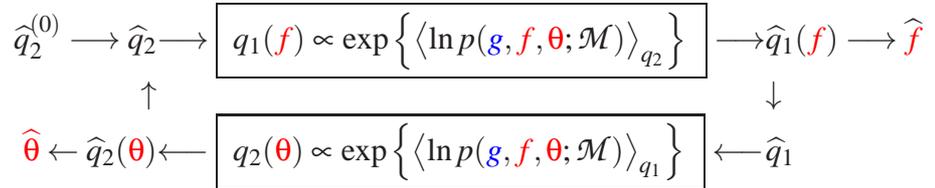


Figure 1: Three different approaches for inferring both unknowns  $f$  and  $\theta$ .

As we will see, the main inconvenient of the first approach is that we are summarizing the joint posterior law  $p(f, \theta|g)$  by only its mode. Also, for obtaining this mode, in general an iterative alternate optimization is used, where at each iteration, only the

values of the estimates at previous iterations are used without accounting for their corresponding uncertainties. In the second approach, first  $\theta$  is estimated and then it is used for the estimation of  $f$ , again without accounting for its uncertainty. In the third approach, as we will see, the estimation of  $f$  depends on the approximated law  $q_2(\theta)$  and the estimation of  $\theta$  depends on the approximated law  $q_1(f)$ , thus accounting for uncertainties in both steps.

A simple prior probability law is often not enough for modeling signals and images, in particular for non stationary signals or non homogeneous images. We then may use hierarchical models with hidden variables  $z$  which may represent, for example, the class labels in mixture models. In those cases, the prior probability model contains two parts  $p(f|z, \theta_2)$  and  $p(z|\theta_3)$  and we will have:

$$p(f, z, \theta|g, \theta_0) \propto p(g|f, \theta_1) p(f|z, \theta_2) p(z|\theta_3) p(\theta|\theta_0) \quad (8)$$

and then again different approaches can be used to infer the unknowns  $f$ ,  $z$  and  $\theta$ .

In this paper, first the general VBA method is detailed for the inference on inverse problems with hierarchical prior models. Then, two particular classes of prior models (Student-t and mixture of Gaussians) are considered and the details of BVA algorithms are given for them.

## BAYESIAN VARIATIONAL APPROXIMATION WITH HIERARCHICAL PRIOR MODELS

When a hierarchical prior model  $p(f|z, \theta)$  is used and when the estimation of the hyper-parameters  $\theta$  has to be considered, the joint posterior law of all the unknowns becomes:

$$p(f, z, \theta|g) \propto p(g|f, \theta_1) p(f|z, \theta_2) p(z|\theta_3) p(\theta) \quad (9)$$

which can also be written as  $p(f, z, \theta|g) = p(f|z, \theta, g) p(z|\theta, g) p(\theta|g)$  where

$$p(f|z, \theta, g) = p(g|f, \theta) p(f|z, \theta) / p(g|z, \theta) \text{ with } p(g|z, \theta) = \int p(g|f, \theta) p(f|z, \theta) df \quad (10)$$

and

$$p(z|\theta, g) = \frac{p(g|z, \theta) p(z|\theta)}{p(g|\theta)} \text{ with } p(g|\theta) = \int p(g|z, \theta) p(z|\theta) dz \quad (11)$$

and finally

$$p(\theta|g) = \frac{p(g|\theta) p(\theta)}{p(g)} \text{ with } p(g) = \int p(g|\theta) p(\theta) d\theta. \quad (12)$$

In general, common choices for  $p(g|f, \theta_1)$  and  $p(f|z, \theta_2)$  are Gaussians and for  $p(z|\theta_3)$  and  $p(\theta)$  are Bernoulli or Binomial (for discrete valued  $z$ ) or Gamma for inverse of the variances. Thus, the first term

$$p(f|z, \theta, g) \propto p(g|f, \theta) p(f|z, \theta) \quad (13)$$

will be easy to handle because it is the product of two Gaussians and so it is a multivariate Gaussian. But the two others are not.

The main idea behind the VBA is to approximate the joint posterior  $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g})$  by a separable one, for example

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f}) q_2(\mathbf{z}) q_3(\boldsymbol{\theta}) \quad (14)$$

and where the expressions of  $q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g})$  is obtained by minimizing the Kullback-Leibler divergence

$$\text{KL}(q : p) = \int q \ln \frac{q}{p} = \left\langle \ln \frac{q}{p} \right\rangle_q. \quad (15)$$

It is then easy to show that  $\text{KL}(q : p) = \ln p(\mathbf{g} | \mathcal{M}) - \mathcal{F}(q)$  where  $p(\mathbf{g} | \mathcal{M})$  is the likelihood of the model

$$p(\mathbf{g} | \mathcal{M}) = \int \int \int p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) d\mathbf{f} d\mathbf{z} d\boldsymbol{\theta} \quad (16)$$

with  $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) = p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z} | \boldsymbol{\theta}) p(\boldsymbol{\theta})$  and  $\mathcal{F}(q)$  is the free energy associated to  $q$  defined as

$$\mathcal{F}(q) = \left\langle \ln \frac{p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta})} \right\rangle_q. \quad (17)$$

So, for a given model  $\mathcal{M}$ , minimizing  $\text{KL}(q : p)$  is equivalent to maximizing  $\mathcal{F}(q)$  and when optimized,  $\mathcal{F}(q^*)$  gives a lower bound for  $\ln p(\mathbf{g} | \mathcal{M})$ .

Without any other constraint than the normalization of  $q$ , an alternate optimization of  $\mathcal{F}(q)$  with respect to  $q_1$ ,  $q_2$  and  $q_3$  results in

$$\begin{cases} q_1(\mathbf{f}) \propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{z})q(\boldsymbol{\theta})} \right\} \\ q_2(\mathbf{z}) \propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{f})q(\boldsymbol{\theta})} \right\} \\ q_3(\boldsymbol{\theta}) \propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{f})q(\mathbf{z})} \right\} \end{cases} \quad (18)$$

Note that these relations represent an implicit solution for  $q_1(\mathbf{f})$ ,  $q_2(\mathbf{z})$  and  $q_3(\boldsymbol{\theta})$  which need, at each iteration, the expression of the expectations in the right hand of exponentials. If  $p(\mathbf{g} | \mathbf{f}, \mathbf{z}, \boldsymbol{\theta}_1)$  is a member of an exponential family and if all the priors  $p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}_2)$ ,  $p(\mathbf{z} | \boldsymbol{\theta}_3)$ ,  $p(\boldsymbol{\theta}_1)$ ,  $p(\boldsymbol{\theta}_2)$ , and  $p(\boldsymbol{\theta}_3)$  are conjugate priors, then it is easy to see that these expressions leads to standard distributions for which the required expectations are easily evaluated. In that case, we may note

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) q_2(\mathbf{z} | \tilde{\mathbf{f}}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) q_3(\boldsymbol{\theta} | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \mathbf{g}) \quad (19)$$

where the tilded quantities  $\tilde{\mathbf{z}}$ ,  $\tilde{\mathbf{f}}$  and  $\tilde{\boldsymbol{\theta}}$  are, respectively functions of  $(\tilde{\mathbf{f}}, \tilde{\boldsymbol{\theta}})$ ,  $(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}})$  and  $(\tilde{\mathbf{f}}, \tilde{\mathbf{z}})$  and where the alternate optimization with respect to  $q_1$ ,  $q_2$  and  $q_3$  becomes alternate updating of the parameters  $(\tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}})$  for  $q_1$ , the parameters  $(\tilde{\mathbf{f}}, \tilde{\boldsymbol{\theta}})$  of  $q_2$  and the parameters  $(\tilde{\mathbf{f}}, \tilde{\mathbf{z}})$  of  $q_3$ .

Finally, we may note that, to monitor the convergence of the algorithm, we may evaluate the free energy

$$\begin{aligned}\mathcal{F}(q) &= \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) \rangle_q + \langle -\ln q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q \\ &= \langle \ln p(\mathbf{g} | \mathbf{f}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{z} | \boldsymbol{\theta}) \rangle_q \\ &\quad + \langle -\ln q(\mathbf{f}) \rangle_q + \langle -\ln q(\mathbf{z}) \rangle_q + \langle -\ln q(\boldsymbol{\theta}) \rangle_q\end{aligned}\quad (20)$$

where all the expectations are with respect to  $q$ .

Other decompositions are also possible:

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = \prod_j q_{1j}(\mathbf{f}_j | \tilde{\mathbf{f}}_{(-j)}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) \prod_j q_{2j}(\mathbf{z}_j | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}_{(-j)}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) \prod_l q_{3l}(\boldsymbol{\theta}_l | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}_{(-l)}, \mathbf{g})\quad (21)$$

or

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) \prod_j q_{2j}(\mathbf{z}_j | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}_{(-j)}, \tilde{\boldsymbol{\theta}}, \mathbf{g}) \prod_l q_{3l}(\boldsymbol{\theta}_l | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}_{(-l)}, \mathbf{g})\quad (22)$$

In the following section, we consider this case and give some more details with the hierarchical model of Infinite Mixture model of Student-t which is used for example for modeling the distributions of sparse signals or images [14].

## JMAP AND BAYESIAN VARIATIONAL APPROXIMATION WITH STUDENT-T PRIORS

The Student-t model is:

$$p(\mathbf{f} | \mathbf{v}) = \prod_j St(\mathbf{f}_j | \mathbf{v}) \text{ with } St(\mathbf{f}_j | \mathbf{v}) = \frac{1}{\sqrt{\pi \mathbf{v}}} \frac{\Gamma((\mathbf{v} + 1)/2)}{\Gamma(\mathbf{v}/2)} (1 + \mathbf{f}_j^2 / \mathbf{v})^{-(\mathbf{v} + 1)/2}\quad (23)$$

Knowing that

$$St(\mathbf{f}_j | \mathbf{v}) = \int_0^\infty \mathcal{N}(\mathbf{f}_j | 0, 1/z_j) \mathcal{G}(z_j | \mathbf{v}/2, \mathbf{v}/2) dz_j\quad (24)$$

we can write this model via the positive hidden variables  $z_j$ :

$$\begin{cases} p(\mathbf{f} | \mathbf{z}) &= \prod_j p(\mathbf{f}_j | z_j) = \prod_j \mathcal{N}(\mathbf{f}_j | 0, 1/z_j) \propto \exp\left\{-\frac{1}{2} \sum_j z_j \mathbf{f}_j^2\right\} \\ p(\mathbf{z}_j | \alpha, \beta) &= \mathcal{G}(z_j | \alpha, \beta) \propto z_j^{(\alpha-1)} \exp\{-\beta z_j\} \text{ with } \alpha = \beta = \mathbf{v}/2\end{cases}\quad (25)$$

The Cauchy model is obtained when  $\mathbf{v} = 1$ .

Now consider this prior model for the unknowns  $\mathbf{f}$  of a linear inverse problem with the linear forward model  $\mathbf{g} = H\mathbf{f} + \boldsymbol{\varepsilon}$  and assign a Gaussian law to the noise  $\boldsymbol{\varepsilon}$  which results to  $p(\mathbf{g} | \mathbf{f}, \mathbf{v}_\boldsymbol{\varepsilon}) = \mathcal{N}(\mathbf{g} | H\mathbf{f}, \mathbf{v}_\boldsymbol{\varepsilon} I)$ . We also assign a prior  $p(\boldsymbol{\tau}_\boldsymbol{\varepsilon} | \alpha_0, \beta_0) = \mathcal{G}(\boldsymbol{\tau}_\boldsymbol{\varepsilon} | \alpha_0, \beta_0)$  to  $\boldsymbol{\tau}_\boldsymbol{\varepsilon} = 1/\mathbf{v}_\boldsymbol{\varepsilon}$ . Figure 2 shows the graphical representation of this model.

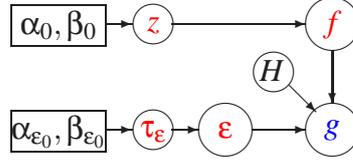


Figure 2: The graphical representation of the proposed model with the Student-t equivalent hierarchical prior.

In the following, we summarize all the equations related to this modeling and inference scheme.

- Forward probability laws:

$$\begin{cases} p(g|f, \tau_\epsilon) = \mathcal{N}(g|Hf, (1/\tau_\epsilon)I), & p(\tau_\epsilon|\alpha_{\epsilon 0}, \beta_{\epsilon 0}) = \mathcal{G}(\tau_\epsilon|\alpha_{\epsilon 0}, \beta_{\epsilon 0}), \\ p(f|z) = \prod_j \mathcal{N}(f_j|0, 1/z_j), & p(z|\alpha_0, \beta_0) = \prod_j \mathcal{G}(z_j|\alpha_0, \beta_0). \end{cases} \quad (26)$$

- Joint posterior laws:

$$\begin{aligned} p(f, z, \tau_\epsilon | g, \alpha_0, \beta_0, \alpha_{\epsilon 0}, \beta_{\epsilon 0}) &\propto p(g|f, \tau_\epsilon) p(f|z) p(z|\alpha_0, \beta_0) p(\tau_\epsilon|\alpha_{\epsilon 0}, \beta_{\epsilon 0}) \\ &\propto \tau_\epsilon^{-M/2} \exp\left\{-\frac{1}{2}\tau_\epsilon \|g - Hf\|^2\right\} \prod_j z_j^{-1/2} \exp\left\{-\frac{1}{2}z_j f_j^2\right\} \\ &\quad \prod_j z_j^{-\alpha_0+1} \exp\{-\beta_0 z_j\} \tau_\epsilon^{-\alpha_{\epsilon 0}+1} \exp\{-\beta_{\epsilon 0}\tau_\epsilon\}. \end{aligned} \quad (27)$$

- Joint MAP alternate maximization algorithm:

The objective of the JMAP optimization is:

$$(\hat{f}, \hat{z}, \hat{\tau}_\epsilon) = \arg \max_{(f, z, \tau_\epsilon)} \{p(f, z, \tau_\epsilon | g, \alpha_0, \beta_0, \alpha_{\epsilon 0}, \beta_{\epsilon 0})\}. \quad (28)$$

The alternate optimization is an iterative optimization, respectively with respect to  $f$ ,  $z$  and  $\tau$ :

$$\begin{cases} \hat{f} = \arg \min_f \left\{ \hat{\tau}_\epsilon \|g - Hf\|^2 + \sum_j \hat{z}_j f_j^2 \right\}, \\ \hat{z} = \arg \min_z \left\{ \frac{N+2\alpha_0-2}{2} \ln z_j + \sum_j z_j \left( \frac{1}{2} \hat{f}_j^2 + \beta_0 \right) \right\}, \\ \hat{\tau}_\epsilon = \arg \min_{\tau_\epsilon} \left\{ \left( \frac{M}{2} + \alpha_{\epsilon 0} - 1 \right) \ln \tau_\epsilon + \left( \frac{1}{2} \|g - H\hat{f}\|^2 + \beta_{\epsilon 0} \right) \right\}. \end{cases} \quad (29)$$

The first optimization can be done either analytically or using any gradient based algorithm. The second and the third optimizations have analytical expressions:

$$\begin{cases} \hat{f} = \hat{\tau}_\epsilon \hat{\Sigma} H' g \text{ with } \hat{\Sigma} = \left( \hat{\tau}_\epsilon H' H + \hat{Z} \right)^{-1} \text{ where } \hat{Z}^{-1} = \text{diag}[\hat{z}], \\ \hat{z}_j = \left( \frac{1}{2} \hat{f}_j^2 + \beta_0 \right) / \left( \frac{M}{2} + \alpha_{\epsilon 0} - 1 \right), \\ \hat{\tau}_\epsilon = \left( \frac{1}{2} \|g - H\hat{f}\|^2 + \beta_{\epsilon 0} \right) / \left( \frac{M}{2} + \alpha_{\epsilon 0} - 1 \right). \end{cases} \quad (30)$$

One iteration of this algorithm is shown in Figure 3.

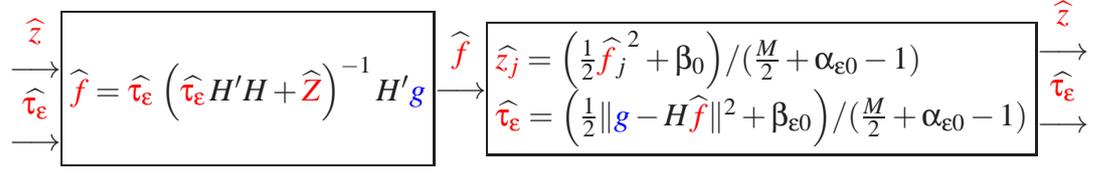


Figure 3: One iteration of the JMAP algorithm.

The main drawback of this method is that the uncertainties of the solution at each step is not accounted for for the next step.

- VBA posterior laws:

$$\begin{cases} q_1(f|\tilde{\mu}, \tilde{\Sigma}) = \mathcal{N}(f|\tilde{\mu}, \tilde{\Sigma}), & \tilde{\mu} = \tilde{\tau} \tilde{\Sigma} H' g, \tilde{\Sigma} = (\tilde{\tau} H' H + \tilde{Z})^{-1} \text{ with } \tilde{Z}^{-1} = \text{diag}[\tilde{z}], \\ q_{2j}(z_j) = \mathcal{G}(z_j|\tilde{\alpha}_j, \tilde{\beta}_j), & \tilde{\alpha}_j = \alpha_0 + \frac{1}{2}, \tilde{\beta}_j = \beta_0 + \langle f_j^2 \rangle / 2, \\ q_3(\tau_\epsilon) = \mathcal{G}(\tau_\epsilon|\tilde{\alpha}_\epsilon, \tilde{\beta}_\epsilon) \\ \tilde{\alpha}_\epsilon = \alpha_{\epsilon 0} + (n+1)/2, & \tilde{\beta}_\epsilon = \beta_{\epsilon 0} + \frac{1}{2}[\|g\|^2 - 2 \langle f \rangle' H' g + H' \langle f f \rangle H]. \end{cases} \quad (31)$$

with

$$\langle f \rangle = \tilde{\mu}, \langle f f \rangle = \tilde{\Sigma} + \tilde{\mu} \tilde{\mu}', \langle f_j^2 \rangle = [\tilde{\Sigma}]_{jj} + \tilde{\mu}_j^2, \tilde{\tau} = \frac{\tilde{\alpha}_{\tau_\epsilon}}{\tilde{\beta}_{\tau_\epsilon}} \text{ and } \tilde{z}_j = \frac{\tilde{\alpha}_j}{\tilde{\beta}_j}. \quad (32)$$

The expression of the free energies can be obtained as follows:

$$\begin{aligned} \mathcal{F}(q) &= \left\langle \ln \frac{p(f, z, \tau, g|M)}{q(f, z, \tau)} \right\rangle = \\ &= \langle \ln p(g|f, z, \tau) \rangle + \langle \ln p(f|z, \tau) \rangle + \langle \ln p(z|\tau) \rangle + \langle -\ln q(f) \rangle + \langle -\ln q(z) \rangle + \langle -\ln q(\tau) \rangle \end{aligned} \quad (33)$$

where

$$\begin{aligned} \langle \ln p(g|f, \tau_\epsilon) \rangle &= \frac{n}{2} (\langle \ln \tau_\epsilon \rangle - \ln(2\pi)) - \frac{1}{2} \{ \langle \lambda \rangle g' g - 2 \langle f \rangle' H' g + H' \langle f f \rangle H \} \\ \langle -\ln p(f|z) \rangle &= -\frac{n+1}{2} \ln(2\pi) - \frac{1}{2} \left\{ \sum_j \langle \ln \alpha_j \rangle \langle \alpha_j \rangle \langle f_j^2 \rangle \right\} \\ \langle -\ln p(z) \rangle &= -(n+1) \alpha_{\epsilon 0} \ln \beta_{\epsilon 0} + (\alpha_{\epsilon 0} - 1) \sum_j \langle \ln \alpha_j \rangle - \beta \langle \alpha_j \rangle - (n+1) \ln \Gamma(\alpha) \\ \langle p(\tau_\epsilon) \rangle &= c \ln d + (c-1) \langle \ln \tau_\epsilon \rangle - d \langle \lambda \rangle - \ln \Gamma(c) \\ \langle -\ln q(f) \rangle &= -\frac{n+1}{2} (1 + \ln(2\pi)) - \frac{1}{2} \ln |\Sigma_f| \\ \langle -\ln q(z) \rangle &= -\sum_j [\tilde{\alpha}_j \ln(\tilde{\beta}_j) + (\tilde{\alpha}_j - 1) \langle \ln \tilde{\alpha}_j \rangle - \tilde{\beta}_j \langle \alpha_j \rangle - \ln \Gamma(\tilde{\alpha}_j)] \\ \langle q(\tau_\epsilon) \rangle &= \tilde{c} \ln \tilde{d} + (\tilde{c} - 1) \langle \ln \tau \rangle - \tilde{d} \langle \lambda \rangle - \ln \Gamma(\tilde{c}). \end{aligned}$$

In these equations,

$$\begin{cases} \langle \ln a_j \rangle = \psi(\tilde{a}_j) - \ln \tilde{b}_j, \\ \langle \ln \tau \rangle = \psi(\tilde{c}) - \ln \tilde{d}, \\ \psi(a) = \frac{\partial \ln \Gamma(a)}{\partial a}. \end{cases} \quad (34)$$

The three steps of this algorithm is shown in Figure 4.

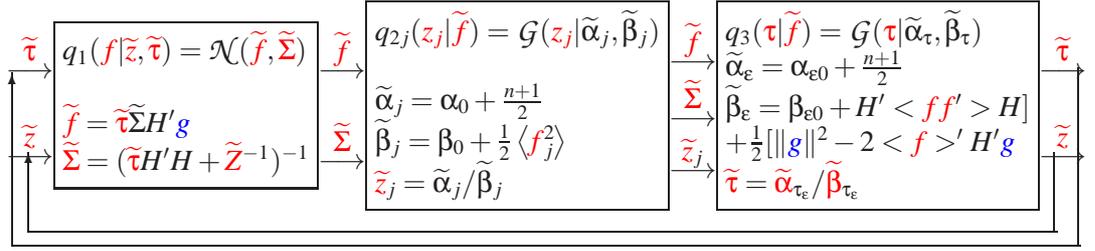


Figure 4: The three steps of the Bayesian Variational Approximation Algorithm.

## BAYESIAN VARIATIONAL APPROXIMATION WITH MIXTURE OF GAUSSIANS PRIORS

The mixture models are also very commonly used as prior models. In particular the Mixture of two Gaussians (MoG2) model:

$$p(f|\lambda, v_1, v_0) = \prod_j (\lambda \mathcal{N}(f_j|0, v_1) + (1 - \lambda) \mathcal{N}(f_j|0, v_0)) \quad (35)$$

which can also be expressed through the binary valued hidden variables  $z_j \in \{0, 1\}$

$$\begin{cases} p(f|z) = \prod_j p(f_j|z_j) = \prod_j \mathcal{N}(f_j|0, v_{z_j}) \propto \exp \left\{ -\frac{1}{2} \sum_j \frac{f_j^2}{v_{z_j}} \right\} \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (36)$$

In general  $v_1 \gg v_0$  and  $\lambda$  measures the sparsity ( $0 < \lambda \ll 1$ ) [11]. In this case also all the equations are very similarly can be obtained. Here, we do not have enough place to write them.

## CONCLUSIONS

In this paper, a VBA method is proposed for doing Bayesian computations for inverse problems where a hierarchical prior modeling is used for the unknowns. In particular, two prior models are considered: the Student-t and the mixture of Gaussian models. In both cases, these priors can be written via hidden variables which gives the model a hierarchical structure which is used to do the factorization. For some applications see for example [19, 7, 4, 6, 5, 10, 8, 16] and two other related papers in this volume.

## REFERENCES

1. M. Beal. *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, Gatsby Computational Neuroscience Unit, University College London, 2003.
2. Sotirios Chatzis and Theodora Varvarigou. Factor analysis latent subspace modeling and robust fuzzy clustering using t-distributionsclassification of binary random patterns. *IEEE Trans. on Fuzzy Systems*, 17:505–517, 2009.
3. R. A. Choudrey. *Variational Methods for Bayesian Independent Component Analysis*. PhD thesis, University of Oxford, 2002.
4. N. Chu, A. Mohammad-Djafari, and J. Picheral. Two robust super-resolution approaches with sparsity constraint and sparse regularization for near-field wideband extended aeroacoustic source imaging. In *Berlin Beamforming Conference 2012 (BeBeC2012)*, number 29, Berlin, Germany, Feb.22-23,2012.
5. N. Chu, A. Mohammad-Djafari, and J. Picheral. Bayesian sparse regularization in near-field wideband aeroacoustic imaging for wind tunnel test. In *11th CongrÃs Franais d'Acoustique and 2012 IOA annual meeting*, Nantes, France, Apr.23-27,2012.
6. N. Chu, A. Mohammad-Djafari, and J. Picheral. A bayesian sparse inference approach in near-field wideband aeroacoustic imaging. In *2012 IEEE International Conference on Image Processing*, Orlando, USA, Sep.30-Oct.4, 2012.
7. N. Chu, J. Picheral, and A. Mohammad-Djafari. A robust super-resolution approach with sparsity constraint for near-field wideband acoustic imaging. In *IEEE International Symposium on Signal Processing and Information Technology*, pages 286–289, Bilbao, Spain, Dec.14-17,2011.
8. Mircea Dumitru and Ali Mohammad-Djafari. Estimating the period of a signal through inverse problem modeling and bayesian inference with sparsity enforcing prior. In *32nd International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, IPP, Garching near Munich, Germany,15-20 July 2012.
9. Aurélia Fraysse and Thomas Rodet. A gradient-like variational Bayesian algorithm. In *SSP 2011*, number S17.5, pages 605–608, Nice, France, jun 2011.
10. Leila Gharsalli, Ali Mohammad-Djafari, Aurélia Fraysse, and Thomas Rodet. Variational bayesian approximation with scale mixture prior for inverse problems: a numerical comparison between three algorithms. In *32nd International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, IPP, Garching near Munich, Germany,15-20 July 2012.
11. J. Rao. H. Ishwaran. Spike and Slab variable selection: Frequentist and Bayesian strategies. *Annals of Statistics*, 2005.
12. L. He, H. Chen, and L. Carin. Tree-Structured Compressive Sensing With Variational Bayesian Analysis. *IEEE Signal. Proc. Let.*, 17(3):233–236, 2010.
13. A. C. Likas and N. P. Galatsanos. A variational approach for bayesian blind image deconvolution. *IEEE Transactions on Signal Processing*, 2004.
14. A Mohammad-Djafari. Bayesian approach with prior models which enforce sparsity in signal and image processing. *EURASIP Journal on Advances in Signal Processing*, Special issue on Sparse Signal Processing (1 Mars 2012):2012:52, 2012.
15. T. Park and G. Casella. The Bayesian Lasso. *Journal of the American Statistical Association*, 2008.
16. R. Pérenon, A. Mohammad-Djafari, E. Sage1, L. Duraffourg, S. Hentz1, A. Brenac, R. Morel, and P. Grangeat. Mcmc-based bayesian estimation algorithm dedicated to NEMS mass spectrometry. In *32nd International Workshop on Bayesian Inference and Maximum Entropy Methods in Science and Engineering*, IPP, Garching near Munich, Germany,15-20 July 2012.
17. M. Tipping. Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 2001.
18. John Winn, Christopher M. Bishop, and Tommi Jaakkola. Variational message passing. *Journal of Machine Learning Research*, 6:661–694, 2005.
19. Sha Zhu, Ali Mohammad-Djafari, Hongqiang Wang, Bin Deng, Xiang Li, and Junjie Mao. Parameter estimation for sar micromotion target based on sparse signal representation. *Eurasip Journal of Signal Processing*, special issue "sparse approximations in signal and image processing", 2012.