

REVIEW

Open Access

Bayesian approach with prior models which enforce sparsity in signal and image processing

Ali Mohammad-Djafari

Abstract

In this review article, we propose to use the Bayesian inference approach for inverse problems in signal and image processing, where we want to infer on sparse signals or images. The sparsity may be directly on the original space or in a transformed space. Here, we consider it directly on the original space (impulsive signals). To enforce the sparsity, we consider the probabilistic models and try to give an exhaustive list of such prior models and try to classify them. These models are either heavy tailed (generalized Gaussian, symmetric Weibull, Student-t or Cauchy, elastic net, generalized hyperbolic and Dirichlet) or mixture models (mixture of Gaussians, Bernoulli-Gaussian, Bernoulli-Gamma, mixture of translated Gaussians, mixture of multinomial, etc.). Depending on the prior model selected, the Bayesian computations (optimization for the joint maximum a posteriori (MAP) estimate or MCMC or variational Bayes approximations (VBA) for posterior means (PM) or complete density estimation) may become more complex. We propose these models, discuss on different possible Bayesian estimators, drive the corresponding appropriate algorithms, and discuss on their corresponding relative complexities and performances.

Keywords: sparsity, Bayesian approach, sparse priors, inverse problems

1 Introduction

In many generic inverse problems in signal and image processing we want to infer on an unknown signal $f(t)$ or an unknown image $f(\mathbf{r})$ with $\mathbf{r} = (x, y)$ through an observed signal $g(s)$ or an observed image $g(\mathbf{s})$ related between them through an operator \mathcal{H} such as convolution $g = h * f$ or any other linear or non linear transformation $g = \mathcal{H}f$. When this relation is linear and we have discretized the problem, we arrive to the relation:

$$\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}, \quad (1)$$

where $\mathbf{f} = [f_1, \dots, f_n]'$ represents the unknowns, $\mathbf{g} = [g_1, \dots, g_m]'$ the observed data, $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_m]'$ the errors of modeling and measurement and \mathbf{H} the matrix of the system response. We may note that, even if the noise could be neglected ($\boldsymbol{\epsilon} = 0$) and the matrix \mathbf{H} invertible ($m = n$), in general, the solution $\hat{\mathbf{f}} = \mathbf{H}^{-1}\mathbf{g}$ is not forcibly the good solution, because this solution may be too sensitive to small changes in the data due to the ill-conditioning of this matrix. For the general case of $m \neq n$, one tries to obtain a regularized solution, for example by

defining it as the optimizer of a two parts criterion

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2\} \quad (2)$$

which is given by $\hat{\mathbf{f}} = [\mathbf{H}\mathbf{H}' + \lambda\mathbf{I}]^{-1}\mathbf{H}'\mathbf{g}$. When the regularization parameter $\lambda = 0$, one gets a generalized inverse $\hat{\mathbf{f}} = [\mathbf{H}\mathbf{H}']^{-1}\mathbf{H}'\mathbf{g}$ and when \mathbf{H} invertible, one gets the normal inverse solution $\hat{\mathbf{f}} = \mathbf{H}^{-1}\mathbf{g}$. The regularization theory has been developed since the pioneer work of Tikhonov [1] and Tikhonov and Arsénine [2] who had introduced a quadratic regularization terms to account for some prior properties of the solution (smoothness). Since that, many different regularization terms have been proposed. In particular, in place of L_2 norm: $L_2(\mathbf{f}) = \|\mathbf{f}\|_2^2 = \sum_j |f_j|^2$, it has been proposed to use the L_0 norm $L_0(\mathbf{f}) = \|\mathbf{f}\|_0 = \sum_j \delta(f_j)$ or the L_1 norm $L_1(\mathbf{f}) = \|\mathbf{f}\|_1 = \sum_j |f_j|$ to enforce the sparsity of the solution [3-11]. Then, due to the fact that $L_0(\mathbf{f})$ is not convex and $L_1(\mathbf{f})$ is convex, but not continuous, the optimization of a criterion with these expressions becomes more difficult than the L_2 norm case. For this reason, there was a great number of works who

Correspondence: djafari@lss.supelec.fr
Laboratoire des signaux et systèmes (L2S), UMR 8506 CNRS-SUPELEC-UNIV
PARIS SUD, SUPELEC, Plateau de Moulon, 91192 Gif-sur-Yvette, France

specialized in proposing algorithms for the optimization of such criteria.

Interestingly, defining the solution of the problem (1) as the optimization of a criterion with two parts can be assimilated to a maximum a posteriori (MAP) solution in a Bayesian approach where the first term of the criterion (2) can be related to the likelihood and the second term to a prior model as we will see in the following where the main objective is to show how the Bayesian approach can go farther than the regularization in at least the following aspects:

- A better account for the noise term characteristics;
- A better and easier way for translating the prior knowledge and in particular the sparsity;
- New tools for assessing the regularization parameter, a great subject of discussion for all those work with regularization theory;
- New solutions and new tools for doing computations (optimizations and integrations).

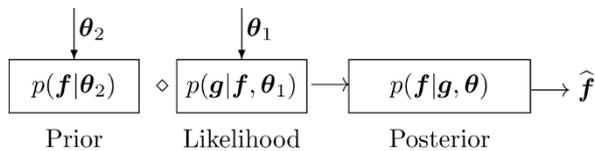
1.1 The Bayesian approach

The Bayesian inference approach is based on the posterior law:

$$p(f | g, \theta_1, \theta_2) = \frac{p(g | f, \theta_1) p(f | \theta_2)}{p(g | \theta_1, \theta_2)} \propto p(g | f, \theta_1) p(f | \theta_2)$$

where the sign \propto stands for “proportional to”, $p(g | f, \theta_1)$ is the likelihood, $p(f | \theta_2)$ the prior model, $\theta = (\theta_1, \theta_2)$ are their corresponding parameters (often called the hyper parameters of the problem) and $p(g | \theta_1, \theta_2)$ is called the evidence of the model.

This general Bayesian approach is illustrated as follows:



In this approach, the likelihood $p(g | f, \theta_1)$ summarizes our knowledge about the noise and the model linking the observed data g to the unknowns f and the prior term $p(f | \theta_2)$ summarizes our incomplete prior knowledge about the unknowns and the posterior law $p(f | g, \theta)$ combines these two terms and contains all our state of knowledge about the unknowns f after accounting for the prior and the observed data.

As a very simple example, when the noise is assumed to be Gaussian, then the MAP solution $\hat{f} = \arg \max_f \{p(f | g, \theta)\}$ is obtained as the optimizer of the criterion $J(f) = \|g - Hf\|^2 + \lambda \Omega(f)$ where the expression of $\Omega(f)$ depends on the prior law. When the prior

knowledge is translated as a Gaussian probability law, then $\Omega(f) = \|f\|_2^2$ and when it is translated as a Laplace probability law, then $\Omega(f) = \|f\|_1$ [12-14].

The first interest of using the Bayesian approach to the regularization approach is to have new tools for handling the hyper parameters [15].

1.2 Full Bayesian approach

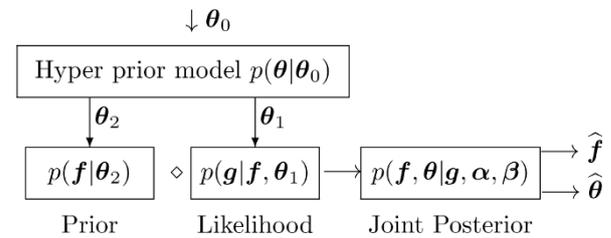
When the parameters θ have to be estimated too, we can assign them a prior $p(\theta | \theta_0)$ with fixed values for θ_0 (often called hyper-hyper-parameters) and express the joint posterior

$$p(f, \theta | g, \theta_0) = \frac{p(g | f, \theta_1) p(f | \theta_2) p(\theta | \theta_0)}{p(g | \theta_0)} \quad (4)$$

and then try to estimate them jointly, for example joint MAP [16]:

$$(\hat{f}, \hat{\theta}) = \arg \max_{(f, \theta)} \{p(f, \theta | g, \theta_0)\} \quad (5)$$

This Full Bayesian approach is illustrated as follows:



One may also first integrate out one of them, for example f to obtain

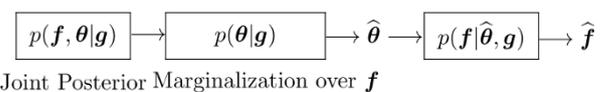
$$p(\theta | g, \theta_0) = \int p(f, \theta | g, \theta_0) df, \quad (6)$$

estimate θ , for example by

$$\hat{\theta} = \arg \max_{\theta} \{p(\theta | g, \theta_0)\} \quad (7)$$

and then use it for the estimation of the other one using $p(f | g, \hat{\theta})$.

This approach (called sometimes type II maximum likelihood) is illustrated as follows:



However, very often this marginalization cannot be done analytically and so the optimization for the estimation of θ cannot be achieved. In such cases, the expectation-maximization (EM) algorithms can be helpful [17]. Considering g as *incomplete data*, f as *hidden variable*, (g, f) as *complete data* and noting $\ln p(g | \theta)$ as

incomplete data log-likelihood and $\ln p(\mathbf{g}, \mathbf{f}|\boldsymbol{\theta})$ complete data log-likelihood, the classical EM algorithm writes:

$$\begin{cases} \text{E - step : } q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)}) = E_{p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)})} \{ \ln p(\mathbf{g}, \mathbf{f}|\boldsymbol{\theta}) \} \\ \text{M - step : } \hat{\boldsymbol{\theta}}^{(k)} = \arg \max_{\boldsymbol{\theta}} \{ q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k-1)}) \} \end{cases} \quad (8)$$

The Bayesian version (Bayesian EM) is not very far and differs only by the introduction of $p(\boldsymbol{\theta})$:

$$\begin{cases} \text{E - step : } q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k)}) = E_{p(\mathbf{f}|\mathbf{g}, \hat{\boldsymbol{\theta}}^{(k)})} \{ \ln p(\mathbf{g}, \mathbf{f}|\boldsymbol{\theta}) + \ln p(\boldsymbol{\theta}) \} \\ \text{M - step : } \hat{\boldsymbol{\theta}}^{(k)} = \arg \max_{\boldsymbol{\theta}} \{ q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}^{(k-1)}) \} \end{cases} \quad (9)$$

This is illustrated as follows:

$$p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) \longrightarrow \text{EM, Bayes EM} \longrightarrow \hat{\boldsymbol{\theta}} \longrightarrow p(\mathbf{f}|\hat{\boldsymbol{\theta}}, \mathbf{g}) \longrightarrow \hat{\mathbf{f}}$$

As we mentioned before, one of the main steps in the Bayesian approach is the prior modeling which has the role of translating our prior knowledge on the unknown signal or image in a probability law. Sparsity is one of the prior knowledge we may translate. The main objective of this article is to see what are the different possibilities.

1.3 Prior modeling

In this article, we propose different prior modeling for signals and images which can be used in a Bayesian inference approach in many inverse problems in signal and image processing where we want to infer on sparse signals or images. The sparsity may be directly on the original space or in a transformed space (see Figures 1, 2, 3, and 4). In this article, we consider the sparsity directly in the original domain.

The prior models discussed are the following:

- generalized Gaussian (GG) with Gaussian (G) and Laplace or double exponential (DE) as particular cases;
- symmetric Weibull (W) with symmetric Rayleigh (R) and again the DE as particular cases;
- Student-t (St) with Cauchy (C) as particular case;

- Elastic net prior model;
- generalized hyperbolic model;
- Dirichlet and symmetric Dirichlet;
- Mixture of two centered Gaussians (MoG2), one with very small and one with a large variances;
- Bernoulli-Gaussian (BG), also called *Spike and slab*;
- Mixture of two Gammas (MoGamm);
- Bernoulli-Gamma (BGamma);
- Mixture of three Gaussians (MoG3), one centered with very small variance and two symmetrically centered on positive and negative axes and large variances;
- Mixture of one Gaussian and two Gammas (MoG-Gammas), and in a more summary the case of
- Bernoulli-Multinomial (BMult) or mixture of Dirichlet (MoD).

Some of these models are well-known [12-14,18-26], some others less. In general, we can classify them into two categories: (i) simple non Gaussian models with heavy tails and (ii) mixture models with hidden variables which result to hierarchical models.

In the Section 2, we give more details about the sparsity and all these prior models which enforce the sparsity.

1.4 Bayesian computation

The second main step in the Bayesian approach is to do the computations. Depending on the prior model selected, the Bayesian computations needed are:

- For simple prior models:
 - Simple optimization of $p(\mathbf{f}|\boldsymbol{\theta}, \mathbf{g})$ for the MAP:

$$p(\mathbf{f}|\boldsymbol{\theta}, \mathbf{g}) \longrightarrow \text{Simple Optimization Algorithm} \longrightarrow \hat{\mathbf{f}}$$

- Joint optimization $p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g})$ for joint MAP:

$$p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) \longrightarrow \text{Joint Optimization Algorithm} \begin{matrix} \longrightarrow \hat{\mathbf{f}} \\ \longrightarrow \hat{\boldsymbol{\theta}} \end{matrix}$$

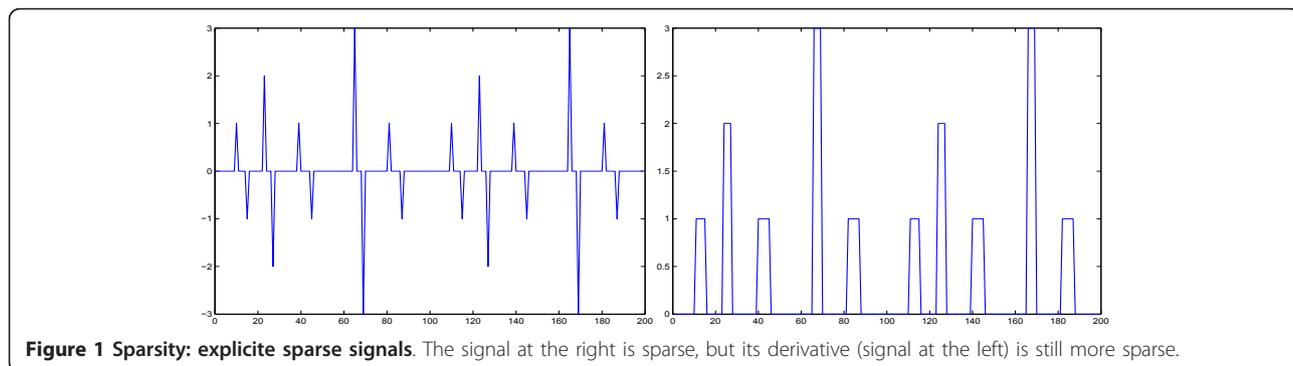
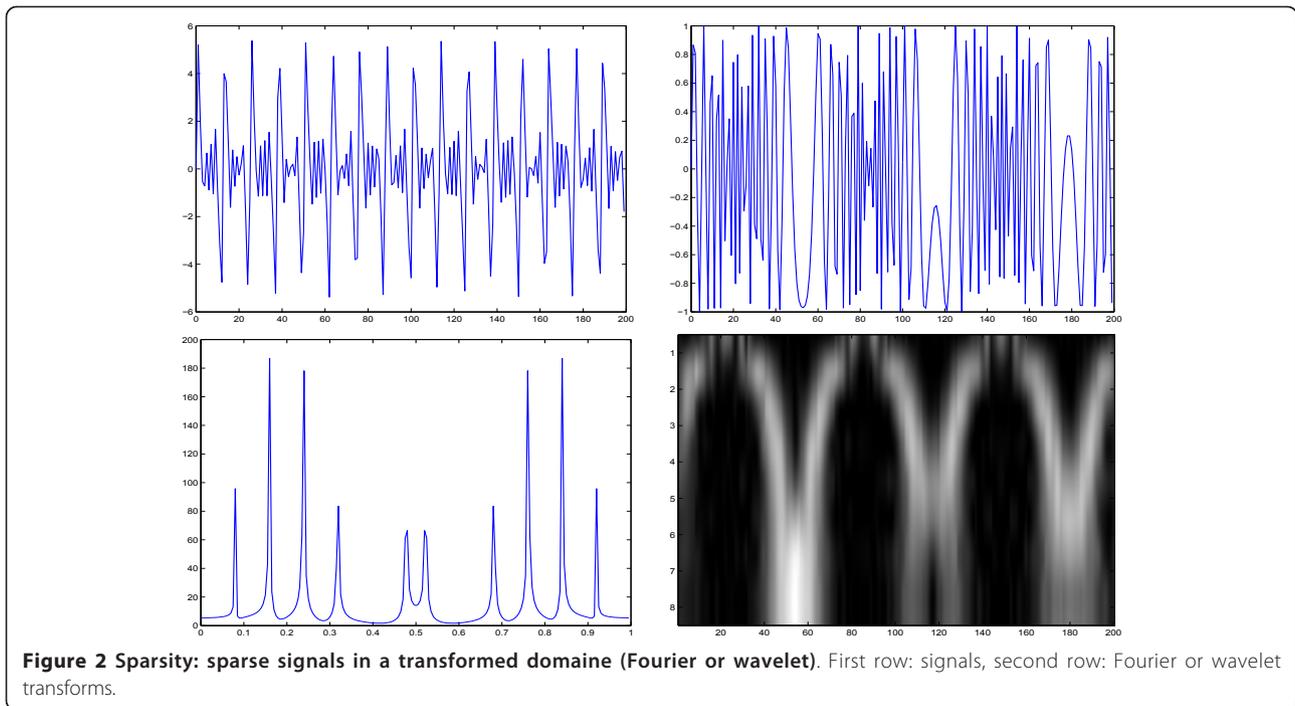
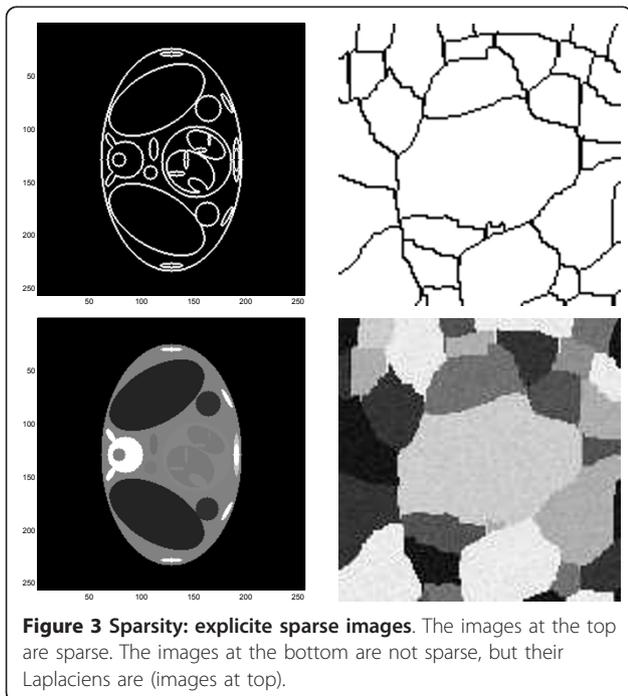
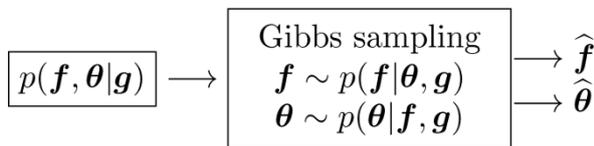


Figure 1 Sparsity: explicite sparse signals. The signal at the right is sparse, but its derivative (signal at the left) is still more sparse.



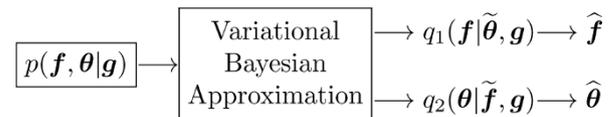
- Generation of samples from the conditionals $p(\mathbf{f}|\boldsymbol{\theta}, \mathbf{g})$ and $p(\boldsymbol{\theta}|\mathbf{f}, \mathbf{g})$ for the MCMC Gibbs sampling methods,



- Variational approximation (VA) of the joint $p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g})$ by a separable

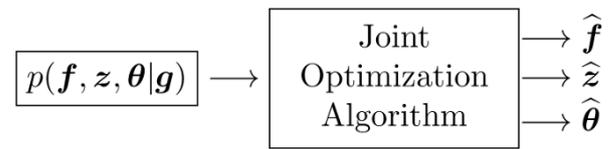
$$q(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) = q_1(\mathbf{f}|\tilde{\boldsymbol{\theta}}, \mathbf{g}) q_2(\boldsymbol{\theta}|\tilde{\mathbf{f}}, \mathbf{g})$$

and then using them for estimation

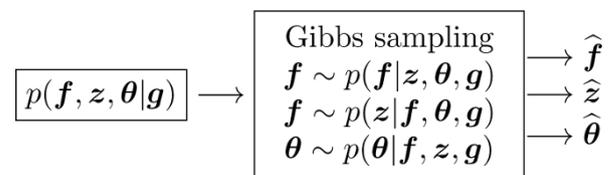


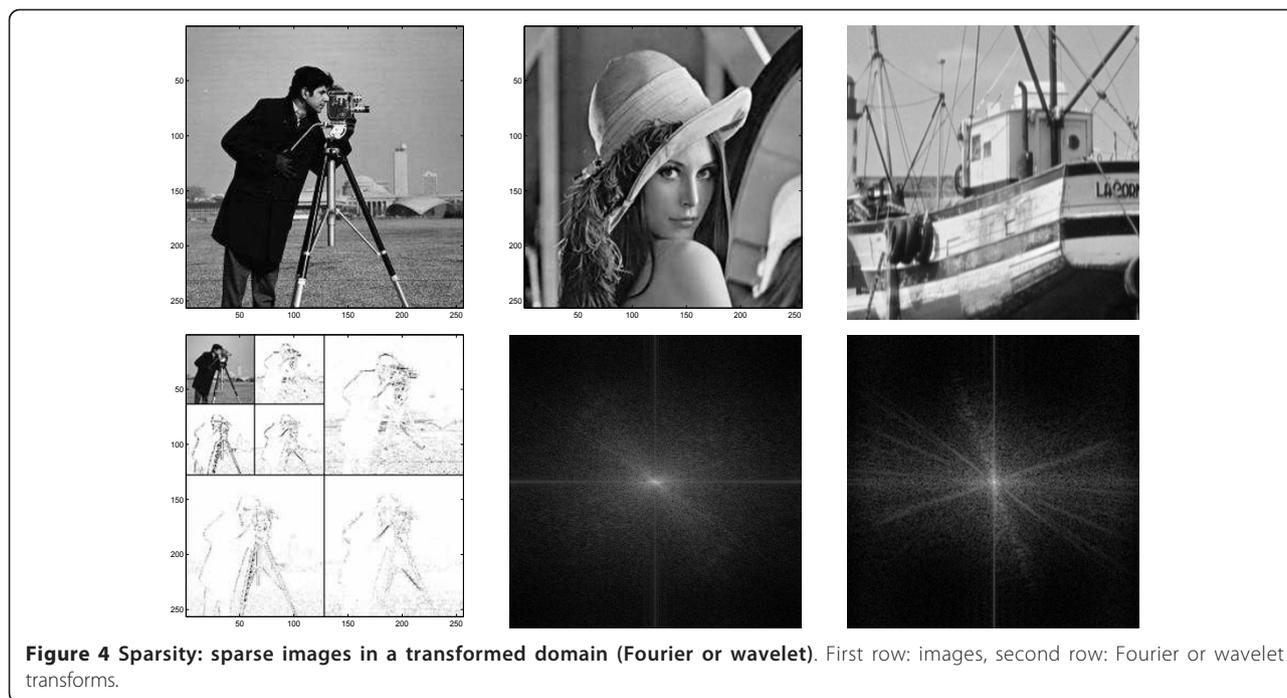
• For hierarchical prior models with hidden variables \mathbf{z} :

- Joint optimization $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}|\mathbf{g})$ for joint MAP,



- Generation of samples from the conditionals $p(\mathbf{f}|\mathbf{z}, \boldsymbol{\theta}, \mathbf{g})$, $p(\boldsymbol{\theta}|\mathbf{f}, \mathbf{z}, \mathbf{g})$ and $p(\mathbf{z}|\mathbf{f}, \boldsymbol{\theta}, \mathbf{g})$ for the MCMC Gibbs sampling methods:

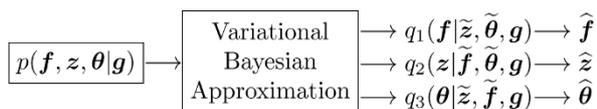




- Variational approximation (VA) of the joint $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g})$ by a separable

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

and then using them for estimation



The second main objective of this article is to discuss on the relative complexities and performances of the algorithms obtained with the proposed prior law.

The rest of the article is organized as follows:

In Section 2, we present in details the proposed prior models and discuss their properties. For example, we will see that the Student-t model can be interpreted as an infinite mixture with a variance hidden variable or that the BG model can be considered as the degenerate case of a MoG2 where one of the variances go to zero. Also, we will examine the less known models of MoG3 and MoGGammas where the heavy tails are obtained by combining a centered Gaussian and two large variance non-centered Gaussians or Gammas.

In Section 3, we examine the expression of the posterior laws that we obtain using these priors and discuss then on complexity of the Bayesian computation of the algorithms. In particular for the mixture models, we give details of the joint estimation of the signal and the

hidden variable as well as the hyper parameters (parameters of the mixtures and the noise) for unsupervised cases.

In Section 4, we give more details on the variational Bayesian approximation method, first for the general case and then for the case of mixture laws and more specifically the case of the Student-t considered as a continuous mixture.

Finally, we present the main conclusions of this article in Section 5.

2 Prior models enforcing sparsity

First, as we mentioned, the sparsity is a property which can be described either directly for the signal itself or after some transformation, for example on the derivative of the signal, or in more general on the coefficients of the projection of the signal on any basis or any set of functions.

Different prior models have been used to enforce sparsity.

2.1 Generalized Gaussian (GG), Gaussian (G) and double exponentials (DE) models

This is the simplest and the most used model (see for example, [27]). Its expression is:

$$p(\mathbf{f} | \gamma, \beta) = \prod_j \mathcal{GG}(f_j | \gamma, \beta) \propto \exp \left\{ -\gamma \sum_j |f_j|^\beta \right\} \quad (10)$$

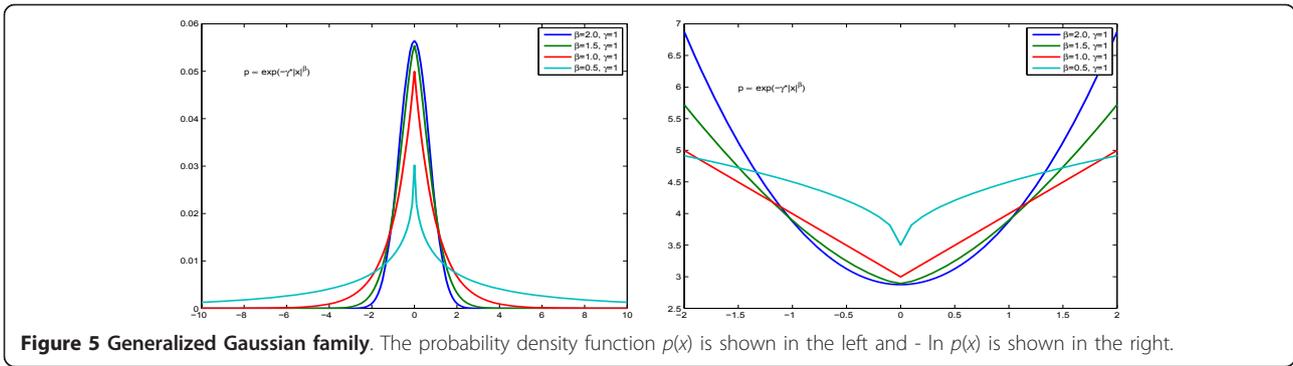


Figure 5 Generalized Gaussian family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

where

$$\mathcal{G}\mathcal{G}(f_j | \gamma, \beta) = \frac{\beta\gamma}{2\Gamma(1/\beta)} \exp\{-\gamma|f_j|^\beta\}. \quad (11)$$

Two particular cases are of importance:

- $\beta = 2$ (Gaussian):

$$p(f | \gamma) = \prod_j \mathcal{N}(f_j | 0, 1/(2\gamma)) \propto \exp\left\{-\gamma \sum_j |f_j|^2\right\} \propto \exp\{-\gamma \|f\|_2^2\} \quad (12)$$

- $\beta = 1$ (double exponential or Laplace):

$$p(f | \gamma) = \prod_j \mathcal{DE}(f_j | \gamma) \propto \exp\left\{-\gamma \sum_j |f_j|\right\} \propto \exp\{-\gamma \|f\|_1\} \quad (13)$$

The general shape of these priors are shown in Figure 5, where the cases $\beta = 1$ and $0 < \beta < 1$, which are of great interest for sparsity enforcing are compared to the Gaussian case $\beta = 2$.

2.2 Symmetric Weibull (W) and symmetric Rayleigh (R) models

The second model we consider is the symmetric Weibull probability density function (pdf):

$$p(f | \gamma, \beta) = \prod_j \mathcal{W}(f_j | \gamma, \beta) \propto \exp\left\{-\gamma \sum_j |f_j|^\beta + (\beta - 1) \log |f_j|\right\} \quad (14)$$

where

$$\mathcal{W}(f_j | \gamma, \beta) = c|f_j|^{(\beta-1)} \exp\{-\gamma|f_j|^\beta\} \quad (15)$$

and where $\gamma > 0$ and $\beta > 0$, and the particular cases of $\beta = 1$ is the double exponential and $\beta = 2$ is the symmetric Rayleigh distribution:

$$p(f | \gamma, \beta) = \prod_j \mathcal{R}(f_j | \gamma) \propto \exp\left\{-\gamma \sum_j |f_j|^2 + \log |f_j|\right\} \quad (16)$$

the cases where $0 < \beta < 1$ are of great interest for sparsity enforcing. This family of models are illustrated on Figure 6.

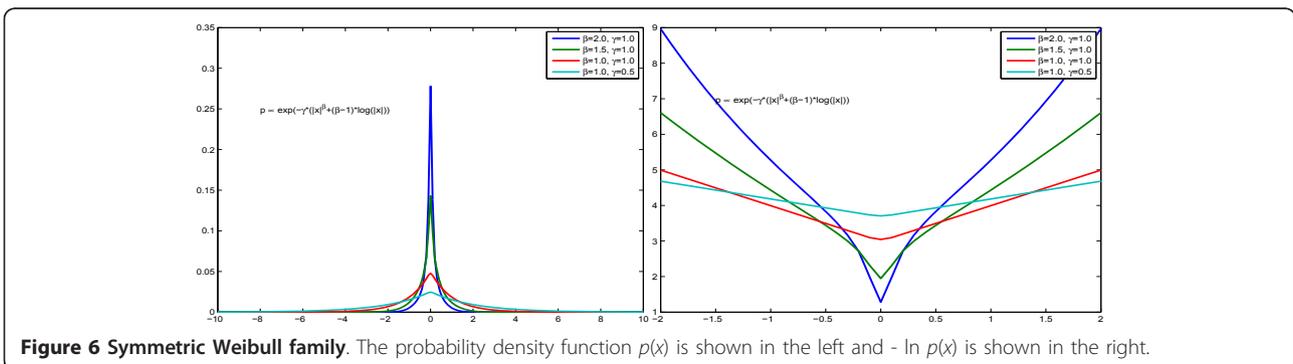


Figure 6 Symmetric Weibull family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

2.3 Student-t (St) and Cauchy (C) models

The second simplest model is the Student-t model:

$$p(f | \nu) = \prod_j St(f_j | \nu) \propto \exp \left\{ -\frac{\nu+1}{2} \sum_j \log(1 + f_j^2/\nu) \right\} \quad (17)$$

where

$$St(f_j | \nu) = \frac{1}{\sqrt{\pi\nu}} \frac{\Gamma((\nu+1)/2)}{\Gamma(\nu/2)} (1 + f_j^2/\nu)^{-(\nu+1)/2} \quad (18)$$

Knowing that

$$St(f_j | \nu) = \int_0^\infty \mathcal{N}(f_j | 0, 1/\tau_j) \mathcal{G}(\tau_j | \nu/2, \nu/2) d\tau_j \quad (19)$$

we can write this model via the positive hidden variables τ_j :

$$p(f, \boldsymbol{\tau}) = \prod_j p(f_j | \tau_j) = \prod_j \mathcal{N}(f_j | 0, 1/\tau_j) \propto \exp \left\{ -\frac{1}{2} \sum_j \tau_j f_j^2 \right\} \quad (20)$$

$$p(\tau_j | a, b) = \mathcal{G}(\tau_j | a, b) \propto \tau_j^{(a-1)} \exp\{-b\tau_j\}$$

with $a = b = \nu/2$

Cauchy model is obtained when $\nu = 1$:

$$p(f) = \prod_j \mathcal{C}(f_j) \propto \exp \left\{ -\sum_j \log(1 + f_j^2) \right\} \quad (21)$$

This family of models are illustrated on Figure 7.

2.4 Elastic Net (EN) prior model

A prior model inspired from elastic net regression literature [28] is:

$$p(f | \nu) = \prod_j \mathcal{EN}(f_j | \nu) \propto \exp \left\{ -\sum_j (\gamma_1 |f_j| + \gamma_2 f_j^2) \right\} \quad (22)$$

where

$$\mathcal{EN}(f_j | \nu) = \mathcal{N}(0, 1/\gamma_1) \mathcal{DE}(\gamma_1) \propto \exp \left\{ -\gamma_1 |f_j| - \gamma_2 f_j^2 \right\} \quad (23)$$

which is a product of a Gaussian and a double exponential pdfs. This family of models are illustrated on Figure 8.

2.5 Generalized hyperbolic (GH) prior model

Another general prior model which can be used is:

$$p(f | \delta, \nu, \beta) = \prod_j (\delta^2 + f_j^2)^{(\nu-1/2)/2} \exp\{\beta x\} K_{\nu-1/2}(\alpha \sqrt{\delta^2 + f_j^2}) \quad (24)$$

where $K_{\nu-1/2}$ is the second kind Bessel function of order $(\nu - 1/2)$. This family of models are illustrated on Figure 9.

2.6 Dirichlet (D) and symmetric Dirichlet (SD) models

When f_j are positive and sums to one, we can use the Dirichlet model

$$\mathcal{D}(f | \boldsymbol{\alpha}) \propto \prod_j f_j^{\alpha_j-1} \quad \text{with } f_j > 0, \quad \sum_j f_j = 1 \quad (25)$$

where $\boldsymbol{\alpha} = \{\alpha_1, \dots, \alpha_N\}$ with $\alpha_j > 0$. The proportionality constant is

$$B(\boldsymbol{\alpha}) = \frac{\prod_j \Gamma(\alpha_j)}{\Gamma(\sum_j \Gamma(\alpha_j))} \quad (26)$$

It is noted that the support of this distribution is $[0,1]^N$ and $\|f\|_1 = \sum_j f_j = 1$.

It is also interesting to note that the domain of the Dirichlet distribution is itself a probability distribution, specifically a N -dimensional discrete distribution and the set of points in the support of a N -dimensional Dirichlet distribution is the open standard $N - 1$ -simplex, which is a generalization of a triangle, embedded in the next-higher dimension.

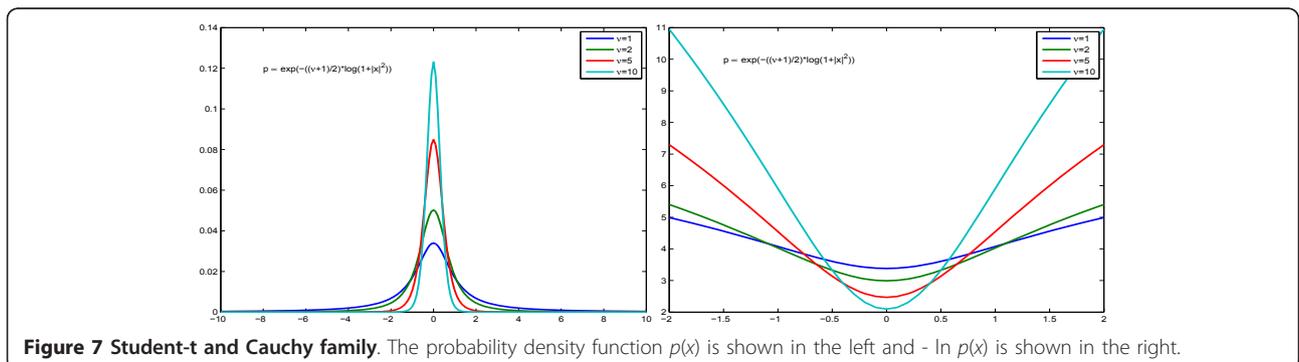


Figure 7 Student-t and Cauchy family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

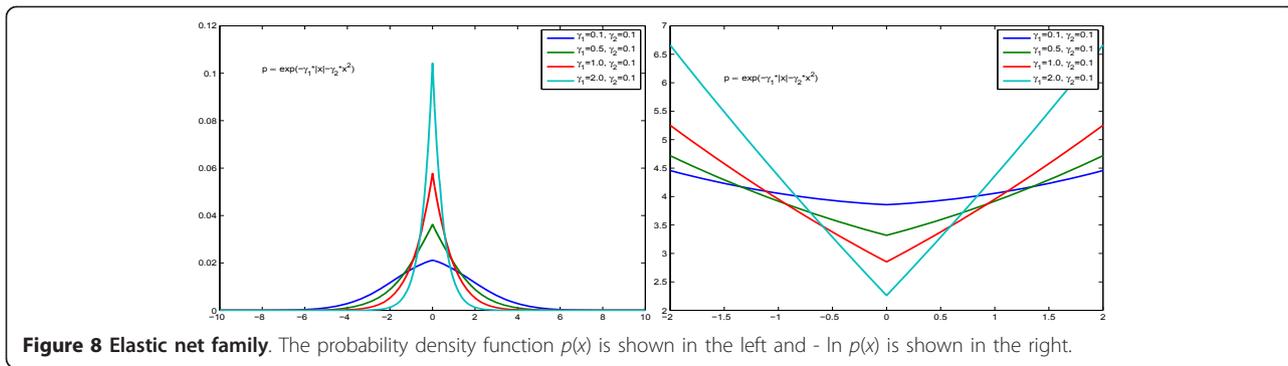


Figure 8 Elastic net family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

A very common special case is the *symmetric Dirichlet* (SD) distribution, where all of the elements making up the parameter vector α have the same value α called the concentration parameter:

$$\mathcal{D}(f|\alpha) \propto \prod_j f_j^{\alpha-1} \text{ with } f_j > 0, \quad \sum_j f_j = 1 \quad (27)$$

When $\alpha > 1$, the symmetric Dirichlet distribution is equivalent to a uniform distribution over the open standard standard $N - 1$ -simplex, i.e., it is uniform over all points in its support. $\alpha > 1$ prefer variants that are dense, evenly-distributed distributions, i.e., all probabilities f_j returned are similar to each other. $\alpha < 1$ prefer

sparse distributions, i.e., most of the probabilities f_j returned will be close to 0, and the vast majority of the mass will be concentrated in a few of them. This is the case on which we are interested. An illustration of this family of models are illustrated on Figure 10.

2.7 Mixture of two Gaussians (MoG2) model

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

$$p(f|\lambda, \nu_1, \nu_0) = \prod_j (\lambda \mathcal{N}(f_j|0, \nu_1) + (1 - \lambda) \mathcal{N}(f_j|0, \nu_0)) \quad (28)$$

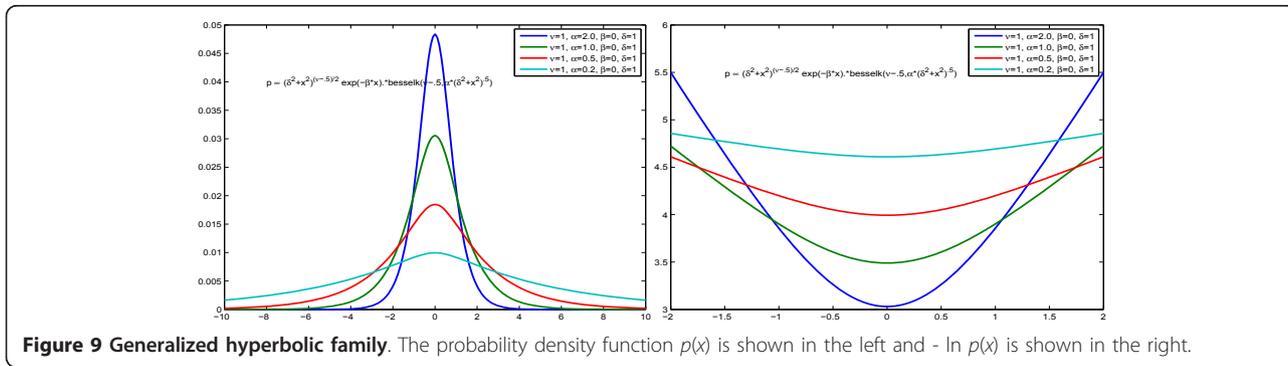


Figure 9 Generalized hyperbolic family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

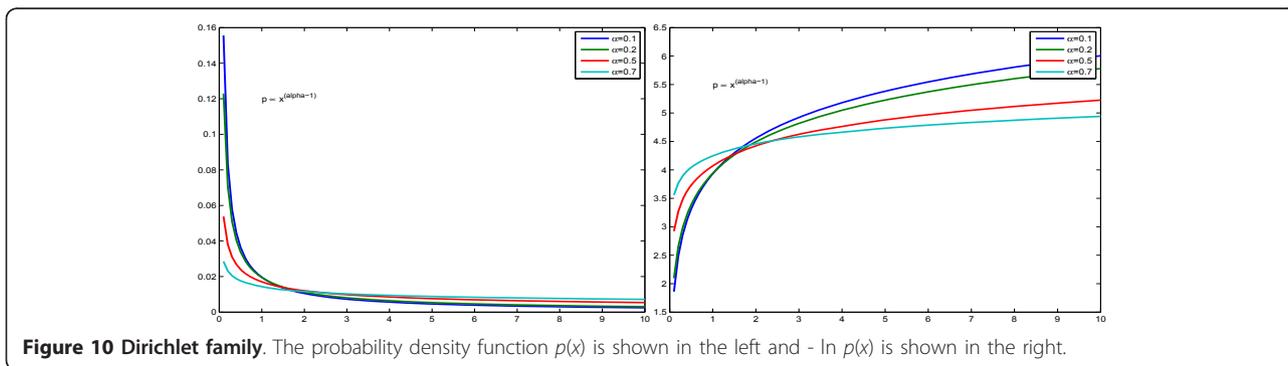


Figure 10 Dirichlet family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

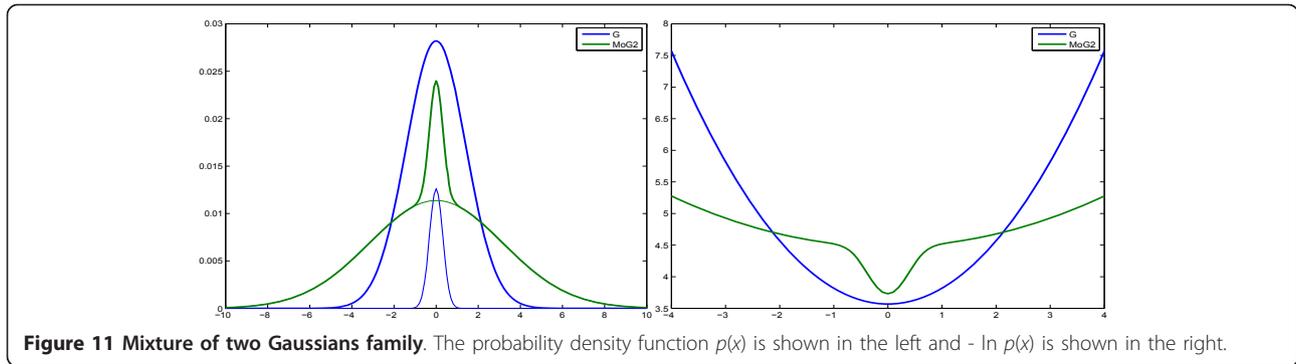


Figure 11 Mixture of two Gaussians family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

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

$$\begin{cases} p(\mathbf{f} | \mathbf{z}) = \prod_j p(f_j | z_j) = \prod_j \mathcal{N}(f_j | 0, v_{z_j}) \\ \propto \exp \left\{ -\frac{1}{2} \sum_j v_{z_j} f_j^2 \right\} \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (29)$$

In general $v_1 \gg v_0$ and λ measures the sparsity ($0 < \lambda \ll 1$). This family of models are illustrated on Figure 11.

2.8 Bernoulli-Gaussian (BG) model

The Bernoulli-Gaussian model can be considered as the particular case of the MoG2 with the particular degenerate case of $v_0 = 0$:

$$p(\mathbf{f} | \lambda, v) = \prod_j p(f_j) = \prod_j (\lambda \mathcal{N}(f_j | 0, v) + (1 - \lambda) \delta(f_j)) \quad (30)$$

which can also be written as

$$\begin{cases} p(\mathbf{f} | \mathbf{z}) = \prod_j p(f_j | z_j) = \prod_j [\mathcal{N}(f_j | 0, v)]^{\delta(z_j)} \prod_j [\delta(f_j)]^{\delta(1-z_j)} \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (31)$$

This model has also been called *spike and slab*. This family of models are illustrated on Figure 12.

2.9 Mixture of three Gaussians (MoG3) model

Another mixture model proposed is using a Mixture of three Gaussians, one centered at zero and two symmetrically placed:

$$p(\mathbf{f} | \lambda, v_0, v_{+1}, v_{-1}, \beta) = \prod_j [(1 - \lambda) \mathcal{N}(f_j | 0, v_0) + (\lambda/2) \mathcal{N}(f_j | +\beta, v_{+1}) + (\lambda/2) \mathcal{N}(f_j | -\beta, v_{-1})] \quad (32)$$

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

$$\begin{cases} p(\mathbf{f} | \mathbf{z}) = \prod_j p(f_j | z_j) = \prod_j \mathcal{N}(f_j | z_j \beta, v_{z_j}) \\ P(z_j = 1) = \lambda/2, \\ P(z_j = -1) = \lambda/2, \\ P(z_j = 0) = 1 - \lambda. \end{cases} \quad (33)$$

In general $v_{+1} = v_{-1} = v \gg v_0$ and λ measures the sparsity ($0 < \lambda \ll 1$). This family of models are illustrated on Figure 13.

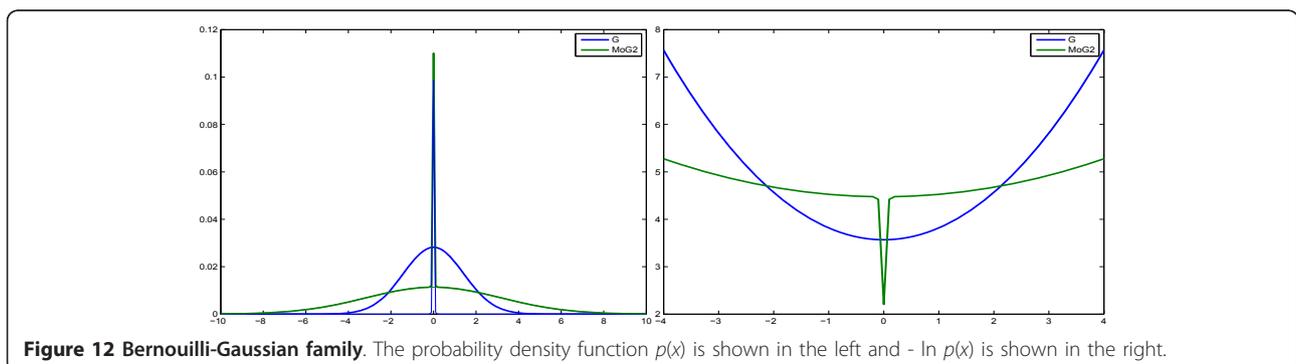


Figure 12 Bernoulli-Gaussian family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

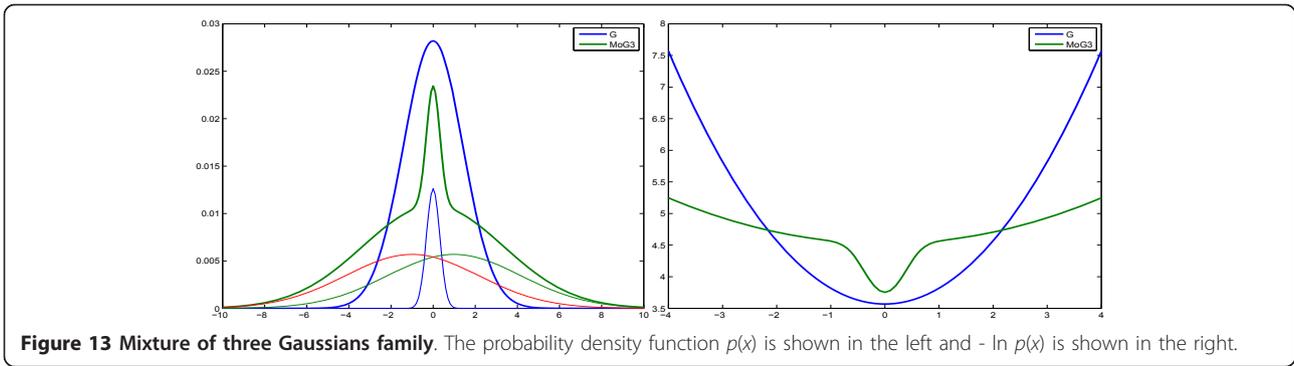


Figure 13 Mixture of three Gaussians family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

2.10 Mixture of one Gaussian and two Gammas (MoGGammas) model

Another mixture model proposed is using a mixture of one central Gaussian and two symmetric Gammas:

$$p(f | \lambda, v_0, \alpha, \beta) = \prod_j [(1 - \lambda)\mathcal{N}(f_j | 0, v_0) + (\lambda/2)\mathcal{G}(f_j | \alpha, \beta) + (\lambda/2)\mathcal{G}(-f_j | \alpha, \beta)] \quad (34)$$

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

$$\begin{cases} p(f | z) = \prod_j (f_j | z_j) = [\mathcal{N}(f_j | 0, v_0)]^{\sum_j \delta(z_j)} \times \\ \quad [\mathcal{G}(f_j | \alpha, \beta)]^{\sum_j \delta(z_j - 1)} \times \\ \quad [\mathcal{G}(-f_j | \alpha, \beta)]^{\sum_j \delta(z_j + 1)} \\ P(z_j = 1) = \lambda/2, \\ P(z_j = -1) = \lambda/2, \\ P(z_j = 0) = 1 - \lambda. \end{cases} \quad (35)$$

This family of models are illustrated on Figure 14.

2.11 Bernoulli-Gamma (BGamma) model

As in the BG model, when we want to enforce both sparsity and positivity, we can use the BGamma model:

$$p(f | \lambda, \alpha, \beta) = \prod_j [\lambda \delta(f_j) + (1 - \lambda)\mathcal{G}(f_j | \alpha, \beta)] \quad (36)$$

or

$$\begin{cases} p(f | z) = \prod_j p(f_j | z_j) = \prod_j [z_j \mathcal{G}(f_j | \alpha, \beta) \\ \quad \prod_j [(1 - z_j) \delta(f_j)] \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (37)$$

A particular case of this model is Bernoulli-exponential (BExponential) which obtained when $\alpha = 1$. These families of models are illustrated on Figure 15 and Figure 16.

2.12 Mixture of Dirichlet (MoD) model

• Mixture of Dirichlet model

$$p(f | \lambda, \alpha_1, \alpha_2) = \lambda \mathcal{D}(f | \alpha_1) + (1 - \lambda) \mathcal{D}(f | \alpha_2) \quad (38)$$

where

$$\mathcal{D}(f | \alpha) \propto \prod_j f_j^{\alpha - 1} \text{ with } f_j > 0, \sum_j f_j = 1 \quad (39)$$

is the symmetric Dirichlet distribution. We need to choose $\alpha_1 > 1$ for dense part and $0 < \alpha_2 < 1$ for the sparse part.

2.13 Bernoulli-multinomial (BMultinomial) model

As in the BG or BGamma model, when we know that the signal is sparse and can only take one of the K

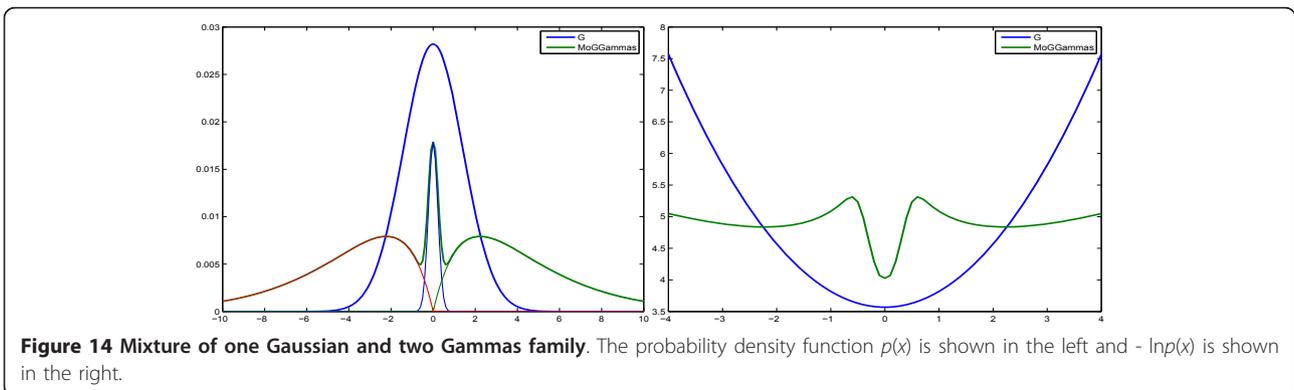


Figure 14 Mixture of one Gaussian and two Gammas family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

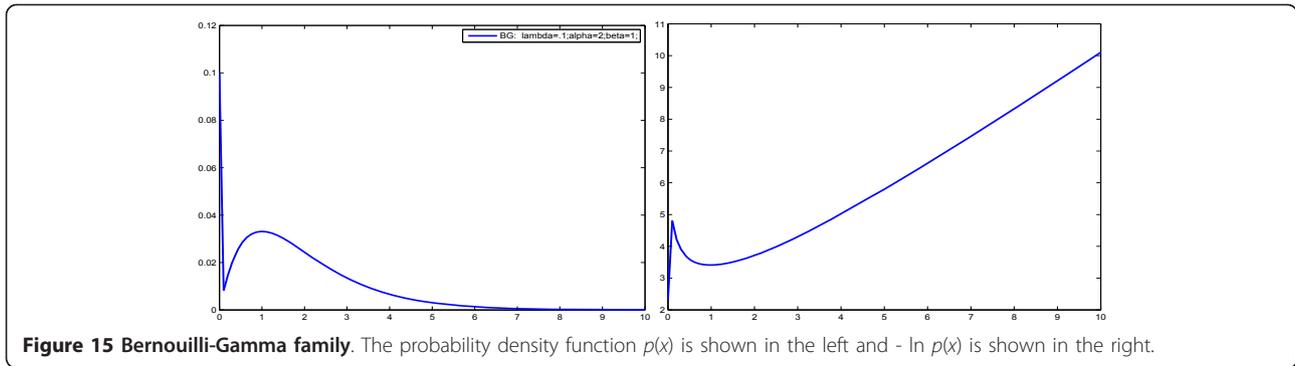


Figure 15 Bernoulli-Gamma family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

discrete values $\{a_1, \dots, a_K\}$, we can use the BMultinomial model:

$$p(\mathbf{f} | \lambda, \mathbf{a}, \boldsymbol{\alpha}) = \prod_j \lambda \mathcal{M}ult(f_j | \mathbf{a}, \boldsymbol{\alpha}) + (1 - \lambda) \delta(f_j) \quad (40)$$

where $\mathbf{a} = \{a_1, \dots, a_K\}$ and $\boldsymbol{\alpha} = \{\alpha_1, \dots, \alpha_K\}$ with $\sum_k \alpha_k = 1$ and

$$\mathcal{M}ult(f_j | \mathbf{a}, \boldsymbol{\alpha}) = \frac{n!}{a_1! \dots a_K!} \prod_k \alpha_k^{a_j}$$

or

$$\begin{cases} p(\mathbf{f} | \mathbf{z}) = \prod_j p(f_j | z_j) = \prod_j [z_j \mathcal{M}ult(f_j | \boldsymbol{\alpha}) \\ \prod_j [(1 - z_j) \delta(f_j)] \\ P(z_j = 1) = \lambda, & P(z_j = 0) = 1 - \lambda \end{cases} \quad (41)$$

3 Bayesian inference with sparsity enforcing priors

The priors proposed can be used in a Bayesian approach to infer on \mathbf{f} given the observed data \mathbf{g} through the posterior law given in Equation (3). First let assume the error ϵ to be centered, Gaussian and white: $\epsilon \sim \mathcal{N}(\epsilon | \mathbf{0}, v_\epsilon \mathbf{I})$. Then, using the forward model (1) we have

$$p(\mathbf{g} | \mathbf{f}) = \mathcal{N}(\mathbf{H}\mathbf{f}, v_\epsilon \mathbf{I}) \propto \exp \left\{ -\frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 \right\} \quad (42)$$

Now, we consider different priors.

3.1 Simple prior models

Given $p(\mathbf{g} | \mathbf{f})$ and any simple prior law $p(\mathbf{f})$, the posterior law is written:

$$p(\mathbf{f} | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}) p(\mathbf{f}) \propto \exp\{J(\mathbf{f})\} \quad (43)$$

with

$$J(\mathbf{f}) = \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \Omega(\mathbf{f}) \quad (44)$$

where $\Omega(\mathbf{f}) = -\ln p(\mathbf{f})$ and so the Maximum A Posteriori (MAP) solution is expressed as the minimizer of this criterion which has two parts: the first part is due to the likelihood and the second part is due to the prior:

$$p(\mathbf{f} | \boldsymbol{\theta}, \mathbf{g}) \rightarrow \left[\text{Optimization of } J(\mathbf{f}) = \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \Omega(\mathbf{f}) \right] \rightarrow \hat{\mathbf{f}}$$

Thus, depending on the choice of the prior we obtain different expressions for $\Omega(\mathbf{f})$. For example for the GG model of (10) we get

$$\Omega(\mathbf{f}) = \gamma \sum_j |f_j|^\beta. \quad (45)$$

For the symmetric Weibull model (14) we get

$$\Omega(\mathbf{f}) = -\gamma \sum_j |f_j|^\beta + (\beta - 1) \log |f_j|. \quad (46)$$

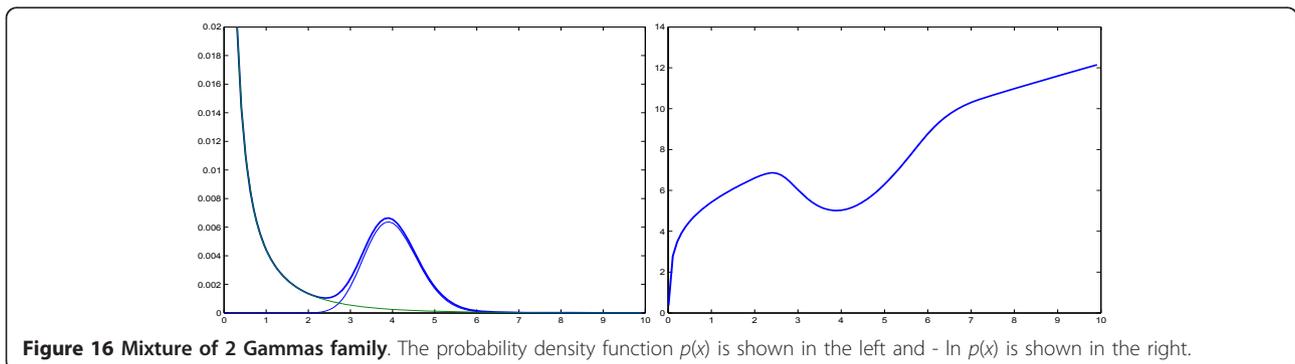


Figure 16 Mixture of 2 Gammas family. The probability density function $p(x)$ is shown in the left and $-\ln p(x)$ is shown in the right.

For the Student-t model (17) we get

$$\Omega(\mathbf{f}) = \frac{\nu + 1}{2} \sum_j \log(1 + f_j^2/\nu). \quad (47)$$

For the elastic net model we get

$$\Omega(\mathbf{f}) = \sum_j [\gamma_1 |f_j| + \gamma_2 f_j^2] \quad (48)$$

and for the Dirichlet model we get

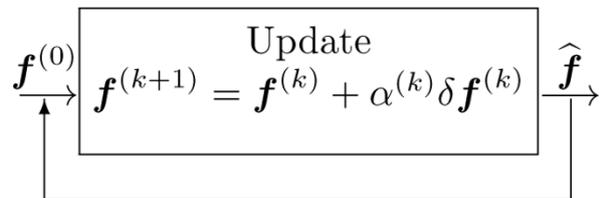
$$\Omega(\mathbf{f}) = \sum_j f_j^{\alpha-1}, \quad f_j > 0, \quad \sum_j f_j = 1. \quad (49)$$

For each of these cases, we may discuss on the unimodality and convexity of the criterion $J(\mathbf{f})$ which depends mainly on its Hessian

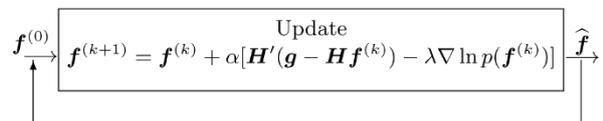
$$\begin{aligned} \Delta J(\mathbf{f}) &= \left[\frac{\partial^2 J(\mathbf{f})}{\partial f_j \partial f_i} \right] = \mathbf{H}'\mathbf{H} + \left[\frac{\partial^2 \Omega(\mathbf{f})}{\partial f_i \partial f_j} \right] \\ &= \mathbf{H}'\mathbf{H} + \left[\frac{\partial^2 \Omega(\mathbf{f})}{\partial f_j^2} \right] \end{aligned} \quad (50)$$

We may look at each case to examine the range of the parameters for which this Hessian matrix is positive definite.

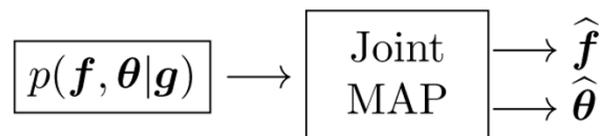
The optimization is done iteratively:



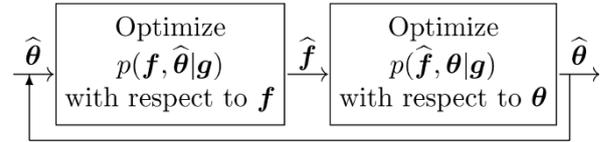
Update operation can be additive, multiplicative or more complex. Updating steps $\alpha^{(k)}$ can be fixed or computed adaptively at each step (steepest descent for example). $\delta \mathbf{f}^{(k)}$ can be, for example proportional to the gradient, in which case, we have



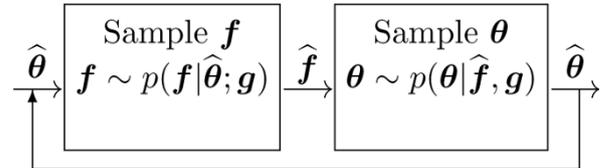
We may also consider to estimate some of these parameters by assigning them appropriate priors and then express the joint $p(\mathbf{f}, \boldsymbol{\theta} | \mathbf{g}, \boldsymbol{\theta}_0)$ as given in Equation (4) and then try to estimate them jointly, for example joint MAP:



or alternate optimization:

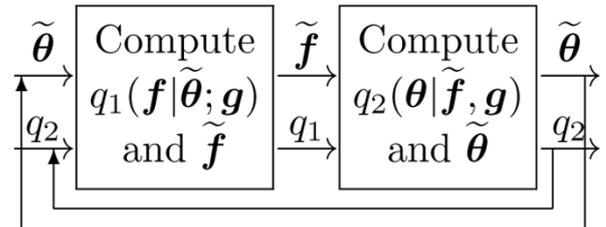


We may also want to explore this joint posterior by generating samples from it. This can be done, for example, through the following Gibbs sampling scheme:



When a great number of samples are thus generated, we may compute their means, variances or any other statistics about them.

Finally, we may try to approximate this joint posterior by a simpler one, for example by a separable $q(\mathbf{f}, \boldsymbol{\theta}) = q_1(\mathbf{f}) q_2(\boldsymbol{\theta})$ using the variational approximation (VA). The main idea and the main basic steps to achieve this is more detailed in the following section. Here, however, we present the result on the following scheme:



To illustrate the differences, we may consider the simple case of a linear forward model and Gaussian priors:

$$\begin{cases} p(\mathbf{g} | \mathbf{f}, v_\varepsilon) = \mathcal{N}(\mathbf{H}\mathbf{f}, v_\varepsilon \mathbf{I}) \\ p(\mathbf{f} | v_f) = \mathcal{N}(0, v_f \mathbf{I}) \end{cases} \quad (51)$$

In this case, if we know $\boldsymbol{\theta} = (v_\varepsilon, v_f)$, then

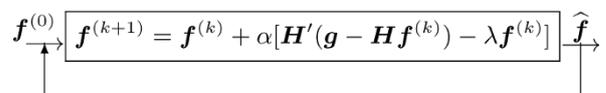
$$p(\mathbf{f} | \mathbf{g}, v_\varepsilon, v_f) = \mathcal{N}(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) \quad (52)$$

with:

$$\hat{\boldsymbol{\mu}} = [\mathbf{H}'\mathbf{H} + \lambda \mathbf{I}]^{-1} \mathbf{H}'\mathbf{g}$$

$$\hat{\boldsymbol{\Sigma}} = [\mathbf{H}'\mathbf{H} + \lambda \mathbf{I}]^{-1}$$

with $\lambda = \frac{v_\varepsilon}{v_f}$. So, we have $\hat{\mathbf{f}} = \hat{\boldsymbol{\mu}}$ which can be computed by optimizing $J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \lambda \|\mathbf{f}\|^2$. A gradient based algorithm is shown below:



Now putting inverse Gamma priors on v_ϵ and v_f or equivalently Gamma priors on $\tau_\epsilon = 1/v_\epsilon$ and $\tau_f = 1/v_f$:

$$\begin{cases} p(\tau_\epsilon | \alpha_{\tau 0}, \beta_{\tau 0}) = \mathcal{G}(\alpha_{\tau 0}, \beta_{\tau 0}) \\ p(\tau_f | \alpha_{\tau f 0}, \beta_{\tau f 0}) = \mathcal{G}(\alpha_{\tau f 0}, \beta_{\tau f 0}) \end{cases} \quad (53)$$

we have

$$\begin{cases} p(\tau_\epsilon | f, \mathbf{g}, \alpha_{\tau 0}, \beta_{\tau 0}) = \mathcal{G}(\hat{\alpha}_\epsilon, \hat{\beta}_\epsilon) \\ p(\tau_f | f, \alpha_{\tau f 0}, \beta_{\tau f 0}) = \mathcal{G}(\hat{\alpha}_f, \hat{\beta}_f) \end{cases} \quad (54)$$

with

$$\begin{cases} \hat{\alpha}_\epsilon = \alpha_{\epsilon 0} + 1/2 \\ \hat{\beta}_\epsilon = \beta_{\epsilon 0} + \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2/2 \\ \hat{\tau}_\epsilon = \hat{\alpha}_\epsilon/\hat{\beta}_\epsilon = \frac{2\alpha_{\epsilon 0} + 1}{2\beta_{\epsilon 0} + \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2} \\ \hat{\alpha}_f = \alpha_{f 0} + 1/2 \\ \hat{\beta}_f = \beta_{f 0} + \|\mathbf{f}\|^2/2 \\ \hat{\tau}_f = \hat{\beta}_f/\hat{\alpha}_f = \frac{2\alpha_{f 0} + 1}{2\beta_{f 0} + \|\mathbf{f}\|^2} \end{cases} \quad (55)$$

and $\hat{\lambda} = \frac{\hat{\tau}_f}{\hat{\tau}_\epsilon}$. Then, the alternate optimization of the JMAP estimate algorithm becomes

$$\begin{array}{c} \lambda \rightarrow \left[\begin{array}{l} \Sigma = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1} \\ \mathbf{f} = \mu = \Sigma\mathbf{H}'\mathbf{g} \end{array} \right] \xrightarrow{\mathbf{f}} \left[\begin{array}{l} \tau_\epsilon = \frac{2\alpha_{\epsilon 0} + 1}{2\beta_{\epsilon 0} + \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2} \\ \tau_f = \frac{2\alpha_{f 0} + 1}{2\beta_{f 0} + \|\mathbf{f}\|^2} \end{array} \right] \xrightarrow{\lambda} \frac{\tau_f}{\tau_\epsilon} \end{array}$$

The Gibbs sampling algorithm becomes

$$\begin{array}{c} \lambda \rightarrow \left[\begin{array}{l} \mathbf{f} \sim \mathcal{N}(\mu, \Sigma) \\ \Sigma = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1} \\ \mu = \Sigma\mathbf{H}'\mathbf{g} \end{array} \right] \xrightarrow{\mathbf{f}} \left[\begin{array}{l} \tau_\epsilon \sim \mathcal{G}(\hat{\alpha}_\epsilon, \hat{\beta}_\epsilon) \\ \tau_f \sim \mathcal{G}(\hat{\alpha}_f, \hat{\beta}_f) \end{array} \right] \xrightarrow{\lambda} \frac{\tau_f}{\tau_\epsilon} \end{array}$$

The VBA algorithm becomes

$$\begin{cases} q(\mathbf{f}) = \mathcal{N}(\hat{\mu}, \hat{\Sigma}) \\ \hat{\Sigma} = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1} \\ \hat{\mu} = \hat{\Sigma}\mathbf{H}'\mathbf{g} \\ q(\tau_\epsilon) = \mathcal{G}(\hat{\alpha}_\epsilon, \hat{\beta}_\epsilon) \\ \hat{\alpha}_\epsilon = \alpha_{\epsilon 0} + 1/2 \\ \hat{\beta}_\epsilon = \beta_{\epsilon 0} + \|\mathbf{g} - \mathbf{H}\langle \mathbf{f} \rangle\|^2/2 \\ q(\tau_f) = \mathcal{G}(\hat{\alpha}_f, \hat{\beta}_f) \\ \hat{\alpha}_f = \alpha_{f 0} + 1/2 \\ \hat{\beta}_f = \beta_{f 0} + \frac{1}{2} \|\langle \mathbf{f}^2 \rangle\| = \beta_{f 0} + \frac{1}{2} \sum_j \langle f_j^2 \rangle \\ \langle \mathbf{f} \rangle = \hat{\mu} \\ \langle \mathbf{f}^2 \rangle = \hat{\mu}^2 + \text{diag}[\hat{\Sigma}] \end{cases} \quad (56)$$

and $\hat{\lambda} = \frac{\hat{\tau}_f}{\hat{\tau}_\epsilon}$:

$$\begin{array}{c} \lambda \rightarrow \left[\begin{array}{l} q(\mathbf{f}) = \mathcal{N}(\hat{\mu}, \hat{\Sigma}) \\ \hat{\Sigma} = [\mathbf{H}'\mathbf{H} + \lambda\mathbf{I}]^{-1} \\ \hat{\mu} = \hat{\Sigma}\mathbf{H}'\mathbf{g} \end{array} \right] \xrightarrow{q(\mathbf{f})} \left[\begin{array}{l} q(\tau_\epsilon) = \mathcal{G}(\hat{\alpha}_\epsilon, \hat{\beta}_\epsilon) \\ q(\tau_f) = \mathcal{G}(\hat{\alpha}_f, \hat{\beta}_f) \end{array} \right] \xrightarrow{\lambda} \frac{\tau_f}{\tau_\epsilon} \end{array}$$

We recently implemented these algorithms for different applications such as: synthetic aperture radar (SAR) Imaging [29], ...

3.2 Mixture models

For the mixture models, and in general for the models which can be expressed via the hidden variables, we want to estimate jointly the original unknowns \mathbf{f} and the hidden variables: τ in Cauchy model, \mathbf{z} in MoG2, BG or BGam models and \mathbf{z} in MoG3 or MoGGammas. Let examine these a little in details.

3.3 Student-t and Cauchy models

In this case the joint prior law can be written as:

$$\begin{aligned} p(\mathbf{f}, \tau) &= \prod_j p(f_j | \tau_j) p(\tau_j) = \prod_j \mathcal{N}(f_j | 0, 1/\tau_j) p(\tau_j) \\ &\propto \exp \left\{ -\frac{1}{2} \sum_j \tau_j f_j^2 + a \ln \tau_j - b \tau_j \right\} \text{ with } a = b = v/2 \end{aligned} \quad (57)$$

such that

$$p(\mathbf{f}, \tau | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}) p(\mathbf{f}, \tau) \propto \exp\{-J(\mathbf{f}, \tau)\} \quad (58)$$

where

$$j(\mathbf{f}, \tau) = \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \sum_j \frac{1}{2} \tau_j f_j^2 - a \ln \tau_j + b \tau_j \quad (59)$$

$$p(\mathbf{f}, \tau | \mathbf{g}) \rightarrow \left[\begin{array}{c} \text{Joint Optimization of} \\ J(\mathbf{f}, \tau) \end{array} \right] \rightarrow \begin{array}{l} \hat{\mathbf{f}} \\ \hat{\tau} \end{array}$$

Joint optimization of this criterion, alternatively with respect to \mathbf{f} (with fixed τ)

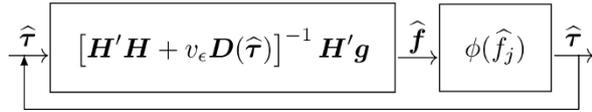
$$\begin{aligned} \hat{\mathbf{f}} &= \arg \min_{\mathbf{f}} \{J(\mathbf{f}, \tau)\} \\ &= \arg \min_{\mathbf{f}} \left\{ \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}\mathbf{f}\|^2 + \sum_j \frac{1}{2} \tau_j f_j^2 \right\} \end{aligned} \quad (60)$$

and with respect to τ (with fixed \mathbf{f})

$$\begin{aligned} \hat{\tau} &= \arg \min_{\tau} \{J(\mathbf{f}, \tau)\} \\ &= \arg \min_{\tau} \left\{ \sum_j \frac{1}{2} \tau_j f_j^2 - a \ln \tau_j + b \tau_j \right\} \end{aligned} \quad (61)$$

results in the following iterative algorithm:

$$\begin{cases} \hat{\mathbf{f}} = [\mathbf{H}'\mathbf{H} + v_\epsilon \mathbf{D}(\hat{\boldsymbol{\tau}})]^{-1} \mathbf{H}'\mathbf{g} \\ \hat{\tau}_j = \phi(\hat{f}_j) = \frac{a}{f_j^2 + b} \\ \mathbf{D}(\hat{\boldsymbol{\tau}}) = \text{diag}[1/\hat{\tau}_j, j = 1, \dots, n] \end{cases} \quad (62)$$



Note that, τ_j is the inverse of a variance and we have $1/\tau_j = \frac{f_j^2 + b}{a}$. We can interpret this as an iterative quadratic regularization inversion followed by the estimation of variances τ_j which are used in the next iteration to define the variance matrix $\mathbf{D}(\boldsymbol{\tau})$.

Here too, we may study the conditions on which the joint criterion is uni-modal and its alternate optimization converges to its unique solution.

We may also consider a Gibbs sampling scheme

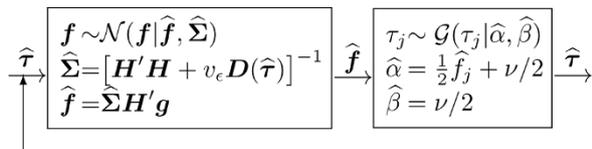
$$\begin{cases} f \sim p(f | \boldsymbol{\tau}, \mathbf{g}) \propto p(\mathbf{g} | f) p(f | \mathbf{u}) = \mathcal{N}(f | \hat{\mathbf{f}}, \hat{\boldsymbol{\Sigma}}) \\ \boldsymbol{\tau} \sim p(\boldsymbol{\tau} | f, \mathbf{g}) \propto p(f | \boldsymbol{\tau}) p(\boldsymbol{\tau}) = \prod_j \mathcal{G}(\tau_j | \hat{\alpha}, \hat{\beta}) \end{cases} \quad (63)$$

where

$$\begin{cases} \hat{\boldsymbol{\Sigma}} = [\mathbf{H}'\mathbf{H} + v_\epsilon \mathbf{D}(\boldsymbol{\tau})]^{-1} \\ \hat{\mathbf{f}} = \hat{\boldsymbol{\Sigma}} \mathbf{H}'\mathbf{g} \end{cases} \quad (64)$$

and

$$\begin{cases} \hat{\alpha} = \frac{1}{2} \hat{f}_j + a = \frac{1}{2} \hat{f}_j + \nu/2 \\ \hat{\beta} = b = \nu/2 \end{cases} \quad (65)$$

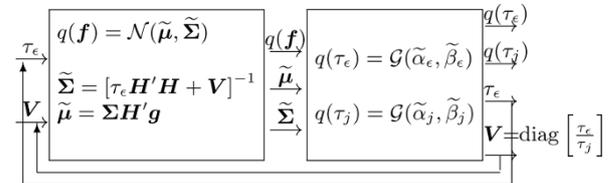


For the VBA, we have

$$\begin{cases} p(\mathbf{g} | f, v_\epsilon) = \mathcal{N}(\mathbf{g} | \mathbf{H}f, v_\epsilon \mathbf{I}), \quad \tau_\epsilon = 1/v_\epsilon \\ p(\tau_\epsilon) = \mathcal{G}(\tau_\epsilon | \alpha_{\epsilon 0}, \beta_{\epsilon 0}) \\ p(f | \mathbf{v}) = \prod_j p(f_j | v_j) = \prod_j \mathcal{N}(f_j | 0, v_j) = \mathcal{N}(f | 0, \mathbf{V}) \\ \mathbf{V} = \text{diag}[v], \quad \tau_j = 1/v_j, \quad \boldsymbol{\tau} = \text{diag}[\boldsymbol{\tau}] = \mathbf{V}^{-1} \\ p(\boldsymbol{\tau}) = \prod_j \mathcal{G}(\tau_j | \alpha_0, \beta_0) \end{cases} \quad (66)$$

$$\begin{cases} \tilde{q}(f) = \mathcal{N}(f | \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \\ \tilde{\boldsymbol{\mu}} = \tilde{\boldsymbol{\Sigma}} \mathbf{H}'\mathbf{g} \\ \tilde{\boldsymbol{\Sigma}} = (\tilde{\boldsymbol{\tau}}_\epsilon \mathbf{H}'\mathbf{H} + \tilde{\mathbf{V}})^{-1}, \text{ with } \tilde{\mathbf{V}} = \text{diag}[\tilde{v}] \end{cases} \quad (67)$$

$$\begin{cases} \tilde{q}(\tau_\epsilon) = \mathcal{G}(\tau_\epsilon | \tilde{\alpha}_\epsilon, \tilde{\beta}_\epsilon), \\ \tilde{\alpha}_\epsilon = \alpha_{\epsilon 0} + (n+1)/2 \\ \tilde{\beta}_\epsilon = \beta_{\epsilon 0} + 1/2 \\ \tilde{\boldsymbol{\tau}} = \tilde{\boldsymbol{\alpha}}_\tau / \tilde{\beta}_\tau \\ \tilde{q}(\tau_j) = \mathcal{G}(\tau_j | \tilde{\alpha}_j, \tilde{\beta}_j) \\ \tilde{\alpha}_j = \alpha_{00} + 1/2 \\ \tilde{\beta}_j = \beta_{00} + \langle f_j^2 \rangle / 2 \\ \tilde{z}_j = \tilde{\beta}_j / \tilde{\alpha}_j \end{cases} \quad (68)$$



3.4 Mixture of two Gaussians (MoG2) model

In this case, following the same arguments, we obtain:

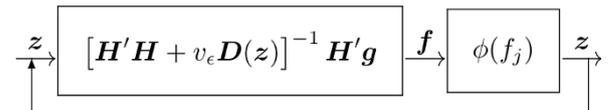
$$p(\mathbf{f}, \mathbf{z} | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}) p(\mathbf{f}, \mathbf{z}) \propto \exp\{-J(\mathbf{f}, \mathbf{z})\} \quad (69)$$

where

$$\begin{aligned} J(\mathbf{f}, \mathbf{z}) &= \frac{1}{2v_\epsilon} \|\mathbf{g} - \mathbf{H}f\|^2 \\ &+ \sum_j \frac{f_j^2}{2v_{z_j}} + z_j \ln \lambda + (1 - z_j) \ln(1 - \lambda) \end{aligned} \quad (70)$$

Again, in this case also, the optimization of this criterion, alternatively with respect to \mathbf{f} and \mathbf{z} results in the following iterative algorithm:

$$\begin{cases} \hat{\mathbf{f}} = [\mathbf{H}'\mathbf{H} + v_\epsilon \mathbf{D}(\mathbf{z})]^{-1} \mathbf{H}'\mathbf{g} \\ \hat{z}_j = \phi(f_j) = \begin{cases} 1, & \text{if } f_j^2 \geq (v_1 - v_0) \ln \frac{1 - \lambda}{\lambda} \\ 0, & \text{if } f_j^2 < (v_1 - v_0) \ln \frac{1 - \lambda}{\lambda} \end{cases} \\ \mathbf{D}(\hat{\mathbf{z}}) = \text{diag}[v_{\hat{z}_j}, j = 1, \dots, n] \end{cases} \quad (71)$$



Here too, we may also consider a Gibbs sampling scheme

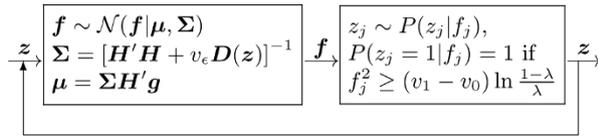
$$\begin{cases} f \sim p(f | \mathbf{z}, \mathbf{g}) \propto p(\mathbf{g} | f) p(f | \mathbf{u}) = \mathcal{N}(f | \hat{\mathbf{f}}, \hat{\boldsymbol{\Sigma}}) \\ \mathbf{z} \sim p(\mathbf{z} | f, \mathbf{g}) \propto p(f | \mathbf{z}) p(\mathbf{z}) = \prod_j P(z_j = k | f_j) \end{cases} \quad (72)$$

where

$$\begin{cases} \hat{\boldsymbol{\Sigma}} = [\mathbf{H}'\mathbf{H} + v_\epsilon \mathbf{D}(\mathbf{z})]^{-1} \\ \hat{\mathbf{f}} = \hat{\boldsymbol{\Sigma}} \mathbf{H}'\mathbf{g} \end{cases} \quad (73)$$

and

$$\begin{cases} P(z_j = 1 | f_j) = 1, & \text{if } f_j^2 \geq (v_1 - v_0) \ln \frac{1 - \lambda}{\lambda} \\ P(z_j = 0 | f_j) = 1, & \text{if } f_j^2 < (v_1 - v_0) \ln \frac{1 - \lambda}{\lambda} \end{cases} \quad (74)$$



3.5 BG model

For the case of BG we have to be more careful, because the joint probability laws are degenerated. Two approaches are then possible:

i) Considering them as the particular case of the MoG models where the variance v_0 is fixed to a small value or reduced gradually during the iterations.

ii) Trying first to integrate out f from the expression of $p(f, z|g)$ to obtain $p(z|g)$ and optimize it with respect to z (detection step) and then use it for the estimation step.

To go further in detail of the second approach, we may remark that for the given z , the expression of $p(f, z|g)$ as a function of f is Gaussian and so it can be easily integrated out and we obtain:

$$\begin{aligned} p(z|g) &\propto p(g|z)p(z) \\ &\propto \mathcal{N}(g|0, H(\text{vdiag}[z_j, j=1, \dots, n])H' + v_e I) \lambda^{\sum_j z_j} (1-\lambda)^{\sum_j (1-z_j)} \end{aligned} \quad (75)$$

Now writing the expression of $\mathcal{L}(z) = -\ln p(z|g)$ and keeping only all terms depending on z we obtain:

$$\mathcal{L}(z) = -g'B^{-1}(z)g - \ln |B(z)| - 2n \ln \frac{1-\lambda}{\lambda} \quad (76)$$

where $B(z) = H(\text{vdiag}[z_j, j=1, \dots, n])H' + v_e I$. We see the complexity of this expression which needs the inversion of the matrix B and its optimization which is a combinatorial optimization needing to evaluate this expression 2^n times.

However, we may also remark that when z obtained, the estimation of f is easy. We have:

$$\hat{f} = H'B^{-1}g. \quad (77)$$

which needs again the inversion of the matrix B .

The exact computations of \hat{z} and \hat{f} are often too costly, one may try to obtain approximate solutions. Many approximations have been proposed. A good overview of these methods can be found in [30, Chap. 5] and also in [31,32].

3.6 BGamma and MoGGammas model

In these cases, it is no more possible to integrate out f analytically as it was the case with Gaussians. One strategy here is to use the MCMC methods to generate samples from the joint posterior. The second approach is to approximate the joint posterior by a simpler one, for example by a separable one on f and the hidden variables z in the BGamma or the MoGGammas cases. Very often then we can do the computations analytically. However, it may happens that, even after these separable approximations, still we need to use the MCMC methods on some of variables. Detailed explanation of these general methods is out of focus of this article. See [30,33,34]. Here, we just give the details for the case of the Gaussian mixtures (MoG2 or MoG3).

4 Variational Bayesian approximation for the case of mixture laws

To start and to be complete as to propose an unsupervised method, we include also the estimation of the parameters θ and write the joint posterior law of all the unknowns:

$$p(f, z, \theta | g) \propto p(g|f, \theta)p(f|z, \theta)p(z|\theta)p(\theta) \quad (78)$$

which can also be written as

$$q(f, z, \theta | g) = p(f|z, \theta; g)p(z|\theta; g)p(\theta | g) \quad (79)$$

where

$$p(f|z, \theta; g) = p(g|f, \theta)p(f|z, \theta)/p(g|z, \theta) \quad (80)$$

with

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

and

$$p(z|\theta; g) = p(g|z, \theta)p(z|\theta)/p(g|\theta) \quad (81)$$

with

$$p(g|\theta) = \int p(g|z, \theta)p(z|\theta) dz$$

or

$$p(g|\theta) = \sum_z p(g|z, \theta)p(z|\theta)$$

when z are discrete valued, and finally

$$p(\theta | g) = p(g|\theta)p(\theta)/p(g) \quad (82)$$

with

$$p(g) = \int p(g|\theta)p(\theta) d\theta$$

One can also write:

$$p(\mathbf{z} | \boldsymbol{\theta}, \mathbf{g}) = \int p(\mathbf{f}, \mathbf{z} | \boldsymbol{\theta}, \mathbf{g}) d\mathbf{f} \quad (83)$$

and

$$p(\boldsymbol{\theta} | \mathbf{g}) = \int \int p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) d\mathbf{f} d\mathbf{z} = \int p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{g}) d\mathbf{z} \quad (84)$$

or

$$p(\boldsymbol{\theta} | \mathbf{g}) = \sum_{\mathbf{z}} \int p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) d\mathbf{f} = \sum_{\mathbf{z}} p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{g}) \quad (85)$$

when \mathbf{z} are discrete valued.

We see that the first term

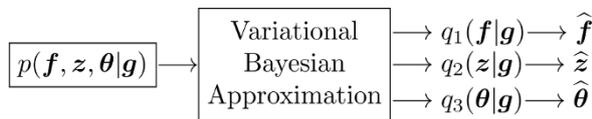
$$p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) \quad (86)$$

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

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \mathbf{g}) q_2(\mathbf{z} | \mathbf{g}) q_3(\boldsymbol{\theta} | \mathbf{g}) \quad (87)$$

illustrated here:



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

$$\text{KL}(q : p) = \int q \ln \frac{q}{p} = \left\langle \ln \frac{q}{p} \right\rangle_q \quad (88)$$

It is then easy to show that

$$\text{KL}(q : p) = \ln p(\mathbf{g} | \mathcal{M}) - \mathcal{F}(q) \quad (89)$$

where $p(\mathbf{g} | \mathcal{M})$ is the likelihood of the model

$$p(\mathbf{g} | \mathcal{M}) = \int \int \int p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) d\mathbf{f} d\mathbf{z} d\boldsymbol{\theta} \quad (90)$$

with

$$p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) = p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) p(\mathbf{z} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (91)$$

and $\mathcal{F}(q)$ is the free energy associated to q defined as

$$\mathcal{F}(q) = \left\langle \ln \frac{p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta})} \right\rangle_q \quad (92)$$

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

$$q_1(\mathbf{f}) \propto \exp \left\{ -\langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{z})q(\boldsymbol{\theta})} \right\}$$

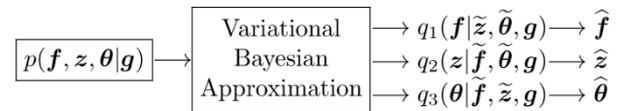
$$q_2(\mathbf{z}) \propto \exp \left\{ -\langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{f})q(\boldsymbol{\theta})} \right\}$$

$$q_3(\boldsymbol{\theta}) \propto \exp \left\{ -\langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q(\mathbf{f})q(\mathbf{z})} \right\}$$

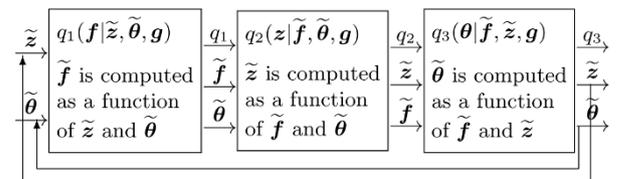
Note that these relations represent an implicit solution for $q_1(\mathbf{f})$, $q_2(\mathbf{z})$, and $q_3(\boldsymbol{\theta})$ which need, at each iteration, the expression of the expectations in the right hand of exponentials. If $p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta}_1)$ is a member of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

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) q_2(\mathbf{z} | \tilde{\mathbf{f}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) q_3(\boldsymbol{\theta} | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}; \mathbf{g}) \quad (93)$$

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 results to 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

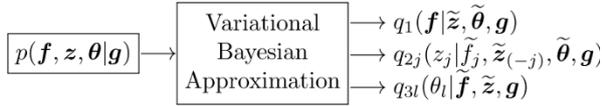
$$\begin{aligned} \mathcal{F}(q) &= \left\langle \ln \frac{p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta})} \right\rangle_q \\ &= \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) | \mathcal{M} \rangle_q + \langle -\ln q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q \quad (94) \\ &= \langle \ln p(\mathbf{g} | \mathbf{f}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{z} | \boldsymbol{\theta}) \rangle_q \\ &\quad + \langle -\ln q(\mathbf{f}) \rangle_q + \langle -\ln q(\mathbf{z}) \rangle_q + \langle -\ln q(\boldsymbol{\theta}) \rangle_q \end{aligned}$$

where all the expectations are with respect to q .

Other decompositions are also possible:

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) \prod_j q_{2j}(z_j | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}_{(-j)}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) \prod_l q_{3l}(\theta_l | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}_{(-l)}; \mathbf{g}) \quad (95)$$

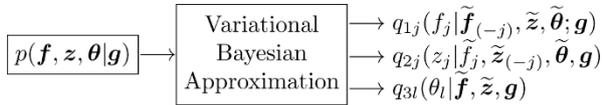
illustrated here:



or even by:

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = \prod_j q_{1j}(f_j | \tilde{\mathbf{f}}_{(-j)}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) \prod_j q_{2j}(z_j | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}_{(-j)}, \boldsymbol{\theta}; \mathbf{g}) \prod_l q_{3l}(\theta_l | \tilde{\mathbf{f}}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}_{(-l)}; \mathbf{g}) \quad (96)$$

illustrated here:



Here, we consider the second case (Equation (95)) and give some more details on it. First to simplify the notations, we write it as:

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}) = q_1(\mathbf{f}) \prod_j q_{2j}(z_j) \prod_l q_{3l}(\theta_l) \quad (97)$$

where it can be shown that:

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

$$q_{2j}(z_j) \propto \exp \left\{ -\langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_1(\mathbf{f})q_3(\boldsymbol{\theta})q_2(\mathbf{z}_{(-j)})} \right\}$$

$$q_{3l}(\theta_l) \propto \exp \left\{ -\langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_1(\mathbf{f})q_2(\mathbf{z})q_3(\boldsymbol{\theta}_{(-l)})} \right\}$$

where $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) = p(\mathbf{g} | \mathbf{f}, \boldsymbol{\theta})p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta})p(\mathbf{z} | \boldsymbol{\theta})p(\boldsymbol{\theta})$ and where $q_2(\mathbf{z}) = \prod_j q_{2j}(z_j)$, $q_3(\boldsymbol{\theta}) = \prod_l q_{3l}(\theta_l)$, $q_2(\mathbf{z}_{(-j)}) = \prod_{i \neq j} q_{2i}(z_i)$, $\langle \cdot \rangle_q$ means expected value with respect to q .

In that case, with appropriate models for the priors (exponential families) and hyper parameters (conjugate priors), we see that $q(\mathbf{f})$ is a multivariate Gaussian $g(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$, $q(\theta_l)$ are either Gaussians (for the means) or Inverse Gammas (for the variances) and $q(z_j)$ are discrete distributions whose expressions can be written easily.

To illustrate this in more detail, we consider the case of the Student-t model.

4.1 Student-t model

In this case, we have the following relations for the forward model and the prior laws:

$$\begin{cases} p(\mathbf{g} | \mathbf{f}, v_\varepsilon) = \mathcal{N}(\mathbf{g} | \mathbf{H}\mathbf{f}, v_\varepsilon \mathbf{I}), & \tau = 1/v_\varepsilon \\ p(\mathbf{f} | \mathbf{z}) = \prod_j p(f_j | z_j) = \prod_j \mathcal{N}(z_j | 0, z_j) = \mathcal{N}(\mathbf{f} | \mathbf{0}, \mathbf{Z}) \\ \mathbf{Z} = \text{diag}[\mathbf{z}], \quad a_j = 1/z_j, \quad \mathbf{A} = \text{diag}[\mathbf{a}] = \mathbf{Z}^{-1} \\ p(\mathbf{a}) = \prod_j \mathcal{G}(a_j | \alpha_0, \beta_0) \\ p(\tau) = \mathcal{G}(\tau | \alpha_{\tau 0}, \beta_{\tau 0}) \end{cases} \quad (98)$$

Then, we obtain the following expressions for the VBA:

$$\begin{cases} \tilde{q}(\mathbf{f}) = \mathcal{N}(\mathbf{f} | \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \\ \tilde{\boldsymbol{\mu}} = \langle \tau \rangle \boldsymbol{\Sigma} \mathbf{H}' \mathbf{g} \\ \tilde{\boldsymbol{\Sigma}} = (\langle \tau \rangle \mathbf{H}' \mathbf{H} + \tilde{\mathbf{Z}})^{-1}, \text{ with } \tilde{\mathbf{Z}} = \tilde{\mathbf{A}}^{-1} = \text{diag}[\tilde{\mathbf{a}}] \end{cases} \quad (99)$$

$$\begin{cases} \tilde{q}(\tau) = \mathcal{G}(\tau | \tilde{\alpha}_\tau, \tilde{\beta}_\tau), \\ \tilde{\alpha}_\tau = \alpha_{\tau 0} + (n+1)/2 \\ \tilde{\beta}_\tau = \beta_{\tau 0} + 1/2 \left[\|\mathbf{g}\|^2 - 2 \langle \mathbf{f} \rangle' \mathbf{H}' \mathbf{g} + \mathbf{H}' \langle \mathbf{f} \mathbf{f}' \rangle \mathbf{H} \right] \\ \tilde{q}(a_j) = \mathcal{G}(a_j | \tilde{\alpha}_j, \tilde{\beta}_j) \\ \tilde{\alpha}_j = \alpha_{00} + 1/2 \\ \tilde{\beta}_j = \beta_{00} + \langle f_j^2 \rangle / 2 \end{cases} \quad (100)$$

where the expressions of the expectations needed are:

$$\begin{cases} \langle \mathbf{f} \rangle = \tilde{\boldsymbol{\mu}} \\ \langle \mathbf{f} \mathbf{f}' \rangle = \boldsymbol{\Sigma} + \boldsymbol{\mu} \boldsymbol{\mu}' \\ \langle f_j^2 \rangle = [\boldsymbol{\Sigma}]_{jj} + \mu_j^2 \\ \langle \tau \rangle = \tilde{\tau} = \tilde{\alpha}_\tau / \tilde{\beta}_\tau \\ \langle a_j \rangle = \tilde{a}_j = \tilde{\alpha}_j / \tilde{\beta}_j \end{cases} \quad (101)$$

We can also express the free energy expression:

$$\begin{aligned} \mathcal{F}(q) &= \left\langle \ln \frac{p(\mathbf{f}, \mathbf{a}, \tau, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{a}, \tau)} \right\rangle \\ &= \langle \ln p(\mathbf{g} | \mathbf{f}, \mathbf{a}, \tau) \rangle + \langle \ln p(\mathbf{f} | \mathbf{a}, \tau) \rangle + \langle \ln p(\mathbf{a} | \tau) \rangle \\ &\quad + \langle -\ln q(\mathbf{f}) \rangle + \langle -\ln q(\mathbf{a}) \rangle + \langle -\ln q(\tau) \rangle \end{aligned} \quad (102)$$

where

$$\begin{aligned} \langle \ln p(\mathbf{g} | \mathbf{f}, \tau) \rangle &= \frac{n}{2} (\langle \ln \tau \rangle - \ln(2\pi)) \\ &\quad - \frac{1}{2} \{ \langle \tau \rangle \mathbf{g}' \mathbf{g} - 2 \langle \mathbf{f} \rangle' \mathbf{H}' \mathbf{g} + \mathbf{H}' \langle \mathbf{f} \mathbf{f}' \rangle \mathbf{H} \} \\ \langle -\ln p(\mathbf{f} | \mathbf{a}) \rangle &= -\frac{n+1}{2} \ln(2\pi) \\ &\quad - \frac{1}{2} \left\{ \sum_j \langle \ln \alpha_j \rangle + \langle \alpha_j \rangle \langle f_j^2 \rangle \right\} \\ \langle -\ln p(\mathbf{a}) \rangle &= -(n+1) \alpha_{\varepsilon_0} \ln(\beta_{\varepsilon_0}) \\ &\quad + (\alpha_{\varepsilon_0} - 1) \sum_j [\langle \ln \alpha_j \rangle - \beta \langle \alpha_j \rangle] - (n+1) \ln \Gamma(\alpha) \\ \langle p(\tau) \rangle &= c \ln d + (c-1) \langle \ln \tau \rangle - d \langle \tau \rangle - \ln \Gamma(c) \end{aligned}$$

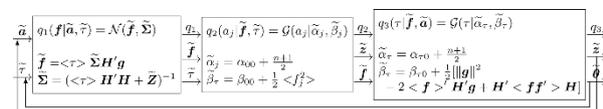
and

$$\begin{aligned} (-\ln q(\mathbf{f})) &= -\frac{n+1}{2}(1 + \ln(2\pi)) - \frac{1}{2} \ln |\Sigma_j| \\ (-\ln q(\mathbf{a})) &= -\sum_j [\tilde{\alpha}_j \ln(\tilde{\beta}_j) + (\tilde{\alpha}_j - 1) \langle \ln \tilde{\alpha}_j \rangle \\ &\quad - \tilde{\beta}_j \langle \alpha_j \rangle - \ln \Gamma(\tilde{\alpha}_j)] \\ (q(\tau)) &= \tilde{c} \ln \tilde{d} + (\tilde{c} - 1) \langle \ln \tau \rangle - \tilde{d} \langle \tau \rangle - \ln \Gamma(\tilde{c}) \end{aligned}$$

In these equations,

$$\begin{cases} \langle \ln a_j \rangle = \psi(\tilde{a}_j) - \ln \tilde{b}_j \\ \langle \ln \tau \rangle = \psi(\tilde{c}) - \ln \tilde{d} \\ \psi(a) = \frac{\partial \ln \Gamma(a)}{\partial a} \end{cases} \quad (103)$$

The resulting algorithm can be summarized as follows



5 Conclusion

The sparsity is a required property in many signal and image processing applications. In this article, first we reviewed the main steps of the Bayesian approach for inverse problems in signal and image processing. Then we presented in a synthetic way the different prior models which can be used to enforce the sparsity. These models have been presented in two categories: simple and hierarchical with hidden variables. For each of these prior models, we discuss their properties and the way to use them in a Bayesian approach resulting to many different inversion algorithms.

We have applied these Bayesian algorithms in many different applications such as X-ray computed tomography [35,36], optical diffraction tomography [37-39], positron emission tomography [40], Microwave imaging [41,42], Sources separation [43-46], spectrometry [47,48], Hyper spectral imaging [49], super resolution [50-52], image fusion [53], image segmentation [54], synthetic aperture radar (SAR) imaging [29]. To save the place and be very synthetic, we did not give here any simulation results or any results on different applications of these methods. These can be found in different articles just referenced.

Acknowledgements

This study had been partially founded by the C5Sys project (Circadian and Cell cycle Clock systems in Cancer) of ERASYSBIO+. <http://www.erasysbio.net/index.php?index=272>

Competing interests

The author declares that they have no competing interests.

Received: 18 January 2012 Accepted: 1 March 2012

Published: 1 March 2012

References

1. A Tikhonov, Regularization of incorrectly posed problems. *Soviet Math Dokl.* **4**, 1624-1627 (1963)
2. A Tikhonov, V Aréline, *Méthodes de Résolution de Problèmes Mal Posés*, (MIR, Moscou, Russia, 1976). Éditions
3. I Daubechies, M Defrise, CD Mol, An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Comm Pure Appl Math.* **57**, 1413-1457 (2004). doi:10.1002/cpa.20042
4. DL Donoho, Compressive sampling. *IEEE Trans Inf Theory.* **52**(4), 1289-1306 (2006)
5. JA Tropp, AC Gilbert, MJ Strauss, Algorithms for simultaneous sparse approximation. Part I: Greedy pursuit. *Signal Processing*, special issue "sparse approximations in signal and image processing". **86**, 572-588 (2006)
6. JA Tropp, Algorithms for simultaneous sparse approximation. Part II: Convex relaxation. *Signal Process* (special issue "Sparse approximations in signal and image processing"). **86**, 589-602 (2006)
7. R Zass, A Shashua, in *Nonnegative Sparse PCA*, vol. 19. (Cambridge, MA: MIT Press, 2007), pp. 1561-1568
8. EJ Candès, M Wakin, S Boyd, Enhancing sparsity by reweighted l1 minimization. *J Fourier Anal Appl.* **14**, 877-905 (2008). doi:10.1007/s00041-008-9045-x
9. DM Witten, R Tibshirani, T Hastie, A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics.* **10**(3), 515-534 (2009). doi:10.1093/biostatistics/kxp008
10. S Vaiter, G Peyré, C Dossal, J Fadili, Robust Sparse Analysis Regularization. Tech rep, preprint Hal-00627452 <http://hal.archives-ouvertes.fr/hal-00627452/> (2011)
11. G Peyré, J Fadili, Learning Analysis Sparsity Priors. *Proc of Sampta'11* <http://hal.archives-ouvertes.fr/hal-00542016/> (2011)
12. P Williams, Bayesian regularization and pruning using a Laplace prior. *Neural Comput.* **7**, 117-143 (1995)
13. T Mitchell, J Beauchamp, Bayesian variable selection in linear regression. *J Am Stat Assoc.* **83**(404), 1023 (1988). doi:10.2307/2290129
14. N Polson, J Scott, Shrink globally, act locally: sparse Bayesian regularization and prediction. *Bayesian Stat.* **9**, 1-24 (2010)
15. A Mohammad-Djafari, On the estimation of hyperparameters in Bayesian approach of solving inverse problems, in *Proc IEEE International Conference on Acoustics, Speech, and Signal Processing ICASSP-93*, vol. 5. (Minneapolis, MN, USA, 1993), pp. 495-498
16. A Mohammad-Djafari, Joint estimation of parameters and hyperparameters in a Bayesian approach of solving inverse problems, in *IEEE Int Conf on Image Processing (ICIP), IEEE ICIP 96*, vol. II. (Lausanne, Suisse, 1996), pp. 473-477
17. R Neal, G Hinton, A view of the EM algorithm that justifies incremental, sparse, and other variants. *Learn Graph Models.* **89**, 355-368 (1998)
18. A Doucet, P Duvaut, Bayesian estimation of state-space models applied to deconvolution of Bernoulli-Gaussian processes. *Signal Process.* **57**(2), 147-161 (1997). doi:10.1016/S0165-1684(96)00192-2
19. T Park, G Casella, The Bayesian Lasso. *J Am Stat Assoc.* **103**(482), 681-686 (2008). doi:10.1198/016214508000000337
20. M Tipping, Sparse Bayesian learning and the relevance vector machine. *J Mach Learn Res.* **1**, 211-244 (2001)
21. C Févotte, S Godsill, A Bayesian approach for blind separation of sparse source. *IEEE Trans Audio Speech Lang Process.* **14**, 2174-2188 (2006)
22. F Caron, A Doucet, Sparse Bayesian nonparametric regression, in *International Conference on Machine Learning*, pp. 88-95 (2008)
23. J Griffin, P Brown, Inference with normal-gamma prior distributions in regression problems. *Bayesian Anal.* **5**, 171-188 (2010)
24. H Snoussi, J Idier, Bayesian blind separation of generalized hyperbolic processes in noisy and underdetermined mixtures. *IEEE Trans Signal Process.* **54**, 3257-3269 (2006)
25. H Ishwaran, JS Rao, Spike and slab variable selection: frequentist and Bayesian strategies. *Ann Stat.* **33**(2), 730-733 (2005). doi:10.1214/009053604000001147
26. S Chatzis, T Varvarigou, Factor analysis latent subspace modeling and robust fuzzy clustering using t-distributionsclassification of binary random patterns. *IEEE Trans Fuzzy Syst.* **17**, 505-517 (2009)

27. CA Bouman, KD Sauer, A generalized Gaussian image model for edge-preserving MAP estimation. *IEEE Trans Image Process.* **2**(3), 296–310 (1993). doi:10.1109/83.236536
28. H Zou, T Hastie, Regularization and variable selection via the elastic net. *J Royal Stat Soc, Ser. B.* **67**(2), 301–320 (2005). doi:10.1111/j.1467-9868.2005.00503.x
29. S Zhu, A Mohammad-Djafari, H Wang, B Deng, X Li, J Mao, Parameter estimation for SAR micromotion target based on sparse signal representation. *Eurasip special issue "Sparse approximations in signal and image processing"*, **13** (2012)
30. J Idier, *Approche Bayésienne Pour les Problèmes Inverses*, (Traité IC2, Série traitement du signal et de l'image, Hermès, Paris, 2001)
31. F Champagnat, Y Goussard, J Idier, Unsupervised deconvolution of sparse spike trains using stochastic approximation. *IEEE Trans Signal Process.* **44**(12), 2988–2998 (1996). doi:10.1109/78.553473
32. D Ge, J Idier, E Le Carpentier, A new MCMC algorithm for blind Bernoulli-Gaussian deconvolution, *Proceedings of EUSIPCO: Septembre 2008*, (Lausanne, Suisse, 2008)
33. JJ Kormylo, JM Mendel, Maximum-likelihood detection and estimation of Bernoulli-Gaussian processes. *IEEE Trans Inf Theory.* **28**, 482–488 (1982). doi:10.1109/TIT.1982.1056496
34. M Lavielle, Bayesian deconvolution of Bernoulli-Gaussian processes. *Signal Process.* **33**, 67–79 (1993)
35. A Mohammad-Djafari, Gauss-Markov-Potts priors for images in computer tomography resulting to joint optimal reconstruction and segmentation. *Int J Tomography Stat.* **11**(W09), 76–92 http://djafari.free.fr/pdf/IJTS_08.pdf (2008)
36. N Gac, A Vabre, A Mohammad-Djafari, F Buyens, GPU implementation of a 3D bayesian CT algorithm and its application on real foam reconstruction, *Proceedings of the first International Conference on Image Formation in X-Ray Computed Tomography*, (Salt Lake City, Utah, USA, 2010)
37. H Ayasso, B Duchne, A Mohammad-Djafari, Bayesian inversion for optical diffraction tomography. *J Modern Opt.* **57**(9), 765–776 (2010). doi:10.1080/09500340903564702
38. H Ayasso, A Mohammad-Djafari, Joint NDT image restoration and segmentation using Gauss-Markov-Potts prior models and variational Bayesian computation. *IEEE Trans. Image Process.* **19**(9), 2265–2277 <http://dx.doi.org/10.1109/TIP.2010.2047902> (2010)
39. H Ayasso, B Duchne, A Mohammad-Djafari, Optical diffraction tomography within a variational Bayesian framework. *Inverse Probl Sci Eng.* **iFirst**, 1–15 (2011)
40. MD Fall, ?É? Barat, C Comtat, T Dautremer, T Montagu, A Mohammad-Djafari, A discrete-continuous Bayesian model for emission tomography, *IEEE International Conference on Image Processing (ICIP)*, (Bruxell, Belgium, 2011), pp. 1401–1404
41. O Féron, B Duchêne, A Mohammad-Djafari, Microwave imaging of inhomogeneous objects made of a finite number of dielectric and conductive materials from experimental data. *Inverse Probl.* **21**(6), 95–115 <http://djafari.free.fr/pdf/> (2005). doi:10.1088/0266-5611/21/6/S08
42. O Féron, B Duchêne, A Mohammad-Djafari, Microwave imaging of piecewise constant objects in a 2D-TE configuration. *Int J Appl Electromag Mech.* **26**(6), 167–174 <http://djafari.free.fr/pdf/jae00905.pdf> (2007)
43. H Snoussi, A Mohammad-Djafari, Fast joint separation and segmentation of mixed images. *J Electron Imag.* **13**(2), 349–361 <http://djafari.free.fr/pdf/> (2004). doi:10.1117/1.1666873
44. H Snoussi, A Mohammad-Djafari, Bayesian unsupervised learning for source separation with mixture of Gaussians prior. *J VLSI Signal Process Syst.* **37**(2/3), 263–279 <http://djafari.free.fr/pdf/VLSIpaper.pdf> (2004)
45. A Mohammad-Djafari, Bayesian source separation: beyond PCA and ICA, in *ESANN 2006*, (Belgium, 2006) <http://djafari.free.fr/pdf/>
46. M Ichir, A Mohammad-Djafari, Hidden Markov models for wavelet-based blind source separation. *IEEE Trans Image Process.* **15**(7), 1887–1899 (2006)
47. A Mohammad-Djafari, J Giovannelli, G Demoment, J Idier, Regularization, maximum entropy and probabilistic methods in mass spectrometry data processing problems. *Int J Mass Spectrom.* **215**(1-3), 175–193 <http://djafari.free.fr/pdf/maspec12013.pdf> (2002). doi:10.1016/S1387-3806(01)00562-0
48. S Moussaoui, D Brie, A Mohammad-Djafari, C Carteret, Separation of non-negative mixture of non-negative sources using a Bayesian approach and MCMC sampling. *IEEE Trans Signal Process.* **54**(11), 4133–4145 (2006)
49. N Bali, A Mohammad-Djafari, Bayesian approach with hidden Markov modeling and mean field approximation for hyperspectral data analysis. *IEEE Trans Image Process.* **17**(2), 217–225 (2008)
50. F Humblot, A Mohammad-Djafari, Super-resolution using hidden Markov model and bayesian detection estimation framework, in *EURASIP J Appl Signal Process*, vol. 16. (Special number on Super-Resolution Imaging: Analysis, Algorithms, and Applications, 2006) <http://www.hindawi.com/GetArticle.aspx?doi=10.1155/ASP/2006/36971>. Article ID 36971
51. A Mohammad-Djafari, Super-resolution: a short review, a new method based on hidden Markov modeling of HR image and future challenges. *Comput J* <http://djafari.free.fr/pdf/bxn005v1.pdf> (2008)
52. M Mansouri, A Mohammad-Djafari, Joint super-resolution and segmentation from a set of low resolution images using a Bayesian approach with a Gauss-Markov-Potts prior. *Int J Signal Imag Syst Eng.* **3**(4), 211–221 (2010). doi:10.1504/IJSISE.2010.038017
53. O Féron, A Mohammad-Djafari, Image fusion and joint segmentation using an MCMC algorithm. *J Electron Imag.* **14**(2) <http://arxiv.org/abs/physics/0403150> (2005). paper no. 023014
54. P Brault, A Mohammad-Djafari, Unsupervised Bayesian wavelet domain segmentation using a Potts-Markov random field modeling. *J Electron Imag.* **14**(4) <http://djafari.free.fr/pdf/> (2005). 043011-1-043011-16

doi:10.1186/1687-6180-2012-52

Cite this article as: Mohammad-Djafari: Bayesian approach with prior models which enforce sparsity in signal and image processing. *EURASIP Journal on Advances in Signal Processing* 2012 **2012**:52.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► springeropen.com