



Bayesian inference methods for sources separation

Ali Mohammad-Djafari

Laboratoire des signaux et systèmes (L2S)

UMR 8506 CNRS-SUPELEC-UNIV PARIS SUD

SUPELEC, Plateau de Moulon, 91192 Gif-sur-Yvette, France

ABSTRACT

The main aim of this paper is first to present the Bayesian inference approach for sources separation where we want to infer on the mixing matrix, the sources and all the hyper-parameters associated to probabilistic modeling (likelihood and priors). For this purpose, the sources separation problem is considered in four steps: i) Estimation of the sources when the mixing matrix is known; ii) Estimation of the mixing matrix when the sources are known; iii) Joint estimation of sources and the mixing matrix; and finally, iv) Joint estimation of sources and the mixing matrix, hidden variables and hyper-parameters. In all cases, one of the main steps is modeling of sources and the mixing matrix prior laws. We propose to use sparsity enforcing probability laws (such as Generalized Gaussian, Student-t and mixture models) both for the sources and the mixing matrix. For algorithmic and computational aspects, we consider either Joint MAP, MCMC Gibbs sampling or Variational Bayesian Approximation tools. For each class of methods we discuss about their relative costs and performances.

1 Introduction

The general sources separation problem can be viewed as an inference problem where first we provide a model linking the observed data (mixed signals) $g(t)$ to unknown sources $f(t)$ through a forward model. In this paper, we only consider the instantaneous mixing model:

$$g(t) = Af(t) + \varepsilon(t), \quad t \in [1, \dots, T] \quad (1)$$

where $\varepsilon(t)$ represents the errors of modeling and measurement. A is called mixing matrix and when it is invertible, its inverse $B = A^{-1}$ is called the separating matrix. The second step is to write down the expression of the $p(f|A, g)$ when the mixing matrix A is known, $p(A|f, g)$ when the sources f are known, or the joint posterior law:

$$p(f, A|g, \theta) = \frac{p(g|f, A, \theta_1) p(f|\theta_2) p(A|\theta_3)}{p(g|\theta_1, \theta_2, \theta_3)} \quad (2)$$

where $p(g|f, A, \theta_1)$ is the likelihood and $p(f|\theta_2)$ and $p(A|\theta_3)$ are the priors on sources and the mixing matrix and $\theta = (\theta_1, \theta_2, \theta_3)$ represent the hyper-parameters of the problem. In real and effective problems, we also have to estimate them. This can be done, by assigning them a prior $p(\theta)$ often separable $p(\theta) = p(\theta_1) p(\theta_2) p(\theta_3)$ and writing the expression of the joint posterior:

$$p(f, A, \theta|g) = \frac{p(g|f, A, \theta_1) p(f|\theta_2) p(A|\theta_3) p(\theta)}{p(g)} \quad (3)$$

In this paper, we will consider these two cases with different prior modeling for sources $p(f|\theta_2)$ and different priors for the mixing matrix $p(A|\theta_3)$. In particular, we consider the Generalized Gaussian (GG), Student-t (St), Elastic net (EN) and Mixture of Gaussians (MoG) models. Some of these models are well-known [2, 5–10], some others less. In general, we can classify them in two categories: i) Simple Non Gaussian models with heavy tails and ii) Mixture models with hidden variables z which result to hierarchical models.

The second main step in the Bayesian approach is to do the computations. The Bayesian computations in general can be:

- Joint optimization of $p(f, A, \theta|g)$ which needs optimization algorithms;
- MCMC Gibbs sampling methods which need generation of samples from the conditionals $p(f|A, \theta, g)$, $p(A|f, \theta, g)$ and $p(\theta|f, A, g)$;
- Bayesian Variational Approximation (BVA) methods which approximate $p(f, A, \theta|g)$ by a separable one

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

and then using them for the estimation [11].

The rest of the paper is organized as follows:

In section 2, we review a few prior models which are frequently used in particular when sparsity has to be enforced. For example, the Generalized Gaussian (GG) with two particular cases of Gaussian (G) and Double Exponential (DE) or Laplace, the Student-t model which can be interpreted as an infinite mixture with a variance hidden variable, Elastic net and the mixture models.

In Section 3, we examine in details the estimation of the sources f when the mixing matrix A is known. In section 4 the estimation of the mixing matrix A when the sources f are known is considered. In section 5 we examine the joint estimation of the mixing matrix A and the sources f , and finally, the more realistic case of joint estimation of the mixing matrix A , the sources f , their hidden variables z and the hyper-parameters θ .

In section 6, we give principal practical algorithms which can be used in real applications, and finally, in section V we show some results and real applications and in particular in SAR imaging [12] and Sonar sources localization [3, 4].

2 Prior models enforcing sparsity

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 [1]). Its expression is:

$$p(\mathbf{f}|\gamma, \beta) = \prod_j \mathcal{GG}(f_j|\gamma, \beta) \propto \exp \left\{ -\gamma \sum_j |f_j|^\beta \right\} \quad (4)$$

where

$$\mathcal{GG}(f_j|\gamma, \beta) = \frac{\beta\gamma}{2\Gamma(1/\beta)} \exp \left\{ -\gamma |f_j|^\beta \right\} \quad (5)$$

and where the particular cases of $\beta = 2$ (Gaussian):

$$p(\mathbf{f}|\gamma) = \prod_j \mathcal{N}(f_j|0, 1/(2\gamma)) \propto \exp \left\{ -\gamma \sum_j |f_j|^2 \right\} \propto \exp \left\{ -\gamma \|\mathbf{f}\|_2^2 \right\} \quad (6)$$

$\beta = 1$ (Double exponential or Laplace):

$$p(\mathbf{f}|\gamma) = \prod_j \mathcal{DE}(f_j|\gamma) \propto \exp \left\{ -\gamma \sum_j |f_j| \right\} \propto \exp \left\{ -\gamma \|\mathbf{f}\|_1 \right\} \quad (7)$$

as well as the cases where $0 < \beta < 1$ are of great interest for sparsity enforcing.

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

The second simplest model is the Student-t model:

$$p(\mathbf{f}|\nu) = \prod_j \mathcal{St}(f_j|\nu) \propto \exp \left\{ -\frac{\nu+1}{2} \sum_j \log(1 + f_j^2/\nu) \right\} \quad (8)$$

where

$$\mathcal{St}(f_j|\nu) = \frac{1}{\sqrt{\pi\nu}} \frac{\Gamma((\nu+1)/2)}{\Gamma(\nu/2)} (1 + f_j^2/\nu)^{-(\nu+1)/2} \quad (9)$$

Knowing that

$$\mathcal{St}(f_j|\nu) = \int_0^\infty \mathcal{N}(f_j|0, 1/\tau_j) \mathcal{G}(\tau_j|\nu/2, \nu/2) d\tau_j \quad (10)$$

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

$$\begin{aligned} p(\mathbf{f}|\tau) &= \prod_j p(f_j|\tau_j) = \prod_j \mathcal{N}(f_j|0, 1/\tau_j) \propto \exp \left\{ -\frac{1}{2} \sum_j \tau_j f_j^2 \right\} \\ p(\tau_j|\alpha, \beta) &= \mathcal{G}(\tau_j|\alpha, \beta) \propto \tau_j^{(\alpha-1)} \exp \left\{ -\beta \tau_j \right\} \text{ with } \alpha = \beta = \nu/2 \end{aligned} \quad (11)$$

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

$$p(\mathbf{f}) = \prod_j \mathcal{C}(f_j) \propto \exp \left\{ -\sum_j \log(1 + f_j^2) \right\} \quad (12)$$

2.3 Elastic Net (EN) prior model

A prior model inspired from Elastic Net regression literature [?] is:

$$p(\mathbf{f}|\mathbf{v}) = \prod_j \mathcal{E}\mathcal{N}(\mathbf{f}_j|\mathbf{v}) \propto \exp \left\{ - \sum_j (\gamma_1 |\mathbf{f}_j| + \gamma_2 \mathbf{f}_j^2) \right\} \quad (13)$$

where

$$\mathcal{E}\mathcal{N}(\mathbf{f}_j|\mathbf{v}) = \mathcal{N}(0, 1/\gamma_1) \mathcal{D}\mathcal{E}(\gamma_1) = c \exp \{ -\gamma_1 |\mathbf{f}_j| - \gamma_2 \mathbf{f}_j^2 \} \quad (14)$$

which is a product of a Gaussian and a Double Exponential pdfs.

2.4 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:

$$p(\mathbf{f}|\lambda, v_1, v_0) = \prod_j (\lambda \mathcal{N}(\mathbf{f}_j|0, v_1) + (1 - \lambda) \mathcal{N}(\mathbf{f}_j|0, v_0)) \quad (15)$$

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

$$\begin{cases} p(\mathbf{f}|\mathbf{z}) = \prod_j p(\mathbf{f}_j|z_j) = \prod_j \mathcal{N}(\mathbf{f}_j|0, v_{z_j}) \propto \exp \left\{ -\frac{1}{2} \sum_j \frac{\mathbf{f}_j^2}{v_{z_j}} \right\} \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (16)$$

In general $v_1 \gg v_0$ and λ measures the sparsity ($0 < \lambda \ll 1$).

2.5 Bernoulli-Gaussian (BG) model

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

$$p(\mathbf{f}|\lambda, v) = \prod_j p(\mathbf{f}_j) = \prod_j (\lambda \mathcal{N}(\mathbf{f}_j|0, v) + (1 - \lambda) \delta(\mathbf{f}_j)) \quad (17)$$

which can also be written as

$$\begin{cases} p(\mathbf{f}|\mathbf{z}) = \prod_j p(\mathbf{f}_j|z_j) = \prod_j [\mathcal{N}(\mathbf{f}_j|0, v)]^{z_j} \prod_j [\delta(\mathbf{f}_j)]^{(1-z_j)} \\ P(z_j = 1) = \lambda, \quad P(z_j = 0) = 1 - \lambda \end{cases} \quad (18)$$

This model has also been called *spike and slab* [6].

3 Bayesian inference for sources f when A is known

First let assume the error ε to be centered, Gaussian and white: $\varepsilon \sim \mathcal{N}(\varepsilon|0, v_\varepsilon I)$. Then, using the forward model (1) we have

$$p(g|f, A) = \mathcal{N}(g|Af, v_\varepsilon I) \propto \exp \left\{ -\frac{1}{2v_\varepsilon} \|g - Af\|^2 \right\} \quad (19)$$

Now, we consider different priors for f .

3.1 Simple prior models

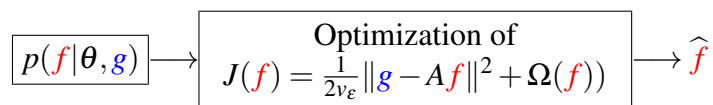
Given $p(g|f, A)$ and any simple prior law $p(f)$, the posterior law is written:

$$p(f|g, A) \propto p(g|f, A) p(f) \propto \exp \{J(f)\} \quad (20)$$

with

$$J(f) = \frac{1}{2v_\varepsilon} \|g - Af\|^2 + \Omega(f) \quad (21)$$

where $\Omega(f) = -\ln p(f)$ 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. The MAP estimate with this prior can be summarized in the following scheme:



Thus, depending on the choice of the prior we obtain different expressions for $\Omega(f)$. For example for the GG model of (4) we get

$$\Omega(f) = \gamma \sum_j |f_j|^\beta. \quad (22)$$

For the Student-t model (8) we get

$$\Omega(f) = \frac{v+1}{2} \sum_j \log(1 + f_j^2/v). \quad (23)$$

For the Elastic Net model we get

$$\Omega(f) = \sum_j [\gamma_1 |f_j| + \gamma_2 f_j^2] \quad (24)$$

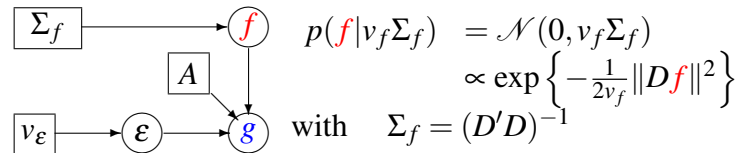
For each of these cases, we may discuss on the unimodality and convexity of the criterion $J(f)$ which depends mainly on its Hessian

$$\Delta J(f) = \left[\frac{\partial^2 J(f)}{\partial f_j \partial f_i} \right] = A'A + \left[\frac{\partial^2 \Omega(f)}{\partial f_i \partial f_j} \right] \quad (25)$$

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

3.2 Gaussian priors

The case of the Gaussian prior is interesting, because, in that case, the posterior is also Gaussian and we have an analytical solution for \hat{f} .



The MAP estimate can be computed as:

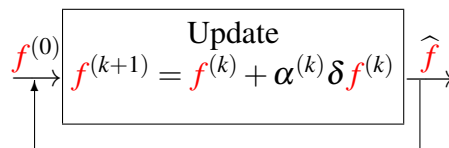
$$\hat{f} = \arg \min_f \{J(f)\} \text{ with } J(f) = \|g - Af\|^2 + \lambda \|Df\|^2 \quad (26)$$

and the solution is

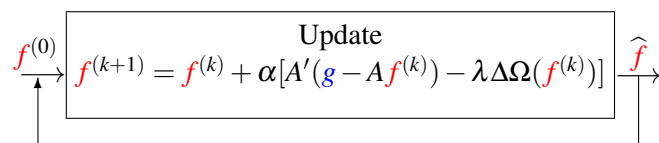
$$\hat{f} = (A'A + \lambda D'D)^{-1} A'g \quad (27)$$

3.3 Non Gaussian priors

For other cases, the optimization is done iteratively. The basics of this algorithm is shown in the following scheme:



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



where $\Delta\Omega$ here means the gradient of $\Omega(f)$.

3.4 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 f and the hidden variables: τ in Cauchy model, z in MoG2 and BG models. Let examine these a little in details.

Student-t and Cauchy models

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

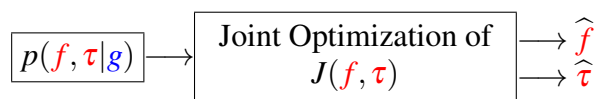
$$\begin{aligned} p(\mathbf{f}, \boldsymbol{\tau}) &= \prod_j p(f_j | \tau_j) p(\tau_j) = \prod_j \mathcal{N}(f_j | 0, 1/\tau_j) p(\tau_j) \\ &\propto \exp \left\{ -\frac{1}{2} \sum_j \tau_j f_j^2 + \alpha \ln \tau_j - \beta \tau_j \right\} \text{ with } \alpha = \beta = \nu/2 \end{aligned} \quad (28)$$

such that

$$p(\mathbf{f}, \boldsymbol{\tau} | \mathbf{g}) \propto p(\mathbf{g} | \mathbf{f}) p(\mathbf{f}, \boldsymbol{\tau}) \propto \exp \{-J(\mathbf{f}, \boldsymbol{\tau})\} \quad (29)$$

where

$$J(\mathbf{f}, \boldsymbol{\tau}) = \frac{1}{2\nu_\varepsilon} \|\mathbf{g} - A\mathbf{f}\|^2 + \sum_j \frac{1}{2} \tau_j f_j^2 - \alpha \ln \tau_j + \beta \tau_j \quad (30)$$



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

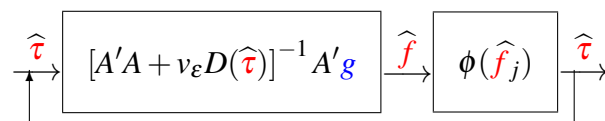
$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f}, \boldsymbol{\tau})\} = \arg \min_{\mathbf{f}} \left\{ \frac{1}{2\nu_\varepsilon} \|\mathbf{g} - A\mathbf{f}\|^2 + \sum_j \frac{1}{2} \tau_j f_j^2 \right\} \quad (31)$$

and with respect to $\boldsymbol{\tau}$ (with fixed \mathbf{f})

$$\hat{\boldsymbol{\tau}} = \arg \min_{\boldsymbol{\tau}} \{J(\mathbf{f}, \boldsymbol{\tau})\} = \arg \min_{\boldsymbol{\tau}} \left\{ \sum_j \frac{1}{2} \tau_j f_j^2 - \alpha \ln \tau_j + \beta \tau_j \right\} \quad (32)$$

results in the following iterative algorithm:

$$\begin{cases} \hat{\mathbf{f}} = [A'A + \nu_\varepsilon D(\hat{\boldsymbol{\tau}})]^{-1} A' \mathbf{g} \\ \hat{\tau}_j = \phi(\hat{f}_j) = \frac{\alpha}{\hat{f}_j^2 + \beta} \\ D(\hat{\boldsymbol{\tau}}) = \text{diag} [1/\hat{\tau}_j, j = 1, \dots, n] \end{cases} \quad (33)$$



Note that τ_j is inverse of a variance and we have $1/\tau_j = \frac{f_j^2 + \beta}{\alpha}$. 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 $D(\boldsymbol{\tau})$.

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

We may also consider a Gibbs sampling scheme

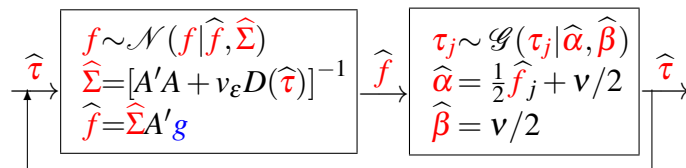
$$\begin{cases} f \sim p(f|\tau, g) \propto p(g|f)p(f|u) = \mathcal{N}(f|\hat{f}, \hat{\Sigma}) \\ \tau \sim p(\tau|f, g) \propto p(f|\tau) p(\tau) = \prod_j \mathcal{G}(\tau_j|\hat{\alpha}, \hat{\beta}) \end{cases} \quad (34)$$

where

$$\begin{cases} \hat{\Sigma} = [A'A + v_\varepsilon D(\tau)]^{-1} \\ \hat{f} = \hat{\Sigma} A' g \end{cases} \quad (35)$$

and

$$\begin{cases} \hat{\alpha} = \frac{1}{2} \hat{f}_j + \alpha = \frac{1}{2} \hat{f}_j + v/2 \\ \hat{\beta} = \beta = v/2 \end{cases} \quad (36)$$



MoG model

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

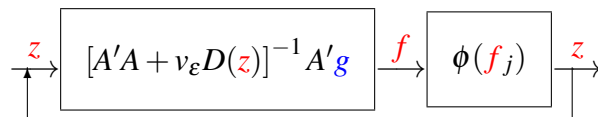
$$p(f, z|g) \propto p(g|f)p(f, z) \propto \exp\{-J(f, z)\} \quad (37)$$

where

$$\begin{aligned} J(f, z) &= \frac{1}{2v_\varepsilon} \|g - Af\|^2 \\ &+ \sum_j \frac{f_j^2}{2v_{z_j}} + z_j \ln \lambda + (1 - z_j) \ln(1 - \lambda) \end{aligned} \quad (38)$$

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

$$\begin{cases} \hat{f} = [A'A + v_\varepsilon D(z)]^{-1} A' g \\ \hat{z}_j = \phi(f_j) = \begin{cases} 1, & \text{if } \hat{f}_j^2 \geq (v_1 - v_0) \ln \frac{1-\lambda}{\lambda} \\ 0, & \text{if } \hat{f}_j^2 < (v_1 - v_0) \ln \frac{1-\lambda}{\lambda} \end{cases} \\ D(\hat{z}) = \text{diag}[v_{\hat{z}_j}, j = 1, \dots, n] \end{cases} \quad (39)$$



Here too, we may also consider a Gibbs sampling scheme

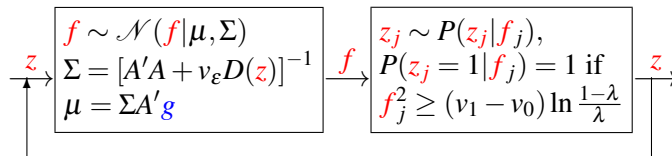
$$\begin{cases} f \sim p(f|z, g) \propto p(g|f)p(f|u) = \mathcal{N}(f|\hat{f}, \hat{\Sigma}) \\ z \sim p(z|f, g) \propto p(f|z) p(z) = \prod_j P(z_j = k|f_j) \end{cases} \quad (40)$$

where

$$\begin{cases} \widehat{\Sigma} = [A'A + v_\varepsilon D(z)]^{-1} \\ \widehat{f} = \widehat{\Sigma} A' g \end{cases} \quad (41)$$

and

$$\begin{cases} P(z_j = 1 | f_j) = 1, & \text{if } f_j^2 \geq (v_1 - v_0) \\ P(z_j = 0 | f_j) = 1, & \text{if } f_j^2 < (v_1 - v_0) \ln \frac{1-\lambda}{\lambda} \end{cases} \quad (42)$$



4 Estimation of the mixing matrix A with known sources f

4.1 Sources separation: a bilinear problem

First, we may note that, the sources separation problem is a bilinear problems: $g = Af = Fa$. To see some details, consider the case of $N = 2$ sources and $M = 2$ mixtures:

$$\begin{bmatrix} g_1 \\ g_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} f_1 \\ f_2 \end{bmatrix} = \begin{bmatrix} f_1 & 0 & f_2 & 0 \\ 0 & f_1 & 0 & f_2 \end{bmatrix} \begin{bmatrix} a_{11} \\ a_{21} \\ a_{12} \\ a_{22} \end{bmatrix}$$

In the following, we use the following notations:

$$F = f \odot I, \quad a = \text{vec}(A)$$

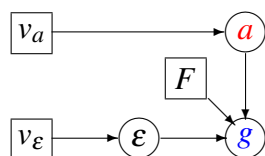
With this introduction and these notations, we may now compare the following cases:

- Estimation of f with known A : $g = Af$;
- Estimation of A with known f : $g = Af = Fa$; and
- Joint Estimation of f and A : $g = Af = Fa$.

We may note that the second and the third problems are more and more under determined.

4.2 The Gaussian case

So, if we follow now the same strategy for $g = Fa$ than the case of $g = Af$, we can summarize the results as follows:



$$p(a | v_f) = \mathcal{N}(0, v_a I) \propto \exp \left\{ -\frac{1}{2v_a} \sum_j |a_j|^2 \right\}$$

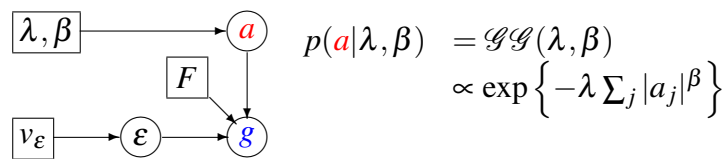
The MAP estimate is given by:

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \{J(\mathbf{a})\} \text{ with } J(\mathbf{a}) = \|\mathbf{g} - F\mathbf{a}\|^2 + \lambda \sum_j |a_j|^2 \quad (43)$$

which results to:

$$\hat{\mathbf{a}} = (F^t F + \lambda I)^{-1} F^t \mathbf{g} \text{ or } \hat{\mathbf{A}} = \mathbf{g} \mathbf{f}^t (\mathbf{f} \mathbf{f}^t + \lambda I)^{-1}. \quad (44)$$

4.3 The Generalized Gaussian case



The MAP estimate becomes

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \{J(\mathbf{a})\} \text{ with } J(\mathbf{a}) = \|\mathbf{g} - F\mathbf{a}\|^2 + \lambda \sum_j |a_j|^\beta \quad (45)$$

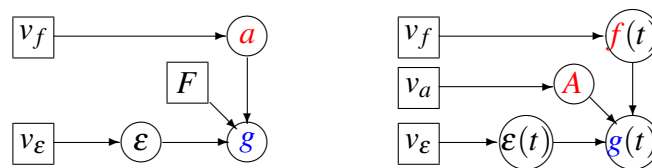
but unfortunately there is no analytic expression for the solution. However, its numerical computation can be done very easily.

5 Joint estimation of the sources f and the mixing matrix A

Now, we consider the main problem, and as in the previous cases, we consider different cases and summarize them as follows:

5.1 Gaussian priors for f and A

Now consider the joint estimation which is illustrated in two forms, one as usual and the second with considering the time index too:



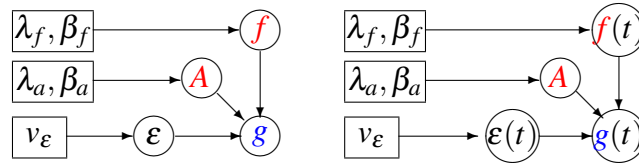
and the equations can be summarized as follows:

$$p(\mathbf{A}, \mathbf{f}(t) | \mathbf{g}(t), v_\epsilon) \propto p(\mathbf{g}(t) | \mathbf{A}, \mathbf{f}(t), v_\epsilon) p(\mathbf{A} | v_a) p(\mathbf{f}(t) | v_f) \quad (46)$$

$$\begin{cases} \hat{\mathbf{f}}(t) = (\hat{\mathbf{A}}' \hat{\mathbf{A}} + \lambda_f I)^{-1} \hat{\mathbf{A}}' \mathbf{g}(t), & \lambda_f = v_\epsilon / v_f \\ \hat{\mathbf{A}} = \sum_t \mathbf{g}(t) \hat{\mathbf{f}}'(t) \left(\sum_t \hat{\mathbf{f}}(t) \hat{\mathbf{f}}'(t) + \lambda_a I \right)^{-1} & \lambda_a = v_\epsilon / v_a \end{cases} \quad (47)$$

5.2 Gaussian prior for f and Generalized Gaussian prior for A

This case also can be summarized as follows:



$$\begin{aligned}
 p(f|v_f) &= \mathcal{N}(0, v_f I) \propto \exp\left\{-\frac{1}{2v_f} \|f\|^2\right\} \\
 p(a|\lambda_a, \beta_a) &= \mathcal{GG}(\lambda_a, \beta_a) \propto \exp\left\{-\lambda_a \sum_j |a_j|^\beta\right\}
 \end{aligned} \tag{48}$$

Joint Posterior:

$$p(A, f|g, \theta) \propto p(g|A, f, v_\epsilon) p(A|v_a) p(f|v_f)$$

Integration over f can be done easily

5.3 Gaussian prior for A and Non Gaussian for sources $f(t)$

Let make a comparison between the Gaussian and Non Gaussian priors for $f(t)$ when trying to implement joint MAP computation iteratively.

- For the **Gaussian** laws for the sources, at each iteration we have:

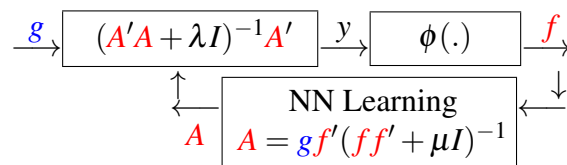
$$\begin{cases} f(t) &= (A'A + \lambda I)^{-1} A' g(t) \\ A &= \sum_t g(t) f'(t) (\sum_t f(t) f'(t) + \mu I)^{-1} \end{cases} \tag{49}$$

with $\lambda = \sigma_\epsilon^2 / \sigma_s^2$ and $\mu = \sigma_\epsilon^2 / \sigma_a^2$.

- For the **Non Gaussian** sources $f(t)$, we have:

$$\begin{cases} y(t) &= (A'A + \lambda I)^{-1} A' g(t) \\ f(t) &= \phi(y(t)) \\ A &= \sum_t g(t) f'(t) (\sum_t f(t) f'(t) + \mu I)^{-1} \end{cases}$$

For both cases, we can see the following scheme:



For Gaussian case $\phi(y) = y$ and a nonlinear function for Non Gaussian case.

5.4 Joint estimation: General case

Let consider now the following general notations which can also be used for accounting for time dependent sources (colored or non stationary sources).

$$\mathbf{g}(t) = \mathbf{A}\mathbf{f}(t) + \boldsymbol{\varepsilon}(t) \longrightarrow \mathbf{g}_{1..T} = \mathbf{A}\mathbf{f}_{1..T} + \boldsymbol{\varepsilon}_{1..T} \quad (50)$$

If we assume that the noise is iid and white, then

$$p(\mathbf{g}_{1..T}|\mathbf{f}_{1..T}) = \prod_t \prod_i p_i(\mathbf{g}_i(t) - [\mathbf{A}\mathbf{f}]_i(t)) \quad (51)$$

and if we assume the same for the sources

$$p(\mathbf{f}_{1..T}) = \prod_t \prod_j p_j(f_j(t)), \quad (52)$$

and

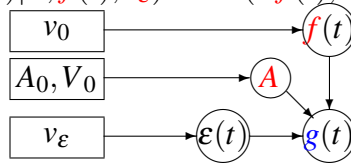
$$p(\mathbf{A}) = \prod_i \prod_j p(a_{ij}) \propto \exp \left\{ -\frac{1}{2\sigma_a^2} \sum_i \sum_j (a_{ij} - a_{0ij})^2 \right\}, \quad (53)$$

then we have:

$$p(\mathbf{A}, \mathbf{f}_{1..T}|\mathbf{g}_{1..T}) \propto \prod_t \prod_i p_i(\mathbf{g}_i(t) - [\mathbf{A}\mathbf{f}]_i(t)) + \prod_t \prod_j p_j(f_j(t)) + \prod_i \prod_j p(a_{ij}) \quad (54)$$

Let consider first the Gaussian priors:

$$\begin{aligned} p(f_j(t)|v_{0j}) &= \mathcal{N}(0, v_{0j}) \\ p(\mathbf{f}(t)|v_0) &\propto \exp \left\{ -\frac{1}{2} \sum_j \mathbf{f}_j^2(t)/v_{0j} \right\} \\ p(\mathbf{A}_{ij}|\mathbf{A}_{0ij}, V_{0ij}) &= \mathcal{N}(\mathbf{A}_{0ij}, V_{0ij}) \\ p(\mathbf{A}|\mathbf{A}_0, V_0) &= \mathcal{N}(\mathbf{A}_0, V_0) \\ p(\mathbf{g}(t)|\mathbf{A}, \mathbf{f}(t), v_\varepsilon) &= \mathcal{N}(\mathbf{A}\mathbf{f}(t), v_\varepsilon \mathbf{I}) \end{aligned} \quad (55)$$



We also assume the following:

- All sources a priori same variance v_f and note $v_0 = [v_f, \dots, v_f]'$.
- All noise terms a priori same variance v_ε and note $v_\varepsilon = [v_\varepsilon, \dots, v_\varepsilon]'$
- $\mathbf{A}_0 = 0$, and $V_0 = v_a \mathbf{I}$

Then, we have:

$$p(\mathbf{f}(t)|\mathbf{g}(t), \mathbf{A}, v_\varepsilon, v_0) = \mathcal{N}(\hat{\mathbf{f}}(t), \hat{\boldsymbol{\Sigma}}) \text{ with } \begin{cases} \hat{\boldsymbol{\Sigma}} = (\mathbf{A}'\mathbf{A} + \lambda_f \mathbf{I})^{-1} \\ \hat{\mathbf{f}}(t) = (\mathbf{A}'\mathbf{A} + \lambda_f \mathbf{I})^{-1} \mathbf{A}'\mathbf{g}(t) \\ \lambda_f = v_\varepsilon/v_f \end{cases} \quad (56)$$

$$p(\mathbf{A}|\mathbf{g}(t), \mathbf{f}(t), v_\varepsilon, \mathbf{A}_0, V_0) = \mathcal{N}(\widehat{\mathbf{A}}, \widehat{\mathbf{V}}) \text{ with } \begin{cases} \widehat{\mathbf{V}} = (\mathbf{F}'\mathbf{F} + \lambda_f \mathbf{I})^{-1} \\ \widehat{\mathbf{A}} = \sum_t \mathbf{g}(t) \mathbf{f}'(t) (\sum_t \mathbf{f}(t) \mathbf{f}'(t) + \lambda_a \mathbf{I})^{-1} \\ \lambda_a = v_\varepsilon / v_a \end{cases} \quad (57)$$

The joint MAP algorithm becomes:

$$\begin{array}{ccc} \mathbf{A}^{(0)} \rightarrow \widehat{\mathbf{A}} \rightarrow & \boxed{(\widehat{\mathbf{A}}'\widehat{\mathbf{A}} + \lambda_f \mathbf{I})^{-1} \widehat{\mathbf{A}}'\mathbf{g}} \rightarrow & \widehat{\mathbf{f}}(t) \\ \uparrow & & \downarrow \\ \widehat{\mathbf{A}} \leftarrow & \boxed{\sum_t \mathbf{g}(t) \widehat{\mathbf{f}}'(t) (\sum_t \widehat{\mathbf{f}}(t) \widehat{\mathbf{f}}'(t) + \lambda_a \mathbf{I})^{-1}} \leftarrow & \widehat{\mathbf{f}}(t) \end{array} \quad (58)$$

6 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 $\boldsymbol{\theta}$ and write the joint posterior law of all the unknowns:

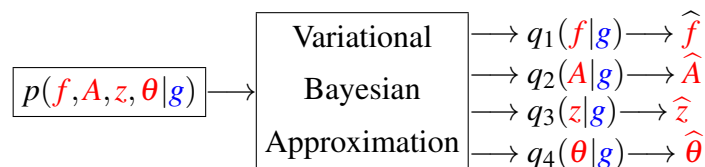
$$p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}|\mathbf{g}) \propto p(\mathbf{g}|\mathbf{f}, \mathbf{A}, \boldsymbol{\theta}_1) p(\mathbf{f}|\mathbf{z}, \boldsymbol{\theta}_2) p(\mathbf{z}|\boldsymbol{\theta}_3) p(\mathbf{A}|\boldsymbol{\theta}_4) p(\boldsymbol{\theta}) \quad (59)$$

where $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \boldsymbol{\theta}_3, \boldsymbol{\theta}_4)$ represents all the hyper-parameters of the problem.

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

$$q(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}|\mathbf{g}) = q_1(\mathbf{f}|\mathbf{g}) q_2(\mathbf{A}|\mathbf{g}) q_3(\mathbf{z}|\mathbf{g}) q_4(\boldsymbol{\theta}|\mathbf{g}) \quad (60)$$

illustrated here:



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

$$\text{KL}(q : p) = \int q \ln p/q = \langle \ln p/q \rangle_q \quad (61)$$

It is then easy to show that

$$\text{KL}(q : p) = \ln p(\mathbf{g}|\mathcal{M}) - \mathcal{F}(q) \quad (62)$$

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

$$p(\mathbf{g}|\mathcal{M}) = \int \int \int \int p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}|\mathcal{M}) d\mathbf{f} d\mathbf{A} d\mathbf{z} d\boldsymbol{\theta} \quad (63)$$

with

$$p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}|\mathcal{M}) = p(\mathbf{g}|\mathbf{f}, \mathbf{A}, \boldsymbol{\theta}_1) p(\mathbf{A}|\boldsymbol{\theta}_2) p(\mathbf{f}|\mathbf{z}, \boldsymbol{\theta}_3) p(\mathbf{z}|\boldsymbol{\theta}_4) p(\boldsymbol{\theta}) \quad (64)$$

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

$$\mathcal{F}(q) = \left\langle \ln \frac{p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta})} \right\rangle_q \quad (65)$$

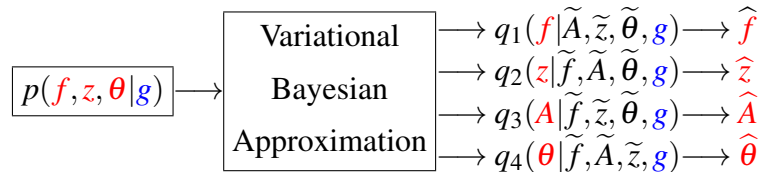
So, for a given model \mathcal{M} , minimizing $\text{KL}(q : p)$ is equivalent to maximizing $\mathcal{F}(q)$ and when optimized, $\mathcal{F}(q^*)$ gives a lower bound for $\ln p(\mathbf{g} | \mathcal{M})$.

An alternate optimization of $\mathcal{F}(q)$ with respect to q_1, q_2, q_3 and q_4 results in

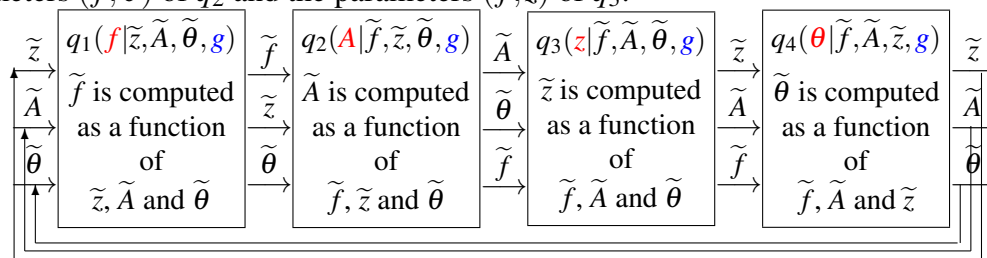
$$\begin{aligned} q_1(\mathbf{f}) &\propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_2(\mathbf{A}) q_3(\mathbf{z}) q_4(\boldsymbol{\theta})} \right\} \\ q_2(\mathbf{A}) &\propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_1(\mathbf{f}) q_3(\mathbf{z}) q_4(\boldsymbol{\theta})} \right\} \\ q_3(\mathbf{z}) &\propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_1(\mathbf{f}) q_2(\mathbf{A}) q_4(\boldsymbol{\theta})} \right\} \\ q_4(\boldsymbol{\theta}) &\propto \exp \left\{ - \langle \ln p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g}) \rangle_{q_1(\mathbf{f}) q_2(\mathbf{A}) q_3(\mathbf{z})} \right\} \end{aligned}$$

Note that these relations represent an implicit solution for $q_1(\mathbf{f}), q_2(\mathbf{A}), q_3(\mathbf{z})$ and $q_4(\boldsymbol{\theta})$. For conjugate conditional distributions we consider, these leads to standard distributions for which the required expectations are easily evaluated. In that case, we may note

$$q(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g}) = q_1(\mathbf{f} | \tilde{\mathbf{A}}, \tilde{\mathbf{z}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) q_2(\mathbf{A} | \tilde{\mathbf{f}}, \tilde{\mathbf{A}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) q_3(\mathbf{z} | \tilde{\mathbf{f}}, \tilde{\mathbf{A}}, \tilde{\boldsymbol{\theta}}; \mathbf{g}) q_4(\boldsymbol{\theta} | \tilde{\mathbf{f}}, \tilde{\mathbf{A}}, \tilde{\mathbf{z}}; \mathbf{g}) \quad (66)$$



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

$$\begin{aligned} \mathcal{F}(q) &= \left\langle \ln \frac{p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M})}{q(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta})} \right\rangle_q \\ &= \langle \ln p(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}, \mathbf{g} | \mathcal{M}) \rangle_q + \langle - \ln q(\mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q \\ &= \langle \ln p(\mathbf{g} | \mathbf{f}, \mathbf{A}, \mathbf{z}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{f} | \mathbf{z}, \mathbf{A}, \boldsymbol{\theta}) \rangle_q + \langle \ln p(\mathbf{z} | \boldsymbol{\theta}) \rangle_q \\ &\quad + \langle - \ln q(\mathbf{f}) \rangle_q + \langle - \ln q(\mathbf{z}) \rangle_q + \langle - \ln q(\boldsymbol{\theta}) \rangle_q \end{aligned} \quad (67)$$

where all the expectations are with respect to q .

Other decompositions are also possible. To illustrate this approach, here we consider only the simple case of sources separation with Gaussian priors. For a more extensive examples in particular for inverse problems, the readers may refer to [7].

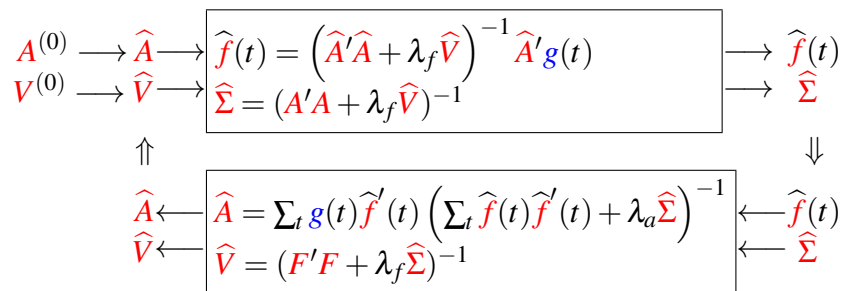
6.1 Variational Bayesian Approximation for sources separation

Here the main objective is to approximate $p(\mathbf{f}_{1..T}, \mathbf{A} | \mathbf{g}_{1..T})$ by $q_1(\mathbf{f}_{1..T} | \mathbf{A}, \mathbf{g}_{1..T}) q_2(\mathbf{A} | \mathbf{f}_{1..T}, \mathbf{g}_{1..T})$.

The solution is:

$$\begin{aligned}
 q_1(\mathbf{f}(t) | \mathbf{g}(t), \mathbf{A}, v_\varepsilon, v_0) &= \mathcal{N}(\hat{\mathbf{f}}(t), \hat{\Sigma}) & \text{with} & \begin{cases} \hat{\Sigma} = (\mathbf{A}'\mathbf{A} + \lambda_f \hat{\mathbf{V}})^{-1} \\ \hat{\mathbf{f}}(t) = (\mathbf{A}'\mathbf{A} + \lambda_f \hat{\mathbf{V}})^{-1} \mathbf{A}'\mathbf{g}(t), \\ \lambda_f = v_\varepsilon / v_f \\ \hat{\mathbf{V}} = (\mathbf{F}'\mathbf{F} + \lambda_f \hat{\Sigma})^{-1} \end{cases} \\
 q_2(\mathbf{A} | \mathbf{g}(t), \mathbf{f}(t), v_\varepsilon, \mathbf{A}_0, \mathbf{V}_0) &= \mathcal{N}(\hat{\mathbf{A}}, \hat{\mathbf{V}}) & \text{with} & \begin{cases} \hat{\mathbf{A}} = \sum_t \mathbf{g}(t) \mathbf{f}'(t) \left(\sum_t \mathbf{f}(t) \mathbf{f}'(t) + \lambda_a \hat{\Sigma} \right)^{-1} \\ \lambda_a = v_\varepsilon / v_a \end{cases}
 \end{aligned} \tag{68}$$

Algorithm:



We may compare this algorithm with the joint MAP in 58. As we can see, in 58 the mixing matrix and the sources are updated without accounting for the uncertainties, but here the uncertainties of each step of the estimation is also updated and used in the next iteration.

7 Conclusion

The main aim of this paper was first to present the Bayesian inference approach for blind sources separation (BSS) where we want to infer on the mixing matrix, the sources and all the hyper-parameters associated to probabilistic modeling. For this purpose, the sources separation problem was considered in a step by step approach: i) Sources estimation when the mixing matrix is known; ii) Mixing matrix estimation when the sources are known; iii) Joint estimation of the mixing matrix and the sources; and finally, iv) Joint estimation of the mixing matrix, the sources and the associated hidden variables and hyper-parameters. One of the main steps in any of these steps is modeling of sources and the mixing matrix prior laws. Here, we focused on the sparsity enforcing probability laws (such as Generalized Gaussian, Student-t and mixture models) both for the sources and the mixing matrix. For algorithmic and computational aspects, we considered either Joint MAP, MCMC Gibbs sampling or Variational Bayesian Approximation (VBA)

tools. For each class of methods we discuss about their relative costs and performances.

References

- [1] C. Bouman and K. Sauer. “A generalized gaussian image model for edge-preserving map estimation.” *IEEE Trans. on Medical Imaging*, MI-2(3), 296–310, 1993.
- [2] S. Chatzis and T. 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] N. CHU, A. M. 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*, page No.18. Berlin, Germany, Feb.22-23,2012.
- [4] 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.
- [5] J. Griffin and P. Brown. “Inference with normal-gamma prior distributions in regression problems.” *Bayesian Analysis*, 2010.
- [6] J. R. H. Ishwaran. “Spike and Slab variable selection: Frequentist and Bayesian strategies.” *Annals of Statistics*, 2005.
- [7] 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, 2012.
- [8] T. Park and G. Casella. “The Bayesian Lasso.” *Journal of the American Statistical Association*, 2008.
- [9] N. Polson and J. Scott. “Shrink globally, act locally: sparse Bayesian regularization and prediction.” *Bayesian Statistics 9*, 2010.
- [10] H. Snoussi and J. Idier. “Bayesian blind separation of generalized hyperbolic processes in noisy and underdeterminate mixtures.” *IEEE Trans. on Signal Processing*, 2006.
- [11] M. Tipping. “Sparse Bayesian learning and the relevance vector machine.” *Journal of Machine Learning Research*, 2001.
- [12] S. Zhu, A. Mohammad-Djafari, H. Wang, B. Deng, X. Li, and J. Mao. “Parameter estimation for sar micromotion target based on sparse signal representation.” *EURASIP Journal on Advances in Signal Processing*, Special issue on Sparse Signal Processing, 2012.