



# Variational Bayesian Approximation methods for inverse problems

Ali Mohammad-Djafari

Laboratoire des Signaux et Systèmes,  
UMR8506 CNRS-SUPELEC-UNIV PARIS SUD 11  
SUPELEC, 91192 Gif-sur-Yvette, France  
<http://lss.supelec.free.fr>

Email: [djafari@lss.supelec.fr](mailto:djafari@lss.supelec.fr)  
<http://djafari.free.fr>

# 1. General inverse problem

$$g(t) = \mathcal{H}f(t) + \epsilon(t), \quad t \in [1, \dots, T]$$

$$g(\mathbf{r}) = \mathcal{H}f(\mathbf{r}) + \epsilon(\mathbf{r}), \quad \mathbf{r} = (x, y) \in R^2$$

- ▶  $f$  unknown quantity (input)
- ▶  $\mathcal{H}$  Forward operator:  
(Convolution, Radon, Fourier or any Linear operator)
- ▶  $g$  observed quantity (output)
- ▶  $\epsilon$  represents the errors of modeling and measurement

Discretization:

$$g = \mathbf{H}f + \epsilon$$

- ▶ Forward operation  $\mathbf{H}f$
- ▶ Adjoint operation  $\mathbf{H}'g$  :  $\langle \mathbf{H}'g, f \rangle = \langle \mathbf{H}f, g \rangle$
- ▶ Inverse operation (if exists)  $\mathbf{H}^{-1}g$

## 2. General Bayesian Inference

- ▶ Bayesian inference:

$$p(\mathbf{f}|\mathbf{g}, \boldsymbol{\theta}) = \frac{p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f}|\boldsymbol{\theta}_2)}{p(\mathbf{g}|\boldsymbol{\theta})}$$

with  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$

- ▶ Point estimators:

Maximum A Posteriori (MAP) or Posterior Mean (PM)  $\rightarrow \hat{\mathbf{f}}$

- ▶ Full Bayesian inference:

- ▶ Simple prior models:  $p(\mathbf{f}|\boldsymbol{\theta}_2)$

$$q(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) \propto p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f}|\boldsymbol{\theta}_2) p(\boldsymbol{\theta})$$

- ▶ Prior models with hidden variables:  $p(\mathbf{f}|\mathbf{z}, \boldsymbol{\theta}_2) p(\mathbf{z}|\boldsymbol{\theta}_3)$

$$q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}|\mathbf{g}) \propto p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f}|\boldsymbol{\theta}_2) p(\mathbf{z}|\boldsymbol{\theta}_3) p(\boldsymbol{\theta})$$

### 3. Sparsity enforcing prior models

- ▶ Simple heavy tailed models:
  - ▶ Generalized Gaussian, Double Exponential
  - ▶ Student-t, Cauchy
  - ▶ Elastic net
  
  - ▶ Symmetric Weibull, Symmetric Rayleigh
  - ▶ Generalized hyperbolic
  
- ▶ Hierarchical mixture models:
  - ▶ Mixture of Gaussians
  - ▶ Bernoulli-Gaussian
  
  - ▶ Mixture of Gammas
  - ▶ Bernoulli-Gamma
  - ▶ Mixture of Dirichlet
  - ▶ Bernoulli-Multinomial

## 4. Simple heavy tailed models

- Generalized Gaussian, Double Exponential

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

$\beta = 1$  Double exponential or Laplace.

$0 < \beta \leq 1$  are of great interest for sparsity enforcing.

- Student-t and Cauchy models

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

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

- Elastic net prior model

$$p(\mathbf{f}|\nu) = \prod \mathcal{EN}(f_j|\nu) \propto \exp \left\{ -\sum (\gamma_1 |f_j| + \gamma_2 f_j^2) \right\}$$

## 5 Mixture models

- Mixture of two Gaussians (MoG2) model

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

- Bernoulli-Gaussian (BG) model

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

- Mixture of Gammas

$$p(\mathbf{f}|\lambda, v_1, v_0) = \prod_j (\lambda \mathcal{G}(f_j|\alpha_1, \beta_1) + (1 - \lambda)\mathcal{G}(f_j|\alpha_2, \beta_2))$$

- Bernoulli-Gamma model

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

## 6. MAP, Joint MAP

- ▶ Inverse problems:  $\mathbf{g} = \mathbf{H}\mathbf{f} + \boldsymbol{\epsilon}$
- ▶ Posterior law:

$$p(\mathbf{f}|\boldsymbol{\theta}, \mathbf{g}) \propto p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f}|\boldsymbol{\theta}_2)$$

- ▶ Examples:

Gaussian noise, Gaussian prior and MAP:

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\} \quad \text{with } J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|_2^2 + \lambda \|\mathbf{f}\|_2^2$$

Gaussian noise, Double Exponential prior and MAP:

$$\hat{\mathbf{f}} = \arg \min_{\mathbf{f}} \{J(\mathbf{f})\} \quad \text{with } J(\mathbf{f}) = \|\mathbf{g} - \mathbf{H}\mathbf{f}\|_2^2 + \lambda \|\mathbf{f}\|_1$$

- ▶ Full Bayesian: Joint Posterior:

$$p(\mathbf{f}, \boldsymbol{\theta}|\mathbf{g}) \propto p(\mathbf{g}|\mathbf{f}, \boldsymbol{\theta}_1) p(\mathbf{f}|\boldsymbol{\theta}_2) p(\boldsymbol{\theta})$$

- ▶ Joint MAP:

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

## 7. Marginal MAP and PM estimates

- ▶ Marginal MAP:  $\hat{\theta} = \arg \max_{\theta} \{p(\theta|g)\}$  where

$$p(\theta|g) = \int p(\mathbf{f}, \theta|g) d\mathbf{f} = \int p(g|\mathbf{f}, \theta_1) p(\mathbf{f}|\theta_2) d\mathbf{f}$$

and then  $\hat{\mathbf{f}} = \arg \max_{\mathbf{f}} \{p(\mathbf{f}|\hat{\theta}, g)\}$

- ▶ Posterior Mean:  $\hat{\mathbf{f}} = \int \mathbf{f} p(\mathbf{f}|\hat{\theta}, g) d\mathbf{f}$

- ▶ EM and GEM Algorithms

- ▶ Variational Bayesian Approximation:

Approximate  $p(\mathbf{f}, \theta|g)$  by  $q(\mathbf{f}, \theta|g) = q_1(\mathbf{f}|g) q_2(\theta|g)$   
and then continue computations.



## 8. Hierarchical models and hidden variables

- ▶ All the mixture models and some of simple models can be modeled via **hidden variables**  $\mathbf{z}$ .
- ▶ Example 1: Student-t model

$$\begin{cases} p(\mathbf{f}|\mathbf{z}) = \prod_j p(f_j|z_j) = \prod_j \mathcal{N}(f_j|0, 1/z_j) \propto \exp \left\{ -\frac{1}{2} \sum_j z_j f_j^2 \right\} \\ p(z_j|a, b) = \mathcal{G}(z_j|a, b) \propto z_j^{(a-1)} \exp \{-bz_j\} \text{ with } a = b = \nu/2 \end{cases}$$

- ▶ Example 2: MoG model:

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

- ▶ With these models we have:

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

## 9. Bayesian Computation and Algorithms

- ▶ Often, the expression of  $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g})$  is complex.
- ▶ Its optimization (for Joint MAP) or its marginalization or integration (for Marginal MAP or PM) is not easy
- ▶ Two main techniques:  
MCMC and Variational Bayesian Approximation (VBA)
- ▶ MCMC:  
Needs the expressions of the conditionals  $p(\mathbf{f} | \mathbf{z}, \boldsymbol{\theta}, \mathbf{g})$ ,  $p(\mathbf{z} | \mathbf{f}, \boldsymbol{\theta}, \mathbf{g})$ , and  $p(\boldsymbol{\theta} | \mathbf{f}, \mathbf{z}, \mathbf{g})$
- ▶ VBA: Approximate  $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta} | \mathbf{g})$  by a separable one

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

and do any computations with these separable ones.

## 10. Bayesian Variational Approximation

- ▶ Objective: Approximate  $p(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}|\mathcal{g})$  by a separable one  $q(\mathbf{f}, \mathbf{z}, \boldsymbol{\theta}|\mathcal{g}) = q_1(\mathbf{f}) q_2(\mathbf{z}) q_3(\boldsymbol{\theta})$

- ▶ Criterion:

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

- ▶ Free energy:  $\text{KL}(q : p) = \ln p(\mathcal{g}|\mathcal{M}) - \mathcal{F}(q)$  where:

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

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

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

- ▶ 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  $\ln p(\mathcal{g}|\mathcal{M})$ .

## 11. BVA with Student-t priors

Scale Mixture Model of Student-t:

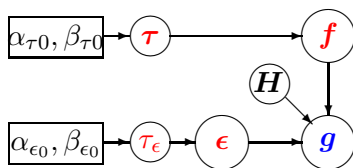
$$St(f_j|\nu) = \int_0^\infty \mathcal{N}(f_j|0, 1/\tau_j) \mathcal{G}(\tau_j|\nu/2, \nu/2) d\tau_j$$

Hidden variables  $\tau_j$ :

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

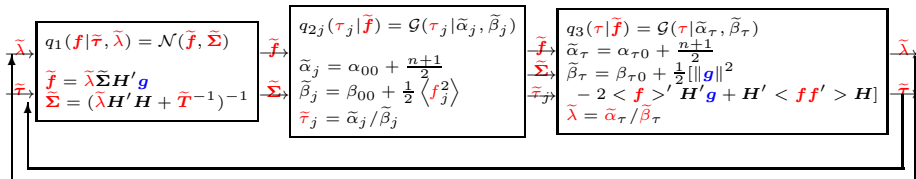
Cauchy model is obtained when  $\nu = 1$ :

► Graphical model:



## 12. BVA with Student-t priors Algorithm

$$\begin{cases}
 p(\mathbf{g}|\mathbf{f}, \tau_\epsilon) = \mathcal{N}(\mathbf{g}|\mathbf{H}\mathbf{f}, (1/\tau_\epsilon)\mathbf{I}) \\
 p(\tau_\epsilon|\alpha_{\tau 0}, \beta_{\tau 0}) = \mathcal{G}(\tau_\epsilon|\alpha_{\tau 0}, \beta_{\tau 0}) \\
 p(\mathbf{f}|\boldsymbol{\tau}) = \prod_j \mathcal{N}(\mathbf{f}_j|0, 1/\tau_j) \\
 p(\boldsymbol{\tau}|\alpha_0, \beta_0) = \prod_j \mathcal{G}(\tau_j|\alpha_0, \beta_0) \\
 \\
 q_1(\mathbf{f}|\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) = \mathcal{N}(\mathbf{f}|\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \\
 \tilde{\boldsymbol{\mu}} = \langle \lambda \rangle \tilde{\boldsymbol{\Sigma}} \mathbf{H}' \mathbf{g} \\
 \tilde{\boldsymbol{\Sigma}} = (\langle \lambda \rangle \mathbf{H}' \mathbf{H} + \tilde{\mathbf{Z}})^{-1}, \\
 \text{with } \tilde{\mathbf{Z}} = \tilde{\mathbf{T}}^{-1} = \text{diag}[\tilde{\tau}]
 \end{cases}
 \begin{cases}
 q_{2j}(\tau_j) = \mathcal{G}(\tau_j|\tilde{\alpha}_j, \tilde{\beta}_j) \\
 \tilde{\alpha}_j = \alpha_{00} + 1/2 \\
 \tilde{\beta}_j = \beta_{00} + \langle \mathbf{f}_j^2 \rangle / 2 \\
 \\
 q_3(\tau_\epsilon) = \mathcal{G}(\tau_\epsilon|\tilde{\alpha}_{\tau_\epsilon}, \tilde{\beta}_{\tau_\epsilon}), \\
 \tilde{\alpha}_{\tau_\epsilon} = \alpha_{\tau 0} + (n+1)/2 \\
 \tilde{\beta}_{\tau_\epsilon} = \beta_{\tau 0} + 1/2[\|\mathbf{g}\|^2 \\
 - 2 \langle \mathbf{f}' \rangle' \mathbf{H}' \mathbf{g} + \mathbf{H}' \langle \mathbf{f} \mathbf{f}' \rangle \mathbf{H}]
 \end{cases}
 \begin{cases}
 \langle \mathbf{f} \rangle = \tilde{\boldsymbol{\mu}} \\
 \langle \mathbf{f} \mathbf{f}' \rangle = \tilde{\boldsymbol{\Sigma}} + \tilde{\boldsymbol{\mu}} \tilde{\boldsymbol{\mu}}' \\
 \langle \mathbf{f}_j^2 \rangle = [\tilde{\boldsymbol{\Sigma}}]_{jj} + \tilde{\mu}_j^2 \\
 \\
 \tilde{\lambda} = \tilde{\alpha}_\tau / \tilde{\beta}_\tau \\
 \\
 \tilde{\tau}_j = \tilde{\alpha}_j / \tilde{\beta}_j
 \end{cases}$$



## 13. Implementation issues

- ▶ In inverse problems, often we do not have access directly to the matrix  $\mathbf{H}$ . But, we can compute:
  - ▶ Forward operator :  $\mathbf{H}\mathbf{f} \rightarrow \mathbf{g}$      $\mathbf{g}=\text{direct}(\mathbf{f}, \dots)$
  - ▶ Adjoint operator :  $\mathbf{H}'\mathbf{g} \rightarrow \mathbf{f}$      $\mathbf{f}=\text{transp}(\mathbf{g}, \dots)$
- ▶ For any particular application, we can always write two programs (`direct` & `transp`) corresponding to the application of these two operators.
- ▶ To compute  $\tilde{\mathbf{f}}$ , we use a gradient based optimization algorithm which will use these operators.
- ▶ We may also need to compute the diagonal elements of  $[\mathbf{H}'\mathbf{H}]$ . We also developed algorithms which computes these diagonal elements with the same programs (`direct` & `transp`)

## 14. Conclusions and Perspectives

- ▶ We proposed a list of **different probabilistic prior models** which can be used for **sparsity enforcing**.
- ▶ We classified these models in two categories: **simple heavy tails** and **hierarchical mixture models**
- ▶ We showed **how to use these models for inverse problems where the desired solutions are sparse**
- ▶ Different algorithms have been developed and their relative performances are compared.
- ▶ We use these models for inverse problems in different signal and image processing applications such as:
  - ▶ **Period estimation in biological time series**
  - ▶ **X ray Computed Tomography,**
  - ▶ **Signal deconvolution in Proteomic and molecular imaging**
  
  - ▶ **Diffraction Optical Tomography**
  - ▶ **Microwave Imaging, Acoustic imaging and sources localization**
  - ▶ **Synthetic Aperture Radar (SAR) Imaging**

# 15. References

1. A. Mohammad-Djafari, "Bayesian approach with prior models which enforce sparsity in signal and image processing," *EURASIP Journal on Advances in Signal Processing*, vol. Special issue on Sparse Signal Processing, (2012).
2. S. Zhu, A. Mohammad-Djafari, H. Wang, B. Deng, X. Li and J. Mao J, "Parameter estimation for SAR micromotion target based on sparse signal representation," *EURASIP Journal on Advances in Signal Processing*, vol. Special issue on Sparse Signal Processing, (2012).
3. N. Chu, J. Picheral and A. Mohammad-Djafari, "A robust super-resolution approach with sparsity constraint for near-field wideband acoustic imaging," *IEEE International Symposium on Signal Processing and Information Technology* pp 286–289, Bilbao, Spain, Dec14-17,2011
4. N. Bali and A. Mohammad-Djafari, "Bayesian Approach With Hidden Markov Modeling and Mean Field Approximation for Hyperspectral Data Analysis," *IEEE Trans. on Image Processing* 17: 2. 217-225 Feb. (2008).
5. J. Griffin and P. Brown, "Inference with normal-gamma prior distributions in regression problems," *Bayesian Analysis*, 2010.
6. N. Polson and J. Scott., "Shrink globally, act locally: sparse Bayesian regularization and prediction," *Bayesian Statistics 9*, 2010.
7. T. Park and G. Casella., "The Bayesian Lasso," *Journal of the American Statistical Association*, 2008.
8. C. Févotte and S. Godsill, "A Bayesian approach for blind separation of sparse source," *IEEE Transactions on Audio, Speech, and Language processing*, 2006.
9. H. Snoussi and J. Idier., "Bayesian blind separation of generalized hyperbolic processes in noisy and underdeterminate mixtures," *IEEE Trans. on Signal Processing*, 2006.
10. J. R. H. Ishwaran, "Spike and slab variable selection: Frequentist and Bayesian strategies," *Annals of Statistics*, 2005.
11. M. Tipping, "Sparse Bayesian learning and the relevance vector machine," *Journal of Machine Learning Research*, 2001.



# 14. References ..

1. Abib. Doucet and P. Duvaut, "Bayesian estimation of state-space models applied to deconvolution of Bernoulli-Gaussian processes," *Signal Processing*, vol. 57, no. 2, 1997.
2. P. Williams, "Bayesian regularization and pruning using a Laplace prior," *Neural Computation*, 1995.
3. M. Lavielle, "Bayesian deconvolution of Bernoulli-Gaussian processes," *Signal Processing*, vol. 33, pp. 67-79, 1993.
4. T. Mitchell and J. Beauchamp, "Bayesian variable selection in linear regression," *Journal of the American Statistical Association*, 1988.
5. J. J. Kormylo and J. M. Mendel, "Maximum-likelihood detection and estimation of Bernoulli-Gaussian processes," vol. 28, pp. 482-488, 1982.
6. H. Snoussi and A. Mohammad-Djafari, " Estimation of Structured Gaussian Mixtures: The Inverse EM Algorithm," *IEEE Trans. on Signal Processing* 55: 7. 3185-3191 July (2007).
7. N. Bali and A. Mohammad-Djafari, "A variational Bayesian Algorithm for BSS Problem with Hidden Gauss-Markov Models for the Sources," in: *Independent Component Analysis and Signal Separation (ICA 2007)* Edited by: M.E. Davies, Ch.J. James, S.A. Abdallah, M.D. Plumbley. 137-144 Springer (LNCS 4666) (2007).
8. N. Bali and A. Mohammad-Djafari, "Hierarchical Markovian Models for Joint Classification, Segmentation and Data Reduction of Hyperspectral Images" *ESANN 2006*, September 4-8, Belgium. (2006)
9. M. Ichir and A. Mohammad-Djafari, "Hidden Markov models for wavelet-based blind source separation," *IEEE Trans. on Image Processing* 15: 7. 1887-1899 July (2005)
10. S. Moussaoui, C. Carteret, D. Brie and A. Mohammad-Djafari, "Bayesian analysis of spectral mixture data using Markov Chain Monte Carlo methods sampling," *Chemometrics and Intelligent Laboratory Systems* 81: 2. 137-148 (2005).
11. H. Snoussi and A. Mohammad-Djafari, "Fast joint separation and segmentation of mixed images" *Journal of Electronic Imaging* 13: 2. 349-361 April (2004)
12. H. Snoussi and A. Mohammad-Djafari, "Bayesian unsupervised learning for source separation with mixture of Gaussians prior," *Journal of VLSI Signal Processing Systems* 37: 2/3. 263-279 June/July (2004)