



Probabilistic models which enforce sparsity

Ali Mohammad-Djafari

Laboratoire des Signaux et Systèmes,
UMR8506 CNRS-SUPELEC-UNIV PARIS SUD 11
SUPELEC, 91192 Gif-sur-Yvette, France

Email: djafari@lss.supelec.fr

<http://djafari.free.fr>

Summary

In this paper, 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. Here we consider it directly on the original space (impulsive signals). These models are either **simple heavy tailed** or **hierarchical mixture models**.

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 and **drive the corresponding appropriate algorithms**, and discuss on their corresponding relative **complexities and performances**.

1. Bayesian inference for inverse problems

- ▶ Inverse problems: $\mathbf{g} = \mathbf{H}\mathbf{f} + \epsilon$
- ▶ 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), Posterior Mean (PM)
- ▶ 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: $p(\mathbf{f}|\mathbf{z}, \boldsymbol{\theta}_2) p(\mathbf{z}|\boldsymbol{\theta}_3)$

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

2. Sparsity enforcing prior models

- ▶ Simple heavy tailed models:
 - ▶ Generalized Gaussian, Double Exponential
 - ▶ Symmetric Weibull, Symmetric Rayleigh
 - ▶ Student-t, Cauchy
 - ▶ Generalized hyperbolic
 - ▶ Elastic net

- ▶ Hierarchical mixture models:
 - ▶ Mixture of Gaussians
 - ▶ Bernoulli-Gaussian
 - ▶ Mixture of Gammas
 - ▶ Bernoulli-Gamma
 - ▶ Mixture of Dirichlet
 - ▶ Bernoulli-Multinomial

3. 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.

- Symmetric Weibull

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

$\beta = 2$ is the Symmetric Rayleigh distribution.

$\beta = 1$ is the Double exponential and

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

- Student-t and Cauchy models

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

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

- Elastic net prior model

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

- Generalized hyperbolic (GH) models

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

4. 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

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

- Mixture of Dirichlets model

$$p(\mathbf{f}|\lambda, \mathbf{a}_1, \boldsymbol{\alpha}_1, \mathbf{a}_2, \boldsymbol{\alpha}_2) = \prod_j \lambda \mathcal{D}(f_j|\mathbf{a}_1, \boldsymbol{\alpha}_1) + (1 - \lambda) \mathcal{D}(f_j|\mathbf{a}_2, \boldsymbol{\alpha}_2)$$

where

$$\mathcal{D}(f_j|\mathbf{a}, \boldsymbol{\alpha}) = \prod_{k=1}^K \frac{\Gamma(\alpha_k)}{\Gamma(\alpha_0)\Gamma(\alpha_K)} a_k^{\alpha_k - 1}, \quad \alpha_k \geq 0, \quad a_k \geq 0$$

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

- Bernoulli-Multinomial (BMultinomial) model

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

5. MAP, Joint MAP

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

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

- ▶ Example: Gaussian noise, Double Exponential prior and MAP:

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

- ▶ 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:

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

6. 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.

7. Hierarchical models and hidden variables

- ▶ All the mixture models and some of simple models can be modeled via **hidden variables** z .
- ▶ Example 1: MoG model:

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

- ▶ Example 2: Student-t model

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

- ▶ With these models we have:

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

8. 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.

9. 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:
 - ▶ **Synthetic Aperture Radar (SAR) Imaging**
 - ▶ **Signal deconvolution in Proteomic and molecular imaging**
 - ▶ **X ray Computed Tomography, Diffraction Optical Tomography, Microwave Imaging, ...**

10. Main references



A. Doucet and P. Duvaut, "Bayesian estimation of state-space models applied to deconvolution of Bernoulli-Gaussian processes," *Signal Processing*, vol. 57, no. 2, 1997.



T. Park and G. Casella., "The Bayesian Lasso," *Journal of the American Statistical Association*, 2008.



M. Tipping, "Sparse Bayesian learning and the relevance vector machine," *Journal of Machine Learning Research*, 2001.



C. Févotte and S. Godsill, "A Bayesian approach for blind separation of sparse source," *IEEE Transactions on Audio, Speech, and Language processing*, 2006.



J. Griffin and P. Brown, "Inference with normal-gamma prior distributions in regression problems," *Bayesian Analysis*, 2010.



N. Polson and J. Scott., "Shrink globally, act locally: sparse Bayesian regularization and prediction," *Bayesian Statistics 9*, 2010.



H. Snoussi and J. Idier., "Bayesian blind separation of generalized hyperbolic processes in noisy and underdeterminate mixtures," *IEEE Trans. on Signal Processing*, 2006.



P. Williams, "Bayesian regularization and pruning using a Laplace prior," *Neural Computation*, 1995.



T. Mitchell and J. Beauchamp, "Bayesian variable selection in linear regression," *Journal of the American Statistical Association*, 1988.



J. R. H. Ishwaran, "Spike and slab variable selection: Frequentist and bayesian strategies," *Annals of Statistics*, 2005.



J. J. Kormylo and J. M. Mendel, "Maximum-likelihood detection and estimation of Bernoulli-Gaussian processes," vol. 28, pp. 482–488, 1982.



M. Lavielle, "Bayesian deconvolution of Bernoulli-Gaussian processes," *Signal Processing*, vol. 33, pp. 67–79, 1993.