## A Bayesian Method for Sparse Regression IMACS 2013

#### Amandine Schreck under the supervision of Gersende Fort and Eric Moulines, joint work with Sylvain Le Corff

Télécom ParisTech

Thursday, July 18th, 2013

A B A A B

The problem Outline

#### Goal : brain imaging - locate activated zones in a brain



Collaboration with Alexandre Gramfort on brain imaging problems.

**The problem** Outline



<ロ> (四) (四) (三) (三) (三)

2

The problem Outline

- Goal: find the active (i.e. non-zero) components of the sparse signal decomposition.
- Difficulty: high dimensional setting, potentially low number of observations, high number of regressors.
- Existing solutions: deterministic methods (e.g. ISTA), transdimensional MCMC methods (Reversible Jump, Metropolised Carlin and Chib).

(日) (同) (日) (日) (日)

The problem Outline

#### Specification of the problem

- The simplified model
- The Bayesian model selection framework

#### 2 The PMALA algorithm

- Two main ingredients
- The algorithm

#### ③ Illustration

Toy examples and simulated data

## 4 Futur directions

The simplified model The Bayesian model selection framework

Simplified model :

$$Y = GX + \sqrt{\tau}E \; ,$$

where

- $Y \in \mathbb{R}^{N \times T}$  is the observed signal
- $G \in \mathbb{R}^{N \times P}$  is the design matrix (known)
- $X \in \mathbb{R}^{P \times T}$  is the emitted signal, directly assumed to be sparse
- $E \in \mathbb{R}^{N imes T}$  is a standard Gaussian noise

For concision of notations: T = 1.

4 日 5 - 4 周 5 - 4 戸 5 - 4 戸 5

- X can be equivalently defined by  $(m, X_m)$  where
  - $m = (m_1, \cdots, m_P) \in \mathcal{M} = \{0, 1\}^P$  is the model, with  $m_i = 0$  iff  $X_i = 0$ ,
  - $X_m \in \mathbb{R}^{|m|}$  collects the active rows of X, where  $|m| = \sum_i m_i$ .

 $\rightarrow$  Sampling set:

$$\Theta = \bigcup_{m \in \mathcal{M}} \left( \{m\} imes \mathbb{R}^{|m|} 
ight) \;.$$

(日) (同) (日) (日)

э

Likelihood and prior distributions:

• 
$$\pi(Y|m, X_m) = (2\pi\tau)^{-N/2} \exp\left(-\frac{1}{\tau} \|Y - G_{\cdot m} X_m\|_2^2\right).$$

• 
$$\pi(X_m|m) = \exp(-\lambda \|X_m\|_1 - |m|\log(c_\lambda))$$
, where  $\lambda \ge 0$ .

• 
$$\pi(m) = w_m$$
, where  $\sum_{m \in \mathcal{M}} w_m = 1$ .

Posterior distribution on  $\Theta = \bigcup_{m \in \mathcal{M}} (\{m\} \times \mathbb{R}^{|m|})$ :

$$\pi(m, X_m | Y) \propto w_m c_{\lambda}^{-|m|} \exp\left(-\frac{1}{2\tau} \|Y - G_{\cdot m} X_m\|_2^2 - \lambda \|X_m\|_1\right)$$

(日) (同) (日) (日)

.

3

The simplified model The Bayesian model selection framework

Goal : propose a transdimensional MCMC method to sample the posterior distribution.

- Robust in high dimensional settings
- Can deal with non-differentiability in the penalization function
- In harmony with sparsity assumption

Two main ingredients The algorithm

Goal of the **Proximal MALA** algorithm (PMALA): build a **Markov** chain converging to a target distribution with density of the form

$$\pi(x) \propto \exp(-g(x) - h(x))$$
,

where

- g: continuously differentiable, convex, such that  $\nabla g$  is  $L_g$ -Lipschitz,
- h: convex.

$$ightarrow$$
 Applied with  $g(x)=rac{1}{2 au}\|Y-{{{{{ G}x}}}}\|_{2}^{2}$  and  $h(x)=\lambda\|x\|_{1}.$ 

Two main ingredients The algorithm

Ingredient 1: The **proximal gradient algorithm** (also known as the Iterative Shrinkage Thresholding Algorithm)

- Goal: minimize g + h where
  - g: continuously differentiable, convex, such that  $\nabla g$  is  $L_g$ -Lipschitz,
  - h: convex

Two main ingredients The algorithm

An iteration of the proximal gradient algorithm starting from  $x^t$ : (1) Define a local approximation of g + h at  $x^t$  by

$$Q_L(x^t, x) = h(x) + g(x^t) + \langle x - x^t, \nabla g(x^t) \rangle + \frac{L}{2} ||x - x^t||_2^2$$
.

(2) Set  $x^{t+1} = \operatorname{argmin}_{x} Q_{L}(x^{t}, x) = \operatorname{prox}_{h/L} \left( x^{t} - \frac{1}{L} \nabla g(x^{t}) \right)$ , where

$$\operatorname{prox}_{\gamma h}(\mathbf{u}) = \operatorname{argmin}_{x} \left( \gamma h(x) + \frac{1}{2} \|x - u\|_{2}^{2} \right)$$

٠

э

(日) (同) (日) (日)

Two main ingredients The algorithm

# Ingredient 2: The Metropolis Adjusted Langevin Algorithm (MALA)

Goal: build a Markov chain converging to a target distribution with density  $\pi(x) \propto \exp(-g(x))$ , where g is differentiable.

Two main ingredients The algorithm

An iteration of MALA starting from  $X^t$ :

(1) **Propose** a new point

$$Y^{t+1} = X^t - \frac{\sigma^2}{2} \nabla g(X^t) + \sigma W^{t+1} ,$$

where  $W^{t+1}$  is a random vector with i.i.d. entries from  $\mathcal{N}(0,1)$ .

(2) Classical Acceptation/Rejection step.

(日) (同) (三) (三)

э

Two main ingredients The algorithm

An iteration of **PMALA** starting from  $X^t$ :

(1) **Propose** a new point

$$Y^{t+1} = \operatorname{prox}_{\sigma^2 \mathbf{h}/2} \left( X^t - \frac{\sigma^2}{2} \nabla g(X^t) + \sigma W^{t+1} \right) ,$$

where  $W^{t+1}$  is a random vector with i.i.d. entries from  $\mathcal{N}(0, 1)$ .

(2) Classical Acceptation/Rejection step, with acceptance probability  $\alpha(x, y) = \frac{\pi(y)q(y,x)}{\pi(x)q(x,y)}$ , where q(x, y) is the density of the proposal distribution (explicitly known).

4 E N 4 E N 4 E N 4 E

Two main ingredients The algorithm

#### Lemma

Let  $\boldsymbol{\mu} \in \mathbb{R}^{P}$  and  $\gamma, \sigma > 0$ . Set  $\boldsymbol{Y} = prox_{\gamma \parallel \cdot \parallel_{1}} (\boldsymbol{\mu} + \sigma \boldsymbol{W})$  where  $\boldsymbol{W} \in \mathbb{R}^{P}$  is a matrix of i.i.d random variables  $\sim \mathcal{N}(0, 1)$ . The distribution of  $\boldsymbol{Y} \in \mathbb{R}^{P}$  is given by

$$\sum_{m \in \mathcal{M}} \left( \prod_{i \notin I_m} p_1(\boldsymbol{\mu}_i) \, \delta_0(\mathrm{d}\boldsymbol{z}_i) \right) \left( \prod_{i \in I_m} f_1(\boldsymbol{\mu}_i, \boldsymbol{z}_i) \mathrm{d}\boldsymbol{z}_i \right) \, ,$$

where for any  $c, z \in \mathbb{R}$ ,

$$p_1(c) = \mathbb{P}\left\{ |c + \xi| \le \gamma \right\} , \text{ with } \xi \sim \mathcal{N}(0, \sigma^2) ,$$
  
$$f_1(c, z) = \left(2\pi\sigma^2\right)^{-1/2} \exp\left(-\frac{1}{2\sigma^2} \left| \left(1 + \frac{\gamma}{|z|}\right) z - c \right|_2^2\right)$$



An iteration of PMALA starting from x is equivalent to:

(i) sample  $m' = (m'_1, \dots, m'_P)$  with  $(m'_i, i \in \{1, \dots, P\})$  i.i.d. and such that  $m'_i$  is a Bernoulli r.v. with success parameter

$$1 - \mathbb{P}\left( \left| \left( x - rac{\sigma^2}{2} 
abla g(x) 
ight)_i + \xi \right|^2 \leq \gamma 
ight) \qquad \xi \sim \mathcal{N}(0, \sigma^2) \; .$$

(ii) sample  $y = (y_i)_{1 \le i \le P}$  in  $\mathbb{R}^{|m'|}$  with independant components such that for any  $i \in I_{m'}$ , the distribution of  $y_i$  is proportional to

$$\exp\left(-\frac{1}{2\sigma^2}\left|\left(1+\frac{\gamma}{|y_i|}\right)y_i-\left(x-\frac{\sigma^2}{2}\nabla g(x)\right)_i\right|^2\right)$$

The data:  $Y = GX + \sqrt{ au}E$ 

- The components of E are samples of  $\mathcal{N}(0,1)$
- $X = (X_i)_{1 \le i \le P}$  with  $X_i = \mathbf{1}_{i \le S}$  with S depending on the example.
- Columns of  $G \in \mathbb{R}^{N \times P}$ :
  - uncorrelated designs: independant Gaussian samples.
  - correlated designs: independant Gaussian samples or linear combinations of other columns plus Gaussian vectors.

Implementation parameters:

- Prior on the models: uniform.
- Starting point: empty model.

医肾管医肾管

Toy examples and simulated data



Figure: Posterior probability of activation in 16 dimensions, with uncorrelated design (i.e.  $\mathbb{P}(x_i \neq 0|M)$  for  $1 \le i \le 16$ ); S=8, P=16, N=100.

Toy examples and simulated data



Figure: Posterior probability of the models in 16 dimensions, with uncorrelated design (i.e.  $\mathbb{P}(m|M)$  for  $1 \le m \le 2^{16}$ ; logarithmic scale); S=8, P=16, N=100.

Toy examples and simulated data



Figure: Dimension 16, uncorrelated design. Right: evolution of the acceptance rate; Left: evolution of the probability of activation.

Toy examples and simulated data



Figure: Posterior probability of activation in 200 dimensions, with uncorrelated design; S=10, P=200, N=100.

B → ((B)

Toy examples and simulated data



Figure: Comparison of posterior probabilities of the models in 16 dimensions, with correlated design; S=8, P=16, N=100.

#### Toy examples and simulated data



Figure: Comparison of mean estimated regression vectors in 300 dimensions, with correlated design; S=16, P=300, N=100.

Toy examples and simulated data



Figure: Comparison of posterior probabilities of activation in 300 dimensions, with correlated design; S=16, P=300, N=100.

-

Futur directions

- partial updating (higher control on the acceptance rate)
- tempering (to deal with multimodality)
- hard thresholding or other thresholding functions (to avoid shrinkage for regression applications)
- theory : geometric ergodicity, ...
- real data (back to brain imaging, regression problems)

・ロト ・ 同ト ・ ヨト ・ 日

### Thank you !

(日) (四) (王) (王)

æ