

# Sparse Bayesian nonparametric regression

François Caron and Arnaud Doucet

Depts of Computer Science & Statistics, UBC

July 7, 2008

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$
- Variable selection, decomposition of a signal over an overcomplete basis, etc.

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$
- Variable selection, decomposition of a signal over an overcomplete basis, etc.
- Numerous models/algorithms

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$
- Variable selection, decomposition of a signal over an overcomplete basis, etc.
- Numerous models/algorithms
  - Spike and Slab

- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$
- Variable selection, decomposition of a signal over an overcomplete basis, etc.
- Numerous models/algorithms
  - Spike and Slab
  - Lasso



- Linear regression model

$$y = X\beta + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2 I_L) \quad (1)$$

- $y \in \mathbb{R}^L$ ,  $X$  design matrix of size  $L \times K$ ,  $\beta \in \mathbb{R}^K$
- Sparse estimate of  $\beta$
- Variable selection, decomposition of a signal over an overcomplete basis, etc.
- Numerous models/algorithms
  - Spike and Slab
  - Lasso
  - Relevance Vector Machine

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$

# Introduction

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Scale-mixture of Gaussians

$$p(\beta_k) = \int \mathcal{N}(\beta_k; 0, \sigma_k^2) p(\sigma_k^2) d\sigma_k^2$$

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Scale-mixture of Gaussians

$$p(\beta_k) = \int \mathcal{N}(\beta_k; 0, \sigma_k^2) p(\sigma_k^2) d\sigma_k^2$$

- Laplace prior  $\rightarrow$  Lasso objective function

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Scale-mixture of Gaussians

$$p(\beta_k) = \int \mathcal{N}(\beta_k; 0, \sigma_k^2) p(\sigma_k^2) d\sigma_k^2$$

- Laplace prior  $\rightarrow$  Lasso objective function
- Normal-Jeffreys

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Scale-mixture of Gaussians

$$p(\beta_k) = \int \mathcal{N}(\beta_k; 0, \sigma_k^2) p(\sigma_k^2) d\sigma_k^2$$

- Laplace prior  $\rightarrow$  Lasso objective function
- Normal-Jeffreys
- Normal-exponential gamma

- Prior distribution  $p(\beta) = \prod_{k=1}^K p(\beta_k)$
- Local minima of the objective function

$$-\frac{1}{2\sigma^2} \|y - X\beta\|^2 - \log p(\beta) \quad (2)$$

- Scale-mixture of Gaussians

$$p(\beta_k) = \int \mathcal{N}(\beta_k; 0, \sigma_k^2) p(\sigma_k^2) d\sigma_k^2$$

- Laplace prior  $\rightarrow$  Lasso objective function
- Normal-Jeffreys
- Normal-exponential gamma
- Find local minimum of Eq. (2) with EM algorithm



# Why this title is too bad

- **Sparse**... but... one model does not lead to 'strictly' sparse estimates!

# Why this title is too bad

- **Sparse**... but... one model does not lead to 'strictly' sparse estimates!
- **Bayesian**... but... EM algorithm

# Why this title is too bad

- **Sparse**... but... one model does not lead to 'strictly' sparse estimates!
- **Bayesian**... but... EM algorithm
- **nonparametric**... but... the number of parameters is finite!

# Why this title is too bad

- **Sparse**... but... one model does not lead to 'strictly' sparse estimates!
- **Bayesian**... but... EM algorithm
- **nonparametric**... but... the number of parameters is finite!
- **regression**... YES IT IS :-)

# Overview

- 1 Introduction
- 2 Models**
- 3 Sparsity properties
- 4 Empirical results
- 5 Conclusion

# Normal-gamma model

- Gamma prior over  $\sigma_k^2$

$$\sigma_k^2 \sim \mathcal{G}\left(\frac{\alpha}{K}, \frac{\gamma^2}{2}\right)$$

# Normal-gamma model

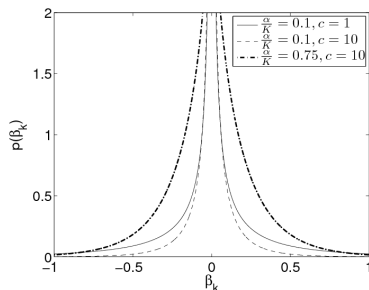
- Gamma prior over  $\sigma_k^2$

$$\sigma_k^2 \sim \mathcal{G}\left(\frac{\alpha}{K}, \frac{\gamma^2}{2}\right)$$

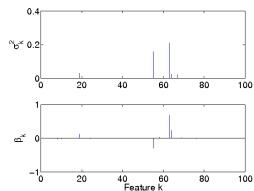
- Marginal distribution over  $\beta_k$

$$p(\beta_k) \propto |\beta_k|^{\frac{\alpha}{K} - \frac{1}{2}} \mathcal{K}_{\frac{\alpha}{K} - \frac{1}{2}}(\gamma|\beta_k|) \quad (3)$$

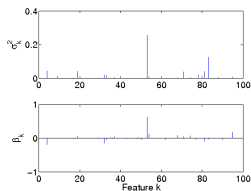
where  $\mathcal{K}_\nu(\cdot)$  is the modified Bessel function of the second kind.



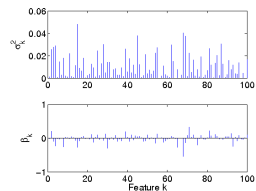
# Normal-gamma model



(a)  $\alpha = 1$



(b)  $\alpha = 5$



(c)  $\alpha = 100$



# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E}\left[\sum_{k=1}^K |\beta_k|\right] = \frac{2\alpha}{\gamma} \quad , \quad \mathbb{E}\left[\sum_{k=1}^K \beta_k^2\right] = \frac{2\alpha}{\gamma^2}.$$

# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E}\left[\sum_{k=1}^K |\beta_k|\right] = \frac{2\alpha}{\gamma}, \quad \mathbb{E}\left[\sum_{k=1}^K \beta_k^2\right] = \frac{2\alpha}{\gamma^2}.$$

- Stick breaking construction for the weights

# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E}\left[\sum_{k=1}^K |\beta_k|\right] = \frac{2\alpha}{\gamma}, \quad \mathbb{E}\left[\sum_{k=1}^K \beta_k^2\right] = \frac{2\alpha}{\gamma^2}.$$

- Stick breaking construction for the weights
  - Let  $\sigma_{(1)}^2 \geq \sigma_{(2)}^2 \geq \dots \geq \sigma_{(K)}^2$  be the order statistics of the sequence  $\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2$ .

# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E}\left[\sum_{k=1}^K |\beta_k|\right] = \frac{2\alpha}{\gamma}, \quad \mathbb{E}\left[\sum_{k=1}^K \beta_k^2\right] = \frac{2\alpha}{\gamma^2}.$$

- Stick breaking construction for the weights
  - Let  $\sigma_{(1)}^2 \geq \sigma_{(2)}^2 \geq \dots \geq \sigma_{(K)}^2$  be the order statistics of the sequence  $\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2$ .
  - $\bar{\sigma}^2 = \left( \frac{\sigma_{(1)}^2}{\sum_k \sigma_{(k)}^2}, \frac{\sigma_{(2)}^2}{\sum_k \sigma_{(k)}^2}, \dots \right)$  and  $\sum_k \sigma_{(k)}^2$  are independent and respectively distributed according to  $PD(\alpha)$  and  $\mathcal{G}(\alpha, \gamma^2/2)$

# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E}\left[\sum_{k=1}^K |\beta_k|\right] = \frac{2\alpha}{\gamma}, \quad \mathbb{E}\left[\sum_{k=1}^K \beta_k^2\right] = \frac{2\alpha}{\gamma^2}.$$

- Stick breaking construction for the weights

- Let  $\sigma_{(1)}^2 \geq \sigma_{(2)}^2 \geq \dots \geq \sigma_{(K)}^2$  be the order statistics of the sequence  $\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2$ .
- $\bar{\sigma}^2 = \left( \frac{\sigma_{(1)}^2}{\sum_k \sigma_{(k)}^2}, \frac{\sigma_{(2)}^2}{\sum_k \sigma_{(k)}^2}, \dots \right)$  and  $\sum_k \sigma_{(k)}^2$  are independent and respectively distributed according to  $PD(\alpha)$  and  $\mathcal{G}(\alpha, \gamma^2/2)$
- Can be recovered from the (Infinite) stick-breaking construction

$$\pi_k = \zeta_k \prod_{j=1}^{k-1} (1 - \zeta_j) \text{ with } \zeta_j \sim \mathcal{B}(1, \alpha) \quad (4)$$

# Asymptotic properties

- Bounded sum of the terms

$$\lim_{K \rightarrow \infty} \mathbb{E} \left[ \sum_{k=1}^K |\beta_k| \right] = \frac{2\alpha}{\gamma}, \quad \mathbb{E} \left[ \sum_{k=1}^K \beta_k^2 \right] = \frac{2\alpha}{\gamma^2}.$$

- Stick breaking construction for the weights

- Let  $\sigma_{(1)}^2 \geq \sigma_{(2)}^2 \geq \dots \geq \sigma_{(K)}^2$  be the order statistics of the sequence  $\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2$ .
- $\bar{\sigma}^2 = \left( \frac{\sigma_{(1)}^2}{\sum_k \sigma_{(k)}^2}, \frac{\sigma_{(2)}^2}{\sum_k \sigma_{(k)}^2}, \dots \right)$  and  $\sum_k \sigma_{(k)}^2$  are independent and respectively distributed according to  $PD(\alpha)$  and  $\mathcal{G}(\alpha, \gamma^2/2)$
- Can be recovered from the (Infinite) stick-breaking construction

$$\pi_k = \zeta_k \prod_{j=1}^{k-1} (1 - \zeta_j) \text{ with } \zeta_j \sim \mathcal{B}(1, \alpha) \quad (4)$$

- Coefficients ( $\beta_k$ ) are the weights (jumps) of the so-called variance gamma process (Brownian motion evaluated at times given by a gamma process)

# Normal-inverse Gaussian model

- Inverse-Gaussian prior over  $\sigma_k^2$

$$\sigma_k^2 \sim \mathcal{IG}\left(\frac{\alpha}{K}, \gamma\right) \quad (5)$$

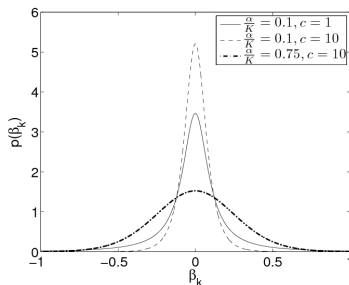
# Normal-inverse Gaussian model

- Inverse-Gaussian prior over  $\sigma_k^2$

$$\sigma_k^2 \sim \text{IG}\left(\frac{\alpha}{K}, \gamma\right) \quad (5)$$

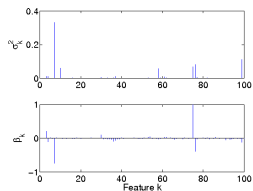
- Marginal pdf of  $\beta_k$

$$p(\beta_k) \propto \left(\frac{\alpha^2}{K^2} + \beta_k^2\right)^{-1/2} \mathcal{K}_1\left(\gamma\sqrt{\frac{\alpha^2}{K^2} + \beta_k^2}\right) \quad (6)$$

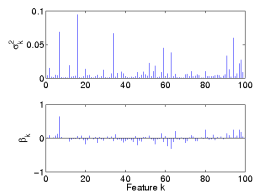




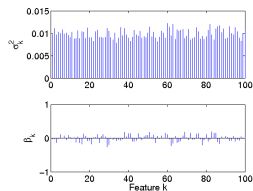
# Normal-inverse Gaussian model



(d)  $\alpha = 1$



(e)  $\alpha = 5$



(f)  $\alpha = 100$

# Extension

- $N$  vectors  $\{y_n\}_{n=1}^N$  where  $y_n \in \mathbb{R}^L$ .

# Extension

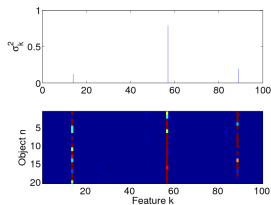
- $N$  vectors  $\{y_n\}_{n=1}^N$  where  $y_n \in \mathbb{R}^L$ .
- For a given  $k$  the random variables  $\{\beta_k^n\}_{n=1}^N$  are statistically dependent and exchangeable

# Extension

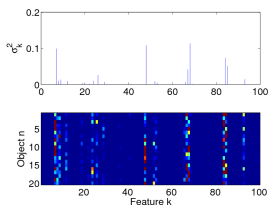
- $N$  vectors  $\{y_n\}_{n=1}^N$  where  $y_n \in \mathbb{R}^L$ .
- For a given  $k$  the random variables  $\{\beta_k^n\}_{n=1}^N$  are statistically dependent and exchangeable
- Hierarchical model

$$\sigma_k^2 \sim \mathcal{G}\left(\frac{\alpha}{K}, \frac{\gamma^2}{2}\right) \text{ or } \sigma_k^2 \sim \text{IG}\left(\frac{\alpha}{K}, \gamma\right)$$

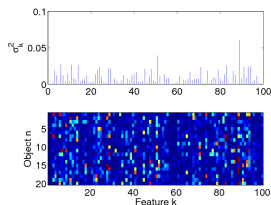
for  $k = 1, \dots, K$  and  $\beta_k^n \sim \mathcal{N}(0, \sigma_k^2)$  for  $n = 1, \dots, N$ .



(m)  $\alpha = 1$



(n)  $\alpha = 5$



(o)  $\alpha = 100$

- As  $K \rightarrow \infty$ , prior distributions over infinite matrices with real-valued entries

- As  $K \rightarrow \infty$ , prior distributions over infinite matrices with real-valued entries
- Complementary to the Indian buffet process and the infinite gamma-Poisson process which are prior distributions over infinite matrices with integer-valued entries.

# Overview

- 1 Introduction
- 2 Models
- 3 Sparsity properties**
- 4 Empirical results
- 5 Conclusion

# Sparsity properties

- Minimize

$$-\sum_{n=1}^N \frac{1}{2\sigma_n^2} \|y^n - X\beta^n\|^2 - \sum_{k=1}^K \text{pen}(\beta_k^{1:N})$$

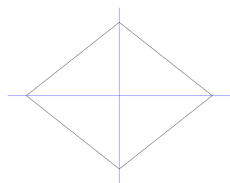
where

	$\text{pen}(\beta_k^{1:N})$
Lasso ( $N = 1$ )	$\gamma \beta_k $
NJ	$N \log(u_k)$
NG	$(\frac{N}{2} - \frac{\alpha}{K}) \log u_k - \log \mathcal{K}_{\frac{\alpha}{K} - \frac{N}{2}}(\gamma u_k)$
NIG	$\frac{N+1}{2} \log(q_k) - \log \mathcal{K}_{\frac{N+1}{2}}(\gamma q_k)$

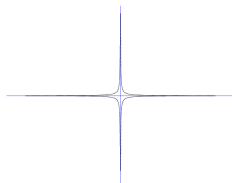
$$\text{where } u_k = \sqrt{\sum_{n=1}^N (\beta_k^n)^2}, \quad q_k = \sqrt{\frac{\alpha^2}{K^2} + u_k^2}$$



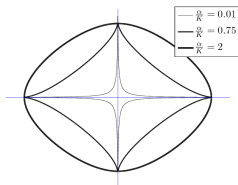
# Sparsity properties



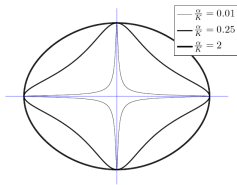
(p) Laplace



(q) Normal-Jeffreys



(r) Normal-gamma



(s) Normal-inverse  
Gaussian

Figure: Contour of constant value of  $pen(\beta_1) + pen(\beta_2)$  for different priors

# Sparsity properties

- The normal-gamma prior is a thresholding rule for  $\alpha/K \leq 1$  and yields sparse estimates
- The normal-inverse Gaussian is not a thresholding rule but it can yield “almost sparse” estimates

# Overview

- 1 Introduction
- 2 Models
- 3 Sparsity properties
- 4 Empirical results**
- 5 Conclusion

- 100 datasets with  $L = 50$  and  $\sigma = 1$
- Correlation between  $X_{k,i}$  and  $X_{k,j}$  is  $\rho^{|i-j|}$  with  $\rho = 0.5$
- True  $\beta = (3 \ 1.5 \ 0 \ 0 \ 2 \ 0 \ 0 \dots)^T \in \mathbb{R}^K$ , where  $K = 20, 60, 100, 200$
- Parameters of the Lasso, NG and NIG are estimated by 5-fold cross-validation

# Empirical results

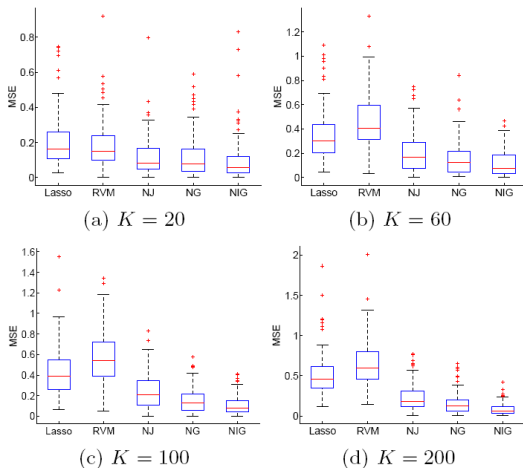


Figure: Box plots of the MSE associated to the simulated data.

# Conclusion

Why this title is not so bad

- Sparse...

# Conclusion

Why this title is not so bad

- Sparse...
  - Two classes of models which lead to sparser estimates

# Conclusion

Why this title is not so bad

- Sparse...
  - Two classes of models which lead to sparser estimates
- Bayesian nonparametric...



# Conclusion

Why this title is not so bad

- Sparse...
  - Two classes of models which lead to sparser estimates
- Bayesian nonparametric...
  - Related to a class of nonparametric Bayesian model when  $K \rightarrow \infty$

# Conclusion

Why this title is not so bad

- Sparse...
  - Two classes of models which lead to sparser estimates
- Bayesian nonparametric...
  - Related to a class of nonparametric Bayesian model when  $K \rightarrow \infty$
- regression...

# Conclusion

Why this title is not so bad

- Sparse...
  - Two classes of models which lead to sparser estimates
- Bayesian nonparametric...
  - Related to a class of nonparametric Bayesian model when  $K \rightarrow \infty$
- regression...
  - Extension to probit regression

- Ongoing work with K. Murphy on graph learning with group sparsity

# Conclusion

- Ongoing work with K. Murphy on graph learning with group sparsity
- How to use the stick-breaking construction?

- Ongoing work with K. Murphy on graph learning with group sparsity
- How to use the stick-breaking construction?
- Marginal distribution?



Barndorff-Nielsen, O. (1997).

Normal inverse Gaussian distributions and stochastic volatility modelling.

*Scandinavian Journal of Statistics*, 24, 1–13.



Figueiredo, M. (2003).

Adaptive sparseness for supervised learning.

*IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25, 1150–1159.



Griffin, J., & Brown, P. (2007).

*Bayesian adaptive lasso with non-convex penalization* (Technical Report).

Dept of Statistics, University of Warwick.