



Bayesian Sparse-Smooth Modeling and Variational Inference

Shunsuke Horii

November 13, 2017

Waseda University

Background

- Learning from **high-dimensional data** is one of the emerging tasks in machine learning
 - In some cases, dimension p is larger than sample size n
- **Sparsity** is an important concept
- Optimization based approach such as **LASSO** has been well studied
 - LASSO can be considered as a method of finding MAP estimate assuming the Laplace prior
- Full information of the posterior distribution is useful for some problems
 - Ex: Bayesian experimental design
- Some methods for sparse Bayesian modeling
 - Relevance Vector Machine, Gibbs sampler (Bayesian LASSO), **Variational Bayes**

Past works

- Some extensions of LASSO:
 - Group LASSO: Group Sparsity
 - Fused LASSO: Sparsity + Smoothness
 - Parameter is sparse and locally constant
- Some Bayesian studies for extended LASSO models:
 - Kyung et al, 2010: Gibbs sampler for Group LASSO, Fused LASSO, Elastic net
 - Babacan et al, 2014: Variational Bayes for Group LASSO

Our work

- Variational Bayes for sparse and smooth model (Fused LASSO model)
 - Variational Bayes is more efficient than Gibbs sampler

Relationship with past works:

	Sparsity only	Group Sparsity	Sparsity + Smoothness
Optimization based	Tibshirani, 1996	Yuan and Lin, 2006	Tibshirani, 2005
MCMC	Park and Casella, 2008	Kyung et al., 2010	Kyung et al., 2010
Variational Bayes	Babacan et al., 2014	Babacan et al., 2014	Our work

Strength of our work:

- Considering sparsity and smoothness: effective for sparse and smooth data
 - Ex: Denoising for sparse image
- Variational Bayes: more efficient than Gibbs sampler
 - Applicable for very high dimensional data analysis

Problem setup

Linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- $\mathbf{y} \in \mathbb{R}^n$: Observation vector
- $\mathbf{X} \in \mathbb{R}^{n \times p}$: Design matrix
- $\boldsymbol{\beta} \in \mathbb{R}^p$: **Unknown parameter**
- $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, s^{-1}\mathbf{I}_p)$: Gaussian noise vector (s : precision parameter)

Assumption:

Sparsity: $\boldsymbol{\beta}$ is sparse (number of nonzero elements is small)

Smoothness: β_i and β_j take similar values for $(i, j) \in E$

- $G = (V, E)$ is a predefined graph

How to statistically model sparsity and smoothness?

Straightforward approach

Assume Laplace priors for $\{\beta_j\}$ and $\{\beta_i - \beta_j\}$

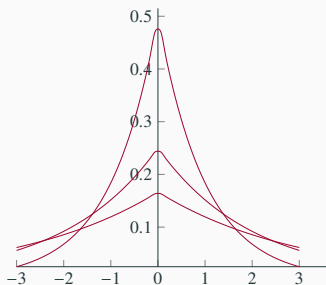
$$p(\boldsymbol{\beta}) \propto \prod_{j=1}^p \exp\left(-\frac{|\beta_j|}{\sqrt{a_{\tau,j}}}\right) \prod_{(j,k) \in E} \exp\left(-\frac{|\beta_j - \beta_k|}{\sqrt{a_{v,jk}}}\right)$$

Sparsity

Smoothness

- MAP estimator corresponds to the Fused LASSO estimator
- Difficult to calculate the posterior distribution
- What values should be set for α ?

⇒ needs another approach



Hierarchical representation of Laplace distribution

Laplace distribution can be expressed as a **scale mixture of Gaussian with exponential mixture**:

$$\frac{\sqrt{\alpha}}{2} \exp(-\sqrt{\alpha}|\beta|) = \int_0^{\infty} p(\beta|\tau)p(\tau|\gamma)d\tau$$

$$p(\beta|\tau) = \mathcal{N}(\beta|0, \tau^{-1})$$

$$p(\tau|\gamma) = \frac{\gamma}{2} \exp\left(-\frac{\gamma}{2}\tau\right)$$

- This property is used in the past Bayesian studies on the sparse modeling
- If we use inverse gamma distribution for the scale mixture, the marginal distribution becomes Student's t distribution

Sparse-Smooth hierarchical modeling

Gaussian distribution for β :

$$p(\beta|\tau, \nu) \propto \prod_{j=1}^p \exp(-\tau_j \beta_j^2) \prod_{(j,k) \in E} \exp(-\nu_{jk} (\beta_j - \beta_k)^2)$$

Exponential (or IG) distribution for τ and ν :

$$p(\tau|\mathbf{a}_\tau) = \prod_{j=1}^p \text{Exp}(\tau_j | a_{\tau,j}), \quad p(\nu|\mathbf{a}_\nu) = \prod_{(j,k) \in E} \text{Exp}(\nu_{jk} | a_{\nu,jk})$$

Gamma distribution for $\mathbf{a}_\tau, \mathbf{a}_\nu$:

$$p(\mathbf{a}_\tau) = \prod_{j=1}^p \text{Gam}(a_{\tau,j} | \theta_\tau, k_\tau), \quad p(\mathbf{a}_\nu) = \prod_{(j,k) \in E} \text{Gam}(a_{\nu,jk} | \theta_\nu, k_\nu)$$

Variational Bayes

Joint distribution:

$$p(\mathbf{y}, \beta, s, \boldsymbol{\tau}, \boldsymbol{\nu}, \mathbf{a}_{\boldsymbol{\tau}}, \mathbf{a}_{\boldsymbol{\nu}}) = p(\mathbf{y}|\beta, s)p(\beta|\boldsymbol{\tau}, \boldsymbol{\nu})p(\boldsymbol{\tau}|\mathbf{a}_{\boldsymbol{\tau}})p(\boldsymbol{\nu}|\mathbf{a}_{\boldsymbol{\nu}})p(\mathbf{a}_{\boldsymbol{\tau}}, \mathbf{a}_{\boldsymbol{\nu}})p(s)$$

We want to calculate **posterior distribution** $p(\beta|\mathbf{y})$

⇒ complex integral calculation is required

Mean field approximation:

For $\boldsymbol{\theta} = (\beta, \boldsymbol{\tau}, \boldsymbol{\nu}, \mathbf{a}_{\boldsymbol{\tau}}, \mathbf{a}_{\boldsymbol{\nu}}, s)$, find an approximate distribution $q(\boldsymbol{\theta})$ which minimizes Kullback-Leibler divergence:

$$q^*(\boldsymbol{\theta}) = \operatorname{argmin}_{q(\boldsymbol{\theta})} \int q(\boldsymbol{\theta}) \ln \frac{q(\boldsymbol{\theta})}{p(\boldsymbol{\theta}|\mathbf{y})} d\boldsymbol{\theta}$$

w.r.t. factorized distribution $q(\boldsymbol{\theta}) = q(\beta)q(\boldsymbol{\tau}, \boldsymbol{\nu})q(\mathbf{a}_{\boldsymbol{\tau}}, \mathbf{a}_{\boldsymbol{\nu}})q(s)$

Variational Bayes

$\langle \cdot \rangle$: Expectation w.r.t. the corresponding distribution

Update equation for $q(\beta)$:

$$\begin{aligned}q^*(\beta) &= \mathcal{N}(\beta | \bar{\beta}, \Sigma_\beta), \\ \bar{\beta} &= \langle s \rangle \Sigma_\beta \mathbf{X}^T \mathbf{y}, \\ \Sigma_\beta &= \left(\langle s \rangle \mathbf{X}^T \mathbf{X} + \langle \mathbf{S}_{\tau, \nu} \rangle \right)^{-1}\end{aligned}$$

Update equation for $q(\tau, \nu)$:

$$q^*(\tau, \nu) = \prod_{j=1}^p \text{Gig} \left(\tau_j \mid \langle a_{\tau,j} \rangle, \langle \beta_j^2 \rangle, \frac{1}{2} \right) \prod_{(j,k) \in E} \text{Gig} \left(\nu_{jk} \mid \langle a_{\nu,jk} \rangle, \langle (\beta_j - \beta_k)^2 \rangle, \frac{1}{2} \right)$$

$\text{Gig}(x|a, b, \rho)$ is the generalized inverse Gaussian:

$$\text{Gig}(x|a, b, \rho) \propto x^{\rho-1} \exp \left(-\frac{1}{2}(ax + bx^{-1}) \right)$$

Update equation for $q(\mathbf{a}_\tau, \mathbf{a}_\nu)$:

$$q^*(\mathbf{a}_\tau, \mathbf{a}_\nu) = \prod_{j=1}^p \text{Gam}(a_{\tau,j} | k_\tau + 1, \theta_\tau + \langle \tau_j \rangle / 2) \cdot \prod_{(j,k) \in E} \text{Gam}(a_{\nu,jk} | k_\nu + 1, \theta_\nu + \langle \nu_{jk} \rangle / 2)$$

Update equation for $q(s)$:

$$q^*(s) = \text{Gam}(s | k_s + n/2, \theta_s + \langle (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^2 \rangle / 2)$$

Assuming that $p(s)$ is $\text{Gam}(s | k_s, \theta_s)$

Experiment on synthetic data

- We considered two cases:

Case 1: True parameter is $\beta_1^* = (\mathbf{0}_{10}, \mathbf{2}_{10}, \mathbf{0}_{10}, \mathbf{2}_{10})$
(sparse and smooth)

Case 2: True parameter is $\beta_2^* = (0, 2, 0, 2, \dots) \in \mathbb{R}^{40}$
(sparse, but not smooth)

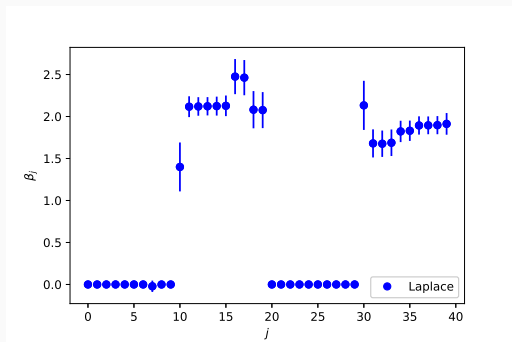
- Randomly generate sample (size = 40)
- Number of experiments is 100
- For $p(\tau, \nu)$, besides exponential distribution, we also consider inverse gamma distribution
 - Referred as **Laplace** when exponential distribution is assumed and **Student** when inverse gamma distribution is assumed

Experiment on synthetic data

Mean squared error:

	Lasso	Laplace	Student
β_1^*	5.75×10^{-1}	8.07×10^{-3}	7.98×10^{-3}
β_2^*	5.53×10^{-1}	1.37	1.11

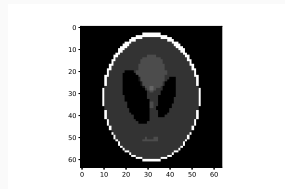
Example of estimation result:



Application for denoising of sparse image

$$\beta \in \mathbb{R}^{60 \times 60}:$$

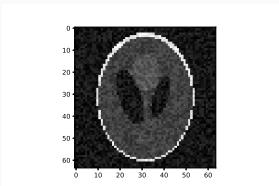
Original sparse image



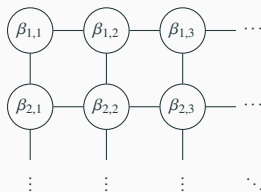
$$\mathbf{y} \in \mathbb{R}^{60 \times 60}:$$

Noisy image

$$\mathbf{y} = \beta + \epsilon$$



Graph G : 2D-grid



Problem: Estimate the original image from the noisy image

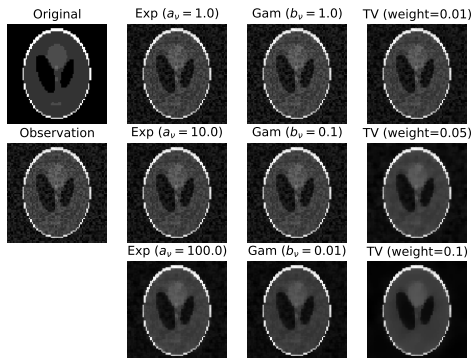
- In this experiment, τ , that controls sparsity, is fixed

Application for denoising of sparse image

PSNR [DB] and SSIM of restored images:

method	Exponential			Gamma			TV		
parameter	$a_v = 1.0$	$a_v = 10.0$	$a_v = 100.0$	$b_v = 1.0$	$b_v = 0.1$	$b_v = 0.01$	$\lambda=0.01$	$\lambda=0.05$	$\lambda=0.1$
PSNR	24.89	24.96	25.90	24.18	24.46	25.04	24.51	24.85	24.49
SSIM	0.5541	0.5680	0.6257	0.5407	0.5620	0.6138	0.5670	0.6119	0.6028

Restored images:



- We proposed **hierarchical model** for the problem of estimating parameter which is **sparse and smooth**.
- We also proposed an approximate estimation algorithm based on the **variational Bayes method**.
- The effectiveness of the proposed method was proved by experiments on synthetic data and real image data