

Modern Bayesian Statistics
Part III: high-dimensional modeling
Example 3: Sparse and time-varying covariance
modeling

HEDIBERT FREITAS LOPES¹
hedibert.org

13^a aMostra de Estatística
IME-USP, October 2018

¹Professor of Statistics and Econometrics at Insper, São Paulo.

Example 3: Sparse and time-varying covariance modeling

Consider the Gaussian linear model

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

where \mathbf{y} is n -length vector and \mathbf{X} is a $n \times q$ design matrix.

Ridge Regression (ℓ_2 penalty on $\boldsymbol{\beta}$)

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = \arg \min_{\boldsymbol{\beta}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right\}, \quad \lambda \geq 0,$$

leading to $\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}'\mathbf{y}$.

LASSO (ℓ_1 penalty on $\boldsymbol{\beta}$)

$$\hat{\boldsymbol{\beta}}_{\text{lasso}} = \arg \min_{\boldsymbol{\beta}} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}, \quad \lambda \geq 0,$$

which can be solved by using quadratic programming techniques such as *coordinate gradient descent*.

Bayesian regularization in linear regression problems

Hierarchical scale mixture of normals:

$$\beta|\psi \sim \mathcal{N}(0, \psi) \quad \text{and} \quad \psi|\theta \sim p(\psi),$$

Maximum a posteriori (MAP): $\arg \max_{\beta} \{p(\mathbf{y}|\beta, \sigma^2)p(\beta|\psi)\}$

A few cases:

Prior	$p(\psi)$	$p(\beta)$
Bayesian Lasso	$\psi \sim \mathcal{E}(\lambda^2/2)$	Laplace
Ridge	$\psi \sim \mathcal{IG}(a, b)$	Scaled-t
Normal-Gamma	$\psi \sim \mathcal{G}(\lambda, 1/(2\gamma^2))$	below

$$p(\beta|\lambda, \gamma^2) = \frac{1}{\sqrt{\pi}2^{\lambda-1/2}\gamma^{\lambda+1/2}\Gamma(\lambda)} |\beta|^{\lambda-1/2} K_{\lambda-1/2}(|\beta|/\gamma),$$

where $\text{Var}(\beta|\lambda, \gamma^2) = 2\lambda\gamma^2$ and excess kurtosis $3/\lambda$.

The Normal-Gamma prior

High mass close to zero and heavy tails

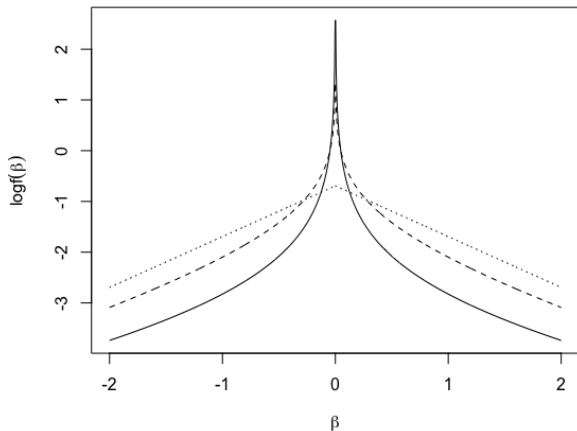


Figure: $\lambda = 0.1$ (dot), $\lambda = 0.33$ (dot-dashed), $\lambda = 1$ (solid).

Spike-and-slab priors

Stochastic search variable selection (SSVS)

SSVS places independent mixture priors directly on the coefficients

$$\beta|J \sim (1 - J) \underbrace{\mathcal{N}(0, \tau^2)}_{\text{spike}} + J \underbrace{\mathcal{N}(0, c^2 \tau^2)}_{\text{slab}},$$

with $c > 1$ large, $\tau > 0$ small and $J \sim \text{Ber}(\omega)$.

SMN representation

$$\beta|\psi \sim \mathcal{N}(0, \psi)$$

$$\psi|J \sim (1 - J)\delta_{\tau^2}(\cdot) + J\delta_{c^2\tau^2}(\cdot)$$

Spike-and-slab priors on the scale parameter ψ

Normal mixture of Inverse-Gammas (NMIG)

$$\begin{aligned}\beta|K &\sim \mathcal{N}(\mathbf{0}, K\tau^2), \\ K|\omega &\sim (1 - \omega)\delta_{v_0}(\cdot) + \omega\delta_{v_1}(\cdot), \quad v_0/v_1 \ll 1, \\ \tau^2 &\sim \mathcal{IG}(\mathbf{a}_\tau, \mathbf{b}_\tau).\end{aligned}\tag{1}$$

SMN representation

$$\beta|\psi \sim \mathcal{N}(\mathbf{0}, \psi)$$

and

$$\psi \sim (1 - \omega)\mathcal{IG}(\mathbf{a}_\tau, v_0\mathbf{b}_\tau) + \omega\mathcal{IG}(\mathbf{a}_\tau, v_1\mathbf{b}_\tau)$$

The resulting marginal distribution of β is a **two component mixture of scaled Student's t distributions**.

Spike-and-slab priors on the scale parameter ψ

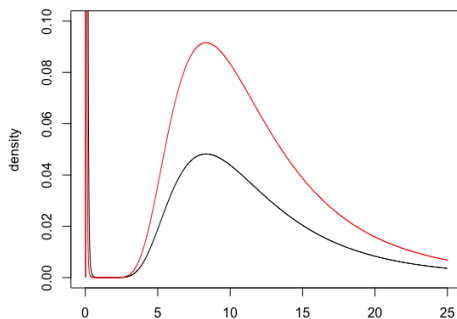


Figure: Conditional density for hypervariance ψ for NMIG mixture prior where $v_0 = 0.005$, $v_1 = 1$, $a_\tau = 5$, $b_\tau = 50$ and (a) $\omega = 0.5$ (black line), (b) $\omega = 0.95$ (red line). Note that as ω has a Uniform prior, (a) also corresponds to the marginal density of ψ . Observe that only the height of the density changes as ω varied.

Spike-and-slab priors: summary

Prior	Spike $\psi J=0$	Slab $\psi J=1$	Marginal $\beta \omega$	Constant c
SSVS	$\psi J=0 = \delta_{rQ}(\cdot)$	$\psi J=1 = \delta_Q(\cdot)$	$\omega\mathcal{N}(0, Q) + (1-\omega)\mathcal{N}(0, rQ)$	1
NMIG	$\mathcal{IG}(\nu, rQ)$	$\mathcal{IG}(\nu, Q)$	$\omega t_{2\nu}(0, Q/\nu) + (1-\omega)t_{2\nu}(0, rQ/\nu)$	$1/(\nu-1)$
Mixture of Laplaces	$\mathcal{E}(1/2rQ)$	$\mathcal{E}(1/2Q)$	$\omega\text{Lap}(\sqrt{Q}) + (1-\omega)\text{Lap}(\sqrt{rQ})$	2
Mixture of Normal-Gammas	$\mathcal{G}(a, 1/2rQ)$	$\mathcal{G}(a, 1/2Q)$	$\omega\mathcal{NG}(\beta_j a, Q) + (1-\omega)\mathcal{NG}(\beta_j a, r, Q)$	$2a$
Laplace-t	$\mathcal{E}(1/2rQ)$	$\mathcal{IG}(\nu, Q)$	$\omega t_{2\nu}(0, Q/\nu) + (1-\omega)\text{Lap}(\sqrt{rQ})$	$c_1 = 2, c_2 = 1/(\nu-1)$

Gaussian dynamic regression problems

- ▶ Consider the univariate Gaussian *dynamic linear model* (DLM) expressed by

$$y_t = \mathbf{F}_t' \boldsymbol{\beta}_t + \nu_t, \quad \nu_t \sim \mathcal{N}(0, V_t) \quad (2)$$

$$\boldsymbol{\beta}_t = \mathbf{G}_t \boldsymbol{\beta}_{t-1} + \boldsymbol{\omega}_t, \quad \boldsymbol{\omega}_t \sim \mathcal{N}(0, \mathbf{W}_t), \quad (3)$$

where $\boldsymbol{\beta}$ is of length q and $\boldsymbol{\beta}_0 \sim \mathcal{N}(\mathbf{m}_0, \mathbf{C}_0)$.

- ▶ **Dynamic regression model:** $\mathbf{F}_t = \mathbf{X}_t$ and $\mathbf{G}_t = \mathbf{I}_q$.
- ▶ **Static regression model:** $\mathbf{W}_t = 0$ for all t .

Shrinkage in dynamic regression problems

- ▶ Two main obstacles:
 1. Time-varying parameters (states), and
 2. A large number of predictors q .

- ▶ Two sources of sparsity:
 1. **horizontal sparsity**: $\beta_{j,t} = 0, \forall t$ for some coefficients j .
 2. **vertical sparsity**: $\beta_{j,t} = 0$ for several js at time t .

Shrinkage in dynamic regression problems

- ▶ Two main obstacles:
 1. Time-varying parameters (states), and
 2. A large number of predictors q .
- ▶ Two sources of sparsity:
 1. **horizontal sparsity**: $\beta_{j,t} = 0, \forall t$ for some coefficients j .
 2. **vertical sparsity**: $\beta_{j,t} = 0$ for several js at time t .

Illustration: $q = 5$ and $T = 12$

	<i>jan</i>	<i>feb</i>	<i>mar</i>	<i>apr</i>	<i>may</i>	<i>jun</i>	<i>jul</i>	<i>aug</i>	<i>sep</i>	<i>oct</i>	<i>nov</i>
x_0	$\beta_{1,1}$	$\beta_{1,2}$	$\beta_{1,3}$	$\beta_{1,4}$	$\beta_{1,5}$	$\beta_{1,6}$	$\beta_{1,7}$	$\beta_{1,8}$	$\beta_{1,9}$	$\beta_{1,10}$	$\beta_{1,11}$
x_1	0	0	0	0	0	0	0	0	0	0	0
x_2	$\beta_{3,1}$	$\beta_{3,2}$	$\beta_{3,3}$	$\beta_{3,4}$	$\beta_{3,5}$	0	0	0	$\beta_{3,9}$	$\beta_{3,10}$	$\beta_{3,11}$
x_3	0	0	$\beta_{4,3}$	$\beta_{4,4}$	$\beta_{4,5}$	$\beta_{4,6}$	$\beta_{4,7}$	$\beta_{4,8}$	$\beta_{4,9}$	$\beta_{4,10}$	$\beta_{4,11}$
x_4	$\beta_{5,1}$	$\beta_{5,2}$	$\beta_{5,3}$	$\beta_{5,4}$	$\beta_{5,5}$	0	0	0	$\beta_{5,9}$	$\beta_{5,10}$	$\beta_{5,11}$

Our contribution: a dynamic spike-and-slab model

Our contribution is defining a a spike-and-slab prior that not only shrinks time-varying coefficients in dynamic regression problems but allows for **dynamic variable selection**.

We use a **non-centered parametrization**:

$$\begin{aligned}y_t &= \mathbf{F}'_t \boldsymbol{\beta}_t + \nu_t, & \nu_t &\sim \mathcal{N}(0, \sigma_t^2) \\ \tilde{\boldsymbol{\beta}}_t &= \mathbf{G}_t \tilde{\boldsymbol{\beta}}_{t-1} + \boldsymbol{\omega}_t, & \boldsymbol{\omega}_t &\sim \mathcal{N}(0, \mathbf{W}_t),\end{aligned}$$

where

$$\begin{aligned}\tilde{\boldsymbol{\beta}}_t &= \left(\frac{\beta_{1,t}}{\sqrt{\psi_{1,t}}}, \dots, \frac{\beta_{q,t}}{\sqrt{\psi_{q,t}}} \right)' \\ \mathbf{G}_t &= \text{diag}(\varphi_1, \dots, \varphi_q) \\ \mathbf{W}_t &= \text{diag}((1 - \varphi_1^2), \dots, (1 - \varphi_q^2)), \\ \mathbf{F}'_t &= (X_{1,t} \sqrt{\psi_{1,t}}, \dots, X_{q,t} \sqrt{\psi_{q,t}}).\end{aligned}$$

Our contribution: a dynamic spike-and-slab model

For shrinking the states β_1, \dots, β_T , for any j coefficient, we place independent priors for each $\psi_t = \tau^2 K_t$ as

$$\begin{aligned}\tau^2 &\stackrel{\text{iid}}{\sim} p(\tau^2 | \boldsymbol{\theta}), \\ (K_t | K_{t-1} = v_i) &\stackrel{\text{ind}}{\sim} \omega_{1,i} \delta_1(\cdot) + (1 - \omega_{1,i}) \delta_r(\cdot), \\ \omega_{1,i} &= p(K_t = 1 | K_{t-1} = v_i),\end{aligned}$$

where $v_i \in \{r, 1\}$, $p(K_1 = r) = p(K_1 = 1) = 1/2$,

$r = \text{Var}_{\text{spike}}(\beta | \boldsymbol{\theta}) / \text{Var}_{\text{slab}}(\beta | \boldsymbol{\theta}) \ll 1$ and $p(\tau^2 | \boldsymbol{\theta})$ is one of priors from the previous Table.

Markov switching process for K_t

That is, K_t is a binary random latent variable that can assume binary values (regimes) $v_0 = r$ or $v_1 = 1$, depending only on the previous value of K_{t-1} and the prior **transition probabilities**

$\{\omega_{0,0}; \omega_{0,1}; \omega_{1,1}; \omega_{1,0}\}$.

Our contribution: direct acyclic graph

