

Advanced Bayesian Learning

Regularization and Variable Selection - Lecture 2

Mattias Villani

Department of Statistics



Lecture overview

- Spike-and-slab variable selection regression
- Polya-Gamma augmentation for logistic regression
- Extensions

Bayesian variable selection

- Linear regression

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon.$$

- Which variables have non-zero coefficients?

- Introduce **variable selection indicators** $\mathbf{z} = (z_1, \dots, z_p)$.

- Example: $\mathbf{z} = (1, 1, 0)$ means that $\beta_1 \neq 0$ and $\beta_2 \neq 0$, but $\beta_3 = 0$, so x_3 drops out of the model.

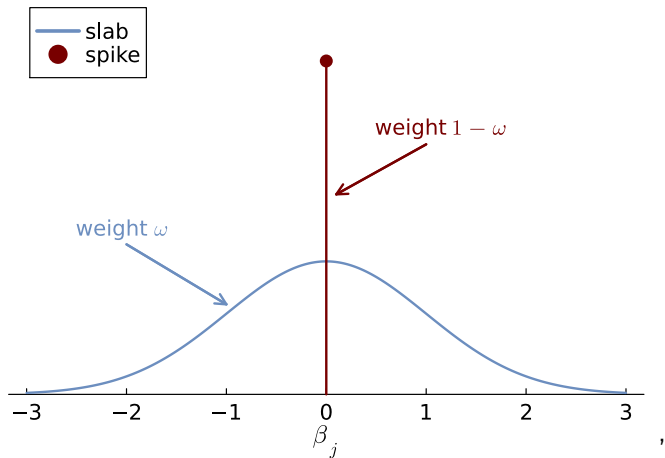
- **Spike-and-slab** prior

$$z_1, \dots, z_p \sim \text{Bernoulli}(\omega)$$

$$\beta_j | z_j \sim \begin{cases} N(0, \tau^2 \sigma^2) & \text{if } z_j = 1 \\ = 0 & \text{if } z_j = 0 \end{cases}$$

- **Prior inclusion probability** ω .

Spike-and-slab prior



Bayesian variable selection

■ Posterior

$$p(\beta, \sigma^2, \mathbf{z} | \mathbf{y}, \mathbf{X}) = p(\beta, \sigma^2 | \mathbf{z}, \mathbf{y}, \mathbf{X}) p(\mathbf{z} | \mathbf{y}, \mathbf{X})$$

$$p(\mathbf{z} | \mathbf{y}, \mathbf{X}) \propto p(\mathbf{y} | \mathbf{X}, \mathbf{z}) \cdot p(\mathbf{z})$$

- The prior $p(\mathbf{z})$ is $z_1, \dots, z_p \sim \text{Bernoulli}(\omega)$ as before.
- Need the **marginal likelihood** $p(\mathbf{y} | \mathbf{X}, \mathbf{z})$ for each model \mathbf{z} .

$$p(\mathbf{y} | \mathbf{X}, \mathbf{z}) = \int p(\mathbf{y} | \beta, \mathbf{X}, \mathbf{z}) p(\beta | \mathbf{X}, \mathbf{z}) d\beta$$

- For **linear Gaussian regression** the marginal likelihood is

$$t_{\nu_{0,z}} \left(\mathbf{y} | 0, \sigma_{0,z}^2 (I_n + \mathbf{X}_z \Omega_{0,z}^{-1} \mathbf{X}_z^\top) \right)$$

where $t_{\nu_{0,z}}$ is the multivariate- t density and \mathbf{X}_z is the matrix of covariates selected by \mathbf{z} .

- Prior hyperparameters ν_0 , σ_0^2 and Ω_0 allowed to depend on \mathbf{z} .

Bayesian variable selection via Gibbs sampling

- But there are 2^P model combinations to go through! *Ouch!*
- ... but most have essentially zero posterior probability. *Phew!*
- **Simulate** from the joint posterior distribution:

$$p(\beta, \sigma^2, \mathbf{z} | \mathbf{y}, \mathbf{X}) = p(\beta, \sigma^2 | \mathbf{z}, \mathbf{y}, \mathbf{X}) p(\mathbf{z} | \mathbf{y}, \mathbf{X})$$

- Simulate from $p(\mathbf{z} | \mathbf{y}, \mathbf{X})$ using **Gibbs sampling**:
 - ▶ Draw $z_1 | z_{-1}, \mathbf{y}, \mathbf{X}$
 - ▶ Draw $z_2 | z_{-2}, \mathbf{y}, \mathbf{X}$
 - ▶ ...
 - ▶ Draw $z_p | z_{-p}, \mathbf{y}, \mathbf{X}$
 - ▶ Draw β, σ^2 from $p(\beta, \sigma^2 | \mathbf{z}, \mathbf{y}, \mathbf{X})$.
- Compute $p(\mathbf{z} | \mathbf{y}, \mathbf{X}) \propto p(\mathbf{y} | \mathbf{X}, \mathbf{z}) \cdot p(\mathbf{z})$ for $z_j = 0$ and for $z_j = 1$, and normalize.
- **Model averaging** in a single simulation run.

Bayesian variable selection algorithm

Gibbs sampling for Bayesian variable selection in regression

Input: $n \times p$ matrix with p covariates as columns \mathbf{X}
vector \mathbf{y} with response observations
slab variance τ^2
prior inclusion probability ω
initial variable indicators $\mathbf{z}^{(0)} = (z_1^{(0)}, \dots, z_p^{(0)})$
number of posterior draws m .

```
for  $j$  in  $1:m$  do
    // Update regression parameters
    Draw  $(\sigma^2)^{(j)} | \mathbf{y}, \mathbf{X}_{\mathbf{z}^{(j-1)}} \sim \text{ScaledInv-}\chi^2(\nu_n, \sigma_n^2)$ 
    Draw  $\boldsymbol{\beta}_{\mathbf{z}^{(j-1)}}^{(j)} | (\sigma^2)^{(j)}, \mathbf{y}, \mathbf{X}_{\mathbf{z}^{(j-1)}} \sim N(\boldsymbol{\mu}_n, (\sigma^2)^{(j)} \Omega_n^{-1})$ 
    Set  $\boldsymbol{\beta}^{(j)}[\mathbf{z}^{(j-1)}] = \boldsymbol{\beta}_{\mathbf{z}^{(j-1)}}^{(j)}$  and  $\boldsymbol{\beta}^{(j)}[\text{Not}(\mathbf{z}^{(j-1)})] = 0$ 

    // Update mixture allocations
    Set  $\tilde{\mathbf{z}} = \mathbf{z}^{(j-1)}$ 
    for  $k$  in  $1:p$  do
        Set  $\tilde{\mathbf{z}}_0$  to  $\tilde{\mathbf{z}}$  but with  $k$ th element equal to 0
        Set  $\tilde{\mathbf{z}}_1$  to  $\tilde{\mathbf{z}}$  but with  $k$ th element equal to 1
        Compute  $\tilde{\omega}_{k,0} \propto (1 - \omega) \cdot p(\mathbf{y} | \mathbf{X}_{\tilde{\mathbf{z}}_0})$ 
        Compute  $\tilde{\omega}_{k,1} \propto \omega \cdot p(\mathbf{y} | \mathbf{X}_{\tilde{\mathbf{z}}_1})$ 
        Normalize  $\tilde{\omega}_{k,0}$  and  $\tilde{\omega}_{k,1}$  to sum to one
        Simulate allocation  $z_k^{(j)} \sim \text{Bernoulli}(\tilde{\omega}_{k,1})$ 
        Update  $\tilde{\mathbf{z}}$  with the new allocation  $z_k^{(j)}$ 
    end
    Set  $\mathbf{z}^{(j)} = \tilde{\mathbf{z}}$ 
end
```

Bayesian variable selection algorithm

```
for j in 1:m do
    // Update regression parameters
    Draw  $(\sigma^2)^{(j)} | \mathbf{y}, \mathbf{X}_{\mathbf{z}^{(j-1)}} \sim \text{ScaledInv-}\chi^2(\nu_n, \sigma_n^2)$ 
    Draw  $\boldsymbol{\beta}_{\mathbf{z}^{(j-1)}}^{(j)} | (\sigma^2)^{(j)}, \mathbf{y}, \mathbf{X}_{\mathbf{z}^{(j-1)}} \sim N(\boldsymbol{\mu}_n, (\sigma^2)^{(j)} \Omega_n^{-1})$ 
    Set  $\boldsymbol{\beta}^{(j)}[\mathbf{z}^{(j-1)}] = \boldsymbol{\beta}_{\mathbf{z}^{(j-1)}}^{(j)}$  and  $\boldsymbol{\beta}^{(j)}[\text{Not}(\mathbf{z}^{(j-1)})] = 0$ 

    // Update mixture allocations
    Set  $\tilde{\mathbf{z}} = \mathbf{z}^{(j-1)}$ 
    for k in 1:p do
        Set  $\tilde{\mathbf{z}}_0$  to  $\tilde{\mathbf{z}}$  but with kth element equal to 0
        Set  $\tilde{\mathbf{z}}_1$  to  $\tilde{\mathbf{z}}$  but with kth element equal to 1
        Compute  $\tilde{\omega}_{k,0} \propto (1 - \omega) \cdot p(\mathbf{y} | \mathbf{X}_{\tilde{\mathbf{z}}_0})$ 
        Compute  $\tilde{\omega}_{k,1} \propto \omega \cdot p(\mathbf{y} | \mathbf{X}_{\tilde{\mathbf{z}}_1})$ 
        Normalize  $\tilde{\omega}_{k,0}$  and  $\tilde{\omega}_{k,1}$  to sum to one
        Simulate allocation  $z_k^{(j)} \sim \text{Bernoulli}(\tilde{\omega}_{k,1})$ 
        Update  $\tilde{\mathbf{z}}$  with the new allocation  $z_k^{(j)}$ 
    end
    Set  $\mathbf{z}^{(j)} = \tilde{\mathbf{z}}$ 
end
```


Simple general Bayesian variable selection

- The previous algorithm only works when we can compute

$$p(\mathbf{z}|\mathbf{y}, \mathbf{X}) = \int p(\boldsymbol{\beta}, \sigma^2, \mathbf{z}|\mathbf{y}, \mathbf{X}) d\boldsymbol{\beta} d\sigma$$

- **MH** - propose $\boldsymbol{\beta}$ and \mathbf{z} jointly from the proposal distribution

$$q(\boldsymbol{\beta}_p|\boldsymbol{\beta}_c, \mathbf{z}_p)q(\mathbf{z}_p|\mathbf{z}_c)$$

- Main difficulty: how to propose the non-zero elements in $\boldsymbol{\beta}_p$?
- Simple approach:
 - ▶ Approximate posterior with **all** variables in the model:

$$\boldsymbol{\beta}|\mathbf{y}, \mathbf{X} \stackrel{approx}{\sim} N\left[\hat{\boldsymbol{\beta}}, J_y^{-1}(\hat{\boldsymbol{\beta}})\right]$$

- ▶ Propose $\boldsymbol{\beta}_p$ from $N\left[\hat{\boldsymbol{\beta}}, J_y^{-1}(\hat{\boldsymbol{\beta}})\right]$, conditional on the zero restrictions implied by \mathbf{z}_p . Formulas are available.

Variable selection in more complex models

Table 1
Posterior summary of the one-component split-t model.^a

| Parameters | Mean | Stdev | Post.Incl. |
|--|---------------|--------------|--------------|
| <i>Location μ</i> | | | |
| Const | 0.084 | 0.019 | – |
| <i>Scale ϕ</i> | | | |
| Const | 0.402 | 0.035 | – |
| LastDay | –0.190 | 0.120 | 0.036 |
| LastWeek | –0.738 | 0.193 | 0.985 |
| LastMonth | –0.444 | 0.086 | 0.999 |
| CloseAbs95 | 0.194 | 0.233 | 0.035 |
| CloseSqr95 | 0.107 | 0.226 | 0.023 |
| MaxMin95 | 1.124 | 0.086 | 1.000 |
| CloseAbs80 | 0.097 | 0.153 | 0.013 |
| CloseSqr80 | 0.143 | 0.143 | 0.021 |
| MaxMin80 | –0.022 | 0.200 | 0.017 |
| <i>Degrees of freedom ν</i> | | | |
| Const | 2.482 | 0.238 | – |
| LastDay | 0.504 | 0.997 | 0.112 |
| LastWeek | –2.158 | 0.926 | 0.638 |
| LastMonth | 0.307 | 0.833 | 0.089 |
| CloseAbs95 | 0.718 | 1.437 | 0.229 |
| CloseSqr95 | 1.350 | 1.280 | 0.279 |
| MaxMin95 | 1.130 | 1.488 | 0.222 |
| CloseAbs80 | 0.035 | 1.205 | 0.101 |
| CloseSqr80 | 0.363 | 1.211 | 0.112 |
| MaxMin80 | –1.672 | 1.172 | 0.254 |
| <i>Skewness λ</i> | | | |
| Const | –0.104 | 0.033 | – |
| LastDay | –0.159 | 0.140 | 0.027 |
| LastWeek | –0.341 | 0.170 | 0.135 |
| LastMonth | –0.076 | 0.112 | 0.016 |
| CloseAbs95 | –0.021 | 0.096 | 0.008 |
| CloseSqr95 | –0.003 | 0.108 | 0.006 |
| MaxMin95 | 0.016 | 0.075 | 0.008 |
| CloseAbs80 | 0.060 | 0.115 | 0.009 |
| CloseSqr80 | 0.059 | 0.111 | 0.010 |
| MaxMin80 | 0.093 | 0.096 | 0.013 |

Model averaging

- Let γ be a quantity with the same interpretation in the two models.
- Example: Prediction $\gamma = (y_{T+1}, \dots, y_{T+h})'$.
- The marginal posterior distribution of γ reads

$$p(\gamma|y) = p(M_1|y)p_1(\gamma|y) + p(M_2|y)p_2(\gamma|y),$$

$p_k(\gamma|y)$ is the marginal posterior of γ conditional on M_k .

- Predictive distribution includes **three sources of uncertainty**:
 - ▶ **Future errors**/disturbances (e.g. the ε 's in a regression)
 - ▶ **Parameter uncertainty** (the predictive distribution has the parameters integrated out by their posteriors)
 - ▶ **Model uncertainty** (by model averaging)

Pólya-Gamma augmentation for logistic regression

■ Logistic regression

$$\Pr(y = y_i | \mathbf{x}_i, \boldsymbol{\beta}) = \frac{\exp(\mathbf{x}_i^\top \boldsymbol{\beta})^{y_i}}{1 + \exp(\mathbf{x}_i^\top \boldsymbol{\beta})} \quad \text{for } y_i \in \{0, 1\}$$

$$\Pr(y_1, \dots, y_n | \mathbf{X}, \boldsymbol{\beta}) = \prod_{i=1}^n \Pr(y = y_i | \mathbf{x}_i, \boldsymbol{\beta})$$

■ The key identity

$$\frac{(e^\psi)^a}{(1 + e^\psi)^b} = 2^{-b} e^{\kappa \psi} \int_0^\infty e^{\omega \psi^2 / 2} p(\omega) d\omega,$$

where $\kappa = a - b/2$ and $p(\omega)$ is the density of the Pólya-Gamma distribution

$$\omega \sim \text{PG}(b, 0)$$

Pólya-Gamma augmentation for logistic regression

- So for each term in the likelihood function:

$$\frac{\exp(\mathbf{x}_i^\top \beta)^{y_i}}{1 + \exp(\mathbf{x}_i^\top \beta)} = \frac{1}{2} e^{\kappa_i \mathbf{x}_i^\top \beta} \int_0^\infty e^{\omega_i (\mathbf{x}_i^\top \beta)^2 / 2} p(\omega_i) d\omega_i,$$

- The likelihood conditional on $\omega = (\omega_1, \dots, \omega_n)$ is

$$\begin{aligned} \prod_{i=1}^n \frac{\exp(\mathbf{x}_i^\top \beta)^{y_i}}{1 + \exp(\mathbf{x}_i^\top \beta)} &\propto \prod_{i=1}^n e^{\kappa_i \mathbf{x}_i^\top \beta} e^{\omega_i (\mathbf{x}_i^\top \beta)^2 / 2} \\ &= \exp \left(\sum_{i=1}^n \kappa_i \mathbf{x}_i^\top \beta + \frac{\omega_i (\mathbf{x}_i^\top \beta)^2}{2} \right) \end{aligned}$$

which is an exponential of a quadratic form in β .

- Hence

$$\beta \sim N(\mu_0, \Sigma_0) \quad \implies \quad \beta | \omega, \mathbf{y}, \mathbf{X} \sim N(\mu_n, \Sigma_n)$$

Pólya-Gamma distribution

Pólya-Gamma distribution

$X \sim \text{PG}(b, c)$ for $X > 0$.

A Pólya-Gamma is defined as a infinite weighted sum (convolution) of iid Gamma distributed variables

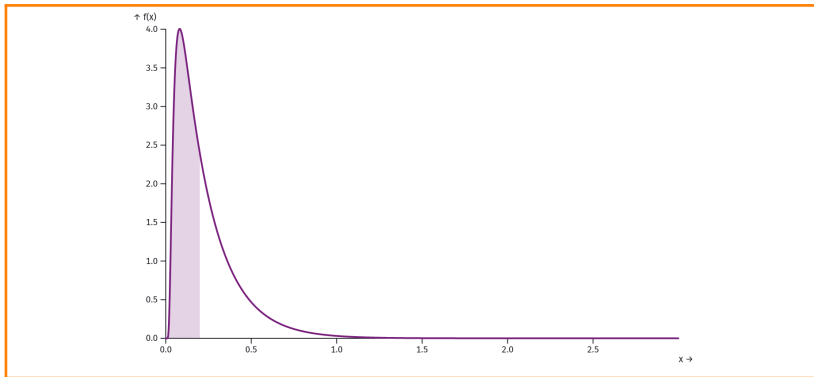
$$X \stackrel{d}{=} \sum_{k=1}^{\infty} v_k Y_k$$

where $\stackrel{d}{=}$ mean equality in distribution, the weights are

$$v_k = \frac{1}{2(k - 1/2)^2 \pi^2 + c^2/2}$$

and $Y_k \stackrel{\text{iid}}{\sim} \text{Gamma}(b, 1)$.

Pólya-Gamma distribution



Pólya-Gamma augmentation for logistic regression

Gibbs sampling for logistic regression using Pólya-Gamma augmentation

Input: response vector $\mathbf{y} = (y_1, \dots, y_n)^\top$
matrix $(n \times p)$ with covariates \mathbf{X}
initial value $\boldsymbol{\beta}^{(0)}$
number of posterior draws m .

$\boldsymbol{\kappa} \leftarrow (y_1 - 1/2, \dots, y_n - 1/2)^\top$

for k in $1:m$ **do**

 // Update Pólya-Gamma variables

for i in $1:n$ **do**

$\omega_i^{(k)} | \boldsymbol{\beta}^{(k-1)}, \mathbf{y}, \mathbf{x}_i \sim \text{PG}(1, \mathbf{x}_i^\top \boldsymbol{\beta}^{(k-1)})$

end

$\boldsymbol{\Omega}^{(k)} \leftarrow \text{Diag}(\omega_1^{(k)}, \dots, \omega_n^{(k)})$

 // Update $\boldsymbol{\beta}$

$\boldsymbol{\Sigma}_n \leftarrow (\mathbf{X}^\top \boldsymbol{\Omega}^{(k)} \mathbf{X} + \boldsymbol{\Sigma}_0^{-1})^{-1}$

$\boldsymbol{\mu}_n \leftarrow \boldsymbol{\Sigma}_n (\mathbf{X}^\top \boldsymbol{\kappa} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0)$

 Draw $\boldsymbol{\beta}^{(k)} | \boldsymbol{\omega}$ from $N(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$

end

Output: m draws $\boldsymbol{\beta}^{(1)}, \dots, \boldsymbol{\beta}^{(m)}$ from the posterior distribution $p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X})$.