

BAYESIAN LEARNING - LECTURE 12

Mattias Villani

**Division of Statistics and Machine Learning
Department of Computer and Information Science
Linköping University**

OVERVIEW

- ▶ Bayesian variable selection
- ▶ Summary and review of course material

BAYESIAN VARIABLE SELECTION

- ▶ Linear regression:

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

- ▶ Which variables have non-zero coefficient? Example of hypotheses:

$$H_0 : \beta_0 = \beta_1 = \dots = \beta_p = 0$$

$$H_1 : \beta_1 = 0$$

$$H_2 : \beta_1 = \beta_2 = 0$$

- ▶ Introduce **variable selection indicators** $\mathcal{I} = (I_1, \dots, I_p)$.
- ▶ Example: $\mathcal{I} = (1, 1, 0)$ means that $\beta_1 \neq 0$ and $\beta_2 \neq 0$, but $\beta_3 = 0$, so x_3 drops out of the model.

BAYESIAN VARIABLE SELECTION, CONT.

- ▶ Model inference, just crank the Bayesian machine:

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

- ▶ The prior $p(\mathcal{I})$ is typically taken to be $I_1, \dots, I_p | \theta \stackrel{iid}{\sim} \text{Bernoulli}(\theta)$.
- ▶ θ is the **prior inclusion probability**.
- ▶ Challenge: Computing the **marginal likelihood** for each model (\mathcal{I})

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

BAYESIAN VARIABLE SELECTION, CONT.

- ▶ Let $\beta_{\mathcal{I}}$ denote the **non-zero** coefficients under \mathcal{I} .
- ▶ Prior:

$$\begin{aligned}\beta_{\mathcal{I}}|\sigma^2 &\sim N\left(0, \sigma^2 \Omega_{\mathcal{I},0}^{-1}\right) \\ \sigma^2 &\sim \text{Inv} - \chi^2\left(\nu_0, \sigma_0^2\right)\end{aligned}$$

- ▶ **Marginal likelihood**

$$p(\mathbf{y}|\mathbf{X}, \mathcal{I}) \propto \left|\mathbf{X}'_{\mathcal{I}}\mathbf{X}_{\mathcal{I}} + \Omega_{\mathcal{I},0}^{-1}\right|^{-1/2} |\Omega_{\mathcal{I},0}|^{1/2} (\nu_0\sigma_0^2 + \text{RSS}_{\mathcal{I}})^{-(\nu_0+n-1)/2}$$

where $\mathbf{X}_{\mathcal{I}}$ is the covariate matrix for the subset given by \mathcal{I} .

- ▶ $\text{RSS}_{\mathcal{I}}$ is (almost) the residual sum of squares under model implied by \mathcal{I}

$$\text{RSS}_{\mathcal{I}} = \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbf{X}_{\mathcal{I}} \left(\mathbf{X}'_{\mathcal{I}}\mathbf{X}_{\mathcal{I}} + \Omega_{\mathcal{I},0}\right)^{-1} \mathbf{X}'_{\mathcal{I}}\mathbf{y}$$

BAYESIAN VARIABLE SELECTION VIA GIBBS SAMPLING

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

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

- ▶ Simulate from $p(\mathcal{I} | \mathbf{y})$ using Gibbs sampling:
 - ▶ Draw $l_1 | \mathcal{I}_{-1}, \mathbf{y}, \mathbf{X}$
 - ▶ Draw $l_2 | \mathcal{I}_{-2}, \mathbf{y}, \mathbf{X}$
 - ▶ ...
 - ▶ Draw $l_p | \mathcal{I}_{-p}, \mathbf{y}, \mathbf{X}$
- ▶ Only need to compute $Pr(l_i = 0 | \mathcal{I}_{-i}, \mathbf{y}, \mathbf{X})$ and $Pr(l_i = 1 | \mathcal{I}_{-i}, \mathbf{y}, \mathbf{X})$.
- ▶ Automatic model averaging, all in one simulation run.
- ▶ If needed, simulate from $p(\beta, \sigma^2 | \mathcal{I}, \mathbf{y}, \mathbf{X})$ for each draw of \mathcal{I} .

PSEUDO CODE FOR BAYESIAN VARIABLE SELECTION

0 Initialize $\mathcal{I}^{(0)} = (I_1^{(0)}, I_2^{(0)}, \dots, I_p^{(0)})$

1 Simulate σ^2 and β from [Note: $\nu_n, \sigma_n^2, \mu_n, \Omega_n$ all depend on $\mathcal{I}^{(0)}$]

▶ $\sigma^2 | \mathcal{I}^{(0)}, \mathbf{y}, \mathbf{X} \sim \text{Inv} - \chi^2 (\nu_n, \sigma_n^2)$

▶ $\beta | \sigma^2, \mathcal{I}^{(0)}, \mathbf{y}, \mathbf{X} \sim N [\mu_n, \sigma^2 \Omega_n^{-1}]$

2.1 Simulate $I_1 | \mathcal{I}_{-1}, \mathbf{y}, \mathbf{X}$ by [define $\mathcal{I}_{prop}^{(0)} = (1 - I_1^{(0)}, I_2^{(0)}, \dots, I_p^{(0)})$]

▶ compute marginal likelihoods: $p(\mathbf{y} | \mathbf{X}, \mathcal{I}^{(0)})$ and $p(\mathbf{y} | \mathbf{X}, \mathcal{I}_{prop}^{(0)})$

▶ Simulate $I_1^{(1)} \sim \text{Bernoulli}(\kappa)$ where

$$\kappa = \frac{p(\mathbf{y} | \mathbf{X}, \mathcal{I}^{(0)}) \cdot p(\mathcal{I}^{(0)})}{p(\mathbf{y} | \mathbf{X}, \mathcal{I}^{(0)}) \cdot p(\mathcal{I}^{(0)}) + p(\mathbf{y} | \mathbf{X}, \mathcal{I}_{prop}^{(0)}) \cdot p(\mathcal{I}_{prop}^{(0)})}$$

2.2 Simulate $I_2 | \mathcal{I}_{-2}, \mathbf{y}, \mathbf{X}$ as in Step 2.1, but $\mathcal{I}^{(0)} = (I_1^{(1)}, I_2^{(0)}, \dots, I_p^{(0)})$

⋮

2.P Simulate $I_p | \mathcal{I}_{-p}, \mathbf{y}, \mathbf{X}$ as in Step 2.1, but $\mathcal{I}^{(0)} = (I_1^{(1)}, I_2^{(1)}, \dots, I_p^{(0)})$

3 Repeat Steps 1-2 many times.

SIMPLE GENERAL BAYESIAN VARIABLE SELECTION

- ▶ The previous algorithm only works when we can integrate out all the model parameters to obtain

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

- ▶ MH - propose β and \mathcal{I} jointly from the proposal distribution

$$q(\beta_p|\beta_c, \mathcal{I}_p)q(\mathcal{I}_p|\mathcal{I}_c)$$

- ▶ Main difficulty: how to propose the non-zero elements in β_p ?
- ▶ Simple approach:
 - ▶ Approximate posterior with all variables in the model:
 $\beta|\mathbf{y}, \mathbf{X} \stackrel{approx}{\sim} N[\hat{\beta}, J_{\mathbf{y}}^{-1}(\hat{\beta})]$
 - ▶ Propose β_p from $N[\hat{\beta}, J_{\mathbf{y}}^{-1}(\hat{\beta})]$, conditional on the zero restrictions implied by \mathcal{I}_p . Formulas are available.