

Bayesian Learning

Lecture 12 - Variable selection

Mattias Villani

Department of Statistics
Stockholm University

Department of Computer and Information Science
Linköping University



 mattiasvillani.com

 [@matvil](https://twitter.com/matvil)



 [mattiasvillani](https://github.com/mattiasvillani)

Overview

- Bayesian variable selection
- Model averaging
- Posterior predictive analysis

Bayesian variable selection

- Linear regression:

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

- Which variables have **non-zero** coefficient?

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

$$H_1 : \beta_1 = 0$$

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

- 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

- Model inference, just crank the Bayesian machine:

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

- 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})

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

Bayesian variable selection

- 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**

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

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

- $\text{RSS}_{\mathcal{I}}$ is (almost) the residual sum of squares for model with \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 have essentially zero posterior probability. *Phew!*
- **Simulate** from the joint posterior distribution:

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

- Simulate from $p(\mathcal{I} | y, X)$ using **Gibbs sampling**:
 - ▶ Draw $l_1 | \mathcal{I}_{-1}, y, X$
 - ▶ Draw $l_2 | \mathcal{I}_{-2}, y, X$
 - ▶ ...
 - ▶ Draw $l_p | \mathcal{I}_{-p}, y, X$
- Note that: $Pr(l_i = 0 | \mathcal{I}_{-i}, y, X) \propto Pr(l_i = 0, \mathcal{I}_{-i} | y, X)$.
- Compute $p(\mathcal{I} | y, X) \propto p(y | X, \mathcal{I}) \cdot p(\mathcal{I})$ for $l_i = 0$ and for $l_i = 1$.
- **Model averaging** in a single simulation run.
- If needed, simulate from $p(\beta, \sigma^2 | \mathcal{I}, y, X)$ for each draw of \mathcal{I} .

Simple general Bayesian variable selection

- The previous algorithm only works when we can compute

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

- **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:

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

- ▶ Propose β_p from $N\left[\hat{\beta}, J_y^{-1}(\hat{\beta})\right]$, conditional on the zero restrictions implied by \mathcal{I}_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)

Posterior predictive analysis

- If $p(y|\theta)$ is a 'good' model, then the data actually observed should not differ 'too much' from simulated data from $p(y|\theta)$.
- Bayesian: simulate data from the **posterior predictive distribution**:

$$p(y^{rep}|y) = \int p(y^{rep}|\theta)p(\theta|y)d\theta.$$

- Difficult to compare y and y^{rep} because of dimensionality.
- Solution: compare **low-dimensional statistic** $T(y, \theta)$ to $T(y^{rep}, \theta)$.
- Evaluates the full probability model consisting of both the likelihood *and* prior distribution.

Posterior predictive analysis

- **Algorithm** for simulating from the posterior predictive density $p[T(y^{rep})|y]$:
 - 1 Draw a $\theta^{(1)}$ from the posterior $p(\theta|y)$.
 - 2 Simulate a data-replicate $y^{(1)}$ from $p(y^{rep}|\theta^{(1)})$.
 - 3 Compute $T(y^{(1)})$.
 - 4 Repeat steps 1-3 a large number of times to obtain a sample from $T(y^{rep})$.
- We may now compare the observed statistic $T(y)$ with the distribution of $T(y^{rep})$.
- **Posterior predictive p-value:** $\Pr[T(y^{rep}) \geq T(y)]$
- Informal graphical analysis.

Posterior predictive analysis - Normal model, max statistic

