# Bayesian variable selection

**Feng Li**
feng.li@cufe.edu.cn

**School of Statistics and Mathematics**
**Central University of Finance and Economics**

**Today we are going to learn...**

# Bayesian variable selection

- Linear regression:

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

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

$$
\begin{aligned}
H_0 &: \quad \beta_0 = \beta_1 = ... = \beta_p = 0 \\
H_1 &: \quad \beta_1 = 0 \\
H_2 &: \quad \beta_1 = \beta_2 = 0
\end{aligned}
$$

- Introduce **variable selection indicators** $\mathcal{I} = (I_1, ..., I_p)$.
- 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 $I_1, ..., I_p|\theta \overset{iid}{\sim} Bernoulli(\theta)$.
▶ $\theta$ is the **prior inclusion probability**.

# 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 $I_1, ..., I_p|\theta \overset{iid}{\sim} Bernoulli(\theta)$.

▶ $\theta$ 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:

$$\beta_{\mathcal{I}}|\sigma^2 \sim N\left(0, \sigma^2 \Omega_{\mathcal{I},0}^{-1}\right)$$

$$\sigma^2 \sim Inv - \chi^2\left(\nu_0, \sigma_0^2\right)$$

- **Marginal likelihood**

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

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

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

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

# 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!

# 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 $I_1|\mathcal{I}_{-1}, \mathbf{y}, \mathbf{X}$
  - ▶ Draw $I_2|\mathcal{I}_{-2}, \mathbf{y}, \mathbf{X}$
  - ▶ ...
  - ▶ Draw $I_p|\mathcal{I}_{-p}, \mathbf{y}, \mathbf{X}$

# 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 $I_1 | \mathcal{I}_{-1}, \mathbf{y}, \mathbf{X}$
    - ▶ Draw $I_2 | \mathcal{I}_{-2}, \mathbf{y}, \mathbf{X}$
    - ▶ ...
    - ▶ Draw $I_p | \mathcal{I}_{-p}, \mathbf{y}, \mathbf{X}$

- ▶ Only need to compute $Pr(I_i = 0 | \mathcal{I}_{-i}, \mathbf{y}, \mathbf{X})$ and $Pr(I_i = 1 | \mathcal{I}_{-i}, \mathbf{y}, \mathbf{X})$.
- ▶ 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 $\mathcal{I}^{(0)} = (I_1^{(0)}, I_2^{(0)} ..., I_p^{(0)})$

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

   - $\sigma^2 | \mathcal{I}^{(0)}, \mathbf{y}, \mathbf{X} \sim Inv - \chi^2 \left( \nu_n, \sigma_n^2 \right)$
   - $\beta | \sigma^2, \mathcal{I}^{(0)}, \mathbf{y}, \mathbf{X} \sim N \left[ \mu_n, \sigma^2 \Omega_n^{-1} \right]$

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

   - compute marginal likelihoods: $p(\mathbf{y}|\mathbf{X}, \mathcal{I}^{(0)})$ and $p(\mathbf{y}|\mathbf{X}, \mathcal{I}_{prop}^{(0)})$
   - Simulate $I_1^{(1)} \sim Bernoulli(\kappa)$ 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 $\mathcal{I}^{(0)} = (I_1^{(1)}, I_2^{(0)}, ..., I_p^{(0)})$

$\vdots$

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

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** $\beta$ 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 $\beta_p$?

▶ Simple approach:
  ▶ Approximate posterior with all variables in the model:
    $\beta|\mathbf{y}, \mathbf{X} \overset{approx}{\sim} N\left[\hat{\beta}, J_{\mathbf{y}}^{-1}(\hat{\beta})\right]$
  ▶ Propose $\beta_p$ from $N\left[\hat{\beta}, J_{\mathbf{y}}^{-1}(\hat{\beta})\right]$, conditional on the zero restrictions implied by $\mathcal{I}_p$. Formulas are available.

# VARIABLE SELECTION IN MORE COMPLEX MODELS

Posterior summary of the one-component split-$t$ model.[a]

| Parameters | Mean | Stdev | Post.Incl. |
|---|---|---|---|
| *Location $\mu$* | | | |
| Const | 0.084 | 0.019 | – |
| | | | |
| *Scale $\phi$* | | | |
| 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 |
| | | | |
| *Degrees of freedom $\nu$* | | | |
| 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 |
| | | | |
| *Skewness $\lambda$* | | | |
| 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 $\gamma$ be a quanitity with an interpretation which stays the same across the two models.

▶ Example: Prediction $\gamma = (y_{T+1}, ..., y_{T+h})'$.

▶ The marginal posterior distribution of $\gamma$ reads

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

where $p_k(\gamma|\mathbf{y})$ is the marginal posterior of $\gamma$ conditional on model $k$.

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

## Variable-selection priors

- The standard modern practice in Bayesian variable-selection problems is to treat variable inclusions as exchangeable Bernoulli trials with common success probability $p$.

- This implies that the **prior probability of a model** is given by

$$p(M_\gamma|p) = p^{k_\gamma}(1-p)^{m-k_\gamma}$$

with $k_\gamma$ representing the number of included variables in the model.

- This indicates that as $m$ grows with the true $k$ remaining fixed, the posterior distribution of $p$ will concentrate near 0. That means **using a fixed $p$ will yield a null model when $m$ is big** (no variable will be selected).

- Selecting $p = 1/2$ does not provide multiplicity correction. Treating $p$ as an unknown parameter to be estimated from the data will, however, yield an automatic multiple-testing penalty.

## Fully Bayesian variable-selection priors

- Assume that $p$ has a Beta distribution, $p \sim \text{Beta}(a, b)$, giving

$$p(M_\gamma) = \frac{\text{Beta}(a + k_\gamma, b + m - k_\gamma)}{\text{Beta}(a, b)}$$

- For the default choice of $a = b = 1$, implying a uniform prior on $p$

$$p(M_\gamma) = \frac{1}{m+1} \binom{m}{k_\gamma}^{-1}$$