



ECNM11060

Bayesian Econometrics

**A Bayesian linear
regression model**

Bayesian analysis of the linear regression model

- ⇒ Now see how general Bayesian theory of overview lecture works in familiar regression model
- ⇒ In lecture, we will focus on multiple regression under classical assumptions (independent errors, homoskedasticity, etc.)
- ⇒ Bayesian methods for freeing up classical assumptions are relatively straightforward (some of them we will discuss in the next lectures)

A regression model

⇒ Assume K explanatory variables, x_{i1}, \dots, x_{iK} , for $i = 1, \dots, N$, and a regression model:

$$y_i = \beta_1 + \beta_2 x_{i2} + \dots + \beta_K x_{iK} + \varepsilon_i$$

⇒ Note x_{i1} is implicitly set to 1 to allow for an intercept

⇒ Matrix notation:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

⇒ ε is $N \times 1$ vector stacked in same way as \mathbf{y}

A regression model

⇒ β is $K \times 1$ vector of regression coefficients

⇒ \mathbf{X} is $N \times K$ matrix of regressors

$$\mathbf{X} = \begin{bmatrix} 1 & X_{12} & \dots & X_{1K} \\ 1 & X_{22} & \dots & X_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{N2} & \dots & X_{NK} \end{bmatrix}$$

⇒ Regression model can be written as:

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

⇒ Classical assumptions imply: ε is $\mathcal{N}(\mathbf{0}_N, h^{-1}\mathbf{I}_N)$ where $h^{-1} = \sigma^2$

How to approach a regression problem as a Bayesian?

⇒ Suppose we are interested in estimating the following AR(1) model:

$$y_t = \rho y_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, h^{-1})$$

⇒ We have data \mathbf{y} and two unknown model parameters (ρ and σ^2)

⇒ We know that y_t is a stationary time series (e.g., a growth rate)

⇒ What else do we know about the data? (e.g., how is the data distributed?)

⇒ Given this information, what prior beliefs can we formulate about the AR(1) parameter ρ and the precision parameter h ?

⇒ In other words, we wish to specify the prior distributions $p(\rho)$ and $p(h)$

→ What do we know about ρ for certain?

→ What do we know about h for certain?

The likelihood function

- ⇒ Likelihood can be derived under the classical assumptions
- ⇒ All elements of \mathbf{X} are fixed (i.e., not random variables)
- ⇒ OLS quantities given by:

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}, \quad \nu = N - K, \quad s^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\beta})' (\mathbf{y} - \mathbf{X}\hat{\beta})}{\nu}$$

Likelihood function can be written in terms of OLS quantities
(see Exercise 10.1, Bayesian Econometric Methods):

$$p(\mathbf{y}|\beta, h) = \frac{1}{(2\pi)^{\frac{N}{2}}} \exp \left[-\frac{h}{2} (\beta - \hat{\beta})' \mathbf{X}'\mathbf{X} (\beta - \hat{\beta}) \right] \times h^{\frac{N}{2}} \exp \left[-\frac{h}{2} \nu s^2 \right]$$

The prior

⇒ Common starting point is natural conjugate prior

⇒ β conditional on h is now **multivariate Normal**:

$$\beta \mid h \sim \mathcal{N}(\underline{\beta}, h^{-1} \underline{\mathbf{V}})$$

⇒ Prior for error precision h is **Gamma**:

$$h \sim \mathcal{G}(\underline{a}, \underline{b})$$

⇒ $\underline{\beta}$, $\underline{\mathbf{V}}$, \underline{a} , and \underline{b} are prior hyperparameter values chosen by the researcher

The joint posterior of β and h

- ⇒ Multiplying likelihood by prior and collecting terms, the posterior can be shown to have the same form as the prior (see Exercise 10.1)
- ⇒ The conditional posterior for β (conditional on h) is **multivariate Normal**:

$$\beta \mid \mathbf{y}, h \sim \mathcal{N}(\bar{\beta}, h^{-1}\bar{\mathbf{V}})$$

with

$$\bar{\mathbf{V}} = (\underline{\mathbf{V}}^{-1} + \mathbf{X}'\mathbf{X})^{-1}$$
$$\bar{\beta} = \bar{\mathbf{V}} (\underline{\mathbf{V}}^{-1}\underline{\beta} + \mathbf{X}'\mathbf{X}\hat{\beta})$$

The joint posterior of β and h (cont.)

⇒ The (marginal) posterior for error precision h is **Gamma**:

$$h \mid \mathbf{y} \sim \mathcal{G}(\bar{a}, \bar{b})$$

with

$$\begin{aligned}\bar{a} &= \underline{a} + \frac{N}{2} \\ \bar{b} &= \underline{b} + \frac{1}{2} \left(\nu s^2 + (\hat{\beta} - \underline{\beta})' \left[\underline{\mathbf{V}} + (\mathbf{X}'\mathbf{X})^{-1} \right]^{-1} (\hat{\beta} - \underline{\beta}) \right) \\ &= \underline{b} + \frac{1}{2} \left(\mathbf{y}'\mathbf{y} + \underline{\beta}'\underline{\mathbf{V}}^{-1}\underline{\beta} - \underline{\beta}'\underline{\mathbf{V}}^{-1}\bar{\beta} \right)\end{aligned}$$

The marginal posterior of β

- ⇒ The marginal posterior for β is obtained by integrating out h
- ⇒ It turns out that β is multivariate t -distributed:

$$\beta \mid \mathbf{y} \sim t \left(\bar{\beta}, \frac{\bar{b} \bar{\mathbf{V}}}{\bar{a}}, 2\bar{a} \right)$$

- ⇒ Useful results for estimation:

$$\mathbb{E}(\beta \mid \mathbf{y}) = \bar{\beta}$$

$$\mathbb{V}(\beta \mid \mathbf{y}) = \frac{\bar{b}}{\bar{a} - 1} \bar{\mathbf{V}}$$

- ⇒ **Intuition:** Posterior mean and variance are weighted average of information in the prior and the data

What does a prior do?

⇒ To show main ideas assume:

→ $h = 1$ is known and β is a scalar

→ $\mathbf{X} = (1, \dots, 1)'$, such that $\mathbf{X}'\mathbf{X} = N$

→ $\underline{V} = \tau$ and $\underline{\beta} = 0$

⇒ In this case the posterior is:

$$\bar{\beta} = \frac{N}{\frac{1}{\tau} + N} \hat{\beta}$$

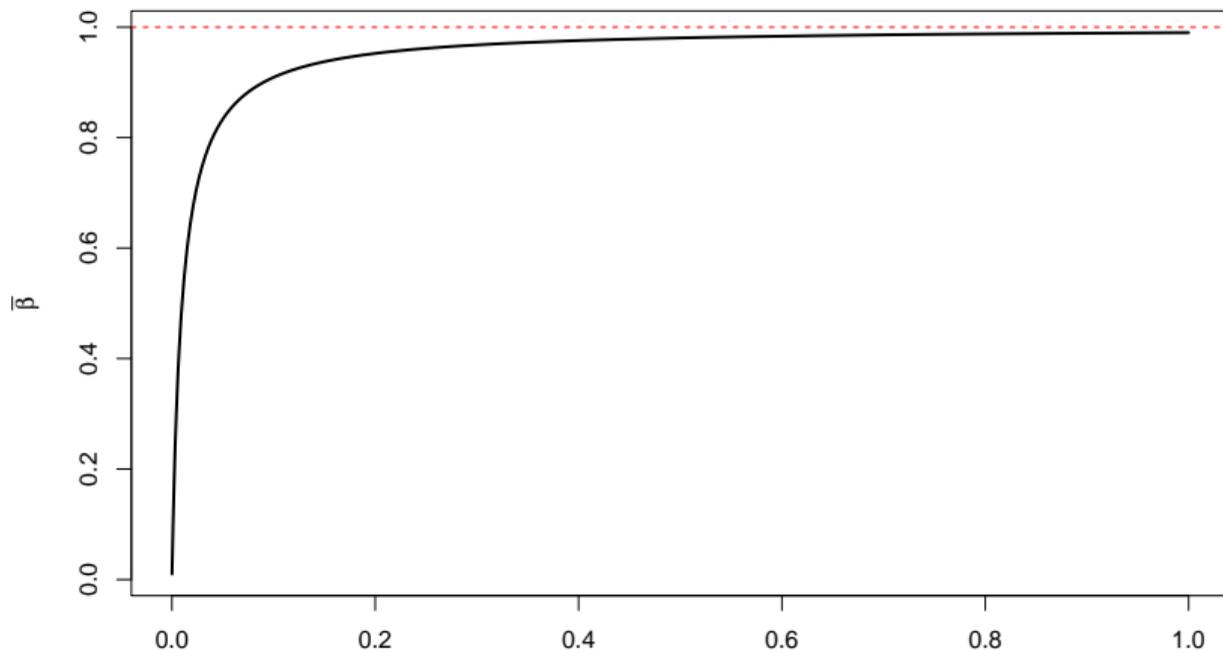
⇒ Think of the extreme cases ($\tau \rightarrow \infty$ and $\tau \rightarrow 0$)

→ If $\tau \rightarrow \infty$, then $\frac{1}{\tau} \rightarrow 0$ and $\frac{N}{\frac{1}{\tau} + N} \rightarrow 1$, so $\bar{\beta} \rightarrow \hat{\beta}$

→ If $\tau \rightarrow 0$, then $\frac{1}{\tau} \rightarrow \infty$ and $\frac{N}{\frac{1}{\tau} + N} \rightarrow 0$, so $\bar{\beta} \rightarrow 0$ (the prior mean)

Implied prior shrinkage by varying τ

$$\bar{\beta} = \frac{N}{1/\tau + N} \hat{\beta}, \text{ for } N = 100 \text{ and } \hat{\beta} = 1$$



Prior shrinkage

- ⇒ Posterior mean is pulled towards zero (“shrinkage”)
- ⇒ Commonly done to avoid over-fitting/over-parameterisation problems
- ⇒ Strength of prior shrinkage controlled through prior variance
 - If τ is small, then strong prior information β is near 0
 - If τ is big, then prior becomes more noninformative
- ⇒ Note: exactly what “small” and “large” means depends on the empirical application and units of measurement of data

Dummy observation view of conjugate shrinkage prior

- ⇒ In macroeconomics, we typically have a maximum of 700/800 observations available
- ⇒ Likely a sufficient amount of information to estimate a single parameter, but what about a large VAR model with thousands of parameters?
- ⇒ In such a case, we typically need additional information! (see Lecture 4)
- ⇒ Why not simply use a **theoretical model** to simulate data and combine it with real observed data?
- ⇒ Bayesians refer to this theoretical model as a **prior model** (our initial belief), as we can simulate observations from it without seeing any data

Dummy observation view of conjugate prior (cont.)

- ⇒ Conjugate prior amounts to pooling real observations, \mathbf{y} and \mathbf{X} , with dummy observations, $\underline{\mathbf{y}}$ and $\underline{\mathbf{X}}$, simulated from a theoretical model
- ⇒ Defining $\bar{\mathbf{y}} = (\mathbf{y}', \underline{\mathbf{y}})'$ and $\bar{\mathbf{X}} = (\mathbf{X}', \underline{\mathbf{X}})'$, the posterior estimator (via OLS) is given by:

$$\bar{\beta} = (\bar{\mathbf{X}}' \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}' \bar{\mathbf{y}}$$

$$\bar{\beta} = \left(\mathbf{X}' \mathbf{X} + \underbrace{\mathbf{X}' \underline{\mathbf{X}}}_{=\underline{\mathbf{V}}^{-1}} \right)^{-1} (\mathbf{X}' \mathbf{y} + \underline{\mathbf{X}}' \underline{\mathbf{y}})$$

$$\bar{\beta} = (\mathbf{X}' \mathbf{X} + \underline{\mathbf{V}}^{-1})^{-1} (\mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} + \underline{\mathbf{X}}' \underline{\mathbf{X}} (\underline{\mathbf{X}}' \underline{\mathbf{X}})^{-1} \underline{\mathbf{X}}' \underline{\mathbf{y}})$$

$$\bar{\beta} = (\mathbf{X}' \mathbf{X} + \underline{\mathbf{V}}^{-1})^{-1} (\mathbf{X}' \mathbf{X} \hat{\beta} + \underline{\mathbf{V}}^{-1} \underline{\beta})$$

A noninformative prior

- ⇒ Noninformative prior sets $\underline{a} = 0$ and \underline{V} is big (implies large prior uncertainty)
- ⇒ But there is not a unique way of doing the latter (see Exercise 10.4 in Bayesian Econometric Methods)
- ⇒ A common way: $\underline{V}^{-1} = \tau^{-1} \mathbf{I}_K$ where τ^{-1} is a scalar and let τ^{-1} go to zero
- ⇒ This noninformative prior is improper and becomes:

$$p(\beta, h) \propto \frac{1}{h}$$

- ⇒ With this choice we get OLS results:

$$\beta \mid h, \mathbf{y} \sim \mathcal{N}(\bar{\beta}, h^{-1} \bar{\mathbf{V}})$$

$$\text{with } \bar{\mathbf{V}} = (\mathbf{X}'\mathbf{X})^{-1} \text{ and } \bar{\beta} = \hat{\beta}$$

$$h \mid \mathbf{y} \sim \mathcal{G}(\bar{a}, \bar{b})$$

$$\text{with } \bar{a} = \frac{N}{2} \text{ and } \bar{b} = \frac{1}{2} (\nu s^2)$$

Model comparison

- ⇒ **Case 1:** M_1 imposes a linear restriction and M_2 does not (nested)
- ⇒ **Case 2:** $M_1 : \mathbf{y} = \mathbf{X}^{(1)}\boldsymbol{\beta}^{(1)} + \boldsymbol{\varepsilon}^{(1)}$ and $M_2 : \mathbf{y} = \mathbf{X}^{(2)}\boldsymbol{\beta}^{(2)} + \boldsymbol{\varepsilon}^{(2)}$, where $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ contain different explanatory variables (non-nested)
- ⇒ Both cases can be handled by defining models as (for $j = 1, 2$):

$$M_j : \mathbf{y} = \mathbf{X}^{(j)}\boldsymbol{\beta}^{(j)} + \boldsymbol{\varepsilon}^{(j)}$$

- ⇒ Nested model comparison defines M_2 as unrestricted regression
- ⇒ M_1 imposes the restriction can involve a redefinition of explanatory and dependent variable

Example: Nested model comparison

⇒ M_2 is unrestricted model

$$\mathbf{y} = \beta_1 + \beta_2 \mathbf{x}_2 + \beta_3 \mathbf{x}_3 + \epsilon$$

⇒ M_1 restricts $\beta_3 = 1$, can be written:

$$\begin{aligned}\mathbf{y} &= \beta_1 + \beta_2 \mathbf{x}_2 + \mathbf{1} \mathbf{x}_3 + \epsilon \\ \mathbf{y} - \mathbf{x}_3 &= \beta_1 + \beta_2 \mathbf{x}_2 + \epsilon\end{aligned}$$

⇒ M_1 has dependent variable $\mathbf{y} - \mathbf{x}_3$ and intercept and \mathbf{x}_2 are explanatory variables

Example: Marginal likelihood computation

⇒ Marginal likelihood is (for $j = 1, 2$):

$$p(\mathbf{y} | M_j) = c^{(j)} \left(\frac{|\overline{\mathbf{V}}^{(j)}|}{|\underline{\mathbf{V}}^{(j)}|} \right)^{\frac{1}{2}} \left(\overline{\mathbf{b}}^{(j)} \right)^{-\frac{\overline{a}^{(j)}}{2}}$$

⇒ $c^{(j)}$ is constant depending on prior hyperparameters, etc.

$$PO_{12} = \frac{c^{(1)} \left(\frac{|\overline{\mathbf{V}}^{(1)}|}{|\underline{\mathbf{V}}^{(1)}|} \right)^{\frac{1}{2}} \left(\overline{\mathbf{b}}^{(1)} \right)^{-\frac{\overline{a}^{(1)}}{2}}}{c^{(2)} \left(\frac{|\overline{\mathbf{V}}^{(2)}|}{|\underline{\mathbf{V}}^{(2)}|} \right)^{\frac{1}{2}} \left(\overline{\mathbf{b}}^{(2)} \right)^{-\frac{\overline{a}^{(2)}}{2}}} \times \frac{p(M_1)}{p(M_2)}$$

⇒ Posterior odds ratio depends on the prior odds ratio and contains rewards for model fit, coherency between prior and data information and parsimony

Model comparison with noninformative priors

Important Rule

- ⇒ **Noninformative priors** acceptable for parameters *common to all models*
- ⇒ **Informative, proper priors** required for all *model-specific* parameters
- ⇒ Prior for $h^{(1)}, h^{(2)}$: noninformative is **fine** (common to both models);
 $\underline{a}^{(1)} = \underline{a}^{(2)} = 0$ still yields sensible PO_{12}
- ⇒ Noninformative priors for $\beta^{(j)}$ **cause problems** when $K_1 \neq K_2$; e.g., $\underline{\mathbf{V}}^{(j)} = \tau \mathbf{I}_{K_j}$,
 $\tau \rightarrow \infty$, since $|\underline{\mathbf{V}}^{(j)}|$ depends on K_j , it does not cancel across models
- ⇒ Consequence:

Condition	Outcome
$K_1 < K_2$	$PO_{12} \rightarrow \infty$ (always favours M_1)
$K_1 > K_2$	$PO_{12} \rightarrow 0$ (always favours M_2)

Prediction

⇒ Want to predict:

$$y^* = \mathbf{x}^{*'}\boldsymbol{\beta} + \varepsilon^*$$

⇒ Remember, prediction is based on:

$$p(y^* | \mathbf{y}) = \int \int p(y^* | \mathbf{y}, \boldsymbol{\beta}, h) p(\boldsymbol{\beta}, h | \mathbf{y}) d\boldsymbol{\beta} dh$$

⇒ The resulting predictive:

$$y^* | \mathbf{y} \sim t \left(\mathbf{x}^{*'}\bar{\boldsymbol{\beta}}, \frac{\bar{b}}{\bar{a}} \{1 + \mathbf{x}^{*'}\bar{\mathbf{V}}\mathbf{x}^*\}, 2\bar{a} \right)$$

⇒ Model comparison, prediction and posterior inference about $\boldsymbol{\beta}$ can all be done analytically

⇒ So no need for posterior simulation in this model

⇒ However, let us illustrate Monte Carlo integration in this model

Monte Carlo integration

⇒ Remember the basic LLN we used for Monte Carlo integration

⇒ Let $\beta^{(s)}$, for $s = 1, \dots, S$, be a random sample from $p(\beta \mid \mathbf{y})$ and $g(\bullet)$ be any function and define

$$\hat{g}_S = \frac{1}{S} \sum_{s=1}^S g(\beta^{(s)})$$

⇒ then \hat{g}_S converges to $\mathbb{E}[g(\beta) \mid \mathbf{y}]$ as S goes to infinity

⇒ How would you write a computer program which did this?

Monte Carlo integration in action

- ⇒ *Step 1:* Take a random draw, $\beta^{(s)}$ from the posterior for β using a random number generator for the multivariate t distribution
- ⇒ *Step 2:* Calculate $g(\beta^{(s)})$ and keep this result
- ⇒ *Step 3:* Repeat Steps 1 and 2 S times
- ⇒ *Step 4:* Take the average of the S draws $g(\beta^{(1)}), \dots, g(\beta^{(S)})$
- ⇒ These steps will yield an estimate of $\mathbb{E}[g(\beta) | \mathbf{y}]$ for any function of interest
- ⇒ Remember: Monte Carlo integration yields only an approximation for $\mathbb{E}[g(\beta) | \mathbf{y}]$ (since you cannot set $S = \infty$)
- ⇒ By choosing S , can control the degree of approximation error

Monte Carlo integration in action (cont.)

In **R** this could look like this (you will do something similar in Matlab)

```
1 tXX <- crossprod(X)           # Compute X'X
2 tXy <- crossprod(X,y)        # Compute X'y
3 tyy <- crossprod(y)          # Compute y'y
4
5 # Posterior moments
6 V_po.inv <- (tXX + V_pr.inv)   # Posterior precision
7 V_po <- solve(V_po.inv)       # Posterior variance-covariance
8 Vchol_po <- t(chol(V_po))     # Cholesky factor of posterior VC
9 beta_po <- V_po%*(tXy + V_pr.inv%*b_pr) # Posterior mean
10
11 a_po <- s_pr + N/2            # Posterior DoF
12 b_po <- S_pr + (tyy + t(beta_pr)%*V_pr.inv%*beta_pr
13   - t(beta_po)%*V_po.inv%*beta_po)/2 # Posterior scaling
```

Monte Carlo integration in action (cont.)

In **R** this could look like this (you will do something similar in Matlab)

```
1 # Monte carlo integration
2 nsave <- 10000 # No. of draws
3
4 # Produce draws for h from the Gamma distribution
5 h_store <- rgamma(nsave, a_po, b_po)
6
7 # Produce draws for beta from the conditional Normal distribution beta|h
8 beta_store <- matrix(NA, nsave, K)
9 for(irep in 1:nsave)
10 {
11   beta_store[irep,] <- beta_po +
12     sqrt(1/h_store[irep])*Vchol_po %*% rnorm(K)
13 }
14
15 # Compute posterior mean of each element in beta
16 beta_pm <- apply(beta_store, 2, mean)
```

⇒ Note: By integrating over h , which is Gamma-distributed, β will be t -distributed

An empirical illustration

Empirical illustration

⇒ Dataset: $N = 546$ houses sold in Windsor, Canada (1987)

⇒ Dependent variable: $y_i =$ sales price (Canadian dollars)

⇒ Regressors:

x_{i2} : lot size (square feet)

x_{i3} : number of bedrooms

x_{i4} : number of bathrooms

x_{i5} : number of storeys

Setup of empirical illustration

- ⇒ Results shown for both **informative** and **noninformative** priors
- ⇒ Our prior implies statements of the form “if we compare two houses which are identical except the first house has one bedroom more than the second, then we expect the first house to be worth \$5,000 more than the second”; this yields prior mean, then choose large prior variance to indicate prior uncertainty
- ⇒ Tables on the next few slides present some empirical results (textbook has lots of discussion of how you would interpret them)
- ⇒ **95% HPDI** = Highest Posterior Density Interval: shortest interval $[a, b]$ such that

$$p(a \leq \beta_j \leq b \mid \mathbf{y}) = 0.95$$

Comparing prior and posterior means for β

	Prior	Posterior	
	Informative	Noninformative prior	Informative prior
β_1	0 (10,000)	-4,009.55 (3,593.16)	-4,035.05 (3,530.16)
β_2	10 (5)	5.43 (0.37)	5.43 (0.37)
β_3	5,000 (2,500)	2,824.61 (1,211.45)	2,886.81 (1,184.93)
β_4	10,000 (5,000)	17,105.17 (1,729.65)	16,965.24 (1,708.02)
β_5	10,000 (5,000)	7,634.90 (1,005.19)	7,641.23 (997.02)

Standard deviations in parentheses

Model comparison involving β

	$p(\beta_j > 0 \mid \mathbf{y})$	95% HPDI	PO for $\beta_j = 0$
Informative Prior			
β_1	0.13	[-10,957; 2,887]	4.14
β_2	1.00	[4.71; 6.15]	2.25×10^{-39}
β_3	0.99	[563.5; 5,210.1]	0.39
β_4	1.00	[13,616; 20,314]	1.72×10^{-19}
β_5	1.00	[5,686; 9,596]	1.22×10^{-11}
Noninformative Prior			
β_1	0.13	[-11,055; 3,036]	—
β_2	1.00	[4.71; 6.15]	—
β_3	0.99	[449.3; 5,200]	—
β_4	1.00	[13,714; 20,497]	—
β_5	1.00	[5,664; 9,606]	—

PO = posterior odds; — not defined under noninformative prior

Posterior results for β_2 (analytical versus Monte Carlo)

Method	Mean	Std. dev.	Numerical std. error
Analytical	5.4316	0.3662	
Monte Carlo (S replications)			
$S = 10$	5.3234	0.2889	0.0913
$S = 100$	5.4877	0.4011	0.0401
$S = 1,000$	5.4209	0.3727	0.0118
$S = 10,000$	5.4330	0.3677	0.0037
$S = 100,000$	5.4323	0.3664	0.0012

As S increases, MC estimates converge to the analytical solution and the numerical standard error shrinks

Summary so far

- ⇒ So far: Normal linear regression with **natural conjugate prior** ⇒ posterior, marginal likelihood and predictive distributions all available in **closed form**
- ⇒ Other priors / more complex models ⇒ **no analytical results**
- ⇒ Next: Extensions of the regression model requiring **Gibbs sampler** for posterior computation
 - Gibbs Sampler is a special case of **Markov Chain Monte Carlo (MCMC)**
 - Keep the Normal linear regression model (under the classical assumptions) as before
 - Likelihood function same as before
 - Parameters of model will still be β and h

The independent Normal-Gamma prior

⇒ Before we had conjugate prior where $p(\beta | h)$ was Normal density and $p(h)$ Gamma density

⇒ Now use similar prior, but assume prior **independence** between β and h

$$p(\beta, h) = p(\beta) p(h)$$

⇒ $p(\beta)$ is Normal and $p(h)$ is Gamma:

$$\beta \sim \mathcal{N}(\underline{\beta}, \underline{\mathbf{V}})$$

$$h \sim \mathcal{G}(\underline{a}, \underline{b})$$

⇒ Key difference: now $\underline{\mathbf{V}}$ is now the prior covariance matrix of β , with conjugate prior we had $\mathbb{V}(\beta | h) = h^{-1} \underline{\mathbf{V}}$

The posterior

- ⇒ Posterior \propto prior \times likelihood
- ⇒ Joint posterior for β and h : **no closed form** \Rightarrow cannot be used directly for inference
- ⇒ But **conditional posterior** for β given h takes a simple form:

$$\beta \mid \mathbf{y}, h \sim \mathcal{N}(\bar{\beta}, \bar{\mathbf{V}})$$

with

$$\bar{\mathbf{V}} = \left(\underline{\mathbf{V}}^{-1} + h\mathbf{X}'\mathbf{X} \right)^{-1}$$

$$\bar{\beta} = \bar{\mathbf{V}} \left(\underline{\mathbf{V}}^{-1} \underline{\beta} + h\mathbf{X}'\mathbf{y} \right)$$

The posterior (cont.)

⇒ **Conditional posterior** for h given β is Gamma:

$$h \mid \mathbf{y}, \beta \sim \mathcal{G}(\bar{a}, \bar{b})$$

with

$$\bar{a} = \underline{a} + \frac{N}{2}, \quad \bar{b} = \underline{b} + (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$$

⇒ Econometrician is interested in **joint posterior** $p(\beta, h \mid \mathbf{y})$ or marginal $p(\beta \mid \mathbf{y})$, **not** in the conditional posteriors

⇒ **Key point:** Since

$$p(\beta, h \mid \mathbf{y}) \neq p(\beta \mid \mathbf{y}, h) p(h \mid \mathbf{y}, \beta)$$

the conditionals **do not directly characterize** the joint posterior

Popular MCMC algorithms

The Gibbs sampler

- ⇒ The **Gibbs sampler** uses conditional posteriors to draw $\beta^{(s)}$ and $h^{(s)}$ for $s = 1, \dots, S \Rightarrow$ averages yield posterior estimates, as with Monte Carlo integration
- ⇒ Powerful and widely-used tool for posterior simulation in econometric models
- ⇒ **General setup:** θ is a p -vector of parameters with likelihood $p(\mathbf{y} \mid \theta)$, prior $p(\theta)$ and posterior $p(\theta \mid \mathbf{y})$
- ⇒ Partition θ into B **blocks**:

$$\theta = \left(\theta'_{(1)}, \theta'_{(2)}, \dots, \theta'_{(B)} \right)'$$

- ⇒ For example, regression model: $B = 2$ with $\theta_{(1)} = \beta$ and $\theta_{(2)} = h$

Intuition of the Gibbs sampler

- ⇒ **Monte Carlo integration:** Draw from $p(\theta | \mathbf{y})$ and average to estimate $\mathbb{E}[g(\theta) | \mathbf{y}]$ for any function $g(\theta)$
- ⇒ **Problem:** Drawing directly from $p(\theta | \mathbf{y})$ is often **not possible** (e.g., not available in closed form)
- ⇒ But it often is easy to draw from **full conditional posteriors:**

$$p(\theta_{(b)} | \mathbf{y}, \{\theta_{(j)}\}_{j \neq b}), \quad b = 1, \dots, B$$

- ⇒ **Gibbs sampler:** Cycle through draws from each full conditional; in turn, produces a sequence $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(S)}$ that can be averaged to estimate $\mathbb{E}[g(\theta) | \mathbf{y}]$, just as with Monte Carlo integration

More motivation for the Gibbs sampler

⇒ Regression model: $B = 2$ blocks, $\theta = \{\beta, h\}$

⇒ Suppose we have one draw $\beta^{(0)}$ from $p(\beta | \mathbf{y})$

⇒ Then since

$$p(\beta, h | \mathbf{y}) = p(h | \mathbf{y}, \beta) p(\beta | \mathbf{y})$$

a draw from $p(h | \mathbf{y}, \beta^{(0)})$ is a valid draw of h ; call this $h^{(1)}$

⇒ Similarly, since

$$p(\beta, h | \mathbf{y}) = p(\beta | \mathbf{y}, h) p(h | \mathbf{y})$$

a draw from $p(\beta | \mathbf{y}, h^{(1)})$ is a valid draw of β ; call this $\beta^{(1)}$

⇒ Hence $\{\beta^{(1)}, h^{(1)}\}$ is a valid draw from $p(\beta, h | \mathbf{y})$ and continuing yields $\{\beta^{(s)}, h^{(s)}\}$ for $s = 1, \dots, S$

More motivation for the Gibbs sampler (cont.)

- ⇒ Hence, if you can successfully find $\beta^{(0)}$, then sequentially drawing $p(h | \mathbf{y}, \beta)$ and $p(\beta | \mathbf{y}, h)$ will give valid draws from posterior
- ⇒ Problem with above strategy is that it is not possible to find such initial $\beta^{(0)}$
- ⇒ If we knew how to easily take random draws from $p(\beta | \mathbf{y})$, we could use this and $p(h | \beta, \mathbf{y})$ to do Monte Carlo integration and have no need for Gibbs sampling
- ⇒ However, it can be shown that subject to weak conditions, the initial draw $\beta^{(0)}$ does not matter: Gibbs sampler will converge to a sequence of draws from $p(\beta, h | \mathbf{y})$
- ⇒ In practice, choose $\beta^{(0)}$ in some manner and then run the Gibbs sampler for S replications
- ⇒ Discard S_0 initial draws (the so-called “**burn-in**”) and remaining S_1 used to estimate $\mathbb{E}[g(\theta) | \mathbf{y}]$

Why is Gibbs sampling so useful?

- ⇒ In Normal linear regression model with independent Normal-Gamma prior, Gibbs sampler is easy
- ⇒ $p(\beta | \mathbf{y}, h)$ is Normal and $p(h | \mathbf{y}, \beta)$ is Gamma (easy to draw from)
- ⇒ Huge number of other models have hard joint posterior, but easy posterior conditionals
- ⇒ Tobit, probit, stochastic frontier model, Markov switching model, threshold autoregressive, smooth transition threshold autoregressive, other regime switching models, state space models, some semiparametric regression models, etc.
- ⇒ What if the full posterior conditionals do not have simple form?
- ⇒ Many other algorithms exist for handling general cases, Metropolis-Hastings algorithm is most popular

The Metropolis-Hastings (MH) algorithm

- ⇒ Another popular MCMC algorithm, useful when **Gibbs sampling is not feasible**
- ⇒ General notation: θ is a vector of parameters with likelihood $p(\mathbf{y} | \theta)$, prior $p(\theta)$ and posterior $p(\theta | \mathbf{y})$
- ⇒ Key idea: Draw from a convenient **candidate generating density** $q(\theta^{(s-1)}; \theta)$, where θ^* denotes a candidate draw whose density depends on $\theta^{(s-1)}$
- ⇒ Since we draw from $q(\theta^{(s-1)}; \theta)$ instead of $p(\theta | \mathbf{y})$, we have to **correct for this** via an **acceptance probability**, with only some candidate draws being accepted

MH algorithm: Steps

1. Choose starting value $\theta^{(0)}$
2. Draw candidate θ^* from $q(\theta^{(s-1)}; \theta)$
3. Compute acceptance probability $\alpha(\theta^{(s-1)}, \theta^*)$
4. Set $\theta^{(s)} = \begin{cases} \theta^* & \text{with probability } \alpha(\theta^{(s-1)}, \theta^*) \\ \theta^{(s-1)} & \text{with probability } 1 - \alpha(\theta^{(s-1)}, \theta^*) \end{cases}$
5. Repeat Steps 2 to 4 for $s = 1, \dots, S(= S_0 + S_1)$; discard S_0 burn-in draws
6. Estimate $\mathbb{E}[g(\theta) \mid \mathbf{y}]$ by averaging $g(\theta^{(S_0+1)}), \dots, g(\theta^{(S)})$

The acceptance probability

⇒ The acceptance probability is:

$$\alpha(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^*) = \min \left[\frac{p(\boldsymbol{\theta}^* | \mathbf{y}) q(\boldsymbol{\theta}^*; \boldsymbol{\theta}^{(s-1)})}{p(\boldsymbol{\theta}^{(s-1)} | \mathbf{y}) q(\boldsymbol{\theta}^{(s-1)}; \boldsymbol{\theta}^*)}, 1 \right]$$

⇒ Intuition: Candidate $\boldsymbol{\theta}^*$ is **always accepted** if it has higher posterior density than $\boldsymbol{\theta}^{(s-1)}$; otherwise accepted with probability < 1 (see textbook pp. 93–94)

⇒ As with Gibbs sampling, discard S_0 **burn-in** draws to remove dependence on starting value $\boldsymbol{\theta}^{(0)}$

⇒ Note: The **Gibbs sampler** is a special case of MH where the **full conditionals** serve as candidate generating densities for each block $\boldsymbol{\theta}_{(b)}$, resulting in $\alpha(\boldsymbol{\theta}_{(b)}^{(s-1)}, \boldsymbol{\theta}_{(b)}^*) = 1$, with every candidate draw being accepted

Candidate generating density: Independence chain

⇒ Candidate density does not depend on $\theta^{(s-1)}$:

$$q(\theta^{(s-1)}; \theta) = q(\theta)$$

⇒ Useful when a convenient **approximation to the posterior** exists

⇒ Acceptance probability simplifies to:

$$\alpha\left(\theta^{(s-1)}, \theta^*\right) = \min \left[\frac{p(\theta^* | \mathbf{y}) q(\theta^{(s-1)})}{p(\theta^{(s-1)} | \mathbf{y}) q(\theta^*)}, 1 \right]$$

Candidate generating density: Random walk chain

⇒ **Random walk chain** popular in DSGE models; useful when no good posterior approximation exists

⇒ Candidate draw:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}^{(s-1)} + \boldsymbol{w}$$

where \boldsymbol{w} is the **increment random variable**

⇒ Acceptance probability simplifies to:

$$\alpha\left(\boldsymbol{\theta}^{(s-1)}, \boldsymbol{\theta}^*\right) = \min \left[\frac{p(\boldsymbol{\theta}^* | \boldsymbol{y})}{p(\boldsymbol{\theta}^{(s-1)} | \boldsymbol{y})}, 1 \right]$$

⇒ Common choice: $\boldsymbol{w} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$, so $q(\boldsymbol{\theta}^{(s-1)}; \boldsymbol{\theta}) = f_N(\boldsymbol{\theta} | \boldsymbol{\theta}^{(s-1)}, \boldsymbol{\Sigma})$

⇒ Researcher must choose $\boldsymbol{\Sigma} = c\boldsymbol{\Omega}$, where $\boldsymbol{\Omega}$ is an estimate of the posterior covariance (e.g., inverse Hessian at posterior mode) and c is a scalar tuned so that the acceptance rate is **neither too high nor too low** (rule of thumb: between 0.2 and 0.4)

Example of independence chain MH

⇒ Consider a simple AR(1) model:

$$y_t = \beta y_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, h^{-1})$$

⇒ We have prior knowledge that y_t is **persistent but stationary**: $\beta \in (0, 1)$

⇒ This may motivate choosing a **Beta prior**: $\beta \sim \mathcal{B}(\underline{c}, \underline{d})$

⇒ **Problem**: Gaussian likelihood \times Beta prior results in a posterior $p(\beta | \mathbf{y})$ with **no closed form**

⇒ In addition, no well-known conditionals available; Gibbs sampling not possible

⇒ Use independence chain MH with a **Uniform candidate density**
 $q(\beta) = \mathcal{U}(0, 1)$

Bayesian model averaging (BMA)

BMA in a nutshell

- ⇒ BMA can be used with any set of models; here use it in a Big Data regression (many explanatory variables)
- ⇒ **Model selection:** Choose a single model and present estimates or forecasts based on it
- ⇒ **Model averaging:** Take a weighted average of estimates or forecasts from all models with weights given by $p(M_r | \mathbf{y})$, incorporating uncertainty about which model generated the data
- ⇒ Allows for a formal treatment of **model uncertainty**
- ⇒ Let M_r , for $r = 1, \dots, R$, denote R competing models, for any parameter ϕ or forecast of interest, the rules of probability imply:

$$p(\phi | \mathbf{y}) = \sum_{r=1}^R p(\phi | \mathbf{y}, M_r) p(M_r | \mathbf{y})$$

Model space of a big data regression

⇒ Let $\mathbf{X}^{(r)}$ be $N \times K_r$ submatrix of X ; each model takes the form:

$$\mathbf{y} = \alpha \mathbf{1}_N + \mathbf{X}^{(r)} \boldsymbol{\beta}^{(r)} + \boldsymbol{\varepsilon}$$

with $\mathbf{1}_N$ being $N \times 1$ vector of ones (intercept included in all models)

⇒ With K potential regressors there are $R = 2^K$ possible models (each regressor is either included or excluded; 2 options per regressor)

⇒ **Computational challenge:** For $K = 41$, $R = 2^{41} \approx 2$ trillion models (at 0.001 seconds per model, estimation would take **over 100 years**)

⇒ Use a **natural conjugate prior** to maximize estimation speed per model (closed-form solutions), combined with MCMC model composition to explore the model space without estimating every model

BMA priors: The g-prior

- ⇒ We want a prior for model M_r that is **informative** (valid marginal likelihoods), **objective** (minimal subjective input), and **automatic** (no individual choice per model)
- ⇒ The **g-prior** (Zellner, 1986) satisfies all three
- ⇒ Prior shrinks coefficients towards zero:

$$\underline{\beta}^{(r)} = \mathbf{0}, \quad \underline{\mathbf{V}}^{(r)} = \left(g \mathbf{X}^{(r)'} \mathbf{X}^{(r)} \right)^{-1}$$

- ⇒ Prior information that $\beta^{(r)} = \mathbf{0}$ takes the **same form** as data information
- ⇒ g is a scalar controlling the **relative weight** of prior vs data information:
 - $g = 1$: equal weight to prior and data; $g = 0.01$: prior receives 1% of the weight of data; g can also be chosen via rules of thumb or treated as unknown
- ⇒ Noninformative prior typically used for h

BMA posterior

- ⇒ With natural conjugate prior, model M_r yields **analytical results**
 - Posterior: Normal-Gamma
 - Marginal likelihood: Closed form
 - Predictive density: Student- t
- ⇒ For each model, everything needed can be computed **quickly**
- ⇒ However, for $K > 20$, evaluating all 2^K models remains **computationally infeasible**

BMA computation: MC³ algorithm

⇒ **MC³** (Markov Chain Monte Carlo Model Composition): Draws models $M^{(s)}$, for $s = 1, \dots, S$, instead of parameters

⇒ For any quantity of interest ϕ :

$$\hat{\phi} = \frac{1}{S} \sum_{s=1}^S \mathbb{E}(\phi \mid \mathbf{y}, M^{(s)}) \xrightarrow{S \rightarrow \infty} \mathbb{E}(\phi \mid \mathbf{y})$$

⇒ Frequencies with which models are drawn converge to **posterior model probabilities**: if M_i drawn A times and M_j drawn B times, A/B converges to the Bayes factor BF_{ij}

⇒ Discard initial draws as **burn-in**

MC³: How are models drawn?

⇒ Given current model $M^{(s-1)}$, propose candidate M^* drawn randomly with equal probability from the following set of models:

→ $M^{(s-1)}$

→ All models deleting one regressor from $M^{(s-1)}$

→ All models adding one regressor to $M^{(s-1)}$

⇒ Accept M^* with probability:

$$\alpha\left(M^{(s-1)}, M^*\right) = \min\left[\frac{p(\mathbf{y} | M^*) p(M^*)}{p(\mathbf{y} | M^{(s-1)}) p(M^{(s-1)})}, 1\right]$$

⇒ Set $M^{(s)} = M^*$ if accepted, else $M^{(s)} = M^{(s-1)}$

⇒ Can be shown MC³ converges to the true **BMA posterior**

Classic BMA application

Determinants of economic growth

- ⇒ Classic cross-country growth regression: Why do some countries grow faster than others?
- ⇒ Dependent variable: Average GDP per capita growth, 1960–1992
- ⇒ $K = 41$ potential explanatory variables (education, investment, governance, institutions, trade, etc.), all standardized
- ⇒ $N = 72$ countries (cross-sectional observations)
- ⇒ This is Big Data: **big K , small N**

Classic BMA application (cont.)

Determinants of economic growth

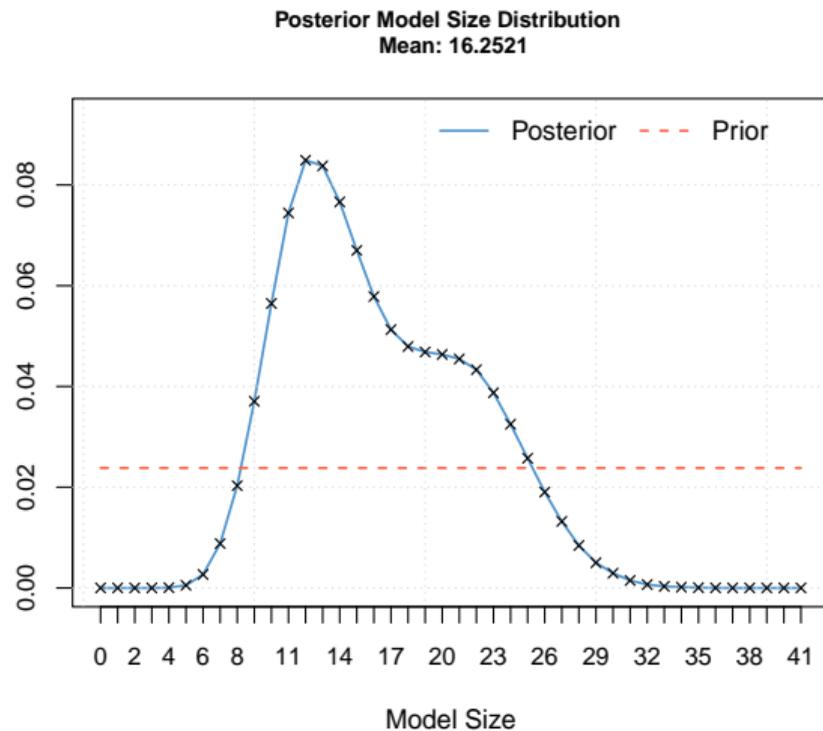
- ⇒ Cross-country growth regression data set with $N = 72$ and $K = 41$
- ⇒ Use common recommendation to set $g = \frac{1}{N}$ (UIP; unit information prior) and $g = \frac{1}{K^2}$ (RIC; risk inflation criteria)
- ⇒ Run MC³ algorithm for 2, 200, 000 draws, discarding first 200, 000 as burn-in
- ⇒ Is this enough draws?
- ⇒ Convergence diagnostic: Calculate posterior model probabilities analytically and using MC³ and compare
- ⇒ Model selection puts all weight on this single model — ignoring huge amount of model uncertainty

Replicating the analysis in R

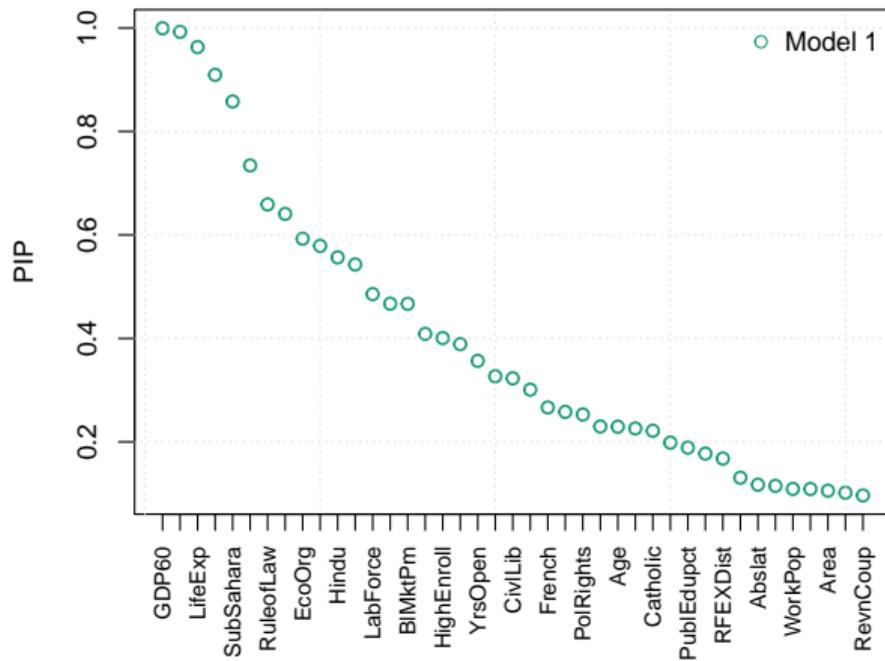
Bayesian Model Averaging can be implemented using the BMS package

```
1 # Load package
2 library(BMS)
3
4 # Load data
5 data(datafls)
6
7 # Estimate BMA model
8 bma.flr <- bms(
9   datafls,
10  burn = 200000,      # burn-in draws
11  iter = 2000000,    # posterior draws
12  g     = "UIP"      # unit information prior
13 )
```

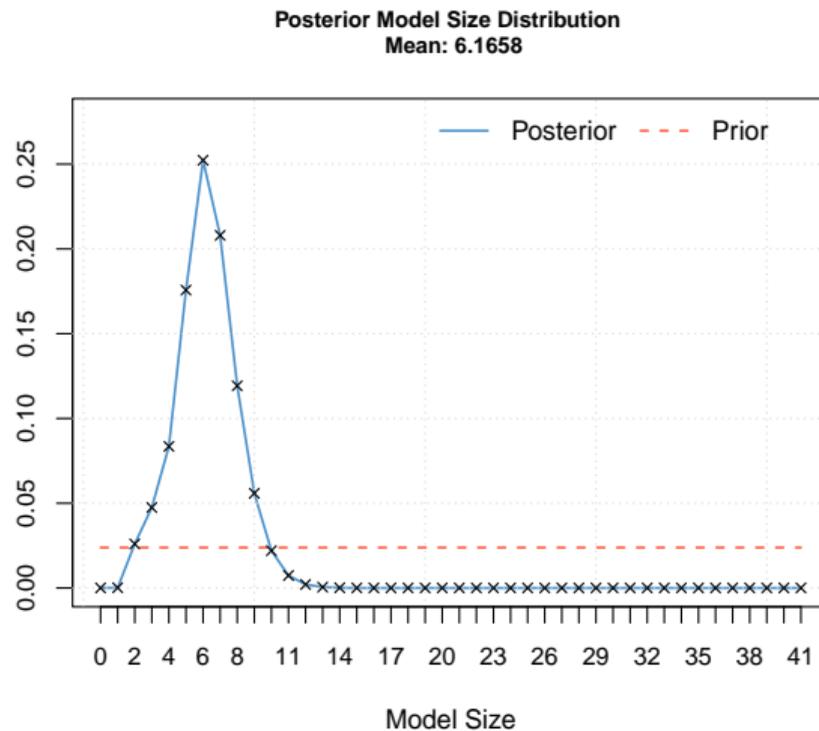
Posterior model size distribution for $g = \text{UIP}$



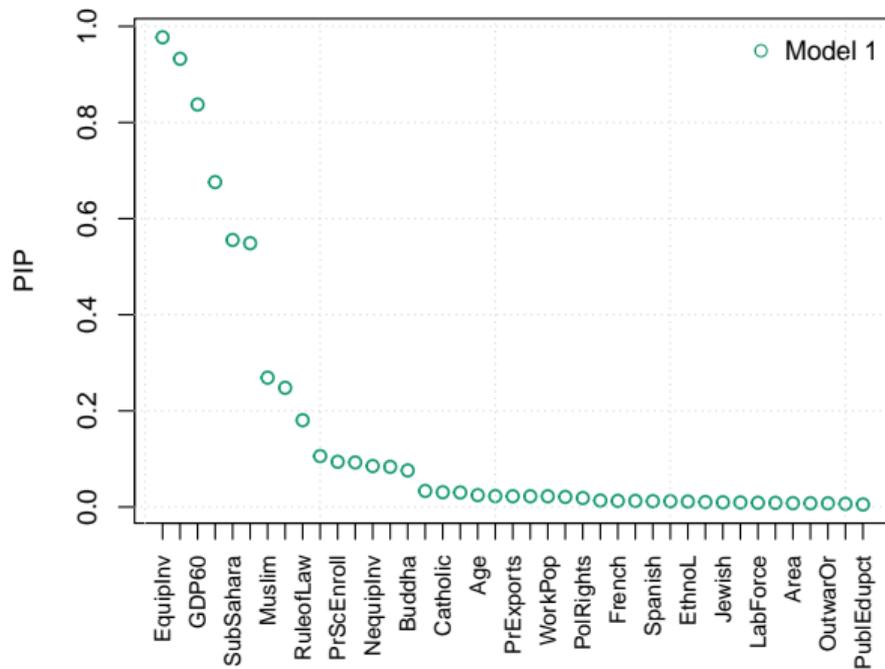
Posterior model size distribution for $g = \text{UIP}$



Posterior model size distribution for $g = \text{RIC}$



Posterior model size distribution for $g = \text{RIC}$



BMA application

- ⇒ Next table presents results:
- ⇒ Posterior mean and standard deviation for each explanatory variable using BMA and BMS
- ⇒ Rule of thumb: If an estimate (posterior mean) more than two standard deviations from zero likely to be important
- ⇒ Column labelled "Prob." = probability that the corresponding explanatory variable should be included (= proportion of models drawn by MC³ which contain the corresponding explanatory variable)
- ⇒ BMS ensures parsimony by choosing 14 variables
- ⇒ By ignoring model uncertainty estimates are more precise (smaller standard dev.)
- ⇒ BMA ensures parsimony by averaging over many small models
- ⇒ Average number of explanatory variables in a model drawn by MC³ is 16.25

Point Estimates and Standard Devs of Regression Coefficients

(Mean and standard deviations multiplied by 100)

Explanatory Variable	BMA			BMS	
	Prob.	Mean.	St. Dev.	Mean	St. Dev.
Primary School Enrolment	0.207	0.104	0.234	0.048	0.018
Life expectancy	0.933	0.961	0.392	0.090	0.020
GDP level in 1960	0.999	-1.425	0.278	-1.463	0.193
Fraction GDP in Mining	0.459	0.147	0.181	0.322	0.108
Degree of Capitalism	0.457	0.151	0.183	0.387	0.094
No. Years Open Economy	0.513	0.260	0.283	0.557	0.138
% Pop. Speaking English	0.069	-0.011	0.047	-	-
% Pop. Speak. For. Lang.	0.068	0.012	0.059	-	-
Exchange Rate Distortions	0.082	-0.017	0.070	-	-
Equipment Investment	0.923	0.552	0.236	0.548	0.128
Non-equipment Investment	0.434	0.136	0.174	0.347	0.099
St. Dev. of Black Mkt. Prem.	0.048	-0.006	0.037	-	-
Outward Orientation	0.037	-0.003	0.029	-	-

Point Estimates and Standard Devs of Regression Coefficients

(Mean and standard deviations multiplied by 100)

Explanatory Variable	BMA			BMS	
	Prob.	Mean.	St. Dev.	Mean	St. Dev.
Black Market Premium	0.179	-0.040	0.097	-	-
Area	0.030	-0.001	0.021	-	-
Latin America	0.215	-0.082	0.191	-	-
Sub-Saharan Africa	0.738	-0.473	0.347	-0.543	0.124
Higher Education Enrolment	0.046	-0.008	0.056	-	-
Public Education Share	0.032	-0.001	0.024	-	-
Revolutions and Coups	0.031	-0.001	0.023	-	-
War	0.075	-0.014	0.062	-	-

Posterior Estimates and Standard Devs of Regression Coefficients					
Explanatory Variable	Bayesian Model Averaging			Single Best Model	
	Prob.	Mean	St. Dev.	Mean	St. Dev.
Political Rights	0.094	-0.028	0.107	-	-
Civil Liberties	0.131	-0.050	0.015	-0.284	0.176
Latitude	0.041	0.001	0.052	-	-
Age	0.085	-0.015	0.058	-	-
British Colony	0.041	-0.003	0.032	-	-
Fraction Buddhist	0.196	0.047	0.109	-	-
Fraction Catholic	0.128	-0.011	0.121	-	-
Fraction Confucian	0.990	0.493	0.127	0.503	0.090
Ethnolinguistic Fractionalization	0.060	0.010	0.056	-	-
French Colony	0.049	0.007	0.040	-	-

Posterior Estimates and Standard Devs of Regression Coefficients					
Explanatory Variable	Bayesian Model Averaging			Single Best Model	
	Prob.	Mean	St. Dev.	Mean	St. Dev.
Fraction Hindu	0.126	-0.035	0.120	-	-
Fraction Jewish	0.037	-0.002	0.028	-	-
Fraction Muslim	0.640	0.025	0.023	0.295	0.093
Primary Exports	0.100	-0.029	0.105	-0.352	0.136
Fraction Protestant	0.455	-0.143	0.178	-0.277	0.098
Rule of Law	0.489	0.244	0.279	0.563	0.134
Spanish Colony	0.058	0.010	0.068	-	-
Population Growth	0.037	0.005	0.048	-	-
Ratio Workers to Population	0.045	-0.005	0.043	-	-
Size of Labor Force	0.075	0.018	0.097	-	-

Summary

- ⇒ This lecture shows how Bayesian ideas work in familiar context (regression model)
- ⇒ Occasionally analytical results are available (no need for posterior simulation)
- ⇒ Usually posterior simulation is required
- ⇒ Monte Carlo integration is simplest, but rarely possible to use it
- ⇒ Gibbs sampling (and related MCMC) methods can be used for estimation and prediction for a wide variety of models
- ⇒ Metropolis-Hastings algorithms popular and can be combined with Gibbs sampling (Metropolis-within-Gibbs)
- ⇒ Note: There are methods for calculating marginal likelihoods using Gibbs sampler output