

Variable Selection for Bayesian Linear Regression Model in a Finite Sample Size

Satoshi KABE¹

Graduate School of Systems and Information Engineering,
University of Tsukuba.

Yuichiro KANAZAWA²

Graduate School of Systems and Information Engineering,
University of Tsukuba.

1 Introduction

Data analysis often involves a comparison of several candidate models. Because true model is seldom known a priori, there is a need for a simple, effective, and objective methods for the selection of the best approximating model. Akaike (1973) proposed an information criterion later known to be Akaike's information criterion (AIC) based on the concept of minimizing the Kullback-Leibler Information, a measure of discrepancy between the true density (or model) and approximating model: The discrepancy between two models or two probability densities is expressed by the expected log-likelihood with respect to the true density. AIC is designed to be an approximately unbiased estimator of the expected log-likelihood under the assumption that "true model" is one of the candidate models being considered.

Hurvich and Tsai (1989) illustrated that AIC can be dramatically biased when the sample size is small by comparing AIC with the finite bias-corrected version of AIC (AIC_C) proposed by Sugiura (1978). This criterion adjusts AIC to be an exact unbiased estimator of the expected log-likelihood. They then extend AIC_C so that it can be employed for autoregressive (AR) model and autoregressive moving average (ARMA) model. Later Bedrick and Tsai (1994) further extend Sugiura (1978) to multivariate regression cases where response variable is expressed by a matrix.

In the 'fully' Bayesian data analysis, the deviance information criterion (DIC) is widely used for the model selection. Spiegelhalter et al. (2002) proposed a Bayesian measure of model complexity (i.e., effective number of parameters p_D) obtained from the difference between the posterior mean of the deviance and the deviance at the posterior mean of the parameters. When the number of data is sufficiently large, DIC is given by adding p_D to the posterior mean of the deviance. Spiegelhalter et al. (2002) gave an asymptotic justification of DIC in cases where the number of observations is large relative to the number of parameters.

In an discussion to Spiegelhalter et al. (2002), Burnham (2002) questioned "[a] lesson that we learned was that, if sample size n is small or the number of estimated parameters p is large relative to n , a modified AIC should be used, such as $AIC_C = AIC + 2p(p + 1)/(n - p - 1)$. I wonder whether DIC needs such a modification or if it really automatically adjusts for a small sample size or large p , relative to n ." We write this article to answer this question, at least partially for a very important case of linear regression. More concretely,

¹E-mail: k0420214@sk.tsukuba.ac.jp.

²E-mail: kanazawa@sk.tsukuba.ac.jp.

we propose a finite-sample bias corrected information criterion for the Bayesian linear regression models with normally distributed error, and then we implement simulation studies when the sample size is relatively small.

The rest of this article is organized as follows: Next section briefly describes the Bayesian linear regression model. In Section 3, we propose our information criterion for the Bayesian linear regression model. Section 4 shows the results of simulation studies to show the validity of our proposed information criterion when the sample size is small. Finally, Section 5 draws some conclusions concerning our proposed criterion.

2 Bayesian Linear Regression Model

We consider the linear regression model as follows

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N) \quad (2.1)$$

where \mathbf{y} is a $N \times 1$ vector and \mathbf{X} is a $N \times K$ non-stochastic matrix. The parameter vector $\boldsymbol{\beta}$ is a $K \times 1$ vector and error term $\boldsymbol{\varepsilon}$ follows a N -dimensional multivariate normal distribution $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$.

We assume that prior distribution of $\boldsymbol{\beta}$ is a K -dimensional multivariate normal distribution and that of σ^{-2} is a gamma distribution:

$$\boldsymbol{\beta} | \sigma^{-2} \sim \mathcal{N}(\mathbf{b}_0, \sigma^2 \mathbf{B}_0) \quad (2.2)$$

$$\sigma^{-2} \sim \mathcal{G}\left(\frac{\nu_0}{2}, \frac{\lambda_0}{2}\right). \quad (2.3)$$

Then posterior distributions of parameters $\boldsymbol{\beta}$ and σ^{-2} are expressed as

$$\boldsymbol{\beta} | \sigma^{-2}, \mathbf{y}, \mathbf{X} \sim \mathcal{N}(\mathbf{b}_1, \sigma^2 \mathbf{B}_1) \quad (2.4)$$

$$\sigma^{-2} | \mathbf{y}, \mathbf{X} \sim \mathcal{G}\left(\frac{\nu_1}{2}, \frac{\lambda_1}{2}\right) \quad (2.5)$$

where $\mathbf{b}_1 = \mathbf{B}_1 (\mathbf{X}'\mathbf{y} + \mathbf{B}_0^{-1}\mathbf{b}_0)$, $\mathbf{B}_1 = (\mathbf{X}'\mathbf{X} + \mathbf{B}_0^{-1})^{-1}$, $\nu_1 = \nu_0 + N + K$, $\lambda_1 = \lambda_0 + (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_N)'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_N) + (\mathbf{b}_0 - \hat{\boldsymbol{\beta}}_N)'[(\mathbf{X}'\mathbf{X})^{-1} + \mathbf{B}_0]^{-1}(\mathbf{b}_0 - \hat{\boldsymbol{\beta}}_N)$ and $\hat{\boldsymbol{\beta}}_N = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.

3 Finite-Sample Bias Correction and Variable Selection Criterion

Let us denote unknown true density as $f_{\mathbf{Y}}(\cdot)$ and approximating model as $g(\cdot | \boldsymbol{\theta})$ with parameter vector $\boldsymbol{\theta}$. Then Kullback-Leibler Information between $f_{\mathbf{Y}}(\cdot)$ and $g(\cdot | \boldsymbol{\theta})$ can be expressed as follows

$$I(f_{\mathbf{Y}}, g(\cdot | \boldsymbol{\theta})) = \int f_{\mathbf{Y}}(\mathbf{z}) \log \left\{ \frac{f_{\mathbf{Y}}(\mathbf{z})}{g(\mathbf{z} | \boldsymbol{\theta})} \right\} d\mathbf{z} \quad (3.1)$$

and (3.1) can be rewritten as

$$I(f_{\mathbf{Y}}, g(\cdot | \boldsymbol{\theta})) = \mathbf{E}_{\mathbf{z}} [\log \{f_{\mathbf{Y}}(\mathbf{z})\}] - \mathbf{E}_{\mathbf{z}} [\log \{g(\mathbf{z} | \boldsymbol{\theta})\}]. \quad (3.2)$$

Even though the true density $f_{\mathbf{Y}}(\cdot)$ is unknown, the first term on the right-hand side of Kullback-Leibler Information in (3.2) can be regarded as a constant since the variable \mathbf{z} is integrated out. Then we select the best approximating model with maximizing the expected log-likelihood.

In Bayesian perspective, parameters follow the posterior distributions estimated by observed data \mathbf{y} . Hence we consider the posterior mean of (3.2):

$$\mathbf{E}_{\theta|\mathbf{y}} [I(f_{\mathbf{Y}}, g(\cdot|\theta))] = \mathbf{E}_{\mathbf{z}} [\log\{f_{\mathbf{Y}}(\mathbf{z})\}] - \mathbf{E}_{\theta|\mathbf{y}} [\mathbf{E}_{\mathbf{z}} [\log\{g(\mathbf{z}|\theta)\}]] \quad (3.3)$$

and as in Spiegelhalter et al. (2002) and Ando (2007), our proposed information criterion is constructed based on the posterior mean of expected log-likelihood.

From the Bayesian linear regression model (2.1), we use \mathbf{y} as observed data of sample size N obtained from the unknown true density $f_{\mathbf{Y}}(\mathbf{y})$ to estimate the posterior distributions of parameters $\boldsymbol{\beta}$ and σ^{-2} , while we also use \mathbf{z} as replicate data of sample size N generated from the unknown true density $f_{\mathbf{Y}}(\mathbf{z})$ to evaluate the goodness of fit of approximating model $g(\mathbf{z}|\mathbf{X}, \boldsymbol{\beta}, \sigma^{-2})$. Then posterior mean of expected log-likelihood \mathcal{T} is defined as

$$\mathcal{T} \equiv \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} [\mathbf{E}_{\mathbf{z}} [\log\{g(\mathbf{z}|\mathbf{X}, \boldsymbol{\beta}, \sigma^{-2})\}]] \quad (3.4)$$

where expectation $\mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}}[\cdot]$ can be calculated by $\mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}}[\mathbf{E}_{\beta|\sigma^{-2}, \mathbf{y}, \mathbf{X}}[\cdot]]$ from the posterior distributions (2.4) and (2.5).

To estimate the posterior mean of expected log-likelihood \mathcal{T} in (3.4), we use the posterior mean of observed log-likelihood $\widehat{\mathcal{T}}_N$:

$$\widehat{\mathcal{T}}_N \equiv \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} [\log\{g(\mathbf{y}|\mathbf{X}, \boldsymbol{\beta}, \sigma^{-2})\}]. \quad (3.5)$$

and the bias-corrected estimator is obtained as $\widehat{\mathcal{T}}_N - \widehat{b}_N$, where \widehat{b}_N is an estimate of bias $b_{\Theta} \equiv \mathbf{E}_{\mathbf{y}}[\widehat{\mathcal{T}}_N - \mathcal{T}] \neq 0$. Then we propose information criterion (IC) of the form

$$\text{IC} \equiv -2\widehat{\mathcal{T}}_N + 2\widehat{b}_N. \quad (3.6)$$

Ignoring the constant term, we can express the log-likelihood function for the replicate data \mathbf{z} such as

$$\log\{g(\mathbf{z}|\mathbf{X}, \boldsymbol{\beta}, \sigma^{-2})\} = \frac{N}{2} \log \sigma^{-2} - \frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{z} - \mathbf{X}\boldsymbol{\beta}) \quad (3.7)$$

where parameters $\boldsymbol{\beta}$ and σ^{-2} follow the posterior distributions estimated by observed data \mathbf{y} and \mathbf{X} . Then the posterior mean of expected log-likelihood \mathcal{T} in (3.4) is rewritten as

$$\begin{aligned} \mathcal{T} &= \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\mathbf{E}_{\mathbf{z}} \left[\frac{N}{2} \log \sigma^{-2} - \frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{z} - \mathbf{X}\boldsymbol{\beta}) \right] \right] \\ &= \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\frac{N}{2} \log \sigma^{-2} - \frac{\sigma^{-2}}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \\ &\quad - \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\mathbf{E}_{\mathbf{z}} \left[\frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{z} - \mathbf{X}\boldsymbol{\beta}) \right] \right] \\ &\quad + \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \end{aligned}$$

$$= \widehat{\mathcal{T}}_N - C_1 + C_2 \quad (3.8)$$

and the bias b_Θ with respect to the true density $f_{\mathbf{Y}}(\mathbf{y})$ is obtained by $b_\Theta \equiv \mathbf{E}_{\mathbf{y}}[\widehat{\mathcal{T}}_N - \mathcal{T}] = \mathbf{E}_{\mathbf{y}}[C_1 - C_2]$.

First we evaluate C_1 in (3.8). However, true density $f_{\mathbf{Y}}(\mathbf{z})$ is seldom known in practice, so that expectation with respect to the true density is not analytically obtained. In the previous studies, Kitagawa (1997) replaced the unknown true density by the prior predictive density to construct the predictive information criterion (PIC) for the Bayesian linear Gaussian model, while Laud and Ibrahim (1995), Gelfand and Ghosh (1998), and Ibrahim et al. (2001) considered using the posterior predictive density to generate the replicate data \mathbf{z} for model assessment. In this article, we use the posterior predictive density to evaluate the expectation with respect to the true density $f_{\mathbf{Y}}(\mathbf{z})$ in C_1 because the prior predictive density is far more sensitive to the selection of prior distribution.

To evaluate C_1 in (3.8), we replace the true density $f_{\mathbf{Y}}(\cdot)$ with a (conditional) posterior predictive density

$$\mathbf{z}|\sigma^{-2}, \mathbf{y}, \mathbf{X} \sim \mathcal{N}(\mathbf{X}\mathbf{b}_1, \sigma^2\boldsymbol{\Sigma}_0), \quad (3.9)$$

where σ^{-2} follows the posterior distribution (2.5) and $\boldsymbol{\Sigma}_0 = \mathbf{I}_N + \mathbf{X}\mathbf{B}_1\mathbf{X}'$. From (3.9), we estimate C_1 in (3.8) as

$$\begin{aligned} \widehat{C}_1 &= \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\mathbf{E}_{\mathbf{z}|\sigma^{-2}, \mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{z} - \mathbf{X}\boldsymbol{\beta}) \right] \right] \\ &= \mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\mathbf{E}_{\mathbf{z}|\sigma^{-2}, \mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\mathbf{b}_1)' (\mathbf{z} - \mathbf{X}\mathbf{b}_1) \right] \right] \\ &\quad + \mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\text{tr} \left\{ \frac{\sigma^{-2}}{2} (\mathbf{X}'\mathbf{X}) \mathbf{E}_{\beta|\sigma^{-2}, \mathbf{y}, \mathbf{X}} [(\boldsymbol{\beta} - \mathbf{b}_1)(\boldsymbol{\beta} - \mathbf{b}_1)'] \right\} \right]. \end{aligned} \quad (3.10)$$

The first term on the right-hand side of (3.10) can be rewritten as follows

$$\begin{aligned} &\mathbf{E}_{\beta, \sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\mathbf{E}_{\mathbf{z}|\sigma^{-2}, \mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\mathbf{b}_1)' (\mathbf{z} - \mathbf{X}\mathbf{b}_1) \right] \right] \\ &= \mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} \text{tr} \left\{ \mathbf{E}_{\mathbf{z}|\sigma^{-2}, \mathbf{y}, \mathbf{X}} [(\mathbf{z} - \mathbf{X}\mathbf{b}_1)(\mathbf{z} - \mathbf{X}\mathbf{b}_1)'] \right\} \right] \\ &= \mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\frac{\sigma^{-2}}{2} \text{tr} \left\{ \sigma^2 \boldsymbol{\Sigma}_0 \right\} \right] \\ &= \frac{1}{2} \text{tr} \{ \mathbf{I}_N + \mathbf{X}\mathbf{B}_1\mathbf{X}' \} \\ &= \frac{N}{2} + \frac{1}{2} \text{tr} \{ (\mathbf{X}'\mathbf{X})\mathbf{B}_1 \}. \end{aligned} \quad (3.11)$$

The second term on the right-hand side of (3.10) is similarly obtained from the posterior distribution (2.4) as follows

$$\begin{aligned} &\mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\text{tr} \left\{ \frac{\sigma^{-2}}{2} (\mathbf{X}'\mathbf{X}) \mathbf{E}_{\beta|\sigma^{-2}, \mathbf{y}, \mathbf{X}} [(\boldsymbol{\beta} - \mathbf{b}_1)(\boldsymbol{\beta} - \mathbf{b}_1)'] \right\} \right] \\ &= \mathbf{E}_{\sigma^{-2}|\mathbf{y}, \mathbf{X}} \left[\text{tr} \left\{ \frac{\sigma^{-2}}{2} (\mathbf{X}'\mathbf{X}) \sigma^2 \mathbf{B}_1 \right\} \right] \end{aligned}$$

$$= \frac{1}{2} \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\}. \quad (3.12)$$

From (3.11) and (3.12), \widehat{C}_1 in (3.10) is evaluated as follows

$$\begin{aligned} \widehat{C}_1 &= \mathbf{E}_{\beta, \sigma^{-2} | y, X} \left[\mathbf{E}_{z | \sigma^{-2}, y, X} \left[\frac{\sigma^{-2}}{2} (\mathbf{z} - \mathbf{X}\mathbf{b}_1)' (\mathbf{z} - \mathbf{X}\mathbf{b}_1) \right] \right] \\ &\quad + \mathbf{E}_{\sigma^{-2} | y, X} \left[\text{tr} \left\{ \frac{\sigma^{-2}}{2} (\mathbf{X}'\mathbf{X}) \mathbf{E}_{\beta | \sigma^{-2}, y, X} [(\boldsymbol{\beta} - \mathbf{b}_1)(\boldsymbol{\beta} - \mathbf{b}_1)'] \right\} \right] \\ &= \frac{N}{2} + \frac{1}{2} \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\} + \frac{1}{2} \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\} \\ &= \frac{N}{2} + \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\}. \end{aligned} \quad (3.13)$$

Then we notice that \widehat{C}_1 does not depend on any data \mathbf{y} , i.e., $\mathbf{E}_y(\widehat{C}_1) = \widehat{C}_1$.

Suppose that interchange of order of integrations is valid, we can rewrite $\mathbf{E}_y(C_2)$ such as

$$\begin{aligned} \mathbf{E}_y(C_2) &= \mathbf{E}_y \left[\mathbf{E}_{\beta, \sigma^{-2} | y, X} \left[\frac{\sigma^{-2}}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \right] \\ &= \mathbf{E}_{\beta, \sigma^{-2}} \left[\mathbf{E}_{y | X, \beta, \sigma^{-2}} \left[\frac{\sigma^{-2}}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right] \right] \end{aligned} \quad (3.14)$$

where $\mathbf{E}_{\beta, \sigma^{-2}}[\cdot]$ is an expectation with respect to the joint prior distribution and $\mathbf{E}_{y | X, \beta, \sigma^{-2}}[\cdot]$ is an expectation with respect to N -dimensional multivariate normal distribution with mean vector $\mathbf{X}\boldsymbol{\beta}$ and variance-covariance matrix $\sigma^2 \mathbf{I}_N$. Since $\mathbf{E}_{y | X, \beta, \sigma^{-2}}[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'] = \sigma^2 \mathbf{I}_N$, $\mathbf{E}_y(C_2)$ in (3.14) can be evaluated as

$$\begin{aligned} \mathbf{E}_y(C_2) &= \mathbf{E}_{\beta, \sigma^{-2}} \left[\mathbf{E}_{y | X, \beta, \sigma^{-2}} \left[\text{tr} \left\{ \frac{\sigma^{-2}}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \right\} \right] \right] \\ &= \frac{1}{2} \text{tr} \{ \mathbf{I}_N \} \\ &= \frac{N}{2}. \end{aligned} \quad (3.15)$$

Therefore \widehat{b}_N is given by

$$\begin{aligned} \widehat{b}_N &= \mathbf{E}_y \left[\widehat{C}_1 - C_2 \right] \\ &= \widehat{C}_1 - \mathbf{E}_y(C_2) \\ &= \frac{N}{2} + \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\} - \frac{N}{2} \\ &= \text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\}. \end{aligned} \quad (3.16)$$

Then \widehat{b}_N in (3.16) can be regarded as a ratio of variance-covariance matrices $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$ and $\sigma^2\mathbf{B}_1$.

Multiplying -2 to the bias-corrected estimator $\widehat{\mathcal{T}}_N - \widehat{b}_N$, our proposed information criterion for variable selection in the Bayesian linear regression model (IC_{BL}) is obtained by

$$\text{IC}_{BL} = -2\widehat{\mathcal{T}}_N + 2\text{tr} \{(\mathbf{X}'\mathbf{X})\mathbf{B}_1\} \quad (3.17)$$

where $\mathbf{B}_1 = (\mathbf{X}'\mathbf{X} + \mathbf{B}_0^{-1})^{-1}$.

For simplicity, let us denote the parameter \mathbf{B}_0 in (2.2) as $\mathbf{B}_0 = \kappa_0 \mathbf{I}_K$ ($\kappa_0 > 0$) and the bias term \widehat{b}_N can be rewritten as

$$\widehat{b}_N = K - \text{tr}\{(\mathbf{X}'\mathbf{X}\mathbf{B}_0 + \mathbf{I}_K)^{-1}\} \quad (3.18)$$

from the matrix inversion lemma³. Then if the sample size $N \rightarrow \infty$, the last term in (3.18) is expected to be zero because $\text{tr}\{(\kappa_0 \mathbf{X}'\mathbf{X}/N + \mathbf{I}_N/N)^{-1}/N\} \rightarrow 0$ (i.e., $\widehat{b}_N \rightarrow K$) when each element of $\mathbf{X}'\mathbf{X}/N$ does not diverge. Furthermore, if κ_0 is sufficiently large (i.e., non-informative prior), we also have $\text{tr}\{(\kappa_0 \mathbf{X}'\mathbf{X} + \mathbf{I}_K)^{-1}\} \rightarrow 0$ in (3.18).

4 Simulation Experiments

4.1 Deviance Information Criterion (DIC)

We conduct two simulation studies to compare small sample performances of our proposed information criterion (IC_{BL}) in (3.17) with the deviance information criterion (DIC) :

$$\text{DIC} = -2\widehat{\mathcal{T}}_N + p_D, \quad (4.1)$$

where Spiegelhalter et al. (2002) termed p_D as the effective number of parameters defined to be $p_D \equiv 2 \log\{g(\mathbf{y}|\mathbf{X}, \bar{\boldsymbol{\beta}}, \bar{\sigma}^{-2})\} - 2\widehat{\mathcal{T}}_N$ evaluated at the posterior means of parameters $\bar{\boldsymbol{\beta}}$ ($= \mathbf{b}_1$) and $\bar{\sigma}^{-2}$ ($= \nu_1/\lambda_1$) in (2.4) and (2.5).

4.2 Simulation studies

In this article, we present two small sample simulation studies: The first one is designed to examine cases where the set of candidate models includes many over-specified models; the second one examines performances with respect to under-specified candidate models.

Case 1

As in Hurvich and Tsai (1989), we consider the nested candidate models by using seven explanatory variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_7$. In our simulation study, \mathbf{x}_1 is a $N \times 1$ vector whose elements are ones (i.e., intercept term) and the other $N \times 1$ vectors \mathbf{x}_i ($2 \leq i \leq 7$) are generated from the uniform distribution $\mathcal{U}(-2, 2)$. These variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_7$ are included into the candidate models in a sequentially nested fashion. The candidate models are linear regression models given by $\mathbf{y} = \beta_1 \mathbf{x}_1 + \beta_2 \mathbf{x}_2 + \dots + \beta_K \mathbf{x}_K + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$. The candidate model with $K = 1$ has only intercept term and with $K = 7$ is the full model. In this simulation study, we determine the number of variables K by using our proposed information criterion (IC_{BL}) in (3.17) and DIC in (4.1) in small sample cases $N = 25, 50, 100$ with informative ($\kappa_0 = 0.1$) and non-informative ($\kappa_0 = 100$) priors. To

³For any matrices \mathbf{A} ($m \times m$), \mathbf{B} ($m \times n$), \mathbf{C} ($n \times m$), and \mathbf{D} ($n \times n$), we have

$$(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} = \mathbf{A}^{-1} + \mathbf{A}^{-1}\mathbf{B}(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}$$

where \mathbf{A} and \mathbf{D} are nonsingular matrices.

examine the small sample performance in the Bayesian linear regression cases, we generate a sample of \mathbf{y} from the true model ($K = 3$):

$$\mathbf{y} = 1.0\mathbf{x}_1 + 2.0\mathbf{x}_2 + 3.0\mathbf{x}_3 + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, 1.0\mathbf{I}_N). \quad (4.2)$$

In Table 1, we examine the performance of IC_{BL} and DIC for the small sample cases ($N = 25, 50, 100$). The parameters of prior distributions in (2.2) and (2.3) are set to be $\mathbf{b}_0 = \mathbf{0}$, $\mathbf{B}_0 = \kappa_0\mathbf{I}_K$ ($\kappa_0 = 0.1$ or 100), and $\nu_0 = \lambda_0 = 0.1$. The simulation considered each combination of $N = 25, 50, 100$ and $\kappa_0 = 0.1, 100$. We draw 50,000 MCMC samples from the posterior distributions in (2.4) and (2.5) to compute the posterior mean of observed log-likelihood $\widehat{\mathcal{T}}_N$ in (3.5). For each combination of (N, κ_0) , we generate 100 observations of IC_{BL} and DIC, and record the number of selected models (i.e., the candidate model with minimum value for the two criteria).

Table 1 shows that our proposed information criterion (IC_{BL}) identifies the true model ($K = 3$) for the small sample cases ($N = 25, 50, 100$) with informative prior ($\kappa_0 = 0.1$) far better than DIC, on the other hand DIC tends to overfit the true model. For the non-informative prior ($\kappa_0 = 100$), both criteria tend to overfit the true model for the sample size $N = 25$ and 50 , but nevertheless our proposed information criterion (IC_{BL}) far outperforms DIC at the sample size $N = 100$.

In Tables 2 and 3, we show the results of average criteria in 100 observations for the small sample cases ($N = 25, 50, 100$) with informative ($\kappa_0 = 0.1$) and non-informative ($\kappa_0 = 100$) priors. Both criteria selected the true model ($K = 3$) in all cases, but the difference between $2\widehat{b}_N$ and p_D becomes more apparent along with an increase in the number of explanatory variables. Hence the effective number of parameters p_D in DIC tends to underestimate the model complexity compared with the bias term $2\widehat{b}_N$ in IC_{BL} .

Case 2

Next we consider the cases where explanatory variables in the candidate models are selected from the subsets of $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$, (i.e., $\{\mathbf{x}_1\}$, $\{\mathbf{x}_2\}$, $\{\mathbf{x}_3\}$, $\{\mathbf{x}_1, \mathbf{x}_2\}$, $\{\mathbf{x}_1, \mathbf{x}_3\}$, $\{\mathbf{x}_2, \mathbf{x}_3\}$ and $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$). Then we specify a true model such as

$$\mathbf{y} = 1.0\mathbf{x}_1 + 2.0\mathbf{x}_2 + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, 1.0\mathbf{I}_N). \quad (4.3)$$

We implement a simulation study to examine the performance of IC_{BL} in (3.17) and DIC in (4.1) for under-specified candidate models in small sample cases.

The parameters of prior distributions in (2.2) and (2.3) are set to be $\mathbf{b}_0 = \mathbf{0}$, $\mathbf{B}_0 = \kappa_0\mathbf{I}_K$ ($\kappa_0 = 0.1$ or 100), and $\nu_0 = \lambda_0 = 0.1$. We draw 50,000 MCMC samples from the posterior distributions in (2.4) and (2.5) and compute the posterior mean of observed log-likelihood $\widehat{\mathcal{T}}_N$ in (3.5). For each combination of (N, κ_0) , we generate 100 observations of IC_{BL} and DIC and record the number of selected models such as Case 1.

Table 4 shows that DIC tends to select a full model (i.e., $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$) as compared with IC_{BL} in small sample cases. On the other hand, under-specified candidate models (i.e., $\{\mathbf{x}_1\}$, $\{\mathbf{x}_2\}$, $\{\mathbf{x}_3\}$, $\{\mathbf{x}_1, \mathbf{x}_3\}$, $\{\mathbf{x}_2, \mathbf{x}_3\}$) are not selected at all. Hence, DIC shows a tendency to overfit the true model in this simulation study.

In Tables 5 and 6, average values of model complexity in IC_{BL} and DIC with respect to the candidate models $\{\mathbf{x}_1, \mathbf{x}_2\}$, $\{\mathbf{x}_1, \mathbf{x}_3\}$, and $\{\mathbf{x}_2, \mathbf{x}_3\}$ are close to each other. However,

true model (i.e., $\{\boldsymbol{x}_1, \boldsymbol{x}_2\}$) has the largest value of posterior mean of observed log-likelihood $\widehat{\mathcal{T}}_N$ among three candidate models. Therefore, average values of IC_{BL} and DIC successfully select the true model (4.3) in all cases.

5 Conclusion and Discussion

In Bayesian data analysis, DIC (Spiegelhalter et al., 2002) is widely used for the model selection, since this criterion is relatively easy to calculate and applicable to a wide range of statistical models. Spiegelhalter et al. (2002) gave an asymptotic justification of DIC in cases where the number of observations is large relative to the number of parameters. However, Burnham (2002) questioned if the DIC needs a modification for small sample size. In this article, we proposed a variable selection criterion IC_{BL} for the Bayesian linear regression models in small sample cases, as Sugiura (1978) proposed AIC_C in frequentist framework. We then examined the performance of our proposed information criterion (IC_{BL}) relative to the DIC for small sample cases.

In our simulation studies, we found that DIC often showed tendency to overfit the true model (see Tables 1 and 4), whereas our proposed information criterion (IC_{BL}) performs well for small sample cases ($N = 25, 50, 100$). We also found that the measure of model complexity $2\widehat{b}_N$ was mostly larger than p_D (see Tables 2 and 3), leading us to conclude bias correction of DIC underestimates the model complexity in small sample cases.

An important directions for the further research would be to extend our information criterion to the several types of Bayesian regression models (e.g., the hierarchical Bayesian regression models, Bayesian regression models with serially correlated error, and Bayesian Markov-switching models) because empirical analysis is generally performed under the limitation of data availability.

Table 1: The number of selected models by IC_{BL} and DIC for small sample cases $N = 25, 50, 100$ with informative ($\kappa_0 = 0.1$) and non-informative ($\kappa_0 = 100$) priors in Case 1 (100 observations).

Model (K)	Informative prior ($\kappa_0 = 0.1$)						Non-informative prior ($\kappa_0 = 100$)					
	$N = 25$		$N = 50$		$N = 100$		$N = 25$		$N = 50$		$N = 100$	
	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC
1	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0	0
3	97	79	95	74	95	76	81	65	89	75	90	75
4	2	9	5	15	5	14	8	10	5	13	7	12
5	1	6	0	7	0	3	7	9	4	8	3	3
6	0	5	0	2	0	5	0	3	0	1	0	6
7	0	1	0	2	0	2	4	13	2	3	0	4

Table 2: Average values of IC_{BL} and DIC in 100 observations for small sample cases $N = 25, 50, 100$ with informative prior ($\kappa_0 = 0.1$) in Case 1.

Model (K)	Informative prior ($\kappa_0 = 0.1$)														
	$N = 25$				$N = 50$				$N = 100$						
	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{\delta}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{\delta}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{\delta}_N$	p_D
1	99.304 (6.596)	99.558 (6.595)	-48.938 (3.298)	1.429 (0.000)	1.683 (0.008)	196.367 (9.404)	196.519 (9.404)	-97.350 (4.702)	1.667 (0.000)	1.819 (0.009)	394.064 (12.972)	394.147 (12.972)	-196.123 (6.486)	1.818 (0.000)	1.902 (0.009)
2	93.196 (5.943)	92.664 (5.945)	-45.134 (2.974)	2.928 (0.077)	2.396 (0.041)	181.946 (7.360)	181.215 (7.362)	-89.275 (3.682)	3.396 (0.030)	2.665 (0.021)	361.471 (9.739)	360.616 (9.737)	-178.897 (4.868)	3.677 (0.012)	2.822 (0.013)
3	55.425 (2.984)	54.115 (2.984)	-25.502 (1.486)	4.420 (0.112)	3.111 (0.065)	86.241 (5.324)	84.633 (5.326)	-40.565 (2.663)	5.111 (0.043)	3.503 (0.046)	136.230 (10.446)	134.440 (10.448)	-65.349 (5.223)	5.531 (0.019)	3.742 (0.038)
4	56.745 (3.030)	54.676 (3.029)	-25.436 (1.516)	5.874 (0.146)	3.805 (0.075)	87.682 (5.385)	85.198 (5.383)	-40.430 (2.693)	6.823 (0.058)	4.338 (0.050)	137.879 (10.578)	135.149 (10.589)	-65.248 (5.289)	7.382 (0.027)	4.653 (0.042)
5	58.004 (3.088)	55.200 (3.075)	-25.357 (1.536)	7.290 (0.208)	4.486 (0.105)	89.188 (5.474)	85.832 (5.477)	-40.333 (2.743)	8.521 (0.063)	5.165 (0.051)	139.577 (10.477)	135.916 (10.480)	-65.171 (5.239)	9.234 (0.031)	5.573 (0.046)
6	59.228 (3.141)	55.703 (3.131)	-25.272 (1.560)	8.684 (0.234)	5.159 (0.119)	90.776 (5.456)	86.568 (5.465)	-40.283 (2.735)	10.211 (0.084)	6.002 (0.052)	141.136 (10.495)	136.540 (10.492)	-65.026 (5.248)	11.083 (0.034)	6.487 (0.041)
7	60.470 (3.094)	56.221 (3.069)	-25.201 (1.530)	10.068 (0.269)	5.820 (0.139)	92.206 (5.446)	87.137 (5.454)	-40.160 (2.729)	11.886 (0.095)	6.817 (0.057)	142.716 (10.512)	137.185 (10.519)	-64.895 (5.256)	12.925 (0.038)	7.395 (0.046)

Standard deviations in parentheses

Table 3: Average values of IC_{BL} and DIC in 100 observations for small sample cases $N = 25, 50, 100$ with non-informative prior ($\kappa_0 = 100$) in Case 1.

Model (K)	Non-informative prior ($\kappa_0 = 100$)														
	$N = 25$			$N = 50$			$N = 100$								
	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D
1	98.416 (6.588)	98.386 (6.589)	-48.208 (3.294)	1.999 (0.000)	1.970 (0.010)	197.876 (8.148)	197.861 (8.147)	-97.938 (4.074)	2.000 (0.000)	1.984 (0.010)	393.852 (12.035)	393.846 (12.034)	-195.926 (6.017)	2.000 (0.000)	1.993 (0.010)
2	93.019 (5.732)	91.952 (5.732)	-44.510 (2.866)	3.999 (1.65E-04)	2.932 (0.011)	182.333 (7.243)	181.298 (7.243)	-89.167 (3.622)	3.999 (4.67E-05)	2.965 (0.011)	360.553 (10.999)	359.537 (10.998)	-178.277 (5.500)	4.000 (1.49E-05)	2.984 (0.011)
3	31.506 (7.097)	29.409 (7.097)	-12.754 (3.549)	5.998 (2.90E-04)	3.901 (0.012)	56.363 (12.093)	54.310 (12.090)	-25.182 (6.047)	5.999 (7.37E-05)	3.946 (0.013)	108.119 (13.496)	106.092 (13.498)	-51.060 (6.748)	5.999 (2.42E-05)	3.973 (0.012)
4	33.519 (7.281)	30.392 (7.282)	-12.761 (3.641)	7.997 (3.70E-04)	4.869 (0.014)	58.419 (12.074)	55.352 (12.074)	-25.210 (6.037)	7.999 (1.05E-04)	4.932 (0.016)	110.077 (13.381)	107.040 (13.381)	-51.039 (6.691)	7.999 (2.78E-05)	4.962 (0.014)
5	35.083 (7.372)	30.928 (7.373)	-12.544 (3.686)	9.996 (4.98E-04)	5.841 (0.015)	60.275 (12.271)	56.190 (12.270)	-25.138 (6.136)	9.998 (1.35E-04)	5.914 (0.015)	111.959 (13.369)	107.913 (13.370)	-50.980 (6.685)	9.999 (3.29E-05)	5.952 (0.016)
6	37.234 (7.633)	32.048 (7.634)	-12.620 (3.817)	11.995 (6.20E-04)	6.809 (0.016)	62.496 (12.272)	57.395 (12.273)	-25.249 (6.136)	11.998 (1.66E-04)	6.897 (0.020)	113.978 (13.766)	108.924 (13.766)	-50.990 (6.883)	11.999 (3.61E-05)	6.945 (0.016)
7	38.753 (8.174)	32.543 (8.177)	-12.379 (4.087)	13.994 (7.31E-04)	7.784 (0.018)	64.278 (12.289)	58.159 (12.289)	-25.140 (6.144)	13.997 (1.85E-04)	7.878 (0.018)	115.613 (14.019)	109.551 (14.020)	-50.807 (7.010)	13.999 (3.84E-05)	7.937 (0.019)

Standard deviations in parentheses

Table 4: The number of selected models by IC_{BL} and DIC for small sample cases $N = 25, 50, 100$ with informative ($\kappa_0 = 0.1$) and non-informative ($\kappa_0 = 100$) priors in Case 2 (100 observations).

Model	Informative prior ($\kappa_0 = 0.1$)						Non-informative prior ($\kappa_0 = 100$)					
	$N = 25$		$N = 50$		$N = 100$		$N = 25$		$N = 50$		$N = 100$	
	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC	IC_{BL}	DIC
$\{\mathbf{x}_1\}$	0	0	0	0	0	0	0	0	0	0	0	0
$\{\mathbf{x}_2\}$	0	0	0	0	0	0	0	0	0	0	0	0
$\{\mathbf{x}_3\}$	0	0	0	0	0	0	0	0	0	0	0	0
$\{\mathbf{x}_1, \mathbf{x}_2\}$	91	80	93	84	93	85	89	79	91	80	96	90
$\{\mathbf{x}_1, \mathbf{x}_3\}$	0	0	0	0	0	0	0	0	0	0	0	0
$\{\mathbf{x}_2, \mathbf{x}_3\}$	0	0	0	0	0	0	0	0	0	0	0	0
$\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$	9	20	7	16	7	15	11	21	9	20	4	10

Table 5: Average values of IC_{BL} and DIC in 100 observations for small sample cases $N = 25, 50, 100$ with informative prior ($\kappa_0 = 0.1$) in Case 2.

Model	Informative prior ($\kappa_0 = 0.1$)														
	$N = 25$				$N = 50$				$N = 100$						
	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D
$\{\mathbf{x}_1\}$	73.331 (5.347)	73.588 (5.347)	-35.951 (2.674)	1.429 (0.000)	1.685 (0.009)	145.125 (7.837)	145.278 (7.837)	-71.729 (3.919)	1.667 (0.000)	1.819 (0.009)	287.602 (10.764)	287.685 (10.765)	-142.892 (5.382)	1.818 (0.000)	1.901 (0.009)
$\{\mathbf{x}_2\}$	49.125 (4.651)	49.324 (4.649)	-23.794 (2.322)	1.537 (0.068)	1.736 (0.041)	89.938 (8.048)	90.054 (8.049)	-44.100 (4.023)	1.738 (0.029)	1.854 (0.022)	173.414 (12.967)	173.476 (12.968)	-85.777 (6.483)	1.860 (0.012)	1.922 (0.016)
$\{\mathbf{x}_3\}$	76.914 (5.308)	77.122 (5.311)	-37.694 (2.656)	1.525 (0.060)	1.733 (0.031)	152.120 (8.618)	152.239 (8.617)	-75.193 (4.309)	1.733 (0.033)	1.852 (0.019)	302.167 (13.281)	302.229 (13.282)	-150.154 (6.641)	1.859 (0.013)	1.921 (0.010)
$\{\mathbf{x}_1, \mathbf{x}_2\}$	39.725 (4.009)	39.186 (4.009)	-18.390 (2.003)	2.944 (0.078)	2.406 (0.043)	64.713 (7.585)	63.985 (7.583)	-30.658 (3.792)	3.396 (0.033)	2.668 (0.031)	112.735 (11.646)	111.877 (11.646)	-54.530 (5.823)	3.675 (0.013)	2.817 (0.023)
$\{\mathbf{x}_2, \mathbf{x}_3\}$	74.736 (5.423)	74.206 (5.425)	-35.903 (2.713)	2.929 (0.067)	2.399 (0.034)	146.744 (7.816)	146.015 (7.817)	-71.676 (3.909)	3.392 (0.036)	2.663 (0.018)	289.677 (10.954)	288.823 (10.952)	-143.001 (5.476)	3.675 (0.013)	2.821 (0.015)
$\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$	50.140 (4.971)	49.553 (4.974)	-23.550 (2.489)	3.041 (0.088)	2.454 (0.047)	91.506 (7.982)	90.742 (7.988)	-44.022 (3.992)	3.463 (0.047)	2.699 (0.032)	175.133 (13.000)	174.255 (12.999)	-85.708 (6.500)	3.717 (0.018)	2.839 (0.019)
$\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$	40.917 (4.183)	39.608 (4.188)	-18.247 (2.094)	4.423 (0.099)	3.115 (0.056)	66.202 (7.639)	64.591 (7.641)	-30.545 (3.825)	5.113 (0.051)	3.502 (0.039)	114.416 (11.596)	112.622 (11.595)	-54.443 (5.797)	5.529 (0.020)	3.736 (0.026)

Standard deviations in parentheses

Table 6: Average values of IC_{BL} and DIC in 100 observations for small sample cases $N = 25, 50, 100$ with non-informative prior ($\kappa_0 = 100$) in Case 2.

Model	Non-informative prior ($\kappa_0 = 100$)														
	$N = 25$				$N = 50$				$N = 100$						
	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D	IC_{BL}	DIC	$\hat{\tau}_N$	$2\hat{b}_N$	p_D
$\{\mathbf{x}_1\}$	73.185 (6.051)	73.156 (6.052)	-35.593 (3.025)	1.999 (0.000)	1.970 (0.009)	145.721 (8.455)	145.705 (8.454)	-71.861 (4.227)	2.000 (0.000)	1.984 (0.010)	288.977 (11.986)	288.969 (11.986)	-143.489 (5.983)	2.000 (0.000)	1.992 (0.010)
$\{\mathbf{x}_2\}$	45.601 (6.308)	45.571 (6.308)	-21.801 (3.154)	1.999 (1.0.E-04)	1.969 (0.009)	86.823 (9.168)	86.809 (9.168)	-42.412 (4.584)	2.000 (4.2.E-05)	1.985 (0.009)	172.104 (11.558)	172.095 (11.559)	-85.052 (5.779)	2.000 (1.5.E-05)	1.991 (0.009)
$\{\mathbf{x}_3\}$	77.356 (7.055)	77.326 (7.054)	-37.678 (3.527)	1.999 (1.1.E-04)	1.969 (0.010)	153.189 (8.989)	153.174 (8.990)	-75.595 (4.495)	2.000 (4.0.E-05)	1.985 (0.008)	303.755 (11.788)	303.749 (11.788)	-150.878 (5.894)	2.000 (1.4.E-05)	1.993 (0.009)
$\{\mathbf{x}_1, \mathbf{x}_2\}$	29.687 (7.380)	28.622 (7.381)	-12.844 (3.690)	3.999 (1.5.E-04)	2.933 (0.010)	55.125 (9.710)	54.092 (9.711)	-25.563 (4.855)	3.999 (5.3.E-05)	2.966 (0.011)	104.365 (13.881)	103.350 (13.883)	-50.183 (6.941)	4.000 (1.6.E-05)	2.984 (0.009)
$\{\mathbf{x}_1, \mathbf{x}_3\}$	75.212 (6.228)	74.147 (6.230)	-35.607 (3.114)	3.999 (1.3.E-04)	2.934 (0.011)	147.694 (8.535)	146.661 (8.537)	-71.848 (4.267)	3.999 (4.3.E-05)	2.966 (0.009)	291.063 (12.005)	290.048 (12.004)	-143.532 (6.003)	4.000 (1.5.E-05)	2.984 (0.011)
$\{\mathbf{x}_2, \mathbf{x}_3\}$	47.295 (6.349)	46.227 (6.349)	-21.648 (3.174)	3.999 (1.8.E-04)	2.931 (0.011)	88.602 (9.404)	87.569 (9.406)	-42.302 (4.702)	3.999 (5.6.E-05)	2.966 (0.010)	174.231 (11.594)	173.213 (11.596)	-85.115 (5.797)	4.000 (2.3.E-05)	2.982 (0.012)
$\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$	31.448 (7.589)	29.349 (7.590)	-12.725 (3.794)	5.998 (2.2.E-04)	3.898 (0.013)	56.957 (10.021)	54.907 (10.023)	-25.479 (5.010)	5.999 (6.9.E-05)	3.949 (0.014)	106.577 (13.769)	104.549 (13.769)	-50.289 (6.884)	5.999 (2.3.E-05)	3.971 (0.013)

Standard deviations in parentheses

References

- Akaike, H. (1973) 'Information theory as an extension of the maximum likelihood principle.' Second International Symposium on Information Theory Akademiai Kiado, Budapest pp. 267–281
- Ando, T. (2007) 'Bayesian predictive information criterion for the evaluation of hierarchical bayesian and empirical bayes models.' *Biometrika* 94(2), 443–458
- Bedrick, E. J., and C. L. Tsai (1994) 'Model selection for multivariate regression in small samples.' *Biometrics* pp. 226–231
- Burnham, K. P. (2002) 'Discussion of a paper by D.J. Spiegelhalter et al.' *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 64(4), 629
- Gelfand, A. E., and S. K. Ghosh (1998) 'Model choice: A minimum posterior predictive loss approach.' *Biometrika* 85(1), 1–11
- Hurvich, C.M., and C.L. Tsai (1989) 'Regression and time series model selection in small samples.' *Biometrika* 76(2), 297–307
- Ibrahim, J. G., M. H. Chen, and D. Sinha (2001) 'Criterion-based methods for bayesian model assessment.' *Statistica Sinica* 11(2), 419–444
- Kitagawa, G. (1997) 'Information criteria for the predictive evaluation of bayesian models.' *Communications in Statistics-Theory and Methods* 26(9), 2223–2246
- Laud, P. W., and J. G. Ibrahim (1995) 'Predictive model selection.' *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 247–262
- Spiegelhalter, D. J., N. G. Best, B. P. Carlin, and A. van der Linde (2002) 'Bayesian measures of model complexity and fit.' *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 64(4), 583–639
- Sugiura, N. (1978) 'Further analysts of the data by akaike's information criterion and the finite corrections.' *Communications in Statistics-Theory and Methods* 7(1), 13–26