

# Methods and Tools for Bayesian Variable Selection and Model Averaging in Univariate Linear Regression

Anabel Forte,

*Department of Statistics and Operations research, University of Valencia*

Gonzalo García-Donato

*Department of Economics and Finance, University of Castilla-La Mancha*

and Mark F.J. Steel\*

*Department of Statistics, University of Warwick*

## Abstract

In this paper we briefly review the main methodological aspects concerned with the application of the Bayesian approach to model choice and model averaging in the context of variable selection in regression models. This includes prior elicitation, summaries of the posterior distribution and computational strategies. We then examine and compare various publicly available R-packages for its practical implementation summarizing and explaining the differences between packages and giving recommendations for applied users. We find that all packages reviewed lead to very similar results, but there are potentially important differences in flexibility and efficiency of the packages.

## 1 Motivation

A very general problem in statistics is where several statistical models are proposed as plausible descriptions for certain observations  $\mathbf{y}$  and the observed data are used to resolve the model uncertainty. This problem is normally known as *model selection* or *model choice* if the aim is to select a single “best” model, but if the model uncertainty is to be formally reflected in the inferential process, we typically use *model averaging*, where inference on issues that are not model-specific (such as prediction or effects of covariates) is averaged over the set of models under consideration.

A particular important model uncertainty problem in practice is *variable selection* where the proposed models share a common functional form (eg. a normal linear regression model) but differ

---

\*Corresponding author: Mark Steel, Department of Statistics, University of Warwick, Coventry, CV4 7AL, UK; email: M.F.Steel@stats.warwick.ac.uk

in which explanatory variables, from a given set, are included to explain the response. The focus in this paper will be on variable selection in the context of normal linear models, a problem formally introduced in Section 2.

Model uncertainty is a classic problem in statistics that has been scrutinized from many different perspectives. Hence, quite often, the main issues for practitioners are to decide which methodology to use and/or how to implement the methodology in practice. One appealing approach is based on the Bayesian paradigm and is considered by many the *formal* Bayesian answer to the problem. This approach is the one based on the posterior probabilities of the models under consideration and results in a coherent and complete analysis of the problem and provides answers to practical questions. For instance, a single model can be selected as that most supported by the data (the model with the highest posterior probability) or inferences can be performed using the posterior model probabilities as weights, normally denoted by Bayesian model averaging (BMA). In this paper we describe how the formal Bayesian method can be implemented in R (R Core Team, 2015), analyzing the different packages that are currently available in CRAN ([cran.r-project.org](http://cran.r-project.org)). Emphasis is placed on comparison but also on putting in perspective the details of the implementations.

As with any Bayesian method, the prior distribution for the unknown parameters needs to be specified. It is well known that this aspect is particularly critical in model uncertainty problem since results are potentially highly sensitive to the priors used (see e.g. Berger and Pericchi, 2001; Ley and Steel, 2009). In this paper, we pay special attention to the family of priors in the tradition started by Jeffreys, Zellner and Siow (Jeffreys, 1961; Zellner and Siow, 1980; Zellner, 1986) and continued by many other authors with important contributions during the last ten years. These types of priors, which we label *conventional*, are introduced in Section 2.1. Bayarri et al. (2012) have recently shown that conventional priors have a number of optimal properties that make them a very appealing choice for dealing with model uncertainty.

## 2 Bayesian variable selection in Linear Models

Consider a Gaussian response variable  $\mathbf{y}$ , size  $n$ , assumed to be explained by an intercept and some subset of  $p$  possible explanatory variables with values grouped in the  $n \times p$  matrix  $\mathbf{X} = (x_1, \dots, x_p)$ . Throughout the paper we suppose that  $n > p$  and that  $\mathbf{X}$  is of full column rank. We define a binary vector  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)^t$  where  $\gamma_i = 1$  if  $x_i$  is included in the model  $M_{\boldsymbol{\gamma}}$  and zero otherwise. This is the variable selection problem, a model uncertainty problem with the following  $2^p$  competing models:

$$M_{\boldsymbol{\gamma}} : \mathbf{y} = \alpha \mathbf{1}_n + \mathbf{X}_{\boldsymbol{\gamma}} \boldsymbol{\beta}_{\boldsymbol{\gamma}} + \boldsymbol{\varepsilon}, \quad (1)$$

where  $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$  and the  $n \times p_{\boldsymbol{\gamma}}$  design matrices  $\mathbf{X}_{\boldsymbol{\gamma}}$  are all possible submatrices of  $\mathbf{X}$ . If we choose the null matrix for  $\mathbf{X}_{\boldsymbol{\gamma}}$ , corresponding to  $\boldsymbol{\gamma} = \mathbf{0}$ , we obtain the null model with only the

intercept

$$M_0 : \mathbf{y} = \alpha \mathbf{1}_n + \varepsilon. \quad (2)$$

Without loss of generality, we assume that columns of  $\mathbf{X}$  have been centered on their corresponding means, which makes the covariates orthogonal to the intercept, and gives the intercept an interpretation that is common to all models. The set of all competing models is called the model space and is denoted as  $\mathcal{M}$ .

Assuming that one of the models in  $\mathcal{M}$  is the true model, the posterior probability of any model is

$$Pr(M_{\gamma^*} | \mathbf{y}) = \frac{m_{\gamma^*}(\mathbf{y})Pr(M_{\gamma^*})}{\sum_{\gamma} m_{\gamma}(\mathbf{y})Pr(M_{\gamma})}, \quad (3)$$

where  $Pr(M_{\gamma})$  is the prior probability of  $M_{\gamma}$  and  $m_{\gamma}$  is the integrated likelihood with respect to the prior  $\pi_{\gamma}$ :

$$m_{\gamma}(\mathbf{y}) = \int p_{\gamma}(\mathbf{y} | \boldsymbol{\beta}_{\gamma}, \alpha, \sigma) \pi_{\gamma}(\boldsymbol{\beta}_{\gamma}, \alpha, \sigma^2) d\boldsymbol{\beta}_{\gamma} d\alpha d\sigma^2, \quad (4)$$

also called the (prior) marginal likelihood. Note that, for  $\gamma = 0$  this integrated likelihood becomes:

$$m_0(\mathbf{y}) = \int p_0(\mathbf{y} | \alpha, \sigma) \pi_0(\alpha, \sigma^2) d\alpha d\sigma^2, \quad (5)$$

An alternative expression for (3) is based on the Bayes factors:

$$Pr(M_{\gamma^*} | \mathbf{y}) = \frac{B_{\gamma^*}(\mathbf{y})Pr(M_{\gamma^*})}{\sum_{\gamma} B_{\gamma}(\mathbf{y})Pr(M_{\gamma})}, \quad (6)$$

where  $B_{\gamma}$  is the Bayes factor of  $M_{\gamma}$  to a fixed model, say  $M_0$  (without any loss of generality) and hence  $B_{\gamma} = m_{\gamma}/m_0$  and  $B_0 = 1$ .

The prior on the model parameters implicitly assigns posterior point mass at zero for those regression coefficients that are not included in  $M_{\gamma}$ , which automatically induces sparsity.

As stated in the introduction, we are mainly interested in software that implements the formal Bayesian answer which implies that we use the posterior distribution in (3). Even with this important characteristic in common there could be substantial differences between R-packages (leaving aside for the moment details on programming and the interface used) due to the following three aspects:

- the priors that the package accommodates, that is,  $\pi_{\gamma}(\boldsymbol{\beta}_{\gamma}, \alpha, \sigma^2)$  and  $Pr(M_{\gamma})$ ,
- the tools provided to summarize the posterior distribution and obtain model averaged inference,
- the numerical methods implemented to compute the posterior distribution.

We now succinctly revise the main methodological proposals for the above points. Emphasis is placed in presenting the revision in a way that accommodates the methods implemented in the different R packages.

| Proposal   | Reference                                      | Name                                 | Label |
|--|--|--------------------------------------|-------|
| <b>Constant <math>g</math></b>                                       |  |                                      |       |
| $g = n$  | Zellner (1986); Kass and Wasserman (1995)      | Unit Information prior (UIP)         | C1    |
| $g = p^2$  | Foster and George (1994)                       | Risk inflation criterion prior (RIC) | C2    |
| $g = \max\{n, p^2\}$   | Fernández et al. (2001)                        | Benchmark prior (BRIC)               | C3    |
| $g = \log(n)$  | Fernández et al. (2001)                        | Hannan-Quinn (HQ)                    | C4    |
| $g_\gamma = \hat{g}_\gamma$  | Liang et al. (2008)                            | Local Empirical Bayes (EBL)          | C5    |
| <b>Random <math>g</math></b>   |  |                                      |       |
| $g \sim IGa(1/2, n/2)$   | Jeffreys (1961); Zellner and Siow (1980, 1984) | Cauchy prior (JZS)                   | R1    |
| $g a \sim \pi(g) \propto (1+g)^{-a/2}$                               | Liang et al. (2008)                            | hyper-g                              | R2    |
| $g a \sim \pi(g) \propto (1+g/n)^{-a/2}$                             | Liang et al. (2008)                            | hyper-g/n                            | R3    |
| $g \sim \pi(g) \propto (1+g)^{-3/2}, g > \frac{1+n}{p_\gamma+1} - 1$ | Bayarri et al. (2012)                          | Robust prior                         | R4    |

Table 1: Specific proposals for the hyperparameter  $g$  in the literature. Column “Label” will be used as convenient reference to particular proposals throughout the paper. For the priors on  $g$ ,  $a > 2$  to ensure a proper prior and  $p_\gamma$  denotes the number of covariates in  $M_\gamma$ .

## 2.1 Prior Specification

The two inputs that are needed to obtain the posterior distribution are  $\pi_\gamma$  and  $Pr(M_\gamma)$ : the  $2^p$  prior distributions for the parameters within each model and the prior distribution over the model space, respectively.

Without loss of generality, the prior distributions  $\pi_\gamma$  can be expressed as

$$\pi_\gamma(\boldsymbol{\beta}_\gamma, \alpha, \sigma^2) = \pi_\gamma(\boldsymbol{\beta}_\gamma | \alpha, \sigma^2) \pi_\gamma(\alpha, \sigma^2).$$

Under the conventional approach (Fernández et al., 2001) the standard Jeffreys’ prior is used for the parameters that are common to all models

$$\pi_\gamma(\alpha, \sigma^2) = \sigma^{-2} \tag{7}$$

and for  $\pi_\gamma(\boldsymbol{\beta}_\gamma | \alpha, \sigma^2)$  we adopt either a normal or mixtures of normal distributions centered on zero (“by reasons of similarity” Jeffreys, 1961) and scaled by  $\sigma^2(\mathbf{X}_\gamma^t \mathbf{X}_\gamma)^{-1}$  (“a matrix suggested by the form of the information matrix” Zellner and Siow, 1980) times a factor  $g$ , normally labelled as “ $g$ -prior”. Recent research has shown that such conventional priors possess a number of optimal properties that can be extended by putting specific priors on the hyperparameter  $g$ . Among these properties are invariance under affine transformations of the covariates, several types of predictive matching and consistency (for details see Bayarri et al., 2012).

The specification of  $g$  has inspired many interesting studies in the literature. Of these, we have collected the most popular ones in Table 1.

Related with the conventional priors is the proposal by Raftery (1995) which is inspired by asymptotically reproducing the popular Bayesian Information Criterion (Schwarz, 1978). Raftery

(1995) proposes using the same covariance matrix as the Unit Information Prior (see Table 1) but with mean the maximum likelihood estimator  $\hat{\beta}_\gamma$  (instead of the zero mean of the conventional prior).

Other priors specifically used in model uncertainty problems are the spike and slab priors, that assume that the components of  $\beta$  are independent, each having a mixture of two distributions: one highly concentrated on zero (the spike) and the other one quite disperse (the slab). There are two different developments of this idea in the literature. In the original proposal by Mitchell and Beauchamp (1988) the spike is a degenerate distribution at zero so this fits with what we have called the formal approach. The proposal by George and McCulloch (1993) in which the spike is a continuous distribution with a small variance also received a lot of attention, perhaps for computational advantages. In this implementation there is no posterior distribution over the model space as every model smaller than the full model has zero probability.

With respect to the priors over the model space  $\mathcal{M}$ , a very popular starting point is

$$Pr(M_\gamma | \theta) = \theta^{p_\gamma} (1 - \theta)^{p - p_\gamma}, \quad (8)$$

where  $p_\gamma$  is the number of covariates in  $M_\gamma$ , and the hyperparameter  $\theta \in (0, 1)$  has the interpretation of the common probability that a given variable is included (independently of all others).

Among the most popular default choices for  $\theta$  are

- Fixed  $\theta = 1/2$ , which assigns equal prior probability to each model, i.e  $Pr(M_\gamma) = 1/2^p$ ;
- Random  $\theta \sim \text{Unif}(0, 1)$ , giving equal probability to each possible number of covariates or model size.

Of course many other choices for  $\theta$  – both fixed and random– have been considered in the literature. In general, fixed values of  $\theta$  have been shown to perform poorly in controlling for multiplicity (the occurrence of spurious explanatory variables as a consequence of performing a large number of tests) and can lead to rather informative priors. This issue can be avoided by using random distributions for  $\theta$  as, for instance, the second proposal above that has been studied in Scott and Berger (2010). Additionally, Ley and Steel (2009) consider the use of  $\theta \sim \text{Beta}(1, b)$  that results in a binomial-beta prior for the number of covariates in the model or the model size,  $W$ :

$$Pr(W = w | b) \propto \binom{p}{w} \Gamma(1 + w) \Gamma(b + p - w), \quad w = 0, 1, \dots, p.$$

Notice that for  $b = 1$  this reduces to the uniform prior on  $\theta$  and also on  $W$ . As Ley and Steel (2009) highlight, this setting is useful to incorporate prior information about the mean model size, say  $w^*$ . This would translate into  $b = (p - w^*)/w^*$ .

## 2.2 Summaries of the posterior distribution and model averaged inference

The simplest summary of the posterior model distribution (3) is its mode

$$\arg \max_{\gamma} Pr(M_{\gamma} | \mathbf{y}).$$

This model is the model most supported by the information (data and prior) and is normally called the HPM (stands for highest posterior model) or MAP (maximum a posteriori) model. Clearly, a measure of uncertainty regarding this summary is reflected by its posterior probability which should always be reported.

When  $p$  is moderate to large, posterior probabilities of individual models can be very small and their interpretation loses appeal. In such situations, posterior inclusion probabilities (normally denoted as PIP) are very useful.

$$Pr(\gamma_i = 1 | \mathbf{y}) = \sum_{x_i \in M_{\gamma}} Pr(M_{\gamma} | \mathbf{y}). \quad (9)$$

These should be understood as the importance of each variable for explaining the response. Interestingly, these probabilities are used to define another summary, namely the median probability model (MPM) which is the model containing the covariates with inclusion probability larger than 0.5. This model is studied in (Barbieri and Berger, 2004) and they show that, in some situations, it is optimal for prediction.

Extending the idea of inclusion probabilities, it is interesting to obtain measures of joint importance of sets of regressors on the response. For instance, we can compute the posterior probability of two (or more) covariates occurring together in the model or the probability that a covariate enters the model given that another covariate is already present (or not). These quantities are known as joint posterior probabilities and conditional posterior probabilities, respectively, and are studied, with other related summaries, in Ley and Steel (2007) (and references therein).

A measure of the model complexity is given by

$$Pr(W = w | \mathbf{y}) = \sum_{M_{\gamma}: p_{\gamma}=w} Pr(M_{\gamma} | \mathbf{y}), \quad (10)$$

which is the posterior probability mass function of the model size.

The posterior distribution easily allows for obtaining model averaged estimates of any quantity of interest  $\Lambda$  (assuming it has the same meaning across all models). Suppose  $\hat{\Lambda}_{\gamma}$  is the estimate of  $\Lambda$  you would use if  $M_{\gamma}$  were the true model. Then the model averaged estimate of  $\Lambda$  is

$$\hat{\Lambda} = \sum_{M_{\gamma}} \hat{\Lambda}_{\gamma} Pr(M_{\gamma} | \mathbf{y}), \quad (11)$$

which has the appeal of incorporating model uncertainty.

When  $\Lambda$  refers to regression coefficients ( $\beta_i$ ) the model averaged estimates should be used and interpreted with caution as they could be potentially misleading since the ‘same’ parameter may have a different meaning in different models (Berger and Pericchi, 2001). Also the posterior distribution of  $\beta_i$  is a discrete mixture and hence summaries like the mean are not natural descriptions.

One particular appealing application of this technique is in predicting new values  $y^*$  of the dependent variable associated with certain values of the covariates. In this case  $\Lambda$  could be the moments of  $y^*$  or even the whole predictive distribution. Apart from their intrinsic interest, predictions can be a very useful tool to run predictive checks (often using score functions) *e.g.* to compare various prior specifications.

### 2.3 Numerical methods

There are two main computational challenges in solving a model uncertainty problem. First is the integral in (4) and second is the sum in the denominator of (3) which involves many terms if  $p$  is moderate or large.

Fortunately, in normal models, conventional priors combine easily with the likelihood, and conditionally on  $g$  lead to closed forms for  $m_\gamma(\mathbf{y})$ . Hence, at most, a univariate integral needs to be computed when  $g$  is taken to be random. Interesting exceptions are the Robust prior of Bayarri et al. (2012) and the prior of Maruyama and George (2011), which despite assuming a hyper prior on  $g$  induce closed form marginals. This is done by making the prior on  $g$  dependent on the size of the model considered, either through the hyperparameters or through truncation.

The second problem, related with the magnitude of the number of models in  $\mathcal{M}$  (i.e.  $2^p$ ), could be a much more difficult one. If  $p$  is small (say,  $p$  in the twenties at most) exhaustive enumeration is possible but if  $p$  is larger, heuristic methods need to be implemented. The question of which method should be used has been studied in Garcia-Donato and Martinez-Beneito (2013) which classify strategies as i) MCMC methods to sample from the posterior (3) in combination with estimates based on frequencies and ii) searching methods looking for ‘good’ models with estimates based on renormalization (i.e with weights defined by the analytic expression of posterior probabilities, cf. (3)). They show that i) is potentially more precise than ii) which could be biased by the searching procedure. Approach i) is the most standard approach but different implementations of ii) have lead to fruitful contributions. The proposals in Raftery et al. (1997) and Fernández et al. (2001) which are based on a Metropolis-Hasting algorithm called  $MC^3$  (originally introduced in Madigan and York (1995)) could be in either class above, while the implementation in Eicher et al. (2011) based on a leaps and bound algorithm proposed by Raftery (1995) is necessarily in (ii), since model visit frequencies are not an approximation to model probabilities in this case.

### 3 CRAN packages screening

In what follows we will write the name of the packages using the font `package`; functions as `function()` and arguments as `argument`.

We seek in CRAN all possible packages that, potentially, could be used to implement the Bayesian approach to variable selection. The key words used to search in CRAN were *Model Selection*, *Variable Selection*, *Bayes Factor* and *Averaging*. The last search was on June 26, 2015 and we found a total of 13 packages: `VarSelectIP`; `spikeslab` (Ishwaran et al., 2013); `spikeslabGAM` (Morey et al., 2015); `ensembleBMA` (Fraley et al., 2015); `dma` (McCormick et al., 2014); `BMA` (Raftery et al., 2015); `mglm` (Katabuchi and Nakamura, 2015); `varbvs` (Carbonetto and Stephens, 2012); `INLABMA` (Bivand et al., 2015); `BayesFactor` (Morey et al., 2015); `BayesVarSel` (Garcia-Donato and Forte, 2015); `BMS` (Zeugner and Feldkircher, 2015) and `mombf` (Rossell et al., 2014).

From these, `VarSelectIP`, appeared as not longer supported and, within the rest, only the last four implement conventional priors described in the previous section to perform variable selection in linear models and hence will be considered for detailed description and comparison in the following sections. Particularly, `BayesVarSel` and `BMS` seem to be specifically conceived for that task, while the main motivation in `BayesFactor` and `mombf` seem different. `BayesFactor` provides many interesting functionalities to carry out  $t$ -tests, ANOVA-type studies and contingency tables using (conventional) Bayes factors with special emphasis on the specification of the hyper parameter  $g$  for certain design matrices. On the other hand, `mombf` focuses on a particular type of priors for the model parameters, namely the non-local priors (Johnson and Rossell, 2010, 2012), applied to either the normal scenario considered here or probit models.

Of the other packages we found, `spikeslab` and `spikeslabGAM`, implement spike and slab priors in the spirit of the approach by George and McCulloch (1993) and hence are not directly comparable with packages that compute the posterior distribution over the model space. Interestingly, the original spike and slab approach by Mitchell and Beauchamp (1988) is used as the base methodology in `varbvs` but with a specific development by Carbonetto and Stephens (2012) with extreme high dimensional problems ( $p \gg$ ) in mind. Finally, `BMA` provides the posterior distribution over the model space, but based on the BIC criterion.

Some other packages consider statistical models that are not of the type studied here (linear regression models). This is the case for `ensembleBMA`, which implements BMA for weather forecasting models and `dma` which focuses on dynamic models.

`INLABMA` interacts with the INLA (Rue et al., 2009) methodology for performing model selection within a given list of models. The priors there used are those in the R package INLA which are not model selection priors.

The package `mglm` is not Bayesian and it uses the Akaike Information Criterion (AIC).



| Package                        | BayesFactor                 | BayesVarSel  | BMS                | mombf   |
|--------------------------------|-----------------------------|--|--------------------|---|
| Commands for model uncertainty | <code>regressionBF()</code> | <code>Bvs()</code> , <code>PBvs()</code> , <code>GibbsBvs</code> | <code>bms()</code> | <code>modelSelection()</code>                       |
| Argument                       | <code>rscaleCont=</code>    | <code>prior.betas=</code>  | <code>g=</code>    | <code>priorCoef=</code>                             |
| Prior                          |                             |  |                    |   |
| C1                             | -                           | "gZellner"   | "UIP"              | <code>zellnerprior(tau=n)</code>                    |
| C2                             | -                           | -  | "RIC"              | <code>zellnerprior(tau=p<sup>2</sup>)</code>        |
| C3                             | -                           | "FLS"  | "BRIC"             | <code>zellnerprior(tau=max(n,p<sup>2</sup>))</code> |
| C4                             | -                           | -  | "HQ"               | <code>zellnerprior(tau=log(n))</code>               |
| C5                             | -                           | -  | "EBL"              | -   |
| R1                             | 1                           | "ZellnerSiow"  | -                  | -   |
| R2                             | -                           | -  | "hyper=a"          | -   |
| R3 (a=3)                       | -                           | "Liangetal"  | -                  | -   |
| R4                             | -                           | "Robust"   | -                  | -   |

Table 2: Priors for the parameters within each model. Main commands and corresponding modifying arguments for the different specifications for the hyper parameter  $g$  (keys in column ‘Prior’ refer to that in Table 1) in conventional prior.

## 4 Selected Packages

The R packages `BayesFactor`, `BayesVarSel`, `BMS` and `mombf` provide functionalities to calculate and study the posterior distribution (3) corresponding to some of the conventional priors described in Table 1. The commands for such calculation are `regressionBF()` in `BayesFactor`; `Bvs()`, `PBvs()` and `GibbsBvs()` (for exhaustive enumeration, distributed enumeration and Gibbs sampling) in `BayesVarSel`; `bms()` in `BMS` and finally `modelSelection()` in the package `mombf`.

**Prior inputs** The different conventional priors available in each package and the corresponding argument for its use are described in Table 2.

The implementation of the conventional priors in `mombf` have certain peculiarities that we now describe. The priors for the common parameters,  $(\alpha, \sigma)$ , in `mombf` do not exactly coincide with (7). In this package, the simplest model  $M_0$  only contains the error term and hence  $\alpha$  is not a common parameter. The more popular problem with fixed intercept examined in this paper (cf. (7)) is handled via the modifying argument `center=TRUE` (given by default) which in turns is equivalent to a prior for  $\alpha$  degenerate at its maximum likelihood estimate. This will, especially if  $n$  is large enough, often lead to very similar results as with a flat prior on  $\alpha$  but small differences could occur because in `mombf` the variability in this parameter is not taken into account. Also, for  $\sigma^2$  this package uses an inverse gamma which has the non informative  $\sigma^{-2}$  as a limiting density. Thus, differences among the two are expected to be negligible if the parameters in the inverse gamma are small (values of 0.01 are given by default). Another logical argument in `modelSelection()` is

| Package                        | BayesFactor                          | BayesVarSel                | BMS                  | mombf                          |
|--------------------------------|--------------------------------------|----------------------------|----------------------|--------------------------------|
| Argument                       | <code>newPriorOdds(BFobject)=</code> | <code>prior.models=</code> | <code>mprior=</code> | <code>priorDelta=</code>       |
| Prior                          |                                      |                            |                      |                                |
| $\theta = 1/2$                 | <code>rep(1,2~ p)</code>             | "constant"                 | "fixed" or "uniform" | <code>modelunifprior()</code>  |
| $\theta \sim \text{Unif}(0,1)$ | -                                    | "ScottBerger"              | "random"             | <code>modelbbprior(1,1)</code> |

Table 3: Most popular default priors over the model space (see (8)) within the selected packages. For more flexible options see the text.

`scale`. If it is set to `TRUE` the  $y$ 's and the  $x$ 's are scaled to have unitary variance. In this article we are fixing it to `scale=FALSE` so that the data that enter in all the main functions exactly coincide.

All four packages are very rich and flexible regarding the choice of the prior over the model space,  $Pr(M_\gamma)$ . The access to the standard approaches is described in Table 3. Apart from these standard priors BMS, following the proposals in Ley and Steel (2009), also allows for the use of a beta distribution for  $\theta$  in (8) by using `mprior="random"` and modifying the argument `mprior.size` to specify the desired expectation for the model prior distribution (the default option is  $p/2$  hence providing the uniform prior on model size). Similarly the `mombf` package provides a beta prior for  $\theta$  with parameter  $(a, b)$  by setting the corresponding argument to `modelbbprior(a,b)`. In `BayesVarSel` particular specifications of prior probabilities are available with `mprior="User"` and a  $p + 1$  dimensional vector defined in `priorprobs` which describes the prior probability,  $Pr(M_\gamma)$ , of a single model of each possible size (models of the same size are assumed to have the same prior probability).

For illustration purposes consider the FLS dataset in Ley and Steel (2009) with  $p = 41$  potential regressors. These authors study the prior (8) with  $\theta \sim \text{Beta}(1, b = (41 - \omega^*)/\omega^*)$  and  $\omega^* = 7$ , reflecting that, a priori, the expected number of regressors is  $\omega^* = 7$ . Such a prior can be implemented in BMS with `mprior="random"`, `mprior.size=7` and in `mombf` with `modelbbprior(1,34/7)`. In `BayesVarSel` the syntax is quite different and we have to specify `prior.models="User"` and `priorprobs = dbetabinom.ab(x = 0 : 41, size = 41, shape1 = 1, shape2 = 34/7)/choose(41, 0 : 41)`.

**Summaries and model averaging** The result of executing the main commands for model uncertainty (see Table 2) is an object describing, with a specific structure depending on the package, the posterior distribution (3). For ease of exposition suppose the object created is called `ob`. We compare here the different possibilities to summarize this distribution under each package. This is illustrated in the Supplementary Material which shows the different ways of summarizing the results for each package using one of the studied data sets.

- In `BayesFactor`, a list of the most probable models and their corresponding Bayes factors (to the null model) can be obtained with the command `head(ob)` or `plot(ob)` over the resulting

object.

- In `mombf`, this list can be obtained with `postprob(ob)` but now best models are displayed with their posterior probabilities. Additionally, inclusion probabilities (9) are contained in `ob$margpp`. In the context of large model spaces, having a list with all the models sampled can be very useful so that the user may program his/her own needs, such as model averaged predictions. Such a list is contained in binary matrix form in `mombf` in `ob$postSample`. To obtain model averaged estimates we also have the command `rnlp` which produces posterior samples of regression coefficients (from which it is easy to obtain any  $\hat{\Lambda}$  in (11) that relates to coefficients).
- In `BayesVarSel` most probable models and their probabilities are viewed printing the object created, `ob`, while `summary(ob)` displays a table with the inclusion probabilities, the HPM and the MPM (see Subsection 2.2). The posterior distribution of the model size (10) is in `ob$postprobdim` which can be graphed with `plotBvs(ob,option="d")`. Plots of several measures of the joint importance of two covariates (e.g. joint inclusion probabilities) can be visualized as an image plot with `plotBvs(ob, option="j")`. All models visited are saved in the matrix `ob$modelslogBF` which, in the last column, have the Bayes factors of each model in log scale.
- In `BMS` the top best models with their probabilities are displayed using `topmodels(ob)`, that can also be plotted with `image(ob)`. A `summary(ob)` of the resulting object also prints the posterior of the model size (10) that can be plotted with the command `plotModelSize(ob)`. Printing `ob` displays a table with model averaged estimates of regression coefficients, namely their expected posterior mean and standard deviation (column Post Mean and Post SD respectively). Interestingly, it is possible to compute predictions with the commands `predict(ob)` (expected predictive mean a posteriori) and `pred.density(ob)` (mixture predictive density based on a selected number of best models). This package does not save all the models visited but only a (necessarily small) list of the best models sampled in `ob$topmod` expressed in hexadecimal code.

**Numerical methods** Exhaustive enumeration can be performed with `BayesFactor`, `BayesVarSel` (command `Bvs()`) and in `BMS` (modifying argument `mcmc="enumerate"`).

When  $p$  is larger, exhaustive enumeration is not feasible and this is solved in `mombf`, `BayesVarSel` and `BMS` by providing specific routines to approximate the posterior distribution in such big model spaces. In summary, all three packages implement the strategy i) briefly described in Section 2.3 with the following peculiarities. The packages `mombf` and `BayesVarSel` implement the same Gibbs sampling scheme. A minor difference between both is that frequency-based estimates of inclusion

probabilities in `mombf` are refined using Rao-Blackwellization. The methods programmed in `BMS` are also MCMC strategies to explore the posterior distribution which can be of the type birth and death (modifying argument `mcmc="bd"`) or a reversible jump (`mcmc="rev.jump"`). There is an important difference between the algorithms in `mombf`, `BayesVarSel` and in `BMS`. While in each MCMC step the inclusion/exclusion of *all*  $p$  covariates is sampled in `mombf` and `BayesVarSel` only one is sampled in `BMS`.

## 5 Performance in selected datasets

To compare the selected packages two different scenarios have been considered:

- Exact scenario: data sets with small  $p$  and hence all the models can be enumerated.
- Sampling scenario: data sets with moderate to large  $p$  where only a small proportion of models can be explored.

As we previously mentioned, `mombf` cannot be considered in the exact scenario nor can `BayesFactor` be considered in the sampling scenario. Ideally, we should compare all possible packages (in each setup) under the same prior. Thus, Table 2 indicates which comparisons are possible. We compared `BayesFactor` with `BayesVarSel` using the Zellner-Siow prior (labelled as R1) while we compared `mombf` and `BMS` and `BayesVarSel` using the UIP (C1). In all cases, the constant prior over the model space was used.

As expected, all four packages produced very similar results in the analyzed datasets. Hence, the question of comparing them reduces basically to comparing computational times and the availability, clarity and organisation of the output.

For the computational comparisons to be fair all the calculations have been done on an iMac computer with Intel Core i5, 2.7 GHz processor. The code used to compute results provided here is publicly available at [www.uv.es/fordela](http://www.uv.es/fordela).

### 5.1 Exact Scenario

We considered two data sets that we briefly describe.

**US Crime Data.** The US Crime data set was first studied by Ehrlich (1973) and is available from R-package `MASS` (Venables and Ripley, 2002). This data set has a total of  $n = 47$  observations (corresponding to states in the US) of  $p = 15$  potential covariates aimed to explain the rate of crimes in a particular category per head of population.

**Returns to schooling** This data set, used by Tobias and Li (2004) and Ley and Steel (2012), concerns returns to education. As these are microeconomic data, the number of potential observations is much larger. In particular we have a response variable: the log of hourly wages recorded for  $n = 1190$  white males in the US in 1990, and a total of  $p = 26$  possible regressors.

For both scenarios we directly compare the time needed to exactly calculate the posterior distribution with `BayesVarSel` and `BMS` using the C1 prior for the parameters and the uniform prior (with fixed  $\theta = 1/2$ ) for the model space. These times are presented in Table 4. The results clearly indicate that `BayesVarSel` is more affected by the sample size since it performs better than `BMS` for the Crime data set ( $n = 47$ ) but not for the returns to schooling application ( $n = 1190$ ). The Bayes factors depend on the data only through the sum of squared errors and we know that `BayesVarSel` computes this statistic from scratch for each model and, thus, the  $n$  matters in that calculation. Hence a likely reason for the differences in computational time between the two packages when  $n$  increases would be that the algorithm in `BMS` exploits reduction by sufficiency and optimally updates when a variable is added/dropped from the current model.

The comparison between `BayesFactor` and `BayesVarSel`, now using the R1 prior, is summarized in the same table for the Crime data set where we can clearly see that `BayesFactor` is outperformed by `BayesVarSel`.

| Data set                      | Prior   | BMS       | BayesVarSel | BayesFactor |
|-------------------------------|---------|-----------|-------------|-------------|
| Crime $p = 15$                | C1 unif | 3.22 secs | 0.35 secs   | -           |
| Returns to schooling $p = 26$ | C1 unif | 1.83 hrs  | 11.24 hrs   | -           |
| Crime $p = 15$                | R1 unif | -         | 1.4 secs    | 12.73 mins  |

Table 4: Computational times in exact scenario.

Table 4 also illustrates the large difference in computational cost between an exhaustive analysis with  $p = 15$  covariates (where  $\mathcal{M}$  has  $2^{15} = 32,768$  models) and  $p = 26$ , leading a model space with 67 million models, which is about 2000 times larger. Computational cost goes up by a factor of about 2000 for `BMS`, which is therefore roughly linear in the size of model space, and thus seems virtually unaffected by the number of observations  $n$ . This is a consequence of how the statistics are computed within each package, as commented above.

## 5.2 Sampling Scenario

We considered here three data sets.

**Ozone.** These data were used by Casella and Moreno (2006), Berger and Molina (2005) and Garcia-Donato and Martinez-Beneito (2013) and contain  $n = 178$  measures of ozone concentration

in the atmosphere with a total of  $p = 35$  covariates. Details on the data can be found in Casella and Moreno (2006).

**GDP growth.** This dataset is larger than Ozone with a total of  $p = 67$  potential drivers for the annual GDP growth per capita between 1960 and 1996 for  $n = 88$  countries. This data set is also used in Sala-I-Martin et al. (2004) and revisited by Ley and Steel (2007).

**Boston housing.** This dataset was used recently in Schäfer and Chopin (2013) and contains  $n = 506$  observations of  $p = 103$  covariates formed by the 13 columns of the original data set, all first order interactions and a squared version of each covariate (except for the binary variable CHAS).

For the Ozone dataset, exact inclusion probabilities, (9), are reported in Garcia-Donato and Martinez-Beneito (2013) for the C1 prior. These are the result of an intensive computational experiment aimed at comparing different searching methods. These numbers allow us to define a simple measure to compare the computational efficiency of the different packages. For a given computational time,  $t$ , we can compute

$$\Delta_t = \max_{i=1, \dots, p} |\widehat{Pr}_t(\gamma_i = 1 | \mathbf{y}) - Pr(\gamma_i = 1 | \mathbf{y})|,$$

where  $\widehat{Pr}_t(\gamma_i = 1 | \mathbf{y})$  is the estimate of the corresponding PIP at time  $t$  provided by the package. Clearly, the faster  $\Delta_t$  approaches zero, the more efficient is the package. In Figure 1 we have plotted  $\Delta_t$  for `mombf` and `BayesVarSel` respectively and the two algorithms in `BMS`.

All four approaches behave quite satisfactorily, providing very reliable estimates with a small computational time (a maximum discrepancy with the exact values of 0.01 in less than 2.5 minutes). It seems that `BayesVarSel` is slightly more efficient than the rest while the reversible jump implemented in `BMS` is less efficient. The apparent constant bias in `mombf` is possibly due to the difference in the prior actually implemented that we have already described.

In the GDP growth and the Boston Housing examples, we cannot compute  $\Delta_t$  simply because the PIP's are unknown. Nevertheless, we observe that for a sufficiently large computational time, all packages converged to almost identical PIP's. Hence, and even in the unlikely case that none of them were capturing the 'truth' it seems that the fairest way to compare the packages is computing time until 'convergence'. This is what we have represented in Figures 2 and 3 where the  $y$ -axes display the difference between estimates at consecutive computational times, *i.e.*

$$\Delta_{t,t-dt} = \max_{i=1, \dots, p} |\widehat{Pr}_t(\gamma_i = 1 | \mathbf{y}) - \widehat{Pr}_{t-dt}(\gamma_i = 1 | \mathbf{y})|,$$

where  $dt = 60$  seconds was used and we have verified that PIPs converge.

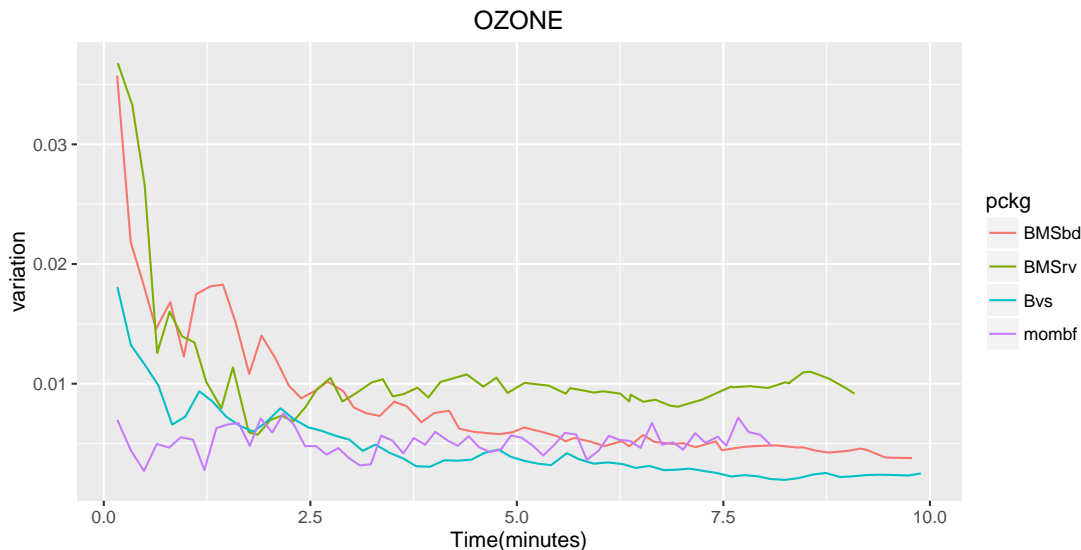


Figure 1: Ozone dataset: maximum difference with the real inclusion probabilities ( $\Delta_t$ ) as a function of computational time. BMSbd (BMSrv) stands for the birth/death (reversible jump) algorithm in BMS; Bvs for BayesVarSel and momfb for momfb.

In the GDP growth data set, we can not find big differences in the performance of all four approaches and all of them behave, again, very satisfactorily. It seems that the procedure implemented by BayesVarSel tends to 0 faster than the rest of algorithms while the performance of momfb manifests more variability, likely due to the Rao-Blackwellization way of computing the results.

In the Boston Housing problem the package BayesVarSel is clearly penalized (with respect to the GDP growth data) by the large number of observations hence having a slower convergence.

For both examples, the inclusion probabilities obtained with each package after 30 minutes of computations (after the burning period) differ, at most, in the second decimal number. Finally, both plots are not affected by the difference in the prior implemented (as each method compares with itself).

## 6 Other Features

Besides the characteristics analysed so far (prior inputs, numerical methods and summaries), there are several other features of the packages that are potentially relevant for the applied user. We list some here under three categories: the interface, extra functionalities and documentation.

**The interface** In general all four packages have simple interfaces with quite intuitive syntaxes. One minor difference is that in BayesVarSel and BayesFactor the dependent and explanatory variables are defined with the use of `formula` (hence inspired by well-known R commands like `lm`) while

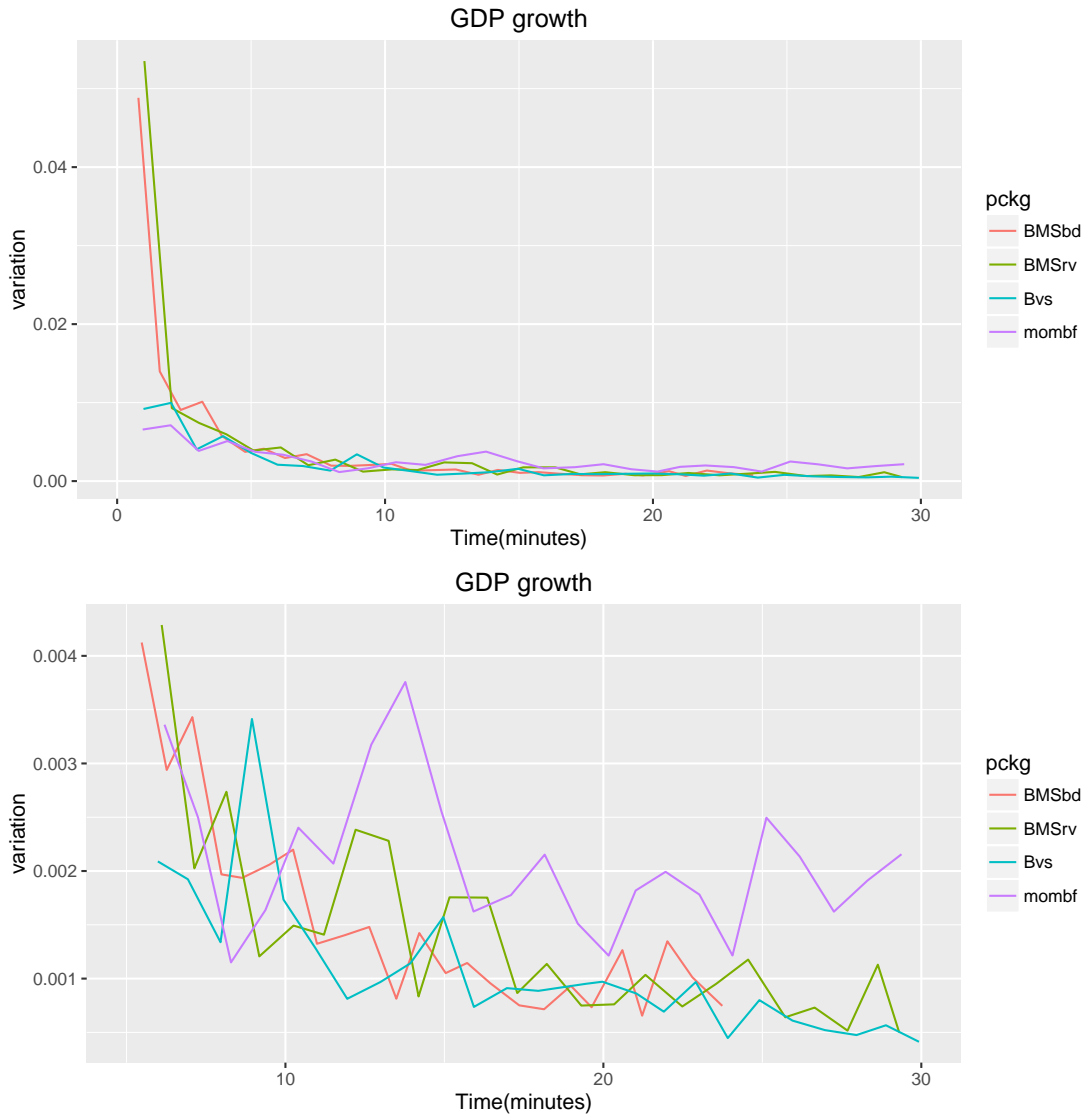


Figure 2: GDP growth data: variations in PIP ( $\Delta_{t,t-dt}$ ) as a function of computational time with  $dt = 60$  seconds (starting after the burning period). Both figures represents the same functions and just differ in that the figure below describes the behaviour after 5 minutes of computation.

in `mombf` these are defined through the arguments `y` and `x`. In `BMS` the dependent variable should be in the first column of the data provided and the rest play the role of explanatory variables.

### Extra functionalities

- *Fixed covariates.* By default only the intercept is included in all the competing models (cf. (1)) in all packages (but recall this is handled in `mombf` via centering). There could be situations where we wish to assume that certain covariates affect the response and these



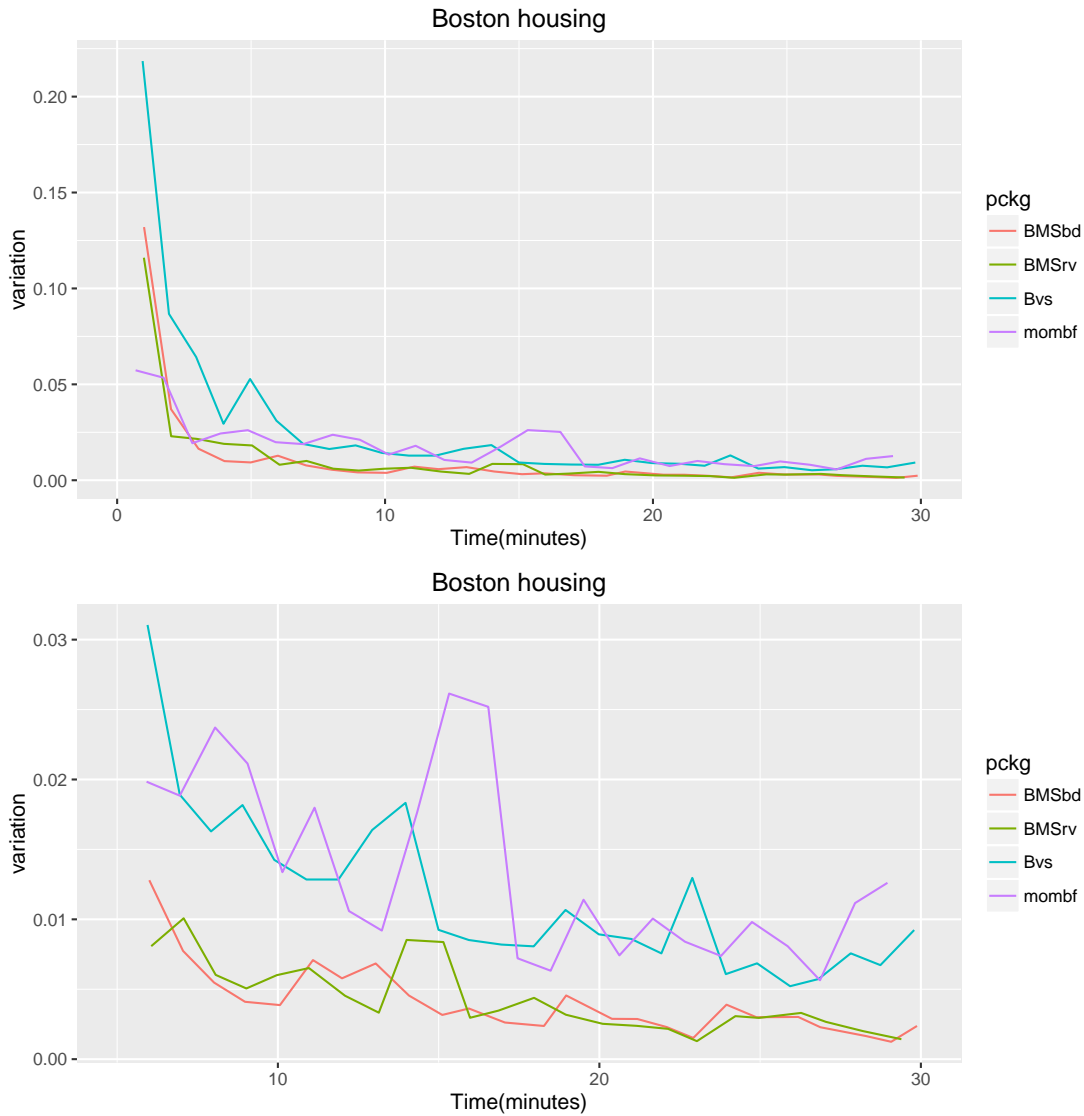


Figure 3: Boston Housing data: variations in PIP ( $\Delta_{t,t-dt}$ ) as a function of computational time with  $dt = 60$  seconds (starting after the burning period). Both figures represents the same functions and just differ in the axis represented (the figure below details the behaviour of the routines after 5 minutes of computation).

should be always included in the analysis (see, for instance Camarero et al., 2015). Both BMS and BayesVarSel include this possibility in their main commands.

- *Main terms and interactions.* On occasion, it is convenient to conserve the hierarchy between the explanatory variables in the way that interactions (or higher polynomial terms) are only included if the main terms are included (Peixoto, 1987). In Chipman et al. (1997) this is called the “heredity principle”. This would translate into a reduction of the model space. The package BMS accommodates this possibility through a modification of the sampling

algorithm.

- *Model comparison.* A complementary tool to the BMA exercise would be comparing separately some of the competing models (e.g. comparing the HPM and the MPM). These type of comparisons can be performed in BMS, BayesVarSel and BayesFactor.
- *Convergence.* BMS includes several interesting tools to analyse the convergence of the sampling methods implemented.
- *Parallel computation.* BMS, BayesVarSel and mombf have facilities to perform computations in parallel.

**Documentation** The four packages come with a detailed help with useful examples. Further, mombf and BMS have a comprehensive *vignette* with additional illustrations and written more pedagogically than the help documentations.

The packages BMS and BayesFactor are documented in the websites associated with Feldkircher and Zeugner (2014) (<http://bms.zeugner.eu>) and Morey (2015) (<http://bayesfactor.blogspot.com.es>), respectively. These sites contain manuals as well as valuable additional information, especially to users less familiar with model uncertainty techniques.

## 7 Conclusions and recommendations

In this paper, we have examined the behaviour and the possibilities of various R-packages available in CRAN for the purpose of Bayesian variable selection in linear regression. In particular, we compare the packages BMS, BayesVarSel, mombf and BayesFactor. It is clear that all packages concerned lead to very similar results, which is reassuring for the user. However, they do differ in the prior choices they allow, the way they present the output and the numerical strategies used. The latter affects CPU times, and, for example means that BayesVarSel is a good choice for small or moderate values of  $n$ , but BMS is preferable when  $n$  is large. The package BayesFactor can not deal with larger values of  $p$  and seems relatively slow, thus is not recommended for general use. mombf uses a slightly different prior from the one we focus on here (and which is the most commonly used), but is relatively competitive and closely approximates the PIPs after a short run time, albeit with slightly more variability than BMS or BayesVarSel.

In practice, users may be interested in specific features, such as always including certain covariates, that will dictate the choice of package. On the basis of its performance, the flexibility of prior choices and the extra features allowed, we would recommend the use of BayesVarSel for small or moderate values of  $n$ , and of BMS when  $n$  is large.

## Acknowledgments

The authors would like to thank David Rossell for valuable comments on a preliminary version of this paper.

## References

- Barbieri, M. M. and J. O. Berger (2004). Optimal predictive model selection. *The Annals of Statistics* 32(3), 870–897.
- Bayarri, M. J., J. O. Berger, A. Forte, and G. García-Donato (2012). Criteria for Bayesian model choice with application to variable selection. *Annals of Statistics* 40, 1550–1577.
- Berger, J. O. and G. Molina (2005). Posterior model probabilities via path-based pairwise priors. *Statistica Neerlandica* 59(1), 3–15.
- Berger, J. O. and L. R. Pericchi (2001). *Objective Bayesian Methods for Model Selection: Introduction and Comparison*, Volume 38 of *Lecture Notes–Monograph Series*, pp. 135–207. Beachwood, OH: Institute of Mathematical Statistics.
- Bivand, R. S., V. Gómez-Rubio, and H. Rue (2015). Spatial data analysis with R-INLA with some extensions. *Journal of Statistical Software* 63(20), 1–31.
- Camarero, M., A. Forte, G. García-Donato, Y. Mendoza, and J. Ordoñez (2015). Variable selection in the analysis of energy consumption-growth nexus. *Energy Economics* 52, Part A, 207–216.
- Carbonetto, P. and M. Stephens (2012). Scalable variational inference for Bayesian variable selection in regression, and its accuracy in genetic association studies. *Bayesian Analysis* 7(1), 73–108.
- Casella, G. and E. Moreno (2006). Objective Bayesian variable selection. *Journal of the American Statistical Association* 101(473), 157–167.
- Chipman, H., M. Hamada, and C. F. J. Wu (1997). A Bayesian variable-selection approach for analyzing designed experiments with complex aliasing. *Technometrics* 39(4), 372–381.
- Ehrlich, I. (1973). Participation in illegitimate activities: a theoretical and empirical investigation. *Journal of Political Economy* 81(3), 521–567.
- Eicher, T., C. Papageorgiou, and A. E. Raftery (2011). Default priors and predictive performance in Bayesian model averaging, with application to growth determinants. *Journal of Applied Econometrics* 26, 30–55.

- Feldkircher, M. and S. Zeugner (2014). R-package BMS Bayesian Model Averaging in R. <http://bms.zeugner.eu>.
- Fernández, C., E. Ley, and M. F. Steel (2001). Benchmark priors for Bayesian model averaging. *Journal of Econometrics* 100, 381–427.
- Foster, D. and E. I. George (1994). The Risk Inflation Criterion for Multiple Regression. *The Annals of Statistics* 22, 381–427.
- Fraley, C., A. E. Raftery, J. M. Slougher, T. Gneiting, and U. of Washington. (2015). *ensembleBMA: Probabilistic Forecasting using Ensembles and Bayesian Model Averaging*. R package version 5.1.2.
- Garcia-Donato, G. and A. Forte (2015). *BayesVarSel: Bayes Factors, Model Choice And Variable Selection In Linear Models*. R package version 1.6.1.
- Garcia-Donato, G. and M. A. Martinez-Beneito (2013). On Sampling strategies in Bayesian variable selection problems with large model spaces. *Journal of the American Statistical Association* 108(501), 340–352.
- George, E. I. and R. E. McCulloch (1993). Variable selection via Gibbs sampling. *Journal of the American Statistical Association* 88(423), 881–889.
- Ishwaran, H., J. Rao, and U. Kogalur (2013). *spikeslab: Prediction and Variable Selection Using Spike and Slab Regression*. R package version 1.1.5.
- Jeffreys, H. (1961). *Theory of Probability* (3rd ed.). Oxford University Press.
- Johnson, V. E. and D. Rossell (2010). On the use of non-local prior densities in Bayesian hypothesis tests. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 72(2), 143–170.
- Johnson, V. E. and D. Rossell (2012). Bayesian model selection in high-dimensional settings. *Journal of the American Statistical Association* 107(498), 649–660.
- Kass, R. E. and L. Wasserman (1995). A reference Bayesian test for nested hypotheses and its relationship to the schwarz criterion. *Journal of the American Statistical Association* 90(431), 928–934.
- Katabuchi, M. and A. Nakamura (2015). *mglm: Model Averaging for Multivariate GLM with Null Models*. R package version 0.0.2.
- Ley, E. and M. F. Steel (2007). Jointness in Bayesian variable selection with applications to growth regression. *Journal of Macroeconomics* 29(3), 476 – 493. Special Issue on the Empirics of Growth Nonlinearities.

- Ley, E. and M. F. Steel (2009). On the effect of prior assumptions in Bayesian model averaging with applications to growth regression. *Journal of Applied Econometrics* 24(4), 651–674.
- Ley, E. and M. F. Steel (2012). Mixtures of g-priors for Bayesian model averaging with economic applications. *Journal of Econometrics* 171(2), 251 – 266. Bayesian Models, Methods and Applications.
- Liang, F., R. Paulo, G. Molina, M. A. Clyde, and J. O. Berger (2008). Mixtures of g-priors for Bayesian variable selection. *Journal of the American Statistical Association* 103(481), 410–423.
- Madigan, D. and J. York (1995). Bayesian graphical models for discrete data. *International Statistical Review* 63, 215–232.
- Maruyama, Y. and E. I. George (2011). Fully Bayes factors with a generalized g-prior. *The Annals of Statistics* 39(5), 2740–2765.
- McCormick, T. H., A. E. Raftery, and D. Madigan (2014). *dma: Dynamic model averaging*. R package version 1.2-2.
- Mitchell, T. and J. Beauchamp (1988). Bayesian variable selection in linear regression. *Journal of the American Statistical Association* 83, 1023–1032.
- Morey, R. (2015). BayesFactor an R package for Bayesian data analysis. <http://bayesfactorppl.r-forge.r-project.org>.
- Morey, R. D., J. N. Rouder, and T. Jamil (2015). *BayesFactor: Computation of Bayes Factors for Common Designs*. R package version 0.9.11-1.
- Peixoto, J. (1987). Hierarchical variable selection in polynomial regression models. *American Statistician* 44(1), 26–30.
- R Core Team (2015). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Raftery, A., J. Hoeting, C. Volinsky, I. Painter, and K. Y. Yeung (2015). *BMA: Bayesian Model Averaging*. R package version 3.18.4.
- Raftery, A. E. (1995). Bayesian model selection in social research. *Sociological Methodology* 25, 111–163.
- Raftery, A. E., D. Madigan, and J. Hoeting (1997). Bayesian model averaging for linear regression models. *Journal of the American Statistical Association* 92, 179–191.

- Rossell, D., J. D. Cook, D. Telesca, and P. Roebuck (2014). *mombf: Moment and Inverse Moment Bayes factors*. R package version 1.5.9.
- Rue, H., S. Martino, and N. Chopin (2009). Approximate Bayesian inference for latent gaussian models using integrated nested Laplace approximations (with discussion). *Journal of the Royal Statistical Society: Series B* 71(2), 319–392.
- Sala-I-Martin, X., G. Doppelhofer, and R. I. Miller (2004). Determinants of long-term growth: A Bayesian averaging of classical estimates (BACE) approach. *American Economic Review* 94(4), 813–835.
- Schäfer, C. and N. Chopin (2013). Sequential monte carlo on large binary sampling spaces. *Statistics and Computing* 23(2), 163–184.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics* 6(2), 461–464.
- Scott, J. G. and J. O. Berger (2010). Bayes and empirical-Bayes multiplicity adjustment in the variable-selection problem. *The Annals of Statistics* 38(5), 2587–2619.
- Tobias, J. L. and M. Li (2004). Returns to schooling and Bayesian model averaging: A union of two literatures. *Journal of Economic Surveys* 18(2), 153–180.
- Venables, W. N. and B. D. Ripley (2002). *Modern Applied Statistics with S* (Fourth ed.). New York: Springer. ISBN 0-387-95457-0.
- Zellner, A. (1986). On assessing prior distributions and Bayesian regression analysis with g-prior distributions. In A. Zellner (Ed.), *Bayesian Inference and Decision techniques: Essays in Honor of Bruno de Finetti*, pp. 389–399. Edward Elgar Publishing Limited.
- Zellner, A. and A. Siow (1980). Posterior odds ratio for selected regression hypotheses. In J. M. Bernardo, M. DeGroot, D. Lindley, and A. F. M. Smith (Eds.), *Bayesian Statistics 1*, pp. 585–603. Valencia: University Press.
- Zellner, A. and A. Siow (1984). *Basic Issues in Econometrics*. Chicago: University of Chicago Press.
- Zeugner, S. and M. Feldkircher (2015). Bayesian model averaging employing fixed and flexible priors: The BMS package for R. *Journal of Statistical Software* 68(4), 1–37.

## Appendix A. Summarizing the output for the Crime data set

### BMS

- Call:

```
bms(X.data=lUScrime, g="UIP", mprior="uniform", nmodel=100,g.stats = F)
```

- print()

```
> Ob
```

|      | PIP       | Post Mean     | Post SD    | Cond.Pos.  | Sign | Idx |
|------|-----------|---------------|------------|------------|------|-----|
| Ineq | 0.9974810 | 1.4165246470  | 0.35866715 | 0.99999993 |      | 12  |
| Ed   | 0.9775864 | 1.9044911340  | 0.61687338 | 0.99999929 |      | 2   |
| Prob | 0.8963338 | -0.2156149890 | 0.11648121 | 0.00000000 |      | 13  |
| M    | 0.8503615 | 1.1652362359  | 0.67546221 | 0.99999999 |      | 1   |
| NW   | 0.6792925 | 0.0666392373  | 0.05770554 | 1.00000000 |      | 8   |
| Po1  | 0.6654873 | 0.6238407271  | 0.52893431 | 0.99999947 |      | 3   |
| U2   | 0.5996084 | 0.2030465034  | 0.21658823 | 0.99898978 |      | 10  |
| Po2  | 0.4215797 | 0.3263306162  | 0.51374655 | 0.94520349 |      | 4   |
| Time | 0.3333490 | -0.0792972600 | 0.15550003 | 0.05684919 |      | 14  |
| Pop  | 0.3301836 | -0.0207565713 | 0.03847876 | 0.00026644 |      | 7   |
| GDP  | 0.3124840 | 0.1830703611  | 0.35290133 | 0.99921131 |      | 11  |
| So   | 0.2306890 | 0.0316629469  | 0.08629093 | 0.97051521 |      | 15  |
| U1   | 0.2082608 | -0.0196768907 | 0.15978060 | 0.40681795 |      | 9   |
| M.F  | 0.1603299 | 0.0007683185  | 0.69992351 | 0.50105177 |      | 6   |
| LF   | 0.1567424 | 0.0445475745  | 0.27607008 | 0.76320873 |      | 5   |

| Mean no. regressors | Draws          | Burnins         | Time            |
|---------------------|----------------|-----------------|-----------------|
| "7.8198"            | "32768"        | "0"             | "3.188557 secs" |
| No. models visited  | Modelspace 2^K | % visited       | % Topmodels     |
| "32768"             | "32768"        | "100"           | "0.31"          |
| Corr PMP            | No. Obs.       | Model Prior     | g-Prior         |
| "NA"                | "47"           | "uniform / 7.5" | "UIP"           |

- summary()

```
> summary(ob)
```

| Mean no. regressors | Draws   | Burnins | Time            |
|---------------------|---------|---------|-----------------|
| "7.8198"            | "32768" | "0"     | "3.188557 secs" |

| No. models visited | Modelspace 2^K | % visited       | % Topmodels |
|--------------------|----------------|-----------------|-------------|
| "32768"            | "32768"        | "100"           | "0.31"      |
| Corr PMP           | No. Obs.       | Model Prior     | g-Prior     |
| "NA"               | "47"           | "uniform / 7.5" | "UIP"       |

- plot() and image() Figure 4

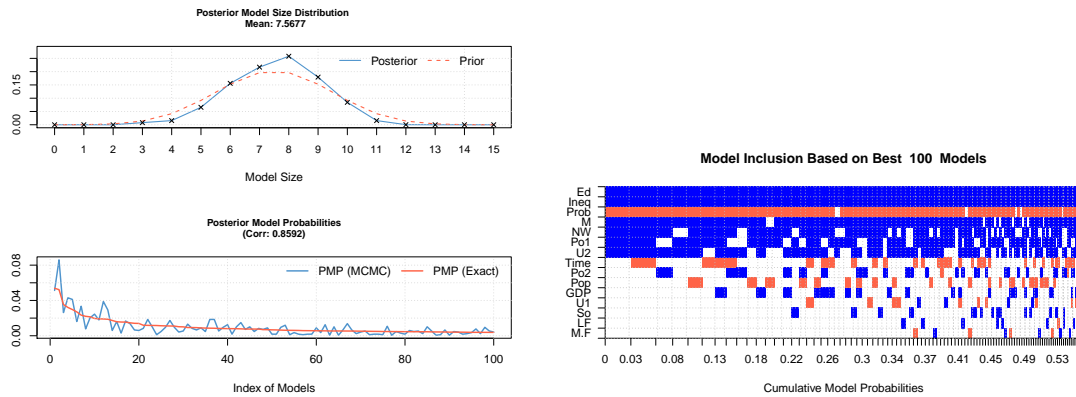


Figure 4: Prior/Posterior probabilities of each dimension and convergence performance plotted using BMS with plot(ob)(top) and Model inclusion probabilities based on the best 100 models using image(ob)(bottom)

- Predictive density plot(pred.density(ob, newdata)) Figure 5

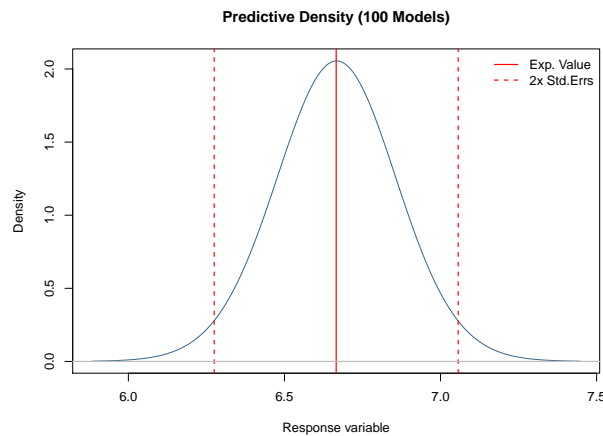


Figure 5: Predictive density using the 100 most probable models



## BayesFactor

- Call:  
`regressionBF(formula=y . ,data=lUScrime, rscaleCont = 1, noSample=TRUE)`
- `head(ob)`

Bayes factor analysis

-----

```
[1] M + Ed + Po1 + NW + U2 + Ineq + Prob + Time      : 23217774828 0%
[2] M + Ed + Po1 + NW + U2 + Ineq + Prob             : 22390162522 0%
[3] M + Ed + Po2 + NW + U2 + Ineq + Prob             : 15146223575 0%
[4] M + Ed + Po1 + Pop + NW + U2 + Ineq + Prob       : 13801956643 0%
[5] M + Ed + Po1 + U2 + Ineq + Prob                  : 12988208736 0%
[6] M + Ed + Po1 + NW + U2 + GDP + Ineq + Prob + Time : 11922310558 0%
```

Against denominator:

Intercept only

----

Bayes factor type: BFlinearModel, JZS

- `plot(head(ob))` Figure 6

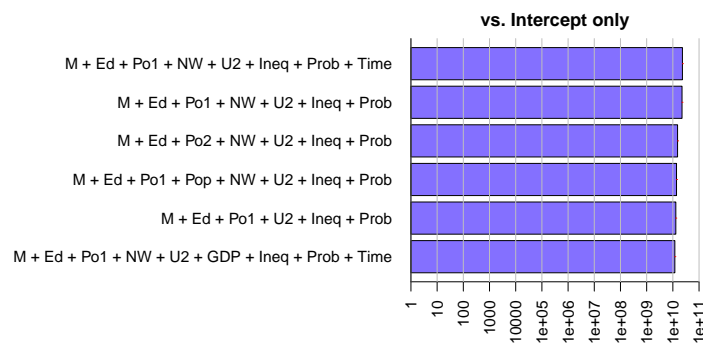


Figure 6: Bayes Factors using BayesFactor

## BayesVarSel

- Call:

```
Bvs(formula="y .", data=lUScrime, prior.betas="gZellner", n.keep=10)
```

- print()

```
>ob
```

```
Call:
```

```
Bvs(formula = "y~.", data = lUScrime, prior.betas = "gZellner",  
     n.keep = 10)
```

The 10 most probable models and their probabilities are:

|    | M | Ed | Po1 | Po2 | LF | M.F | Pop | NW | U1 | U2 | GDP | Ineq | Prob | Time | So | prob        |
|----|---|----|-----|-----|----|-----|-----|----|----|----|-----|------|------|------|----|-------------|
| 1  | * | *  | *   |     |    |     |     |    | *  | *  |     | *    | *    |      |    | 0.024695841 |
| 2  | * | *  | *   |     |    |     |     |    | *  | *  |     | *    | *    | *    |    | 0.023987471 |
| 3  | * | *  |     | *   |    |     |     |    | *  | *  |     | *    | *    |      |    | 0.016258755 |
| 4  | * | *  | *   |     |    |     |     |    |    | *  |     | *    | *    |      |    | 0.014728156 |
| 5  | * | *  | *   |     |    |     | *   | *  | *  | *  |     | *    | *    |      |    | 0.013640810 |
| 6  | * | *  | *   |     |    |     |     | *  |    |    |     | *    | *    | *    |    | 0.012415831 |
| 7  | * | *  | *   |     |    |     |     | *  | *  | *  | *   | *    | *    | *    |    | 0.010720692 |
| 8  | * | *  |     | *   |    |     |     | *  | *  |    |     | *    | *    | *    |    | 0.010106903 |
| 9  | * | *  |     | *   |    |     |     |    |    | *  |     | *    | *    |      |    | 0.009834356 |
| 10 | * | *  | *   |     |    |     | *   | *  |    |    |     | *    | *    |      |    | 0.008994454 |

- summary()

```
> summary(ob)
```

```
Call:
```

```
Bvs(formula = "y~.", data = lUScrime, prior.betas = "gZellner",  
     n.keep = 10)
```

Inclusion Probabilities:

|     | Incl.prob. | HPM | MPM |
|-----|------------|-----|-----|
| M   | 0.8504     | *   | *   |
| Ed  | 0.9776     | *   | *   |
| Po1 | 0.6655     | *   | *   |
| Po2 | 0.4216     |     |     |

|      |        |   |   |
|------|--------|---|---|
| LF   | 0.1567 |   |   |
| M.F  | 0.1603 |   |   |
| Pop  | 0.3302 |   |   |
| NW   | 0.6793 | * | * |
| U1   | 0.2083 |   |   |
| U2   | 0.5996 | * | * |
| GDP  | 0.3125 |   |   |
| Ineq | 0.9975 | * | * |
| Prob | 0.8963 | * | * |
| Time | 0.3333 |   |   |
| So   | 0.2307 |   |   |

---

Code: HPM stands for Highest posterior Probability Model and  
MPM for Median Probability Model.

- `plotBvs()` The `BayesVarSel` package has its own plot function named as `plotBvs()` which shows different information depending of the value of the argument `option`. If `option` is `'d'` a Barplot of the posterior probabilities of each model dimension is printed. If we set `option='j'` we obtain the posterior probability of every two covariates being together in the model. `option='c'` present the posterior conditional probability of a variable (in rows) given an other variable (in columns) is already in the model. Finally `option='n'` plots the posterior conditional probability of a variable (in rows) being in the model given than other variable (in columns) is not. See Figure 7.

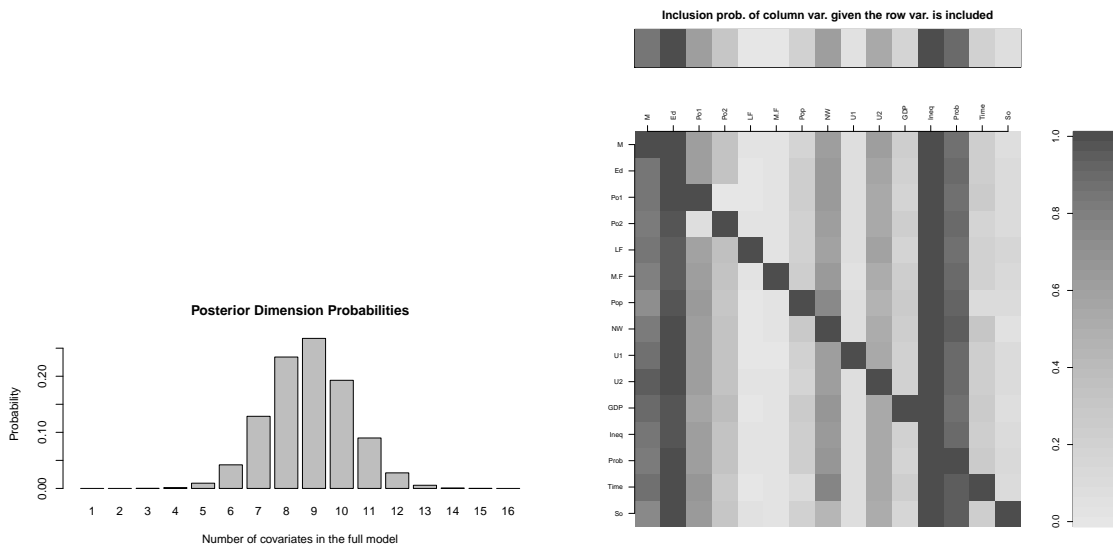


Figure 7: Posterior probabilities of each dimension plotted using BayesVarSel with `plotBvs(ob, option='d')`(left) and posterior conditional probabilities `plotBvs(ob, option='c')`(right)

## mombf

- Call:  
`modelSelection(y=lUScrime$y, x=lUScrime[,-1], priorCoef=zellnerprior(tau=n), center=TRUE, scale=FALSE, priorDelta=modelunifprior(),priorVar=igprior(alpha=.001,lambda=.001))`
- `print()` and `summary()`

```
> ob
```

```
msfit object with 15 variables
```

```
Use postProb() to get posterior model probabilities
```

```
Elements $margpp, $postMode, $postSample and $coef contain further information
```

```
(see help('msfit') and help('modelSelection') for details)
```

```
> mombf.crime$margpp #posterior inclusion probabilities
```

```
[1] 0.8573872 0.9807889 0.6853461 0.4037062 0.1563102 0.1611739 0.3349094 0.6960262
```

```
[9] 0.2103670 0.6125000 0.3171652 0.9977052 0.9081704 0.3449367 0.2321784
```

```
> mombf.crime$postMode
```

```
[1] 1 1 1 0 0 0 0 1 0 1 0 1 1 1 0
```

```
> round(mombf.crime$coef,4)
```

```
[1] 1.4472 2.1749 0.8347 0.0000 0.0000 0.0000 0.0000 0.1066 0.0000 0.2827
```

```
[11] 0.0000 1.2120 -0.3039 -0.2806 0.0000
```