

Mode jumping MCMC for Bayesian variable selection in GLMM

Aliaksandr Hubin ¹

Department of Mathematics, University of Oslo

and

Geir Storvik

Department of Mathematics, University of Oslo

Abstract

Generalized linear mixed models (GLMM) are used for inference and prediction in a wide range of different applications providing a powerful scientific tool. An increasing number of sources of data are becoming available introducing a variety of hypothetical explanatory variables for these models to be considered. Selection of an optimal combination of these variables is thus becoming crucial. In a Bayesian setting, the posterior distribution of the models, based on the observed data, can be viewed as a relevant measure for the model evidence. The model space increases exponential in the number of candidate variables and has numerous local extrema. To resolve these issues a novel MCMC algorithm for the search through the model space via efficient mode jumping for GLMMs is introduced. The algorithm is based on that marginal likelihoods can be efficiently calculated within each model. It

¹The corresponding author, Aliaksandr Hubin, is a PhD candidate at the University of Oslo, 0851 Moltke Moes vei 35 Oslo, Norway. Email: aliaksah@math.uio.no, tel: +4745171361.

The paper has supplementary materials consisting of:

R package: *R* package *EMJMCMC* to perform the efficient mode jumping MCMC described in the paper. (EMJMCMC_1.2.tar.gz; GNU zipped tar file).

Data and code: Data (simulated and real) and *R* code for MJMCMC algorithm, post-processing and creating figures wrapped together into a reference based EMJMCMC class. (code-and-data.zip; zip file containing the data, code and a read-me file (readme.pdf))

Details and Pseudo code: Proofs of the ergodicity of MJMCMC procedure, pseudo codes for MJMCMC and local combinatorial optimizers, parallelization strategies and some supplementary tables for the experiments. (appendix.pdf)

is recommended that either exact expressions or precise approximations of marginal likelihoods are applied. The suggested algorithm is further applied to some simulated data, the famous U.S. crime data, protein activity data and epigenetic data and compared to several existing approaches.

Keywords: Bayesian variable selection; Bayesian model averaging; Generalized linear mixed models; Auxiliary variables; Combinatorial optimization; High performance computations.

1. Introduction

In this paper we study variable selection in generalized linear mixed models (GLMM) addressed in a Bayesian setting. Being one of the most powerful modeling tools in modern statistical science (Stroup, 2013) these models have proven to be efficient in numerous applications including the simple banking scoring problems (Grossi and Bellini, 2006) and insurance claims modeling (David, 2015), studies on the course of illness in schizophrenia and linking diet with heart diseases (Skrondal and Rabe-Hesketh, 2003), analyzing sophisticated astrophysical data (de Souza et al., 2015), and genomics data (Lobbraux and Melodelima, 2015). In many of these applications, the number of candidate explanatory variables (covariates) is large, making variable selection a difficult problem, both conceptually and numerically. In this paper we will focus on efficient Markov chain Monte Carlo algorithms for such variable selection problems. Our focus will be on posterior model probabilities although other model selection criteria can also easily be adopted within the algorithm.

Algorithms for variable selection in the Bayesian settings have been previously addressed, but primarily in the combined space of models *and* parameters. George and McCulloch (1997) describe and compare various hierarchical mixture prior formulations for Bayesian variable selection in normal linear regression models. Then they describe computational methods including Gray Code sequencing and standard MCMC for posterior evaluation and exploration of the space of models. They also describe the infeasibility of exhaustive exploration of the space of models for moderately large problems as well as the inability of standard MCMC techniques to escape from local optima efficiently. Ghosh (2015) also addresses MCMC algorithms to estimate the posterior distribution over models. However, she mentions that estimates of posterior probabilities of individual models based on MCMC

output are often not reliable because the number of MCMC samples is typically by far smaller than the size of the model space. The authors show that their algorithm can, under some conditions, outperform standard MCMC. [Song and Liang \(2015\)](#) address the case when there is by far more explanatory variables than observations and suggest a split and merge Bayesian model selection algorithm that first splits the set of covariates into a number of subsets, then finds relevant variables from these subsets and in the second stage merges these relevant variables and performs a new selection from the merged set. This algorithm in general cannot guarantee convergence to a global optimum or find the true posterior distribution of the models, however under some strict regularity conditions it does so asymptotically. [Al-Awadhi et al. \(2004\)](#) considered using several MCMC steps within a new model to obtain good proposals within the combined parameter and model domain while [Yeh et al. \(2012\)](#) proposed local annealing approaches. Multiple try MCMC methods with local optimization have been described by [Liu et al. \(2000\)](#). These methods fall into the category of generating auxiliary states for proposals ([Storvik, 2011](#); [Chopin et al., 2013](#)). Yet another approach for Bayesian model selection is addressed by [Bottolo et al. \(2011\)](#), who propose the moves of MCMC between local optima through a permutation based genetic algorithm that has a pool of solutions in a current generation suggested by the parallel tempered chains, which allows to achieve a reasonably good mixing of the chains and escape from local modes at a reasonable rate. A similar idea is considered by [Frommlet et al. \(2012\)](#).

For an increasing number of models, marginal likelihoods for specific models can be efficiently calculated, making the exploration of models far much easier. [Bivand et al. \(2014\)](#) combine approximations of marginal likelihood with Bayesian model averaging within spatial models. [Clyde et al. \(2011\)](#) suggest a Bayesian adaptive sampling (BAS) algorithm as an alternative to MCMC allowing for perfect sampling without replacement. [Bové and Held \(2011\)](#) consider an MCMC algorithm within the model space in cases where marginal likelihoods are available, but only allow local moves.

Different approaches for calculation of marginal likelihoods are available. For linear models with conjugate priors analytic expressions are available ([Clyde et al., 2011](#)). In more general settings, MCMC algorithms combined with e.g. Chib’s method ([Chib, 1995](#)) can be applied, although computational expensive. See also [Friel and Wyse \(2012\)](#) for alternative MCMC based methods. For Gaussian latent variables, the computational task can be efficiently solved through the integrated nested Laplace approximation

(INLA) approach (Rue et al., 2009). Hubin and Storvik (2016) compare INLA with MCMC based methods, showing that INLA based approximations are extremely accurate and require much less computational effort than the MCMC approaches for within model calculations.

In this paper we introduce a novel MCMC algorithm for the search through the model space, the mode jumping MCMC (MJMCMC). The focus will be on Gaussian latent variable models, for which efficient approximations to marginal likelihoods are available. The algorithm is based on the idea of mode jumping within MCMC - resulting in an MCMC algorithm which manages to efficiently explore the model space by means of mode jumping, applicable through large jumps combined with local optimization. Mode jumping MCMC methods within a continuous space setting were first suggested by Tjelmeland and Hegstad (1999). We modify the algorithm to the discrete space of possible models, requiring both new ways of making large jumps and of performing local optimization. We include mixtures of proposal distributions and parallelization for further improving the performance of the algorithm. A valid acceptance probability within the Metropolis-Hastings setting is constructed based on the use of backward kernels.

2. The generalized linear mixed model

We consider the following generalized linear mixed model:

$$Y_i|\mu_i \sim \mathbf{f}(y|\mu_i), \quad \mu_i = g^{-1}(\eta_i), \quad (1)$$

$$\eta_i = \beta_0 + \sum_{j=1}^p \gamma_j \beta_j x_{ij} + \delta_i, \quad (2)$$

$$\boldsymbol{\delta} = (\delta_1, \dots, \delta_n) \sim N_n(\mathbf{0}, \boldsymbol{\Sigma}_b). \quad (3)$$

Here Y_i is the response variable while $x_{ij}, j = 1, \dots, p$ are the covariates. We assume $\mathbf{f}(y|\mu)$ is a density/distribution from the exponential family with corresponding link function $g(\cdot)$. The latent indicators $\gamma_j \in \{0, 1\}, j = 1, \dots, p$ define if covariate x_{ij} is included into the model ($\gamma_j = 1$) or not ($\gamma_j = 0$) while $\beta_j \in \mathbb{R}, j \in \{0, \dots, p\}$ are the corresponding regression coefficients. We are also addressing the unexplained variability of the responses and the correlation structure between them through random effects δ_i with a specified parametric and sparse covariance matrix structure defined through $\boldsymbol{\Sigma}_b = \boldsymbol{\Sigma}_b(\boldsymbol{\psi}) \in \mathbb{R}^n \times \mathbb{R}^n$ where $\boldsymbol{\psi}$ are parameters describing the correlation structure.

In order to put the model into a Bayesian framework, we assume

$$\gamma_j|q \sim \text{Binom}(1, q), \quad j = 1, \dots, p, \quad (4)$$

$$q \sim \text{Beta}(\alpha_q, \beta_q), \quad (5)$$

where q is the prior probability of including a covariate into the model. For $(\boldsymbol{\beta}, \boldsymbol{\psi})$ different priors are possible, see the applications in section 4.

Let $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$, which uniquely defines a specific model. Assuming the constant term β_0 is always included, there are $L = 2^p$ different models. We want to find a set of the best models with respect to posterior model probabilities $p(\boldsymbol{\gamma}|\mathbf{y})$, where $\mathbf{y} = (y_1, \dots, y_n)$. We assume that marginal likelihoods $p(\mathbf{y}|\boldsymbol{\gamma})$ are available for a given $\boldsymbol{\gamma}$, and then use MCMC to explore $p(\boldsymbol{\gamma}|\mathbf{y})$. By Bayes formula

$$p(\boldsymbol{\gamma}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\gamma})p(\boldsymbol{\gamma})}{\sum_{\boldsymbol{\gamma}' \in \Omega} p(\mathbf{y}|\boldsymbol{\gamma}')p(\boldsymbol{\gamma}')}. \quad (6)$$

In order to calculate $p(\boldsymbol{\gamma}|\mathbf{y})$ we have to iterate through the whole model space Ω , which becomes computationally infeasible for large p . We aim at approximating $p(\boldsymbol{\gamma}|\mathbf{y})$ by means of searching for some subspace \mathbb{V} of Ω giving rise to the approximation

$$\hat{p}(\boldsymbol{\gamma}|\mathbf{y}) = \frac{\mathbb{1}(\boldsymbol{\gamma} \in \mathbb{V})p(\mathbf{y}|\boldsymbol{\gamma})p(\boldsymbol{\gamma})}{\sum_{\boldsymbol{\gamma}' \in \mathbb{V}} p(\mathbf{y}|\boldsymbol{\gamma}')p(\boldsymbol{\gamma}')}, \quad (7)$$

where $\mathbb{1}(\cdot)$ is the indicator function. [Clyde et al. \(2011\)](#) named these the renormalized (RM) model estimates. Models with high values of $p(\mathbf{y}|\boldsymbol{\gamma})$ are important to be addressed. This means that modes and near modal values of marginal likelihoods are particularly important for construction of reasonable $\mathbb{V} \subset \Omega$ and missing them can dramatically influence our estimates. In this context the denominator of (7), which we would like to be as high as possible, becomes an extremely relevant measure for the quality of the search in terms of being able to capture whether the algorithm visits all of the modes, whilst the size of \mathbb{V} should be low in order to save computational time.

An alternative to (7) is the ordinary MCMC based estimate

$$\tilde{p}(\boldsymbol{\gamma}|\mathbf{y}) = \frac{\sum_{i=1}^W \mathbb{1}(\boldsymbol{\gamma}^{(i)} = \boldsymbol{\gamma})}{W} \xrightarrow[W \rightarrow \infty]{d} p(\boldsymbol{\gamma}|\mathbf{y}), \quad (8)$$

where W is the total number of MCMC samples. Although (8) is asymptotically consistent, (7) will often be preferable estimators since convergences of the MCMC based approximation (8) is much slower, see [Clyde et al. \(2011\)](#).

The posterior marginal inclusion probability $p(\gamma_j = 1|\mathbf{y})$ can be approximated by

$$\widehat{p}(\gamma_j = 1|\mathbf{y}) = \sum_{\boldsymbol{\gamma}' \in \mathbb{V}} \mathbb{1}(\gamma'_j = 1) \widehat{p}(\boldsymbol{\gamma}'|\mathbf{y}), \quad (9)$$

giving a measure for assessing importance of the covariates. Other parameters can be estimated similarly.

Algorithms for estimating \mathbb{V} are described in section 3. In practice $p(\mathbf{y}|\boldsymbol{\gamma})$ may not be available analytically. We then rely on some precise approximations $\widehat{p}(\mathbf{y}|\boldsymbol{\gamma})$. Such approximations introduce additional errors in (7) and (9), but we assume them to be small enough to be ignored. This is further discussed in section 3.4.

3. Mode jumping Markov chain Monte Carlo

MCMC algorithms (Robert and Casella, 2005) have been extremely popular for the exploration of model spaces for model selection, being capable of providing samples from the posterior distribution of the models. In our setting, the most important aspect becomes building a method to explore the model space in a way to efficiently switch between potentially sparsely located modes, whilst avoiding visiting models with a low $p(\mathbf{y}|\boldsymbol{\gamma})$ too often.

3.1. Standard Metropolis-Hastings

Metropolis-Hastings algorithms (Robert and Casella, 2005) are a class of MCMC methods for drawing from a complicated target distribution living on some space Ω , which in our setting will be $\pi(\boldsymbol{\gamma}) = p(\boldsymbol{\gamma}|\mathbf{y})$. Given some proposal distribution $q(\boldsymbol{\gamma}^*|\boldsymbol{\gamma})$, the Metropolis-Hastings algorithm accepts the proposed $\boldsymbol{\gamma}^*$ with probability

$$r_{mh}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^*) = \min \left\{ 1, \frac{\pi(\boldsymbol{\gamma}^*)q(\boldsymbol{\gamma}|\boldsymbol{\gamma}^*)}{\pi(\boldsymbol{\gamma})q(\boldsymbol{\gamma}^*|\boldsymbol{\gamma})} \right\}, \quad (10)$$

and otherwise remains in the old state $\boldsymbol{\gamma}$. This will generate a Markov chain which, given the chain is irreducible and aperiodic, will have π as stationary distribution. Theoretical results related to convergence of MCMC based estimates can be found in e.g. Tierney (1996). Note that the discrete finite space of models make these results easily applicable in our case.

Given that the γ_j 's are binary, changes correspond to swaps between the values 0 and 1. One can address various options for generating proposals. A

Type	Proposal	Label
1	$\frac{\prod_{i \in \{i_1, \dots, i_S\}} \rho_i}{\binom{p}{S}(\eta - \zeta + 1)}$	<i>Random change with random size of the neighborhood</i>
2	$\frac{\prod_{i \in \{i_1, \dots, i_S\}} \rho_i}{\binom{p}{S}}$	<i>Random change with fixed size of the neighborhood</i>
3	$\frac{1}{\binom{p}{S}(\eta - \zeta + 1)}$	<i>Swap with random size of the neighborhood</i>
4	$\binom{p}{S}^{-1}$	<i>Swap with fixed size of the neighborhood</i>
5	$\frac{1 - \mathbb{1}(\sum_i^p \gamma_i = p)}{p - \sum_i^p \gamma_i + \mathbb{1}(\sum_i^p \gamma_i = p)}$	<i>Uniform addition of a covariate</i>
6	$\frac{1 - \mathbb{1}(\sum_i^p \gamma_i = 0)}{\sum_i^p \gamma_i + \mathbb{1}(\sum_i^p \gamma_i = 0)}$	<i>Uniform deletion of a covariate</i>

Table 1: Types of proposals suggested for moves between models during an MCMC procedure. Here S is either a deterministic or random ($S \sim \text{Unif}\{\zeta, \dots, \eta\}$) size of the neighborhood; ρ_i is the probability of inclusion of variable γ_i .

simple proposal is to first select the number of components to change, e.g. $S \sim \text{Unif}\{\zeta, \dots, \eta\}$, followed by a sample of size S without replacement from $\{1, \dots, p\}$. This implies that in (10) the proposal probability for switching from γ to γ^* becomes symmetric, simplifying calculation of the acceptance probability. Other possibilities for proposals are summarized in Table 1, allowing, among others, different probabilities of swapping for the different components. Such probabilities can for instance be associated with marginal inclusion probabilities from a preliminary MCMC run.

3.2. MJMCMC - the mode jumping MCMC

The main problem with the standard Metropolis-Hastings algorithms is the trade-off between possibilities of large jumps (by which we understand proposals with a large neighborhood) and high acceptance probabilities. Large jumps will typically result in proposals with low probabilities. In a continuous setting, Tjelmeland and Hegstad (1999) solved this by introducing local optimization after large jumps. We adapt this approach to the discrete model selection setting by the following algorithm:

Algorithm 1 Mode jumping MCMC

- 1: Generate a large jump \mathbf{x}_0^* according to a proposal distribution $q_l(\mathbf{x}_0^*|\gamma)$.
- 2: Perform a local optimization, defined through $\mathbf{x}_k^* \sim q_o(\mathbf{x}_k^*|\mathbf{x}_0^*)$.
- 3: Perform a small randomization to generate the proposal $\gamma^* \sim q_r(\gamma^*|\mathbf{x}_k^*)$.
- 4: Generate backwards auxiliary variables $\mathbf{x}_0 \sim q_l(\mathbf{x}_0|\gamma^*)$, $\mathbf{x}_k \sim q_o(\mathbf{x}_k|\mathbf{x}_0)$.
- 5: Put

$$\gamma' = \begin{cases} \gamma^* & \text{with probability } r_{mh}(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*); \\ \gamma & \text{otherwise,} \end{cases}$$

where

$$r_{mh}^*(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) = \min \left\{ 1, \frac{\pi(\gamma^*)q_r(\gamma|\mathbf{x}_k)}{\pi(\gamma)q_r(\gamma^*|\mathbf{x}_k^*)} \right\}. \quad (11)$$

Here by the local optimization we understand some combinatorial optimization algorithm with a small neighborhood (compared to the one of the large jump).

The procedure is illustrated in Figure 1 where the backward sequence $\gamma^* \rightarrow \mathbf{x}_0 \rightarrow \mathbf{x}_k \rightarrow \gamma$, needed for calculating the acceptance probability, is included. For this algorithm, three proposals need to be specified, $q_l(\cdot|\cdot)$ specifying the first large jump, $q_o(\cdot|\cdot)$ specifying the local optimizer, and $q_r(\cdot|\cdot)$ specifying the last randomization.

π -invariance of the MJMCMC procedures is given by the following theorem (based on similar arguments as in [Storvik, 2011](#); [Chopin et al., 2013](#)):

Theorem 1. *Assume $\gamma \sim \pi(\cdot)$ and γ' is generated according to Algorithm 1. Then $\gamma' \sim \pi(\cdot)$.*

Proof. Since $\gamma \sim \pi(\cdot)$ and $(\mathbf{x}_0^*, \mathbf{x}_k^*) \sim q_l(\mathbf{x}_0^*|\gamma)q_o(\mathbf{x}_k^*|\mathbf{x}_0^*)$ we have that

$$(\gamma, \mathbf{x}_0^*, \mathbf{x}_k^*) \sim \pi(\gamma)q_l(\mathbf{x}_0^*|\gamma)q_o(\mathbf{x}_k^*|\mathbf{x}_0^*) \equiv \bar{\pi}(\gamma, \mathbf{x}_0^*, \mathbf{x}_k^*).$$

We may now consider $(\gamma^*, \mathbf{x}_0, \mathbf{x}_k)$ as a proposal in the extended space, generated according to the distribution $q_r(\gamma^*|\mathbf{x}_k^*)q_l(\mathbf{x}_0|\gamma^*)q_o(\mathbf{x}_k|\mathbf{x}_0)$. An ordinary Metropolis-Hastings iteration with respect to $\bar{\pi}(\gamma, \mathbf{x}_0^*, \mathbf{x}_k^*)$ is then to accept $(\gamma^*, \mathbf{x}_0, \mathbf{x}_k)$ with probability $r_{mh}^* = \min\{1, \alpha_{mh}^*\}$ where

$$\begin{aligned} \alpha_{mh}^* &= \frac{\bar{\pi}(\gamma^*, \mathbf{x}_0, \mathbf{x}_k)q_r(\gamma|\mathbf{x}_k)q_l(\mathbf{x}_0^*|\gamma)q_o(\mathbf{x}_k^*|\mathbf{x}_0^*)}{\bar{\pi}(\gamma, \mathbf{x}_0^*, \mathbf{x}_k^*)q_r(\gamma^*|\mathbf{x}_k^*)q_l(\mathbf{x}_0|\gamma^*)q_o(\mathbf{x}_k|\mathbf{x}_0)} \\ &= \frac{\pi(\gamma^*)q_l(\mathbf{x}_0|\gamma^*)q_o(\mathbf{x}_k|\mathbf{x}_0)q_r(\gamma|\mathbf{x}_k)q_l(\mathbf{x}_0^*|\gamma)q_o(\mathbf{x}_k^*|\mathbf{x}_0^*)}{\pi(\gamma)q_l(\mathbf{x}_0^*|\gamma)q_o(\mathbf{x}_k^*|\mathbf{x}_0^*)q_r(\gamma^*|\mathbf{x}_k^*)q_l(\mathbf{x}_0|\gamma^*)q_o(\mathbf{x}_k|\mathbf{x}_0)} = \frac{\pi(\gamma^*)q_r(\gamma|\mathbf{x}_k)}{\pi(\gamma)q_r(\gamma^*|\mathbf{x}_k^*)}, \end{aligned}$$

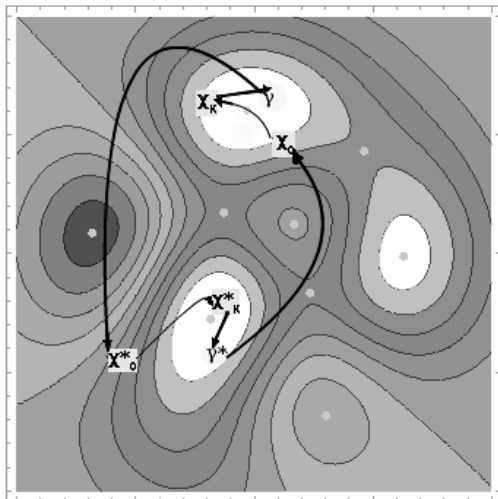


Figure 1: Illustration of intra mode mixing in MJMCMC.

proving the algorithm has $\bar{\pi}(\cdot)$ as invariant distribution. Since this distribution has $\pi(\cdot)$ as marginal distribution it follows that $\gamma' \sim \pi(\cdot)$. \square

Note that neither the large jump distribution $q_l(\cdot)$ nor the optimization distribution $q_o(\cdot)$ (which can be both deterministic and stochastic) are involved in the acceptance probability. This gives great flexibility in the choice of these distributions.

Large jumps are not performed at each iteration, but rather through a composition of standard Metropolis-Hastings steps with local moves and large jumps. As a rule of thumb, based on suggestions of [Tjelmeland and Hegstad \(1999\)](#) and our own experience, we recommended that in not more than 5% of the cases large jumps are performed. This is believed to provide the global Markov chain with both good mixing between the modes and accurate exploration of the regions around the modes. This in turn induces good performance of the algorithm in terms of the captured posterior mass for a given number of iterations. However, some tuning might well be required for the particular practical applications.

The mode jumping MCMC steps can be modified to include random choices of different proposal kernels q_l , q_o , and q_r and parallelized using the multiple try MCMC idea. Technical details are given in [Appendix A](#).

An illustrative example. Assume 10 covariates x_1, \dots, x_{10} and thus 1024 possible models. We generated $Y \sim N(1+10x_1+0.89x_8+1.43x_5, 1)$ with correlated

	Forward		Backward	
	Model	$\log(p(\mathbf{y} \gamma))$	Model	$\log(p(\mathbf{y} \gamma))$
Initial mode	$\gamma = 1010110111$	1606.21	$\gamma^* = 1101100001$	1612.27
Large jump	$\chi_0^* = 100\textcolor{red}{1}1100\textcolor{red}{0}1$	1541.51	$\chi_0 = 11\textcolor{red}{1}0100\textcolor{red}{1}11$	1608.55
Optimize	$\chi_k^* = 11\textcolor{red}{0}1\textcolor{blue}{1}0000\textcolor{blue}{0}$	1616.16	$\chi_k = 10\textcolor{red}{1}0100\textcolor{blue}{1}1\textcolor{blue}{0}$	1612.00
Randomize	$\gamma^* = 11\textcolor{red}{0}1\textcolor{blue}{1}0000\textcolor{green}{1}$	1612.27	$\gamma = 10\textcolor{red}{1}0\textcolor{blue}{1}10\textcolor{green}{1}1\textcolor{green}{1}$	1606.21
Acceptance probability: $\min\{1, 541.11\}$, accept $\gamma' = \gamma^* = \mathbf{1101100001}$				

Table 2: Illustration of a typical MJMCMC step with locally optimized proposals. The red components correspond to components swapped in the large jumps, the blue components to the ones changed in the optimizer, the green components of γ to the randomization step.

binary covariates (see supplementary material for details) and 1000 observations. We used a Gaussian linear regression with a Zellner’s g-prior (Zellner, 1986) with $g = 1000$. This model has tractable marginal likelihoods described in detail in section 4. We consider an MJMCMC step with a large jump swapping randomly 4 components of γ and a local greedy search as optimization routine. The last randomization changes each component of γ independently with probability equal to 0.1. A typical MJMCMC step with locally optimized proposals is illustrated in Table 2.

Large jumps. A change is defined by the components that are to be swapped. A simple choice is to give all components an equal probability ρ to be swapped and independence between components, in which case

$$q_l(\chi_0^*|\gamma) = \prod_{j=1}^p \rho^{I_j} (1 - \rho)^{1-I_j} = \rho^S (1 - \rho)^{p-S},$$

where I_j is a binary variable equal to 1 if component γ_j is to be swapped and $S = \sum_{j=1}^p I_j$ is the number of components to be swapped. An alternative is to first draw the number of components, S , to swap according to a distribution $q_S(\cdot)$ and thereafter choose (uniformly) among the possibilities. Table 1 describes different ways of jumping from γ to χ_0^* .

Optimization. In order to increase the quality of proposals and consequently both improve the acceptance ratio and increase the probability of escaping from local optima, the large jump is followed by local optimization. Typically, $q_o(\cdot)$ contains many iterations, generating intermediate states $\chi_0^* \rightarrow \chi_1^* \rightarrow \dots \rightarrow \chi_k^*$ but none of these intermediate states are needed for the final evaluation. Different local learning and optimization routines can be applied

for the generation of $\boldsymbol{\chi}_k^*$, both deterministic and stochastic ones. We will consider several feasible computationally options: local greedy optimization, local simulated annealing (SA) optimization, and local MCMC methods.

Randomization. A last randomization step defined through q_r is needed in order to make the move back from $\boldsymbol{\gamma}^*$ to $\boldsymbol{\gamma}$ feasible. We typically use randomizing kernels with a high mass on a small neighborhood around the mode but with a positive probability for any change. The two possible appropriate kernels from Table 1 are the random change of either random $S \sim \text{Unif}\{1, \dots, p\}$ or deterministic $S = p$ number of components with reasonably small but positive probabilities $0 < \rho_i \ll 1$. This guarantees that the MJMCMC procedure is irreducible in Ω .

In order for the acceptance probability to be high, it is crucial that the auxiliary variables in the reverse sequence $\boldsymbol{\chi} = (\boldsymbol{\chi}_0, \boldsymbol{\chi}_k)$ make $\boldsymbol{\gamma}$ plausible ($q_r(\boldsymbol{\gamma}|\boldsymbol{\chi}_k)$ should be large in (11)). This may be difficult to achieve because the backwards large jump has no guarantee to be close to the current state. One way to achieve this is to choose $q_l(\boldsymbol{\chi}_0^*|\boldsymbol{\gamma})$ to be symmetric, increasing the probability of returning close to the initial mode in the reverse step. The symmetry is achieved by swapping the same set of $\boldsymbol{\gamma}_j$'s in the large jumps in the forward simulation as in the backwards simulation. We record the components I that have been swapped. In our current implementation we require that only the components that do not correspond to I can be changed in optimization transition kernels.

Algorithm 2 Mode jumping MCMC with symmetric backwards jump

- 1: Generate a large jump χ_0^* by first generating a set $I \subset \{1, \dots, p\} \sim q_I(\cdot)$ defining the components to be swapped.
- 2: Perform a local optimization, defined through $\chi_k^* \sim q_o(\chi_k^*|\chi_0^*)$.
- 3: Perform a small randomization to generate the proposal $\gamma^* \sim q_r(\gamma^*|\chi_k^*)$.
- 4: Define the backwards large jump χ_0 through swapping the components I in γ^* .
- 5: Generate $\chi_k \sim q_o(\chi_k|\chi_0)$.
- 6: Put

$$\gamma' = \begin{cases} \gamma^* & \text{with probability } r_m(\gamma, \gamma^*; \chi_k, \chi_k^*); \\ \gamma & \text{otherwise,} \end{cases}$$

where

$$r_{mh}^*(\gamma, \gamma^*; \chi_k, \chi_k^*) = \min \left\{ 1, \frac{\pi(\gamma^*)q_r(\gamma|\chi_k)}{\pi(\gamma)q_r(\gamma^*|\chi_k^*)} \right\}. \quad (12)$$

The following theorem shows that also this algorithm has $\pi(\gamma)$ as invariant distribution:

Theorem 2. Assume $\gamma \sim \pi(\cdot)$ and γ' is generated according to Algorithm 2. Then $\gamma' \sim \pi(\cdot)$.

Proof. The stochastic auxiliary components are now I, χ_k^* and χ_k where χ_0^* and χ_0 are deterministic functions of (γ, I) and (γ^*, I) , respectively. We have

$$(\gamma, I, \chi_k^*) \sim \pi(\gamma)q_I(I)q_o(\chi_k^*|\chi_0^*) \equiv \bar{\pi}(\gamma, I, \chi_k^*).$$

We may now consider (γ^*, I, χ_k) as a proposal in the extended space, generated according to the distribution $q_r(\gamma^*|\chi_k^*)q_o(\chi_k|\chi_0)$. An ordinary Metropolis-Hastings iteration with respect to $\bar{\pi}(\gamma, I, \chi_k^*)$ is then to accept (γ^*, I, χ_k) with probability $r_{mh}^* = \min\{1, \alpha_{mh}^*\}$ where

$$\begin{aligned} \alpha_{mh}^* &= \frac{\bar{\pi}(\gamma^*, I, \chi_k)q_r(\gamma|\chi_k)q_o(\chi_k^*|\chi_0^*)}{\bar{\pi}(\gamma, I, \chi_k^*)q_r(\gamma^*|\chi_k^*)q_o(\chi_k|\chi_0)} \\ &= \frac{\pi(\gamma^*)q_I(I)q_o(\chi_k|\chi_0)q_r(\gamma|\chi_k)q_o(\chi_k^*|\chi_0^*)}{\pi(\gamma)q_I(I)q_o(\chi_k^*|\chi_0^*)q_r(\gamma^*|\chi_k^*)q_o(\chi_k|\chi_0)} = \frac{\pi(\gamma^*)q_r(\gamma|\chi_k)}{\pi(\gamma)q_r(\gamma^*|\chi_k^*)}, \end{aligned}$$

proving the algorithm has $\bar{\pi}(\cdot)$ as invariant distribution. Since this distribution has $\pi(\cdot)$ as marginal distribution it follows that $\gamma' \sim \pi(\cdot)$. \square

3.3. Delayed acceptance

The most computationally demanding parts of the MJMCMC algorithms are the forward and backward optimizations. In many cases, the proposal generated through the forward optimization may lead to a very small value of $\pi(\gamma^*)$ resulting in a low acceptance probability regardless of the way the backwards auxiliary variables are generated. In such cases, one would like to reject directly without the need for performing the backward optimization. Such a scheme can be constructed by the use of the delayed acceptance procedure (Christen and Fox, 2005; Banterle et al., 2015). We then have:

Theorem 3. Assume $\gamma \sim \pi(\cdot)$ and assume γ^* is generated according to either Algorithm 1 or Algorithm 2. Accept γ^* if both

1. γ^* is preliminary accepted with a probability $\min\{1, \frac{\pi(\gamma^*)}{\pi(\gamma)}\}$
2. and is finally accepted with a probability $\min\{1, \frac{q_r(\gamma|\mathbf{x}_k)}{q_r(\gamma^*|\mathbf{x}_k^*)}\}$.

Then also $\gamma \sim \pi(\cdot)$.

Proof. We have that

$$\alpha_{mh}^*(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) = \alpha_{mh}^1(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) \times \alpha_{mh}^2(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*)$$

where

$$\alpha_{mh}^1(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) = \frac{\pi(\gamma^*)}{\pi(\gamma)}, \quad \alpha_{mh}^2(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) = \frac{q_r(\gamma|\mathbf{x}_k)}{q_r(\gamma^*|\mathbf{x}_k^*)}$$

Since $\alpha_{mh}^j(\gamma, \gamma^*; \mathbf{x}_k, \mathbf{x}_k^*) = [\alpha_{mh}^j(\gamma^*, \gamma; \mathbf{x}_k^*, \mathbf{x}_k)]^{-1}$ for $j = 1, 2$, it follows by the general results in Banterle et al. (2015) that we obtain an invariant kernel with respect to $\bar{\pi}$. \square

In general the total acceptance rate will be smaller than without delayed acceptance (Banterle et al., 2015, remark 1), but the gain by avoiding a backwards optimization step if not accepted in the preliminary step can compensate on this.

3.4. Calculation of marginal densities

In practice exact calculation of the marginal density can only be performed in simple models such as linear Gaussian ones, so alternatives need to be considered. One approach is to use estimators that are accurate enough to neglect the approximation errors involved. Such approximative approaches have been used in various settings of Bayesian variable selection and Bayesian model averaging. Laplace’s method (Tierney and Kadane, 1986) has been widely used, but is based on rather strong assumptions. The Harmonic mean estimator (Newton and Raftery, 1994) is an easy to implement MCMC based method but can give high variability in the estimates. Chib’s method (Chib, 1995), and its extension (Chib and Jeliazkov, 2001), have gained increasing popularity and can be very accurate provided enough MCMC iterations are performed. Approximate Bayesian Computation (Marin et al., 2012) has also been considered in this context, being much faster than MCMC alternatives, but also giving cruder approximations. Variational methods (Jordan et al., 1999) provide lower bounds for the marginal likelihoods and have been used for model selection in e.g. mixture models (McGrory and Titterton, 2007). Integrated nested Laplace approximation (INLA, Rue et al., 2009) provides accurate estimates of marginal likelihoods within the class of latent Gaussian models. In the context of generalized linear models, BIC type approximations can be used.

An alternative is to insert unbiased estimates of $\pi(\gamma)$ into (10). Andrieu and Roberts (2009) name this the *pseudo-marginal* approach and show that this leads to exact algorithms (in the sense of converging to the right distribution). Importance sampling (Beaumont, 2003) and particle filter (Andrieu et al., 2010) are two approaches that can be used within this setting. In general, the convergence rate will depend on the amount of Monte Carlo effort that is applied. Doucet et al. (2015) provide some guidelines.

Our implementation of the MJMCMC algorithm allows for all of the available possibilities for calculation of marginal likelihoods and assumes that the approximation error can be neglected. For the experiments in section 4 we have applied exact evaluations in the case of linear Gaussian models, approximations based on the assumed informative priors in case of generalized linear models (Clyde et al., 2011), and INLA (Rue et al., 2009) in the case of latent Gaussian models. Bivand et al. (2015) also apply INLA within an MCMC setting, but then concentrating on hyperparameters that (currently) can not be estimated within the INLA framework. Friel and Wyse (2012) performed comparison of some of the mentioned approaches for calculation of marginal

likelihoods, including Laplace’s approximations, harmonic mean approximations, Chib’s method and others. [Hubin and Storvik \(2016\)](#) reported some comparisons of INLA and other methods for approximating marginal likelihood. There it is demonstrated that INLA provides extremely accurate approximations on marginal likelihoods in a fraction of time compared to Monte Carlo based methods. [Hubin and Storvik \(2016\)](#) also demonstrated that by means of adjusting tuning parameters within the algorithm (the grid size and threshold values within the numerical integration procedure, [Rue et al., 2009](#)) one can often make the difference between INLA and unbiased methods of estimating of the marginal likelihood arbitrary small.

3.5. Parallelization and tuning parameters of the search

With large number of potential explanatory variables it is important to be able to utilize multiple cores and GPUs of either local machines or clusters in parallel. General principles of utilizing multiple cores in local optimization are provided in [Eksioglu et al. \(2002\)](#). At every step of the local optimization within the large jump steps one can simultaneously draw several proposals with respect to a certain transition kernel during the optimization procedure and then sequentially calculate the transition probabilities as the proposed models are evaluated by the corresponding CPUs, GPUs or clusters in the order they are returned. In those iterations where no large jumps are performed, we are utilizing multiple cores by means of addressing multiple try MCMC to explore the solutions around the current mode. The parallelization strategies are described in detail in [Appendix A](#).

In practice, tuning parameters of the local optimization routines such as the choice of the neighborhood, generation of proposals within it, the cooling schedule for *simulated annealing* ([Michiels et al., 2010](#)) or number of steps in greedy optimization also become crucially important and it yet remains unclear whether we can optimally tune them before or during the search. Random selection of proposals for from [Table 1](#) and of optimizer at each iteration is also possible. Tuning the probabilities of addressing these different options can be beneficial. Such tuning is a sophisticated mathematical problem, which we are not trying to resolve optimally within this paper, however we suggest a simple practical idea for obtaining reasonable solutions. Within the BAS algorithm, an important feature was to utilize the marginal inclusion probabilities of different covariates. We have introduced this in our algorithms as well by allowing insertion of estimates of the ρ_i ’s in proposals given in [Table 1](#) based on some burn-in period. They then correspond to the

marginal inclusion probabilities after burn-in shifted with some small ϵ from 0 and 1 if necessary in order to guarantee irreducibility. Additional literature review on search parameter tuning can be found in [Luo \(2016\)](#).

4. Experiments

In this section we are going to apply the MJMCMC algorithm to different data sets and analyze the results in relation to other algorithms. Linear regression is addressed through the U.S. Crime Data ([Raftery et al., 1997](#)) and a protein activity data ([Clyde et al., 1998](#)). Logistic regression is considered in a simulated example based on a data set and through Arabidopsis epigenetic data. The Arabidopsis example also includes random effects.

We compare the performance of our approach to competing MCMC methods such as MCMC model composition (MC³, [Madigan et al., 1995](#); [Raftery et al., 1997](#)) and the random-swap (RS) algorithm ([Clyde et al., 2011](#)) as well as the BAS algorithm ([Clyde et al., 2011](#)). Both MC³ and RS are simple MCMC procedures based on the standard Metropolis-Hastings algorithm with proposals chosen correspondingly as an inversion or a random change of one coordinate in γ at a time ([Clyde et al., 2011](#)). BAS carries out sampling without repetition from the space of models with respect to the adaptively updated marginal inclusion probabilities. For one of the examples, also a comparison with the ESS++ software (evolutionary stochastic search [Bottolo et al., 2011](#)) is made. For the cases when full enumeration of the model space is possible we additionally compare all of the aforementioned approaches to the benchmark TOP method that consists of the best quantile of models in terms of the posterior probability for the corresponding number of addressed models $\|\mathbb{V}\|$ and can not by any chance be outperformed in terms of the posterior mass captured.

Following [Clyde et al. \(2011\)](#), approximations for model probabilities (7) and marginal inclusion probabilities (9) based on a subspace of models are further referred to as RM (renormalized) approximations, whilst the corresponding MCMC based approximations (8) are referred to as MC approximations. The validation criteria addressed include root mean squared errors and bias of parameters of interest based on multiple replications of each algorithm, similar to [Clyde et al. \(2011\)](#). In addition to marginal inclusion

probabilities, we also include a global measure

$$C(\gamma) = \frac{\sum_{\gamma' \in \mathbb{V}} p(\mathbf{y}|\gamma')p(\gamma')}{\sum_{\gamma' \in \Omega} p(\mathbf{y}|\gamma')p(\gamma')}, \quad (13)$$

describing the fraction of probability mass contained in the subspace \mathbb{V} . This measure allows to address how well the search worked in terms of capturing posterior mass within a given model space. By formula (7) maximization of $C(\gamma)$ automatically induces minimization of the bias in terms of posterior marginal model probabilities, which vanishes gradually as long as $C(\gamma) \rightarrow 1$.

Mixtures of different proposals from Table 1 and local optimizers mentioned in section 3.2 were used in the studied examples in MJMCMC algorithm. A validation of the gain in using such mixtures is given in example 4.1, where we address both MJMCMC with mixtures and without them. The details on the choices and frequencies of different proposals are given in Tables B.7, B.8, B.9, and B.10 in Appendix B. The choice is based on some tuning based on the simulated data example, reported in the . Further arbitrary adaptations for the size of the problem were made in some of the examples. Generally speaking, we can not claim that the choice of the tuning parameters is optimal. It is rather some arbitrary and subjectively rational choice. Settings of the combinatorial optimizers for the addressed examples are reported in Appendix B, Table B.6.

4.1. Example 1

We address a real U.S. Crime data set, first introduced by Vandaele (1978) and stated to be a test bed for evaluation of methods for model selection (Raftery et al., 1997). The data set consists of 47 observations on 15 covariates and the responses, which are the corresponding crime rates. We will compare performance of the algorithms based on a linear Bayesian regression model with a Zellner’s g-prior, $g = 47$, making the marginal likelihood to be of the following form:

$$p(\mathbf{y}|\gamma) \propto (1 + g)^{(T-p-1)/2} (1 + g[1 - R_\gamma^2])^{-(T-1)/2}, \quad (14)$$

where R_γ^2 is the usual coefficient of determination of a linear regression model; with this scaling, the marginal likelihood of the null model (the model containing no covariates) is 1.0.

This is a sophisticated example with several local modes, which results in that all simple MCMC methods easily get stuck and have extremely poor performances in terms of the captured mass and precision of both the marginal

posterior inclusion probabilities and the posterior model probabilities. Table 3 shows the RMSE (scaled by 10^2) for different parameters over 100 repeated runs for each algorithm. The True column contains the true marginal inclusion probabilities (obtained from full enumeration) while the TOP column shows the RMSE results based on the 3276 models with highest posterior probabilities. The MJMCMC columns show the results based on using mixtures of proposals and optimizers while the MJMCMC* results are based on one choice of proposal with swaps of only 2 components at a time and a greedy optimizer.

For this example MJMCMC gives a much better performance than other MCMC methods in terms of both MC and RM based estimations with respect to the posterior mass captured ($C(\gamma)$). Using 3276 iterations, BAS slightly outperforms MJMCMC. However, when running MJMCMC so that the number of *unique* models visited ($\|\mathbb{V}\|$) are comparable with BAS, MJMCMC gives better results. compared to BAS in terms of posterior mass captured, biases and root mean squared errors for both posterior model probabilities and marginal inclusion probabilities (Table 3).

BAS has the property of never revisiting the same solutions, whilst all MCMC based procedures tend to do such revisiting with respect to the corresponding posterior probabilities. When generating a proposal is much cheaper than estimation of the model (which is usually the case, also in this example) and we are storing the results for the already estimated models, having generated a bit more models by MJMCMC does not seem to be a serious issue. Those unique models that are visited have a higher posterior mass than those suggested by BAS (for the same number of models visited). Furthermore MJMCMC (like BAS) can escape from local modes.

Also the results based on no mixture of proposals (MJMCMC* in the table) gives much better results than standard MCMC methods, however the results obtained by the MJMCMC algorithm with a mixture of proposals were even better. We have tested this on some other examples too and the use of mixtures was always beneficial and thus recommended. For this reason in other experiments only the cases with mixtures of proposals are addressed.

4.2. Example 2

In this example we are considering a new simulated data set for logistic regression. We generated $p = 20$ covariates as a mixture of binary and continuous variables. The correlation structure is shown in Figure 2 while the full details of how the data was generated is given in Appendix B.1.

Par	True	TOP	MJMCMC				BAS	MC ³			RS		MJMCMC*	
Δ	π_j	-	RM	MC	RM	MC	RM	MC	RM	MC	RM	MC	RM	MC
γ_8	0.16	3.51	6.57	10.68	5.11	10.29	5.21	6.49	3.49	5.87	3.31	6.23	9.06	
γ_{13}	0.16	3.34	7.46	10.54	5.60	10.19	6.26	8.62	3.39	8.83	3.05	6.38	10.54	
γ_{14}	0.19	3.24	8.30	12.43	6.30	12.33	6.20	6.58	2.55	6.22	2.46	7.15	10.91	
γ_{12}	0.22	3.27	6.87	13.61	5.57	13.64	3.10	5.81	6.23	4.93	5.27	5.29	10.93	
γ_5	0.23	2.56	6.30	13.45	4.59	13.65	1.84	6.07	13.05	5.13	12.77	5.39	10.90	
γ_9	0.23	3.27	9.49	16.21	7.40	16.21	9.27	5.99	2.99	5.70	2.60	7.68	11.06	
γ_7	0.29	2.31	4.37	13.63	3.45	12.73	2.28	4.74	9.61	3.46	9.70	3.91	10.10	
γ_4	0.30	1.57	6.18	19.22	3.79	17.31	0.99	13.24	21.84	13.53	21.48	4.63	13.22	
γ_6	0.33	1.92	8.61	19.71	6.14	19.49	3.11	10.19	7.47	10.99	7.12	5.87	15.43	
γ_1	0.34	2.51	11.32	22.68	7.29	20.50	8.43	22.89	25.19	23.63	24.71	7.58	12.97	
γ_3	0.39	0.43	3.95	11.13	2.38	6.99	5.02	21.48	30.24	21.39	29.94	2.99	12.66	
γ_2	0.57	1.58	5.92	13.21	3.82	9.03	13.78	30.81	37.57	29.27	37.15	5.11	14.04	
γ_{11}	0.59	0.58	3.57	13.49	2.37	15.94	4.04	11.88	21.79	11.16	21.31	2.77	12.77	
γ_{10}	0.77	3.25	7.62	7.28	5.97	4.78	15.45	21.83	19.18	20.53	19.65	6.41	14.27	
γ_{15}	0.82	3.48	9.23	4.45	6.89	5.85	14.50	69.68	76.81	69.19	76.30	6.75	14.76	
$C(\gamma)$	1.00	0.86	0.58	0.58	0.71	0.71	0.66	0.10	0.10	0.10	0.10	0.60	0.60	
Eff	2^{15}	3276	1909	1909	3237	3237	3276	829	829	1071	1071	3264	3264	
Tot	2^{15}	3276	3276	3276	5936	5936	3276	3276	3276	3276	3276	4295	4295	

Table 3: Average root mean squared error (RMSE) over the 100 repeated runs of every algorithm on the Crime data [example 1]; the values reported in the table are $\text{RMSE} \times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$. Tot is the total number of generated proposals, while Eff is the number of unique models visited during the iterations of the algorithms (for the TOP column all 2^{15} models were visited but the RMSE are based on the best 3276 models). RM corresponds to using the renormalization procedure (7) while MC corresponds to using the MC procedure (8). The two runs of MJMCMC are based on different Eff. The corresponding biases are reported in [Appendix C](#) in Table C.11.

A total of $2^{20} = 1\,048\,576$ potential models need to be considered in this case. Additionally, in this example $n = 2000$, which makes estimation of a single model significantly slower than in the previous examples. For γ we use the Binomial prior (4) with $q = 0.057$. We are in this case using the BIC-approximation for the marginal likelihood,

$$\log \hat{p}(\mathbf{y}|\gamma) = \log \hat{p}(\mathbf{y}|\hat{\beta}_\gamma) - \frac{n}{2} \log(|\hat{\beta}_\gamma|), \quad (15)$$

where $\hat{\beta}_\gamma$ is the maximum likelihood (or MAP) estimate for the β_j 's involved and $|\hat{\beta}_\gamma|$ is the number of parameters. This choice was made in order to compare the results with implementations of BAS, RS and MC³ available in the supplementary to [Clyde et al. \(2011\)](#), where this approximation is considered. In that way, the model search procedures are compared based on the same selection criterion.

Some of the covariates involved have large correlations. This induces both multimodality of the space of models and sparsity of the locations of the modes and creates an interesting example for comparison of different

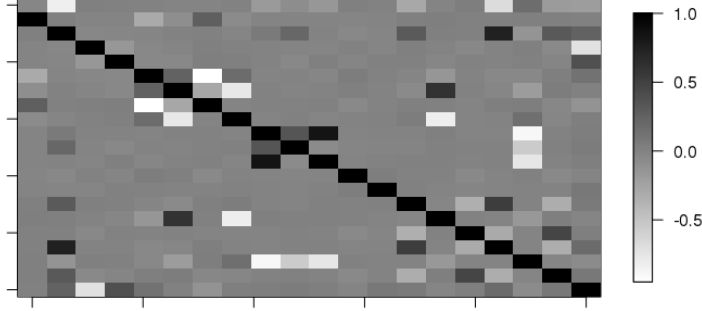


Figure 2: Correlation structure of the covariates in example 3.

Par	True	TOP	MJMCMC				BAS	MCBAS	RS	
Δ	π_j	-	RM	MC	RM	MC	RM	RM	RM	MC
γ_6	0.29	0.00	7.38	15.54	4.54	16.62	6.47	3.67	6.01	2.11
γ_8	0.31	0.00	6.23	15.50	3.96	16.94	5.58	3.02	5.37	2.55
γ_{12}	0.35	0.00	4.86	14.62	2.78	13.66	4.22	2.12	3.91	2.37
γ_{15}	0.35	0.00	4.55	15.24	2.56	15.45	4.66	1.64	3.40	2.56
γ_2	0.36	0.00	4.90	16.52	2.92	17.39	5.42	2.45	3.65	2.61
γ_{20}	0.37	0.00	4.82	14.35	2.66	14.08	3.32	1.80	4.15	2.18
γ_3	0.40	0.00	9.25	20.93	5.65	22.18	9.75	4.82	6.76	2.83
γ_{14}	0.44	0.00	3.14	17.54	1.58	16.24	3.73	1.30	1.33	2.93
γ_{10}	0.44	0.00	4.60	18.73	2.29	17.90	4.87	1.30	1.51	2.42
γ_5	0.46	0.00	3.10	17.17	1.53	16.97	4.06	1.51	1.09	2.85
γ_9	0.61	0.00	3.68	16.29	1.63	13.66	3.89	1.39	2.19	2.35
γ_4	0.88	0.00	5.66	6.70	3.74	6.26	6.60	5.57	7.61	2.15
γ_{11}	0.91	0.00	5.46	6.81	3.95	6.90	4.66	3.14	4.32	1.57
γ_1	0.97	0.00	1.90	1.74	1.35	1.34	2.43	1.96	2.30	1.1
γ_{13}	1.00	0.00	0.00	0.43	0.00	0.32	0.00	0.00	0.00	0.37
γ_7	1.00	0.00	0.00	0.57	0.00	0.41	0.00	0.00	0.00	0.33
γ_{16}	1.00	0.00	0.00	0.41	0.00	0.33	0.00	0.00	0.00	0.23
γ_{17}	1.00	0.00	0.00	0.43	0.00	0.39	0.00	0.00	0.00	0.23
γ_{18}	1.00	0.00	0.00	0.47	0.00	0.35	0.00	0.00	0.00	0.24
γ_{19}	1.00	0.00	0.00	0.52	0.00	0.36	0.00	0.00	0.00	0.41
$C(\gamma)$	1.00	1.00	0.72	0.72	0.85	0.85	0.74	0.85	0.68	0.68
Eff	2^{20}	10000	5148	5148	9988	9988	10000	10000	1889	1889
Tot	2^{20}	10000	9998	9998	19849	19849	10000	10000	10000	10000

Table 4: Average root mean squared error (RMSE) from the 100 repeated runs of every algorithm on the simulated logistic regression data (example 2); the values reported in the table are $\text{RMSE} \times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$. See the caption of Table C.14 for further details. The corresponding biases are reported in the appendix Appendix C in Table C.12.

search strategies. As one can see in Table 4, MJMCMC based on the same number of unique models outperforms both pure BAS by far in terms of posterior mass captured and root mean square errors of marginal inclusion probabilities and model probabilities. MJMCMC based on the same number of total models outperformed RS as well. The latter got stuck in some local modes and for the 10000 models visited only could search through 1889 unique models. We could not reach 10000 unique models for RS algorithm within a reasonable time for this examples (most likely the algorithm could not escape from local extrema), hence such a scenario is not reported. Even for almost two times less originally visited models in \mathbb{V} , comparing to BAS, MJMCMC gives almost the same results in terms of the posterior mass captured and errors. MJMCMC, for the given number of unique models visited, did not outperform a combination of MCMC and BAS, which is recommended by Clyde et al. (2011) for larger model spaces; both of them gave approximately identical results.

4.3. Example 3

This experiment is based on a much larger model space in comparison to all of the other examples. We address the protein activity data (Clyde et al., 1998) and consider all main effects together with the two-way interactions and quadratic terms of the continuous covariates resulting in 88 covariates in total. This corresponds to a model space of cardinality 2^{88} . This model space is additionally multimodal, which is the result of having high correlations between numerous of the addressed covariates (17 pairs of covariates have correlations above 0.95). We analyze the data set using the Bayesian linear regression with the Binomial prior (4) with $q = 0.5$ for γ and a Zellner's g-prior with $g = 96$ for β (the data has 96 observations). We then compare the performance of MJMCMC, BAS and RST. For this example we have also addressed the ESS algorithm (Bottolo et al., 2011).

The reported RST results are based on the RS algorithm run for 88×2^{20} iterations and a thinning rate of $\frac{1}{88}$. BAS was run with several choices of initial sampling probabilities such as uniformly distributed within the model space one, eplogp adjusted (Clyde et al., 2011), and those based on RM and MC approximations obtained by the RST algorithm. For the first two initial sampling probabilities BAS was run for 2^{20} iterations, whilst for the latter two first RST was run for 88×2^{19} iterations providing 2^{19} models for estimating initial sampling probabilities and then BAS was run for the other 2^{19} iterations. MJMCMC was run until 2^{20} unique models were obtained.

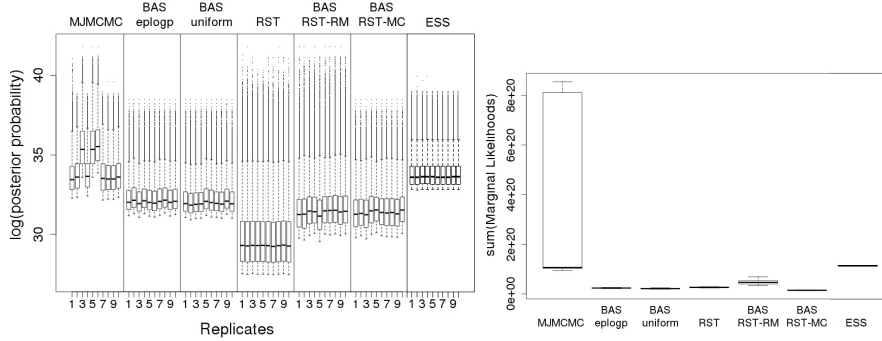


Figure 3: Comparisons of the log marginal likelihood in the protein data of the top 100000 models (left) and box-plots of the posterior mass captured (right) obtained by MJMCMC, BAS-eplogp, BAS-uniform, thinned version of Random Swap (RST), BAS with Monte Carlo estimates of inclusion probabilities from the RST samples (BAS-RST-MC), BAS re-normalized estimates of inclusion probabilities (BAS-RST-RM) from the RST samples, and ESS.

ESS was run with default search settings until 2^{20} unique models are visited. All of the algorithms were run on 10 replications.

In Figure 3 box-plots of the best 100 000 models captured by the corresponding replications of the algorithms as well as posterior masses captured by them are displayed. BAS with both uniform and eplogp initial sampling probabilities perform rather poorly in comparison to other methods, whilst BAS combined with RM approximations from RST does slightly better. ESS as well as MJMCMC show the most promising results. BAS with RM initial sampling probabilities usually manages to find models with the highest posterior probabilities, however MJMCMC in general captures by far higher posterior mass within the same amount of unique models addressed. Marginal inclusion probabilities obtained by the best run of MJMCMC with a mass (denominator of (7)) equal to 8.56×10^{20} are reported in Figure 4, whilst those obtained by other methods can be found in Clyde et al. (2011). Since MJMCMC obtains the highest posterior mass, we expect that the corresponding RM estimates of the marginal inclusion probabilities are the least biased, moreover they perfectly agree with the MC approximations. Although MJMCMC in all of the obtained replications outperformed most of the competitors in terms of the posterior mass captured, it itself exhibits significant variation between the runs (right panel of Figure 3). The latter issue can be explained by that we are only allowing visiting $3.39 \times 10^{-19}\%$ of the total model space in the addressed replications, which might be not enough to always converge to the same posterior mass captured. Note however that the variability in

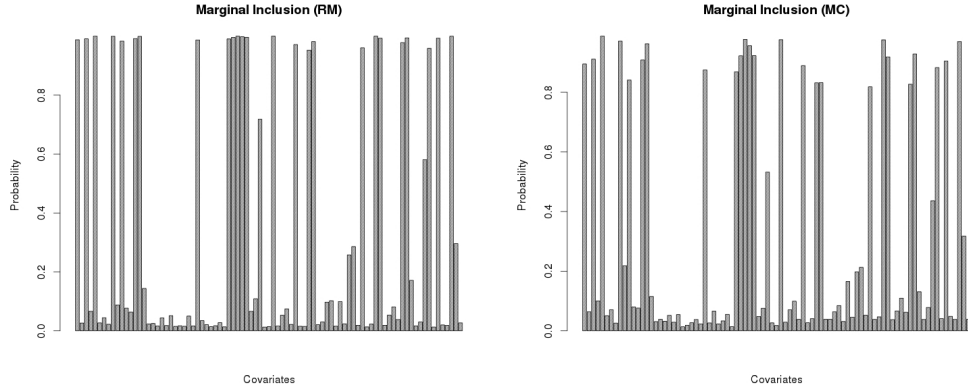


Figure 4: Comparisons of RM (left) and MC (right) estimates of marginal posterior inclusion probabilities obtained by the best run of MJMCMC with $8.56e + 20$ posterior mass captured.

the results obtained from different runs of MJMCMC clearly indicates that more iterations are needed, while the other methods may indicate (wrongly) that sufficient iterations are performed.

4.4. Example 4

In this example we illustrate how MJMCMC works for GLMM models. As illustration, we address genomic and epigenomic data on Arabidopsis. Arabidopsis is a plant model organism with a lot of genomic/epigenomic data easily available (Becker et al., 2011). At each position on the genome, a number of reads are allocated. At locations with a nucleotide of type cytosine (C), reads are either methylated or not. Our focus will be on modeling the amount of methylated reads through different covariates including (local) genomic structures, gene classes and expression levels. The studied data was obtained from the NCBI GEO archive (Barrett et al., 2013).

We model the number of methylated reads $Y_i \in \{1, \dots, R_i\}$ per loci $i \in \{1, \dots, n\}$, where $R_i \in \mathbb{N}$ is the number of reads, through (1)-(3) by a Poisson distribution for the response. Since in general the ratio of methylated bases is low, we have preferred the Poisson distribution of the responses to the binomial. The mean η_i is modeled via the log link to the chosen covariates, including an offset defined by R_i per location, and a spatially correlated random effect δ_i which is modeled via an $AR(1)$ process with parameter $\rho \in \mathbb{R}$, namely $\delta_i = \rho\delta_{i-1} + \epsilon_i \in \mathbb{R}$ with $\epsilon_i \sim N(0, \sigma_\epsilon^2)$, $i \in \{1, \dots, n\}$. Thus, we take into account spatial dependence structures of methylation rates along the genome as well as the variance of the observations not explained by the

covariates. We use the binomial prior (4) with $q = 0.5$ for γ and the Gaussian prior

$$\beta|\gamma \sim N_{p_\gamma}(\mu_{\beta_\gamma}, \Sigma_{\beta_\gamma})$$

for the regression coefficients. For the parameters within the random effects, we first reparametrize to $\psi_1 = \log \frac{1}{\sigma_{\epsilon,t}^2}(1 - \rho^2)$, $\psi_2 = \log \frac{1+\rho}{1-\rho}$ and assume

$$\psi_1 \sim \text{logGamma}(1, 5 \times 10^{-5}), \quad (16)$$

$$\psi_2 \sim N(0, 0.15^{-1}). \quad (17)$$

Marginal likelihoods for this example are calculated through the INLA package (www.r-inla.org).

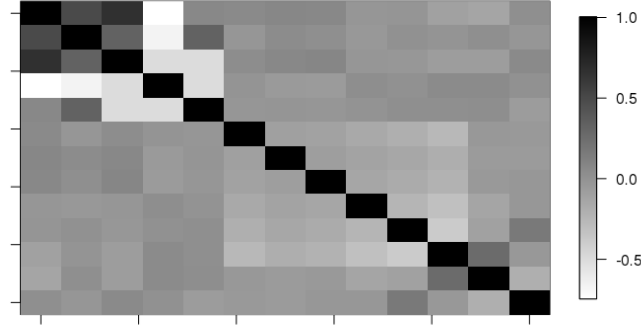


Figure 5: Correlation structure of the covariates in example 5.

We have addressed $p = 13$ different covariates in addition to the intercept. Among these covariates we address a factor with 3 levels corresponding to whether a location belongs to a CGH, CHH or CHG genetic region, where H is either A, C or T and thus generating two covariates X_1 and X_2 corresponding to whether a location is CGH or CHH. The second group of factors indicates whether a distance to the previous cytosine nucleobase (C) in DNA is 1, 2, 3, 4, 5, from 6 to 20 or greater than 20 inducing the binary covariates $X_3 - X_8$. The third factor corresponds to whether a location belongs to a gene from a particular group of genes of biological interest, these groups are indicated as M_α , M_γ , M_δ or M_0 inducing 3 additional covariates $X_9 - X_{11}$. Finally, we have two binary covariates X_{12} and X_{13} represented by expression levels exceeding 3000 and 10000, respectively. The cardinality of our search space

Par	True	TOP	MJMCMC		RS MCMC	
Δ	π_j	RM	RM	MC	RM	MC
γ_4	0.0035	0.0005	0.0022	2.0416	0.0198	1.9768
γ_6	0.0048	0.0006	0.0051	2.0899	0.0257	1.9352
γ_7	0.0065	0.0006	0.0056	2.3459	0.0353	0.6887
γ_3	0.0076	0.0007	0.0017	3.3660	0.0353	1.2374
γ_8	0.0076	0.0007	0.0079	2.3279	0.0344	1.6163
γ_5	0.0096	0.0007	0.0075	2.3342	0.0455	1.7170
γ_{11}	0.0813	0.0007	0.0200	3.6851	0.1679	2.8022
γ_{12}	0.0851	0.0006	0.0134	2.7179	0.0766	1.9136
γ_9	0.1185	0.0008	0.0184	3.3149	0.1773	3.0463
γ_{10}	0.3042	0.0006	0.0071	9.4926	0.1106	3.7344
γ_{13}	0.9827	0.0002	0.0063	2.5350	0.0638	1.5681
γ_1	1.0000	0.0007	0.0000	4.7091	0.0000	1.2258
γ_2	1.0000	0.0000	0.0000	2.7343	0.0000	0.9971
$C(\gamma)$	1.0000	1.0000	0.9998	0.9998	0.9977	0.9977
Eff	8192	385	1758	1758	155	155
Tot	8192	385	3160	3160	10000	10000

Table 5: Average root mean squared error (RMSE) from the 100 simulated runs of MJMCMC on the epigenetic data (example 4); the values reported in the table are $\text{RMSE} \times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$.

Ω is $2^{13} = 8192$ for this example. The correlation structure between these 13 covariates is represented in Figure 5.

As seen from Table 5 (TOP column) within just the 385 best unique models (2.35% of the total model space) we are able to capture almost full posterior mass for this problem. The model space, as shown in Figure 6, has very few sparsely located modes in a quite large model space. In this example we compare MJMCMC and a simple RS MCMC algorithm, the latter is allowed to only swap one component per iteration. This example does contain most of the mass in just two closely located models as can be seen in Figure 6. This is why a simple MCMC can capture essentially most of the mass after 10 000 iterations. At the same time there are a few small modes that lie a bit further from the region of the high concentration of mass, which the simple MCMC algorithm did not capture. Essentially, RS MCMC was staying within a few modes for most of the time, never being able to travel to the more remote parts of the model space and generating very few (155 on average) unique models. This number is here very low compared to the total number of models visited (10 000). If there were more

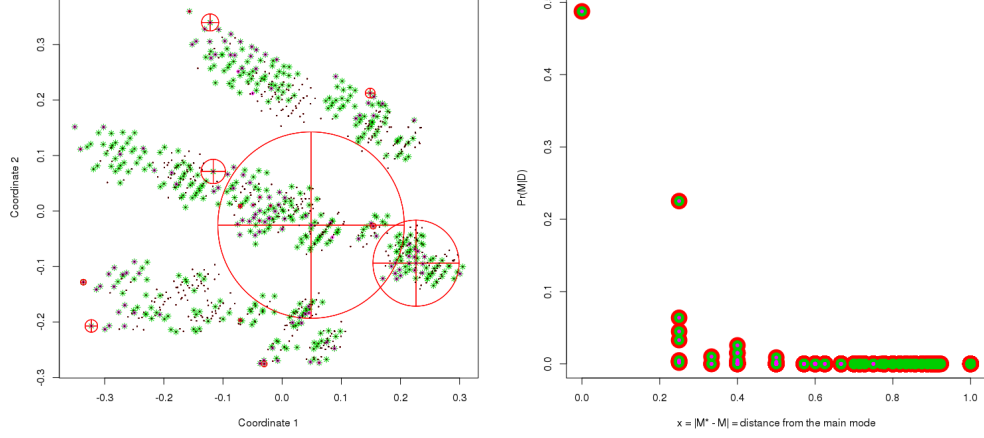


Figure 6: Left: Multidimensional scale plot (Rohde, 2002) of the best 1024 models in terms of posterior model probability in the space of models (black dots are centers of the models red circles proportional to the posterior probabilities of models, green stars - models visited by MJMCMC, purple stars - models visited by MCMC). Right: A plot of posterior probabilities with respect to distance from the global mode (red circles correspond to all the models, the green circles - models visited by MJMCMC, the purple circles - by simple MCMC).

sparsely located remote modes the simple MCMC algorithm would run into the problems similar to those discussed in the previous examples and miss a significant amount of mass. For MJMCMC, we run the algorithm until 3160 models were visited, resulting in 1758 unique models. As a result, MJMCMC captures the mass from the remote small modes, adding a bit to the captured mass, which allows it to slightly outperform the simple MCMC algorithm. As can be seen in Table 5, MJMCMC outperforms the simple MCMC algorithm in terms of the errors of marginal model probabilities. Marginal inclusion probabilities in terms of RM are also more precise when MJMCMC is used, but the MC based approximations of them are slightly better for the MCMC case.

According to marginal inclusion probabilities (π_j column in Table 5), factors of whether the location is CGH or CHH (γ_1 and γ_2) are both extremely significant, as well as the higher cut off for the level of expression (γ_{13}). Additionally factors for M_α and M_δ groups of genes have non-zero marginal inclusion probabilities and reasonably high significance. In future it would be of interest to obtain additional covariates such as whether a nucleobase belongs to a particular part of the gene like promoter or a coding region. Furthermore, it is of interest to address factors whether a base is located

within a CpG island or whether it belongs to a transposone. Moreover interactions of these covariates may be interesting. Alternative choices of the response distributions (e.g. binomial or negative binomial) and/or type of random effects ($AR(k)$, $ARMA(l, k)$) might also be of an interest.

5. Summary and discussion

In this paper we have introduced the mode jumping MCMC (MJMCMC) approach for estimating posterior model probabilities and performing Bayesian model averaging and selection. The algorithm incorporates the ideas of MCMC with the possibility of large jumps combined with local optimizers to generate proposals in the discrete space of models. Unlike standard MCMC methods applied to variable selection, the developed procedure avoids getting stuck in local modes and manages to iterate through all of the important models much faster. It also in many cases outperforms Bayesian Adaptive Sampling (BAS), having the tendency to capture a higher posterior mass within the same amount of unique models visited. This can be explained by that for problems with numerous covariates BAS requires good initial marginal inclusion probabilities to perform well. [Clyde et al. \(2011\)](#) demonstrated that estimates of marginal inclusion probabilities obtained from preliminary MCMC runs could largely improve BAS. A combination of MJMCMC with BAS could possibly improve both algorithms even further.

The *EMJMCMC* R-package is developed and currently available from the Git Hub repository: <http://aliaksah.github.io/EMJMCMC2016/>. The methodology depends on the possibility of calculating marginal likelihoods within models accurately. The developed package gives a user high flexibility in the choice of methods to obtain marginal likelihoods. Whilst the default choice for marginal likelihood calculations is based on INLA ([Rue et al., 2009](#)), we also have adopted efficient C based implementations for exact calculations in Bayesian linear regression and approximate calculations in Bayesian logistic and Poisson regressions in combination with g-priors as well as other priors. Several model selection criteria for the class of methods are also addressed. Extensive parallel computing for both MCMC moves and local optimizers is available within the developed package; in particular, with a default option a user specifies how many threads are addressed within the in-build *mclapply* function or *snow* based parallelization, however an advanced user can specify his own function to parallelize computations on

both the MCMC and local optimization levels tapping, for instance, modern graphical processing units - GPUs, which in turn allows additional efficiency and flexibility.

Whilst estimators (7) for marginal inclusion and posterior model probabilities based on Bayes formula and obtained by MJMCMC, as noticed by [Clyde et al. \(2011\)](#), are Fisher consistent, they remain generally speaking biased; although their bias reduces to zero asymptotically. MCMC based estimators such as (8), which is both consistent and unbiased, are also available through our procedure; these estimators however tend to have a much higher variance than the aforementioned ones. As one of the further developments it would be of an interest to combine knowledge available from both groups of estimators to adjust for bias and variance, which is vital for higher dimensional problems.

Another aspect that requires being discussed is the model selection criterion. Different criteria can sometimes disagree about the results of model selection. In order to avoid confusion, the researcher should be clear about the stated goals. If the goal is prediction rather than inference one should adjust for that and use AIC, WAIC ([Watanabe, 2009](#)) or DIC ([Spiegelhalter et al., 2002](#)) rather than BIC or posterior model probability as selection criterion in MJMCMC. These choices are possible within the *EMJMCMC* package as well.

Based on several experiments, we can claim MJMCMC to be a rather competitive algorithm that is addressing the wide class of Generalized Linear Mixed Models (GLMM). In particular for this class of models one can incorporate a random effect, which both models the variability unexplained by the covariates and introduces dependence between observations, creating additional modeling flexibility. Estimations of parameters of such models and Bayesian inference within them becomes significantly harder in comparison to simple GLM. This creates the necessity to address parallel computing extensively. We have enabled the latter within our package by means of combining methods for calculating marginal likelihoods, such as the INLA methodology, and parallel MJMCMC algorithm.

Mode jumping MCMC suggested by [Tjelmeland and Hegstad \(1999\)](#) was not the only idea of adapting MCMC algorithm for exploration of model spaces with multiple sparse modes. Mainly the other approaches can be divided into two groups - those, based on exploration of the tempered target distributions (allowing to flatten or increase multimodality for different temperatures) and those based on utilization of the local gradient (the orig-

inal mode jumping approach hence belongs to the second group). The first group was initialized with the parallel tempering approach (Geyer, 1991), which further had numerous modifications Liang (2010); Miasojedow et al. (2013); Salakhutdinov (2009). One of the most prominent extensions is the equi-energy sampling approach (Kou et al., 2006), which utilizes the physical duality between temperature and energy. This approach targets directly the former to flatten or spike the parameter spaces. Another extension is the multi domain sampling approach (Zhou, 2011), which first uses the target distribution tempering idea to find the set of local modes and then uses local MCMC to explore the regions around them for further global inference. The second group of the algorithms includes such methods as (Neal et al., 2011; Chen et al., 2014; Sengupta et al., 2016) and many others. Like the original mode jumping MCMC (Tjelmeland and Hegstad, 1999) they are using local gradients of the target distribution to move between local extrema. Both groups of the algorithms were mainly developed for exploration of continuous parameter spaces. Hence it is not possible to directly compare them with the proposed in this paper approach. But generally speaking temperature (energy) based methods are easier to adapt to discrete parameter spaces, where there exists no gradient in a classical sense. For example ESS algorithm (Bottolo et al., 2011) addressed in the protein activity data examples relies upon this idea. The gradient based methods are slightly trickier, since instead of utilizing the gradient local combinatorial optimization has to be performed in the discrete parameter spaces. Interestingly in our approach we combine both of the group of methods, since generally the approach is based on mode jumping MCMC idea, however for local combinatorial optimization we use (among others) simulated annealing algorithm, which heavily utilizes the tempering (or rather annealing) idea based on the Boltzmann distribution during the search. Anyways in future it would be of a particular interest to adapt other mentioned above approaches to the model selection problem and compare them to the suggested MJMCMC algorithm.

Currently we use decision variables only on the level of choice of covariates, however the mode jumping procedure can be easily extended to more general cases. In future it would be of interest to extend the procedure to model selection and model averaging jointly across covariates, link functions, random effect structures and response distributions. Such extensions will require even more accurate tuning of control parameters of the algorithm introducing another important direction for further research.

ACKNOWLEDGMENTS

The authors gratefully acknowledge the *CELS project at the University of Oslo*, <http://www.mn.uio.no/math/english/research/groups/cels/index.html>, for giving the opportunity, inspiration and motivation to write this paper. The authors also thank the editor, the associate editor, and the two referees for helpful comments and suggestions which significantly improved the manuscript.

References

- Al-Awadhi, F., Hurn, M. and Jennison, C. (2004), ‘Improving the acceptance rate of reversible jump MCMC proposals’, *Statistics and Probability Letters* **69**(2), 189 – 198.
- Andrieu, C., Doucet, A. and Holenstein, R. (2010), ‘Particle Markov chain Monte Carlo methods’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **72**(3), 269–342.
- Andrieu, C. and Roberts, G. O. (2009), ‘The pseudo-marginal approach for efficient Monte Carlo computations’, *The Annals of Statistics* (2), 697–725.
- Banterle, M., Grazian, C., Lee, A. and Robert, C. P. (2015), ‘Accelerating Metropolis-Hastings algorithms by delayed acceptance’, *arXiv preprint arXiv:1503.00996*.
- Barrett, T., Wilhite, S. E., Ledoux, P., Evangelista, C., Kim, I. F., Tomashevsky, M., Marshall, K. A., Phillippy, K. H., Sherman, P. M., Holko, M. et al. (2013), ‘NCBI GEO: archive for functional genomics data setsupdate’, *Nucleic acids research* **41**(D1), D991–D995.
- Beaumont, M. A. (2003), ‘Estimation of population growth or decline in genetically monitored populations’, *Genetics* **164**(3), 1139–1160.
- Becker, C., Hagemann, J., Müller, J., Koenig, D., Stegle, O., Borgwardt, K. and Weigel, D. (2011), ‘Spontaneous epigenetic variation in the Arabidopsis thaliana methylome’, *Nature* **480**(7376), 245–249.
- Bivand, R., Gómez-Rubio, V., Rue, H. et al. (2015), ‘Spatial Data Analysis with R-INLA with Some Extensions’, *Journal of Statistical Software* **63**(i20).
- Bivand, R. S., Gómez-Rubio, V. and Rue, H. (2014), ‘Approximate Bayesian inference for spatial econometrics models’, *Spatial Statistics* **9**, 146–165.
- Blum, C. and Roli, A. (2003), ‘Metaheuristics in combinatorial optimization: Overview and conceptual comparison’, *ACM Computing Surveys (CSUR)* **35**(3), 268–308.
- Bottolo, L., Chadeau-Hyam, M., Hastie, D. I., Langley, S. R., Petretto, E., Turet, L., Tregouet, D. and Richardson, S. (2011), ‘ESS++: a C++ objected-oriented algorithm for Bayesian stochastic search model exploration’, *Bioinformatics* **27**(4), 587–588.

- Bové, D. S. and Held, L. (2011), ‘Bayesian fractional polynomials’, *Statistics and Computing* **21**(3), 309–324.
- Chen, T., Fox, E. and Guestrin, C. (2014), Stochastic gradient hamiltonian monte carlo, in ‘International Conference on Machine Learning’, pp. 1683–1691.
- Chib, S. (1995), ‘Marginal likelihood from the Gibbs output’, *Journal of the American Statistical Association* **90**(432), 1313–1321.
- Chib, S. and Jeliazkov, I. (2001), ‘Marginal likelihood from the Metropolis–Hastings output’, *Journal of the American Statistical Association* **96**(453), 270–281.
- Chopin, N., Jacob, P. E. and Papaspiliopoulos, O. (2013), ‘SMC2: an efficient algorithm for sequential analysis of state space models’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **75**(3), 397–426.
- Christen, J. A. and Fox, C. (2005), ‘Markov chain Monte Carlo using an approximation’, *Journal of Computational and Graphical statistics* **14**(4), 795–810.
- Clyde, M. A., Ghosh, J. and Littman, M. L. (2011), ‘Bayesian adaptive sampling for variable selection and model averaging’, *Journal of Computational and Graphical Statistics* **20**(1), 80–101.
- Clyde, M., Parmigiani, G. and Vidakovic, B. (1998), ‘Multiple shrinkage and subset selection in wavelets’, *Biometrika* **85**(2), 391–401.
- David, M. (2015), ‘Auto insurance premium calculation using generalized linear models’, *Procedia Economics and Finance* **20**, 147 – 156.
- de Souza, R., Cameron, E., Killedar, M., Hilbe, J., Vilalta, R., Maio, U., Biffi, V., Ciardi, B. and Riggs, J. (2015), ‘The overlooked potential of generalized linear models in astronomy, i: Binomial regression’, *Astronomy and Computing* **12**, 21 – 32.
- Doucet, A., Pitt, M. K., Deligiannidis, G. and Kohn, R. (2015), ‘Efficient implementation of Markov chain Monte Carlo when using an unbiased likelihood estimator’, *Biometrika* **102**(2), 295.
- Eksioglu, S. D., Pardalos, P. M. and Resende, M. G. (2002), Parallel metaheuristics for combinatorial optimization, in R. Corra, I. Dutra, M. Fiallos and F. Gomes, eds, ‘Models for Parallel and Distributed Computation’, Vol. 67 of *Applied Optimization*, Springer US, pp. 179–206.
- Friel, N. and Wyse, J. (2012), ‘Estimating the evidence a review’, *Statistica Neerlandica* **66**(3), 288–308.
- Frommlet, F., Ljubic, I., Arnardttr Helga, B. and Bogdan, M. (2012), ‘QTL Mapping Using a Memetic Algorithm with Modifications of BIC as Fitness Function’, *Statistical Applications in Genetics and Molecular Biology* **11**(4), 1–26.

- George, E. I. and McCulloch, R. E. (1997), ‘Approaches for Bayesian variable selection’, *Statistica Sinica* **7**(2), 339–373.
- Geyer, C. J. (1991), ‘Markov chain monte carlo maximum likelihood’.
- Ghosh, J. (2015), ‘Bayesian model selection using the median probability model’, *Wiley Interdisciplinary Reviews: Computational Statistics* **7**(3), 185–193.
- Grossi, L. and Bellini, T. (2006), Credit risk management through robust generalized linear models, in S. Zani, A. Cerioli, M. Riani and M. Vichi, eds, ‘Data Analysis, Classification and the Forward Search’, Studies in Classification, Data Analysis, and Knowledge Organization, Springer Berlin Heidelberg, pp. 377–386.
- Hubin, A. and Storvik, G. (2016), ‘Estimating the marginal likelihood with Integrated nested Laplace approximation (INLA)’. arXiv:1611.01450v1.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S. and Saul, L. K. (1999), ‘An introduction to variational methods for graphical models’, *Machine learning* **37**(2), 183–233.
- Kou, S., Zhou, Q. and Wong, W. H. (2006), ‘Discussion paper equi-energy sampler with applications in statistical inference and statistical mechanics’, *The annals of Statistics* pp. 1581–1619.
- Liang, F. (2010), ‘A double metropolis–hastings sampler for spatial models with intractable normalizing constants’, *Journal of Statistical Computation and Simulation* **80**(9), 1007–1022.
- Liu, J. S., Liang, F. and Wong, W. H. (2000), ‘The multiple-try method and local optimization in Metropolis sampling’, *Journal of the American Statistical Association* **95**(449), 121–134.
- Lobraux, S. and Melodelima, C. (2015), ‘Detection of genomic loci associated with environmental variables using generalized linear mixed models’, *Genomics* **105**(2), 69 – 75.
- Luo, G. (2016), ‘A review of automatic selection methods for machine learning algorithms and hyper-parameter values’, *Network Modeling Analysis in Health Informatics and Bioinformatics* **5**(1), 1–16.
- Madigan, D., York, J. and Allard, D. (1995), ‘Bayesian graphical models for discrete data’, *International Statistical Review/Revue Internationale de Statistique* pp. 215–232.
- Marin, J.-M., Pudlo, P., Robert, C. P. and Ryder, R. J. (2012), ‘Approximate Bayesian computational methods’, *Statistics and Computing* **22**(6), 1167–1180.
- McGrory, C. A. and Titterton, D. (2007), ‘Variational approximations in Bayesian model selection for finite mixture distributions’, *Computational Statistics & Data Analysis* **51**(11), 5352–5367.

- Miasojedow, B., Moulines, E. and Vihola, M. (2013), ‘An adaptive parallel tempering algorithm’, *Journal of Computational and Graphical Statistics* **22**(3), 649–664.
- Michiels, W., Aarts, E. and Korst, J. (2010), *Theoretical Aspects of Local Search*, 1st edn, Springer Publishing Company, Incorporated.
- Neal, R. M. et al. (2011), ‘Mcmc using hamiltonian dynamics’, *Handbook of Markov Chain Monte Carlo* **2**(11).
- Newton, M. A. and Raftery, A. E. (1994), ‘Approximate Bayesian inference with the weighted likelihood bootstrap’, *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 3–48.
- Raftery, A. E., Madigan, D. and Hoeting, J. A. (1997), ‘Bayesian model averaging for linear regression models’, *Journal of the American Statistical Association* **92**(437), 179–191.
- Robert, C. P. and Casella, G. (2005), *Monte Carlo Statistical Methods*, Springer Texts in Statistics, Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- Rohde, D. L. T. (2002), ‘Methods for binary multidimensional scaling’, *Neural Comput.* **14**(5), 1195–1232.
- Rue, H., Martino, S. and Chopin, N. (2009), ‘Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations’, *Journal of the Royal Statistical Society* **71**(2), 319–392.
- Salakhutdinov, R. R. (2009), Learning in markov random fields using tempered transitions, in ‘Advances in neural information processing systems’, pp. 1598–1606.
- Sengupta, B., Friston, K. J. and Penny, W. D. (2016), ‘Gradient-based mcmc samplers for dynamic causal modelling’, *NeuroImage* **125**, 1107–1118.
- Skrondal, A. and Rabe-Hesketh, S. (2003), ‘Some applications of generalized linear latent and mixed models in epidemiology: repeated measures, measurement error and multilevel modeling’, *Norwegian Journal of Epidemiology* **13**(2), 265–278.
- Song, Q. and Liang, F. (2015), ‘A split-and-merge Bayesian variable selection approach for ultrahigh dimensional regression’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **77**(5), 947–972.
- Spiegelhalter, D. J., Best, N. G., Carlin, B. P. and Van Der Linde, A. (2002), ‘Bayesian measures of model complexity and fit’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **64**(4), 583–639.
- Storvik, G. (2011), ‘On the flexibility of Metropolis-Hastings acceptance probabilities in auxiliary variable proposal generation’, *Scandinavian Journal of Statistics* **38**, 342–358.

- Stroup, W. W. (2013), *Generalized linear mixed models : modern concepts, methods and applications*, CRC Press, Taylor & Francis, Boca Raton.
- Tierney, L. (1996), ‘Introduction to general state-space Markov chain theory’, *Markov chain Monte Carlo in practice* pp. 59–74.
- Tierney, L. and Kadane, J. B. (1986), ‘Accurate approximations for posterior moments and marginal densities’, *Journal of the american statistical association* **81**(393), 82–86.
- Tjelmeland, H. and Hegstad, B. K. (1999), ‘Mode jumping proposals in MCMC’, *Scandinavian journal of statistics* **28**, 205–223.
- Vandaele, W. (1978), ‘Participation in illegitimate activities: Ehrlich revisited’, *Deterrence and Incapacitation* **1**, 270–335.
- Watanabe, S. (2009), ‘An introduction to algebraic geometry and statistical learning theory’.
- Yeh, Y. T., Yang, L., Watson, M., Goodman, N. and Hanrahan, P. (2012), ‘Synthesizing open worlds with constraints using locally annealed reversible jump MCMC’, *ACM Transactions on Graphics* **31**(4), 56–58.
- Zellner, A. (1986), ‘On assessing prior distributions and Bayesian regression analysis with g-prior distributions’, *Bayesian inference and decision techniques: Essays in Honor of Bruno De Finetti* **6**, 233–243.
- Zhou, Q. (2011), ‘Multi-domain sampling with applications to structural inference of bayesian networks’, *Journal of the American Statistical Association* **106**(496), 1317–1330.

Appendix A. Details of the MJMCMC algorithm

Appendix A.1. Multiple try MCMC algorithm

In addition to ordinary MCMC steps and mode jump MCMC, also multiple-try Metropolis (Liu et al., 2000) is considered. Multiple-try Metropolis is a sampling method that is a modified form of the Metropolis-Hastings method, designed to be able to properly parallelize the original Metropolis-Hastings algorithm. The idea of the method is to allow generating S trial proposals $\chi_1^*, \dots, \chi_S^*$ in parallel from a proposal distribution $q(\cdot|\gamma)$. Then, $\gamma^* \in \{\chi_1^*, \dots, \chi_S^*\}$ is selected with probabilities proportional to some importance weights $w(\gamma, \chi_i^*) = \pi(\gamma)q(\chi_i^*|\gamma)\lambda(\chi_i^*, \gamma)$ where $\lambda(\chi_i^*, \gamma) = \lambda(\gamma, \chi_i^*)$. In the reversed move $\chi_1, \dots, \chi_{S-1}$ are generated from the proposal $q(\chi|\gamma^*)$ while $\chi_S = \gamma$. Finally, the move is accepted with probability

$$r_m(\gamma, \gamma^*) = \min \left\{ 1, \frac{w(\chi_1^*, \gamma) + \dots + w(\chi_S^*, \gamma)}{w(\chi_1, \gamma^*) + \dots + w(\chi_S, \gamma^*)} \right\}. \quad (\text{A.1})$$

In the implementation of the algorithm, ordinary MCMC is considered as a special case of multiple try MCMC with $S = 1$. We recommend ordinary or multiple try MCMC steps are used in at least 95% of the iterations with proposals of large jumps for the remaining 5%.

Appendix A.2. Choice of proposal distributions

The implementation of MJMCMC allows for great flexibility in the choices of proposal distributions for the large jumps, the local optimization and the last randomization.

- Table 1 lists the current possibilities for drawing indexes to swap in the first large jump. One would choose distributions where a large number of components are swapped.
- An important ingredient of the MJMCMC algorithm is the choice of local optimizer. In the current implementation of the algorithm, several choices are possible; simulated annealing, greedy optimizers based on best neighbor optimization or first improving neighbor (Blum and Roli, 2003) which is another variant of greedy local search accepting the first randomly selected solution better than the current. For each alternative the neighbors are defined through swapping a few of the γ_j 's in the current model.
- For the last randomization, again Table 1 lists the possibilities, but in this case a small number of swaps will be preferable.

Different possibilities to combine the optimizers and proposals in a hybrid setting are also possible. Then, at each iteration, which proposal distributions and which optimizer to use are randomly drawn from the set of possibilities, see Robert and Casella (2005, sec 10.3) for the validity of such procedures.

Appendix A.3. Parallel computing in local optimizers

General principles of utilizing multiple cores in local optimization are provided in [Eksioglu et al. \(2002\)](#). Given a current state χ^* , one can simultaneously draw several proposals χ_1, \dots, χ_K with respect to a certain transition kernel $s_o(\cdot|\gamma)$ and, if necessary, calculate the transition probabilities as the proposed models are evaluated. This step can be performed by parallel CPUs, GPUs or clusters. Consider an optimizer with the acceptance probability function $r_o^t(\chi_j; \chi^*)$, $j \in 1, \dots, K$, which either changes over the time (iterations) t or remains unchanged. For the greedy local search $r_o^t(\chi; \chi^*) = \mathbb{1}\{\pi(\chi) \geq \pi(\chi^*)\}$, $t \in 1, 2, \dots$. For the implemented version of the simulated annealing algorithm we consider $r_o^t(\chi; \chi^*) = \min\left\{1, \exp\left(\frac{\log \pi(\chi) - \log \pi(\chi^*)}{T_t}\right)\right\}$, $i \in 1, \dots, N$, where T_t is the SA temperature ([Blum and Roli, 2003](#)) parameter at iteration t . The proposed parallelization strategy is given in detail in Algorithm 3.

Algorithm 3 Parallel optimization

```

1: procedure OPTIMIZE(N)
2:    $\chi^* \leftarrow \chi_0$ 
3:   for  $i = 1, \dots, N$  do
4:      $\chi_{i,1}, \dots, \chi_{i,K} \sim s_o(\cdot|\chi^*)$  ▷ make  $K$  proposals in parallel
5:     for  $j = 1, \dots, K$  do
6:        $r \leftarrow r_o^i(\chi_{i,j}; \chi^*)$  ▷ calculate acceptance probability
7:       if  $\text{Unif}[0; 1] \leq r$  then
8:          $\chi^* \leftarrow \chi_{i,j}$  ▷ accept the transition
9:       end if
10:    end for
11:     $\chi_i^* \leftarrow \chi^*$ 
12:  end for
13:  return  $\chi_N^*$ 
14: end procedure

```

Appendix A.4. Parallel MJMCMC with a mixture of proposals

Here we described the full version of our algorithm based on a combination of Algorithm 2 and the multiple try idea. The suggested MJMCMC approach allows to both mix between local modes efficiently and explore the solutions around the modes simultaneously whilst keeping the desired ergodicity of the MJMCMC procedure. This implementation allows for the mixtures of both local optimizers and proposals addressed within MJMCMC. Both the local optimization and the multiple try steps utilize multiple CPUs and GPUs of a single machine or a cluster of nodes. The pseudo-code of the algorithm is given in Algorithm 4 below. In this pseudo-code we consider the following notation:

- ϱ - the probability for a large jump;

- $P_o(\cdot)$ - the distribution for the choice of the local optimizers, a discrete distribution over a finite number of possibilities;
- $P_l(\cdot)$ - the distribution for the choice of large jump transition kernel, a discrete distribution over the possibilities in Table 1 with high probabilities on a large number of swaps;
- $P_r(\cdot)$ - the distribution for the choice of the randomizing kernel, a discrete distribution over a finite number of possibilities;
- $P_g(\cdot)$ - the distribution for the choice of the multiple try MCMC, a discrete distribution over the possibilities in Table 1 proposals with a high probability on a small number of swaps.

Algorithm 4 Mode jumping MCMC

```

1: procedure MJMCMC( $Numit$ )
2:    $\gamma \leftarrow \gamma_0$  ▷ define the initial state
3:   for  $t = 1, \dots, Numit$  do
4:     if  $\text{Unif}[0; 1] \leq \varrho$  then ▷ large jump with local optimization
5:        $q_l \sim P_l(\cdot)$  ▷ choose large jump kernel
6:        $q_o \sim P_o(\cdot)$  ▷ choose local optimizer
7:        $q_r \sim P_r(\cdot)$  ▷ choose randomization kernel
8:        $I \sim q_l(\cdot | \gamma)$  ▷ Indices for large jump
9:        $\chi_0^* \leftarrow \text{SWAP}(\gamma, I)$  ▷ large jump
10:       $\chi_k^* \sim q_o(\cdot | \chi_0^*)$  ▷ local optimization
11:       $\gamma^* \sim q_r(\cdot | \chi_k^*)$  ▷ randomization around the mode
12:       $\chi_0 \leftarrow \text{SWAP}(\gamma^*, I)$  ▷ reverse large jump
13:       $\chi_k \sim q_o(\cdot | \chi_0)$  ▷ local optimization
14:       $r \leftarrow r_m(\chi, \gamma; \chi^*, \gamma^*)$  ▷ from (12)
15:    else ▷ ordinary proposal
16:       $q_g \sim P_g(\cdot)$  ▷ choose multiple try proposal kernel
17:       $\gamma^* \sim q_g(\cdot | \gamma)$  ▷ proposed solution
18:       $r \leftarrow r_m(\gamma, \gamma^*)$  ▷ from (A.1)
19:    end if
20:    if  $\text{Unif}[0; 1] \leq r$  then
21:       $\gamma \leftarrow \gamma^*$  ▷ accept the move
22:    end if
23:  end for
24: end procedure

```

The essential ingredients of the parallel version of the MJMCMC with a mixture of proposals (Algorithm 4) are as follows:

- Multiple try MCMC steps are performed for the steps with no mode jumps;
- At the iterations with mode jumps the large jump proposals $q_l \sim P_l(\zeta)$, the optimization proposals $q_o \sim P_o(\zeta)$, and the randomizing kernels $q_r \sim P_r(\zeta)$ are chosen randomly;
- At the iterations with no mode jumps the proposal is chosen randomly as $q_g \sim P_g(\zeta)$;
- The optimization steps are parallelized as described in [Appendix A.3](#).
- The multiple-try steps are parallelized.

Appendix B. Supplementary materials for the experiments

Table B.6 describes some of the tuning parameters used for the different examples. The remaining tuning parameters, describing the mixture distributions P_o, P_l and P_r are specified in tables B.7 (example 1), B.8 (example 2), B.9 (example 3) and B.10 (example 4).

Example	CPU	SA				Greedy			MTMCMC	
No	Num	S_t	Δt	t_0	t_f	S	LS	FI	Size	Steps
1	4	4	3	10	14×10^{-5}	15	F	T	4	15
2	2	5	3	10	14×10^{-5}	20	F	T	2	20
3	10	18	3	10	14×10^{-5}	88	F	T	10	88
4	1	3	3	10	14×10^{-5}	13	F	T	2	13
S.1	4	4	3	10	14×10^{-5}	15	F	T	4	15

Table B.6: Tuning parameters of the blocks of MJMCMC in the examples (Example No); CPU (Num) - the number of CPUs utilized within the examples; S_t - number of iterations per temperature in SA algorithm; Δt - cooling factor of the cooling schedule of SA algorithm; t_0 - initial temperature of SA algorithm; t_f - final temperature of SA algorithm; S - number of iterations in Greedy algorithm (per run); LS - if local stop is allowed in Greedy algorithm; FI - if the first improving neighbor strategy is applied in Greedy algorithm; Size - number of proposals per step in MTMCMC algorithm; Steps - number of MTMCMC steps (only makes sense when MTMCMC is used as an optimizer).

Proposal	Optimizer	Frequency	Type 1	Type 4	Type 3	Type 5	Type 6	Type 2
q_g	-	$\varrho = 0.9836$	0.1176	0.3348	0.2772	0.0199	0.2453	0.0042
S	-	-	$\{2, 2\}$	2	$\{2, 2\}$	1	1	15
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_l	-	0.0164	0	1	0	0	0	0
S	-	-	-	4	-	-	-	-
ρ_i	-	-	-	-	-	-	-	-
q_o	SA	0.5553	0.0788	0.3942	0.1908	0.1928	0.1385	0.0040
q_o	GREEDY	0.2404	0.0190	0.3661	0.2111	0.2935	0.1046	0.0044
q_o	MTMCMC	0.2043	0.2866	0.1305	0.2329	0.1369	0.2087	0.0040
S	-	-	$\{2, 2\}$	2	$\{2, 2\}$	1	1	15
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_r	-	-	0	0	0	0	0	1
S	-	-	-	-	-	-	-	15
ρ_i	-	-	-	-	-	-	-	0.0010

Table B.7: Other tuning parameters of MTMCMC for all proposal types (q_g , g_l , q_o , and q_r) in example 1; see Table C.15 for details. Notice that for MJMCMC* reported in the example only proposals of type 4 are used.

Proposal	Optimizer	Frequency	Type 1	Type 4	Type 3	Type 5	Type 6	Type 2
q_g	-	$\varrho = 0.9820$	0.1179	0.3357	0.2779	0.0200	0.2459	0.0021
S	-	-	$\{1, 1\}$	1	$\{1, 1\}$	1	1	20
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_l	-	0.0180	0	1	0	0	0	0
S	-	-	-	5	-	-	-	-
ρ_i	-	-	-	-	-	-	-	-
q_o	SA	0.5042	0.0636	0.3249	0.1571	0.2288	0.2246	0.0009
q_o	GREEDY	0.2183	0.0160	0.3085	0.1779	0.2474	0.2493	0.0007
q_o	MTMCMC	0.2774	0.2879	0.3016	0.1582	0.1107	0.1401	0.0013
S	-	-	$\{1, 1\}$	1	$\{1, 1\}$	1	1	20
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_r	-	-	0	0	0	0	0	1
S	-	-	-	-	-	-	-	20
ρ_i	-	-	-	-	-	-	-	0.0010

Table B.8: Other tuning parameters of MTMCMC for all proposal types (q_g , g_l , q_o , and q_r) in example 2; see Table C.15 for details.

Proposal	Optimizer	Frequency	Type 1	Type 4	Type 3	Type 5	Type 6	Type 2
q_g	-	$\varrho = 0.9816$	0.0932	0.2654	0.2197	0.0158	0.1944	0.2116
S	-	-	$\{1, 3\}$	3	$\{1, 3\}$	1	1	88
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_l	-	0.0164	0	1	0	0	0	0
S	-	-	-	20	-	-	-	-
ρ_i	-	-	-	-	-	-	-	-
q_o	SA	0.5553	0.0633	0.3165	0.1532	0.1548	0.1112	0.2011
q_o	GREEDY	0.2404	0.0149	0.2871	0.1656	0.2302	0.0820	0.2201
q_o	MTMCMC	0.2043	0.2310	0.1052	0.1877	0.1103	0.1682	0.1980
S	-	-	$\{1, 3\}$	3	$\{1, 3\}$	1	1	88
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_r	-	-	0	0	0	0	0	1
S	-	-	-	-	-	-	-	88
ρ_i	-	-	-	-	-	-	-	0.0010

Table B.9: Other tuning parameters of MTMCMC for all proposal types (q_g , q_l , q_o , and q_r) in example 3; see Table C.15 for details.

Proposal	Optimizer	Frequency	Type 1	Type 4	Type 3	Type 5	Type 6	Type 2
q_g	-	$\varrho = 0.9615$	0.1662	0.3323	0.1662	0.1662	0.1662	0.0029
S	-	-	$\{1, 1\}$	1	$\{1, 1\}$	1	1	13
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_l	-	0.0385	0	1	0	0	0	0
S	-	-	-	4	-	-	-	-
ρ_i	-	-	-	-	-	-	-	-
q_o	SA	0.5000	0.0657	0.3281	0.1588	0.2247	0.2209	0.0019
q_o	GREEDY	0.2500	0.0160	0.3083	0.1778	0.2472	0.2491	0.0014
q_o	MTMCMC	0.2500	0.2875	0.3012	0.1580	0.1105	0.1398	0.0026
S	-	-	$\{1, 1\}$	1	$\{1, 1\}$	1	1	13
ρ_i	-	-	$\hat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\hat{p}(\gamma_i \mathbf{y})$
q_r	-	-	0	0	0	0	0	1
S	-	-	-	-	-	-	-	13
ρ_i	-	-	-	-	-	-	-	0.0010

Table B.10: Other tuning parameters of MTMCMC for all proposal types (q_g , q_l , q_o , and q_r) in example 4; see Table C.15 for details.

Appendix B.1. Details on example 2

In the addressed data set the true regression parameters were chosen to be $\beta_0 = 99$ for the intercept, and for the slope coefficients

$$\beta = (-4, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1.2, 0, 37.1, 0, 0, 50, -0.00005, 10, 3, 0).$$

What concerns the covariates, X_1 and X_3 are factors from a group with 3 levels, X_4 and X_6 are correlated to them factors from another group with 3 levels, X_7 and X_8

are two jointly dependent through copulas exponentially distributed variables with rate 0.3, X_9, X_{10} and X_{11} are all uniformly distributed with range from -1 to 10 and also jointly dependent through copulas, X_{12}, X_{13}, X_{14} and X_{15} are multivariate normal with a zero mean, standard deviation of 0.2 and some covariance structure, X_{16} represents some seasonality incorporated by the sinus transformation of the radiant representation of some angle equal to the corresponding ordering numbers of observations, X_{17} is the quadratic trend associated to the squared value of positions of observations, $X_{19} = (-4 + 5X_1 + 6X_3)X_{15}$ and $X_{20} = (-4 + 5X_1 + 6X_3)X_{11}$, finally to avoid over specification 2 layers from the mentioned above groups of factors were replaced with some auxiliary covariates $X_2 = (X_{10} + X_{14}) \times X_9$ and $X_5 = (X_{11} + X_{15}) \times X_{12}$. The linear predictor is drawn as $\eta \sim N(\beta'X, 0.5)$, whilst the observations Y are independent Bernoulli variables with the probability of success modeled by a logit transformation of the linear predictor, namely $Y \sim \text{Bernoulli}\left(p = \frac{\exp(\eta)}{1+\exp(\eta)}\right)$.

Appendix C. Further results

In tables C.11 (example 1), C.12 (example 2) and C.13 (example 4) the estimated biases, corresponding to the RMSE estimates given in tables 3, 4 and 5, are reported. In addition, an extra simulation experiment on linear regression based on simulated data is reported in Appendix C.1.

Par	True	TOP	MJMCMC				BAS	MC ³		RS		MJMCMC*	
Δ	π_j	-	RM	MC	RM	MC	RM	MC	RM	MC	RM	RM	MC
γ_8	0.16	-3.51	-6.54	-10.28	-5.09	-9.64	-5.19	5.37	-3.20	4.96	-3.06	6.23	9.06
γ_{13}	0.16	-3.34	-7.44	-10.12	-5.57	-9.94	-6.25	7.46	2.86	8.06	2.65	6.38	10.54
γ_{14}	0.19	-3.24	-8.27	-11.69	-6.28	-11.93	-6.19	5.27	-1.86	5.37	-2.03	7.15	10.91
γ_{12}	0.22	-3.27	-6.82	-12.91	-5.54	-13.15	-3.08	3.00	-5.82	3.76	-5.06	5.29	10.93
γ_5	0.23	-2.56	-6.21	-12.71	-4.55	-13.35	-1.80	-4.79	-12.98	-4.28	-12.72	5.39	10.90
γ_9	0.23	-3.27	-9.45	-15.67	-7.35	-16.11	-9.26	4.53	-2.45	4.33	-2.10	7.68	11.06
γ_7	0.29	-2.31	-4.15	-12.04	-3.41	-12.36	-2.24	-0.47	-9.41	-1.00	-9.56	3.91	10.10
γ_4	0.30	-1.57	-5.82	-18.74	-3.67	-17.10	0.85	-12.67	-21.79	-13.24	-21.45	4.63	13.22
γ_6	0.33	-1.92	-8.49	-19.07	-6.09	-18.84	-3.06	8.99	7.16	10.09	6.81	5.87	15.43
γ_1	0.34	-2.51	-11.25	-21.94	-7.25	-20.29	-8.42	22.36	25.10	23.32	24.63	7.58	12.97
γ_3	0.39	-0.43	3.51	-7.20	2.09	-4.43	4.98	-21.11	-30.20	-21.13	-29.92	2.99	12.66
γ_2	0.57	1.58	5.66	-8.73	3.71	-7.51	13.73	-30.41	-37.52	-29.05	-37.12	5.11	14.04
γ_{11}	0.59	0.58	2.86	11.75	2.13	15.32	-3.95	10.67	21.68	10.29	21.23	2.77	12.77
γ_{10}	0.77	3.25	7.50	-2.57	5.91	2.33	15.42	-21.22	-19.06	-20.01	-19.55	6.41	14.27
γ_{15}	0.82	3.48	9.17	0.22	6.85	3.65	14.50	-69.61	-76.81	-69.14	-76.30	6.75	14.76
$C(\gamma)$	1.00	0.86	0.58	0.58	0.71	0.71	0.66	0.10	0.10	0.10	0.10	0.60	0.60
Eff	2^{15}	3276	1909	1909	3237	3237	3276	829	829	1071	1071	3264	3264
Tot	2^{15}	3276	3276	3276	5936	5936	3276	3276	3276	3276	3276	4295	4295

Table C.11: Bias for the 100 simulated runs of every algorithm on the Crime data (example 1); the values reported in the table are Bias $\times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$. See the caption of Table C.16 for further details.

Par	True	TOP	MJMCMC				BAS	MCBAS	RS	
Δ	π_j	-	RM	MC	RM	MC	RM	RM	RM	MC
γ_6	0.29	0.00	-7.23	-14.89	-4.48	-16.40	-6.46	-3.59	-5.96	0.23
γ_8	0.31	0.00	-5.97	-13.94	-3.89	-16.57	-5.57	-2.85	-5.28	-0.35
γ_{12}	0.35	0.00	-4.07	-8.12	-2.56	-11.65	-4.20	-1.82	-3.80	0.06
γ_{15}	0.35	0.00	-3.66	-8.85	-2.21	-12.04	-4.58	-1.35	-3.25	-0.28
γ_2	0.36	0.00	-4.60	-14.71	-2.81	-16.80	-5.39	-2.19	-3.51	0.04
γ_{20}	0.37	0.00	-4.16	-8.38	-2.46	-12.03	-3.30	-1.75	-4.07	-0.12
γ_3	0.40	0.00	-8.99	-19.22	-5.58	-21.72	-9.73	-4.63	-6.69	0.23
γ_{14}	0.44	0.00	1.08	7.12	0.51	7.63	3.68	-0.62	-0.99	0.22
γ_{10}	0.44	0.00	-2.68	-7.62	-1.68	-11.89	-4.79	-0.29	-1.19	0.13
γ_5	0.46	0.00	-1.74	-10.78	-0.88	-12.29	-3.93	0.57	0.55	-0.23
γ_9	0.61	0.00	0.32	-2.29	0.00	-1.24	3.78	0.22	1.99	-0.11
γ_4	0.88	0.00	5.61	6.20	3.71	6.13	6.60	5.54	7.58	-0.45
γ_{11}	0.91	0.00	5.36	6.47	3.87	6.84	4.64	3.01	4.29	-0.28
γ_1	0.97	0.00	1.86	0.98	1.32	1.17	2.43	1.94	2.28	-0.31
γ_{13}	1.00	0.00	0.00	-0.33	0.00	-0.29	0.00	0.00	0.00	-0.3
γ_7	1.00	0.00	0.00	-0.41	0.00	-0.36	0.00	0.00	0.00	-0.27
γ_{16}	1.00	0.00	0.00	-0.33	0.00	-0.31	0.00	0.00	0.00	-0.17
γ_{17}	1.00	0.00	0.00	-0.38	0.00	-0.35	0.00	0.00	0.00	-0.17
γ_{18}	1.00	0.00	0.00	-0.37	0.00	-0.32	0.00	0.00	0.00	-0.19
γ_{19}	1.00	0.00	0.00	-0.40	0.00	-0.32	0.00	0.00	0.00	-0.34
$C(\gamma)$	1.00	1.00	0.72	0.72	0.85	0.85	0.74	0.85	0.68	0.68
Eff	2^{20}	10000	5148	5148	9988	9988	10000	10000	1889	1889
Tot	2^{20}	10000	9998	9998	19849	19849	10000	10000	10000	10000

Table C.12: Bias for the 100 simulated runs of every algorithm on the simulated data of experiment 2; the values reported in the table are Bias $\times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$.

Appendix C.1. Example S.1

In this experiment we compare MJMCMC to BAS and competing MCMC methods (MC³, RS) using simulated data following the same linear Gaussian regression model as [Clyde et al. \(2011\)](#) with $p = 15$ and $n = 100$. All columns of the design matrix except for the ninth were generated from independent standard normal random variables and then centered. The ninth column was constructed so that its correlation with the second column was approximately 0.99. The regression parameters were chosen as $\beta_0 = 2$, $\beta = (-0.48, 8.72, -1.76, -1.87, 0, 0, 0, 0, 4, 0, 0, 0, 0, 0, 0)$ while the precision $\phi = 1$. For the parameters in each model Zellner's g-prior with $g = T$ is used. This leads to the marginal likelihood of the model to be proportional to (14). To complete the prior specification, we use (4) with $q = 0.5$. This leads to a rather simple example with two main modes in the model space. Simple approaches are expected to work well in this case. The exact posterior model probabilities may be obtained by enumeration of the model space in this case, making comparison with the truth possible.

In the BAS algorithm 3276 models unique were visited (about 10% of the total number of models). When running the MCMC algorithms approximately the same number of

Par	True	TOP	MJMCMC		RS	
Δ	π_j	RM	RM	MC	RM	MC
γ_4	0.0035	-0.0005	-0.0019	1.7361	-0.0189	1.6397
γ_6	0.0048	-0.0006	-0.0041	1.8155	-0.0241	1.5437
γ_7	0.0065	-0.0006	-0.0045	1.9763	-0.0338	0.2191
γ_3	0.0076	-0.0007	-0.0014	2.9714	-0.0339	0.5167
γ_8	0.0076	-0.0007	-0.0066	1.8370	-0.0326	1.1101
γ_5	0.0096	-0.0007	-0.0055	1.5439	-0.0430	1.1780
γ_{11}	0.0813	-0.0007	-0.0131	-0.7623	-0.1060	1.0394
γ_{12}	0.0851	-0.0006	-0.0042	-0.4290	-0.0637	0.3118
γ_9	0.1185	-0.0008	-0.0121	-1.3414	-0.1277	-0.4439
γ_{10}	0.3042	-0.0006	-0.0036	-8.4912	-0.0501	2.6866
γ_{13}	0.9827	-0.0002	0.0051	-1.6177	0.0607	-1.0082
γ_1	1.0000	0.0007	0.0000	-4.4528	0.0000	-1.0018
γ_2	1.0000	0.0000	0.0000	-2.3865	0.0000	-0.7782
$C(\gamma)$	1.0000	1.0000	0.9998	0.9998	0.9977	0.9977
Eff	8192	385	1758	1758	155	155
Tot	8192	385	3160	3160	10000	10000

Table C.13: Bias of the mean squared error (BIAS) from the 100 simulated runs of MJMCMC on the epigenetic data (example 4); the values reported in the table are $\text{BIAS} \times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$.

iterations were used. For the MJMCMC algorithm calculation of marginal likelihoods of models are stored making it unnecessary to recompute these when a model is revisited. Therefore, for MJMCMC also a number of iterations giving the number of *unique* models visited comparable with BAS was included. For each algorithm 100 replications were performed.

Table C.14, showing the root mean squared errors for different quantities, demonstrate that MJMCMC is outperforming simpler MCMC methods in terms of RM approximations of marginal posterior inclusion probabilities, individual model probabilities and the total captured mass. However the MC approximations seem to be slightly poorer for this example. Whenever both MC and RM approximations are available one should address the latter since they always have less noise. Comparing MJMCMC results to RM approximations provided by BAS (MC are not available for this method), MJMCMC is performing slightly worse when we have 3276 proposals (but 1906 unique models visited). However MJMCMC becomes equivalent to BAS when we consider 6046 proposals with 3212 unique models visited in MJMCMC (corresponding to similar computational time as BAS). In this example we are not facing a really multiple mode issue having just two modes. All MCMC based methods tend to revisit the same states from time to time and for such a simple example one can hardly ever beat BAS, which never revisits the same solutions and

Par	True	TOP	MJMCMC				BAS	MC ³			RS	
Δ	π_j	-	RM	MC	RM	MC	RM	MC	RM	MC	MC	RM
γ_{12}	0.09	0.29	2.11	5.31	1.19	5.73	1.23	2.77	4.27	2.14	3.83	
γ_{14}	0.10	0.28	2.13	6.99	1.13	6.25	1.14	2.92	4.31	2.59	3.95	
γ_{10}	0.11	0.28	2.31	7.41	1.31	7.74	1.15	3.06	4.31	2.40	4.07	
γ_8	0.12	0.27	1.97	6.44	1.09	7.80	0.97	2.77	4.01	2.23	3.87	
γ_6	0.13	0.25	2.25	8.87	1.27	8.46	1.05	3.12	4.74	2.72	4.31	
γ_7	0.14	0.25	2.06	7.75	1.29	8.51	1.05	3.45	4.52	2.50	4.17	
γ_{13}	0.15	0.24	2.42	9.98	1.36	8.79	1.15	3.50	4.87	2.44	4.38	
γ_{11}	0.16	0.24	2.36	9.38	1.22	8.31	1.13	3.64	4.71	3.01	4.52	
γ_{15}	0.17	0.23	1.96	9.38	1.08	9.73	0.78	3.92	4.27	3.32	3.84	
γ_5	0.48	0.00	1.22	15.66	0.50	12.90	0.27	3.69	1.41	4.35	1.59	
γ_9	0.51	0.10	1.15	16.35	0.38	12.92	0.37	16.70	5.62	6.93	2.08	
γ_2	0.54	0.07	1.46	20.69	0.58	15.38	0.39	16.56	5.25	6.91	1.46	
γ_1	0.74	0.18	2.15	6.43	1.06	5.97	1.20	4.10	3.55	4.51	3.90	
γ_3	0.91	0.25	1.61	3.03	0.92	3.33	1.57	2.96	3.66	3.42	4.10	
γ_4	1.00	0.01	0.00	6.08	0.00	2.66	0.00	0.01	0.01	0.17	0.01	
$C(\gamma)$	1.00	0.99	0.89	0.89	0.95	0.95	0.95	0.72	0.72	0.74	0.74	
Eff	2^{15}	3276	1906	1906	3212	3212	3276	400	400	416	416	
Tot	2^{15}	3276	3276	3276	6046	6046	3276	3276	3276	3276	3276	

Table C.14: Average root mean squared error (RMSE) over the 100 repeated runs of every algorithm on the simulated linear regression data; the values reported in the table are $\text{RMSE} \times 10^2$ for $p(\gamma_j = 1|\mathbf{y})\Delta = \gamma_j$. Tot is the total number of generated proposals, while Eff is the number of unique models visited during the iterations of the algorithms (for the TOP column all 2^{15} models were visited but the RMSE are based on the best 3276 models). RM corresponds to using the renormalization procedure (7) while MC corresponds to using the MC procedure (8). The two runs of MJMCMC are based on different Eff. The corresponding biases are reported in Table C.16 in Appendix C.2.

simultaneously draws the models to be estimated in a clever adaptive way with respect to the current marginal posterior inclusion probabilities of individual covariates.

Appendix C.2. Example S.1

Proposal	Optimizer	Frequency	Type 1	Type 4	Type 3	Type 5	Type 6	Type 2
q_g	-	$\varrho = 0.9836$	0.1176	0.3348	0.2772	0.0199	0.2453	0.0042
S	-	-	$\{2, 2\}$	2	$\{2, 2\}$	1	1	15
ρ_i	-	-	$\widehat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\widehat{p}(\gamma_i \mathbf{y})$
q_l	-	0.0164	0	1	0	0	0	0
S	-	-	-	4	—	—	—	—
ρ_i	-	-	-	-	-	-	-	-
q_o	SA	0.5553	0.0788	0.3942	0.1908	0.1928	0.1385	0.0040
q_o	GREEDY	0.2404	0.0190	0.3661	0.2111	0.2935	0.1046	0.0044
q_o	MTMCMC	0.2043	0.2866	0.1305	0.2329	0.1369	0.2087	0.0040
S	-	-	$\{2, 2\}$	2	$\{2, 2\}$	1	1	15
ρ_i	-	-	$\widehat{p}(\gamma_i \mathbf{y})$	-	-	-	-	$\widehat{p}(\gamma_i \mathbf{y})$
q_r	-	-	0	0	0	0	0	1
S	-	-	-	-	—	—	—	15
ρ_i	-	-	-	-	-	-	-	0.0010

Table C.15: Other tuning parameters of MTMCMC for all proposal types (q_g, q_l, q_o , and q_r) in example 1; Optimizer - to which optimizer the proposal belongs (if not relevant "-"); Frequency - the frequency at which the proposal is addressed (ϱ for q_g and $1 - \varrho$ for q_l) and the frequency within the set of local optimizers (P_o for local optimizers); Type X - the frequency of proposal of type X Table 1; S - maximal allowed size of the neighborhood for the corresponding proposal; ρ_i - probability of change of component i of the current solution (if applicable to the proposal).

Par	True	TOP	MJMCMC				BAS	MC ³		RS	
Δ	π_j	-	RM	MC	RM	MC	RM	MC	RM	MC	RM
γ_{12}	0.09	-0.29	-2.11	-4.95	-1.19	-5.47	-1.23	-0.14	-4.21	0.35	-3.80
γ_{14}	0.10	-0.28	-2.12	-6.58	-1.12	-6.07	-1.14	-0.23	-4.23	0.05	-3.89
γ_{10}	0.11	-0.28	-2.30	-6.89	-1.30	-7.64	-1.14	-0.10	-4.23	0.11	-4.02
γ_8	0.12	-0.27	-1.96	-6.16	-1.08	-7.69	-0.97	0.36	-3.94	-0.51	-3.81
γ_6	0.13	-0.25	-2.24	-8.03	-1.26	-8.33	-1.05	-0.65	-4.64	0.06	-4.24
γ_7	0.14	-0.25	-2.05	-7.45	-1.28	-8.37	-1.04	-0.13	-4.41	0.08	-4.12
γ_{13}	0.15	-0.24	-2.39	-9.62	-1.35	-8.62	-1.15	-0.49	-4.76	0.28	-4.32
γ_{11}	0.16	-0.24	-2.33	-8.69	-1.21	-7.95	-1.13	-0.38	-4.59	-0.10	-4.44
γ_{15}	0.17	-0.23	-1.93	-7.64	-1.06	-9.59	-0.78	-0.58	-4.15	-0.19	-3.74
γ_5	0.48	0.00	-1.15	-14.18	-0.47	-11.97	-0.25	-0.29	-0.94	0.46	-1.17
γ_9	0.51	-0.10	0.78	13.11	0.23	11.96	-0.32	-1.79	-2.20	-0.22	-1.53
γ_2	0.54	-0.07	-1.21	-18.43	-0.50	-14.64	0.34	1.73	0.29	0.35	-0.25
γ_1	0.74	0.18	2.12	4.88	1.04	3.99	1.19	-0.23	3.39	0.41	3.69
γ_3	0.91	0.25	1.60	-1.79	0.91	0.03	1.56	-0.40	3.59	-0.14	4.00
γ_4	1.00	0.01	0.00	-5.94	0.00	-2.49	0.00	0.01	0.01	-0.02	0.01
$C(\gamma)$	1.00	0.99	0.89	0.89	0.95	0.95	0.95	0.72	0.72	0.74	0.74
Eff	2^{15}	3276	1906	1906	3212	3212	3276	400	400	416	416
Tot	2^{15}	3276	3276	3276	6046	6046	3276	3276	3276	3276	3276

Table C.16: Bias for the 100 simulated runs of every algorithm on the simulated data; the values reported in the table are Bias $\times 10^2$ for $p(\gamma_j = 1|\mathbf{y})$. $C(\gamma)$ is as defined in (13). Tot is the total number of generated proposals, while Eff is the number of unique models visited during the iterations of the algorithms. RM corresponds to using the re-normalization procedure (9) while MC corresponds to using the MC procedure based on (8).