

# FAST BAYESIAN WHOLE-BRAIN FMRI ANALYSIS WITH SPATIAL 3D PRIORS

PER SIDÉN, ANDERS EKLUND, DAVID BOLIN AND MATTIAS VILLANI

**ABSTRACT.** Spatial whole-brain Bayesian modeling of task-related functional magnetic resonance imaging (fMRI) is a great computational challenge. Most of the currently proposed methods therefore do inference in subregions of the brain separately or do approximate inference without comparison to the true posterior distribution. A popular such method, which is now the standard method for Bayesian single subject analysis in the SPM software, is introduced in Penny et al. (2005b). The method processes the data slice-by-slice and uses an approximate variational Bayes (VB) estimation algorithm that enforces posterior independence between activity coefficients in different voxels. We introduce a fast and practical Markov chain Monte Carlo (MCMC) scheme for exact inference in the same model, both slice-wise and for the whole brain using a 3D prior on activity coefficients. The algorithm exploits sparsity and uses modern techniques for efficient sampling from high-dimensional Gaussian distributions, leading to speed-ups without which MCMC would not be a practical option. Using MCMC, we are for the first time able to evaluate the approximate VB posterior against the exact MCMC posterior, and show that VB can lead to spurious activation. In addition, we develop an improved VB method that drops the assumption of independent voxels a posteriori. This algorithm is shown to be much faster than both MCMC and the original VB for large datasets, with negligible error compared to the MCMC posterior.

## 1. INTRODUCTION

Over the past fifteen years, there has been much work devoted to Bayesian spatial modeling of task-related functional magnetic resonance imaging (fMRI) data. The motivation to this line of work has been to develop an extension to the classical general linear model (GLM) approach (Friston et al., 1995). The idea is to replace pre-smoothing of data and post-correction of multiple hypothesis testing – an approach that was recently shown to be problematic for cluster inference (Eklund et al., 2016) – with a proper spatial model. A cornerstone in the field is a series of papers (Penny et al., 2003, 2005a,b, 2007; Penny and Flandin, 2005) upon which the Bayesian spatial single-subject analysis method in the Statistical Parametric Mapping (SPM) software (SPM, 2002) is built. The method extends the classical GLM approach to a Bayesian framework with a spatial Gaussian Markov Random Field (GMRF) prior on the activity parameters and the temporal noise parameters, encouraging them to

---

Per Sidén: *Division of Statistics and Machine Learning, Dept. of Computer and Information Science, Linköping University, SE-581 83 Linköping, Sweden. E-mail: per.siden@liu.se. Corresponding author.*

Anders Eklund: *Division of Statistics and Machine Learning, Dept. of Computer and Information Science, Division of Medical Informatics, Dept. of Biomedical Engineering and Center for Medical Image Science and Visualization (CMIV), Linköping University, SE-581 83 Linköping, Sweden. E-mail: anders.eklund@liu.se.*

David Bolin: *Division of Mathematical Statistics, Dept. of Mathematical Sciences, Chalmers and University of Gothenburg, SE-412 96 Göteborg, Sweden. E-mail: david.bolin@chalmers.se.*

Mattias Villani: *Division of Statistics and Machine Learning, Dept. of Computer and Information Science, Linköping University, SE-581 83 Linköping, Sweden. E-mail: mattias.villani@liu.se.*

vary smoothly across the brain. The prior is formulated such that the optimal amount of smoothness can be estimated from the data. A fast variational Bayes (VB) algorithm is used for inference, but makes the assumptions that (i) the posterior factorizes over different types of parameters and (ii) the posterior for the activity parameters and the temporal noise model parameters factorizes over voxels. Assumption (i) is validated in the univariate (single-voxel) case in Penny et al. (2003) by comparing the VB posterior to the exact posterior obtained from MCMC sampling, but the error from assumption (ii) has to our knowledge never been properly examined. The VB framework also allows Bayesian model comparison based on the model evidence lower bound (Penny et al., 2007), but the computation of determinants of spatial precision matrices limits these types of analyses to be performed slice by slice or to sub-volumes containing  $\lesssim 10000$  voxels. Even without the model comparison, the SPM method would be considered too time consuming for the full 3D brain analysis for most practitioners.

A number of extensions of the SPM method have been developed. Harrison et al. (2008a) replace the stationary spatial prior in Penny et al. (2005b) with a non-stationary prior which is more adaptive and find evidence for this in data from an fMRI study of the auditory system. Groves et al. (2009) note that the spatial prior in Penny et al. (2005b) actually performs simultaneous smoothing and shrinkage of activity parameters in a rather non-flexible way. They separate these two effects using a Gaussian Process (GP) prior with a squared exponential kernel function, and infer the kernel length scale using Evidence Optimization (EO). The same issue is targeted in Yue et al. (2014), who use the GMRF representation of a Matérn field (Lindgren et al., 2011) and perform estimation using Integrated Nested Laplace Approximation (INLA) (Rue et al., 2009).

Even though all of these extensions seem like improvements from a modeling perspective they all struggle with computational complexity and inference is only performed slice by slice. A computationally attractive approach towards spatial 3D modeling is to first partition the brain into sub-volumes or parcels (Thirion et al., 2014), and then estimate a separate 3D model for each parcel, which can also be done in parallel, see for example Harrison et al. (2008b); Vincent et al. (2010); Musgrove et al. (2016). The inability to model dependencies between parcels is a bit unnatural, however, and Harrison et al. (2008b) notice discontinuities in posterior estimates along partition boundaries when using their model. Full 3D brain analysis with reasonable speed is achieved by Harrison and Green (2010) who approximate the zero mean spatial prior in Penny et al. (2005b) with a non-zero mean empirical prior, with the cost of no longer having a true generative model.

An alternative route to a spatial model is to view selecting active voxels as a variable selection problem, with a spatial prior on the probability of activation in a certain voxel, see for example Smith and Fahrmeir (2007) and Zhang et al. (2014). A comparison between estimation by VB and MCMC for this kind of model is available in the recent paper by Zhang et al. (2016). In a related line of work (Vincent et al., 2010; Risser et al., 2011; Chaari et al., 2013), voxel activations are modeled using a spatial mixture model of Gaussian distributions, where one of the mixture components has mean zero, corresponding to non-active voxels. This framework has the benefit of allowing for simultaneous estimation of the hemodynamic

response function (HRF), using both MCMC and variational methods which are applied parcel-wise. Another possibility is to model not only the activity parameters, but also the noise as spatially dependent (see for example Woolrich et al. (2004)), an assumption that seems natural but which comes with computational trouble because of the spatio-temporal structure.

A typical goal with these types of analyses is to compute Posterior Probability Maps (PPMs) (Friston and Penny, 2003), that is, brain maps showing the marginal probability of activation under a certain stimuli in each voxel. It is also common that these PPMs are thresholded at some level, for example at  $p > 0.9$ , to display active voxels or regions. However, such thresholding implicitly defines a hypothesis test in each voxel and these multiple tests must be corrected for. Given a spatially dependent posterior, a natural way to do this is using theory on excursion sets (Bolin and Lindgren, 2015) as in Yue et al. (2014) who define the term *joint PPM* based on the joint posterior, as opposed to *marginal PPMs*.

Our paper makes a number of contributions. First, we propose a fast and practical MCMC algorithm for slice-by-slice and whole brain task-fMRI analysis with spatial priors on the activity fields and the autoregressive noise parameters. The algorithm makes efficient use of sparsity and sampling from high-dimensional Gaussian distributions using preconditioned conjugate gradient (PCG) methods. These efficiency improvements reduce the computational complexity by several orders of magnitude. Collectively, they make it possible to perform posterior whole-brain analysis with spatial priors on problems where MCMC was previously simply not a practical option. Second, we develop a very fast VB approach that maintains the weaker independence assumption (i) in SPM’s VB method, but drops the stronger assumption (ii). Letting go of the assumption of spatially a posteriori independent voxels is non-trivial from a computational standpoint, and we employ several numerical techniques that together make the non-factorized VB approximation a very fast alternative to MCMC. The approximation errors from this non-factorized VB are shown to be essentially negligible for practical applications. Third, we demonstrate that the completely factorized VB in SPM12 can lead to spurious activations via a hitherto unexplored channel. Factorized VBs are well known to underestimate posterior variances, but we highlight and explain why the factorization over voxels can also result in a quite distorted smoothing of the activations.

The paper is organized as follows. We give a short background on the spatial model in SPM and the VB method used there to estimate the model. We then introduce the MCMC algorithm and the improved VB method and put extra emphasis on how to make these methods computationally efficient. In the next section we show results for simulated and real experiment data, and illustrate how the parameter estimates and PPMs differ between the different methods. The speed and convergence properties across methods are also compared. The last section contains a discussion and our conclusions. Derivation of results and implementation details can be found in the appendices.

The new methods presented in this article have been implemented as an extension to the SPM software, available at <http://www.fil.ion.ucl.ac.uk/spm/ext/#BFAST3D>.

## 2. BACKGROUND

We will consider single-subject fMRI-data containing  $T$  volumes with  $N$  voxels ordered in a  $T \times N$  matrix  $\mathbf{Y}$ . The experiment is described by the  $T \times K$  design matrix  $\mathbf{X}$ , with  $K$  regressors. The model in Penny et al. (2007) is written as

$$\mathbf{Y} = \mathbf{X}\mathbf{W} + \mathbf{E}, \quad (2.1)$$

where  $\mathbf{W}$  is a  $K \times N$  matrix of regression coefficients and  $\mathbf{E}$  is a  $T \times N$  matrix of error terms. Since this is the same model that will be considered here, we will use the same notation throughout, unless stated otherwise. Voxel-specific, normally distributed  $P$ th order AR models are assumed for the error terms, but for a clearer presentation here we will consider the special case  $P = 0$  and handle the more general case in Appendix A. The likelihood then becomes

$$p(\mathbf{Y}|\mathbf{W}, \lambda) = \prod_{n=1}^N \mathcal{N}(\mathbf{Y}_{:,n}; \mathbf{X}\mathbf{W}_{:,n}, \lambda_n^{-1} \mathbf{I}_T), \quad (2.2)$$

with  $\mathbf{Y}_{:,n}$  and  $\mathbf{W}_{:,n}$  denoting the  $n$ th column of  $\mathbf{Y}$  and  $\mathbf{W}$ ,  $\lambda_n$  as the noise precision of voxel  $n$  and  $\mathbf{I}_T$  a  $T \times T$  identity matrix. The likelihood factorizes over voxels, which is an assumption of non-spatial measurement noise that is made because the opposite assumption would be very computationally challenging. Instead, the spatial part of the model will enter via the prior on the regression coefficients

$$\begin{aligned} \mathbf{W}'_{k,\cdot} | \alpha_k &\sim \mathcal{N}(0, \alpha_k^{-1} \mathbf{D}_w^{-1}), \\ p(\mathbf{W}|\alpha) &= \prod_{k=1}^K p(\mathbf{W}'_{k,\cdot} | \alpha_k), \end{aligned} \quad (2.3)$$

where  $\mathbf{W}'_{k,\cdot}$  denotes the transposed  $k$ th row of  $\mathbf{W}$ ,  $\mathbf{D}_w$  is a fixed spatial  $N \times N$  precision matrix and  $\alpha = (\alpha_1, \dots, \alpha_K)'$  are hyperparameters to be estimated from the data. There are several possible choices for  $\mathbf{D}_w$ , but we will here focus on the SPM12 default choice, the unweighted graph-Laplacian (UGL) (called  $\mathbf{L}$  in Penny et al. (2005b)) which for each voxel has the number of adjacent voxels on the diagonal and  $\mathbf{D}_w(i, j) = -1$  if  $i$  and  $j$  are adjacent. For voxels in the interior part of the brain we will thus have 6's on the diagonal when modeling the whole 3D brain, and 4's when modeling each 2D slice separately. The main focus in Penny et al. (2005b) is on a different prior defined as  $\mathbf{D}_w = \mathbf{L}'\mathbf{L}$ . It is straightforward to use that prior also in our framework, but our experience is that it leads to too smooth posteriors, which is probably the reason why it is not the default option in SPM12. It also leads to slower inference since  $\mathbf{D}_w$  is less sparse when using this prior. The assumption of a sparse  $\mathbf{D}_w$  is the key to fast inference for this type of model using any method.

The hyperparameter  $\alpha_k$  will be estimated for each regressor to put the appropriate weight on the prior, depending on what smoothness is supported by the data. A higher  $\alpha_k$  brings the regression coefficient in each voxel closer to the mean of the coefficients in neighboring voxels (more smoothness), and globally all coefficients closer to zero (more shrinkage). The precision parameters  $\alpha_k$  and  $\lambda_n$  are assigned conjugate Gamma priors. For all the details on the generative model and the priors, see Figure 1 and Appendix A in Penny et al. (2007)

which also gives the values of the prior parameters that are default in SPM12 which are also used here (except  $r_1 = 10000$  in SPM12). That is, we use  $q_1 = u_1 = 10$ ,  $r_1 = 10000$  and  $q_2 = r_2 = u_2 = 0.1$ .

The VB algorithm for inference presented in Penny et al. (2005b) makes two independence assumptions for the joint posterior distribution. Firstly, the posterior is assumed to factorize over the different kinds of parameters, that is

$$q(\mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\lambda} | \mathbf{Y}) = q(\mathbf{W} | \mathbf{Y}) q(\boldsymbol{\alpha} | \mathbf{Y}) q(\boldsymbol{\lambda} | \mathbf{Y}), \quad (2.4)$$

with  $q$  denoting VB posteriors. Secondly, the regression parameters (and the AR-parameters in the general case) are assumed to factorize over voxels

$$q(\mathbf{W} | \mathbf{Y}) = \prod_{n=1}^N q(\mathbf{W}_{:,n} | \mathbf{Y}). \quad (2.5)$$

The second assumption is possibly the strongest and most counter-intuitive since it is clear that the spatial prior will generate dependence between voxels in the posterior. In the following section we will present an efficient MCMC (Gibbs) algorithm that performs inference without any of these two assumptions and an improved VB algorithm that drops only the second one. For an introduction to MCMC and VB, see Penny et al. (2003), who introduce these methods for the one-voxel case.

### 3. THEORY

Penny et al. (2005b) motivate the second independence assumption by noting that the posterior distribution for  $\mathbf{W}$  will otherwise have a full covariance matrix of size  $KN \times KN$ , which is prohibitively large. This is certainly true, and our algorithms are therefore designed to never compute the covariance matrix explicitly, but work with precision matrices instead. The posterior precision matrix is also of size  $KN \times KN$ , but it is sparse and can therefore be stored and used for computations quite cheaply. For example, the full conditional posterior distribution for  $\mathbf{w}_r = \text{vec}(\mathbf{W}')_r$  will be multivariate normal and can be characterized by (see Appendix A for the derivation)

$$\begin{aligned} p(\mathbf{w}_r | \mathbf{Y}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) &\propto \exp\left(-\frac{1}{2} \mathbf{w}_r' \tilde{\mathbf{B}} \mathbf{w}_r + \mathbf{b}_w' \mathbf{w}_r\right), \\ \mathbf{b}_w &= \text{vec}(\text{diag}(\boldsymbol{\lambda}) \mathbf{Y}' \mathbf{X}), \\ \tilde{\mathbf{B}} &= \mathbf{X}' \mathbf{X} \otimes \text{diag}(\boldsymbol{\lambda}) + \text{diag}(\boldsymbol{\alpha}) \otimes \mathbf{D}_w, \end{aligned} \quad (3.1)$$

where  $\tilde{\mathbf{B}}$  is a precision matrix which will maximally have  $K + 6$  non-zero elements in any row when the UGL prior is used in 3D. In other words,  $\mathbf{w}_r | \mathbf{Y}, \boldsymbol{\lambda}, \boldsymbol{\alpha} \sim \mathcal{N}(\tilde{\mathbf{B}}^{-1} \mathbf{b}_w, \tilde{\mathbf{B}}^{-1})$ .

As all the matrix inversions are avoided, the computational bottleneck of the algorithms developed in this paper will instead be to generate random samples from this and similar normal distributions, and to solve equation systems of the form  $\tilde{\mathbf{B}} \mathbf{x} = \mathbf{b}$ , given some  $KN \times 1$  vector  $\mathbf{b}$ . We will approach these problems in two ways, using Cholesky decomposition-based exact methods and preconditioned conjugate gradient (PCG) approximate methods.



### 3.1. GMRF sampling.

Rue and Held (2005) give a nice introduction to inference in latent Gaussian spatial models with sparse precision matrices, so called Gaussian Markov Random Fields (GMRFs), and also provide a good list of historical references, for example Besag (1974) and Woods (1972). In particular, they give computationally effective algorithms for sampling and equation solving for GMRFs based on first computing the Cholesky factor of the precision matrix, that is for a precision matrix  $\tilde{\mathbf{B}}$  compute a lower triangular matrix  $\mathbf{L}$  such that  $\mathbf{L}\mathbf{L}' = \tilde{\mathbf{B}}$ . Using so called reordering methods, such as the approximate minimum degree permutation (Amestoy et al., 1996), one can find a way to reorder the rows and columns of a sparse  $\tilde{\mathbf{B}}$  such that  $\mathbf{L}$  will be reasonably sparse as well, which will be important for speed. We can sample from the distribution in equation (3.1) and calculate  $E(\mathbf{w}_r | \mathbf{Y}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = \boldsymbol{\mu}_w = \tilde{\mathbf{B}}^{-1} \mathbf{b}_w$  using Algorithm 1. A key notion here is that, algorithmically, solving equations involving  $\mathbf{L}$  or  $\mathbf{L}'$  (called forward or backward substitution) is much faster than solving equations involving  $\tilde{\mathbf{B}}$  directly. In practice, this algorithm works well for GMRFs of dimension up to  $\approx 10000$ , but then starts to become too slow for practical use, the bottleneck being the Cholesky decomposition. This means that it can be used for slice- or parcel-wise inference in many cases, but not for whole 3D brain inference.

Papandreou and Yuille (2010) offer a method to sample from the posterior for  $\mathbf{W}$  that avoids the Cholesky decomposition, see Algorithm 2. To use it in our setting we first have to rewrite the prior density for  $\mathbf{W}'_{k,\cdot}$  according to

$$p(\mathbf{W}'_{k,\cdot} | \alpha_k) \propto \mathcal{N}(\sqrt{\alpha_k} \mathbf{G}_w \mathbf{W}'_{k,\cdot}; 0, \mathbf{I}_{N_{G_w}}) \propto \exp\left(-\frac{1}{2} \mathbf{W}_{k,\cdot} \alpha_k \mathbf{G}'_w \mathbf{G}_w \mathbf{W}'_{k,\cdot}\right), \quad (3.2)$$

for some  $N_{G_w} \times N$  matrix  $\mathbf{G}_w$ . For this prior to be equal to the one in equation (2.3),  $\mathbf{G}_w$  needs to be chosen such that  $\mathbf{D}_w = \mathbf{G}'_w \mathbf{G}_w$ . Fortunately, many priors are naturally specified through  $\mathbf{G}_w$  directly. In our case, we construct  $\mathbf{G}_w$  based on the interpretation of the UGL as a priori saying that the differences between adjacent voxels are i.i.d. normal, that is

$$W_{k,i} - W_{k,j} \stackrel{iid}{\sim} \mathcal{N}(0, \alpha_k^{-1}), \quad (3.3)$$

for all adjacent  $i$  and  $j$ . Thus, we can construct  $\mathbf{G}_w$  as having one row for every pair of adjacent voxels  $i$  and  $j$  with 1 in column  $i$  and  $-1$  in column  $j$ . A similar construction is possible for the main prior in Penny et al. (2005b). In cases where the prior is instead specified through  $\mathbf{D}_w$ , one way to obtain a  $\mathbf{G}_w$  is always available as the Cholesky factor of  $\mathbf{D}_w$ , as long as this is computable. Further, we construct  $\mathbf{B}_{data} = \text{diag}(\boldsymbol{\lambda}) \otimes \mathbf{X}'\mathbf{X}$  and  $\mathbf{L}_{data}$  as its Cholesky factor which is cheap to compute since  $\mathbf{B}_{data}$  will be banded with bandwidth  $K$  and also block diagonal. We define  $\mathbf{H}_w$  as in Penny et al. (2007) as the permutation matrix such that  $\text{vec}(\mathbf{W}) = \mathbf{H}_w \text{vec}(\mathbf{W}')$ .

The last piece of the second sampling method is to use PCG for solving equations of the form  $\tilde{\mathbf{B}}\mathbf{x} = \mathbf{b}$  approximately (Barrett et al., 1994; Manteuffel, 1980), with a computationally cheap incomplete Cholesky pre-conditioner  $\mathbf{M}$ . The efficiency of the PCG method increases if given a starting value  $\mathbf{x}^{start}$  that is close to the solution. The PCG algorithm will iterate until the relative residual  $\|\tilde{\mathbf{B}}\mathbf{x} - \mathbf{b}\| / \|\mathbf{b}\|$  becomes less than a user specified tolerance level  $\delta$ , e.g.

$10^{-8}$ , why PCG can be set to approximate the true solution arbitrarily well. This means that the approximation error that comes from PCG will be negligible in practice, which is exemplified in Figure 4.1b below, where the difference in posterior mean for the MCMC method using  $\delta = 10^{-6}$  and  $\delta = 10^{-8}$  is very close to zero. PCG can be used both as in Algorithm 2 for sampling or simply to compute the mean  $\mu_w = \tilde{\mathbf{B}}^{-1}\mathbf{b}_w$ .

---

**Algorithm 1** Cholesky based sampling from  $p(\mathbf{w}_r|\mathbf{Y}, \alpha, \lambda)$ 


---

**Require:**  $\tilde{\mathbf{B}}, \mathbf{b}_w$

- 1: Compute reordering based on  $\tilde{\mathbf{B}}$  and reorder  $\tilde{\mathbf{B}}$  and  $\mathbf{b}_w$  accordingly
  - 2: Compute  $\mathbf{L}$  as the Cholesky factor of  $\tilde{\mathbf{B}}$
  - 3: Solve  $\mathbf{L}\mathbf{x} = \mathbf{b}_w$
  - 4: Solve  $\mathbf{L}'\mu_w = \mathbf{x}$
  - 5: Sample  $\mathbf{z} \sim N(0, \mathbf{I}_{KN})$
  - 6: Solve  $\mathbf{L}'\mathbf{v} = \mathbf{z}$
  - 7: Compute  $\mathbf{w}_r = \mu_w + \mathbf{v}$
  - 8: Reorder  $\mathbf{w}_r$  and  $\mu_w$  using the inverse reordering computed in step 1
  - 9: **return**  $\mathbf{w}_r, \mu_w$
- 

---

**Algorithm 2** PCG based sampling from  $p(\mathbf{w}_r|\mathbf{Y}, \alpha, \lambda)$ 


---

**Require:**  $\tilde{\mathbf{B}}, \mathbf{b}_w, \mathbf{G}_w, \mathbf{B}_{data}, \mathbf{w}_r^{start}$

- 1: Compute  $\mathbf{L}_{data}$  as the Cholesky factor of  $\mathbf{B}_{data}$
  - 2: Sample  $\mathbf{z}_1 \sim \mathcal{N}(0, \mathbf{I}_{N_{G_w}})$
  - 3: Sample  $\mathbf{z}_2 \sim \mathcal{N}(0, \mathbf{I}_{KN})$
  - 4: Compute  $\mathbf{b} = \left( \text{blkdiag} \left[ \sqrt{\alpha_k} \mathbf{G}_w \right]_{k \in \{1, \dots, K\}} \right)' \mathbf{z}_1 + \mathbf{H}_w' \mathbf{L}_{data} \mathbf{H}_w \mathbf{z}_2 + \mathbf{b}_w$
  - 5: Compute reordering based on  $\tilde{\mathbf{B}}$  and reorder  $\tilde{\mathbf{B}}$  and  $\mathbf{b}$  accordingly
  - 6: Compute  $\mathbf{M}$  as the Incomplete Cholesky factor of  $\tilde{\mathbf{B}}$
  - 7: Solve  $\tilde{\mathbf{B}}\mathbf{w}_r = \mathbf{b}$  approximately using PCG with preconditioner  $\mathbf{M}$  and starting value  $\mathbf{w}_r^{start}$
  - 8: Reorder  $\mathbf{w}_r$  using the inverse reordering computed in step 5
  - 9: **return**  $\mathbf{w}_r$
- 

Figure 3.1 shows the average time it takes to produce a single sample from the conditional posterior of  $\mathbf{W}$  for the simulated data presented in Section 4.1 ( $K = 5$ ) for the Cholesky and PCG based sampling algorithms, as a function of the number of voxels. The Cholesky algorithm runs out of memory (32GB RAM) for  $N = 10^5$  (which is another issue for the Cholesky approach), but extending the seemingly linear behavior on the log-log-scale indicates that the PCG method would be roughly 100 times faster in this case. By providing a good starting value, as in the SVB algorithm presented below, we have observed the PCG timings to decrease with an additional factor in the range  $[2, 15]$ , a factor that increases as the algorithm converges.

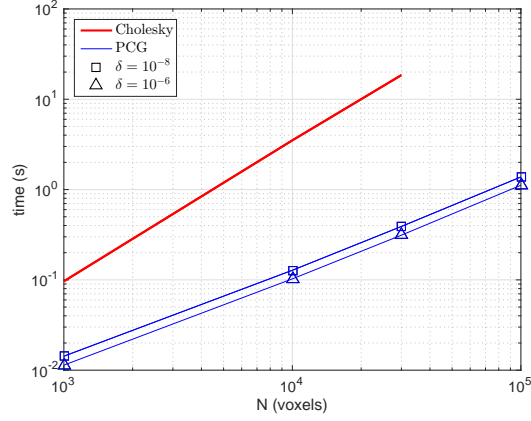


FIGURE 3.1. Average sampling times using Cholesky and PCG based sampling for GMRFs of size  $KN$  for the simulated data ( $K = 5$ ), with two different PCG tolerance levels  $\delta$ . In each case the processing time was computed as an average over 100 samples. Note that the presented PCG timings can be reduced by an additional factor in the range  $[2, 15]$  when a good starting value is available. The Cholesky algorithm runs out memory for  $N = 10^5$ .

### 3.2. MCMC algorithm.

In order to evaluate the true posterior of the model without any independence assumptions between parameters, we develop an algorithm for MCMC sampling. Since all priors are conjugate, we obtain closed form expressions for all full conditionals and can therefore perform Gibbs sampling. The full conditional posterior for  $\mathbf{W}$  is given in (3.1) and the full conditionals for  $\lambda$  and  $\alpha$  are given by

$$\begin{aligned}\lambda_n | \mathbf{Y}, \mathbf{W}, \alpha &\sim \text{Ga}(\tilde{u}_{1n}, \tilde{u}_2), \\ \alpha_k | \mathbf{Y}, \mathbf{W}, \lambda &\sim \text{Ga}(\tilde{q}_{1k}, \tilde{q}_2),\end{aligned}\tag{3.4}$$

with

$$\begin{aligned}\frac{1}{\tilde{u}_{1n}} &= \frac{1}{2} (\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n} - 2\mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n}) + \frac{1}{u_1}, \\ \tilde{u}_2 &= \frac{T}{2} + u_2, \\ \frac{1}{\tilde{q}_{1k}} &= \frac{1}{2} \mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}'_{k,\cdot} + \frac{1}{q_1}, \\ \tilde{q}_2 &= \frac{N}{2} + q_2.\end{aligned}\tag{3.5}$$

See Appendix A for the derivation of these and the corresponding full conditionals in the  $P > 0$  case.

The Gibbs algorithm returns  $N_{iter}$  samples from the joint posterior of all parameters, which can be used for posterior inference about any subset of parameters. For example, given samples  $\mathbf{W}^{(1:N_{iter})}$  we can compute the marginal PPM for any  $K \times 1$  contrast vector  $\mathbf{c}$  as

$$P(\mathbf{c}' \mathbf{W}_{\cdot,n} > \gamma | \mathbf{Y}) \approx \frac{1}{N_{iter}} \sum_{j=1}^{N_{iter}} I(\mathbf{c}' \mathbf{W}_{\cdot,n}^{(j)} > \gamma),\tag{3.6}$$



for voxel  $n$  and some activity threshold  $\gamma$ . Furthermore, since the MCMC posterior does not factorize over voxels, it is also meaningful and straightforward to compute the joint probability of activation for any set of voxels  $E$  as

$$P\left(\mathbf{c}'\mathbf{W}_{:,E} > \gamma \mathbf{1}'_{|E|} | \mathbf{Y}\right) \approx \frac{1}{N_{iter}} \sum_{j=1}^{N_{iter}} I\left(\mathbf{c}'\mathbf{W}_{:,E}^{(j)} > \gamma \mathbf{1}'_{|E|}\right), \quad (3.7)$$

where  $\mathbf{1}_{|E|}$  is a vector of ones of length  $|E|$ . Using the theory on excursion sets developed by Bolin and Lindgren (2015), we can thereby compute the joint PPMs introduced in Yue et al. (2014) that avoid the problem of multiple hypothesis testing.

### 3.3. Spatial VB.

Out of the two posterior independence assumptions in Penny et al. (2005b), we view the second one as the strongest, that is the assumption that the posterior for  $\mathbf{W}$  factorizes over voxels. The developed MCMC algorithm relieves us from both independence assumptions, but has its limitations in terms of speed and memory. We therefore seek to develop an improved VB algorithm that maintains the efficiency gain from the first assumption of independence between the different types of parameters, but drops the second assumption and models the joint posterior of  $\mathbf{W}$ . We will refer to this algorithm as Spatial Variational Bayes (SVB) and to SPM's factorized VB algorithm as Independent Variational Bayes (IVB).

The SVB posterior will be computed iteratively, just as in the IVB algorithm, for one parameter at a time given the approximate posterior of the others. If we denote  $\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\lambda}, \boldsymbol{\alpha}\}$ , then (Bishop, 2006)

$$\log q(\theta_j) = E_{\theta_{-j}}[\log p(\mathbf{Y}, \boldsymbol{\theta})] + \text{const}, \quad (3.8)$$

with the expectation taken with respect to the VB posterior  $q(\theta_{-i})$ . This means that for  $q(\mathbf{W})$  we get (see Appendix B)

$$\begin{aligned} \log q(\mathbf{w}_r) &= E_{\boldsymbol{\lambda}, \boldsymbol{\alpha}}[\log p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\lambda}) + \log p(\mathbf{W}|\boldsymbol{\alpha})] + \text{const} \\ &= -\frac{1}{2}\mathbf{w}_r' \tilde{\mathbf{B}}^{SVB} \mathbf{w}_r + \mathbf{b}_w^{SVB'} \mathbf{w}_r + \text{const}, \\ \mathbf{b}_w^{SVB} &= \text{vec}(\text{diag}(E_{\boldsymbol{\lambda}}[\boldsymbol{\lambda}]) \mathbf{Y}' \mathbf{X}), \\ \tilde{\mathbf{B}}^{SVB} &= \mathbf{X}' \mathbf{X} \otimes \text{diag}(E_{\boldsymbol{\lambda}}[\boldsymbol{\lambda}]) + \text{diag}(E_{\boldsymbol{\alpha}}[\boldsymbol{\alpha}]) \otimes \mathbf{D}_w, \end{aligned} \quad (3.9)$$

that is, we get exactly the same expression as for the full conditionals only the values for  $\boldsymbol{\lambda}$  and  $\boldsymbol{\alpha}$  are replaced by their expectations with respect to their variational posteriors  $q(\boldsymbol{\lambda})$  and  $q(\boldsymbol{\alpha})$ . This expression is simple because  $\boldsymbol{\lambda}$  and  $\boldsymbol{\alpha}$  enter linearly, which will not be the case in general. For  $q(\boldsymbol{\alpha})$  we get (note that the likelihood  $p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\lambda})$  does not depend on  $\boldsymbol{\alpha}$ )

$$\begin{aligned} \log q(\boldsymbol{\alpha}) &= E_{\mathbf{W}, \boldsymbol{\lambda}}[\log p(\mathbf{Y}|\mathbf{W}, \boldsymbol{\lambda}) + \log p(\mathbf{W}|\boldsymbol{\alpha}) + \log p(\boldsymbol{\alpha})] + \text{const} \\ &= \left(\frac{N}{2} + q_2 - 1\right) \sum_{k=1}^K \log \alpha_k - \sum_{k=1}^K \alpha_k \left[\frac{1}{2} E_{\mathbf{W}}[\mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}'] + \frac{1}{q_1}\right] + \text{const}. \end{aligned} \quad (3.10)$$

So, just as for the full conditionals,  $\alpha_k$  will be Gamma distributed a posteriori with parameters

$$\begin{aligned}\frac{1}{\tilde{q}_{1k}^{SVB}} &= \frac{1}{2}E_{\mathbf{W}} [\mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}'] + \frac{1}{q_1}, \\ \tilde{q}_2^{SVB} &= \frac{N}{2} + q_2.\end{aligned}\tag{3.11}$$

The problem here is the expectation of the quadratic form

$$E_{\mathbf{W}} [\mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}'] = E_{\mathbf{W}} [\mathbf{W}_{k,\cdot}] \mathbf{D}_w E_{\mathbf{W}} [\mathbf{W}_{k,\cdot}'] + \text{tr}(\mathbf{D}_w \text{Cov}[\mathbf{W}_{k,\cdot}, \mathbf{W}_{k,\cdot}]),\tag{3.12}$$

where the second term requires inversion (or at least partial inversion) of the posterior precision for  $\mathbf{W}$ ,  $\tilde{\mathbf{B}}^{SVB}$ , which is computationally infeasible in the general 3D case. To avoid this, we adopt a Monte Carlo (MC) sampling based approach to compute the expectation. The PCG sampling method provides an efficient way to generate a number of ( $N_s$ ) samples  $\mathbf{W}^{(1:N_s)}$  from the VB posterior  $q(\mathbf{W})$ , by simply replacing  $\tilde{\mathbf{B}}$  and  $\mathbf{b}_w$  in Algorithm 2 with  $\tilde{\mathbf{B}}^{SVB}$  and  $\mathbf{b}_w^{SVB}$ . These samples can be used to approximate the expectation as

$$E_{\mathbf{W}} [\mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}'] \approx \frac{1}{N_s} \sum_{j=1}^{N_s} \mathbf{W}_{k,\cdot}^{(j)} \mathbf{D}_w \mathbf{W}_{k,\cdot}^{(j)'}.\tag{3.13}$$

Similar MC approximations will be used in the VB update equations of  $\lambda$  and for the other parameters when  $P > 0$ , see Appendix B.

The SVB algorithm is similar to the MCMC algorithm in that the computational bottleneck will be the sampling of  $\mathbf{W}$ , but there are some important differences. The MCMC algorithm runs for a large number of iterations ( $N_{iter}$ , thousands) producing one sample from  $p(\mathbf{W}|\mathbf{Y}, \lambda, \alpha)$  in each iteration. For SVB it is enough to run for a much smaller number of iterations (tens), but each iteration draws a larger number of samples ( $N_s$ , tens or hundreds) from  $q(\mathbf{W})$ . When PCG based sampling is used SVB is advantageous because (i) the same pre-conditioner  $\mathbf{M}$  can be used for all samples in each VB iteration (ii) the same random seeds (the same  $\mathbf{z}_1$  and  $\mathbf{z}_2$  in Algorithm 2) can be used across VB iterations, making the previous iteration samples very good starting values for the PCG (iii) since the samples are independent, the sampling can be fully parallelized within each iteration by running Algorithm 2 on separate cores. All these points contribute to that SVB can be run much faster than MCMC in general, with the price being the first assumption of independence between different kinds of parameters.

The MC approximation introduced in equation (3.13) adds a stochastic approximation error to the already approximate VB posterior. Figure 4.1b below quantifies the size of this error with respect to the number of samples  $N_s$  on simulated data. Convergence results for stochastic VB methods (Kingma and Welling, 2014; Gunawan et al., 2016) and for stochastic variants of the related expectation-maximization (EM) algorithm (Chan and Ledolter, 1995; Delyon et al., 1999) are available in the literature. However, these results do not apply to our setting since they build on repeated sampling, while the SVB algorithm presented here only draws  $\mathbf{z}_1$  and  $\mathbf{z}_2$  in Algorithm 2 in the first iteration and then re-uses those same random numbers at the subsequent iterations. Re-using the random numbers speeds up SVB (typically by a factor between 2 and 15) since it allows us to use good PCG starting values,

and the additional approximation error resulting from fixed random numbers is small in our applications. See Appendix C for more details about the convergence.

In theory, all SVB posterior statistics, such as PPMs, can be computed from the approximate posterior  $q(\mathbf{W})$  in equation (3.9). However, since it is parametrized using the precision matrix even basic marginal statistics, as variances and posterior probabilities, are not trivially obtained since this requires the inversion (or at least the Cholesky factor) of the precision matrix. We can instead use the sample variance of the samples  $\mathbf{W}^{(1:N_s)}$  to get a fast approximation (Papandreou and Yuille, 2011), which can be used to compute the marginal PPMs in each voxel. For contrast PPMs, we similarly use the sample covariance matrix for each voxel. While a small value of  $N_s$  ( $\leq 100$ ) seems sufficient for convergence of the SVB algorithm, covariance estimates based on the same number of samples will be quite noisy (see Figure 4.2 below). A straightforward strategy to reduce the noise would be to generate additional samples from  $q(\mathbf{W})$  in a post-processing step, to further improve the covariance estimates. Such a step could be time consuming, however, and we are currently working on a different, more efficient way to compute covariances for a given sparse precision matrix.

#### 4. RESULTS

In this section we present results comparing the three methods (IVB, SVB and MCMC) on the same data. As MCMC is exact (in the sense of being simulation consistent with a small and controllable error), we can view this as the ground truth when evaluating the other methods. IVB is run using SPM12. The SVB algorithm is implemented as an extension to the original IVB algorithm by manipulating the original SPM12 Matlab code, while the MCMC algorithm is implemented in a separate Matlab function. The overhead time is low compared to the main computation steps for all three methods, so all timing comparisons should be considered fair. We perform the comparison on simulated data with the focus on computational efficiency as a function of data size, and on two different real data sets with the focus on the resulting posteriors and PPMs.

For any comparison to be fair, both regarding computing time and estimated posteriors, we need the different algorithms to reach the same level of convergence. In Appendix C we discuss some details on the convergence of the different methods and show that the SPM12 default setting of 4 VB iterations is usually not sufficient. In the results below, all methods including IVB are run until convergence. In Appendix C we also provide some practical details about the implementation of respective method and about the computers used to perform the analyses.

##### 4.1. Simulated data.

We simulate synthetic data from the model, with  $K = 5$  and  $P = 1$  and pick parameter values for the voxel intercept and noise standard deviation that approximately match those of the face repetition data described below. Several different values of  $\alpha$  were used to simulate conditions with varying informativeness, see Appendix D for details. We run the algorithms on simulations of  $N = 10^3, 10^4$  and  $10^5$  voxels to get an idea how they scale and compare with the number of voxels. Figure 4.1a shows the processing time until convergence, defined as the time until the estimated posterior mean of  $\alpha$  reaches within 1% of its final value, for

respective algorithm. In Figure 4.1b, the accuracy of respective algorithm is evaluated. This is based on the root-mean-squared-error (RMSE) of the marginal posterior mean of activation coefficients  $\mathbf{W}$  compared to the posterior mean from MCMC with  $\delta = 10^{-8}$ .

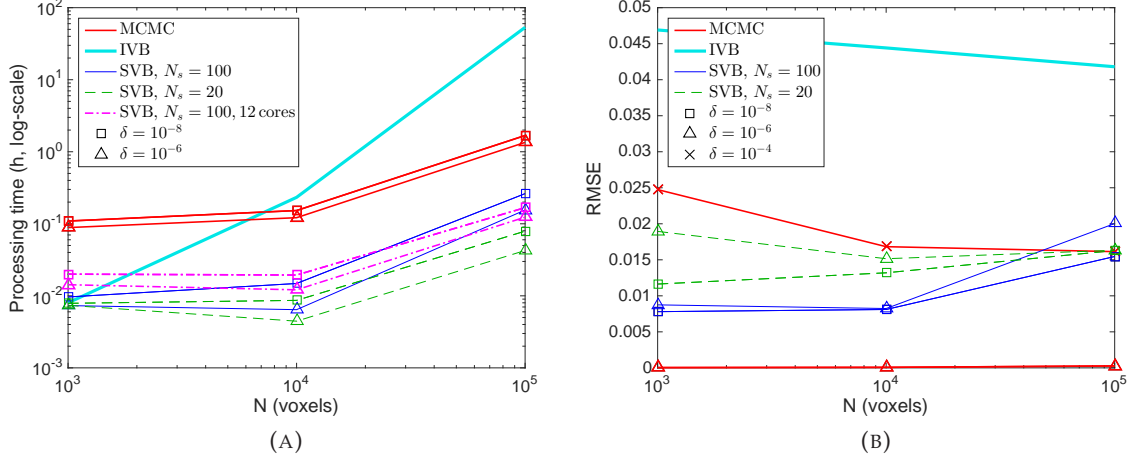


FIGURE 4.1. (a) Processing time as a function of the number of voxels for the different algorithms and different PCG tolerance levels  $\delta$ . For SVB, the number of MC samples  $N_s$  and computing cores are also varied. The computations were made on simulated data with  $K = 5$  using the 3D prior. (b) RMSE of activation coefficient posterior mean for the different algorithms as compared to the MCMC posterior with  $\delta = 10^{-8}$ , based on the first 4 regressors.

Figure 4.1 indicates that the SVB and MCMC algorithms scale much better with the number of voxels than IVB does, and also provide a higher accuracy. Lowering the PCG tolerance  $\delta$  from  $10^{-8}$  to  $10^{-6}$  gives some speedup while seemingly sacrificing little in accuracy. However, for  $\delta = 10^{-4}$  the RMSE for MCMC increases and for SVB it becomes as high as 0.25 with very noisy convergence times (not shown), so this tolerance must not be set too high. The speedup achieved by lowering  $N_s$  seems almost linear, but results in a lower accuracy. The speed gain in parallelizing is small due to overhead costs, but we expect greater speed-ups for larger data sets when each VB iteration requires more time, which is seen for the real data in Table 1 below. Note that the timing results in Figure 4.1 (and also in Table 1) need to be interpreted with caution since they are based on single runs of the stochastic MCMC and SVB methods, but these graphs provide insight about how the timing largely compares for the different methods and settings.

#### 4.2. Real data.

Two real task-fMRI data sets are considered, the face repetition data used in Penny et al. (2005b) and data from a visual object recognition experiment from the OpenfMRI database (Poldrack et al., 2013), see Appendix D for more details on these data sets.

Approximate processing times for these data are shown in Table 1, both for slice-wise analysis using the 2D prior, and whole-brain analysis using the 3D prior. The IVB method scales much worse with the number of voxels and is hence slow in the 3D case, while MCMC

and SVB are not necessarily slower in 3D than in 2D. For whole-brain inference with the 3D prior, SVB is generally the fastest option, especially when run in parallel.

	Face repetition		Object recognition	
	2D prior	3D prior	2D prior	3D prior
IVB	4.9	190	1.9	26
MCMC	110	150	230	76
SVB	5.3	11	22	20
SVB, 4 cores	2.8	3.8	8.7	7.6

TABLE 1. Approximate processing times (h) for the different real data sets and algorithms using the 2D/3D prior.

We first consider the face repetition data. Figure 4.2 shows the posterior mean and standard deviation of the activation contrast, estimated using all three methods for the 2D prior. Comparing IVB to MCMC, we see bias both in the estimated IVB mean (the maximum error is 1.4 across all voxels) and standard deviation (maximum error 25%). A separate experiment with  $\alpha$  fixed to the same value for both methods showed the well-known systematic underestimation of standard deviations by IVB (roughly by 8%, results not shown). However, IVB tends to also underestimate  $\alpha$ , as shown in Figure 4.6 below, leading to less shrinkage and VB errors in standard deviation that go in both directions. Comparing SVB to MCMC, we see that the estimated SVB mean is much more correct than the IVB (maximum error 0.2). However, the SVB standard deviation estimates are quite inaccurate (maximum error 26%), but this is mainly due to the noisy covariance estimates discussed in Section 3.3.

Computed marginal PPMs for various settings, thresholded at 0.9, are shown in Figure 4.3. The first row shows estimates for IVB after 4 VB iterations for the 2D prior (which is the SPM12 default), and after convergence for both the 2D and 3D prior. In the second and third row one can see the corresponding PPMs estimated using MCMC and SVB, and also the MCMC based joint PPMs which were computed using the R package `excursions`. The 2D joint PPM is computed based only on voxels within this slice while the 3D joint PPM is computed based on all voxels in the brain, so these maps are not directly comparable. The differences between the non-converged and converged SPM results seem rather small for this data set. Overall, the PPMs from IVB agree quite well with the PPMs from MCMC; an exception is a small cluster of activity in IVB which is lacking in the activation map from MCMC. The MCMC and SVB PPMs are hard to distinguish. Larger differences are found when comparing results based on the 2D vs. 3D prior and when comparing marginal vs. joint PPMs.

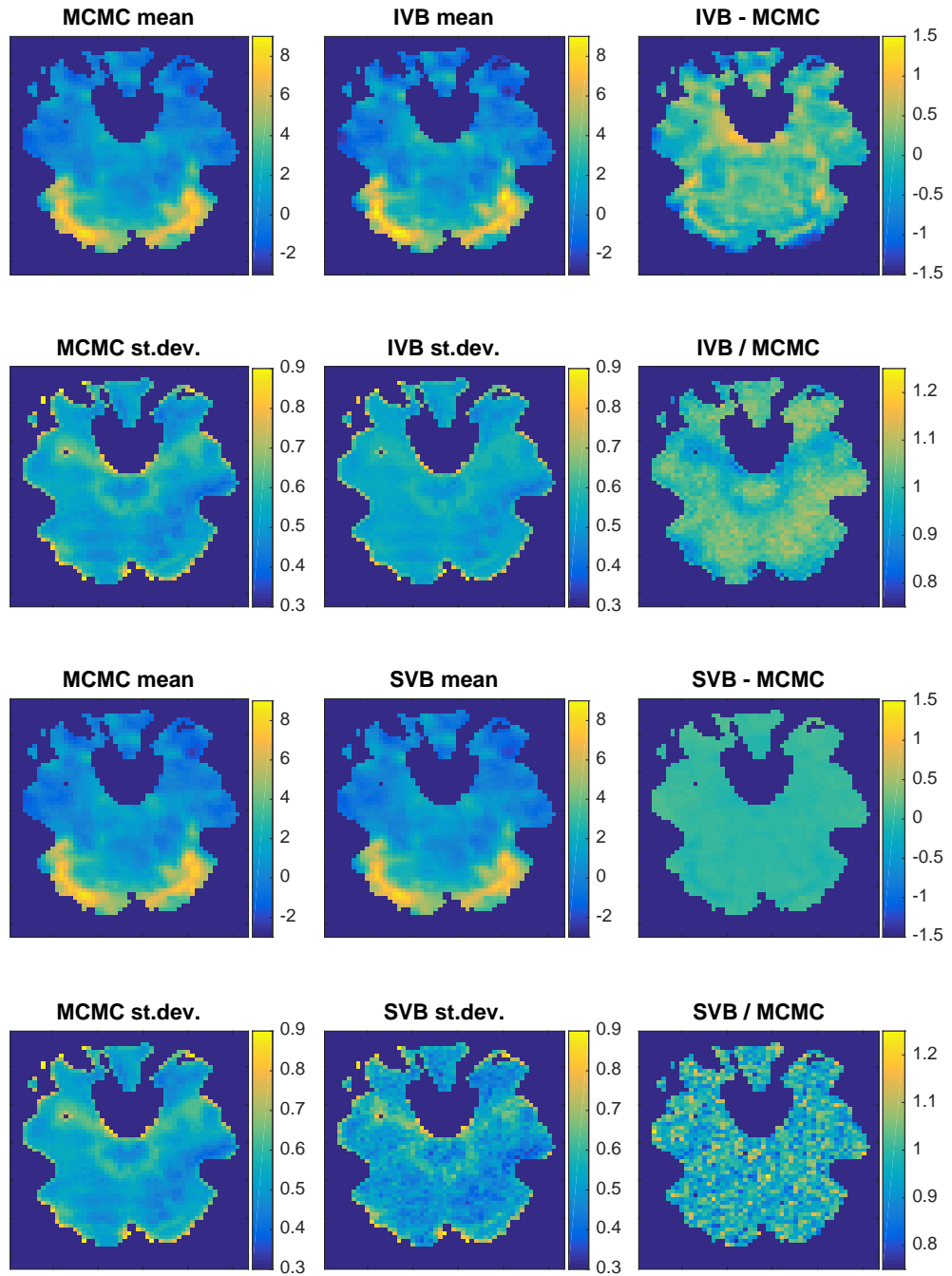


FIGURE 4.2. Posterior mean (top row) and standard deviation (second row) estimated using MCMC (left) and IVB (middle) for the contrast (mean effect of faces) using the 2D prior. The right column shows the differences in mean and standard deviation ratio of the estimated posteriors. Row three and four show the corresponding results for SVB.



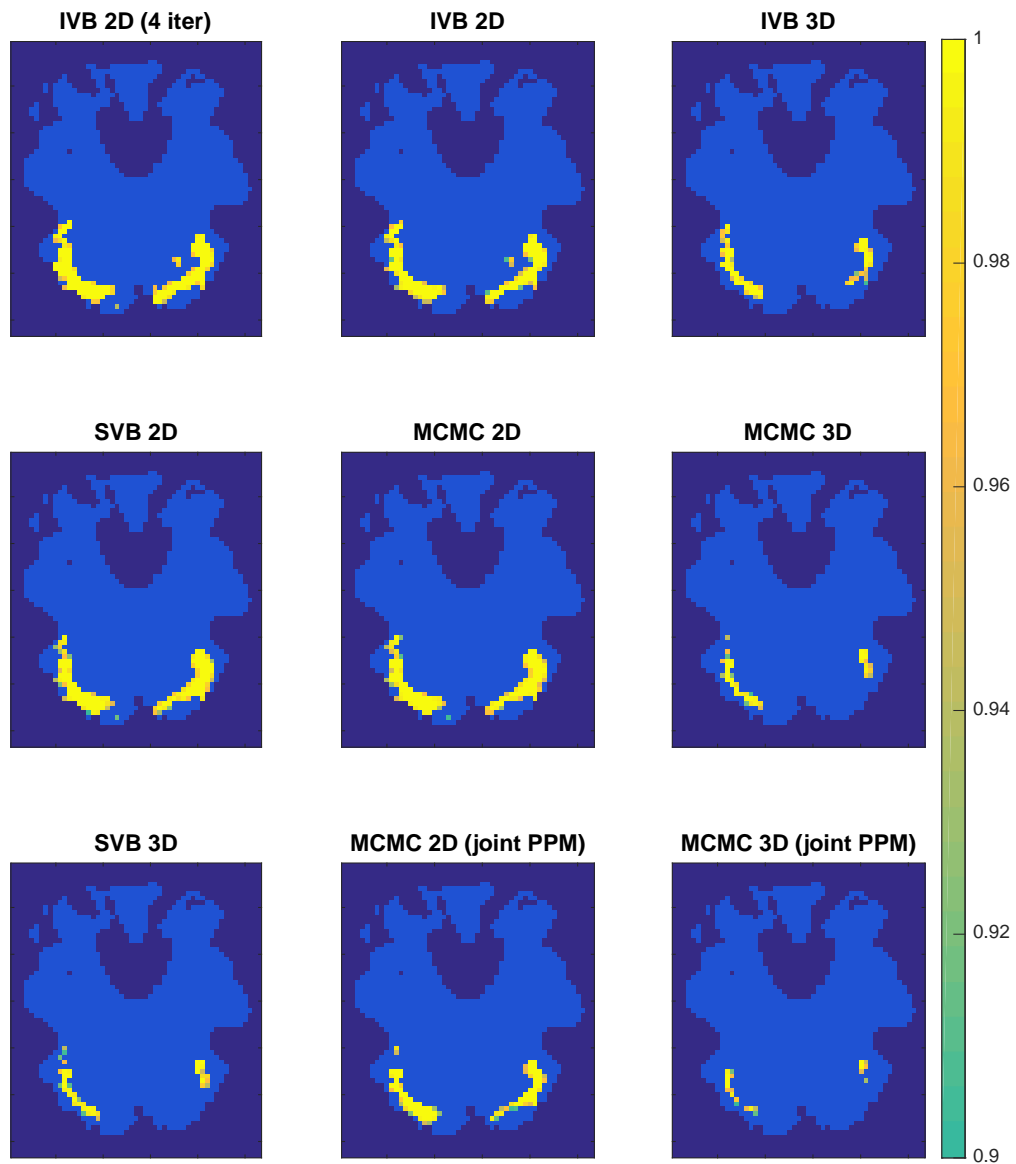


FIGURE 4.3. PPMs estimated for the face repetition data with different methods and 2D or 3D spatial priors. The PPMs show probabilities of the contrast (mean effect of faces) exceeding 1% of the global mean signal, thresholded at 0.9.

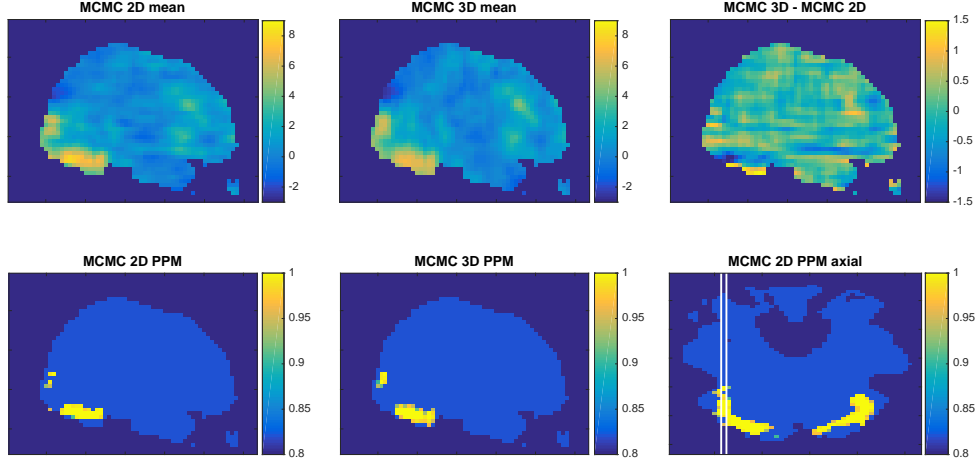


FIGURE 4.4. Posterior mean (top row) and marginal PPMs (bottom row) for the face repetition data for the MCMC method when using the 2D and 3D spatial priors for the sagittal slice indicated in the bottom right figure. The PPMs show probabilities of the contrast (mean effect of faces) exceeding 1% of the global mean signal, thresholded at 0.9.

The smaller activity regions from the 3D prior compared to the 2D prior is due to shrinkage towards a larger number of non-active voxels nearby in the  $z$ -dimension, which is depicted in Figure 4.4 that compares the posterior mean and PPM for the two priors in a sagittal slice. It is clear that the posterior mean is generally lower in the active regions when using the 3D prior, which can only be explained with the assumed dependence with non-active voxels in the  $z$ -dimension. For this particular slice, the effect is nevertheless strong enough for these voxels to be classified as active in the PPMs, but for other slices the higher smoothness can bring posterior probabilities below the 0.9-threshold, for voxels classified as active when using the 2D prior. At the same time, for the 2D prior we observe discontinuous effects between slices, for example in the inferior part of the largest active blob, while the 3D prior lends strength to the voxels below, classifying them as active. It should be noted that the 2D and 3D prior lead to rather different models (the 2D prior implies a different covariance structure and has many more parameters and is therefore more flexible), so the results are not directly comparable and greatly data dependent. To decide which of these priors (or perhaps a more flexible 3D prior or parcel based method) is best is a model selection problem, which we see as important future work.

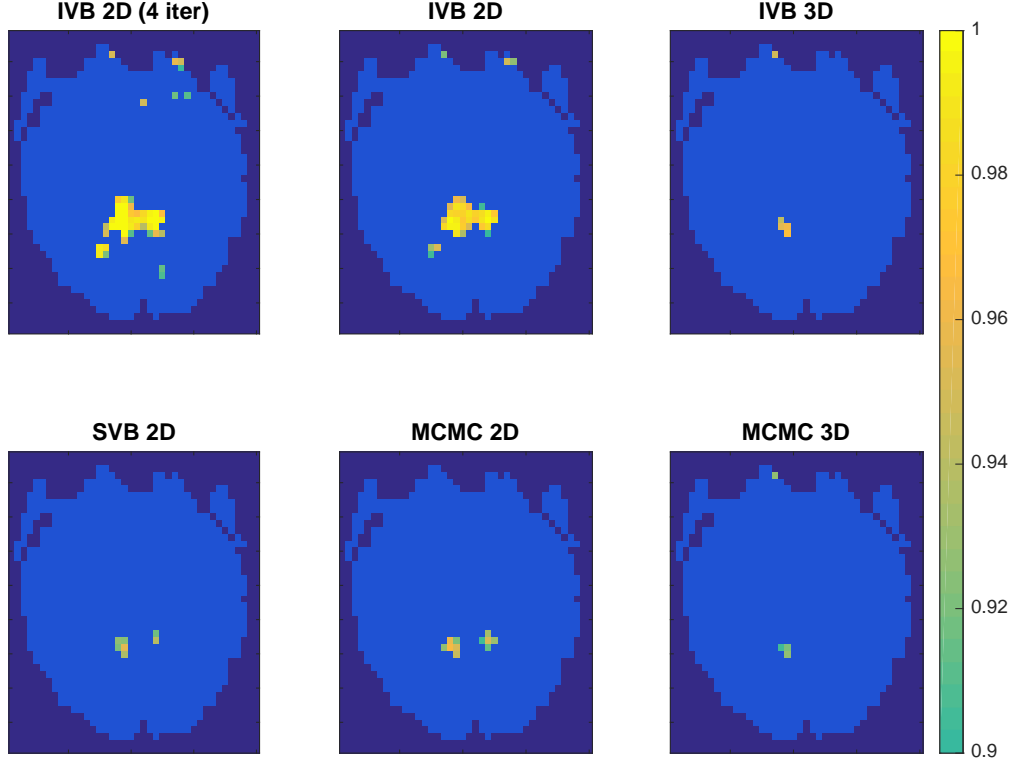


FIGURE 4.5. Marginal PPMs estimated for the object recognition data with different methods and 2D or 3D spatial priors. The PPMs show probabilities of the contrast (houses vs. faces) exceeding 0.5% of the global mean signal, thresholded at 0.9.

Although we observe some errors in the posterior mean and standard deviation for IVB as compared to MCMC, this error is not big enough to make much impact on the PPMs for the face repetition data set. Our second dataset, the object recognition data, shows a situation where this error can in fact be severe also for the PPMs. Figure 4.5 shows computed PPMs for one slice of the object recognition data, for some different methods and using the 2D/3D prior. It is clear that the independence assumption in IVB can lead to severely distorted activation maps, and that SVB is a much more accurate approximation for this dataset.

Much of the differences in brain maps between IVB and the other methods can be attributed to the underestimation of hyperparameters. Figure 4.6 shows the estimated posteriors of the spatial hyperparameters  $\alpha$  and  $\beta$  for the main regressors and AR coefficient for the different methods and data sets. When data are informative, as for the intercept, the VB methods generally approximate the hyperparameter posteriors well, but as data become less informative (this is clearly seen for the simulated data) the approximate posteriors from IVB underestimate both the location and dispersion of the posterior. The underestimation of posterior dispersion is a well known issue with VB quite generally (Bishop, 2006), but the inability of IVB to correctly approximate the posterior location is unusual. The answer comes

from the assumption of posterior independent voxels. Since the noise is assumed to be spatially independent, the posterior dependence between voxel activations comes solely from the spatial prior. IVB therefore pushes the  $\alpha_k$ 's to lower values in an attempt to reduce the influence of the prior. Indeed, SVB is much better at finding the correct location, but, like most VB methods, tends to underestimate the dispersion.

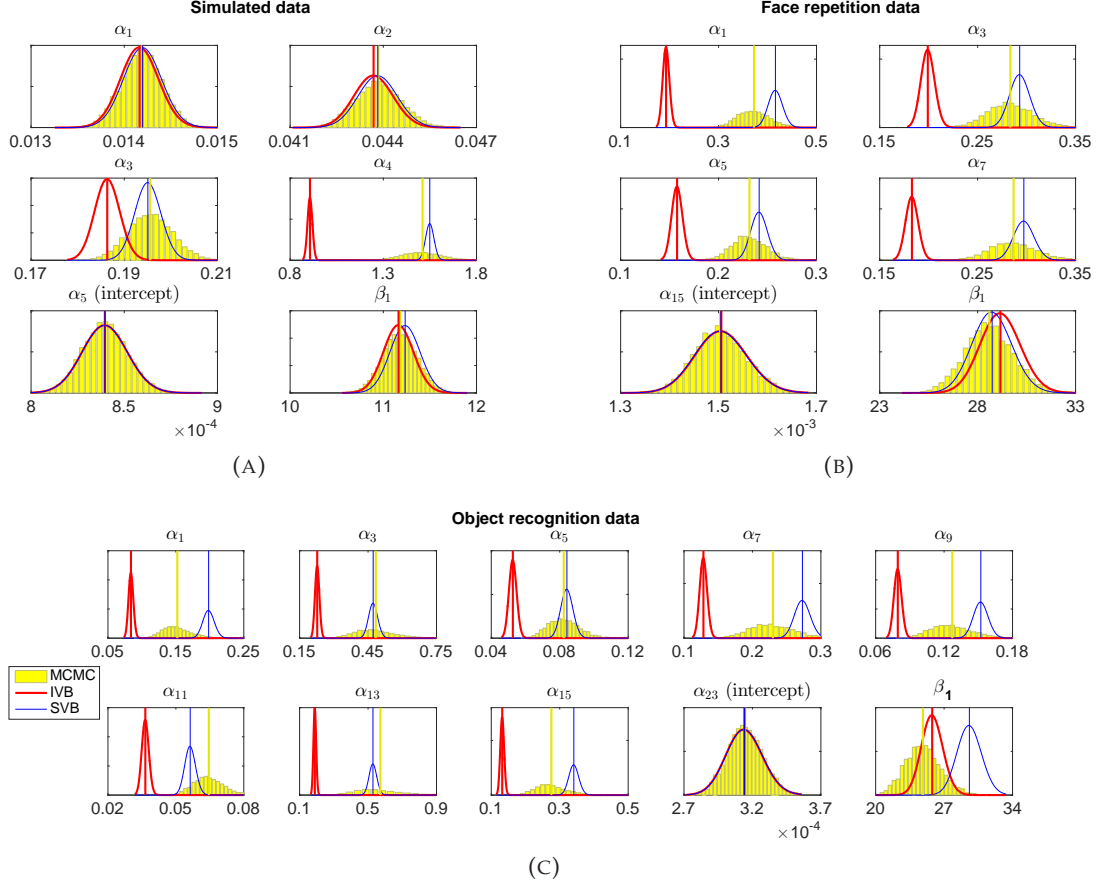


FIGURE 4.6. Estimated posteriors for the spatial hyperparameters for the main regressors and the first AR coefficient for the simulated (a), face repetition (b) and object recognition (c) data by MCMC (histogram), IVB (red) and SVB (blue).

## 5. DISCUSSION AND FUTURE WORK

The PCG based GMRF sampling provides a fast way to perform exact inference in high-dimensional spatial models such as the one considered here for task-fMRI. The presented results show that our methods scale better than SPM's IVB, which is simultaneously shown to generate erroneous results for certain data.

*VB approximation error.* The IVB error comes from two sources. Firstly, factorized VB is known to underestimate posterior variances (Bishop, 2006) and we have seen that fixing  $\alpha$  to the same value in both SPM's VB and the MCMC algorithms results in underestimated  $\mathbf{W}$  posterior standard deviation for IVB for the face repetition data. The same type of posterior

variance underestimation can be seen for the SVB hyperparameters in Figure 4.6, a behavior that is theoretically motivated in Rue et al. (2009) (Appendix A). Secondly, we have seen that IVB tends to underestimate also the mean of  $\alpha$  for many regressors, resulting in the wrong level of smoothing/shrinkage.

The SVB method seems to approximate the exact  $\mathbf{W}$  posterior well in most cases, but the underestimation of hyperparameter variance is occasionally quite large (sometimes also the mean is slightly wrong) which could motivate dropping the VB assumptions entirely and instead optimizing these parameters and perhaps using a Gaussian approximation for the uncertainty.

*Computational improvements and model extensions.* There is room for further improvement of the SVB method, by investigation of the sensitivity to settings like the PCG tolerance  $\delta$  and the number of MC samples  $N_s$  on a larger number of data sets. Different pre-conditioners could be used, for example the robust incomplete Cholesky (Ajiz and Jennings, 1984), and Bolin et al. (2014) discuss other approximations of traces like the one in equation (3.12). The evaluation should be in relation to the output of interest, for example PPMs, to find the optimal balance between accuracy and processing time. In addition, a better online criterion for convergence would be beneficial. In future work we will also investigate if graphics processing units (GPUs) can be used to reduce the processing time further (Eklund et al., 2013, 2014).

Even though the MCMC algorithm is probably not fast enough for the everyday practitioner to run on whole-brain data sets with many conditions, it is orders of magnitude faster than what would be the case without PCG sampling, in which case it would be impractical to run at all. It fulfills an important purpose as the ground truth when evaluating approximate methods such as the VB methods in this paper, and could practically be used on sub-volumes of interest.

For the algorithms presented in this paper we have seen that the convergence rate is determined by the spatial hyperparameters for the least informative regressors. Hence, a small model change that would increase the convergence rates could be to drop the spatial prior on regressors that are very non-informative (for example motion regressors, as previously suggested by Groves et al. (2009)).

In this work, we have for brevity only considered the spatial UGL prior and focused on exact inference using this approximate model. For example, the UGL prior is stationary, isotropic and cannot separate shrinkage and smoothing. As mentioned in the introduction, many alternative, less approximate, priors have been proposed which would be interesting to adopt to the PCG sampling framework for 3D inference. Many of these would be straightforward to implement, for example the anatomically motivated tissue-type AR-priors in Penny et al. (2007), while others would require solving additional computational issues. A particularly interesting alternative would be the Matérn kernel which is a standard choice in spatial statistics. The Matérn kernel defines a Gaussian field (GF) that in general is not Markov, however, Lindgren et al. (2011) give an explicit link between GFs and GM-RFs for the Matérn class such that it can be used in a sparsity exploiting manner, and they

also provide a possible non-stationary Matérn model. Other parallel model improvements would be to include a spatial model also for the noise precision  $\lambda$ , to add a spatial prior also for the probability of activation and to include spatial dependence in the likelihood, which is motivated by Kriegeskorte et al. (2008) who, using a phantom, demonstrate that noise from echoplanar imaging is naturally spatial.

While the voxels are normally equally sized in the  $x$  and  $y$  direction of each slice, this is not necessarily the case relative to the  $z$  direction and one might worry that this leads to more anisotropic data that cannot be accounted for by the UGL prior. A simple solution to this is to resample the voxels to have the same size during the pre-processing, which is in fact what we did for the face repetition data, but not for the object recognition data. The same solution is possible for the second data set, but another straightforward solution would be to replace the UGL prior with a weighted graph-Laplacian (WGL) prior, that is to redefine  $\mathbf{D}_w = \mathbf{G}_w' \mathbf{C} \mathbf{G}_w$  instead of  $\mathbf{D}_w = \mathbf{G}_w' \mathbf{G}_w$  for a diagonal weight matrix  $\mathbf{C}$  that can be chosen based on the Euclidian distances between neighboring voxels as in Harrison et al. (2008a). This solution would also be simply adapted to the PCG framework by just multiplying row  $i$  in  $\mathbf{G}_w$  with  $\sqrt{C_{ii}}$  when using Algorithm 2. Another solution would be change the prior precision matrix in equation (2.3) from  $\alpha_k \mathbf{G}_w' \mathbf{G}_w$  to  $\mathbf{G}_w' \mathbf{C}_k \mathbf{G}_w$  with

$$\mathbf{C}_k = \begin{bmatrix} \alpha_{x,k} \mathbf{I} & & \\ & \alpha_{y,k} \mathbf{I} & \\ & & \alpha_{z,k} \mathbf{I} \end{bmatrix},$$

with  $\alpha_{x,k}$ ,  $\alpha_{y,k}$  and  $\alpha_{z,k}$  being random parameters (to be inferred) that applies to neighboring pairs of voxels in the  $x$ ,  $y$  and  $z$  direction respectively. Adopting such a model to the MCMC or SVB frameworks would however be difficult for computational reasons.

For inference using this or other more advanced priors than the UGL, Gibbs sampling is often not feasible. In these situations, or when the mixing of the Gibbs sampling chain is dissatisfactory, one can instead attempt to perform MCMC using Metropolis-Hastings (MH) steps or Hamiltonian Monte Carlo (HMC) (Duane et al., 1987; Neal, 2011). These methods require the computation of acceptance probabilities based on the joint posterior ratio, as demonstrated for our model in Appendix A. These become computable in practice because the hyperparameters  $\alpha_k$  can be factored out of the determinant of the prior precision matrix as in equation (A.14), but for a more advanced prior the computation of this determinant would be problematic. In a recent work, independent to ours, Teng et al. (2016) use HMC to do inference in the same model as the one in our article. They also make comparisons to SPM's VB for the face repetition data using the 3D prior, and obtain similar results. Their reported computing times seem fast, but since they are using a much smaller model ( $K = 5$ ) and a different implementation (C++), another number of samples, etc., their computing times cannot be directly compared to ours. Nevertheless, HMC seems like a competitive method to the proposed PCG based Gibbs sampling method for MCMC in these models. In a different, recent and independent article, Rad et al. (2016) use a similar PCG based Gibbs sampling method to ours in a similar model for other kinds of neuroscientific data. However,



they use a different spatial edge-preserving Laplace prior, which would be interesting to apply also to fMRI data within this framework.

In order to be able to properly motivate any of the above mentioned model improvements, however, a model selection criterion is required. Penny et al. (2007) used the model evidence lower bound, which is a good alternative, but requires the computation of determinants of precision matrices of size  $KN \times KN$ , which is infeasible in the 3D case. Good approximations of such determinants would therefore be an eligible direction for future research. Model selection criteria based on cross validation or the marginal likelihood could also be explored.

*Multiple comparisons.* There is currently no consensus on how to control for multiple comparisons in Bayesian spatial models for fMRI data. We showed how joint PPMs based on excursion sets can be computed for the MCMC method, but the large data size is currently preventing us from computing the same for the SVB method. The joint PPMs control the family-wise error rate given the spatial model and a threshold, but a separate objective could be to instead control the false discovery rate (FDR) for which a unified framework within large-scale spatial models is provided by Sun et al. (2015). An interesting area of future research would thus be to adopt both these approaches to the 3D spatial modeling of fMRI data when using both the MCMC and SVB method.

*Group analysis.* Another important area of future work is to extend this single subject analysis to the group level. The simplest way to do this is to consider the posterior mean maps from the single subject analysis as spatially processed and use these as input to a voxel-wise Bayesian regression, which is basically what is done in SPM’s Bayesian second level analysis. This approach has several drawbacks, one being that mis-registration between subjects causes activation to be located in different voxels, and can therefore be missed when averaging across subjects. The classic GLM approach “handles” this by using smoothing as a pre-processing step. A more elegant, Bayesian solution is presented in Xu et al. (2009) which explicitly models population level activation centers. A second drawback is that such a procedure discards the posterior uncertainty from the subject level analyses. The most natural way to do the group level analysis would instead be using a Bayesian hierarchical model with a common group level activation working as a latent prior for each subject, and to also model the spatial hyperparameters at the group level. However, estimating such a model would intuitively be very demanding in terms of memory and speed, and an approximate idea is to instead target the mean of the subject-level activations, as in Yue et al. (2014).

### 5.1. In conclusion.

A fast and practical MCMC scheme for exact whole-brain spatial inference in task-fMRI is suggested and implemented. Also, a non-factorizing VB method is developed and shown to give practically the same results, but in shorter time. The methods are compared to the popular factorizing VB method in SPM and shown to scale better with problem size. The comparison with the exact MCMC estimates gives evidence that SPM’s VB can produce false activity estimates in some settings.

## ACKNOWLEDGEMENTS

We thank Will Penny and three anonymous reviewers for helpful comments. This work was funded by Swedish Research Council (Vetenskapsrådet) grant no 2013-5229. David Bolin was also supported by the Knut and Alice Wallenberg foundation.

## REFERENCES

- (2002). SPM. Wellcome Department of Imaging Neuroscience, Available at <http://www.fil.ion.ucl.ac.uk/spm/software>.
- Ajiz, M. A. and Jennings, A. (1984). A robust incomplete Cholesky-conjugate gradient algorithm. *Int. J. Numer. Meth. Engng.*, 20:949–966.
- Amestoy, P. R., Davis, T. A., and Duff, I. S. (1996). An Approximate Minimum Degree Ordering Algorithm. *SIAM. J. Matrix Anal. & Appl.*, 17(4):886–905.
- Barrett, R., Berry, M. W., Chan, T. F., Demmel, J., Donato, J., Dongarra, J., Eijkhout, V., Pozo, R., Romine, C., and Van der Vorst, H. (1994). *Templates for the solution of linear systems: building blocks for iterative methods*, volume 43. Siam.
- Besag, J. (1974). Spatial Interaction and the Statistical Analysis of Lattice Systems. *Journal of the Royal Statistical Society. Series B (Methodological)*, 36(2):192–236.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Bolin, D. and Lindgren, F. (2015). Excursion and contour uncertainty regions for latent Gaussian models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 77(1):85–106.
- Bolin, D., Wallin, J., and Lindgren, F. (2014). Multivariate latent Gaussian random field mixture models. *Preprint Dept. Math. Sci., Chalmers University of Technology and Göteborg University 2014:1*.
- Chaari, L., Vincent, T., Forbes, F., Dojat, M., and Ciuciu, P. (2013). Fast Joint Detection-Estimation of Evoked Brain Activity in Event-Related fMRI Using a Variational Approach Lotfi. *IEEE Transactions on Medical Imaging*, 32(5):821–837.
- Chan, K. S. and Ledolter, J. (1995). Monte Carlo EM Estimation for Time Series Models Involving Counts. *Journal of the American Statistical association*, 90(429):242–252.
- Delyon, B., Lavielle, M., and Moulines, E. (1999). Convergence of a stochastic approximation version of the EM algorithm. *Annals of Statistics*, 27(1):94–128.
- Duane, S., Kennedy, A. D., Pendleton, B. J., and Roweth, D. (1987). Hybrid Monte Carlo. *Physics letters B*, 195(2):216–222.
- Eklund, A., Dufort, P., Forsberg, D., and LaConte, S. (2013). Medical image processing on the GPU - Past, present and future. *Medical Image Analysis*, 17:1073–1094.
- Eklund, A., Dufort, P., Villani, M., and LaConte, S. (2014). BROCCOLI: Software for fast fMRI analysis on many-core CPUs and GPUs. *Frontiers in Neuroinformatics*, 8:24.
- Eklund, A., Nichols, T. E., and Knutsson, H. (2016). Cluster failure: why fMRI inferences for spatial extent have inflated false positive rates. *Proceedings of the National Academy of Sciences*, 113(28):7900–7905.

- Friston, K. J., Holmes, a. P., Worsley, K. J., Poline, J.-P., Frith, C. D., and Frackowiak, R. S. J. (1995). Statistical parametric maps in functional imaging: A general linear approach. *Human Brain Mapping*, 2(4):189–210.
- Friston, K. J. and Penny, W. (2003). Posterior probability maps and SPMs. *NeuroImage*, 19(3):1240–1249.
- Groves, A. R., Chappell, M. A., and Woolrich, M. W. (2009). Combined spatial and non-spatial prior for inference on MRI time-series. *NeuroImage*, 45(3):795–809.
- Gunawan, D., Tran, M., and Kohn, R. (2016). Fast Inference for Intractable Likelihood Problems using Variational Bayes. *preprint: <http://hdl.handle.net/2123/14594>*.
- Hanson, S. J., Matsuka, T., and Haxby, J. V. (2004). Combinatorial codes in ventral temporal lobe for object recognition: Haxby (2001) revisited: is there a "face" area? *NeuroImage*, 23(1):156–166.
- Harrison, L. M. and Green, G. G. R. (2010). A Bayesian spatiotemporal model for very large data sets. *NeuroImage*, 50(3):1126–1141.
- Harrison, L. M., Penny, W., Daunizeau, J., and Friston, K. J. (2008a). Diffusion-based spatial priors for functional magnetic resonance images. *NeuroImage*, 41(2):408–423.
- Harrison, L. M., Penny, W., Flandin, G., Ruff, C. C., Weiskopf, N., and Friston, K. J. (2008b). Graph-partitioned spatial priors for functional magnetic resonance images. *NeuroImage*, 43(4):694–707.
- Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P. (2001). Distributed and Overlapping Representations of Faces and Objects in Ventral Temporal Cortex. *Science*, 293(5539):2425–2430.
- Henson, R., Shallice, T., Gorno-Tempini, M. L., and Dolan, R. (2002). Face repetition effects in implicit and explicit memory tests as measured by fMRI. *Cerebral Cortex*, 12:178–186.
- Kingma, D. P. and Welling, M. (2014). Auto-Encoding Variational Bayes. *arXiv:1312.6114v10*.
- Kriegeskorte, N., Bodurka, J., and Bandettini, P. (2008). Artifactual time-course correlations in echo-planar fMRI with implications for studies of brain function. *International Journal of Imaging Systems and Technology*, 18(5-6):345–349.
- Lindgren, F., Rue, H., and Lindström, J. (2011). An explicit link between Gaussian fields and Gaussian Markov random fields: The SPDE approach. *Journal of the Royal Statistical Society Series B*, 73(4):423–498.
- Manteuffel, T. a. (1980). An incomplete factorization technique for positive definite linear systems. *Mathematics of Computation*, 34(150):473–473.
- Musgrove, D. R., Hughes, J., and Eberly, L. E. (2016). Fast, fully Bayesian spatiotemporal inference for fMRI data. *Biostatistics*, 17(2):291–303.
- Neal, R. M. (2011). MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, pages 113–162.
- O’Toole, A. J., Jiang, F., Abdi, H., and Haxby, J. V. (2005). Partially Distributed Representations of Objects and Faces in Ventral Temporal Cortex. *Journal of Cognitive Neuroscience*, 17(4):580–590.
- Papandreou, G. and Yuille, A. (2010). Gaussian sampling by local perturbations. *Advances in Neural Information Processing Systems* 23, 90(8):1858–1866.

- Papandreou, G. and Yuille, A. L. (2011). Efficient variational inference in large-scale Bayesian compressed sensing. *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, pages 1332–1339.
- Penny, W. and Flandin, G. (2005). Bayesian analysis of fMRI data with spatial priors. In *Proceedings of the Joint Statistical Meeting (JSM). American Statistical Association*.
- Penny, W., Flandin, G., and Trujillo-Barreto, N. (2007). Bayesian comparison of spatially regularised general linear models. *Human Brain Mapping*, 28(4):275–293.
- Penny, W., Kiebel, S., and Friston, K. (2003). Variational Bayesian inference for fMRI time series. *NeuroImage*, 19(3):727–741.
- Penny, W. D., Trujillo-Barreto, N., and Flandin, G. (2005a). Bayesian analysis of single-subject fMRI data: SPM implementation. *Technical report, Wellcome Department of Imaging Neuroscience*.
- Penny, W. D., Trujillo-Barreto, N. J., and Friston, K. J. (2005b). Bayesian fMRI time series analysis with spatial priors. *NeuroImage*, 24(2):350–362.
- Poldrack, R. A., Barch, D. M., Mitchell, J. P., Wager, T. D., Wagner, A. D., Devlin, J. T., Cumba, C., Koyejo, O., and Milham, M. P. (2013). Toward open sharing of task-based fMRI data: the OpenfMRI project. *Frontiers in neuroinformatics*, 7:12.
- Rad, K. R., Machado, T. A., and Paninski, L. (2016). Robust and scalable Bayesian analysis of spatial neural tuning function data. *arXiv:1606.07845v1*.
- Risser, L., Vincent, T., Forbes, F., Idier, J., and Ciuciu, P. (2011). Min-max Extrapolation Scheme for Fast Estimation of 3D Potts Field Partition Functions. Application to the Joint Detection-Estimation of Brain Activity in fMRI. *Journal of Signal Processing Systems*, 65(3):325–338.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*. CRC Press.
- Rue, H., Martino, S., and Chopin, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximation. *Journal of the Royal Statistical Society, Series B*, 71(2):319–392.
- Smith, M. and Fahrmeir, L. (2007). Spatial Bayesian variable selection with application to functional magnetic resonance imaging. *Journal of the American Statistical Association*, 102:417–431.
- Sun, W., Reich, B. J., Tony Cai, T., Guindani, M., and Schwartzman, A. (2015). False discovery control in large-scale spatial multiple testing. *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 77(1):59–83.
- Teng, M., Johnson, T., and Nathoo, F. (2016). A Comparison of Variational Bayes and Hamiltonian Monte Carlo for Bayesian fMRI Time Series Analysis with Spatial Priors. *arXiv:1609.02123v1*.
- Thirion, B., Varoquaux, G., Dohmatob, E., and Poline, J. B. (2014). Which fMRI clustering gives good brain parcellations? *Frontiers in Neuroscience*, 8:167.
- Vincent, T., Risser, L., and Ciuciu, P. (2010). Spatially Adaptive Mixture Modeling for Analysis of fMRI Time Series. *IEEE transactions on medical imaging*, 29(4):1059–1074.

- Woods, J. W. (1972). Two-Dimensional Discrete Markovian Fields. *IEEE Transactions on Information Theory*, 18(2):232–240.
- Woolrich, M. W., Jenkinson, M., Brady, J. M., and Smith, S. M. (2004). Fully Bayesian spatio-temporal modeling of fMRI data. *IEEE transactions on medical imaging*, 23(2):213–31.
- Xu, L., Johnson, T. D., Nichols, T. E., and Nee, D. E. (2009). Modeling inter-subject variability in fMRI activation location: A bayesian hierarchical spatial model. *Biometrics*, 65(4):1041–1051.
- Yue, Y. R., Lindquist, M. a., Bolin, D., Lindgren, F., Simpson, D., and Rue, H. (2014). A Bayesian General Linear Modeling Approach to Slice-wise fMRI Data Analysis. *Preprint*.
- Zhang, L., Guindani, M., Versace, F., Engelmann, J. M., and Vannucci, M. (2016). A spatiotemporal nonparametric Bayesian model of multi-subject fMRI data. *Ann. Appl. Stat.*, 10(2):638–666.
- Zhang, L., Guindani, M., Versace, F., and Vannucci, M. (2014). A spatio-temporal non-parametric Bayesian variable selection model of fMRI data for clustering correlated time courses. *NeuroImage*, 95:162–175.

## APPENDIX A. DERIVATION OF FULL CONDITIONAL POSTERiors FOR THE MCMC ALGORITHM

This section is divided into two pieces, the first handling the case when the noise in each voxel is modeled as i.i.d. over time and the second handling the case with auto-regressive temporal noise. The first part is shorter and basically contains all the concepts needed for the second part. It should be noted that adding a temporal model does not add much to the time complexity as long as  $P < K$ , which is normally the case. We also provide an expression for computing the joint posterior ratio.

### A.1. i.i.d. noise model.

*The likelihood.* The likelihood in (2.2) can be expressed in logs as

$$\begin{aligned} \log p(\mathbf{Y}|\mathbf{W}, \lambda) &= \frac{T}{2} \sum_{n=1}^N \log(\lambda_n) \\ &\quad - \frac{1}{2} \sum_{n=1}^N \lambda_n [\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n} - 2\mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n}] + \text{const}, \end{aligned} \tag{A.1}$$

where we have omitted everything that is constant with respect to the parameters. Since  $\mathbf{Y}$  and  $\mathbf{X}$  are data that will not change during the MCMC algorithm, quantities as  $\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n}$ ,  $\mathbf{Y}'_{\cdot,n} \mathbf{X}$  and  $\mathbf{X}' \mathbf{X}$  can be effectively pre-computed, removing the time dimension from the likelihood which leads to significant speed up. This is similar to what is done in Penny et al. (2005a).

Full conditional posterior of  $\mathbf{W}$ .

$$\begin{aligned}
\log p(\mathbf{W}|\mathbf{Y}, \lambda, \alpha) &= \log p(\mathbf{Y}|\mathbf{W}, \lambda) + \log p(\mathbf{W}|\alpha) + \text{const} \quad (\text{A.2}) \\
&= -\frac{1}{2} \left[ \sum_{n=1}^N \lambda_n (\mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n} - 2 \mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}_{\cdot,n}) + \sum_{k=1}^K \mathbf{W}_{k,\cdot} \alpha_k \mathbf{D}_w \mathbf{W}'_{k,\cdot} \right] + \text{const} \\
&= -\frac{1}{2} \mathbf{w}'_r \tilde{\mathbf{B}} \mathbf{w}_r + \mathbf{b}'_w \mathbf{w}_r + \text{const}, \\
\mathbf{b}_w &= \text{vec}(\text{diag}(\lambda) \mathbf{Y}' \mathbf{X}), \\
\tilde{\mathbf{B}} &= \mathbf{X}' \mathbf{X} \otimes \text{diag}(\lambda) + \text{diag}(\alpha) \otimes \mathbf{D}_w,
\end{aligned}$$

or equivalently  $\mathbf{w}_r|\mathbf{Y}, \lambda, \alpha \sim \mathcal{N}(\tilde{\mathbf{B}}^{-1} \mathbf{b}_w, \tilde{\mathbf{B}}^{-1})$ , where  $\mathbf{w}_r = \text{vec}(\mathbf{W}')$ .

Full conditional posterior of  $\lambda$ .

$$\begin{aligned}
\log p(\lambda|\mathbf{Y}, \mathbf{W}, \alpha) &= \log p(\mathbf{Y}|\mathbf{W}, \lambda) + \log p(\lambda) + \text{const} \quad (\text{A.3}) \\
&= \frac{T}{2} \sum_{n=1}^N \log(\lambda_n) - \frac{1}{2} \sum_{n=1}^N \lambda_n [\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n} - 2 \mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}'_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n}] \\
&\quad + (u_2 - 1) \sum_{n=1}^N \log(\lambda_n) - \sum_{n=1}^N \frac{\lambda_n}{u_1} + \text{const} \\
&= (\tilde{u}_2 - 1) \sum_{n=1}^N \log(\lambda_n) - \sum_{n=1}^N \frac{\lambda_n}{\tilde{u}_{1n}} + \text{const}, \\
\frac{1}{\tilde{u}_{1n}} &= \frac{1}{2} (\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n} - 2 \mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}'_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n}) + \frac{1}{u_1}, \\
\tilde{u}_2 &= \frac{T}{2} + u_2,
\end{aligned}$$

so  $\lambda_n|\mathbf{Y}, \mathbf{W} \sim \text{Ga}(\tilde{u}_{1n}, \tilde{u}_2)$  for all  $n$ .

Full conditional posterior of  $\alpha$ .

$$\begin{aligned}
\log p(\alpha|\mathbf{Y}, \mathbf{W}, \lambda) &= \log p(\mathbf{Y}|\mathbf{W}, \lambda) + \log p(\mathbf{W}|\alpha) + \log p(\alpha) + \text{const} \quad (\text{A.4}) \\
&= \frac{N}{2} \sum_{k=1}^K \log(\alpha_k) - \frac{1}{2} \sum_{k=1}^K \mathbf{W}_{k,\cdot} \alpha_k \mathbf{D}_w \mathbf{W}'_{k,\cdot} \\
&\quad + (q_2 - 1) \sum_{k=1}^K \log(\alpha_k) - \sum_{k=1}^K \frac{\alpha_k}{q_1} + \text{const} \\
&= (\tilde{q}_2 - 1) \sum_{k=1}^K \log(\alpha_k) - \sum_{k=1}^K \frac{\alpha_k}{\tilde{q}_{1k}} + \text{const}, \\
\frac{1}{\tilde{q}_{1k}} &= \frac{1}{2} \mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}'_{k,\cdot} + \frac{1}{q_1}, \\
\tilde{q}_2 &= \frac{N}{2} + q_2,
\end{aligned}$$

so  $\alpha_k|\mathbf{W} \sim \text{Ga}(\tilde{q}_{1k}, \tilde{q}_2)$  for all  $k$ .

## A.2. Temporal noise model.

The derivations of the full conditionals for the temporal noise model follows the same pattern



as for the i.i.d. model, why we leave out some steps for brevity. Also note that the form of the posterior for the AR coefficients  $\mathbf{A}$  will be very similar to that of the regression coefficients  $\mathbf{W}$ , and the hyperparameter  $\beta$  has the same form as  $\alpha$ . The permutation matrices  $\mathbf{H}_w$  and  $\mathbf{H}_a$  are defined as in Penny et al. (2007) such that  $\text{vec}(\mathbf{W}) = \mathbf{H}_w \text{vec}(\mathbf{W}')$  and  $\text{vec}(\mathbf{A}) = \mathbf{H}_a \text{vec}(\mathbf{A}')$ .

*The likelihood.* Using the temporal noise model, the likelihood can be expressed as

$$p(\mathbf{Y}|\mathbf{W}, \mathbf{A}, \lambda) \propto \prod_{t=P+1}^T \prod_{n=1}^N \mathcal{N}\left(Y_{tn} - \mathbf{x}_{t,\cdot} \mathbf{W}_{\cdot,n}; (\mathbf{d}_{tn} - \tilde{\mathbf{X}}_t \mathbf{W}_{\cdot,n})' \mathbf{A}_{\cdot,n}, \lambda_n^{-1}\right), \quad (\text{A.5})$$

where  $\tilde{\mathbf{X}}_t$  is a  $P \times K$  matrix containing the  $P$  rows of the design matrix prior to time point  $t$  and  $\mathbf{d}_{tn}$  is  $P \times 1$  and similarly contains the  $P$  values of  $\mathbf{Y}_{\cdot,n}$  just before  $t$ . Note that we condition on the first  $P$  time points for simplicity. The non-constant part of the log-likelihood can be written as

$$\begin{aligned} \log p(\mathbf{Y}|\cdot) &= \frac{T-P}{2} \sum_{n=1}^N \log(\lambda_n) - \frac{1}{2} \sum_{t=P+1}^T \sum_{n=1}^N \left[ (Y_{tn} - \mathbf{x}_{t,\cdot} \mathbf{W}_{\cdot,n}) - (\mathbf{d}_{tn} - \tilde{\mathbf{X}}_t \mathbf{W}_{\cdot,n})' \mathbf{A}_{\cdot,n} \right]' \lambda_n \\ &\quad \times \left[ (Y_{tn} - \mathbf{x}_{t,\cdot} \mathbf{W}_{\cdot,n}) - (\mathbf{d}_{tn} - \tilde{\mathbf{X}}_t \mathbf{W}_{\cdot,n})' \mathbf{A}_{\cdot,n} \right] + \text{const} \\ &= \frac{T-P}{2} \sum_{n=1}^N \log(\lambda_n) - \frac{1}{2} \sum_{n=1}^N \left[ (\mathbf{Y}_{\cdot,n} - \mathbf{X} \mathbf{W}_{\cdot,n}) - (\mathbf{d}_n - \tilde{\mathbf{X}} \mathbf{W}_{\cdot,n})' \mathbf{A}_{\cdot,n} \right]' \lambda_n \\ &\quad \times \left[ (\mathbf{Y}_{\cdot,n} - \mathbf{X} \mathbf{W}_{\cdot,n}) - (\mathbf{d}_n - \tilde{\mathbf{X}} \mathbf{W}_{\cdot,n})' \mathbf{A}_{\cdot,n} \right] + \text{const}. \end{aligned} \quad (\text{A.6})$$

In the last expression the sums and indexing with respect to  $t$  has been removed. Since none of the parameters depend on the time dimension, this expression can be rewritten so that sums and matrix multiplications over the time dimension can be isolated to the data  $(\mathbf{Y}, \mathbf{X}, \mathbf{d}, \tilde{\mathbf{X}})$ , so that they can be pre-computed outside the Gibbs algorithm, leading to a higher computational efficiency as in Penny et al. (2005a). Note that the size of  $\mathbf{d}_n$  is  $P \times (T-P)$  and the size of  $\tilde{\mathbf{X}}$  is  $P \times (T-P) \times K$ . Matrix multiplications including the 3-dimensional matrix  $\tilde{\mathbf{X}}$  and other tensors will be carried out over the appropriate dimension in what follows, even if this is not stated explicitly. The log-likelihood can now be rewritten as

$$\begin{aligned} \log p(\mathbf{Y}|\cdot) &= \frac{T-P}{2} \sum_{n=1}^N \log(\lambda_n) \\ &\quad - \frac{1}{2} \sum_{n=1}^N \lambda_n \left[ \mathbf{Y}_{\cdot,n}' \mathbf{Y}_{\cdot,n} - 2 \mathbf{Y}_{\cdot,n}' \mathbf{X} \mathbf{W}_{\cdot,n} + \mathbf{W}_{\cdot,n}' \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n} - 2 \mathbf{Y}_{\cdot,n}' \mathbf{d}_n' \mathbf{A}_{\cdot,n} \right. \\ &\quad \left. + \mathbf{A}_{\cdot,n}' \mathbf{d}_n \mathbf{d}_n' \mathbf{A}_{\cdot,n} + \mathbf{W}_{\cdot,n}' \mathbf{B}_n' \mathbf{A}_{\cdot,n} + \mathbf{A}_{\cdot,n}' \mathbf{B}_n \mathbf{W}_{\cdot,n} - \mathbf{W}_{\cdot,n}' (\mathbf{R} \mathbf{A}_{\cdot,n} + (\mathbf{R} \mathbf{A}_{\cdot,n})') \mathbf{W}_{\cdot,n} \right. \\ &\quad \left. - \mathbf{A}_{\cdot,n}' (\mathbf{D}_n \mathbf{W}_{\cdot,n} + (\mathbf{D}_n \mathbf{W}_{\cdot,n})') \mathbf{A}_{\cdot,n} + \mathbf{W}_{\cdot,n}' (\mathbf{A}_{\cdot,n}' \mathbf{S} \mathbf{A}_{\cdot,n}) \mathbf{W}_{\cdot,n} \right] + \text{const}, \end{aligned} \quad (\text{A.7})$$

where

$$\begin{aligned}\mathbf{B}_n &= \mathbf{Y}'_{\cdot,n} \tilde{\mathbf{X}} + \mathbf{d}_n \mathbf{X} \quad \text{is of size } P \times K, \\ \mathbf{R} &= \mathbf{X}' \tilde{\mathbf{X}} \quad \text{is of size } K \times K \times P, \\ \mathbf{D}_n &= \mathbf{d}_n \tilde{\mathbf{X}} \quad \text{is of size } P \times K \times P, \\ \mathbf{S} &= \tilde{\mathbf{X}} \tilde{\mathbf{X}} \quad \text{is of size } P \times K \times K \times P.\end{aligned}$$

Full conditional posterior of  $\mathbf{W}$ .

$$\log p(\mathbf{W}|\mathbf{Y}, \cdot) = -\frac{1}{2} \mathbf{w}'_r \tilde{\mathbf{B}} \mathbf{w}_r + \mathbf{b}'_w \mathbf{w}_r + \text{const}, \quad (\text{A.8})$$

$$\begin{aligned}\mathbf{b}_w &= \text{vec} \left( \begin{bmatrix} \vdots \\ \lambda_n (\mathbf{Y}'_{\cdot,n} \mathbf{X} - \mathbf{A}'_{\cdot,n} \mathbf{B}_n + \mathbf{A}'_{\cdot,n} \mathbf{D}_n \mathbf{A}_{\cdot,n}) \\ \vdots \end{bmatrix}_{n \in \{1, \dots, N\}} \right), \\ \tilde{\mathbf{B}} &= \mathbf{H}'_w \text{blkdiag}_{n \in \{1, \dots, N\}} [\lambda_n (\mathbf{X}' \mathbf{X} - \mathbf{R} \mathbf{A}_{\cdot,n} - (\mathbf{R} \mathbf{A}_{\cdot,n})' + \mathbf{A}'_{\cdot,n} \mathbf{S} \mathbf{A}_{\cdot,n})] \mathbf{H}_w \\ &\quad + \text{diag}(\boldsymbol{\alpha}) \otimes \mathbf{D}_w,\end{aligned}$$

where  $\text{blkdiag}_{n \in \{1, \dots, N\}} [\mathbf{C}_n]$  is a  $KN \times KN$  block diagonal matrix with the  $K \times K$  matrix  $\mathbf{C}_n$  as the  $n$ th block. So  $\mathbf{w}_r | \mathbf{Y}, \cdot \sim \mathcal{N}(\tilde{\mathbf{B}}^{-1} \mathbf{b}_w, \tilde{\mathbf{B}}^{-1})$ .

Full conditional posterior of  $\mathbf{A}$ .

$$\log p(\mathbf{A}|\mathbf{Y}, \cdot) = -\frac{1}{2} \mathbf{a}'_r \tilde{\mathbf{J}} \mathbf{a}_r + \mathbf{b}'_a \mathbf{a}_r + \text{const}, \quad (\text{A.9})$$

$$\begin{aligned}\mathbf{b}_a &= \text{vec} \left( \begin{bmatrix} \vdots \\ \lambda_n (\mathbf{Y}'_{\cdot,n} \mathbf{d}'_n - \mathbf{W}'_{\cdot,n} \mathbf{B}'_n + \mathbf{W}'_{\cdot,n} \mathbf{R} \mathbf{W}_{\cdot,n}) \\ \vdots \end{bmatrix}_{n \in \{1, \dots, N\}} \right), \\ \tilde{\mathbf{J}} &= \mathbf{H}'_a \text{blkdiag}_{n \in \{1, \dots, N\}} [\lambda_n (\mathbf{d}_n \mathbf{d}'_n - \mathbf{D}_n \mathbf{W}_{\cdot,n} - (\mathbf{D}_n \mathbf{W}_{\cdot,n})' + \mathbf{W}'_{\cdot,n} \mathbf{S} \mathbf{W}_{\cdot,n})] \mathbf{H}_a \\ &\quad + \text{diag}(\boldsymbol{\beta}) \otimes \mathbf{D}_a,\end{aligned}$$

so  $\mathbf{a}_r | \mathbf{Y}, \cdot \sim \mathcal{N}(\tilde{\mathbf{J}}^{-1} \mathbf{b}_a, \tilde{\mathbf{J}}^{-1})$ .

Full conditional posterior of  $\lambda$ .

$$\log p(\lambda|\mathbf{Y}, \cdot) = (\tilde{u}_2 - 1) \sum_{n=1}^N \log(\lambda_n) - \sum_{n=1}^N \frac{\lambda_n}{\tilde{u}_{1n}} + \text{const}, \quad (\text{A.10})$$

$$\begin{aligned}\frac{1}{\tilde{u}_{1n}} &= \frac{1}{u_1} + \frac{1}{2} [\mathbf{Y}'_{\cdot,n} \mathbf{Y}_{\cdot,n} - 2 \mathbf{Y}'_{\cdot,n} \mathbf{X} \mathbf{W}_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{X}' \mathbf{X} \mathbf{W}_{\cdot,n} - 2 \mathbf{Y}'_{\cdot,n} \mathbf{d}'_n \mathbf{A}_{\cdot,n} \\ &\quad + \mathbf{A}'_{\cdot,n} \mathbf{d}_n \mathbf{d}'_n \mathbf{A}_{\cdot,n} + \mathbf{W}'_{\cdot,n} \mathbf{B}'_n \mathbf{A}_{\cdot,n} + \mathbf{A}'_{\cdot,n} \mathbf{B}_n \mathbf{W}_{\cdot,n} - \mathbf{W}'_{\cdot,n} (\mathbf{R} \mathbf{A}_{\cdot,n} + (\mathbf{R} \mathbf{A}_{\cdot,n})') \mathbf{W}_{\cdot,n} \\ &\quad - \mathbf{A}'_{\cdot,n} (\mathbf{D}_n \mathbf{W}_{\cdot,n} + (\mathbf{D}_n \mathbf{W}_{\cdot,n})') \mathbf{A}_{\cdot,n} + \mathbf{W}'_{\cdot,n} (\mathbf{A}'_{\cdot,n} \mathbf{S} \mathbf{A}_{\cdot,n}) \mathbf{W}_{\cdot,n}], \\ \tilde{u}_2 &= \frac{T - P}{2} + u_2,\end{aligned}$$

so  $\lambda_n | \mathbf{Y}, \cdot \sim Ga(\tilde{u}_{1n}, \tilde{u}_2)$  for all  $n$ .

Full conditional posterior of  $\alpha$ .

$$\begin{aligned} \log p(\alpha | \cdot) &= (\tilde{q}_2 - 1) \sum_{k=1}^K \log(\alpha_k) - \sum_{k=1}^K \frac{\alpha_k}{\tilde{q}_{1k}} + const, \\ \frac{1}{\tilde{q}_{1k}} &= \frac{1}{2} \mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}' + \frac{1}{q_1}, \\ \tilde{q}_2 &= \frac{N}{2} + q_2, \end{aligned} \quad (\text{A.11})$$

so  $\alpha_k | \cdot \sim Ga(\tilde{q}_{1k}, \tilde{q}_2)$  for all  $k$ . This is exactly the same as for the i.i.d. case.

Full conditional posterior of  $\beta$ .

$$\begin{aligned} \log p(\beta | \cdot) &= (\tilde{r}_2 - 1) \sum_{p=1}^P \log(\beta_p) - \sum_{p=1}^P \frac{\beta_p}{\tilde{r}_{1p}} + const, \\ \frac{1}{\tilde{r}_{1p}} &= \frac{1}{2} \mathbf{A}_{p,\cdot} \mathbf{D}_a \mathbf{A}_{p,\cdot}' + \frac{1}{r_1}, \\ \tilde{r}_2 &= \frac{N}{2} + r_2, \end{aligned} \quad (\text{A.12})$$

so  $\beta_p | \cdot \sim Ga(\tilde{r}_{1p}, \tilde{r}_2)$  for all  $p$ .

*Joint posterior ratio.* The ratio of the joint posterior  $p(\mathbf{W}, \mathbf{A}, \lambda, \alpha, \beta | \mathbf{Y})$  evaluated for two different sets of parameter values,  $\{\mathbf{W}, \mathbf{A}, \lambda, \alpha, \beta\}$  and  $\{\mathbf{W}^*, \mathbf{A}^*, \lambda^*, \alpha^*, \beta^*\}$ , can be used to compare the posterior density in different points, even when the normalized joint posterior itself is not available in closed form. The ratio can be computed as the ratio of the unnormalized joint posterior  $\tilde{p}(\mathbf{W}, \mathbf{A}, \lambda, \alpha, \beta | \mathbf{Y})$  defined by

$$\begin{aligned} \log \tilde{p}(\mathbf{W}, \mathbf{A}, \lambda, \alpha, \beta | \mathbf{Y}) &= \log \tilde{p}(\mathbf{Y} | \cdot) + \log \tilde{p}(\mathbf{W} | \alpha) + \log \tilde{p}(\mathbf{A} | \beta) + \log p(\lambda) \\ &\quad + \log p(\alpha) + \log p(\beta), \end{aligned} \quad (\text{A.13})$$

where  $\log \tilde{p}(\mathbf{Y} | \cdot)$  is defined as in equation (A.7) without the constant part,

$$\begin{aligned} \log \tilde{p}(\mathbf{W} | \alpha) &= \sum_{k=1}^K \left[ \log(|\alpha_k \mathbf{D}_w|^{1/2}) - \frac{1}{2} \alpha_k \mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}' \right] + const \\ &= \sum_{k=1}^K \left[ \frac{N}{2} \log(\alpha_k) - \frac{1}{2} \alpha_k \mathbf{W}_{k,\cdot} \mathbf{D}_w \mathbf{W}_{k,\cdot}' \right] \end{aligned} \quad (\text{A.14})$$

and  $\log \tilde{p}(\mathbf{A} | \beta)$  is defined correspondingly.

## APPENDIX B. DERIVATION OF SPATIAL VB POSTERiors

From equation (3.8) one sees that the SVB approximate posteriors can be derived from the full conditionals in Appendix A, since

$$\log q(\theta_j) = E_{\theta_{-j}} [\log p(\mathbf{Y}, \theta)] + const = E_{\theta_{-j}} [\log p(\theta_j | \mathbf{Y}, \theta_{-j})] + const. \quad (\text{B.1})$$

The only difference is that we must take the expectation with respect to all other parameters. It turns out that the dependencies on  $\lambda, \alpha$  and  $\beta$  are always linear, so we can just replace these

with their expectations  $\bar{\lambda}$ ,  $\bar{\alpha}$  and  $\bar{\beta}$  under respective SVB gamma posterior. For example,

$$\bar{\alpha}_k = E_{\alpha} [\alpha_k] = \bar{q}_{1k}^{SVB} \cdot \bar{q}_2^{SVB}. \quad (\text{B.2})$$

For dependencies on  $\mathbf{W}$  and  $\mathbf{A}$  we use MC approximations as in (3.13) to compute the expectations, by simulating  $N_s$  samples of  $\mathbf{W}$  and  $\mathbf{A}$  respectively. For brevity, we do not derive all approximate posteriors here, only  $q(\mathbf{W})$  for the temporal model as an example.

$$\begin{aligned} \log q(\mathbf{w}_r) &= E_{\mathbf{A}, \lambda, \alpha, \beta} [\log p(\mathbf{W} | \mathbf{Y}, \cdot)] + \text{const} \\ &= -\frac{1}{2} \mathbf{w}_r' E_{\mathbf{A}, \lambda, \alpha, \beta} [\tilde{\mathbf{B}}] \mathbf{w}_r + E_{\mathbf{A}, \lambda, \alpha, \beta} [\mathbf{b}_w] \mathbf{w}_r + \text{const}, \end{aligned} \quad (\text{B.3})$$

with  $\tilde{\mathbf{B}}$  and  $\mathbf{b}_w$  from equation (A.8). The expectations are computed as

$$\mathbf{b}_w^{SVB} = E_{\mathbf{A}, \lambda, \alpha, \beta} [\mathbf{b}_w] \quad (\text{B.4})$$

$$\begin{aligned} &= \frac{1}{N_s} \sum_{j=1}^{N_s} \text{vec} \left( \begin{bmatrix} \bar{\lambda}_n \left( \mathbf{Y}_{:,n}' \mathbf{X} - \mathbf{A}_{:,n}^{(j)'} \mathbf{B}_n + \mathbf{A}_{:,n}^{(j)'} \mathbf{D}_n \mathbf{A}_{:,n}^{(j)} \right) \\ \vdots \\ \bar{\lambda}_n \left( \mathbf{Y}_{:,n}' \mathbf{X} - \mathbf{A}_{:,n}^{(j)'} \mathbf{B}_n + \mathbf{A}_{:,n}^{(j)'} \mathbf{D}_n \mathbf{A}_{:,n}^{(j)} \right) \end{bmatrix}_{n \in \{1, \dots, N\}} \right) \\ \tilde{\mathbf{B}}^{SVB} &= E_{\mathbf{A}, \lambda, \alpha, \beta} [\tilde{\mathbf{B}}] \\ &= \frac{1}{N_s} \sum_{j=1}^{N_s} \mathbf{H}_w' \text{blkdiag}_{n \in \{1, \dots, N\}} \left[ \bar{\lambda}_n \left( \mathbf{X}' \mathbf{X} - \mathbf{R} \mathbf{A}_{:,n}^{(j)} - \left( \mathbf{R} \mathbf{A}_{:,n}^{(j)} \right)' + \mathbf{A}_{:,n}^{(j)'} \mathbf{S} \mathbf{A}_{:,n}^{(j)} \right) \right] \mathbf{H}_w \\ &\quad + \text{diag}(\bar{\alpha}) \otimes \mathbf{D}_w. \end{aligned} \quad (\text{B.5})$$

$$\text{So } q(\mathbf{W}) \sim \mathcal{N} \left( (\tilde{\mathbf{B}}^{SVB})^{-1} \mathbf{b}_w^{SVB}, (\tilde{\mathbf{B}}^{SVB})^{-1} \right).$$

#### APPENDIX C. CONVERGENCE AND IMPLEMENTATION DETAILS

Our experience from running the MCMC and VB methods is that the convergence of the algorithms is largely governed by that of the spatial hyperparameters  $\alpha$  and  $\beta$ . Figure C.1a shows the relative error by iteration number for the hyperparameters for IVB, for the presented slice of the face repetition data. The relative error is here defined as compared to the final value after a long run (200 iterations) and if  $\bar{\alpha}_k^{[j]}$  denotes the value of  $\bar{\alpha}_k$  after  $j$  iterations, then the relative error  $\epsilon_k = \left| \bar{\alpha}_k^{[j]} / \bar{\alpha}_k^{[200]} - 1 \right|$  for the  $k$ th parameter. We see that  $\bar{\alpha}_k$  converges the fastest for the intercept, a bit slower for regressors connected to the HRF and its temporal derivative and for the AR coefficients, and the slowest for the head motion nuisance regressors.

To understand why the parameters differ in convergence speed one has to consider how the VB algorithm works, updating the approximate posterior for  $\mathbf{W}$  given  $\alpha$  and vice versa. Since  $\alpha$  controls the smoothness and shrinkage of  $\mathbf{W}$ , there will be much dependence between the estimated posteriors of these two parameters, leading to slower convergence. In general, the more informative the data are the faster is seemingly the convergence. For example, for  $k$  corresponding to the intercept in every voxel,  $\mathbf{W}_{k,\cdot}$  is quite well defined from the data, so it will depend relatively little on  $\alpha_k$ . Given the smoothness of  $\mathbf{W}_{k,\cdot}$ ,  $\alpha_k$  is also well determined, so convergence will be quick. On the other hand, for  $k$  corresponding to

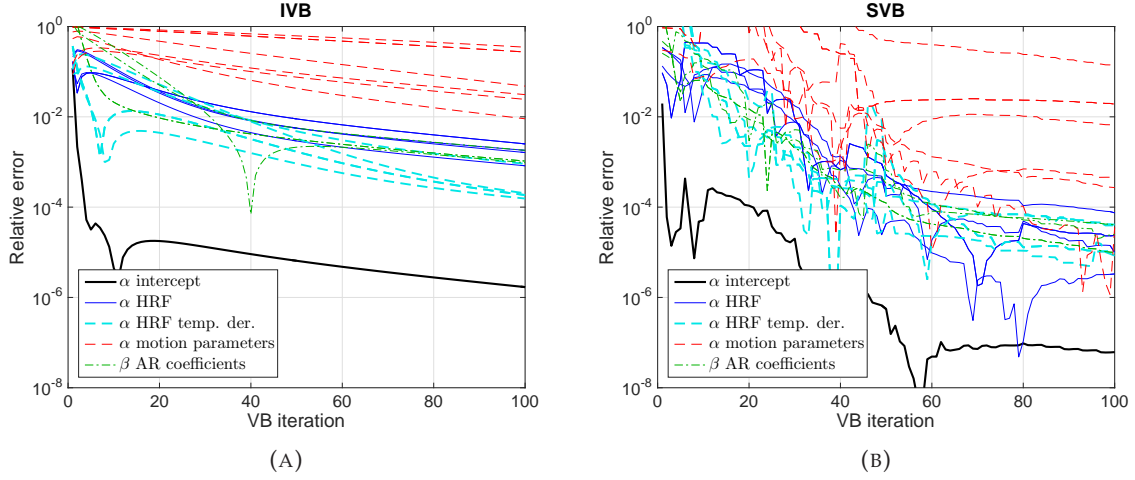


FIGURE C.1. Spatial hyperparameter relative errors by VB iteration number, relative to the value after 200 iterations, for the face repetition data. The errors for IVB are shown in (a), while the errors for SVB are shown in (b). The regression coefficient hyperparameters  $\alpha$  are divided by regressor type and  $\beta$  are the hyperparameters for the AR coefficients.

head motion,  $\mathbf{W}_{k_r}$ , will be non-significant in most voxels, giving more dependence with  $\alpha$  and slower convergence.

It takes almost 50 iterations for all the HRF regressor  $\alpha_k$ :s (which are the most interesting ones for PPMs) to reach a relative error smaller than 1% for the IVB algorithm. This is interesting, since the SPM12 default setting is to use 4 VB iterations and the other available stopping criterion, based on the model evidence lower bound, results in 8 iterations. Even though this might seem as too few, the effect of this on the PPMs is not necessarily so large (see Figure 4.3 and 4.5), so keeping the number of iterations low might well be a reasonable strategy in order to reduce the processing time.

The SVB convergence plot would look similar to that of IVB, if implemented directly as it is presented above. However, we added some ad hoc steps to the SVB algorithm to speed up convergence. In contrast to IVB, the posterior update step for  $q(\mathbf{W})$  only depends on the other parameters (in particular on  $\alpha$ ) and the data, but not on the previous iteration value of  $q(\mathbf{W})$ . Therefore, if there is some faster way to approach the optimal value of  $\bar{\alpha}$ , than the SVB update equation for  $q(\alpha)$ , that speedup will transfer to  $q(\mathbf{W})$  as well. In the SVB implementation we use three ad hoc tricks to speed up convergence, which we have noticed work well in practice:

- (1) In every other VB iteration, instead of accepting the  $\tilde{q}_{1k}^{SVB[j]}$  (and hence  $\bar{\alpha}_k^{[j]} = \tilde{q}_{1k}^{SVB[j]}$  .  $\tilde{q}_2^{SVB}$ ) from the VB update equation, we fit a quadratic function to the values of the last three iterations,  $\bar{\alpha}_k^{[j-2]}$ ,  $\bar{\alpha}_k^{[j-1]}$  and  $\bar{\alpha}_k^{[j]}$ , as a function of the iteration number  $j$ . If the vertex of the quadratic function occurs after  $j$ , we take the value of the function at the vertex as a prediction of what  $\bar{\alpha}_k$  will converge to. Otherwise, we assume that we are far from the final value and set the prediction to  $\bar{\alpha}_k^{[j-1]} + 20(\bar{\alpha}_k^{[j]} - \bar{\alpha}_k^{[j-1]})$ . We

- set  $\bar{\alpha}_k^{[j]}$  to the predicted value but only allow changes in alpha up to a factor 5 from the value proposed by the VB update equation. This methodology will make the SVB algorithm behave less stable, but allow for larger steps when  $\bar{\alpha}_k$  is far from the final value which leads to faster convergence, especially when data are not so informative.
- (2) We use the prior mean as starting values for all parameters, which is a little different from what SPM's IVB does, but leads to faster convergence for the SVB algorithm.
  - (3) In the first 10 iterations of the SVB algorithm we use  $N_s = 5$ , which will make these iterations quicker, but not as exact.

The effect of using these tricks can be seen in Figure C.1b, showing faster convergence for the hyperparameters using SVB as compared to IVB.

Different random seeds in the MC approximation (equation (3.13)) make the SVB algorithm converge to slightly different posteriors. We reran SVB several times with  $N_s = 100$  and different seeds, and found very small differences in the results for the simulated data, but for some real data sets we found that the differences could be slightly larger. SVB appears to sometimes get stuck in local modes of the model evidence lower bound that is implicitly optimized by the VB algorithm. The resulting SVB posteriors for all of these modes were however always closer to the exact MCMC posterior than the IVB posterior was. Since the model evidence lower bound is computationally intractable, we instead compared the different modes using the joint posterior ratio in equation (A.13) evaluated in the SVB posterior mean to decide which results to present for these problematic data sets. The same multi-modal issues were not observed when rerunning MCMC with different seeds or IVB with different starting values.

The MCMC Gibbs algorithm suffers from the same problem with the dependence between  $\mathbf{W}$  and  $\alpha$ , slowing down convergence. Inspecting trace plots and the estimated inefficiency factor  $IF = 1 + 2 \sum_{j=1}^{\infty} \rho_j$ , where  $\rho_j$  is the autocorrelation function of the MCMC chain, for different data sets shows that the convergence for the parameters  $\mathbf{W}$ ,  $\mathbf{A}$  and  $\lambda$  is generally excellent. For example, the maximum IF was less than 1.5 across all regression coefficients  $\mathbf{W}$  for both the real data sets when using the 3D prior, except for the head motion nuisance regressors. With 20000 post-burnin iterations and thinning factor 5, this means that we have at least  $4000/1.5 \approx 2700$  effective samples to base the PPMs on, and Monte Carlo standard deviations less than  $\sqrt{\frac{0.9 \cdot 0.1}{2700}} \approx 0.0058$  for posterior probabilities larger than 90%. We also looked at posterior mean maps and PPMs for the main contrast for the different data sets and compared to the same maps computed based only on the first 5000, 10000 and 15000 samples respectively, and saw small differences in general, much smaller than for example when comparing to the VB maps. The spatial hyperparameters  $\alpha$  and  $\beta$  can however mix poorly, especially when the data are non-informative. It would be tempting to improve the mixing using a collapsed Gibbs sampling step for  $\alpha$  (and  $\beta$ ), but this would require the computation of the precision matrix determinant  $|\tilde{\mathbf{B}}|$ , which would be too time consuming in general. For the main parameters of interest, the  $\alpha_k$ 's belonging to the HRF regressors, and also for the hyperparameters belonging to the intercept and to the first AR coefficients, the



convergence rates are acceptable in general. A 1000 iteration burnin is usually sufficient to reach stationarity for these hyperparameters.

In the timing comparisons in the results section, all methods were run for a long time (200 iterations for IVB, 50 iterations for SVB and 10000 (simulated data) / 20000 (real data) iterations with thinning factor 5 after 1000 burnin samples for MCMC (an exception was MCMC in 3D for the face repetition data, which required 3500 burnin samples)) and we compute the time until the estimated posterior mean of  $\alpha_k$  reaches within 1% of its final values for respective algorithm. For VB, this is the same as based on the relative error defined above, and for MCMC this is based on the relative error of the cumulative mean of MCMC samples. For the simulated data, this is based on all  $\alpha_k$ , but for the real data we only consider the  $\alpha_k$  corresponding to the intercept and HRF regressors. This is because the convergence is sometimes extremely slow for  $\alpha_k$ :s corresponding to the head motion and HRF temporal derivative regressors, while these have negligible effects on the results. For the real data, we use  $\delta = 10^{-8}$  and  $N_s = 100$  throughout.

The simulated data were analyzed on a computing cluster, using two 8-core (16 threads) Intel Xeon E5-2660 processors at 2.2GHz. The real data were mainly analyzed on the same cluster, but some demanding runs were carried out on a faster workstation with a 4-core (8 threads) Intel Xeon E5-1620 processor at 3.5GHz. The workstation ran 46% faster for a large SVB estimation test and therefore the timings using this computer were multiplied with a factor 1/0.54 in this report, hence the word “approximate” in Table 1. The operating system was Linux in both cases.

#### APPENDIX D. SIMULATED AND REAL DATA DETAILS

The synthetic data are simulated from the model in Section 2, with  $K = 5$  and  $P = 1$  and parameter values that are similar to those of the pre-processed face repetition data. The design matrix is set to have the first 4 columns equal to the standard canonical HRF regressors from the paradigm in the face repetition data (so  $T = 351$ ) and the fifth column corresponds to the intercept. The intercept of each voxel is sampled i.i.d. from a  $\mathcal{N}(900, 130^2)$ -distribution.  $\lambda_n^{-1} = 100$  for each voxel and the other hyperparameters are set as  $\alpha_1 = 10^{-4}$ ,  $\alpha_2 = 5 \cdot 10^{-4}$ ,  $\alpha_3 = 2 \cdot 10^{-3}$ ,  $\alpha_4 = 10^{-2}$ ,  $\beta_1 = 10$ , which generate reasonable values of for the  $\mathbf{W}$  and  $\mathbf{A}$ , with varying levels of informativeness.  $\mathbf{W}_{k\cdot}$  (and similarly  $\mathbf{A}$ ) are sampled independently for  $k \in \{1, 2, 3, 4\}$  from the UGL prior using PCG with  $\mathbf{B}_{data} = \mathbf{0}$  and  $\mathbf{w}_r^{start} = \mathbf{0}$ . Even though the UGL prior is improper this works well and generates samples with mean close to zero. Conditioned on the parameters the simulation of fMRI-data  $\mathbf{Y}$  is straightforward using the model. The data are simulated in a big rectangular block of size  $53 \times 63 \times 46$  and we follow the SPM default to scale down the value in each voxel to get mean 100. We then use centered masks of size  $10 \times 10 \times 10$ ,  $25 \times 20 \times 20$  and  $50 \times 50 \times 40$  to obtain data of size  $N = 10^3$ ,  $10^4$  and  $10^5$  that we test the methods on.

The real data from the face repetition experiment (Henson et al., 2002), previously used in Penny et al. (2005b) and available at SPM’s homepage ([http://www.fil.ion.ucl.ac.uk/spm/data/face\\_rep/](http://www.fil.ion.ucl.ac.uk/spm/data/face_rep/)), was pre-processed using the same steps as in Penny et al. (2005b) using SPM12. After masking away voxels outside the brain, this

results in a data set with  $N = 57535$  voxels of size  $3 \times 3 \times 3$  mm and  $T = 351$  volumes. Using the same design matrix based on the canonical HRF and its temporal derivative, 6 head motion parameters and an intercept, we have  $K = 15$ . The contrast considered is the main effect of faces which is the mean across the four regressors corresponding to the HRF of each condition.  $P = 3$  AR parameters are used for all estimations on real data. The presented PPMs for this data set show axial slice 12, which is approximately the same region shown in Penny et al. (2005b).

The real data from the visual object recognition experiment (Haxby et al., 2001; Hanson et al., 2004; O'Toole et al., 2005) was obtained from the OpenfMRI database (<http://openfmri.org/>). Its accession number is ds000105 and we consider only subject001, run001. The data was pre-processed using the same SPM pipeline as for the face repetition data, except no slice time correction was performed (since the slice order information is missing) neither was it normalized to a standard brain. This data set has  $N = 31241$  voxels of size  $3.125 \times 3.125 \times 4$  mm and  $T = 121$ . There are 8 conditions and using the same structure for the design matrix as for the face experiment data, we have  $K = 23$ . The contrast considered is the difference between seeing houses and faces ( $c_1 = 0.5$  and  $c_{15} = -0.5$ ). The axial slice presented in the results is number 30.