

Multi-Level Linear Modelling for FMRI Group Analysis Using Bayesian Inference

FMRIB Technical Report TR03MW1

(A related paper has been accepted for publication in Neuroimage)

Mark W. Woolrich *, **Timothy E.J. Behrens ***, **Christian F. Beckmann**, **Mark Jenkinson** and **Stephen M. Smith**

Oxford Centre for Functional Magnetic Resonance Imaging of the Brain (FMRIB),
Department of Clinical Neurology, University of Oxford, John Radcliffe Hospital,
Headley Way, Headington, Oxford, UK

* First two authors contributed equally to this work

Corresponding author is Mark Woolrich: woolrich@fmrib.ox.ac.uk

Abstract

Functional magnetic resonance imaging studies often involve the acquisition of data from multiple sessions and/or multiple subjects. A hierarchical approach can be taken to modelling such data with a General Linear Model at each level of the hierarchy introducing different random effects variance components. Inferring on these models is non-trivial with frequentist solutions being unavailable. A solution is to use a Bayesian framework. One important ingredient in this is the choice of prior on the variance components and top-level regression parameters. Due to the typically small numbers of sessions or subjects in neuro-imaging the choice of prior is critical. To alleviate this problem we introduce to neuro-image modelling the approach of reference priors, which drives the choice of prior such that it is non-informative in an information-theoretic sense. We propose two inference techniques at the top-level for multi-level hierarchies (a fast approach and a slower more accurate approach). We also demonstrate that we can infer on the top-level of multi-level hierarchies by inferring on the levels of the hierarchy separately and passing summary statistics of a non-central multivariate t-distribution between them.

1 Introduction

Functional magnetic resonance imaging studies are typically used to address questions about activation effects in populations of subjects. This generally involves a multi-subject and/or multi-session approach where data are analysed in such a way as to allow for hypothesis tests at the group level (15; 28), e.g. in order to assess whether the observed effects are common and stable across or between groups of interest.

Calculating the level and probability of brain activation for a single subject is typically achieved using a linear model of the signal together with a Gaussian noise model for the residuals. This model is commonly referred to as the General Linear Model (GLM) and much attention to date has been focussed on ways of modelling and fitting the (time-series) signal and residual noise at the individual single-session level (4; 25; 27).

In order to be able to generate results that extend to the population, we also need to account for the fact that the individual subjects themselves are sampled from the population and thus are random quantities with associated variances. It is exactly this step that marks the transition from a simple *fixed-effects model* to a *mixed-effects model*¹ and it is imperative to formulate a model at the group-level that allows for the explicit modelling of these additional variance terms (15; 7).

¹Note that in the FMRI literature this has often been referred to as a *random-effects model*. Within this paper, however, the separate fixed-effects and random-effects contributions to the mixed-effects variance are considered, thus making a clear distinction between “random-” and “mixed-effects” important.

We can formulate the problem of group statistics in neuro-imaging as being hierarchical (11; 1). For example, the different levels of the hierarchy could be separate GLMs for a session-level, subject-level and group-level. In this paper we attempt to deal with inference on these multi-level GLM hierarchies by utilising a fully Bayesian framework. Typically, the most important inference is at the top-level of the hierarchy, for example we may be looking for significance of a group mean. Whether we are looking to infer at the top-level with the within-session fMRI time-series data (11) or with summary statistic results from the level below (15; 28), a fully Bayesian approach provides us with the means to assess the full uncertainty in the parameter of interest (contrasts of regression parameters) at the top-level; taking into account all of the unknown variance components (fixed and random) in the model.

Bayesian statistics provides the only generic tool for inferring model parameter probability distribution functions from data. It provides strict rules for the rational and consistent adjustment of belief (in the form of probability density functions) in the presence of new information (5), which are not available in the frequentist literature. The major consequences of this are twofold. First, we may make inference about the absolute value of the parameters of interest. i.e. we may ask questions of our parameters such as, “What is the probability that θ lies in the interval $[\theta_0, \theta_1]$?”, a question unavailable to any frequentist technique. Frequentist statistics is typically limited to posing questions of the data under the “Null hypothesis” that the parameter value is zero. Inference in a frequentist framework is then limited to the simple acceptance or rejection of this null hypothesis without being able to make any statement about the parameter values. Second, Bayesian statistics gives us a tool for inferring on *any* model we choose, and *guarantees* that uncertainty will be handled correctly. Only in certain special cases (not including the model presented here) is it possible to derive analytical forms for the null distributions required by frequentist statistics. In their absence, frequentist solutions rely on null distributions derived from the data (e.g. permutation tests), losing the statistical power gained from educated assumptions about, for example, the distribution of the noise.

These features of Bayesian analysis mean that we may make inference on physiological parameters of the haemodynamic response in the complex non-linear balloon model (9) or on spatial noise relationships in multivariate spatial auto-regressive models of fMRI data (24) or, in this paper, on higher level statistics in the presence of multiple variance components.

One important ingredient in a Bayesian approach is the choice of prior on the variance components and top-level regression parameters. Due to the typically small numbers of observations in neuro-imaging above the first-level (e.g. small numbers of subjects), this choice of prior is critical. To solve this problem we introduce to neuro-image modelling the approach of reference priors, which drives the choice of prior such that it is non-informative in an information-theoretic sense. For GLMs where a frequentist solution is available, reference analysis gives the same inference as a frequentist approach. Importantly, reference analysis allows us to perform inference when frequentist solutions are unavailable.

Using fully Bayesian reference analysis we propose two approaches to inferring at the top-level; these are a fast approximation to the marginal posterior, and a slower approach utilising Markov Chain Monte Carlo (MCMC) followed by a multivariate non-central t-distribution fit to the MCMC chains.

In (11) the hierarchical model is solved “all-in-one” using the within-session fMRI time-series data as input. However, in neuro-imaging, where the human and computational costs involved in data analysis are relatively high it is desirable to be able to make top level inferences using the *results* of separate lower-level analyses without the need to re-analyse any of the lower-level data; an approach commonly referred to as the summary statistics approach to fMRI analysis (15). Within such a summary statistic split-level approach, group parameters of interest can easily be refined as more data become available.

In (15), when inferring at the top level, this summary statistic split-level approach is shown to be equivalent to inferring all-in-one under certain conditions (e.g. the approach in (15) requires balanced designs). (1) show that top-level inference using the split-level summary statistics approach can be made equivalent to the all-in-one approach with no restrictions, if we pass up the correct summary statistics (in particular, the covariances from previous levels). Furthermore, (1) demonstrate that by taking into account lower-level covariances and heterogeneity a substantial increase in higher-level z-statistic is possible. However, (1) only show that this is the case when all variance components are *known*. Independently, in this paper, using the fully Bayesian approach, we show this equivalence for when the variance components (excluding autocorrelation) are *unknown*. The equivalence relies on the assumption that the summary statistics, which correspond to the marginal distributions of the GLM regression parameters, can be represented as a multivariate non-central t-distributions. Between the first-level (within session) and the second-level this can be shown analytically. For summary statistics at higher-levels this is an assumption which we test empirically using artificial data.

In summary, there are three main contributions presented in this paper. Firstly, we introduce reference analysis to neuro-imaging. Secondly, we propose two inference techniques at the top-level for multi-level hierarchies (a fast approach and a slower more accurate approach). Thirdly, we demonstrate that we can infer on the top-level of multi-level hierarchies by inferring on the split levels separately and passing summary statistics between them.

1.1 Paper Overview

We start in section 2 by considering the traditional two-level model. In section 3, using the reference analysis fully Bayesian framework, we show how inference on the two-level model can be split into separate inference on the two levels with the summary statistics of a multivariate non-central t-distribution being passed between the two-levels of inference. We then propose two approaches to inferring at the top level. In section 4 we discuss how we can extend the split model inference approach to higher level models than the two-level model. In section 5 we also discuss how we can deal with multiple group variances under certain conditions. In section 6 we validate the crucial assumption of the marginal distribution of the GLM regressions parameters being a multivariate non-central t-distribution at levels higher than the first using artificial data. Finally, in section 7 we go on to show results on FMRI data.

2 Model

To begin with we consider the familiar two-level univariate GLM for FMRI. For example, the model that in the first level deals with individual sessions for individual subjects, relating time-series to activation, and in the second level deals with a group of subjects or sessions (or both), relating the combined individual activation estimates to some group parameter, such as mean activation level. Note that all models and inference in this paper are mass-univariate, i.e. each voxel is modelled and processed independently of the others in the data.

2.1 Two-level GLM

Consider an experiment where there are N_K first-level sessions and that for each first-level session, k , the preprocessed FMRI data is a $T \times 1$ vector Y_k , the $T \times P_K$ design matrix is X_k , and β_k is a $P_K \times 1$ vector of parameter estimates ($k = 1, \dots, N_K$). The preprocessed FMRI data, Y_k , is assumed to have been prewhitened (4; 25). An individual GLM relates first-level parameters to the N_k individual data sets:

$$Y_k = X_k \beta_k + \epsilon_k, \tag{1}$$

where $\epsilon_k \sim N(0, \sigma_k^2 I)$. In this paper we consider the variance components as unknown with the exception of the first-level FMRI time-series autocorrelation. The residuals ϵ_k are assumed to be prewhitened data and as a result are uncorrelated. This inherently means that we assume that the autocorrelation is known with no uncertainty, an assumption which is commonly made in FMRI time-series analysis (4; 25; 8). Note that the first level design matrices, X_k , do not need to be the same for all k .

Using the block diagonal forms, i.e. with

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_{N_K} \end{bmatrix}, \quad X = \begin{bmatrix} X_1 & 0 & \cdots & 0 \\ 0 & X_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & X_{N_K} \end{bmatrix}, \quad \beta_K = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{N_K} \end{bmatrix} \quad \text{and} \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_{N_K} \end{bmatrix}$$

the two-level model is

$$Y = X \beta_K + \epsilon_K \tag{2}$$

$$\beta_K = X_g \beta_g + \epsilon_g \tag{3}$$

where X_g is the $N_K \times P_G$ second-level design matrix (e.g. separating controls from normals or modelling different sessions for subjects), β_g is the $P_G \times 1$ vector of second-level parameters, and $\epsilon_g \sim N(0, \sigma_g^2 I)$ and where $\epsilon_K \sim N(0, V_K)$ with V denoting the diagonal form of first-level covariance matrices $\sigma_k^2 I$. We call σ_g^2 the random effects variance.

3 Inference

There are no solutions in the frequentist literature to this model when the variance components are unknown. Furthermore, inference is highly sensitive to any assumptions made, due to the low number of observations typically available at the subject level in fMRI.

(11) have proposed an approximate Bayesian solution for the model all-in-one, by assuming that the posterior over the regression parameters is multivariate Normal. However, this does not fully incorporate the full uncertainty of the variance components into the parameters of interest (the regression parameters) at the top level. Indeed, the marginal posterior over the regression parameters turns out to be multivariate t-distributed.

In this section we start by introducing the Bayesian inference framework. However, when using a Bayesian framework, we also need to choose priors for the parameters in our model. In particular, we need to choose priors on the top-level regression and variance parameters. Hence, in the next part of this section we describe how we can use reference priors as noninformative priors.

We could proceed to infer on the full model all-in-one. Instead, by using the fully Bayesian approach with reference priors, we go on to show how we can use summary statistics (from inferring on the first-level model in isolation) as the input into a second level. We show that this gives the same inference as we would obtain from using the full model all-in-one.

3.1 Bayesian Inference

The two rules at the heart of Bayesian learning techniques are conceptually very simple. The first tells us how (for a model \mathcal{M}) we should use the data, \mathbf{Y} , to update our *prior* belief in the values of the parameters Θ , $p(\Theta|\mathcal{M})$ to a *posterior* distribution of the parameter values $p(\Theta|\mathbf{Y}, \mathcal{M})$. This is known as Bayes' rule:

$$p(\Theta|\mathbf{Y}, \mathcal{M}) = \frac{p(\mathbf{Y}|\Theta, \mathcal{M})p(\Theta|\mathcal{M})}{p(\mathbf{Y}|\mathcal{M})} \quad (4)$$

Unfortunately, calculating this *posterior pdf* is seldom straightforward. The denominator in equation 4 is:

$$p(\mathbf{Y}|\mathcal{M}) = \int_{\Theta} p(\mathbf{Y}|\Theta, \mathcal{M})p(\Theta|\mathcal{M})d\Theta \quad (5)$$

an integral which is often not tractable analytically. To make matters worse, this *joint posterior pdf* on all parameters is often not the distribution which we are most interested in. We are often interested in the *posterior pdf* on a single parameter, or an interesting subset of parameters. Obtaining these *marginal* distributions again involves performing large integrals,

$$p(\Theta_I|\mathbf{Y}, \mathcal{M}) = \int_{\Theta_{-I}} p(\Theta|\mathbf{Y}, \mathcal{M})d\Theta_{-I} \quad (6)$$

where Θ_I are the parameters of interest and Θ_{-I} are all other parameters. Again these integrals are seldom tractable analytically.

One solution is to use approximations to the marginal distributions. This is the approach we take in section 3.5. Another solution is to draw samples in parameter space from the joint posterior distribution, implicitly performing the integrals numerically. For example, we may repetitively choose random sets of parameter values and choose to accept or reject these *samples* according to a criterion based on the value of the numerator in equation 4. It can be shown (e.g (13)) that a correct choice of this criterion will result in the *accepted* samples being distributed according to the joint posterior pdf (equation 4). Schemes such as this are *rejection sampling* and *importance sampling* which generate independent samples from the posterior. Any marginal distributions may then be generated by examining the samples from only the parameters of interest. However, these kinds of sampling schemes tend to be very slow, particularly in high dimensional parameter spaces, as samples are proposed at random, and thus each has a very small chance of being accepted.

Markov Chain Monte Carlo (MCMC) (see (13) and (12) for texts on MCMC) is a sampling technique which addresses this problem by proposing samples preferentially in areas of high probability. Samples drawn from the posterior are no longer independent of one another, but the high probability of accepting samples, allows for many samples to be drawn and, in many cases, for the posterior pdf to be built in a *relatively* short period of time. This is the approach we take in section 3.6.

3.2 Priors and Reference Analysis

In the fully Bayesian framework the choice of prior is critical to the inference we perform. In group statistics for fMRI the number of observations we have is typically so small as to make the influence of the priors significant. As we have no prior information we want the priors we use to be in some sense “non-informative”, i.e. we want to “let the data speak for itself”. Reference priors are priors which attempt to reflect such prior ignorance. For an overview see (17; 2).

An intuitive approach would be to choose the prior of θ to be $\pi(\theta) = 1$. However, the resulting posteriors can change significantly depending on the parameterisation used. This is because a constant prior for one parameter will not typically transform into a constant prior for another. To overcome this reparameterisation problem the Jeffreys prior was introduced for one-dimensional problems (17):

$$\pi(\theta) \propto \det(H(\theta))^{1/2} \quad (7)$$

where $H(\theta)$ is the Fisher information. However, this has difficulties dealing with multi-dimensional problems, i.e. $\Theta = (\theta_1 \dots \theta_m)$. The Berger-Bernardo method (2) of reference analysis overcomes this by determining reference priors using information-theoretical ideas which maximise the amount of expected “information” from the data. See appendix 10.5 for the derivation of the reference priors used in this paper.

The use of reference priors can be justified by consideration of the information theory that underpins them (2). However, whilst in this paper the null hypothesis frequentist inference is generally unknown, it is interesting to note that the Berger-Bernardo reference priors give the same inference as frequentist null hypothesis testing for cases of GLM inference for which the frequentist null hypothesis test *is* known. For example, in frequentist inference on the GLM we typically examine the probability of attaining statistics for linear combinations (contrasts) of regression parameters under the null hypothesis. In cases where such null hypothesis frequentist inference for the GLM is known, the Berger-Bernardo reference priors give the same probabilities when we test the probability that a contrast is greater than zero.

3.3 First-level

Here we consider the first-level in isolation and derive the marginal posterior distribution for β_k , the vector of GLM height parameters for the *first-level* model fit. Equation 1 gives us the likelihood for a first-level model in isolation, $p(Y_k | \beta_k, \sigma_k^2)$. The joint posterior on all parameters in this model is then:

$$p(\beta_k, \sigma_k^2 | Y_k) \propto p(Y_k | \beta_k, \sigma_k^2) p(\beta_k, \sigma_k^2) \quad (8)$$

where $p(\beta_k, \sigma_k^2)$ is the prior distribution on the regression and variance parameters. We use the Berger-Bernardo reference prior (see section 3.2), which is:

$$p(\beta_k, \sigma_k^2) = 1/\sigma_k^2. \quad (9)$$

However, equation 8 does not give the distribution of interest for inference. We would like to infer on the posterior distribution on the activation height parameters β_k when the effect of estimating σ_k^2 is accounted for, i.e. we would like to infer on $p(\beta_k | Y)$. To get this distribution, we must marginalise the joint posterior (equation 8) over the parameter of no interest σ_k^2 . This integral gives a multivariate non-central t-distribution for the posterior distribution on β_k (18):

$$\begin{aligned} p(\beta_k | Y_k) &\propto \int p(Y_k | \beta_k, \sigma_k^2) / \sigma_k^2 d\sigma_k^2 \\ &= \mathcal{T}(\beta_k; \mu_{\beta_k}, \sigma_{\beta_k}^2, \Sigma_{\beta_k}, \nu_{\beta_k}). \end{aligned} \quad (10)$$

where

$$\begin{aligned} \mu_{\beta_k} &= (X_k^T X_k)^{-1} X_k^T Y_k \\ \sigma_{\beta_k}^2 &= (Y_k - X_k \mu_{\beta_k})^T (Y_k - X_k \mu_{\beta_k}) / (T - P_K) \\ \Sigma_{\beta_k} &= (X_k^T X_k)^{-1} \\ \nu_{\beta_k} &= T - P_K. \end{aligned} \quad (11)$$

Note that if inference is performed in the frequentist framework, the null distribution on β_k is the multivariate *central* t-distribution with the *exact* same covariance structure $\sigma_{\beta_k}^2 \Sigma_{\beta_k}$ and degrees of freedom, ν_{β_k} , and the maximum likelihood estimate for β_k is *exactly* μ_{β_k} , the mean of the posterior distribution in the Bayesian framework.

3.4 Two-level

Here we consider the full two-level model laid out in equations 2 and 3, applying the same ideas as in the previous section to infer on the *second-level* GLM height parameters β_g . We will substitute into the posterior for the full two-level model the summary result of the first-level model derived in the previous section. This will provide us with the way of inferring on the full two-level model using just the summary result of the first-level, i.e. without re-using the data Y .

Considering equations 2 and 3. The full joint posterior for the two-level model is:

$$p(\beta_g, \sigma_g^2, \beta_{\mathbf{K}}, \sigma_{\mathbf{K}}^2 | Y) \propto \prod_k \{p(Y_k | \beta_k, \sigma_k^2)\} p(\beta_{\mathbf{K}} | \beta_g, \sigma_g^2) p(\beta_g, \sigma_g^2, \sigma_{\mathbf{K}}^2), \quad (12)$$

where $\sigma_{\mathbf{K}}^2$ is the $(K \times 1)$ vector of first level variances σ_k^2 , and $\beta_{\mathbf{K}}$ is the $(K \times 1)$ vector of first level regression parameters β_k (for $k = 1 \dots K$). We set the prior to be the Berger-Bernardo reference prior for this full two-level model (see section 3.2):

$$p(\beta_g, \sigma_g^2, \sigma_{\mathbf{K}}^2) = \frac{1}{\sigma_g^2} \prod_k \frac{1}{\sigma_k^2}. \quad (13)$$

Note that this model specification gives the posterior distribution, not only on the second-level parameters (β_g, σ_g^2) but also on the parameters from **all** of the first-level models $(\beta_{\mathbf{K}}, \sigma_{\mathbf{K}}^2)$. However, if we are only interested in the top/second-level parameters, we may substitute the summary result from the first level into this two-level model and marginalise over $(\beta_{\mathbf{K}}, \sigma_{\mathbf{K}}^2)$ (see appendix 10.6), showing that the *marginal* distribution on (β_g, σ_g^2) does not depend on the original data, but on the summary parameters from the first level, i.e. μ_{β_k} and $\sigma_{\beta_k}^2 \Sigma_{\beta_k}$:

$$p(\beta_g, \sigma_g^2, \tau_{\mathbf{K}} | Y) \propto \prod_k \{ \mathcal{N}(\mu_{\beta_k}; X_{gk} \beta_g, (\sigma_{\beta_k}^2 \Sigma_{\beta_k} / \tau_k) + \sigma_g^2 I) \Gamma(\tau_k; \nu_{\beta_k} / 2, \nu_{\beta_k} / 2) \} 1 / \sigma_g^2 \quad (14)$$

where X_{gk} is the k^{th} row vector of the second-level design matrix X_g , and $\tau_{\mathbf{K}}$ is a $(K \times 1)$ vector of latent variables τ_k for $k = 1 \dots K$ introduced for mathematical convenience (see appendix 10.6).

A special case of equation 14 is when the variances, $\sigma_{\beta_k}^2 \Sigma_{\beta_k}$ on the first-level GLM parameters are known with very high degrees of freedom ($\nu_k \rightarrow \infty$). This is equivalent to $p(\beta_k | Y_k)$ in equation 10 being a Normal distribution instead of a t-distribution. In this case, the prior distribution on $\tau_{\mathbf{K}}$ reduces to a delta function centered on $\tau_{\mathbf{K}} = \mathbf{1}$ and the joint posterior distribution on the second-level parameters reduces to:

$$p(\beta_g, \sigma_g^2 | Y) \propto \prod_k \{ \mathcal{N}(\mu_{\beta_k}; X_{gk} \beta_g, (\sigma_{\beta_k}^2 \Sigma_{\beta_k}) + \sigma_g^2 I) \} 1 / \sigma_g^2. \quad (15)$$

Equation 14 (or, in the special case, equation 15) gives us the joint posterior distributions of β_g, σ_g^2 and $\tau_{\mathbf{K}}$. However, as in the first-level model, we are actually interested in inferring upon the marginal distribution over the GLM height parameters, β_g . This marginal posterior $p(\beta_g | Y)$ cannot be obtained analytically. Therefore, we consider two approaches, a fast posterior approximation and a slower but more accurate approach using Markov Chain Monte Carlo (MCMC) sampling. Crucially, in both approaches we are going to assume that $p(\beta_g | Y)$ is a multivariate non-central t-distribution:

$$p(\beta_g | Y) \propto \int p(\beta_g, \sigma_g^2, \tau_{\mathbf{K}} | Y) d\sigma_g^2 d\tau_{\mathbf{K}} \quad (16)$$

$$\approx \mathcal{T}(\beta_g; \mu_{\beta_g}, \sigma_{\beta_g}^2 \Sigma_{\beta_g}, \nu_{\beta_g}) \quad (17)$$

This assumption is crucial to the idea of being able to split hierarchies into inference on different levels for higher order models as we shall see in section 4. We shall test the validity of this assumption later. The fast posterior approximation or MCMC approaches are the means by which we get the distribution parameters $\mu_{\beta_g}, \sigma_{\beta_g}^2 \Sigma_{\beta_g}, \nu_{\beta_g}$.

3.5 Fast Posterior Approximation

Here we propose a fast but approximate approach for estimating the distribution parameters, $\mu_{\beta_g}, \sigma_{\beta_g}^2 \Sigma_{\beta_g}, \nu_{\beta_g}$, in equation 16. First, we assume high degrees of freedom at the first-level, i.e. $\tau_k = 1$ for all k . We then obtain a point estimate of σ_g^2 , and use this point estimate to compute a point estimate of β_g . For the details of how we obtain these point estimates $\widehat{\sigma}_g^2$ and μ_{β_g} , see appendix 10.7.

We then make the assumption that the effect of uncertainty in σ_g^2 is the same as the effect of uncertainty in σ_k^2 in a first-level model. This means that $p(\beta_g|Y)$ is a multivariate non-central t-distribution:

$$\mathcal{N}(\beta_g; \widehat{\beta}_g, (X_G^T U^{-1} X_G)^{-1}, \nu). \quad (18)$$

where U is a diagonal matrix with the k^{th} diagonal element given by $S_k = (\sigma_{\beta_k}^2 \Sigma_{\beta_k} / \tau_k) + \sigma_g^2 I$. However, we do not know the degrees of freedom (DOF), ν . We might expect the DOF to be within the range, $N_K - P_G \leq \nu \leq \infty$. In the validation section we will look at using $\nu = N_K - P_G$ (lower estimate) and $\nu = \infty$ (upper estimate). The accuracy of these assumptions are examined with simulations in section 6.

3.6 MCMC

Here, we use Markov Chain Monte Carlo (MCMC) to sample from the full joint posterior distribution given in equation 14. This also automatically provides us with samples from the marginal posterior distribution, $p(\beta_g|Y)$.

We use single-component Metropolis-Hastings jumps (i.e. we propose separate jumps for each of the parameters in turn) for all parameters. We use separate Normal proposal distributions for each parameter, with the mean fixed on the current value, and with a scale parameter σ_p for the p^{th} parameter that is updated every 30 jumps. At the j^{th} update σ_p is updated according to:

$$\sigma_p^{j+1} = \sigma_p^j \tilde{R} \frac{(1 + A + R)}{(1 + R)} \quad (19)$$

where A and R are the number of accepted and rejected jumps since the last σ_p update respectively, \tilde{R} is the desired rejection rate, which we fix at 0.5.

We require a good initialisation of the parameters in the model purely to reduce the required burn-in of the MCMC chains (the burn-in is the part of the MCMC chain which is used to ensure that the chain has converged to be sampling from the true distribution). To initialise we use the fast approximation approach described in section 3.5.

3.7 BIDE T

MCMC can be used to directly obtain samples from $p(\beta_g|Y)$. However, we would need to get lots of samples well into the tail of the distribution, and MCMC sampling is computationally intensive. Hence, we avoid the need for large numbers of samples by assuming that $p(\beta_g|Y)$ is a multivariate non-central t-distribution. Recall that assuming a multivariate non-central t-distribution is also important to the idea of being able to split hierarchies into inference on different levels. Therefore, we clean up the samples of the posterior using Bayesian Inference with Distribution Estimation using a T-fit (BIDE T).

BIDE T fits a multivariate non-central t-distribution to the MCMC samples of $p(\beta_g|Y)$ as described in appendix 10.4. Figure 1 shows the result of using the multivariate non-central t-distribution fit to an MCMC chain obtained (see section 3.6) on a voxel in Dataset 2 described in section 6.

3.8 Contrasts

Whether from the fast approximation approach or from MCMC plus BIDE T, the output from the analysis at any level in the hierarchy gives us a multivariate non-central t-distribution (equation 16). As in the frequentist framework we can ask questions about linear combinations (or contrasts) $c^T \beta_g$ of the parameters in β_g .

If c is a $P \times 1$ vector representing a t-contrast, we can use equation 16 to give us the univariate non-central t-distribution over $c^T \beta_g$

$$p(c^T \beta_g | Y) = \mathcal{T}(c^T \beta_g; c^T \mu_{\beta_g}, \sigma_{\beta_g}^2 c^T \Sigma_{\beta_g} c, \nu_{\beta_g}) \quad (20)$$

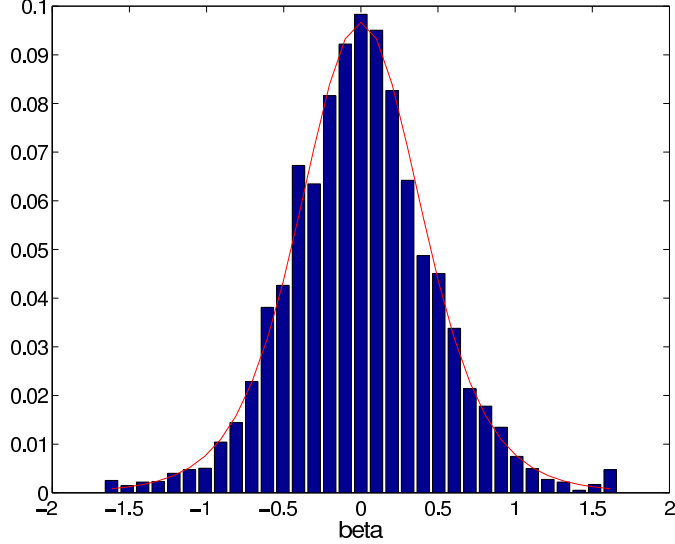


Figure 1: The (in this case 1-dimensional) t-fit obtained on the MCMC samples from a voxel in dataset 1.

We can then look at the $p(c^T \beta_g > 0 | Y)$. Note that this is equal to the probability of getting a t-value greater than the t-statistic:

$$t = c^T \mu_{\beta_g} / \sqrt{c^T \sigma_{\beta_g}^2 \Sigma_{\beta_g} c} \quad (21)$$

under a central t-distribution with degrees of freedom ν_{β_g} .

4 Higher-level Models

An increasing number of studies have three levels, in particular: a within session level, a session level, and a subject level. With multiple sessions for multiple subjects it becomes possible to model the between-session variance separately from the between-subject variance, and hence one can benefit from the improvements in sensitivity (due to heterogeneity of variance) this produces.

In section 3.4 we showed that we could infer on the full two-level model using just the summary result of the first level without using the data Y . We can use a similar argument to show that we can infer on a full three-level model using the summary result of the two-level model (given by equation 16) without using the data Y . The resulting distribution is similar to that in equation 14. Hence, we similarly assume that the marginal posterior is a multivariate non-central t-distribution equivalent to equation 16, and again we can use the fast posterior approximation or MCMC approaches to get the distribution parameters.

Higher-level models can be considered using exactly the same argument. This is because after the first-level, outputs and inputs, for subsequent levels can be summarised as a multivariate non-central t-distribution. Hence, the assumption that the marginal distribution in equation 14 is a multivariate non-central t-distribution is integral to the idea of being able to split inference on multiple-level models into inference on the different levels. We shall test the validity of this assumption later.

5 Multiple group variances

We can use the framework we have described to work with multiple group variances at any level after the first-level. An example of when this would be useful is when we might expect different between-subject variances for a patient group and a control group. We can easily deal with such multiple group variances if we limit ourselves to design matrices which are “separable” with respect to the variance groupings.

We define a sub-design matrix as the part of the design matrix belonging to a group of observations for which we want to have a separate variance group. A design matrix would be “separable” with respect to the variance groupings if the sub-design matrices could be inferred upon using separate GLMs to give the same result as inferring on one GLM using the full design matrix.

We define a “group-regressor” as that part of an regressor that belongs to a particular group variance:

Variance Group	regressor 1	regressor 2
1	1	0
1	1	0
1	1	0
2	0	1
2	0	1
2	0	1

The group-regressor for regressor 1, and for group 1 is $[1, 1, 1]^T$. The group-regressor for regressor 1, group 2 is $[0, 0, 0]^T$.

We can check if our design matrix is “separable” by checking that within each regressor only one group-regressor has non-zero values in it. An example of a design matrix which violates this is:

Variance Group	regressor 1	regressor 2
1	1	1
1	1	1
1	1	1
2	1	-1
2	1	-1
2	1	-1

Simulations have shown that if this constraint is not met then the resulting β_g vector is not generally multivariate t-distributed. Whilst MCMC could deal with it, this violation prohibits the use of BIDEt. This would require the use of longer MCMC chains and would also prohibit carrying the output to higher levels as the output from a level with these properties could not be summarised as a multivariate t-distribution. Hence, we need in practice to ensure that our designs are “separable” with respect to the variance groupings.

These “separable” multiple group variance designs can then be implemented by inferring on separate GLMs using the fast approximation or MCMC plus BIDEt. The results for different variance groups are pooled into one multivariate t-distribution for β_g . We can then proceed to the contrast stage and ask questions within- or across-variance groupings.

6 Artificial Data

6.1 Methods

In section 3.4 we showed that the two-level model can be inferred upon using the summary statistics of the first-level model inference (equation 14). This means that all-in-one and split-level inference is equivalent when we infer on the top-level regression parameters. Here we use four different null artificial datasets from the two-level model for 400 voxels to validate the fast approximation and MCMC/BIDEt inference we perform on equation 14.

6.1.1 Inference approaches

The different inference approaches are all different ways of obtaining a z-statistic for the t-contrast of interest. The different inference approaches considered are:

- *[MCMC]* We sample from $p(\beta_g | Y)$ to get an MCMC chain of 200,000 samples (with a burnin of 1000 samples) using the approach described in section 3.6, we directly calculate the $p(c^T \beta_g > 0 | Y)$ from the MCMC samples of the marginal posterior of $p(c^T \beta_g | Y)$. We can then use a p-to-z transform to calculate a z-statistic at each voxel.

- *[BIDET]* We fit a non-central t-distribution to an MCMC chain of 200,000 samples (with a burnin of 1000 samples) using the approach described in section 3.7. We can then use a t-to-p-to-z transform to calculate a z-statistic at each voxel.
- *[LOWER]* We use the lower bound from the fast approximation approach described in section 3.5 to get an approximate non-central t-distribution. The lower bound is obtained when we assume DOF, $\nu = N_g - P_G$. We can then use a t-to-p-to-z transform to calculate a z-statistic at each voxel.
- *[UPPER]* We use the upper bound from the fast approximation approach described in section 3.5 to get an approximate non-central t-distribution. The upper bound is obtained when we assume DOF, $\nu = \infty$.
- *[OLS]* This is the standard frequentist approach (described at the start of section 3) of estimating the total mixed effects variance. This ignores $\sigma_{\beta_k}^2$. Using the total mixed effects variance estimate, frequentist theory gives that the normalised OLS estimate of $c^T \beta_g$ is t-distributed with DOF, $\nu = N_g - P_G$. We can then use a t-to-p-to-z transform to calculate a z-statistic at each voxel.

6.1.2 z-statistics

We want to be able to compare the resulting inference of these different approaches. It is difficult to compare different t-statistics with different DOF. Therefore, for each of the different inference approaches we convert to the probability of the contrast being greater than zero. This provides us with a measure which we can compare directly between the different approaches. We represent this probability as a z-statistic by ensuring that the area under one tail of a standardised (zero mean and standard deviation of one) Normal distribution corresponds to that probability. In section 6.3 we will explore the possibility of using these z-statistics to mimic null-hypothesis frequentist inference.

6.1.3 Relating the [MCMC] approach to [OLS]

It is important to appreciate that there are two different ways in which the z-statistic can be changed between [OLS] and [MCMC]. The first was demonstrated in (1), in that by taking into account lower-level covariances and their heterogeneity, a substantial increase in higher-level z-statistic is possible. This is due to the fact that the heterogeneity of the lower-level covariances is effectively used to weight the summary statistic data to give more efficient estimates (resulting in reduced top-level regression parameter variance). This is analogous to the way in which prewhitening is used in first-level analyses to weight the regression parameter estimation to give more efficient estimators (25).

(1) were unable to demonstrate the second way in which the z-statistic can be changed between [OLS] and [MCMC], because they assumed that variances were known. In this paper, when we estimate the higher-level variances, they are constrained to be positive. This overcomes the well-known “negative variance” problem in OLS (19), by forcing the total variance to be greater than it would be in the OLS case. This increased variance translates into lower z-statistics in voxels which would have suffered from this problem.

In summary, we have two ways in which z-statistics can change between [OLS] and [MCMC]. Firstly, they can increase due to increased efficiency from using lower-level variance heterogeneity. Secondly, they can decrease due to the higher-level variance being constrained to be positive.

6.1.4 Datasets

To avoid unnecessary consideration of first-level design matrices, and because we are only looking to validate the inference on equation 14, we do not generate artificial first-level data Y . Instead, we directly generate second-level summary “data”, μ_{β_k} , via equation 14. To do this we specify that we want null data by setting $\beta_g = \mathbf{0}$, and then choose values for σ_{β_k} and σ_g . As a result, μ_{β_k} , ν_k and σ_{β_k} form the summary statistic data we then use in the second-level inference.

To generate artificial data we need to decide on our values for σ_{β_k} (for $k=1 \dots K$) and σ_g . Our choice is governed by the variance ratios we want between the top level and the lower levels. In section 6.1.3 we discussed two ways in which we would expect differences between [OLS] and [MCMC] inference. However, we would expect this difference in z-statistics to be less and less substantial as the top-level variance dominates over the lower-level variance. (1) demonstrated that at a 10:1 ratio of between-session/subject variance to within-session variance, the increase in

higher-level z-statistic (due to taking into account variance heterogeneity) is negligible. One of our datasets ([Dataset 4]) utilises a 10:1 variance ratio to explore if the combination of the two possible effects discussed in section 6.1.3 shows any difference in z-statistics between [OLS] and [MCMC].

However, we also consider variance ratios of the order of 1:1. The widely reported existence of the negative variance problem in FMRI (28; 19) along with the effects seen in the real FMRI data later in this paper, demonstrate that such low group to first-level variance ratios do exist in FMRI data. We need such a ratio to reproduce data which will suffer from the well reported “negative variance” problem when using traditional OLS estimation (19). Furthermore, we need to consider the case of three level hierarchies, which are popular in neuro-imaging studies (e.g. hierarchies containing within session levels, session levels and subject levels). When one uses the summary statistics from the output of the second level to infer on the third level, the variance ratio we are concerned with is between session variance to between subject variance, for which a ratio of the order of 1:1 is realistic.

The four datasets are:

- [Dataset 1] A group mean design, with $N_K = 8$ subjects and $\sigma_{\beta_k}^2 \approx 0$, $\sigma_g^2 = 1$. The second level design matrix is: $X_g = [1, 1, 1, 1, 1, 1, 1, 1]^T$ with t-contrast $c = [1]$.
- [Dataset 2] A group mean design, with $N_K = 8$ subjects and $\sigma_{\beta_k}^2 \sim \text{Uniform}(0.1, 1.9)$, $v_k = 8$ and random effects variance $\sigma_g^2 = 1$. The design matrix and t-contrast is the same as for dataset 1.
- [Dataset 3] A paired t-test design with 5 subjects under two conditions (giving $N_K = 10$) and $\sigma_{\beta_k}^2 \sim \text{Uniform}(0.1, 1.9)$, $v_k = 8$ and random effects variance $\sigma_g^2 = 0.5$. The second level design matrix is:

$$X_g = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ -1 & 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

with t-contrast $c = [1, 0, 0, 0, 0, 0]^T$.

- [Dataset 4] A group mean design, with $N_K = 8$ subjects and $\sigma_{\beta_k}^2 \approx \sim \text{Uniform}(0.1, 1.9)$, $\sigma_g^2 = 10$. The second level design matrix is: $X_g = [1, 1, 1, 1, 1, 1, 1, 1]^T$ with t-contrast $c = [1]$.

6.2 Results

Figure 2 show boxplots of the difference in z-statistics between those obtained from a long MCMC chain of 200,000 samples and those obtained from the different inference approaches considered. The intention is to consider the inference from a very long MCMC time series as a “gold standard”. To help validate this assumption the first boxplot (labelled [MCMC]) compares this “gold standard” inference with another equally long MCMC chain but with a different random seed. This allows us to assess the inaccuracies in the “gold standard” due to the finite length of the MCMC chain. In all four datasets the difference in z-statistics for this is of the order of 0.01.

The second boxplot (labelled [BIDET]) compares our “gold standard” to the inference obtained when we fit the non-central multivariate t-distribution to the long MCMC chain with a different random seed. This allows us to validate one of the strongest assumptions that we make in this paper. That is that the marginal posterior in equations 14 are a non-central multivariate t-distribution. This is crucial to the idea of being able to split hierarchies into inference on different levels. By making this distributional assumption it also allows us to infer on shorter MCMC chains, and gives us some basis for the fast approximation approach. This assumption is well supported by these [BIDET] boxplots with the difference in z-statistics being of the order of 0.01 for all four datasets.

Figure 2 also shows boxplots for the fast approximation approaches. We show boxplots for the upper bound (labelled [UPPER]) and lower bound (labelled [LOWER]). Of particular interest is how good these bounds are

at actually bounding the “gold standard” [MCMC]. Hence, a third boxplot (labelled [BOUND]) shows the how far outside the bound the “gold standard” is. This shows a z-statistic difference of up to 0.2 for dataset 2. This z-statistic difference of up to 0.2 between the fast approximation bounds and the [MCMC] “gold standard” will be used later as part of the [HYBRID] inference approach (see section 7.1.1).

The final boxplot shows the traditional inference approach of ignoring the known fixed effects variance estimating the total mixed effects variance, and using OLS to perform inference (labelled [OLS]). Because this ignores the fixed effects variance this makes this approach the “gold standard” for Dataset 1, in which $\sigma_{\beta_k}^2 \approx 0$. Indeed this is supported by the boxplot. However, for Datasets 2 and 3, $\sigma_{\beta_k}^2 > 0$ and varies over k . For these datasets OLS will give unbiased statistics, but very inefficient statistics as the $\sigma_{\beta_k}^2$ information is ignored. These boxplots illustrate the difference in z-statistics between OLS and the “gold standard” due to this inefficiency. In Dataset 4, $\sigma_{\beta_k}^2$ is sufficiently small compared to σ_g^2 so that the differences between [OLS] and [MCMC] are negligible.

Figure 3 shows the z-statistics obtained for 20 voxels from the 3 Datasets for the inference approaches of [UPPER],[LOWER],[BIDET], and [OLS]. For Dataset 1 the correspondance of [OLS], [LOWER] and [BIDET] is reiterated. For Datasets 2 and 3 the difference between [BIDET] and [OLS] is illustrated, as is the small inaccuracy of the [UPPER] and [LOWER] fast approximation approaches compared with [BIDET].

Figure 4 shows the histograms for the four different datasets of the degrees of freedom (DOF) obtained at each voxel from fitting the non-central t-distribution to an MCMC chain of 200,000 samples from the marginal posterior, $p(c^T \beta_g | Y)$, as part of [BIDET]. For Dataset 1, we know that the OLS solution is the correct one and that the DOF, $\nu = 7$. In Dataset 4 σ_{β_k} is sufficiently small compared to σ_g so that the differences between [OLS] and [MCMC] are negligible and the range of DOF match those found in Dataset 1. Figure 4 shows that [BIDET] correctly finds the DOF as being 7 for the majority of voxels in Dataset 1. However, for Datasets 2 and 3 the OLS DOF will be $\nu = 7$ and $\nu = 4$ respectively. We should not expect [BIDET] to have the same DOF values as this. Indeed the histograms show that the DOF obtained from [BIDET] varies from about these OLS DOF values to values up to about 60 or 70 DOF. Without using [BIDET] there would be no way of knowing, for a particular voxel, the required DOF.

Figure 5 shows boxplots of the difference in z-statistics between those obtained from a long MCMC chain of 200,000 samples and those obtained from using [BIDET] on MCMC chains of varying sample sizes. This illustrates the need for an MCMC chain of at least 20,000 samples to achieve accuracies of the order of 0.02 in z-statistics.

6.3 Relating Fully Bayesian Inference to Frequentist Inference

We have a number of choices for how we use the posterior distribution $p(c^T \beta_g | Y)$. We could simply use the posterior, $p(c^T \beta_g | Y)$, to build up posterior probability maps representing the probability of activation at each voxel (10). Another possibility is the use of (spatial) mixture modelling (6; 14; 23) to classify voxels as activating and non-activating. We do not attempt to explore or discuss the relative merits of these approaches in this paper. Here, we consider another possibility of the inference produced if we mimic null-hypothesis frequentist inference (i.e. controlling a False Positive Rate (FPR)) by assuming that under the null hypothesis the z-statistics, that the fully Bayesian [BIDET] approach produces, are standardised (zero mean and standard deviation of one) Normally distributed.

To examine this possibility, figure 6 shows the log probability-log probability plots for the four different datasets for [BIDET] and [OLS]. These are plots of the nominal/theoretical frequentist FPR against the probabilities obtained empirically from our four null artificial datasets. For all four datasets [OLS] does, as expected, produce a log probability plot that matches the nominal/theoretical frequentist FPR. However, this is not true for the [BIDET] approach.

Datasets 1 and 4 with small σ_{β_k} compared to σ_g gives close to the same inference using [BIDET] as when using [OLS]. Hence, we would expect the log probability that [BIDET] produces to match the nominal/theoretical frequentist FPR. Figure 6 demonstrates that this is true.

However, for Datasets 2 and 3 ($\sigma_{\beta_k}^2$ is of the same order as σ_g) [BIDET] produces different results to [OLS]. The empirical log probabilities are lower than the nominal/theoretical FPR in figure 6(b) and (c). Recall from section 6.1.3, that we have two ways in which we expect z-statistics to change between [OLS] and [MCMC]. Firstly, they can increase due to increased efficiency from using lower-level variance heterogeneity. Secondly, they can decrease due to the higher-level variance being constrained to be positive. The first of these effects will introduce no bias into the p-p plots. Hence, only the second of these effects will be visible and the p-p plots for datasets 2 and 3 in figure 6 are consistent with this.

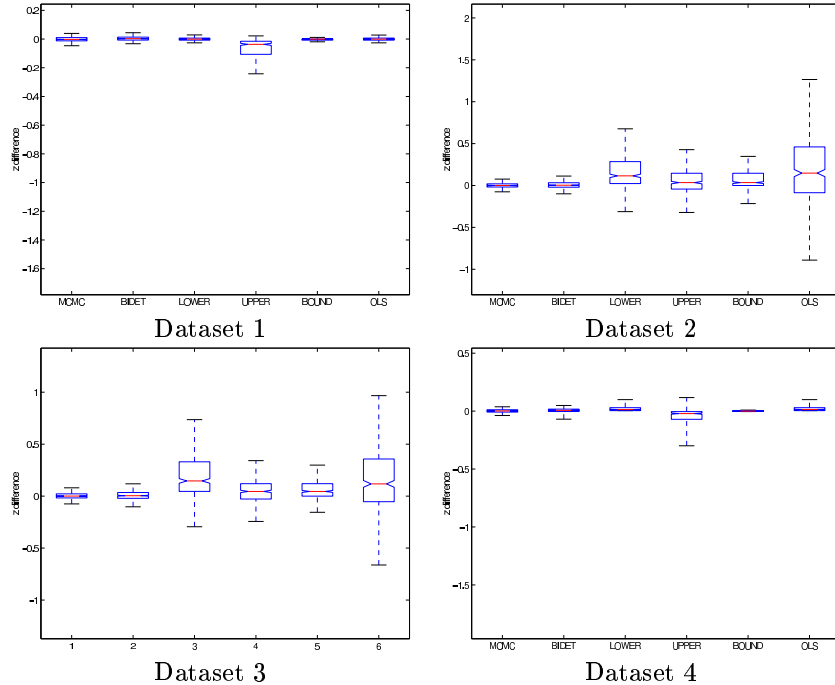


Figure 2: Boxplots over 400 voxels showing the z-statistics obtained from a long MCMC chain of 200,000 samples minus the z-statistics obtained from the different inference approaches considered. The box has lines at the lower quartile, median, and upper quartile values. The length of the whiskers is 1.5 times the Inter Quartile Range. The boxplots labelled [MCMC] correspond to the difference in z-statistics between those obtained from the 200,000 sample MCMC chain and those obtained from another 200,000 sample MCMC chain with a different random seed. The boxplots labelled [BOUND] correspond to how far outside the fast approximation bound (described as [UPPER] and [LOWER]) the z-statistics obtained from the 200,000 sample MCMC chain are.

This means that whilst we produce more accurate estimates of the total mixed effects variance, it also means that the z-statistics resulting from [BIDE] are *not* standardised Normally distributed under the null hypothesis. This is not a problem if we just report posterior probability maps or use mixture modelling.

However, if we do choose to proceed with assuming that the z-statistics from [BIDE] are standardised Normally distributed, since the empirical log probabilities are *lower* than the nominal/theoretical frequentist FPR, then the validity of our statistics will not be violated. In other words, the z-statistics from [BIDE] are, on average conservative. The disadvantage of this is that we will lose some sensitivity when compared with using the unknown, correct null distribution. The advantage is that we can utilise cluster based inference techniques on the z-statistic maps, such as Gaussian Random Field Theory (26; 21).

7 FMRI data

7.1 Methods

Here, we consider two different FMRI datasets, both of which are simple motor tasks:

- [INDEX] index finger vs. rest tapping task.
- [SEQUENTIAL] sequential finger tapping vs index finger tapping.

Each dataset consists of single sessions for eight different subjects. In both datasets the overall aim is to infer the group means at the top-level.

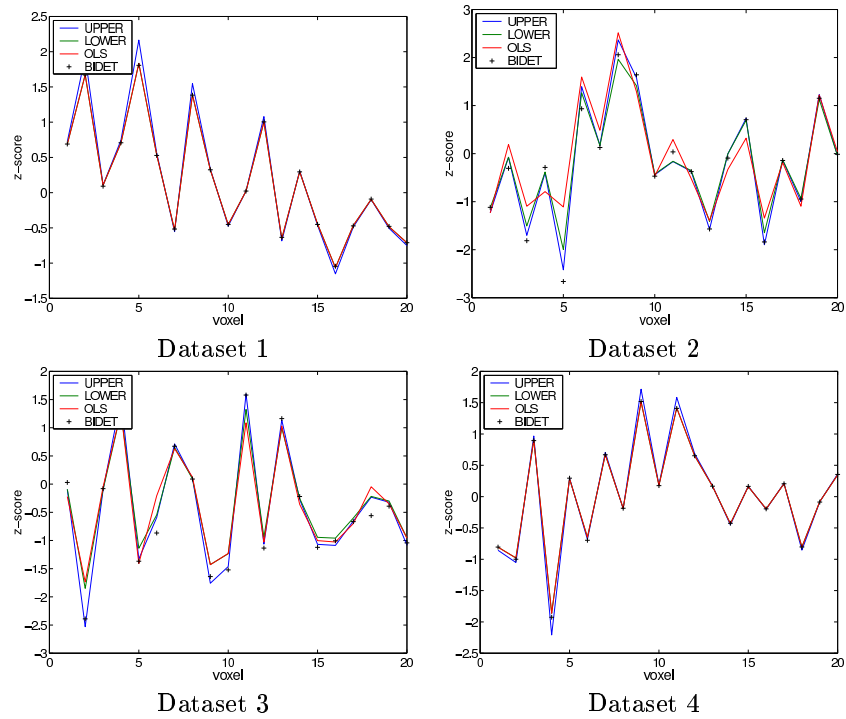


Figure 3: Plots showing the z-statistics for 20 voxels obtained from different inference approaches for the 3 different artificial datasets.

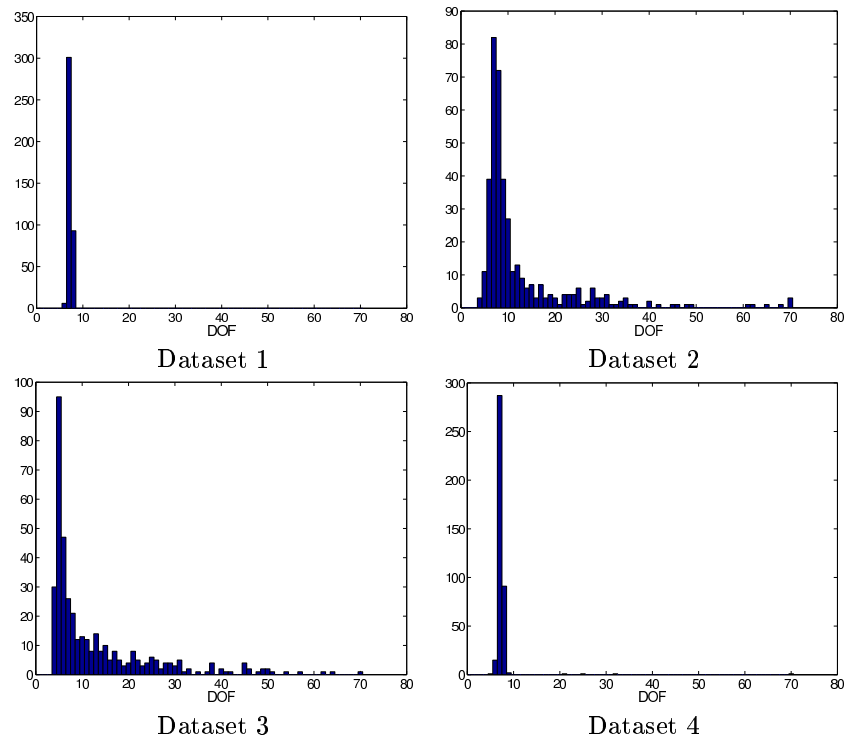


Figure 4: Histograms over 400 voxels of the DOF estimated by [BIDET] for the different datasets for the 3 different artificial datasets.

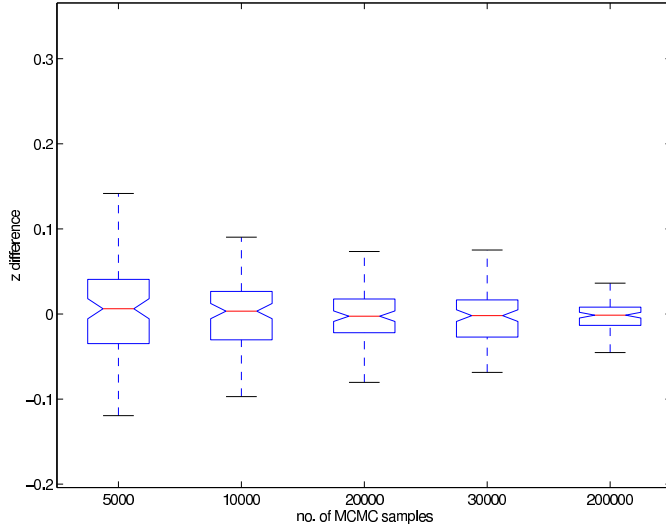


Figure 5: Boxplots over 400 voxels showing the difference in z-statistics between those obtained from a long MCMC chain of 200,000 samples and those obtained from using BIDET on MCMC chains of varying sample sizes on Dataset 1. The box has lines at the lower quartile, median, and upper quartile values. The length of the whiskers is 1.5 times the Inter Quartile Range.

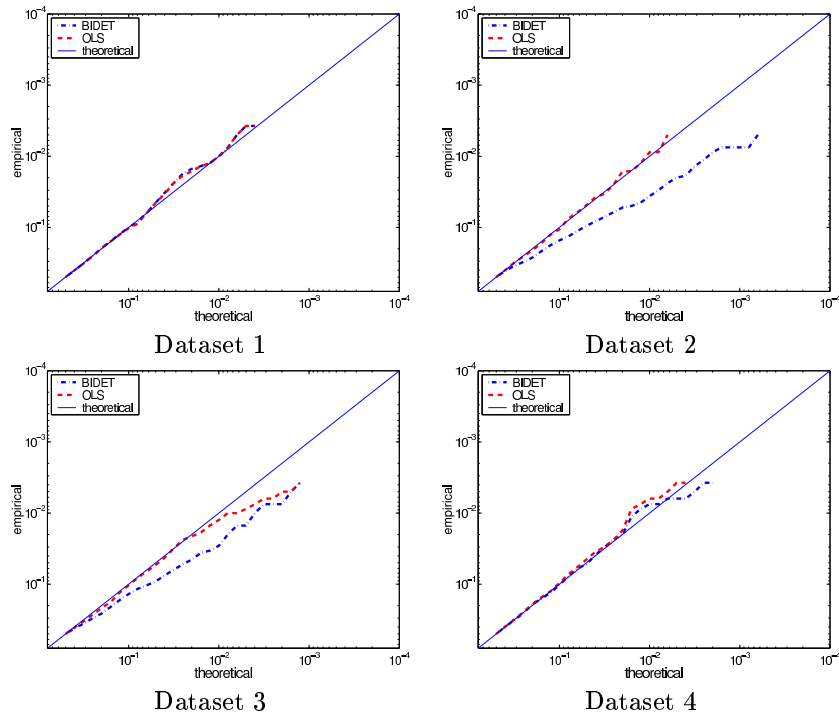


Figure 6: Log probability-log probability plots over 400 voxels for the four different datasets for [BIDET] and [OLS]. These show plots of (nominal/theoretical) FPR against that obtained experimentally from our 3 null artificial datasets. The straight diagonal line shows the result for what would be a perfect match.

For each subject, echo planar images (EPI) were acquired using a 3 Tesla system with TR=3 seconds, time to echo (TE) = 30ms, in-plane resolution 4mm and slice thickness 7mm. The first 4 scans were removed and the data was motion corrected using MCFLIRT (16) and high-pass filtered as described in (25).

The overall model for this group experiment consists of two levels. The first-level is a standard FMRI GLM with a design matrix for subject k , X_k , containing regressors modelling the response to the task within each subject’s dataset. The second-level is a GLM which models the group mean of the individual subject’s responses to the tasks, via a design matrix $X_g = [1, 1, 1, 1, 1, 1, 1]^T$.

To infer on this two-level model we utilise the summary statistic approach we have laid out in this paper. To do this we firstly produce the multivariate non-central t-distribution summary statistics of equation 16 using a first-level analyses of standard generalised least squares (GLS). This GLS analysis was performed using FEAT (FSL). FEAT performs voxel-wise time-series statistical analysis using local autocorrelation estimation to prewhiten the data (25).

To infer the group mean, we now need to infer on the marginal posterior, $p(\beta_g|Y)$, using the multivariate non-central t-distribution summary statistics obtained from these first-level analyses (equation 14).

To do this, we use two different approaches. Firstly, the [OLS] approach as described in section 6. Secondly, a hybrid approach which provides a compromise between the fast posterior approximation approach (section 3.5) and the slower but more accurate approach of using Markov Chain Monte Carlo (MCMC) sampling and the fitting of a non-central multivariate t-distribution (BIDET, section 3.7). The [HYBRID] approach is now described in detail. It is this which is implemented as the FLAME (FMRI’s Local Analysis of Mixed Effects) C++ program used for higher-level analyses in FEAT (part of FSL v3.1).

7.1.1 Hybrid inference approach

Firstly, we can determine bounds on the accuracy of the fast approximation’s z-statistic bounds by using artificial data with “worst case scenario” variance components by comparing the [LOWER] and [UPPER] inference approaches with [BIDET] (as described in section 6). For the design matrices we are using here, the corresponding artificial dataset we need to use is Dataset 1 from section 6.

We can then run the fast approximation approach on our real FMRI data first, and subsequently only run the computationally expensive MCMC sampling (with 30,000 samples and a burnin of 1000 samples) and the fitting of a non-central multivariate t-distribution (BIDET, section 3.7) on voxels at which the desired z threshold lies within the estimated bounds.

This hybrid approach takes approximately 1 hour (for the datasets considered here) on a 2GHz Intel PC on a full volume.

7.1.2 Thresholding

Using [HYBRID], we obtain the marginal posterior, $p(\beta_g|Y)$, as a multivariate non-central t-distribution. We can then use a contrast $c = 1$ to produce $p(c^T\beta_g|Y)$. As discussed in section 6.3 we have a number of choices as to how we infer on this posterior distribution. Here we take the option of performing a t-to-p-to-z transform and mimicking a null-hypothesis frequentist inference (i.e. controlling a FPR) by assuming that under the null hypothesis the z-statistics produced are standardised Normally distributed (see section 6.3). One advantage of doing this is that we can utilise Gaussian Random Field Theory (GRFT) (26; 21). Here we use GRFT to threshold the z-statistic maps and generate activation clusters determined by $z > 2.3$ with a significance threshold of $p = 0.01$.

7.2 Results

Figures 7 and 8 show cluster-thresholded ($z > 2.3, p < 0.01$) group activation for the two motor tasks. Table 9 shows the number of suprathreshold voxels and the maximum z-statistics for the two tasks. Figure 7 shows the results from index finger tapping against rest ([INDEX] dataset). There is a general decrease in z-statistics in potentially activating voxels. This demonstrates the dominance of one of the two possible effects of incorporating first-level variances into the second level estimation process - that is we get an increase in estimated group variance, σ_g , due to it being constrained to be positive. Figure 8 shows the results of a contrast of sequential finger tapping vs index finger tapping ([SEQUENTIAL] dataset). There is a general increase in z-statistics in potentially activating voxels. This demonstrates the dominance of the other possible effect of incorporating first-level variances into the second

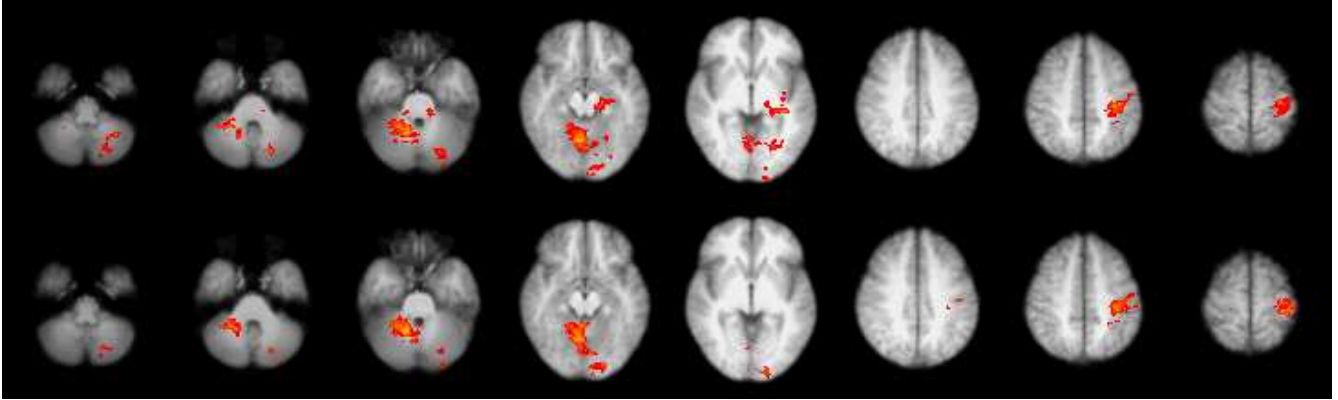


Figure 7: Cluster thresholded ($z > 2.3, p < 0.01$) group activation from the [INDEX] dataset. [top] the [OLS] approach, and [bottom] the [HYBRID] approach.

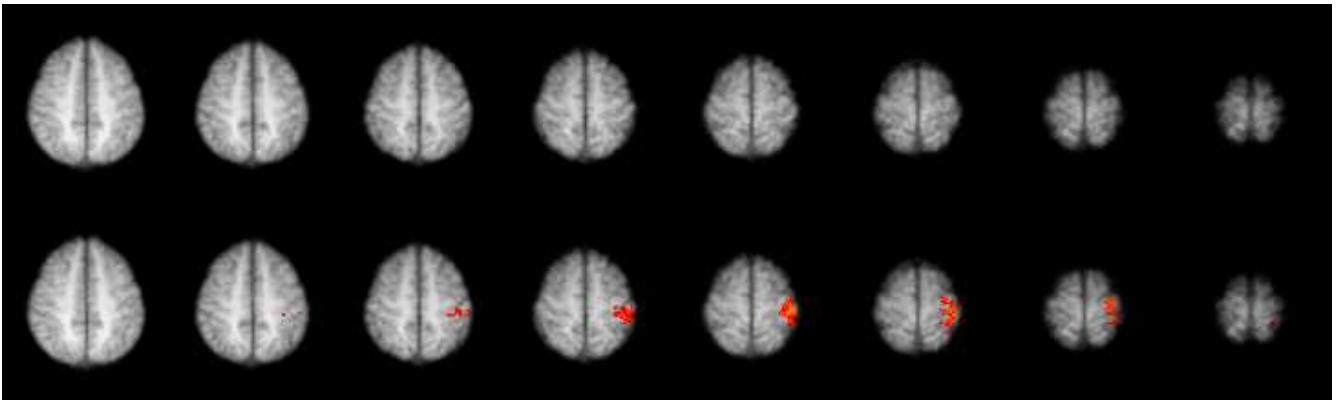


Figure 8: Cluster thresholded ($z > 2.3, p < 0.01$) group activation from the [SEQUENTIAL] dataset. [top] the [OLS] approach, and [bottom] the [HYBRID] approach.

level estimation process - that is we get increased efficiency in parameter estimation due to the use of lower-level variance heterogeneity.

8 Conclusions

We have shown how multi-level hierarchical GLM inference can be split into different levels with the summary statistics of a multivariate non-central t-distribution being passed between the levels. This was achieved by formulating the model in a fully Bayesian framework and using reference analysis to drive our crucial choice of priors (sections 3.3 and 3.4). Using this framework we have proposed two approaches to inferring at the top-level. A fast approximation to the marginal posterior, and a slower approach utilising Markov Chain Monte Carlo (MCMC) followed by a multivariate non-central t-distribution fit to the MCMC chains. These inference approaches are applicable whether we are attempting to infer using the all-in-one approach or the summary statistic split-model approach. We have validated the crucial assumption of the marginal distribution of the GLM regressions parameters being a multivariate non-central t-distribution at levels higher than the first using artificial data. The artificial data also demonstrates the difference between a standard OLS approach and the approach proposed in this paper. We have also shown results on fMRI data.

No. suprathreshold voxels			Max. z		
	[INDEX]	[SEQUENTIAL]		[INDEX]	[SEQUENTIAL]
[OLS]	5206	0	[OLS]	5.20	4.02
[HYBRID]	3861	657	[HYBRID]	4.69	4.22

Figure 9: [left] Number of suprathreshold voxels, and [right] Maximum z-statistic, from the [INDEX] and [SEQUENTIAL] fMRI datasets using the two different inference techniques [OLS] and [HYBRID].

9 Discussion

When we attempt to infer on mixed effects models, we need to deal with the fact that the variance components are unknown. Classically, variance components tend to be estimated separately using iterative estimation schemes employing Ordinary Least Squares (OLS), Expectation Maximisation (EM) or Restricted Maximum Likelihood (ReML), see (22) for details. As an example of a non-Bayesian approach, (Worsley) estimates variance components at each split-level of the model separately. At higher than first levels, they propose EM for estimation of the random effects variance contribution, in order to reduce bias in the variance estimation - a potential problem in higher-level analyses if simple OLS were used. Positivity of the random-effects variance, avoiding what is known as the ‘negative variance problem’ (where mixed-effects variance estimates are lower than fixed-effects variances implying negative random-effects variance (19)), is partially addressed but not strictly enforced.

However, only in certain special cases (not including the model presented here) is it possible to derive analytical forms for the null distributions required by frequentist statistics. In the absence of analytical forms, frequentist solutions rely on null distributions derived from the data using such techniques as permutation tests (20). However, these lose the statistical power gained from educated assumptions about, for example, the distribution of the noise, and limit inference to the number of available points in the empirical null distribution. Bayesian statistics gives us a tool for inferring on *any* model we choose, and *guarantees* that uncertainty will be handled correctly.

(11) have proposed a parametric empirical-Bayesian (PEB) approach for estimation of the all-in-one multi-level model. Unlike (Worsley) they relate the parameters of interest to the full set of original data, i.e. they do not utilise the ‘summary statistics’ approach. Conditional posterior point estimates are generated using EM which give rise to posterior probability maps.

Working in a fully Bayesian reference analysis framework we have the capacity to infer either using the summary statistic split-level (Worsley) approach, or the all-in-one (11) approach. However, all-in-one inference is not part of this paper and is an area of future work. The difference between an all-in-one inference based on the work described in this paper, and the work PEB of (11), is that they assume a multivariate Gaussian marginal posterior for the regression parameters (and then heuristically convert it to a t-statistic), whereas we work in a fully Bayesian framework using reference priors which we can validate as giving a multivariate t-distribution with certain degrees of freedom using MCMC. Without reference priors, (11) have nothing principled to drive the important choice of prior at the top-level and as a result assume flat priors.

Importantly, one of the results demonstrated in this paper is that the inference we would obtain at the *top* level will be approximately the same *regardless* of whether we infer using the summary statistic split-level (Worsley) or the all-in-one approaches (11) (assuming that first-level temporal autocorrelations are effectively known). However, it is very important to realise that there *will* be a difference if we look to infer at intermediate levels in the model. This is because in the all-in-one approach, the regression parameters at these intermediate levels will be regularised by the levels above in the hierarchy, whereas in the split-level approach they will not. Whether or not an experimenter would like to infer on, for example, a subject in isolation, or on a subject in the context of the group of which it is a member, is a choice for the experimenter to make.

Acknowledgements The authors would like to acknowledge support from the UK MRC, EPSRC and GSK. Thanks also to Heidi Johansen-Berg for the data used in this work.

10 Appendix

10.1 Gamma Distribution

x has a two-parameter gamma distribution, denoted by $Ga(a, b)$, with parameters a and b , if its density is given by:

$$\Gamma(x; a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx} \quad (22)$$

where $\Gamma(a)$ is the single-parameter Gamma function. Note, that a two-parameter gamma distribution has mean = a/b and variance = a/b^2 .

10.2 Multivariate Normal distribution

x is a $P \times 1$ random vector and has a multivariate normal distribution, denoted by $N(\mu, \sigma^2 \Sigma)$, if its density is given by:

$$\mathcal{N}(x; \mu, \sigma^2 \Sigma) = \frac{1}{(2\pi)^{P/2} |\sigma^2 \Sigma|^{1/2}} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) \quad (23)$$

The multivariate normal distribution has mean = μ and covariance = $\sigma^2 \Sigma$.

10.3 Multivariate Non-Central t distribution

x is a $P \times 1$ random vector and has a multivariate non-central t distribution, denoted by $t(\mu, \sigma^2 \Sigma, \nu)$, if its density is given by:

$$\mathcal{T}(x; \mu, \sigma^2 \Sigma, \nu) = \frac{\Gamma[(\nu + P)/2]}{(\pi\nu)^{P/2} |\sigma^2 \Sigma|^{1/2} \Gamma[\nu/2]} \left(1 + \frac{(x - \mu)^T \Sigma^{-1} (x - \mu)}{\sigma^2 \nu}\right)^{-(\nu+P)/2} \quad (24)$$

where $\Gamma(a)$ is the single-parameter Gamma function. The non-central t distribution has mean = μ and covariance = $\sigma^2 \Sigma \nu / (\nu - 2)$ for $\nu > 2$.

We can represent a multivariate non-central t distribution using a two-parameter gamma distribution and a multivariate Normal distribution in a Bayesian framework. If we introduce a variable τ , and specify a joint posterior over x and τ as:

$$\begin{aligned} p(\tau, x | \mu, \sigma^2 \Sigma, \nu) &\propto p(x | \tau, \mu, \sigma^2 \Sigma) p(\tau | \nu) \\ x | \tau, \mu, \sigma^2 \Sigma &\sim N(\mu, (\sigma^2 \Sigma / \tau)) \\ \tau | \nu &\sim Ga(\nu/2, \nu/2) \end{aligned} \quad (25)$$

then the marginal posterior for x is a multivariate non-central t distribution, i.e.:

$$\begin{aligned} p(x | \mu, \sigma^2 \Sigma, \nu) &= \int p(\tau, x | \mu, \sigma^2 \Sigma, \nu) d\tau \\ x | \mu, \sigma^2 \Sigma, \nu &\sim t(\mu, \sigma^2 \Sigma, \nu) \end{aligned} \quad (26)$$

10.4 Multivariate Non-central t-distribution fit

In this section we describe how the multivariate non-central t-distribution fit is performed in BIDEt.

Assume that we have $P \times N_J$ matrix, x , with elements, (x_{jp}) , where $j = 1 \dots N_J$ indexes samples and $p = 1 \dots P$ indexes parameters. The task is to fit to these samples a multivariate non-central t-distribution, $t(\mu, \sigma^2 \Sigma, \nu)$ (as described in appendix 10.3).

In BIDEt we constrain the mean of the multivariate non-central t-distribution, μ_{β_g} , to be equal to that from the fast posterior approximation for μ_{β_g} described in section 3.5. If we are not using this constraint then we can set the mean μ to the sample mean, i.e:

$$\mu_p = \frac{1}{N_J} \sum_j x_{jp} \quad (27)$$

We can also directly estimate the normalised covariance Σ using the sample covariance, $\widehat{\Sigma}$:

$$\begin{aligned}\tilde{\Sigma} &= \widehat{\Sigma}/|\widehat{\Sigma}|^{1/P} \\ \widehat{\Sigma} &= (x - M)(x - M)^T / (N_J - 1)\end{aligned}\tag{28}$$

where $M = \{\mu, \mu, \dots, \mu\}^T$.

We still need to estimate σ^2 and ν . Fortunately, we can represent a multivariate non-central t-distribution using a two-parameter gamma distribution and a multivariate Normal distribution in a Bayesian framework, by introducing hidden variables τ_i (see appendix 10.3). With hidden variables we can use the Expectation-Maximisation (EM) algorithm. In the E-step we obtain the expected value of the hidden variables, τ_j :

$$E_{\tau_j | \nu_{(t)}, \sigma_{(t)}^2, x}[\tau_j] = \frac{\sigma_{(t)}^2 (\nu_{(t)} + P)}{\nu_{(t)} \sigma_{(t)}^2 + s_j}\tag{29}$$

where:

$$s_j = (x_j - \mu_j)^T \tilde{\Sigma}^{-1} (x_j - \mu_j)\tag{30}$$

and then in the M-step we can minimise the joint posterior over ν, σ^2 given $\tau_{j(t)} = E_{\tau_j | \nu_{(t)}, \sigma_{(t)}^2, x}[\tau_j]$ to get updates for ν, σ^2 as:

$$\begin{aligned}\sigma_{(t+1)}^2 &= \frac{1}{N_J P} \sum_j \tau_{j(t)} s_j \\ \nu_{(t+1)} &= \frac{2}{1 - \sigma_{(t)}^2 / (\frac{1}{N_J - 1} \sum_j s_j)}\end{aligned}\tag{31}$$

Convergence normally occurs after about 10 iterations. To be conservative we therefore use 50 iterations.

10.5 Determining Reference Priors

Here we show how we determine the reference prior for a vector of parameters θ for a model with likelihood $p(y|\theta)$. This is taken from section 5.4.5. of (2):

The Fisher information matrix, $\mathbf{H}(\theta)$, is given by:

$$H(\theta) = -E_{y|\theta} \left\{ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log p(y|\theta) \right\}\tag{32}$$

For the models in this paper, the Fisher information matrix, $\mathbf{H}(\theta)$, is block diagonal:

$$H(\theta) = \begin{bmatrix} h_{11}(\theta) & 0 & \dots & 0 \\ 0 & h_{22}(\theta) & 0 & \vdots \\ \vdots & 0 & \ddots & 0 \\ 0 & \dots & 0 & h_{mm}(\theta) \end{bmatrix}\tag{33}$$

and we can separate out the block $h_{jj}(\theta)$ as being the product:

$$\{h_{jj}(\theta)\}^{1/2} = f_j(\theta_j) g_j(\theta_{-j})\tag{34}$$

where $f_j(\theta_j)$ is a function depending only on θ_j and $g_j(\theta_{-j})$ does not depend on θ_j . The Berger-Bernardo reference prior is then given by:

$$\pi(\theta) \propto \prod_j^m f_j(\theta_j)\tag{35}$$

Note that this approach yields the Jeffreys prior in one-dimensional problems.

10.6 Marginalising over $(\beta_{\mathbf{K}}, \sigma_{\mathbf{K}}^2)$ in the two-level model

From the two-level model the full joint posterior distribution is (equation 12):

$$p(\beta_g, \sigma_g^2, \beta_{\mathbf{K}}, \sigma_{\mathbf{K}}^2 | Y) \propto \prod_k \{p(Y_k | \beta_k, \sigma_k^2)\} p(\beta_{\mathbf{K}} | \beta_g, \sigma_g^2) p(\beta_g, \sigma_g^2, \sigma_{\mathbf{K}}^2), \quad (36)$$

where the prior is the reference prior for this full two-level model (equation 13):

$$p(\beta_g, \sigma_g^2, \sigma_{\mathbf{K}}^2) = \frac{1}{\sigma_g^2} \prod_k \frac{1}{\sigma_k^2}. \quad (37)$$

If we marginalise out $\sigma_{\mathbf{K}}^2$ then we get:

$$p(\beta_g, \sigma_g^2, \beta_{\mathbf{K}} | Y) \propto \prod_k \left\{ \int p(Y_k | \beta_k, \sigma_k^2) / \sigma_k^2 d\sigma_k^2 \right\} \mathcal{N}(\beta_{\mathbf{K}}; X_g \beta_g, \sigma_g^2 I) 1 / \sigma_g^2 \quad (38)$$

and then substitute in the summary result of the first-level model in isolation (equation 10):

$$p(\beta_g, \sigma_g^2, \beta_{\mathbf{K}} | Y) \propto \prod_k \{ \mathcal{T}(\beta_k; \mu_{\beta_k}, \sigma_{\beta_k}^2 \Sigma_{\beta_k}, \nu_{\beta_k}) \} \mathcal{N}(\beta_{\mathbf{K}}; X_g \beta_g, \sigma_g^2 I) 1 / \sigma_g^2. \quad (39)$$

We can represent a multivariate non-central t-distribution using a two-parameter Gamma distribution and a multivariate Normal distribution (see appendix 10.3). This is achieved by introducing a parameter τ_k for each vector β_k :

$$p(\beta_g, \sigma_g^2, \beta_{\mathbf{K}}, \tau_{\mathbf{K}} | Y) \propto \prod_k \{ \mathcal{N}(\beta_k; \mu_{\beta_k}, (\sigma_{\beta_k}^2 \Sigma_{\beta_k} / \tau_k)) f_{Ga}(\tau_k; \nu_{\beta_k} / 2, \nu_{\beta_k} / 2) \} \mathcal{N}(\beta_{\mathbf{K}}; X_g \beta_g, \sigma_g^2 I) 1 / \sigma_g^2. \quad (40)$$

Writing $\mathcal{N}(\beta_{\mathbf{K}}; X_g \beta_g, \sigma_g^2 I) = \prod_k \mathcal{N}(\beta_k; X_{gk} \beta_g, \sigma_g^2 I)$, where X_{gk} is the k^{th} row vector of the second-level design matrix X_g , we can now easily integrate out β_k for all k to give:

$$p(\beta_g, \sigma_g^2, \tau_{\mathbf{K}} | Y) \propto \prod_k \{ \mathcal{N}(\mu_{\beta_k}; X_{gk} \beta_g, (\sigma_{\beta_k}^2 \Sigma_{\beta_k} / \tau_k) + \sigma_g^2 I) \Gamma(\tau_k; \nu_{\beta_k} / 2, \nu_{\beta_k} / 2) \} 1 / \sigma_g^2 \quad (41)$$

where $\tau_{\mathbf{K}}$ is a $(K \times 1)$ vector of the variables τ_k for $k = 1 \dots K$.

10.7 Fast Approximation Point Estimates

Here we describe how we obtain the point estimates of σ_g^2 and β_g for use in the fast approximation approach described in section 3.5.

We start by rewriting equation 14 as:

$$\begin{aligned} p(\beta_g, \sigma_g^2, \tau_{\mathbf{K}} | Y) &= N(\mu_{\beta_{\mathbf{K}}}; X_G \beta_g, U) 1 / \sigma_g^2 \\ U &= \begin{bmatrix} S_1 & 0 & \cdots & 0 \\ 0 & S_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & S_N \end{bmatrix} \\ S_k &= (\sigma_{\beta_k}^2 \Sigma_{\beta_k} / \tau_k) + \sigma_g^2 I \end{aligned} \quad (42)$$

Point estimate of σ_g^2 We get a point estimate of σ_g^2 by finding the maximum a posteriori (MAP) over the marginal posterior distribution $p(\sigma_g^2, \tau_{\mathbf{K}} | Y)$. If we marginalise out β_g then the marginal posterior is:

$$p(\sigma_g^2, \tau_{\mathbf{K}} | Y) = |U^{-1}|^{1/2} |X_G^T U^{-1} X_G|^{-1/2} \exp \left\{ -\frac{1}{2} \left(\mu_{\beta_{\mathbf{K}}}^T U^{-1} \mu_{\beta_{\mathbf{K}}} - \tilde{\beta}_g^T X_G^T U^{-1} X_G \tilde{\beta}_g \right) \right\} 1 / \sigma_g^2 \quad (43)$$

where

$$\tilde{\beta}_g = (X_G^T U^{-1} X_G)^{-1} X_G^T U^{-1} \mu_{\beta_{\mathbf{K}}} \quad (44)$$

We then assume $\tau_k = 1$ and look to find the MAP for σ_g^2 . However, there is a question of parameterisation. The mode we get will depend on the parametrisation we use. For example, we could look to maximise with respect to σ_g^2 , or σ_g , or $\log(\sigma_g^2)$, or $\phi_g = 1/\sigma_g^2$ etc., all of which will give us different MAPs. Note that as we reparameterise, the reference prior might change but the reference posterior always stays the same, see (2). Hence, a natural way to reparameterise such that the parameter we use gives us a uniform reference prior.

The parameterisation which gives us a uniform reference prior is $\theta = \log(\sigma_g^2)$. Hence, we need to solve:

$$\hat{\theta} = \arg \max_{\theta} p(\sigma_g^2 | Y, \tau_{\mathbf{K}} = \mathbf{1}) \quad (45)$$

where $p(\sigma_g^2 | Y, \tau_{\mathbf{K}} = \mathbf{1})$ is the marginal in equation 43 with $\tau_{\mathbf{K}} = \mathbf{1}$. To solve for $\hat{\theta}$ using this equation we use Brent's algorithm (3). We can then easily convert from $\hat{\theta}$ to $\hat{\sigma}_g^2$.

Point estimate of β_g We approximate $\hat{\beta}_g$ using the point estimate $\hat{\sigma}_g^2$ and $\tau_{\mathbf{K}} = \mathbf{1}$:

$$\hat{\beta}_g = \arg \max_{\beta_g} p(\beta_g | Y, \sigma_g^2 = \hat{\sigma}_g^2, \tau_{\mathbf{K}} = \mathbf{1}) \quad (46)$$

where $p(\beta_g | Y, \sigma_g^2 = \hat{\sigma}_g^2, \tau_{\mathbf{K}} = \mathbf{1})$ is equation 15 with $\sigma_g^2 = \hat{\sigma}_g^2$ and $\tau_{\mathbf{K}} = \mathbf{1}$. The solution to this is:

$$\hat{\beta}_g = (X_G^T U^{-1} X_G)^{-1} X_G^T U^{-1} \mu_{\beta_{\mathbf{K}}} \quad (47)$$

with U as in equation 15, but with $S_k = (\sigma_{\beta_k}^2 \Sigma_{\beta_k}) + \hat{\sigma}_g^2 I$.

References

- [1] Beckmann, C., Jenkinson, M., and Smith, S. (2003). General multi-level linear modelling for group analysis in fMRI. *NeuroImage*, 20:1052–1063. first two authors contributed equally.
- [2] Bernardo, J. and Smith, A. (2000). *Bayesian Theory*. Wiley, New York, USA.
- [3] Brent, R. (1973). *Algorithms for Minimization without Derivatives*. Prentice Hall, Englewood Cliffs, NJ, USA.
- [4] Bullmore, E., Brammer, M., Williams, S., Rabe-Hesketh, S., Janot, N., David, A., Mellers, J., Howard, R., and Sham, P. (1996). Statistical methods of estimation and inference for functional MR image analysis. *Magnetic Resonance in Medicine*, 35(2):261–277.
- [5] Cox, R. (1946). Probability, frequency and reasonable expectation. *American Journal of Physics*, 14:1–13.
- [6] Everitt, B. and Bullmore, E. (1999). Mixture model mapping of brain activation in functional magnetic resonance images. *Human Brain Mapping*, 7:1–14.
- [7] Frison, L. and Pocock, S. (1992). Repeated measures in clinical trials: Analysis using mean summary statistics and its implications for design. *Statistics in Medicine*, 11:1685–1704.
- [8] Friston, K., Josephs, O., Zarahn, E., Holmes, A., Rouquette, S., and Poline, J.-B. (2000). To smooth or not to smooth? *NeuroImage*, 12:196–208.
- [9] Friston, K. J. (2002). Bayesian estimation of dynamical systems: An application to fMRI. *NeuroImage*, 16:513–530.
- [10] Friston, K. J. and Penny, W. (2003). Posterior probability maps and spms. *NeuroImage*, 19(3):1240–1249.
- [11] Friston, K. J., Penny, W., Phillips, C., Kiebel, S., Hinton, G., and Ashburner, J. (2002). Classical and Bayesian inference in neuroimaging: Theory. *NeuroImage*, 16:465–483.
- [FSL] FSL. <http://www.fmrib.ox.ac.uk/fsl>.
- [12] Gamerman, D. (1997). *Markov Chain Monte Carlo*. Chapman and Hall, London.
- [13] Gilks, W., Richardson, S., and Spiegelhalter, D. (1996). *Markov Chain Monte Carlo in Practice*. Chapman and Hall, London.
- [14] Hartvig, N. and Jensen, J. (2000). Spatial mixture modelling of fMRI data. *Human Brain Mapping*, 11(4):233–248.
- [15] Holmes, A. and Friston, K. (1998). Generalisability, random effects & population inference. In *Fourth Int. Conf. on Functional Mapping of the Human Brain*, *NeuroImage*, volume 7, page S754.
- [16] Jenkinson, M., Bannister, P., Brady, J., and Smith, S. (2002). Improved optimisation for the robust and accurate linear registration and motion correction of brain images. *NeuroImage*, 17(2):825–841.
- [17] Kass, R. and Wasserman, L. (1996). Formal rules for selecting prior distributions: A review and annotated bibliography. *jasa*, 91:1343–70.
- [18] Lee, P. (1997). *Bayesian Statistics*. Arnold, London, U.K.
- [19] Leibovici, D. and Smith, S. (2001). Min-max filter for multi-subject analysis. In *Seventh Int. Conf. on Functional Mapping of the Human Brain*, page 10217.
- [20] Nichols, T. E. and Holmes, A. P. (2001). Nonparametric permutation tests for functional neuroimaging: A primer with examples. *Human Brain Mapping*, 15:1–25.
- [21] Poline, J.-B., Worsley, K., Evans, A., and Friston, K. (1997). Combining spatial extent and peak intensity to test for activations in functional imaging. *NeuroImage*, 5:83–96.

- [22] Searle, S., Casella, G., and McCulloch, C. (1992). *Variance Components*. Wiley, New York, USA.
- [23] Woolrich, M., Behrens, T., Beckmann, C., and Smith, S. (2004a). Mixture models with adaptive spatial regularisation for segmentation with an application to fMRI data. in submission.
- [24] Woolrich, M., Jenkinson, M., Brady, J., and Smith, S. (2004b). Fully Bayesian spatio-temporal modelling of fMRI data. *IEEE Trans. on Medical Imaging*, 22(2).
- [25] Woolrich, M., Ripley, B., Brady, J., and Smith, S. (2001). Temporal autocorrelation in univariate linear modelling of fMRI data. *NeuroImage*, 14(6):1370–1386.
- [Worsley] Worsley, K. chapter 14, pages 251–270.
- [26] Worsley, K., Evans, A., Marrett, S., and Neelin, P. (1992). A three-dimensional statistical analysis for CBF activation studies in human brain. *Journal of Cerebral Blood Flow and Metabolism*, 12:900–918.
- [27] Worsley, K. and Friston, K. (1995). Analysis of fMRI time series revisited - again. *NeuroImage*, 2:173–181.
- [28] Worsley, K., Liao, C., Aston, J., Petre, V., Duncan, G., Morales, F., and Evans, A. (2002). A general statistical analysis for fMRI data. *NeuroImage*, 15(1):1–15.