

# Fast methods for posterior inference of two-group normal-normal models

Philip Greengard\*, Jeremy Hoskins†, Charles C. Margossian‡,  
Andrew Gelman§, Aki Vehtari¶

2 Oct 2021

## Abstract

We describe a class of algorithms for evaluating posterior moments of certain Bayesian linear regression models with a normal likelihood and a normal prior on the regression coefficients. The proposed methods can be used for hierarchical mixed effects models with partial pooling over one group of predictors, as well as random effects models with partial pooling over two groups of predictors. We demonstrate the performance of the methods on two applications, one involving U.S. opinion polls and one involving the modeling of COVID-19 outbreaks in Israel using survey data. The algorithms involve analytical marginalization of regression coefficients followed by numerical integration of the remaining low-dimensional density. The dominant cost of the algorithms is an eigendecomposition computed once for each value of the outside parameter of integration. Our approach drastically reduces run times compared to state-of-the-art Markov chain Monte Carlo (MCMC) algorithms. The latter, in addition to being computationally expensive, can also be difficult to tune when applied to hierarchical models.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Mathematical apparatus</b>	<b>5</b>
<b>3</b>	<b>Normalizing constant</b>	<b>6</b>
<b>4</b>	<b>Expectations of <math>\beta</math></b>	<b>8</b>

---

\*Department of Statistics, Columbia University, Corresponding author, Email: pg2118@columbia.edu

†Department of Statistics, University of Chicago

‡Department of Statistics, Columbia University

§Department of Statistics and Political Science, Columbia University

¶Department of Computer Science, Aalto University

<b>5</b>	<b>Covariance of <math>\beta</math></b>	<b>9</b>
<b>6</b>	<b>Numerical implementation</b>	<b>10</b>
<b>7</b>	<b>Special case of two-group model</b>	<b>12</b>
7.1	Numerical apparatus . . . . .	14
<b>8</b>	<b>Mixed effects</b>	<b>15</b>
<b>9</b>	<b>A simple example: Hierarchical linear model</b>	<b>17</b>
<b>10</b>	<b>Application: COVID-19 symptom survey</b>	<b>18</b>
<b>11</b>	<b>Application: Public opinion on abortion policies</b>	<b>23</b>
<b>12</b>	<b>Conclusions and generalizations</b>	<b>25</b>
<b>13</b>	<b>Acknowledgements</b>	<b>26</b>
<b>A</b>	<b>Integral with respect to <math>\rho</math></b>	<b>26</b>

## 1 Introduction

Advances over the last decade in statistical methods and their implementation in open-source, user-friendly software have drastically simplified statistical modeling for applied researchers. For example, with probabilistic programming languages such as Stan [Carpenter et al., 2017] a user can specify and sample from a very general choice of posterior density with flexible language and an easy-to-use interface. For its primary tool of inference, Stan (as well as other probabilistic programming languages) samples from the posterior distribution via dynamic Hamiltonian Monte Carlo sampler (HMC) [Betancourt, 2018, Hoffman and Gelman, 2014]. HMC is a gradient-based sampling method that has become ubiquitous in statistics over the last decade due to its being flexible, reliable, and general.

Despite its widespread use, HMC, as well as other Markov chain Monte Carlo (MCMC) methods, have a substantial drawback in statistical problems with large amounts of data – they can be prohibitively slow (and difficult to tune [e.g. Betancourt et al., 2015]). For example, in the case of a linear regression with  $n$  observations and  $k$  predictors, evaluation of the posterior density requires  $O(nk)$  operations with straightforward implementation. To make matters worse, MCMC methods require large numbers of evaluations of the posterior density, and in the case of HMC, the posterior’s gradient.

Alternative methods for inference have been proposed for problems where MCMC is impractical. These approaches typically involve a suitable approximation of the posterior density with a function with desirable properties. Laplace approximation methods [e.g. Margossian et al., 2020] and variational inference [Blei et al., 2017] are two examples. More generally, there is extensive literature on efficient computational tools and analysis of posterior densities, and there are various software packages devoted to their implementation [see, e.g. Rue et al., 2017, Kristensen et al., 2016].

While these packages, and indeed most of the literature, are devoted to general tools for a wide range of posterior densities, in this paper we introduce an efficient algorithm for computing posterior expectations for two particular classes of Bayesian regression models—two-group normal-normal models and mixed-effects models. These classes of models find a broad range of applications in, for example, social sciences, epidemiology, biochemistry, and environmental sciences [Gelman et al., 2013, Gelman and Hill, 2006, Greenland, 2000, Merlo et al., 2005, Bardini et al., 2017]. Furthermore, in the broader context of model development, these regression models can serve as template models [Gelman et al., 2020].

Using general MCMC methods for sampling from these posteriors can be exceedingly slow for problems with large amounts of data. By specializing to this particular family of models, we leverage their structure to create customized algorithms for fast and accurate inference.

The two Bayesian linear regression models we consider are:

1. **Two group normal-normal:** We define the two-group normal-normal model by

$$\begin{aligned} y &\sim \text{normal}(X_1\beta_1 + X_2\beta_2, \sigma_3) \\ \beta_1 &\sim \text{normal}(0, \sigma_1) \\ \beta_2 &\sim \text{normal}(0, \sigma_2), \end{aligned} \tag{1}$$

where  $X_1$  is a  $n \times k_1$  matrix,  $\beta_1 \in \mathbb{R}^{k_1}$  is a vector of regression coefficients,  $X_2$  is a  $n \times k_2$  matrix, and  $\beta_2 \in \mathbb{R}^{k_2}$  is a vector of regression coefficients. For Bayesian inference, we assume priors on the scale parameters  $\sigma_1, \sigma_2, \sigma_3$ . The performance of the algorithm is largely independent to the choice of these priors. In the models that we use in this paper, we assign independent weakly informative  $\text{normal}^+(0, 1)$  priors on the variance parameters  $\sigma_1, \sigma_2, \sigma_3$  (assuming  $y$  and the columns of  $X$  have been normalized to have standard deviation 1).

2. **Mixed effects:** The mixed-effects model differs slightly from the two-group normal-normal model. Instead of modeling the scale parameter  $\sigma_2$ , fixed scale parameters are assigned to the normal priors on  $\beta_2$ . The mixed-effects model is defined by

$$\begin{aligned} y &\sim \text{normal}(X_1\beta_1 + X_2\beta_2, \sigma_3) \\ \beta_1 &\sim \text{normal}(0, \sigma_1) \\ \beta_{2,i} &\sim \text{normal}(0, \sigma_{2,i}), \end{aligned} \tag{2}$$

where  $\sigma_{2,i}$  is the fixed scale parameter prior on each regression coefficient  $\beta_{2,i}$  for  $i = 1, \dots, k_2$  where  $\beta_2 \in \mathbb{R}^{k_2}$ . We will assume priors on the scale parameters  $\sigma_1, \sigma_3$ .

The models we discuss in this paper are standard models of Bayesian statistics and appear when seeking to model an outcome,  $y$ , as a linear combination of two (or more) distinct groups of predictors. The Gaussian prior on the predictors enable various strategies commonly used in statistical modeling and machine learning; notably regularization and partial pooling between various sources of data. We demonstrate these models on three applications.

1. **COVID-19:** Due to a lack of reliable, fast, and widespread testing, an online survey initiative was created in Israel [Rossman et al., 2020] for tracking and predicting COVID-19 outbreaks. We constructed a mixed-effects model for estimating geographic and age effects on the spread of the virus. With tens of thousands of responses, straightforward implementation of MCMC methods takes hours. Using the methods of this paper, we obtain accurate posterior inference in seconds.
2. **Rat growth:** We demonstrate the efficiency of our two-group algorithm on the classical two-group model for rat growth [Gelfand et al., 1990], which estimates the growth rates of a population of rats over the first few weeks of life.
3. **Public opinion on abortion:** We use 2018 results of the annual Cooperative Congressional Election Study (CCES) to estimate geographic and demographic effects on attitudes towards abortion. The CCES contains nearly 100,000 responses, and performing inference via MCMC sampling can be prohibitively slow. We use the mixed-effects algorithm introduced in this paper to perform posterior inference in seconds.

The computational methods we introduce for the two-group normal-normal model and the mixed-effects models are closely related. In fact, the mixed-effects model is a special case of the two-group model. We organize this paper by first describing our algorithm for the two-group normal-normal model in detail, and then outline the minor modifications that allow for efficient evaluation of mixed-effects models.

The unnormalized density corresponding to the two-group model is given by

$$q(\beta, \sigma_1, \sigma_2, \sigma_3) = \frac{e^{-\sigma_1^2/2 - \sigma_2^2/2 - \sigma_3^2/2}}{\sigma_1^n \sigma_2^{k_1} \sigma_3^{k_2}} e^{-\frac{1}{2\sigma_1^2} \|X\beta - y\|^2} e^{-\frac{1}{2\sigma_2^2} \|\beta_1\|^2} e^{-\frac{1}{2\sigma_3^2} \|\beta_2\|^2}, \quad (3)$$

where  $\beta = (\beta_1, \beta_2)$  with  $\beta_1 \in \mathbb{R}^{k_1}$ ,  $\beta_2 \in \mathbb{R}^{k_2}$ ,  $\beta \in \mathbb{R}^k$ , and  $y \in \mathbb{R}^n$ . For convenience, we will be denoting by  $\sigma$  the vector of scale parameters  $(\sigma_1, \sigma_2, \sigma_3) \in \mathbb{R}^3$ .

In the methods of this paper, we compute posterior moments of  $q$  by analytically reducing the calculation of moments from integrals over  $k + 3$  dimensions to 3-dimensional integrals. We then integrate the remaining 3-dimensional integrals with a tensor product of Gaussian nodes. For example, we evaluate the normalizing constant  $C$  and posterior means for the regression coefficients,  $\beta$ , via

$$C = \int_0^\infty \int_{-\infty}^\infty q(\beta, \sigma_1, \sigma_2, \sigma_3) d\beta d\sigma \approx \sum_{i=1}^n f(\sigma_i) w_i$$

$$E[\beta_j] = \int_0^\infty \int_{-\infty}^\infty \beta_j q(\beta, \sigma_1, \sigma_2, \sigma_3) d\beta d\sigma \approx \frac{1}{C} \sum_{i=1}^n f_j(\sigma_i) w_i,$$

where  $\sigma_i \in \mathbb{R}^3$  and  $w_i \in \mathbb{R}$  are three-dimensional Gaussian nodes and weights [Trefethen, 2020] and

$$f(\sigma_i) = \int_{\mathbb{R}^k} q(\beta, \sigma_i) d\beta \quad (4)$$

$$f_j(\sigma_i) = \int_{\mathbb{R}^k} \beta_j q(\beta, \sigma_i) d\beta. \quad (5)$$

Integrals (4) and (5) can be evaluated analytically via well-known equations [Lindley and Smith, 1972], but a straightforward implementation of those equations results in a computational cost of  $O(m^3k^3)$  operations where  $m$  is the number of discretization nodes needed in each dimension. In the methods of this paper, we improve the computational cost of those integrals to  $O(mk^3 + m^2k^2 + m^3)$  operations after a change of variables, allowing for the rapid evaluation of marginals.

The tools used in the algorithm of this paper are a generalization of the approach proposed by Greengard et al. [2021] and generalize to higher-dimensional multilevel and higher-dimensional multigroup posterior distributions. Since we integrate the marginal density using a tensor product of Gaussian nodes, the cost of the integration scales like  $O(m^d)$  where  $m$  is the number of discretization nodes in each direction and  $d$  is the dimension of the marginalized integral (where  $d = 3$  in the models of this paper). As a result, higher dimensional problems require evaluation of marginal integrals via sampling-based algorithms and cannot rely solely on Gaussian quadrature. We leave the analysis and description of numerical tools for such models to a subsequent publication.

The structure of this paper is as follows. In the following section we describe the change of variables of (3) and provide mathematical analysis that will be used in the algorithm of this paper. Section 3 includes formulas that will allow for the evaluation of the normalizing constant of (3). In Section 4 we describe analysis that will be used in computing expectations and in Section 5 we describe formulas for second moments. In Section 6 we discuss details of the implementation of the algorithm. Section 7 contains an algorithm for a special case of the two-group normal-normal model in which one group is much smaller than the other. The mixed effects algorithm, or rather the modification of the two-group normal-normal algorithm, is contained in Section 8. In Section 10, Section 9 and Section 11 we apply the algorithms of this paper to applications. Conclusions and generalizations of the algorithm of this paper are presented in Section 12.

## 2 Mathematical apparatus

Over the next several sections, we describe a numerical algorithm for computing expectations and second moments of the density  $q : \mathbb{R}^{k+3} \rightarrow \mathbb{R}^+$  defined by

$$q(\beta, \sigma_1, \sigma_2, \sigma_3) = \frac{e^{-\sigma_1^2/2 - \sigma_2^2/2 - \sigma_3^2/2} e^{-\frac{1}{2\sigma_1^2} \|X\beta - y\|^2} e^{-\frac{1}{2\sigma_2^2} \|\beta_1\|^2} e^{-\frac{1}{2\sigma_3^2} \|\beta_2\|^2}}{\sigma_1^n \sigma_2^{k_1} \sigma_3^{k_2}},$$

where  $\sigma_1, \sigma_2, \sigma_3 > 0$ ,  $\beta = (\beta_1, \beta_2)$  with  $\beta_1 \in \mathbb{R}^{k_1}, \beta_2 \in \mathbb{R}^{k_2}, \beta \in \mathbb{R}^k$ , and  $y \in \mathbb{R}^n$ . We begin by introducing notation that will be used throughout the numerical sections of this paper.

Let  $y = X\tilde{\beta} + d$ , where  $X^t d = 0$ , so that  $\tilde{\beta}$  is the least-squares solution to the linear system  $X\beta = y$  and  $d$  is the residual. Let  $I_1$  be the diagonal  $k \times k$  matrix with ones in the first  $k_1$  places on the diagonal, and zeroes in the remaining  $k_2$  places. Similarly, let  $I_2$  be the diagonal  $k \times k$  matrix with zeroes in the first  $k_1$  places on the diagonal, and one in the remaining  $k_2$  places.

We perform a change of variables in  $\sigma_1, \sigma_2$ , and  $\sigma_3$ , defining  $\rho, \theta$ , and  $\phi$  implicitly by

$$\begin{aligned}\sigma_1 &= \rho \cos \phi, \\ \sigma_2 &= \rho \sin \phi \cos \theta, \\ \sigma_3 &= \rho \sin \phi \sin \theta.\end{aligned}$$

This corresponds to changing to spherical coordinates in the  $\sigma$  variables. With these substitutions, and with some minor abuse of notation, we obtain

$$f(\beta, \rho, \theta, \phi) = \frac{e^{-\rho^2/2}}{\rho^{n+k} \cos^n(\phi) \sin^k \phi \cos^{k_1} \theta \sin^{k_2} \theta} \exp \left[ -\frac{1}{2\rho^2} \left( \frac{1}{\cos^2 \phi} \|X(\beta - \tilde{\beta})\|^2 + \frac{\|d\|^2}{\cos^2 \phi} + \frac{\beta^t \left( \frac{I_1}{\cos^2 \theta} + \frac{I_2}{\sin^2 \theta} \right) \beta}{\sin^2 \phi} \right) \right].$$

The differentials transform as follows:

$$d\sigma_1 d\sigma_2 d\sigma_3 = \rho^2 \sin \phi d\rho d\theta d\phi.$$

Moreover, the condition that  $\sigma_1, \sigma_2, \sigma_3 > 0$ , is equivalent to  $0 < \phi, \theta < \pi/2$ .

### 3 Normalizing constant

In this section we describe the computation of the integral of  $f$  over its domain. If we denote this quantity by  $I_0$ , then  $f/I_0$  is a probability density on  $\mathbb{R}^+ \times \mathbb{R}^+ \times \mathbb{R}^+ \times \mathbb{R}^k$ . We begin by making a change of variables in  $\beta$ , setting  $\beta = (I_1 \cos \theta + I_2 \sin \theta)z$ . Similarly, we define  $\tilde{z}$  implicitly by  $\tilde{\beta} = (I_1 \cos \theta + I_2 \sin \theta)\tilde{z}$ . Then,

$$f(z, \rho, \theta, \phi) = \frac{e^{-\rho^2/2}}{\rho^{n+k} \sin^k \phi \cos^n \phi \cos^{k_1} \theta \sin^{k_2} \theta} \exp \left[ -\frac{1}{2\rho^2} \left( \frac{1}{\cos^2 \phi} \|X_\theta(z - \tilde{z})\|^2 + \frac{\|d\|^2}{\cos^2 \phi} + \frac{\|z\|^2}{\sin^2 \phi} \right) \right],$$

where  $X_\theta = X(I_1 \cos \theta + I_2 \sin \theta)$ . The differentials transform as follows

$$d\beta_1 \dots d\beta_k = \cos^{k_1} \theta \sin^{k_2} \theta dz_1 \dots dz_k.$$

To diagonalize the quadratic form appearing in the exponent, we perform an eigendecomposition of on  $X_\theta^t X_\theta$  obtaining

$$X_\theta^t X_\theta = V_\theta D_\theta V_\theta^t, \tag{6}$$

where  $D_\theta \in \mathbb{R}^{k \times k}$  is diagonal and positive, and  $V_\theta \in \mathbb{R}^{k \times k}$  is a unitary matrix. We have assumed here that  $n \geq k$ . If the converse is true then a small modification is required. In the following, for notational convenience we denote the diagonal entries of  $D_\theta$  by  $\lambda_i(\theta)$ .

Next we set  $z = V_\theta w$  and  $\tilde{z} = V_\theta \tilde{w}$ . In terms of the original variables,  $\beta = (I_1 \cos \theta + I_2 \sin \theta)V_\theta w$  and  $\tilde{\beta} = (I_1 \cos \theta + I_2 \sin \theta)V_\theta \tilde{w}$ . In particular,

$$d\beta_1 \dots d\beta_k = \cos^{k_1} \theta \sin^{k_2} \theta dw_1 \dots dw_k.$$

After making these substitutions, we obtain

$$\begin{aligned} & \int \dots \int f(\beta, \sigma_1, \sigma_2, \sigma_3) d\beta_1 \dots d\beta_k \\ &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \int \dots \int \exp \left[ -\frac{1}{2\rho^2} \sum_{i=1}^k \left( \frac{\lambda_i (w_i - \tilde{w}_i)^2}{\cos^2 \phi} + \frac{w_i^2}{\sin^2 \phi} \right) \right] dw_1 \dots dw_k \\ &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i=1}^k \int \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_i (w_i - \tilde{w}_i)^2}{\cos^2 \phi} + \frac{w_i^2}{\sin^2 \phi} \right) \right] dw_i. \end{aligned}$$

Thus, the integrals over the  $\beta$  variables have been reduced to the product of  $k$  one-dimensional Gaussian integrals. Using the identity

$$\int_{\mathbb{R}} e^{-\frac{a}{2}(s-s_0)^2 - \frac{b}{2}s^2} ds = \sqrt{\frac{2\pi}{a+b}} e^{-\frac{ab s_0^2}{2(a+b)}},$$

we find that

$$\begin{aligned} & \int \dots \int f(\beta, \sigma_1, \sigma_2, \sigma_3) d\beta_1 \dots d\beta_k \\ &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i=1}^k \sqrt{\frac{2\pi\rho^2}{\frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}} \exp \left[ -\frac{1}{2\rho^2} \tilde{w}_i^2 \frac{\lambda_i}{\cos^2 \phi \sin^2 \phi \left( \frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \right] \\ &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^n \cos^{n-k} \phi} \prod_{i=1}^k \sqrt{\frac{2\pi}{\lambda_i \sin^2 \phi + \cos^2 \phi}} \exp \left[ -\frac{1}{2\rho^2} \frac{\lambda_i \tilde{w}_i^2}{\lambda_i \sin^2 \phi + \cos^2 \phi} \right]. \end{aligned}$$

Next, we define the functions  $\alpha : (0, \pi/2)^2 \rightarrow \mathbb{R}^+$  and  $\beta : (0, \pi/2)^2 \rightarrow \mathbb{R}^+$  by

$$\begin{aligned} \alpha(\phi, \theta) &= \prod_{i=1}^k \sqrt{\frac{2\pi}{\lambda_i(\theta) \sin^2 \phi + \cos^2 \phi}}, \\ \beta(\phi, \theta) &= \sum_{i=1}^k \frac{\lambda_i(\theta) \tilde{w}_i^2(\theta)}{\lambda_i(\theta) \sin^2 \phi + \cos^2 \phi}. \end{aligned}$$

Then

$$\int \dots \int f(\beta, \sigma_1, \sigma_2, \sigma_3) d\beta_1 \dots d\beta_k = \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^n \cos^{n-k} \phi} \alpha(\phi, \theta) e^{-\frac{1}{2\rho^2} \beta(\phi, \theta)}.$$

For a fixed  $\theta$ , the vector  $\tilde{w}$  and the eigenvalues values  $\lambda_i$  are independent of  $\rho$  and  $\phi$ , and hence need to be recomputed only when  $\theta$  is changed. Moreover, for fixed  $\theta$  and  $\phi$ , the above

integral can be computed in  $O(1)$  floating operations for each new value of  $\rho$ . In particular, the normalization constant can be computed efficiently via the formula

$$I_0 = \int_0^{\frac{\pi}{2}} \int_0^{\frac{\pi}{2}} \frac{\alpha(\phi, \theta) \sin \phi}{\cos^{n-k} \phi} \int_0^\infty \frac{e^{-\frac{\rho^2}{2} - \frac{1}{2\rho^2} \left( \frac{\|d\|^2}{\cos^2 \phi} + \beta(\phi, \theta) \right)}}{\rho^n} \rho^2 d\rho d\phi d\theta. \quad (7)$$

## 4 Expectations of $\beta$

In this section we describe how to compute moments in  $\beta$  of the distribution  $f$ . We begin by observing that

$$\beta_\ell = Q_\ell(\theta) \sum_{\ell=1}^k (V_\theta)_{\ell,j} w_j,$$

where  $Q_\ell(\theta) = \cos \theta$  if  $1 \leq \ell \leq k_1$ , and  $\sin \theta$  if  $k_1 < \ell \leq k$ . Let  $M_j(\rho, \phi, \theta)$  be defined by

$$\begin{aligned} M_j(\rho, \phi, \theta) &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i \neq j}^k \int \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_i (w_i - \tilde{w}_i)^2}{\cos^2 \phi} + \frac{w_i^2}{\sin^2 \phi} \right) \right] dw_i \\ &\quad \times \int w_j \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_j (w_j - \tilde{w}_j)^2}{\cos^2 \phi} + \frac{w_j^2}{\sin^2 \phi} \right) \right] dw_j. \end{aligned}$$

Using the results of the previous section we can perform the integrals to obtain

$$\begin{aligned} M_j(\rho, \phi, \theta) &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i \neq j}^k \sqrt{\frac{2\pi\rho^2}{\frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}} \exp \left[ -\frac{1}{2\rho^2} \tilde{w}_i^2 \frac{\lambda_i}{\cos^2 \phi \sin^2 \phi \left( \frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \right] \\ &\quad \times \frac{\lambda_j \tilde{w}_j}{\cos^2 \phi \left( \frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \sqrt{\frac{2\pi\rho^2}{\frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}} \exp \left[ -\frac{1}{2\rho^2} \tilde{w}_j^2 \frac{\lambda_j}{\cos^2 \phi \sin^2 \phi \left( \frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \right] \\ &= \left( \frac{\lambda_j \tilde{w}_j}{\lambda_j \sin^2 \phi + \cos^2 \phi} \right) \sin^2 \phi \left[ \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^n \cos^{n-k} \phi} \alpha(\phi, \theta) e^{-\frac{1}{2\rho^2} \beta(\phi, \theta)} \right]. \end{aligned}$$

We remark that the second factor in the above expression is the same as the one arising in the computation of  $I_0$ . The first factor depends on  $j$ ,  $\phi$ , and  $\theta$  but not on  $\rho$ . For ease of exposition, let us define  $\tilde{M}_j$  by

$$\tilde{M}_j(\phi, \theta) = \left( \frac{\lambda_j \tilde{w}_j}{\lambda_j \sin^2 \phi + \cos^2 \phi} \right). \quad (8)$$

Then

$$E[\beta_\ell] = \sum_{\ell,j} \frac{1}{I_0} \int_0^{\frac{\pi}{2}} Q_\ell(\theta) (V_\theta)_{\ell,j} \int_0^{\frac{\pi}{2}} \frac{\alpha(\phi, \theta) \sin^3 \phi}{\cos^{n-k} \phi} \tilde{M}_j(\phi, \theta) \int_0^\infty \frac{e^{-\frac{\rho^2}{2} - \frac{1}{2\rho^2} \left( \frac{\|d\|^2}{\cos^2 \phi} + \beta(\phi, \theta) \right)}}{\rho^n} \rho^2 d\rho d\phi d\theta. \quad (9)$$



In particular, if  $I_1^{(j)}$  is defined via the formula

$$I_1^{(j)}(\theta) = \int_0^{\frac{\pi}{2}} \frac{\alpha(\phi, \theta) \sin^3 \phi}{\cos^{n-k} \phi} \tilde{M}_j(\phi, \theta) \int_0^\infty \frac{e^{-\frac{\rho^2}{2} - \frac{1}{2\rho^2} \left( \frac{\|d\|^2}{\cos^2 \phi} + \beta(\phi, \theta) \right)}}{\rho^n} \rho^2 d\rho d\phi, \quad (10)$$

then

$$E[\beta_\ell] = \int_0^{\frac{\pi}{2}} \sum_{\ell, j} Q_i(\theta) (V_\theta)_{\ell, j} I_1^{(j)}(\theta) d\theta, \quad (11)$$

and hence all  $k$  moments can be computed simultaneously with an integrand requiring  $O(k^2)$  operations to compute (assuming the number of quadrature nodes required to achieve a fixed precision is more or less independent of  $k$ ).

## 5 Covariance of $\beta$

In this section, we describe formulas for computing the posterior covariance matrix of  $\beta$ . We use the identity

$$-\frac{a}{2}(s - s_0)^2 - \frac{b}{2}s^2 = -\frac{a+b}{2} \left( s - \frac{as_0}{-(a+b)} \right)^2 + \frac{-abs_0^2}{2(a+b)} \quad (12)$$

to compute second moments of  $w_j$ . That is, letting  $P_j(\rho, \phi, \theta)$  be defined by

$$\begin{aligned} P_j(\rho, \phi, \theta) &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i \neq j}^k \int \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_i (w_i - \tilde{w}_i)^2}{\cos^2 \phi} + \frac{w_i^2}{\sin^2 \phi} \right) \right] dw_i \\ &\quad \times \int w_j^2 \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_j (w_j - \tilde{w}_j)^2}{\cos^2 \phi} + \frac{w_j^2}{\sin^2 \phi} \right) \right] dw_j \\ &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i \neq j}^k \int \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_i (w_i - \tilde{w}_i)^2}{\cos^2 \phi} + \frac{w_i^2}{\sin^2 \phi} \right) \right] dw_i \\ &\quad \times \int w_j^2 \exp \left[ -\frac{1}{2\rho^2} \left( \frac{\lambda_j (w_j - \tilde{w}_j)^2}{\cos^2 \phi} + \frac{w_j^2}{\sin^2 \phi} \right) \right] dw_j. \end{aligned}$$

we use (12) to obtain

$$\begin{aligned} P_j(\rho, \phi, \theta) &= \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^{n+k} \cos^n \phi \sin^k \phi} \prod_{i \neq j}^k \sqrt{\frac{2\pi\rho^2}{\frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}} \exp \left[ -\frac{1}{2\rho^2} \tilde{w}_i^2 \frac{\lambda_i}{\cos^2 \phi \sin^2 \phi \left( \frac{\lambda_i}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \right] \\ &\quad \times \frac{\lambda_j \tilde{w}_j}{\cos^2 \phi \left( \frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \sqrt{\frac{2\pi\rho^2}{\frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}} \exp \left[ -\frac{1}{2\rho^2} \tilde{w}_j^2 \frac{\lambda_j}{\cos^2 \phi \sin^2 \phi \left( \frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi} \right)} \right] \\ &= \left( \frac{\rho^2}{\frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}} \right) \sin^2 \phi \left[ \frac{e^{-\frac{\rho^2}{2} - \frac{\|d\|^2}{2\rho^2 \cos^2 \phi}}}{\rho^n \cos^{n-k} \phi} \alpha(\phi, \theta) e^{-\frac{1}{2\rho^2} \beta(\phi, \theta)} \right] - M_j(\rho, \phi, \theta)^2. \end{aligned}$$

Defining  $\tilde{P}_j$  by

$$\tilde{P}_j(\phi, \theta) = \frac{\rho^2}{\frac{\lambda_j}{\cos^2 \phi} + \frac{1}{\sin^2 \phi}}$$

and defining  $I_2^{(j)}$  via the formula

$$I_2^{(j)}(\theta) = \int_0^{\frac{\pi}{2}} \frac{\alpha(\phi, \theta) \sin^3 \phi}{\cos^{n-k} \phi} \tilde{P}_j(\phi, \theta) \int_0^\infty \frac{e^{-\frac{\rho^2}{2} - \frac{1}{2\rho^2} \left( \frac{\|d\|^2}{\cos^2 \phi} + \beta(\phi, \theta) \right)}}{\rho^n} \rho^2 d\rho d\phi, \quad (13)$$

we observe that

$$E[\beta\beta^t]_{i,m} = \int_0^{\frac{\pi}{2}} (CI_2(\theta)C^t)_{i,j} d\theta$$

where  $C$  is the  $k \times k$  matrix defined by

$$C_{i,j} = \sum_m Q_i(\theta)(V_\theta)_{m,j}.$$

We then compute the posterior covariance of  $\beta$  via

$$E[(\beta - E[\beta])(\beta - E[\beta])^t] = E[\beta\beta^t] - E[\beta]E[\beta]^t, \quad (14)$$

where  $E[\beta]$  is obtained via (11).

## 6 Numerical implementation

We now describe a numerical approach for computing the normalizing constant  $I_0$ , see (7), and moments of  $q$  (see (3)). The quadrature rules used provide arbitrary user-specified precision for both the normalizing constant as well as moments. Our integration scheme is a tensor product of Gaussian nodes in the  $\theta$ ,  $\phi$ , and  $\rho$  directions.

In the remainder of this section we describe how to adaptively determine integration bounds in the  $\phi$  and  $\rho$  directions as well as the number of nodes to use in  $\theta$ .

### Integrals with respect to $\rho$

We first describe an approach for integrating the inner integral,

$$\int_0^\infty \frac{e^{-\frac{\rho^2}{2} - \frac{1}{2\rho^2} \left( \frac{\|d\|^2}{\cos^2 \phi} + \beta(\phi, \theta) \right)}}{\rho^{n-2}} d\rho. \quad (15)$$

We do this by finding upper and lower integration bounds and then performing Gaussian quadrature with 80 nodes. The maximum of the integrand of (15), which we denote  $\rho_{max}$ , is achieved at

$$\rho_{max} = \frac{1}{\sqrt{2}} \left( \sqrt{4 \frac{\|d\|^2}{\cos^2 \phi} + 4\beta(\phi, \theta) + (n-2)^2} - n + 2 \right)^{1/2}. \quad (16)$$

See Appendix A for details. Furthermore, for all  $\rho > \rho_{max}$ , the integrand is monotonically decreasing and for all  $\rho < \rho_{max}$ , the integrand is monotonically increasing. We choose our upper integration bound,  $\rho_1$ , to be the value  $\rho_1 > \rho$  such that the integrand of (15) is smaller than its maximum by a factor of  $10^{20}$ . We evaluate  $\rho_1$  with bisection with initial bounds of

$$\rho_{max} \quad \text{and} \quad \rho_{max} + \frac{40}{\sqrt{-\sigma}}, \quad (17)$$

where  $\sigma$  is the second derivative of the integrand of (15) with respect to  $\rho$  evaluated at  $\rho_{max}$  (see Appendix A). We also find the lower bound of integration using bisection where our starting bounds for the bisection are 0 and  $\rho_{max}$ . After finding the bounds of integration, we evaluate (15) using Gaussian quadrature. Using 80 nodes was sufficient for full double precision accuracy in our examples.

### Integral with respect to $\phi$

We evaluate the integral  $I_1^{(j)}(\theta)$  (see (10)) using Gaussian quadrature where we determine the upper bound of integration adaptively. That is, we determine the upper bound of integration,  $\phi_1$ , by sequentially evaluating the integrand of (10) at order 100 Gaussian nodes until the integrand is smaller than the maximum observed value of the integrand by a factor of  $10^{20}$ . We then declare the first such point,  $\phi_1$ , to be the upper integration bound and evaluate (10) with 80 Gaussian nodes on the interval  $(0, \phi_1)$ .

### Integral with respect to $\theta$

In this section, we describe a technique for evaluating the outer integral of (11). For each evaluation of the integrand of (11), we perform eigendecomposition (6). As a result, each evaluation of the integrand requires  $O(k^3)$  operations and the total computational time for evaluating moments of posterior (3) is roughly linear in the number of evaluations of the integrand of (11).

We compute integral (11) using Gaussian quadrature. We found that for a general set of problems, the integrand of (11) was sufficiently smooth, that using around 10 nodes was sufficient for several digits of accuracy. To check accuracy of an  $m$ -point Gaussian quadrature we compute the same integral with  $2m$  nodes and check the difference, which is a proxy for the accuracy of the  $m$ -point Gaussian quadrature.

For problems in which the number of evaluations of the integrand of (11) needs to be reduced, we can use the following approach. Approximate the integral using two Chebyshev nodes (see, e.g., Trefethen [2020]) and four Chebyshev nodes and compute the difference between the two values. That difference provides an estimate for the error of the approximation using two nodes. If the difference between the two approximations is larger than desired, then double the number of nodes and again compute the difference. Continue to double the number of nodes until the desired accuracy is achieved.

By using practical Chebyshev nodes in  $\theta$  the order  $n$  nodes are a subset of the order  $2n$  nodes. This saves  $n$  evaluations of the integrand when approximating the integral with  $2n$

nodes.

---

**Algorithm 1:** *Two-group normal-normal models: Evaluation of posterior expectations.*

---

- 1 Construct  $\theta_1, \dots, \theta_m$ , Gaussian nodes and weights in  $\theta$  on  $[0, 2\pi]$
  - 2     For each  $\theta_i$ :
  - 3         Compute eigendecomposition (6)
  - 4         Evaluate upper integration bound in  $\phi$
  - 5         Construct  $\phi_1, \dots, \phi_m$ , Gaussian nodes in  $\phi$
  - 6         For each  $\phi_\ell$ :
  - 7             Evaluate integration bounds for  $\rho$  integral
  - 8             Compute integral using Gaussian nodes
  - 9             Evaluate  $\tilde{M}_j(\phi_\ell, \theta_i)$  of (8) for  $j = 1, \dots, k$
  - 10         Convert expectations back to  $\beta$
- 

## 7 Special case of two-group model

In this section we consider the special case where the posterior density  $f$  corresponds to an intercept model and one of the two groups is much smaller than the other. That is,  $k_1$  is small relative to  $k_2$ , and each row of  $X$  contains two non-zero entries—one in the first  $k_1$  columns and another in the next  $k_2$  columns. This model appears in applications in which we seek to model an outcome,  $y$ , as a combination of two unrelated factors.

For the corresponding posterior, we introduce a fast algorithm that analytically marginalizes the normal-normal parameters  $\beta$  by evaluating a determinant and solving a linear system.

The terms dependent on  $\beta$  of the exponential of  $f$  can be written as

$$\frac{1}{\sigma_1^2}(\beta - \beta_0)^t X^t X (\beta - \beta_0) + \beta^t R \beta = (\beta - \bar{\beta}(\sigma))^t \left( \frac{1}{\sigma_1^2} X^t X + R \right) (\beta - \bar{\beta}(\sigma)) + C(\sigma),$$

where

$$\sigma_1^2 \left( \frac{1}{\sigma_1^2} X^t X + R \right) \bar{\beta}(\sigma) = X^t X \beta_0 = X^t y \quad (18)$$

and

$$\begin{aligned} C(\sigma) &= \frac{1}{\sigma_1^2} (\beta_0 - \bar{\beta}(\sigma))^t X^t X \beta_0 \\ &= \frac{1}{\sigma_1^2} (\beta_0 - \bar{\beta}(\sigma))^t X^t y \\ &= \frac{1}{\sigma_1^2} (y - y_r - X \bar{\beta}(\sigma))^t y \\ &= \frac{1}{\sigma_1^2} (y_p - X \bar{\beta}(\sigma))^t y. \end{aligned}$$

Before proceeding further, we introduce a change of variables. Let

$$(\sigma_1, \sigma_2, \sigma_3) = \rho(\cos \theta, \nu_2, \sin \theta)$$

and note that

$$d\sigma_1 d\sigma_2 d\sigma_3 = \rho^2 d\rho d\nu_2 d\theta.$$

Then  $R$  of (18) satisfies

$$R = \frac{1}{\rho^2} \begin{bmatrix} \frac{1}{\nu_2^2} I_{k_1} & 0 \\ 0 & \frac{1}{\sin^2 \theta} I_{k_2} \end{bmatrix}.$$

Let  $\tilde{R} = \rho^2 \cos^2 \theta R$ ; this is a function only of  $\nu_2$  and  $\theta$ . It follows that

$$f(\beta, \rho, \nu_2, \theta) = \frac{1}{\rho^{n+k} \nu_2^{k_1} \cos^n \theta \sin^{k_2} \theta} e^{-\frac{\rho^2}{2}(1+\nu_2^2) - \frac{1}{2\rho^2 \cos^2 \theta} (\beta - \bar{\beta})^t (X^t X + \tilde{R}) (\beta - \bar{\beta}) - \frac{1}{2\rho^2 \cos^2 \theta} [y_p^t y - \bar{\beta} X^t y + \|y_r\|^2]}.$$

For ease of exposition, let  $\beta$  be the function defined by

$$\beta(\nu_2, \theta) = \frac{1}{\cos^2 \theta} [y_p^t y - \bar{\beta} X^t y + \|y_r\|^2].$$

The following lemma will be used in Theorem 1 and will provide formulas for computing posterior moments of  $f$ .

**Lemma 1** (moments of a Gaussian). *Let  $C$  be a symmetric positive definite  $k \times k$  matrix. Then, for any vector  $\tilde{\beta} \in \mathbb{R}^k$ ,*

1.

$$\int_{\mathbb{R}^k} e^{-(\beta - \tilde{\beta})^t C (\beta - \tilde{\beta})} d\beta = \frac{\pi^{\frac{k}{2}}}{\sqrt{\det C}}$$

2.

$$\int_{\mathbb{R}^k} \beta e^{-(\beta - \tilde{\beta})^t C (\beta - \tilde{\beta})} d\beta = \frac{\pi^{\frac{k}{2}}}{\sqrt{\det C}} \tilde{\beta}$$

3.

$$\int_{\mathbb{R}^k} \beta \beta^t e^{-(\beta - \tilde{\beta})^t C (\beta - \tilde{\beta})} d\beta = \frac{\pi^{\frac{k}{2}}}{\sqrt{\det C}} \left( \frac{1}{2} C^{-1} + \tilde{\beta} \tilde{\beta}^t \right)$$

*Proof.* The proof of the first two is immediate. For the second, let  $z = \sqrt{C}(\beta - \tilde{\beta})$ . Then, the integral becomes

$$\begin{aligned} & \int_{\mathbb{R}^k} [\sqrt{C}^{-1} z + \tilde{\beta}] [\sqrt{C}^{-1} z + \tilde{\beta}]^t e^{-z^2} \frac{1}{\det \sqrt{C}} dz \\ &= \int_{\mathbb{R}^k} \left( \sqrt{C}^{-1} z z^t \sqrt{C}^{-1} + \tilde{\beta} \tilde{\beta}^t \right) e^{-z^2} \frac{1}{\det \sqrt{C}} dz \\ &= \frac{\pi^{\frac{k}{2}}}{\sqrt{\det C}} \left( \frac{1}{2} C^{-1} + \tilde{\beta} \tilde{\beta}^t \right). \end{aligned}$$

□

The following theorem follows immediately from the previous lemma and provides formulas that will be used to compute posterior moments of  $f$ .

**Theorem 1.** *Let  $f$  be the unnormalized probability density defined above. Then*

1.

$$\begin{aligned} f_0 &:= \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^k} f(\beta, \sigma) \, d\beta \, d\sigma \\ &= (2\pi)^{\frac{k}{2}} \int_{\mathbb{R}^+} \frac{1}{\cos^{n-k} \theta \sin^{k_2} \theta} \int_{\mathbb{R}^+} \frac{1}{\nu_2^{k_1}} \frac{1}{\sqrt{\det X^t X + \tilde{R}}} \int_{\mathbb{R}^+} \frac{e^{-\frac{\rho^2(1+\nu_2^2)}{2} - \frac{1}{2\rho^2}\beta}}{\rho^{n-2}} \, d\rho \, d\nu_2 \, d\theta \end{aligned}$$

2.

$$\begin{aligned} f_j &:= \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^k} \beta_j f(\beta, \sigma) \, d\beta \, d\sigma \\ &= (2\pi)^{\frac{k}{2}} \int_{\mathbb{R}^+} \frac{1}{\cos^{n-k} \theta \sin^{k_2} \theta} \int_{\mathbb{R}^+} \frac{1}{\nu_2^{k_1}} \frac{\bar{\beta}_j(\nu_2, \theta)}{\sqrt{\det X^t X + \tilde{R}}} \int_{\mathbb{R}^+} \frac{e^{-\frac{\rho^2(1+\nu_2^2)}{2} - \frac{1}{2\rho^2}\beta}}{\rho^{n-2}} \, d\rho \, d\nu_2 \, d\theta \end{aligned}$$

3.

$$\begin{aligned} f_{jk} &:= \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^+} \int_{\mathbb{R}^k} \beta_j \beta_k f(\beta, \sigma) \, d\beta \, d\sigma \\ &= (2\pi)^{\frac{k}{2}} \int_{\mathbb{R}^+} \frac{1}{\cos^{n-k} \theta \sin^{k_2} \theta} \int_{\mathbb{R}^+} \frac{1}{\nu_2^{k_1}} \frac{\frac{1}{2}(X^t X + \tilde{R})_{jk}^{-1} + \bar{\beta}_j \bar{\beta}_k}{\sqrt{\det X^t X + \tilde{R}}} \int_{\mathbb{R}^+} \frac{e^{-\frac{\rho^2(1+\nu_2^2)}{2} - \frac{1}{2\rho^2}\beta}}{\rho^{n-2}} \, d\rho \, d\nu_2 \, d\theta. \end{aligned}$$

## 7.1 Numerical apparatus

In this section, we describe a numerical method for exploiting the structure of  $X^t X$  and  $\tilde{R}$  to obtain a computationally efficient strategy for computing moments of  $f$ .

We use the following lemma to compute the determinant found in each integral of Theorem 1. The lemma follows immediately from the Schur complement formula [Trefethen and Bau, 1997].

**Lemma 2.** *Let  $X^t X$  be the matrix given by*

$$X^t X = \begin{bmatrix} D_1 & B^t \\ B & D_2 \end{bmatrix},$$

and  $\tilde{R}$  be the matrix given by

$$\tilde{R} = \begin{bmatrix} \frac{\cos^2 \theta}{\nu_2^2} I_{k_1} & 0 \\ 0 & \frac{1}{\tan^2 \theta} I_{k_2} \end{bmatrix}.$$

Then,

$$\det X^t X + \tilde{R} = \det \left( D_2 + \frac{1}{\tan^2 \theta} \right) \det \left( D_1 + \frac{\cos^2 \theta}{\nu_2^2} - B^t \left( D_2 + \frac{1}{\tan^2 \theta} \right)^{-1} B \right).$$

In the following lemma, which follows immediately from elementary linear algebra operations, we show that  $X^t X + \tilde{R}$  can be written as  $D + UU^t$  where  $U$  is a  $k \times 2k_1$  matrix and  $D$  is diagonal. This will be used to solve the linear system (18) via the Sherman-Morrison-Woodbury formula.

**Lemma 3.** *Suppose  $X^t X$  and  $\tilde{R}$  are the same matrices as defined previously. Then*

$$X^t X + \tilde{R} = \begin{bmatrix} D_1 + \frac{\cos^2 \theta}{\nu_2^2} & 0 \\ 0 & D_2 + \frac{1}{\tan^2 \theta} \end{bmatrix} + \begin{bmatrix} I_{k_1} & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} 0 & I_{k_1} \\ B & 0 \end{bmatrix}^t.$$

Moreover,

$$\begin{bmatrix} 0 & I_{k_1} \\ B & 0 \end{bmatrix}^t \begin{bmatrix} \left(D_1 + \frac{\cos^2 \theta}{\nu_2^2}\right)^{-1} & 0 \\ 0 & \left(D_2 + \frac{1}{\tan^2 \theta}\right)^{-1} \end{bmatrix} \begin{bmatrix} I_{k_1} & 0 \\ 0 & B \end{bmatrix} = \begin{bmatrix} 0 & B^t \left(D_2 + \frac{1}{\tan^2 \theta}\right)^{-1} B \\ \left(D_1 + \frac{\cos^2 \theta}{\nu_2^2}\right)^{-1} & 0 \end{bmatrix}$$

The following theorem allows us to efficiently solve the linear system (18) and follows immediately from the combination of the previous lemma with the Sherman-Morrison-Woodbury formula [Trefethen and Bau, 1997].

**Theorem 2.** *Let  $X^t X$  and  $\tilde{R}$  be as defined above. Let  $W(\nu_3)$  be the  $k_1 \times k_1$  matrix defined by*

$$W(\nu_3) = B^t \left( D_2 + \frac{1}{\tan^2 \theta} \right)^{-1} B.$$

Then

$$= - \begin{bmatrix} P_1 & 0 \\ 0 & P_2 \end{bmatrix} \begin{bmatrix} I_{k_1} & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} I & W \\ P_1 & I \end{bmatrix}^{-1} \begin{bmatrix} 0 & I_{k_1} \\ B & 0 \end{bmatrix}^t \begin{bmatrix} P_1 & 0 \\ 0 & P_2 \end{bmatrix},$$

where  $P_1 = \left(D_1 + \frac{\cos^2 \theta}{\nu_2^2}\right)^{-1}$ , and  $P_2 = \left(D_2 + \frac{1}{\tan^2 \theta}\right)^{-1}$ . Moreover,

$$\begin{bmatrix} I & W \\ P_1 & I \end{bmatrix}^{-1} = \begin{bmatrix} (I - P_1 W)^{-1} & -W(I - P_1 W)^{-t} \\ -P(I - P_1 W)^{-1} & (I - P_1 W)^{-t} \end{bmatrix}.$$

Using the preceding theorems and the quadrature rule described in Section 6, we compute posterior moments of the special case of  $f$  described in this section.

## 8 Mixed effects

From a computational standpoint, the mixed effects model is nearly identical to the two-group normal-normal model, however they differ in one key respect. In the two-group normal-normal model (see (1)), the scale parameters  $\sigma_1$  and  $\sigma_2$  are treated as unknowns that are fit to the data and assigned priors—they're treated as modeled coefficients, also called "random effects." In the mixed effects model,  $\sigma_1$  is still a random effect, however instead of treating

$\sigma_2$  as a random effect, each regression coefficient in the second group of predictors,  $\beta_{2,i}$  is given a normal prior with fixed scale parameter. The corresponding Bayesian model is

$$\begin{aligned} y &\sim \text{normal}(X_1\beta_1 + X_2\beta_2, \sigma_3) \\ \beta_1 &\sim \text{normal}(0, \sigma_1) \\ \beta_{2,i} &\sim \text{normal}(0, \sigma_{2,i}), \end{aligned} \tag{19}$$

where  $\sigma_{2,i}$  is the fixed scale parameter prior on each regression coefficient  $\beta_{2,i}$  for  $i = 1, \dots, k_2$  where  $\beta_2 \in \mathbb{R}^{k_2}$ . For the purposes of demonstrating the algorithm in this paper, we assign the priors

$$\begin{aligned} \sigma_1 &\sim \text{normal}^+(0, 1) \\ \sigma_3 &\sim \text{normal}^+(0, 1). \end{aligned}$$

The choice of priors on  $\sigma_1$  and  $\sigma_3$  is somewhat arbitrary. The algorithm described allows for general choices for these priors. The unnormalized posterior density for the mixed effects model,  $f : \mathbb{R}^{k+2} \rightarrow \mathbb{R}$ , is given by

$$f(\beta, \sigma_1, \sigma_2) = \frac{e^{-\sigma_1^2/2 - \sigma_2^2/2}}{\sigma_1^n \sigma_2^{k_1}} e^{-\frac{1}{2\sigma_1^2} \|X\beta - y\|^2} e^{-\frac{1}{2\sigma_2^2} \|\beta_1\|^2} e^{-\sum_{i=1}^{k_2} \frac{\beta_{2,i}^2}{2\sigma_{3,i}}},$$

where  $\sigma_3 \in \mathbb{R}^{k_2}$  is a vector of fixed scale parameter priors for the regression coefficients  $\beta_2 \in \mathbb{R}^{k_2}$ . Now our density is over  $k+2$  dimensions—the  $k$  regression coefficients,  $\beta = (\beta_1, \beta_2)$ , and the two scale parameters,  $\sigma_1$  and  $\sigma_2$ .

With a change of variables we convert the posterior density,  $f$ , to the posterior density of a two-group normal-normal model where one scale parameter is fixed. That is, we now convert  $f(\beta, \sigma_1, \sigma_2)$  to  $q(\beta, \sigma_1, \sigma_2, 1)$  where  $q$  is the posterior density of a two-group normal-normal model (see (3)).

We first scale the last  $k_2$  columns of  $X$  by  $\sigma_3^2$  and define  $\hat{X}$  to be resulting matrix:

$$\hat{X}_{i,j} = X_{i,k_1+j} \sigma_{3,j}^2$$

for  $i = 1, \dots, n$  and for  $j = 1, \dots, k_2$ . We define  $\hat{\beta}_2$  to be the vector

$$\hat{\beta}_{2,i} = \beta_{2,i} \sigma_i^2,$$

and define  $\hat{f}$  by the unnormalized density

$$\hat{f}(\hat{\beta}, \sigma_1, \sigma_2) = \frac{e^{-\sigma_1^2/2 - \sigma_2^2/2}}{\sigma_1^n \sigma_2^{k_1}} e^{-\frac{1}{2\sigma_1^2} \|\hat{X}\hat{\beta} - y\|^2} e^{-\frac{1}{2\sigma_2^2} \|\beta_1\|^2} e^{-\frac{1}{2} \|\hat{\beta}_2\|^2}.$$

It follows that posterior means and standard deviations of  $\beta_2$  become

$$E_f[\beta_{2,i}] = \sigma_i^2 E_{\hat{f}}[\hat{\beta}_{2,i}], \quad E_f[(\beta_{2,i} - E_f[\beta_{2,i}])^2] = \sigma_i E_{\hat{f}}[(\hat{\beta}_{2,i} - E_{\hat{f}}[\hat{\beta}_{2,i}])^2],$$

and all other posterior first and second moments are unchanged under density  $\hat{f}$ . We've now reduced the problem of finding moments of  $f$  to finding moments of  $\hat{f}$ , which is equal to  $q(\hat{\beta}, \sigma_1, \sigma_2, 1)$  where  $q$  is defined in (3).



At this point, we rely on the analysis and numerical tools of the two-group normal-normal model for evaluating posterior moments. The only difference between evaluation of moments of the two-group model is that in the two group model the marginal density is a 3-dimensional density over  $(\sigma_1, \sigma_2, \sigma_3)$  whereas in the mixed effects model, the marginal density is over two dimensions,  $(\sigma_1, \sigma_2)$  and  $\sigma_3 = 1$  is fixed. As a result, we perform the change of variables from  $(\sigma_1, \sigma_2, 1)$  to polar coordinates

$$\begin{aligned}\sigma_1 &= \rho \cos(\phi) \\ \sigma_2 &= \rho \sin(\phi) \cos(\theta) \\ 1 &= \rho \sin(\phi) \sin(\theta),\end{aligned}$$

or equivalently

$$\begin{aligned}\rho &= \sqrt{\sigma_1^2 + \sigma_2^2 + 1} \\ \theta &= \text{atan}(1/\sigma_2) \\ \phi &= \text{acos}\left(\frac{\sigma_1}{\sqrt{\sigma_1^2 + \sigma_2^2 + 1}}\right).\end{aligned}$$

The differentials become

$$d\sigma_1 d\sigma_2 = |\gamma|^{-1} d\theta d\phi,$$

where

$$\gamma = \frac{-1}{1 + \sigma_2^2} \left( \frac{1}{\sqrt{\alpha}} - \frac{\sigma_1^2}{\alpha^{3/2}} \right) \left( 1 - \frac{\sigma_1^2}{\alpha} \right)^{-1/2}$$

and

$$\alpha = \sigma_1^2 + \sigma_2^2 + 1.$$

We can now evaluate the posterior of  $\hat{f}$  with Algorithm 1, where the integral with respect to  $\rho$  is replaced with  $\rho = \sqrt{\sigma_1^2 + \sigma_2^2 + 1}$ .

## 9 A simple example: Hierarchical linear model

We demonstrate the algorithm on a hierarchical linear model describing the growth of a group of young rats over a period of several weeks; this is a small example that has been used in the statistical literature [Gelfand et al., 1990]. In the experiment, the weight of each rat is measured at regular time intervals. Regression coefficients are computed for each rat; that is, for the  $j^{\text{th}}$  rat, we estimate an intercept  $\alpha_j$  and a linear coefficient  $\beta_j$ . We assign a normal prior on both parameters and estimate the prior scale. The full model is as follows:

$$\begin{aligned}y_i &\sim \text{normal}(X_i^1 \alpha + X_i^2 \beta, \sigma_1) \\ \alpha_j &\sim \text{normal}(0, \sigma_2) \\ \beta_j &\sim \text{normal}(0, \sigma_3) \\ \sigma_k &\sim \text{normal}^+(0, 10) \text{ for } k = 1, 2, 3,\end{aligned}\tag{20}$$

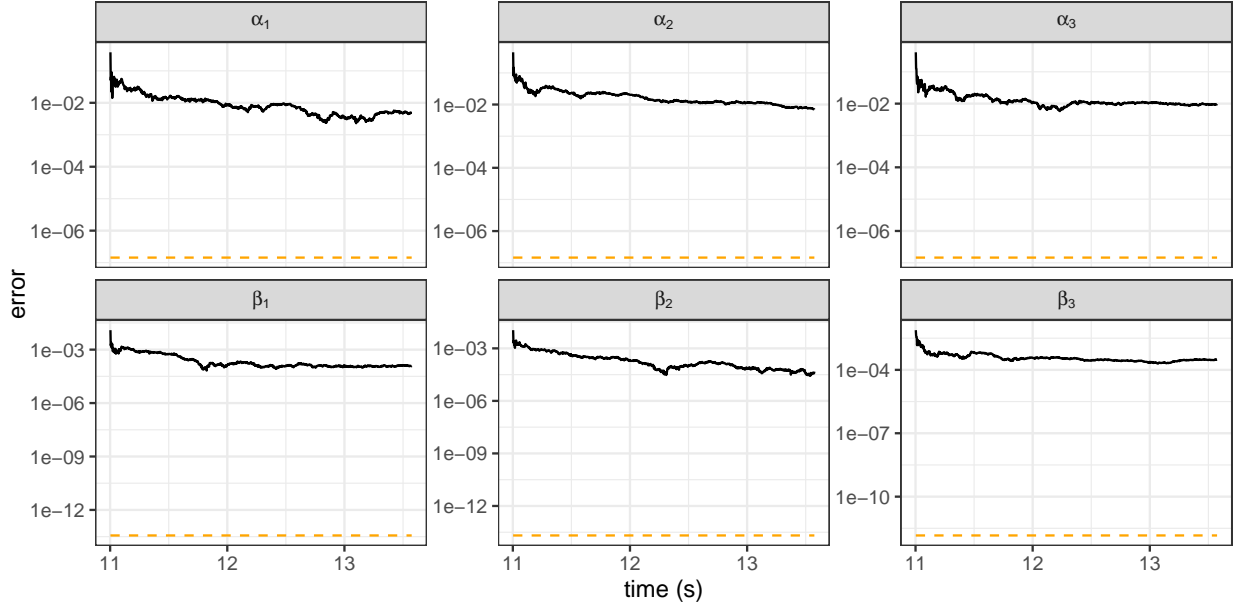


Figure 1: *Error of MCMC estimates via Stan as a function of run time. The horizontal orange line is the error of Algorithm 1.*

where  $X_1$  is an indicator matrix indicating to which rat each observation (weighing) corresponds;  $X_2$  is that same indicator matrix multiplied by  $w - \bar{w}$ , where  $w$  is the observation week and  $\bar{w}$  the mean observation week. In other words, we have an intercept and a slope parameter for each rat. The data is centered at 0 and the priors on the scale parameters are weakly informative.

We demonstrate the efficiency of the two-group normal-normal Algorithm 1 on evaluating posterior means and standard deviations of the rats model on randomly generated data. We assume an experiment with 100 rats and 20 weighing times and randomly generated data for each weighing. As a result, matrix  $X_1$  and  $X_2$  of model (20) are  $2000 \times 100$  matrices.

Because the data size is relatively small and the data matrices have a friendly structure, running MCMC with Stan (4 chains in parallel, each with 1,000 warmup iterations and 1,000 sampling iterations) only takes 13.9s. Our algorithm takes 1.6s and achieves significantly smaller errors than MCMC estimates. Figure 1 shows the error of the MCMC estimates as a function of time and the accuracy achieved by our algorithm. For problems with larger data, the difference in time scale becomes important.

## 10 Application: COVID-19 symptom survey

As of the writing of this article, the coronavirus pandemic is still raging in many countries and stressing healthcare systems around the world. A challenge at the start of the pandemic was tracking its spread, especially in locations where reliable testing was not widely available. Having accurate estimates of infection rates across geographical regions can be extremely helpful. For example, reliable estimates allow hospital systems to allocate resources efficiently, they can alert residents of the need to take extra precaution in their daily rou-

tines, and they can facilitate better policy from local governments. In order to get improved estimates of infection rates in the absence of widespread testing, initiatives were deployed in early 2020 in several countries that allowed individuals to report symptoms via publicly available surveys [e.g., Segal et al., 2020].

One country where these surveys provided valuable information was Israel [Rossman et al., 2020], where demographic and health data was provided by tens of thousands of respondents across the country. The large amount of data collected from survey respondents provided data scientists and policy-makers with a great resource, however at the same time, large amounts of data turns the computational aspect of statistical modeling into a substantial challenge.

In this section, we present an exploratory model used to analyze data from the COVID-19 survey conducted in Israel [Rossman et al., 2020]. Using straightforward MCMC with Stan [Carpenter et al., 2017] was inconvenient; using the full data set resulted in runtimes of several hours. Using the algorithms of this paper, we were able to evaluate posterior moments in seconds without loss of accuracy.

## Multilevel regression and poststratification procedure

The respondents are anonymous, but several of their features are recorded, including their age and the city in which they live. We can use the data to identify regions in which the average symptom score seems unusually high.

A first exploratory model uses an intercept, age group, and population density in the respondent’s city, as covariates,  $X$ , and an indicator matrix  $Z$  for city:

$$y_i \sim \text{normal}(X\beta + Zu, \sigma_1),$$

with a hierarchical prior on the city parameters,

$$u_j \sim \text{normal}(0, \sigma_2),$$

and weakly informative priors on the other coefficients,

$$\beta_j \sim \text{normal}(0, 1).$$

This unit prior is weakly informative if the outcome  $y_i$  has been standardized and the continuous predictors (in this case, population density) has also been standardized to be on unit scale.

In addition, we put weakly informative half-normal priors (standard normal distributions restricted to the non-negative reals) on the hyperparameters  $\sigma_1$  and  $\sigma_2$ :

$$\begin{aligned} \sigma_1 &\sim \text{normal}^+(0, 1) \\ \sigma_2 &\sim \text{normal}^+(0, 1). \end{aligned}$$

This corresponds to a two-group normal-normal model with an additional covariate. In cities where  $u$  cannot be well estimated due to a low response rate, we can rely on the rest of the model, that is a regression model based on age and population density.

Only a fraction of the population responds to the survey, which raises questions about biases. This is notably a concern because different age groups behave differently: not only do their chances of contracting and spreading the disease vary, their susceptibility to the disease also changes. In multilevel regression and poststratification (MRP), we adjust for these biases by using estimates of the proportion of people in each city that belong to each age group. For this model, the proportions are estimated using census data. This leads to a corrected estimate for the expected symptom score of an individual in city  $i$ :

$$\tilde{u}_i = u_i + \beta_0 + \beta_{\text{density}}d_i + \sum_{j=1}^n a_j^i \beta_{\text{age},j},$$

where  $\beta_0$  is the intercept,  $\beta_{\text{density},i}$  is the regression coefficient of the population density covariate,  $d_i$  denotes the density of city  $i$ ,  $a_j^i$  is the proportion of individuals in the  $j^{\text{th}}$  age group in the  $i^{\text{th}}$  city, and  $\beta_{\text{age},j}$  is the regression coefficient of age group  $j$ .

Using the means and covariances of  $u, \beta_0, \beta_{\text{density}}$ , and  $\beta_{\text{age}}$  we compute the posterior mean and variance for  $\tilde{u}$ , per the following formulas. Given a linear combination of random variables,  $Y = \sum_i \delta_i Z_i$ , we have

$$EY = \sum_i \delta_i EZ_i,$$

and

$$\text{Var}Y = \sum_i \delta_i^2 \text{Var}Z_i + 2 \sum_{i < j} \delta_i \delta_j \text{Cov}(Z_i, Z_j).$$

Moreover, variance and thence standard deviations of  $\tilde{u}$  can be computed, provided we also evaluate the relevant posterior covariances.

## Comparison of our algorithm to MCMC

We analyze the data collected over the two weeks between April 15<sup>th</sup> and 30<sup>th</sup> 2020, across 351 cities. These are cities for which we know, through census data, the population density and the age distribution. The total number of responses is 135,501.

Our proposed algorithm returns the posterior mean and standard deviation for all variables of interest and takes  $\sim 7$ s to run.

We next fit the model in Stan using the default dynamic HMC sampler. After warming up the sampler for 500 iterations, we compute another 500 draws, using 4 chains computed in parallel, for a total of 2,000 sampling iterations. The wall time for this procedure is  $\sim 12,000$ s ( $> 3$  hours). For each city, we computed the Monte Carlo mean. Figure 2 plots the posterior mean and standard deviation of  $\tilde{u}$  for all cities, computed by both methods. Figure 3 shows the difference between our algorithm and the Monte Carlo estimate, as a function of computation time. While it takes on the order of hours to get accurate results with MCMC, our algorithm achieves better results within seconds.

## Limitations of the model and our numerical method

We believe the presented model offers an improvement on the analysis conducted on the survey data [Rossman et al., 2020], because (i) it uses full Bayesian inference to quantify

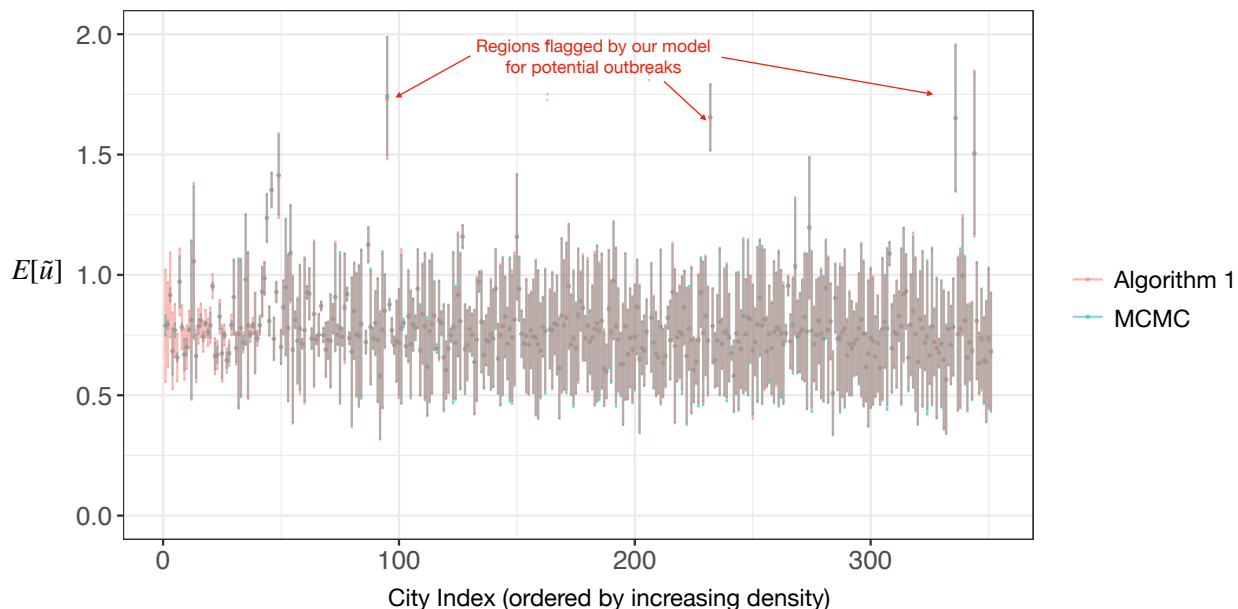


Figure 2: Posterior mean and standard deviation for  $\tilde{u}$  computed using Algorithm 1 and MCMC. The points represent the estimated mean and the “error bars” span two standard deviations.

Regression coefficient	MCMC Accuracy ( $\sim 12,000$ s)	Algorithm 1 Accuracy (7 s)
$\beta_0$	$3e-02$	$1e-05$
$u_1$	$1e-02$	$4e-04$
$u_2$	$1e-02$	$3e-04$
$u_3$	$2e-02$	$6e-05$
$\beta_{\text{age},1}$	$3e-02$	$7e-06$
$\beta_{\text{age},2}$	$3e-02$	$8e-06$
$\beta_{\text{age},3}$	$3e-02$	$9e-08$
$\beta_{\text{density}}$	$2e-03$	$2e-05$

Table 1: Accuracy of the approximation of posterior means for several regression coefficients with both MCMC and Algorithm 1. For MCMC, the total time for the approximation was  $\sim 12,000$ s. Total time for Algorithm 1 was 7s.

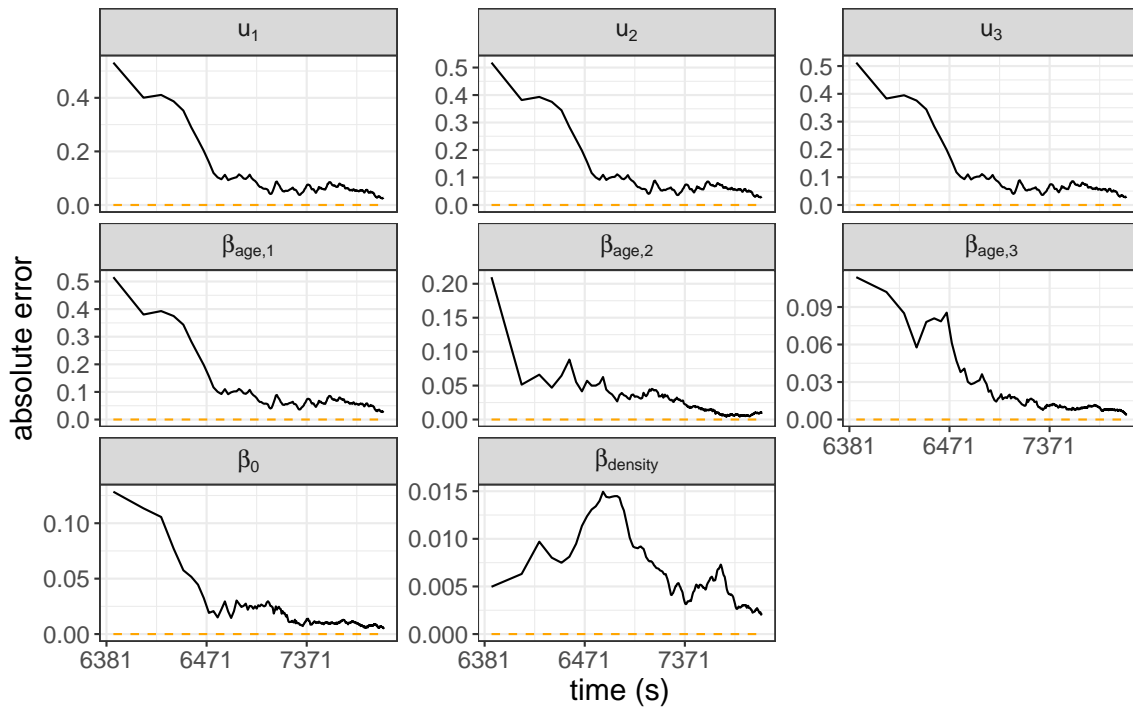


Figure 3: *Error of MCMC estimates via Stan as a function of run time. As a benchmark, we use the estimate returned by Algorithm 1. The horizontal orange line is the error of our method, which took 7 seconds to run. MCMC in Stan required over 6,000 seconds of warmup before error can be measured. Total run time for Stan was  $\sim 12,000$  seconds.*

uncertainty and (ii) it corrects sampling biases using a poststratification step. A more careful quantification of uncertainty would use posterior intervals, rather than posterior variance. Such an interval can be estimated using MCMC draws. Extending our numerical scheme into a sampling scheme to estimate such intervals is a direction we are actively pursuing.

For the model of this paper we only used a fraction of the available covariates, that is, the data collected in survey responses. As a result, the model can be extended to include more than two groups. Estimates tend to be noisy because the studied covariates can be strongly correlated with the outcome. For example, age is correlated with intensity of symptoms. The marginal correlation however is weak. This, and other considerations, suggest that it might be beneficial from a modeling standpoint to build a more sophisticated model, which might be outside of the scope of application of the methods of this paper. Nevertheless, the model considered here is an important step in the development of a better model.

## 11 Application: Public opinion on abortion policies

We next apply our method to a hierarchical linear regression used to model attitudes on abortion policies as they vary across states, ethnicity, age groups, and education levels. Modeling this heterogeneity requires partitioning an initially large data set into small groups. Furthermore, we must address biases that can arise in our survey and correct them using more comprehensive surveys, such as census data. As in Section 10, we use MRP to do inference for small slices of big data and correct biases in our survey.

We analyze data from the 2018 Cooperative Congressional Election Study (CCES) using, as in the case study of [Lopez-Martin et al., 2020], a random subset of 5,000 respondents. Respondents express support or opposition on six abortion policies, for example “Ban abortion after the 20<sup>th</sup> week of pregnancy” or “Allow employers to decline coverage of abortion in insurance plan.” These policies are intended to restrict access to abortion. Each respondent is given a support score,  $y$ , ranging from 0 to 6, indicating the number of supported policies.

We use a normal likelihood with the following covariates, recorded for each respondent: state, ethnicity, age group, education level, and sex. We use the proportion of votes for the Republican party in the state in 2016 as an additional predictor, denoted as `repvote`. The model also admits an intercept term. The statistical formulation of the model is the following:

$$y_i \sim \text{normal}(\beta_0 + X_i^{\text{state}}\beta_{\text{state}} + X_i^{\text{ethnicity}}\beta_{\text{ethnicity}} + X_i^{\text{age}}\beta_{\text{age}} + X_i^{\text{sex}}\beta_{\text{sex}} + X_i^{\text{education}}\beta_{\text{education}} + X_i^{\text{repvote}}\beta_{\text{repvote}}, \sigma_1)$$

The difficult parameters to estimate here are the state coefficients, to which we give normal( $0, \sigma_2$ ) priors. Because the model includes `repvote`, the partial pooling is done toward the prediction of the state based on its previous vote, not toward the national mean.

Table 2 summarizes the performance of the algorithm on this model. The posterior mean and standard deviation of the MRP estimates for each state can be computed as in Section 10 and are plotted in Figure 4.

Figure 4 shows that the expected support score increases with the level of support for the Republican party, bearing some fluctuations. The large posterior standard deviations

$n$	$k_1$	$k_2$	max error	total time (s)
5000	50	19	$1.2 \times 10^{-8}$	0.05

Table 2: *Computation time and accuracy of Algorithm 1 applied to a model of support/opposition for abortion policies. The column “max error” shows the maximum error of posterior means and standard deviations of regression coefficients and scale parameters.*

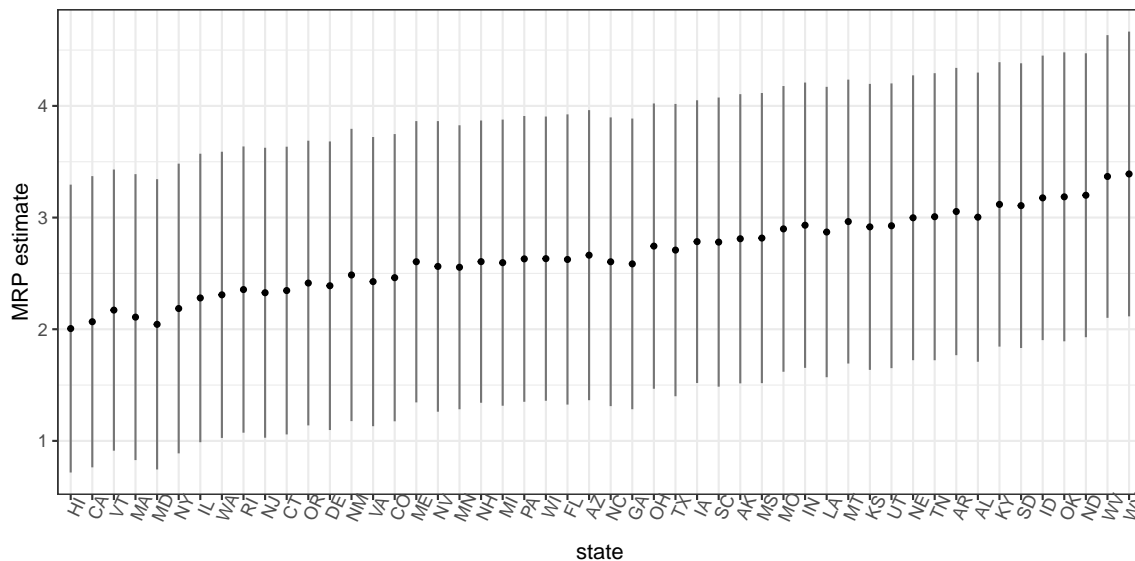


Figure 4: *MRP estimate of the expected level of support for anti-abortion policies in each state. The point represents the posterior mean, and the bars span two posterior standard deviations. The states are ordered based on Republican vote share in the 2016 presidential election.*

indicate there is quite a bit of heterogeneity within each state. For further insight, we may examine how groups other than states, e.g. ethnic groups, age groups, etc. behave.

The present model has certain limitations. First, one could consider interaction terms. This seems sensible since, for instance, white males with no college education likely behave differently than white males with a college degree. The numerical method presented in this paper can handle interaction terms. Computing the posterior standard deviation of the MRP estimate however requires some data wrangling. We plan to create an R package with routines that seamlessly implement these MRP calculations, making it straightforward for modelers to experiment with different covariates and interaction terms.

There is also interest in nonlinear models with non-normal likelihoods. Lopez-Martin et al. [2020] consider an item-response or ideal-point logistic regression. This sort of model can better capture certain characteristics of the data, such as dependence among different survey responses. For such models, we cannot use the proposed integration scheme. This presents us with a tradeoff: the proposed algorithm takes a fraction of a second to run, while fitting the ideal point model with Stan’s MCMC takes hundreds of seconds. The difference is more severe if, rather than fitting a subset of 5,000 respondents, we use all 60,000 respondents in the survey. The modeler then needs to assess how useful it is to use



a non-normal likelihood. Even then, the normal likelihood model can be a fast way to do model exploration, by for example examining various covariates and interaction terms.

## 12 Conclusions and generalizations

In this paper we describe a class of fast algorithms for evaluating the posterior moments of two Bayesian linear regression models:

1. Two-group normal-normal: The two-group normal-normal model is used to model a continuous outcome with two groups of parameters:

$$\begin{aligned} y &\sim \text{normal}(X_1\beta_1 + X_2\beta_2, \sigma_3) \\ \beta_1 &\sim \text{normal}(0, \sigma_1) \\ \beta_2 &\sim \text{normal}(0, \sigma_2) \end{aligned} \tag{21}$$

where  $X_1$  is a  $n \times k_1$  matrix of predictors,  $\beta_1 \in \mathbb{R}^{k_1}$  is a vector of regression coefficients,  $X_2$  is a  $n \times k_2$  matrix of predictors, and  $\beta_2 \in \mathbb{R}^{k_2}$  is a vector of regression coefficients.

2. Mixed effects model: The mixed-effects model is a slight variant of the two-group normal-normal model. In the mixed-effects model we model the scale parameter on one group of coefficients and assign fixed scale parameters to the priors on all other coefficients:

$$\begin{aligned} y &\sim \text{normal}(X_1\beta_1 + X_2\beta_2, \sigma_3) \\ \beta_1 &\sim \text{normal}(0, \sigma_1) \\ \beta_{2,i} &\sim \text{normal}(0, \sigma_{2,i}) \end{aligned} \tag{22}$$

where  $\sigma_{2,i}$  is the fixed scale parameter prior on each regression coefficient  $\beta_{2,i}$  for  $i = 1, \dots, k_2$  where  $\beta_2 \in \mathbb{R}^{k_2}$ .

The algorithms of this paper allow for assigning a general choice of priors on the scale parameters. We demonstrated the performance of our algorithm for posterior inference on two applications. In Section 10 we used COVID-19 symptom survey data to model geographic and age effects. We also used the mixed-effects model with public opinion survey data to estimate geographic and demographic impacts on attitudes towards abortion. These are both existing applications that have been fit with MCMC; by allowing these models to be fit much faster, our algorithm can facilitate a workflow in which users can fit and explore many more models in real time.

The algorithms of this paper provide substantial improvements over standard MCMC methods in both computation time and accuracy in approximating posterior moments. These improvements rely on analytically integrating the regression coefficients, which make up the bulk of the posterior dimensions, and then numerically integrating the remaining low-dimensional density with Gaussian quadrature.

Many of the techniques and analysis used in this paper generalize to multilevel and multi-group models with more than two-groups. For an  $m$  group model, the numerical integration

of our algorithm is computed over a  $m + 1$  dimensional density. For models with large  $m$  (large number of groups) the analytic marginalization of this paper can still be applied, however, integration via a tensor product of Gaussian nodes will not be feasible. On the other hand, using MCMC or other integration schemes can be used on the  $m + 1$  dimensional marginal density.

Bayesian models with more than two groups and non-Gaussian likelihoods are directions of future research.

## 13 Acknowledgements

The authors are grateful to Hagai Rossman and Ayya Keshet for useful discussions and their contribution to the COVID-19 model.

## A Integral with respect to $\rho$

In this section we describe analytical properties of the integrand of the inner integral of (3) that are used in the evaluation of the integral.

Let  $\psi : \mathbb{R}^{+3} \rightarrow \mathbb{R}$  be defined by the formula

$$\psi(\rho, c, n) = \frac{e^{-\frac{\rho^2}{2} - \frac{c}{2\rho^2}}}{\rho^{n-2}}. \quad (23)$$

We seek, for fixed  $c$  and  $n$ , the value for  $\rho$  that maximizes  $\psi(\rho, c, n)$ . We observe that

$$\begin{aligned} \frac{\partial \psi}{\partial \rho} &= e^{-\frac{\rho^2}{2} - \frac{c}{2\rho^2}} \left( \rho^{1-n}(2-n) + \left( \frac{c}{\rho^3} - \rho \right) \rho^{2-n} \right) \\ &= e^{-\frac{\rho^2}{2} - \frac{c}{2\rho^2}} \rho^{1-n} \left( 2-n + \frac{c}{\rho^2} - \rho^2 \right) \end{aligned} \quad (24)$$

Setting

$$2 - n + \frac{c}{\rho^2} - \rho^2 = 0 \quad (25)$$

and rearranging terms, we have

$$\rho^4 + \rho^2(n-2) - c = 0. \quad (26)$$

Then setting

$$\rho_{max} = \frac{1}{\sqrt{2}} \left( \sqrt{4c + (n-2)^2} - n + 2 \right)^{1/2} \quad (27)$$

we observe

$$\frac{\partial \psi}{\partial \rho}(\rho_{max}) = 0. \quad (28)$$

That is, for fixed  $c, n$ , we have  $\psi$  achieves its maximum at

$$\rho = \frac{1}{\sqrt{2}} \left( \sqrt{4c + (n-2)^2} - n + 2 \right)^{1/2}. \quad (29)$$

Furthermore,

$$\frac{\partial}{\partial \rho} \log(\psi(\rho, c, n)) = -\rho + \frac{c}{\rho^3} - \frac{n-2}{\rho} \quad (30)$$

and

$$\frac{\partial^2}{\partial \rho^2} \log(\psi(\rho, c, n)) = -1 - \frac{3c}{\rho^4} + \frac{n-2}{\rho^2}. \quad (31)$$

## References

- R. Bardini, G. Politano, A. Benso, and S. Di Carlo. Multi-level and hybrid modelling approaches for systems biology. *Computational and Structural Biotechnology Journal*, 15: 396–402, 2017. ISSN 2001-0370. doi: <https://doi.org/10.1016/j.csbj.2017.07.005>. URL <https://www.sciencedirect.com/science/article/pii/S2001037017300314>.
- Michael Betancourt. A conceptual introduction to Hamiltonian Monte Carlo, 2018.
- Michael Betancourt, Simon Byrne, and Mark Girolami. Optimizing the integrator step size for hamiltonian monte carlo. *arXiv*, stat/1411.6669, 2015.
- David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- Bob Carpenter, Andrew Gelman, Matthew D. Hoffman, Daniel Lee, Ben Goodrich, Michael Betancourt, Marcus Brubaker, Jiqiang Guo, Peter Li, and Allen Riddell. Stan: A probabilistic programming language. *Journal of Statistical Software, Articles*, 76(1):1–32, 2017.
- Alan E. Gelfand, Susan E. Hills, Amy Racine-Poon, and Adrian F. M. Smith. Illustration of Bayesian inference in normal data models using Gibbs sampling. *Journal of the American Statistical Association*, 85(412):972–985, 1990. ISSN 01621459. URL <http://www.jstor.org/stable/2289594>.
- A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian Data Analysis*. Chapman and Hall/CRC, New York, NY, 3rd edition, 2013.
- Andrew Gelman and Jennifer Hill. *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Analytical Methods for Social Research. Cambridge University Press, 2006. doi: 10.1017/CBO9780511790942.
- Andrew Gelman, Aki Vehtari, Daniel Simpson, Charles C. Margossian, Bob Carpenter, Yuling Yao, Lauren Kennedy, Jonah Gabry, Paul-Christian Bürkner, and Martin Modrák. Bayesian workflow. *arXiv*, stat/2011.01808, 2020.
- P. Greengard, A. Gelman, and A. Vehtari. A Fast Linear Regression via SVD and Marginalization. *Computational Statistics*, 2021. doi: 10.1007/s00180-021-01135-x. URL <https://doi.org/10.1007/s00180-021-01135-x>.
- Sander Greenland. Principles of multilevel modelling. *International Journal of Epidemiology*, 29(1):158–167, 02 2000. doi: 10.1093/ije/29.1.158. URL <https://doi.org/10.1093/ije/29.1.158>.

- Matthew D. Hoffman and Andrew Gelman. The No-U-Turn sampler: Adaptively setting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning Research*, 15(47):1593–1623, 2014. URL <http://jmlr.org/papers/v15/hoffman14a.html>.
- Kasper Kristensen, Anders Nielsen, Casper W. Berg, Hans Skaug, and Bradley M. Bell. Tmb: Automatic differentiation and laplace approximation. *Journal of Statistical Software, Articles*, 70(5):1–21, 2016. doi: 10.18637/jss.v070.i05. URL <https://www.jstatsoft.org/v070/i05>.
- D.V. Lindley and A.F.M. Smith. Bayes estimates for the linear model. *Journal of the Royal Statistical Society. Series B (Methodological)*, 34(1):1–41, 1972. ISSN 00359246. URL <http://www.jstor.org/stable/2985048>.
- Juan Lopez-Martin, Justin H. Phillips, and Andrew Gelman. Multilevel regression and poststratification case studies. 2020. URL <https://juanlopezmartin.github.io/>.
- Charles C. Margossian, Aki Vehtari, Daniel Simpson, and Raj Agrawal. Hamiltonian Monte Carlo using an adjoint-differentiated Laplace approximation: Bayesian inference for latent Gaussian models and beyond. In *Advances in Neural Information Processing Systems*, 2020.
- Juan Merlo, Basile Chaix, Min Yang, John Lynch, and Lennart Rastam. A brief conceptual tutorial of multilevel analysis in social epidemiology: linking the statistical concept of clustering to the idea of contextual phenomenon. *Journal of epidemiology and community health*, 59(6):443–449, 2005. URL <https://doi.org/10.1136/jech.2004.023473>.
- Hagai Rossman, Ayya Keshet, Smadar Shilo, Amir Gavrieli, Tal Bauman, Ori Cohen, Esti Shelly, Ran Balicer, Benjamin Geiger, Yuval Dor, and Eran Segal. A framework for identifying regional outbreak and spread of COVID-19 from one-minute population-wide surveys. *Nature Medicine*, 26(5):634 – 638, 2020. doi: 10.1038/s41591-020-0857-9.
- Håvard Rue, Andrea Riebler, Sigrunn H. Sørbye, Janine B. Illian, Daniel P. Simpson, and Finn K. Lindgren. Bayesian computing with INLA: A review. *Annual Review of Statistics and Its Application*, 4(1):395–421, 2017. doi: 10.1146/annurev-statistics-060116-054045. URL <https://doi.org/10.1146/annurev-statistics-060116-054045>.
- Eran Segal, Feng Zhang, Xihong Lin, Gary King, Ophir Shalem, Smadar Shilo, William E. Allen, Faisal Alquaddoomi, Han Altae-Tran, Simon Anders, Ran Balicer, Tal Bauman, Ximena Bonilla, Gisel Booman, Andrew T. Chan, Ori Cohen, Silvano Coletti, Natalie Davidson, Yuval Dor, David A. Drew, Olivier Elemento, Georgina Evans, Phil Ewels, Joshua Gale, Amir Gavrieli, Benjamin Geiger, Yonatan H. Grad, Casey S. Greene, Iman Hajirasouliha, Roman Jerala, Andre Kahles, Olli Kallioniemi, Ayya Keshet, Ljupco Kocarev, Gregory Landua, Tomer Meir, Aline Muller, Long H. Nguyen, Matej Oresic, Svetlana Ovchinnikova, Hedi Peterson, Jana Prodanova, Jay Rajagopal, Gunnar Rättsch, Hagai Rossman, Johan Rung, Andrea Sboner, Alexandros Sigaras, Tim Spector, Ron Steinherz, Irene Stevens, Jaak Vilo, and Paul Wilmes. Building an international consortium for tracking coronavirus health status. 26(8):1161–1165, 2020. doi: 10.1038/s41591-020-0929-x.

L. N. Trefethen. *Approximation Theory and Approximation Practice: Extended Edition*. SIAM, Philadelphia, PA, 2020.

L. N. Trefethen and D. Bau. *Numerical Linear Algebra*. SIAM, New York, NY, 1997.