

# Modeling and Analysis for Categorical Response Data

*Siddhartha Chib*

## 1. Introduction

In this chapter we discuss how Bayesian methods are used to model and analyze categorical response data. As in other areas of statistics, the growth of Bayesian ideas in the categorical data setting has been rapid, aided by developments in *Markov chain Monte Carlo* (MCMC) methods (Gelfand and Smith, 1990; Smith and Roberts, 1993; Tanner and Wong, 1987; Tierney, 1994; Chib and Greenberg, 1995), and by the framework for fitting and criticizing categorical data models developed in Albert and Chib (1993). In this largely self-contained chapter we summarize the various possibilities for the Bayesian modeling of categorical response data and provide the associated inferential techniques for conducting the prior–posterior analyses.

The rest of the chapter is organized as follows. We begin by including a brief overview of MCMC methods. We then discuss how the output of the MCMC simulations can be used to calculate the marginal likelihood of the model for the purpose of comparing models via Bayes factors. After these preliminaries, the chapter turns to the analysis of categorical response models in the cross-section setting, followed by extensions to multivariate and longitudinal responses.

### 1.1. Elements of Markov chain Monte Carlo

Suppose that  $\boldsymbol{\psi} \in \mathfrak{R}^d$  denotes a vector of random variables of interest (typically consisting of a set of parameters  $\boldsymbol{\theta}$  and other variables) in a particular Bayesian model and let  $\pi(\boldsymbol{\psi}|\mathbf{y}) \propto \pi(\boldsymbol{\psi})p(\mathbf{y}|\boldsymbol{\psi})$  denote the posterior density, where  $\pi(\boldsymbol{\psi})$  is the prior density and  $p(\mathbf{y}|\boldsymbol{\psi})$  is the sampling density. Suppose that we are interested in calculating the posterior mean  $\boldsymbol{\eta} = \int_{\mathfrak{R}^d} \boldsymbol{\psi} \pi(\boldsymbol{\psi}|\mathbf{y}) d\boldsymbol{\psi}$  and that this integral cannot be computed analytically, or numerically, because the dimension of the integration precludes the use of quadrature-based methods. To tackle this problem we can rely on Monte Carlo sampling methods. Instead of concentrating on the computation of the above integral we can proceed in a more general way. We construct a method to produce a sequence of draws from the posterior density, namely

$$\boldsymbol{\psi}^{(1)}, \dots, \boldsymbol{\psi}^{(M)} \sim \pi(\boldsymbol{\psi}|\mathbf{y}).$$

Provided  $M$  is large enough, this sample from the posterior density can be used to summarize the posterior density. We can estimate the mean by taking the average of these simulated draws. We can estimate the quantiles of the posterior density from the quantiles of the sampled draws. Other summaries of the posterior density can be obtained in a similar manner. Under suitable laws of large numbers these estimates converge to the posterior quantities as the simulation-size becomes large.

The sampling of the posterior distribution is, therefore, the central concern in Bayesian computation. One important breakthrough in the use of simulation methods was the realization that the sampled draws need not be independent, that simulation-consistency can be achieved with correlated draws. The fact that the sampled variates can be correlated is of immense practical and theoretical importance and is the defining characteristic of Markov chain Monte Carlo methods, popularly referred to by the acronym MCMC, where the sampled draws form a Markov chain. The idea behind these methods is simple and extremely general. In order to sample a given probability distribution, referred to as the target distribution, a suitable Markov chain is constructed with the property that its limiting, invariant distribution, is the target distribution. Once the Markov chain has been constructed, a sample of draws from the target distribution is obtained by simulating the Markov chain a large number of times and recording its values. Within the Bayesian framework, where both parameters and data are treated as random variables, and inferences about the parameters are conducted conditioned on the data, the posterior distribution of the parameters provides a natural target for MCMC methods.

Markov chain sampling methods originate with the work of [Metropolis et al. \(1953\)](#) in statistical physics. A vital extension of the method was made by [Hastings \(1970\)](#) leading to a method that is now called the Metropolis–Hastings algorithm (see [Tierney, 1994](#) and [Chib and Greenberg, 1995](#) for detailed summaries). This algorithm was first applied to problems in spatial statistics and image analysis ([Besag, 1974](#)). A resurgence of interest in MCMC methods started with the papers of [Geman and Geman \(1984\)](#) who developed an algorithm, a special case of the Metropolis method that later came to be called the Gibbs sampler, to sample a discrete distribution, [Tanner and Wong \(1987\)](#) who proposed a MCMC scheme involving data augmentation to sample posterior distributions in missing data problems, and [Gelfand and Smith \(1990\)](#) where the value of the Gibbs sampler was demonstrated for general Bayesian problems with continuous parameter spaces.

The Gibbs sampling algorithm is one of the simplest Markov chain Monte Carlo algorithms and is easy to describe. Suppose that  $\psi_1$  and  $\psi_2$  denote some grouping (blocking) of  $\psi$  and let

$$\pi_1(\psi_1|y, \psi_2) \propto p(y|\psi_1, \psi_2)\pi(\psi_1, \psi_2), \quad (1.1)$$

$$\pi_2(\psi_2|y, \psi_1) \propto p(y|\psi_1, \psi_2)\pi(\psi_1, \psi_2) \quad (1.2)$$

denote the associated conditional densities, often called the full conditional densities. Then, one cycle of the Gibbs sampling algorithm is completed by sampling each of the full conditional densities, using the most current values of the conditioning block. The Gibbs sampler in which each block is revised in fixed order is defined as follows.

ALGORITHM (*Gibbs Sampling*).

- 1 Specify an initial value  $\boldsymbol{\psi}^{(0)} = (\boldsymbol{\psi}_1^{(0)}, \boldsymbol{\psi}_2^{(0)})$ :
- 2 Repeat for  $j = 1, 2, \dots, n_0 + G$ .
  - Generate  $\boldsymbol{\psi}_1^{(j)}$  from  $\pi_1(\boldsymbol{\psi}_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)})$
  - Generate  $\boldsymbol{\psi}_2^{(j)}$  from  $\pi_2(\boldsymbol{\psi}_2 | \mathbf{y}, \boldsymbol{\psi}_1^{(j)})$
- 3 Return the values  $\{\boldsymbol{\psi}^{(n_0+1)}, \boldsymbol{\psi}^{(n_0+2)}, \dots, \boldsymbol{\psi}^{(n_0+G)}\}$ .

In some problems it turns out that the full conditional density cannot be sampled directly. In such cases, the intractable full conditional density is sampled via the Metropolis–Hastings (M–H) algorithm. For specificity, suppose that the full conditional density  $\pi(\boldsymbol{\psi}_1 | \mathbf{y}, \boldsymbol{\psi}_2)$  is intractable. Let

$$q_1(\boldsymbol{\psi}_1, \boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2)$$

denote a proposal density that generates a candidate  $\boldsymbol{\psi}'_1$ , given the data and the values of the remaining blocks. Then, in the first step of the  $j$ th iteration of the MCMC algorithm, given the values  $(\boldsymbol{\psi}_1^{(j-1)}, \boldsymbol{\psi}_2^{(j-1)})$  from the previous iteration, the transition to the next iterate of  $\boldsymbol{\psi}_1$  is achieved as follows.

ALGORITHM (*Metropolis–Hastings step for sampling an intractable  $\pi_1(\boldsymbol{\psi}_1 | \mathbf{y}, \boldsymbol{\psi}_2)$* ).

- 1 Propose a value for  $\boldsymbol{\psi}_1$  by drawing:

$$\boldsymbol{\psi}'_1 \sim q_1(\boldsymbol{\psi}_1^{(j-1)}, \cdot | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)})$$

- 2 Calculate the probability of move  $\alpha(\boldsymbol{\psi}_1^{(j-1)}, \boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)})$  given by

$$\min \left\{ 1, \frac{\pi(\boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)}) q_1(\boldsymbol{\psi}'_1, \boldsymbol{\psi}_1^{(j-1)} | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)})}{\pi(\boldsymbol{\psi}_1^{(j-1)} | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)}) q_1(\boldsymbol{\psi}_1^{(j-1)}, \boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)})} \right\}.$$

- 3 Set

$$\boldsymbol{\psi}_1^{(j)} = \begin{cases} \boldsymbol{\psi}'_1 & \text{with prob } \alpha(\boldsymbol{\psi}_1^{(j-1)}, \boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)}), \\ \boldsymbol{\psi}_1^{(j-1)} & \text{with prob } 1 - \alpha(\boldsymbol{\psi}_1^{(j-1)}, \boldsymbol{\psi}'_1 | \mathbf{y}, \boldsymbol{\psi}_2^{(j-1)}). \end{cases}$$

A similar approach is used to sample  $\boldsymbol{\psi}_2$  if the full conditional density of  $\boldsymbol{\psi}_2$  is intractable. These algorithms extend straightforwardly to problems with more than two blocks.

### 1.2. Computation of the marginal likelihood

Posterior simulation by MCMC methods does not require knowledge of the normalizing constant of the posterior density. Nonetheless, if we are interested in comparing alter-

native models, then knowledge of the normalizing constant is essential. This is because the formal Bayesian approach for comparing models is via *Bayes factors*, or ratios of *marginal likelihoods*. The marginal likelihood of a particular model is the normalizing constant of the posterior density and is defined as

$$m(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta})\pi(\boldsymbol{\theta}) \, d\boldsymbol{\theta}, \quad (1.3)$$

the integral of the likelihood function with respect to the prior density. If we have two models  $\mathcal{M}_k$  and  $\mathcal{M}_l$ , then the Bayes factor is the ratio

$$B_{kl} = \frac{m(\mathbf{y}|\mathcal{M}_k)}{m(\mathbf{y}|\mathcal{M}_l)}. \quad (1.4)$$

Computation of the marginal likelihood is, therefore, of some importance in Bayesian statistics (DiCiccio et al., 1997, Chen and Shao, 1998, Roberts, 2001). It is possible to avoid the direct computation of the marginal likelihood by the methods of Green (1995) and Carlin and Chib (1995), but this requires formulating a more general MCMC scheme in which parameters and models are sampled jointly. We do not consider these approaches in this chapter. Unfortunately, because MCMC methods deliver draws from the posterior density, and the marginal likelihood is the integral with respect to the prior, the MCMC output cannot be used directly to average the likelihood. To deal with this problem, a number of methods have appeared in the literature. One simple and widely applicable method is due to Chib (1995) which we briefly explain as follows.

Begin by noting that  $m(\mathbf{y})$  by virtue of being the normalizing constant of the posterior density can be expressed as

$$m(\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta}^*)\pi(\boldsymbol{\theta}^*)}{\pi(\boldsymbol{\theta}^*|\mathbf{y})}, \quad (1.5)$$

for any given point  $\boldsymbol{\theta}^*$  (generally taken to be a high density point such as the posterior mode or mean). Thus, provided we have an estimate  $\hat{\pi}(\boldsymbol{\theta}^*|\mathbf{y})$  of the posterior ordinate, the marginal likelihood can be estimated on the log scale as

$$\log m(\mathbf{y}) = \log p(\mathbf{y}|\boldsymbol{\theta}^*) + \log \pi(\boldsymbol{\theta}^*) - \log \hat{\pi}(\boldsymbol{\theta}^*|\mathbf{y}). \quad (1.6)$$

In the context of both single and multiple block MCMC chains, good estimates of the posterior ordinate are available. For example, when the MCMC simulation is run with  $B$  blocks, Chib (1995) employs the marginal–conditional decomposition

$$\pi(\boldsymbol{\theta}^*|\mathbf{y}) = \pi(\boldsymbol{\theta}_1^*|\mathbf{y}) \times \cdots \times \pi(\boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*) \times \cdots \times \pi(\boldsymbol{\theta}_B^*|\mathbf{y}, \boldsymbol{\Theta}_{B-1}^*), \quad (1.7)$$

where the typical term is of the form

$$\begin{aligned} &\pi(\boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*) \\ &= \int \pi(\boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\theta}^{i+1}, \mathbf{z})\pi(\boldsymbol{\theta}^{i+1}, \mathbf{z}|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*) \, d\boldsymbol{\theta}^{i+1} \, d\mathbf{z}. \end{aligned}$$

In the latter expression  $\boldsymbol{\Theta}_i = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_i)$  and  $\boldsymbol{\Theta}^i = (\boldsymbol{\theta}_i, \dots, \boldsymbol{\theta}_B)$  denote the list of blocks up to  $i$  and the set of blocks from  $i$  to  $B$ , respectively, and  $\mathbf{z}$  denotes any latent

data that is included in the sampling. This is the *reduced conditional ordinate*. It is important to bear in mind that in finding the reduced conditional ordinate one must integrate only over  $(\boldsymbol{\Theta}^{i+1}, \mathbf{z})$  and that the integrating measure is conditioned on  $\boldsymbol{\Theta}_{i-1}^*$ .

Consider first the case where the normalizing constant of each full conditional density is known. Then, the first term of (1.7) can be estimated by the Rao–Blackwell method. To estimate the typical reduced conditional ordinate, one conducts a MCMC run consisting of the full conditional distributions

$$\begin{aligned} & \{ \pi(\boldsymbol{\theta}_i | \mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1}, \mathbf{z}); \dots; \pi(\boldsymbol{\theta}_B | \mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\theta}_i, \dots, \boldsymbol{\theta}_{B-1}, \mathbf{z}); \\ & \pi(\mathbf{z} | \mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^i) \}, \end{aligned} \tag{1.8}$$

where the blocks in  $\boldsymbol{\Theta}_{i-1}$  are set equal to  $\boldsymbol{\Theta}_{i-1}^*$ . By MCMC theory, the draws on  $(\boldsymbol{\Theta}^{i+1}, \mathbf{z})$  from this run are from the distribution  $\pi(\boldsymbol{\Theta}^{i+1}, \mathbf{z} | \mathbf{y}, \boldsymbol{\Theta}_{i-1}^*)$  and so the reduced conditional ordinate can be estimated as the average

$$\hat{\pi}(\boldsymbol{\theta}_i^* | \mathbf{y}, \mathcal{M}, \boldsymbol{\Theta}_{i-1}^*) = M^{-1} \sum_{j=1}^M \pi(\boldsymbol{\theta}_i^* | \mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1,(j)}, \mathbf{z}^{(j)})$$

over the simulated values of  $\boldsymbol{\Theta}^{i+1}$  and  $\mathbf{z}$  from the reduced run. Each subsequent reduced conditional ordinate that appears in the decomposition (1.7) is estimated in the same way though, conveniently, with fewer and fewer distributions appearing in the reduced runs. Given the marginal and reduced conditional ordinates, the marginal likelihood on the log scale is available as

$$\begin{aligned} & \log \hat{m}(\mathbf{y} | \mathcal{M}) \\ & = \log p(\mathbf{y} | \boldsymbol{\theta}^*, \mathcal{M}) + \log \pi(\boldsymbol{\theta}^*) - \sum_{i=1}^B \log \hat{\pi}(\boldsymbol{\theta}_i^* | \mathbf{y}, \mathcal{M}, \boldsymbol{\Theta}_{i-1}^*), \end{aligned} \tag{1.9}$$

where  $p(\mathbf{y} | \boldsymbol{\theta}^*)$  is the density of the data marginalized over the latent data  $\mathbf{z}$ .

Consider next the case where the normalizing constant of one or more of the full conditional densities is not known. In that case, the posterior ordinate is estimated by a modified method developed by Chib and Jeliazkov (2001). If sampling is conducted in one block by the M–H algorithm with proposal density  $q(\boldsymbol{\theta}, \boldsymbol{\theta}' | \mathbf{y})$  and probability of move

$$\alpha(\boldsymbol{\theta}^{(j-1)}, \boldsymbol{\theta}' | \mathbf{y}) = \min \left\{ 1, \frac{\pi(\boldsymbol{\theta}' | \mathbf{y}) q(\boldsymbol{\theta}', \boldsymbol{\theta}^{(j-1)} | \mathbf{y})}{\pi(\boldsymbol{\theta}^{(j-1)} | \mathbf{y}) q(\boldsymbol{\theta}^{(j-1)}, \boldsymbol{\theta}' | \mathbf{y})} \right\}$$

then it can be shown that the posterior ordinate is given by

$$\pi(\boldsymbol{\theta}^* | \mathbf{y}) = \frac{E_1 \{ \alpha(\boldsymbol{\theta}, \boldsymbol{\theta}^* | \mathbf{y}) q(\boldsymbol{\theta}, \boldsymbol{\theta}^* | \mathbf{y}) \}}{E_2 \{ \alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta} | \mathbf{y}) \}},$$

where the numerator expectation  $E_1$  is with respect to the distribution  $\pi(\boldsymbol{\theta} | \mathbf{y})$  and the denominator expectation  $E_2$  is with respect to the proposal density  $q(\boldsymbol{\theta}^*, \boldsymbol{\theta} | \mathbf{y})$ . This

leads to the simulation consistent estimate

$$\hat{\pi}(\boldsymbol{\theta}^*|\mathbf{y}) = \frac{M^{-1} \sum_{g=1}^M \alpha(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^*|\mathbf{y}) q(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^*|\mathbf{y})}{J^{-1} \sum_{j=1}^M \alpha(\boldsymbol{\theta}^*, \boldsymbol{\theta}^{(j)}|\mathbf{y})}, \quad (1.10)$$

where  $\boldsymbol{\theta}^{(g)}$  are the given draws from the posterior distribution while the draws  $\boldsymbol{\theta}^{(j)}$  in the denominator are from  $q(\boldsymbol{\theta}^*, \boldsymbol{\theta}|\mathbf{y})$ , given the fixed value  $\boldsymbol{\theta}^*$ .

In general, when sampling is done with  $B$  blocks, the typical reduced conditional ordinate is given by

$$\begin{aligned} \pi(\boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*) \\ = \frac{E_1\{\alpha(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1}, \mathbf{z}) q_i(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1}, \mathbf{z})\}}{E_2\{\alpha(\boldsymbol{\theta}_i^*, \boldsymbol{\theta}_i|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1}, \mathbf{z})\}}, \end{aligned} \quad (1.11)$$

where  $E_1$  is the expectation with respect to  $\pi(\boldsymbol{\Theta}^i, \mathbf{z}|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*)$  and  $E_2$  that with respect to the product measure  $\pi(\boldsymbol{\Theta}^{i+1}, \mathbf{z}|\mathbf{y}, \boldsymbol{\theta}_i^*) q_i(\boldsymbol{\theta}_i^*, \boldsymbol{\theta}_i|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1})$ . The quantity  $\alpha(\boldsymbol{\theta}_i, \boldsymbol{\theta}_i^*|\mathbf{y}, \boldsymbol{\Theta}_{i-1}^*, \boldsymbol{\Theta}^{i+1}, \mathbf{z})$  is the M–H probability of move. The two expectations are estimated from the output of the reduced runs in an obvious way.

## 2. Binary responses

Suppose that  $y_i$  is a binary  $\{0, 1\}$  response variable and  $\mathbf{x}_i$  is a  $k$  vector of covariates. Suppose that we have a random sample of  $n$  observations and the Bayesian model of interest is

$$\begin{aligned} y_i|\boldsymbol{\beta} &\sim \Phi(\mathbf{x}_i'\boldsymbol{\beta}), \quad i \leq n, \\ \boldsymbol{\beta} &\sim \mathcal{N}_k(\boldsymbol{\beta}_0, \mathbf{B}_0), \end{aligned}$$

where  $\Phi$  is the cdf of the  $N(0, 1)$  distribution. We consider the case of the logit link below. On letting  $p_i = \Phi(\mathbf{x}_i'\boldsymbol{\beta})$ , the posterior distribution of  $\boldsymbol{\beta}$  is proportional to

$$\pi(\boldsymbol{\beta}|\mathbf{y}) \propto \pi(\boldsymbol{\beta}) \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{(1-y_i)},$$

which does not belong to any named family of distributions. To deal with this model, and others involving categorical data, [Albert and Chib \(1993\)](#) introduce a technique that has led to many applications. The Albert–Chib algorithm capitalizes on the simplifications afforded by introducing latent or auxiliary data into the sampling.

Instead of the specification above, the model is re-specified as

$$\begin{aligned} z_i|\boldsymbol{\beta} &\sim \mathcal{N}(\mathbf{x}_i'\boldsymbol{\beta}, 1), \\ y_i &= I[z_i > 0], \quad i \leq n, \\ \boldsymbol{\beta} &\sim \mathcal{N}_k(\boldsymbol{\beta}_0, \mathbf{B}_0). \end{aligned} \quad (2.1)$$

This specification is equivalent to the binary probit regression model since

$$\Pr(y_i = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \Pr(z_i > 0 | \mathbf{x}_i, \boldsymbol{\beta}) = \Phi(\mathbf{x}'_i \boldsymbol{\beta}),$$

as required.

Albert and Chib (1993) exploit this equivalence and propose that the latent variables  $\mathbf{z} = \{z_1, \dots, z_n\}$ , one for each observation, be included in the MCMC algorithm along with the regression parameter  $\boldsymbol{\beta}$ . In other words, they suggest using MCMC methods to sample the joint posterior distribution

$$\pi(\boldsymbol{\beta}, \mathbf{z} | \mathbf{y}) \propto \pi(\boldsymbol{\beta}) \prod_{i=1}^N \mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \Pr(y_i | z_i, \boldsymbol{\beta}).$$

The term  $\Pr(y_i | z_i, \boldsymbol{\beta})$  is obtained by reasoning as follows:  $\Pr(y_i = 0 | z_i, \boldsymbol{\beta})$  is one if  $z_i < 0$  (regardless of the value taken by  $\boldsymbol{\beta}$ ) and zero otherwise – but this is the definition of  $I(z_i < 0)$ . Similarly,  $\Pr(y_i = 1 | z_i, \boldsymbol{\beta})$  is one if  $z_i > 0$ , and zero otherwise – but this is the definition of  $I(z_i > 0)$ . Collecting these two cases we have therefore that

$$\pi(\mathbf{z}, \boldsymbol{\beta} | \mathbf{y}) \propto \pi(\boldsymbol{\beta}) \prod_{i=1}^N \mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \{I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i}\}.$$

The latter posterior density is now sampled by a two-block Gibbs sampler composed of the full conditional distributions

$$\mathbf{z} | \mathbf{y}, \boldsymbol{\beta}; \boldsymbol{\beta} | \mathbf{y}, \mathbf{z}.$$

Even though the parameter space has been enlarged, the introduction of the latent variables simplifies the problem considerably. The second conditional distribution, i.e.,  $\boldsymbol{\beta} | \mathbf{y}, \mathbf{z}$ , is the same as the distribution  $\boldsymbol{\beta} | \mathbf{z}$  since  $\mathbf{z}$  implies  $\mathbf{y}$ , and hence  $\mathbf{y}$  has no additional information for  $\boldsymbol{\beta}$ . The distribution  $\boldsymbol{\beta} | \mathbf{z}$  is easy to derive by standard Bayesian results for a continuous response. The first conditional distribution, i.e.,  $\mathbf{z} | \mathbf{y}, \boldsymbol{\beta}$ , factors into  $n$  distributions  $z_i | y_i, \boldsymbol{\beta}$  and from the above is seen to be truncated normal given the value of  $y_i$ . Specifically,

$$p(z_i | y_i, \boldsymbol{\beta}) \propto \mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \{I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i}\}$$

which implies that if  $y_i = 0$ , the posterior density of  $z_i$  is proportional to  $\mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) I[z < 0]$ , a truncated normal distribution with support  $(-\infty, 0)$ , whereas if  $y_i = 1$ , the density is proportional to  $\mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) I[z > 0]$ , a truncated normal distribution with support  $(0, \infty)$ . These truncated normal distributions are simulated by applying the *inverse cdf* method (Ripley, 1987). Specifically, it can be shown that the draw

$$\mu + \sigma \Phi^{-1} \left[ \Phi \left( \frac{a - \mu}{\sigma} \right) + U \left( \Phi \left( \frac{b - \mu}{\sigma} \right) - \Phi \left( \frac{a - \mu}{\sigma} \right) \right) \right], \tag{2.2}$$

where  $\Phi^{-1}$  is the inverse cdf of the  $\mathcal{N}(0, 1)$  distribution and  $U \sim \text{Uniform}(0, 1)$ , is from a  $\mathcal{N}(\mu, \sigma^2)$  distribution truncated to the interval  $(a, b)$ .

The Albert–Chib algorithm for the probit binary response model may now be summarized as follows.

ALGORITHM 1 (*Binary probit link model*).

1 Sample

$$z_i | y_i, \boldsymbol{\beta} \propto \mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \{ I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i} \}, \quad i \leq n$$

2 Sample

$$\boldsymbol{\beta} | \mathbf{z} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \mathbf{x}_i z_i \right), \mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}'_i \right)^{-1} \right)$$

3 Go to 1

EXAMPLE. Consider the data in Table 1, taken from Fahrmeir and Tutz (1994), which is concerned with the occurrence or nonoccurrence of infection following birth by cesarean section. The response variable  $y$  is one if the cesarean birth resulted in an infection, and zero if not. The available covariates are three indicator variables:  $x_1$  is an indicator for whether the cesarean was nonplanned;  $x_2$  is an indicator for whether risk factors were present at the time of birth and  $x_3$  is an indicator for whether antibiotics were given as a prophylaxis. The data in the table contains information from 251 births. Under the column of the response, an entry such as 11/87 means that there were 98 deliveries with covariates (1, 1, 1) of whom 11 developed an infection and 87 did not.

Let us model the binary response by a *probit* model, letting the probability of infection for the  $i$ th birth be given as

$$\Pr(y_i = 1 | \mathbf{x}_i, \boldsymbol{\beta}) = \Phi(\mathbf{x}'_i \boldsymbol{\beta}),$$

where  $\mathbf{x}_i = (1, x_{i1}, x_{i2}, x_{i3})'$  is the covariate vector,  $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \beta_3)$  is the vector of unknown coefficients and  $\Phi$  is the cdf of the standard normal random variable. Let us assume that our prior information about  $\boldsymbol{\beta}$  can be represented by a multivariate normal density with mean centered at zero for each parameter, and variance given by  $5\mathbf{I}_4$ , where  $\mathbf{I}_4$  is the four-dimensional identity matrix.

Algorithm 1 is run for 5000 cycles beyond a burn-in of 100 iterations. In Table 2 we provide the prior and posterior first two moments, and the 2.5th (lower) and 97.5th (upper) percentiles, of the marginal densities of  $\boldsymbol{\beta}$ . Each of the quantities is computed from the posterior draws in an obvious way, the posterior standard deviation in the table is

Table 1  
Cesarean infection data

| $y$ (1/0) | $x_1$ | $x_2$ | $x_3$ |
|-----------|-------|-------|-------|
| 11/87     | 1     | 1     | 1     |
| 1/17      | 0     | 1     | 1     |
| 0/2       | 0     | 0     | 1     |
| 23/3      | 1     | 1     | 0     |
| 28/30     | 0     | 1     | 0     |
| 0/9       | 1     | 0     | 0     |
| 8/32      | 0     | 0     | 0     |



Table 2  
 Cesarean data: Prior–posterior summary based on 5000 draws  
 (beyond a burn-in of 100 cycles) from the Albert–Chib algorithm

|           | Prior |         | Posterior |         |        |        |
|-----------|-------|---------|-----------|---------|--------|--------|
|           | Mean  | Std dev | Mean      | Std dev | Lower  | Upper  |
| $\beta_0$ | 0.000 | 3.162   | -1.100    | 0.210   | -1.523 | -0.698 |
| $\beta_1$ | 0.000 | 3.162   | 0.609     | 0.249   | 0.126  | 1.096  |
| $\beta_2$ | 0.000 | 3.162   | 1.202     | 0.250   | 0.712  | 1.703  |
| $\beta_3$ | 0.000 | 3.162   | -1.903    | 0.266   | -2.427 | -1.393 |

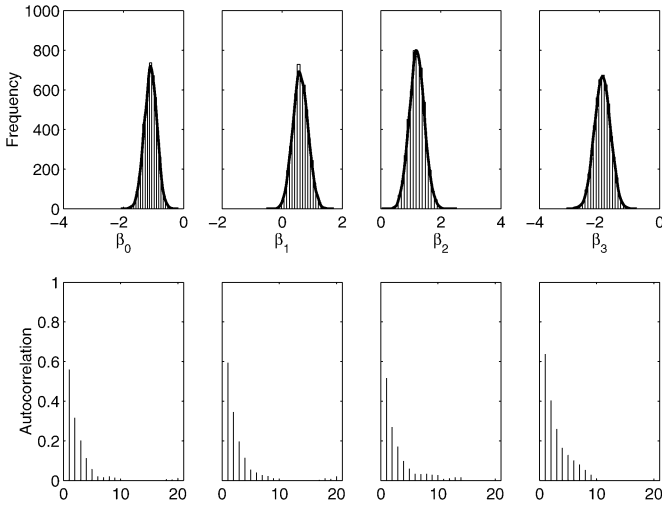


Fig. 1. Cesarean data with Albert–Chib algorithm: Marginal posterior densities (top panel) and autocorrelation plot (bottom panel).

the standard deviation of the sampled variates and the posterior percentiles are just the percentiles of the sampled draws. As expected, both the first and second covariates increase the probability of infection while the third covariate (the antibiotics prophylaxis) reduces the probability of infection.

To get an idea of the form of the posterior density we plot in Figure 1 the four marginal posterior densities. The density plots are obtained by smoothing the histogram of the simulated values with a Gaussian kernel. In the same plot we also report the autocorrelation functions (correlation against lag) for each of the sampled parameter values. The autocorrelation plots provide information of the extent of serial dependence in the sampled values. Here we see that the serial correlations decline quickly to zero indicating that the sampler is mixing well.

### 2.1. Marginal likelihood of the binary probit

The marginal likelihood of the binary probit model is easily calculated by the method of Chib (1995). Let  $\beta^* = E(\beta|y)$ , then from the expression of the marginal likelihood

in (1.6), we have that

$$\ln m(\mathbf{y}) = \ln \pi(\boldsymbol{\beta}^*) + \ln p(\mathbf{y}|\boldsymbol{\beta}^*) - \ln \pi(\boldsymbol{\beta}^*|\mathbf{y}),$$

where the first term is  $\ln \mathcal{N}_k(\boldsymbol{\beta}^*|\boldsymbol{\beta}_0, \mathbf{B}_0)$ , the second is

$$\sum_{i=1}^n y_i \ln \Phi(\mathbf{x}'_i \boldsymbol{\beta}^*) + (1 - y_i)(1 - \ln \Phi(\mathbf{x}'_i \boldsymbol{\beta}^*))$$

and the third, which is

$$\ln \int \pi(\boldsymbol{\beta}|\mathbf{z})\pi(\mathbf{z}|\mathbf{y}) \, d\mathbf{z}$$

is estimated by averaging the conditional normal density of  $\boldsymbol{\beta}$  in Step 2 of Algorithm 1 at the point  $\boldsymbol{\beta}^*$  over the simulated values  $\{z^{(g)}\}$  from the MCMC run. Specifically, letting  $\hat{\pi}(\boldsymbol{\beta}^*|\mathbf{y})$  denote our estimate of the ordinate at  $\boldsymbol{\beta}^*$  we have

$$\hat{\pi}(\boldsymbol{\beta}^*|\mathbf{y}) = M^{-1} \sum_{i=1}^M \mathcal{N}_k\left(\boldsymbol{\beta}^*|\mathbf{B}_n\left(\mathbf{B}_0^{-1}\boldsymbol{\beta}_0 + \sum_{i=1}^n \mathbf{x}_i z_i^{(g)}\right), \mathbf{B}_n\right).$$

### 2.2. Other link functions

The preceding algorithm is for the probit link function. Albert and Chib (1993) showed how it can be readily extended to other link functions that are generated from the scale mixture of normals family of distributions. For example, we could let

$$\begin{aligned} z_i|\boldsymbol{\beta}, \lambda_i &\sim \mathcal{N}(\mathbf{x}'_i \boldsymbol{\beta}, \lambda_i^{-1}), \\ y_i &= I[z_i > 0], \quad i \leq n, \\ \boldsymbol{\beta} &\sim \mathcal{N}_k(\boldsymbol{\beta}_0, \mathbf{B}_0), \\ \lambda_i &\sim \mathcal{G}\left(\frac{\nu}{2}, \frac{\nu}{2}\right) \end{aligned} \tag{2.3}$$

in which case  $\Pr(y_i = 1|\boldsymbol{\beta}) = F_{t,\nu}(\mathbf{x}'_i \boldsymbol{\beta})$ , where  $F_{t,\nu}$  is the cdf of the standard- $t$  distribution with  $\nu$  degrees of freedom. Albert and Chib (1993) utilized this setup with  $\nu = 8$  to approximate the logit link function. In this case, the posterior distribution of interest is given by

$$\begin{aligned} \pi(\mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\beta}|\mathbf{y}) &\propto \pi(\boldsymbol{\beta}) \prod_{i=1}^n \mathcal{N}(z_i|\mathbf{x}'_i \boldsymbol{\beta}, \lambda_i^{-1}) \mathcal{G}\left(\lambda_i \left| \frac{\nu}{2}, \frac{\nu}{2} \right.\right) \\ &\quad \times \{I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i}\} \end{aligned}$$

which can be sampled in two-blocks as

$$(\mathbf{z}, \boldsymbol{\lambda})|\mathbf{y}, \boldsymbol{\beta}; \boldsymbol{\beta}|\mathbf{y}, \mathbf{z}, \boldsymbol{\lambda},$$

where  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$ . The first block is sampled by the method of composition by sampling  $z_i|y_i, \boldsymbol{\beta}$  marginalized over  $\lambda_i$  from the truncated student- $t$  density

$\mathcal{T}_v(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \{ I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i} \}$ , followed by sampling  $\lambda_i | z_i, \boldsymbol{\beta}$  from an updated Gamma distribution. The second distribution is normal with parameters obtained by standard Bayesian results.

ALGORITHM 2 (*Student-t binary model*).

1 Sample

(a)

$$z_i | y_i, \boldsymbol{\beta} \propto \mathcal{T}_v(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \{ I(z_i < 0)^{1-y_i} + I(z_i > 0)^{y_i} \}$$

(b)

$$\lambda_i | z_i, \boldsymbol{\beta} \sim \mathcal{G}\left(\frac{v+1}{2}, \frac{v + (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2}{2}\right), \quad i \leq n$$

2 Sample

$$\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\lambda} \sim \mathcal{N}_k\left(\mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \lambda_i \mathbf{x}_i z_i \right), \mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}'_i \right)^{-1}\right)$$

3 Goto 1

Chen et al. (1999) present a further extension of this idea to model a skewed link function while Basu and Mukhopadhyay (2000) and Basu and Chib (2003) model the distribution of  $\lambda_i$  by a Dirichlet process prior leading to a semiparametric Bayesian binary data model. The logistic normal link function is considered by Allenby and Lenk (1994) with special reference to data arising in marketing.

### 2.3. Marginal likelihood of the student-t binary model

The marginal likelihood of the student-t binary model is also calculated by the method of Chib (1995). It is easily seen from the expression of the marginal likelihood in (1.6) that the only quantity that needs to be estimated is  $\pi(\boldsymbol{\beta}^* | \mathbf{y})$ , where  $\boldsymbol{\beta}^* = E(\boldsymbol{\beta}^* | \mathbf{y})$ . An estimate of this ordinate is obtained by averaging the conditional normal density of  $\boldsymbol{\beta}$  in Step 2 of Algorithm 2 at the point  $\boldsymbol{\beta}^*$  over the simulated values  $\{z^{(g)}\}$  and  $\{\lambda_i^{(g)}\}$  from the MCMC run. Specifically, letting  $\hat{\pi}(\boldsymbol{\beta}^* | \mathbf{y})$  denote our estimate of the ordinate at  $\boldsymbol{\beta}^*$  we have

$$\hat{\pi}(\boldsymbol{\beta}^* | \mathbf{y}) = M^{-1} \sum_{i=1}^M \mathcal{N}_k\left(\mathbf{B}_n^{(g)} \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \lambda_i^{(g)} \mathbf{x}_i z_i^{(g)} \right), \mathbf{B}_n^{(g)}\right),$$

where

$$\mathbf{B}_n^{(g)} = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \lambda_i^{(g)} \mathbf{x}_i \mathbf{x}'_i \right)^{-1}.$$

### 3. Ordinal response data

Albert and Chib (1993) extend their algorithm to the ordinal categorical data case where  $y_i$  can take one of the values  $\{0, 1, \dots, J\}$  according to the probabilities

$$\Pr(y_i \leq j | \boldsymbol{\beta}, \mathbf{c}) = F(c_j - \mathbf{x}'_i \boldsymbol{\beta}), \quad j = 0, 1, \dots, J, \quad (3.1)$$

where  $F$  is some link function, say either  $\Phi$  or  $F_{t,v}$ . In this model the  $\{c_j\}$  are category specific cut-points such that  $c_0 = 0$  and  $c_J = \infty$ . The remaining cut-points  $\mathbf{c} = (c_1, \dots, c_{J-1})$  are assumed to satisfy the order restriction  $c_1 \leq \dots \leq c_{J-1}$  which ensures that the cumulative probabilities are nondecreasing. Given data  $y_1, \dots, y_n$  from this model, the likelihood function is

$$p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{c}) = \prod_{j=0}^J \prod_{i: y_i=j} [F(c_j - \mathbf{x}'_i \boldsymbol{\beta}) - F(c_{j-1} - \mathbf{x}'_i \boldsymbol{\beta})] \quad (3.2)$$

and the posterior density is proportional to  $\pi(\boldsymbol{\beta}, \mathbf{c}) p(\mathbf{y} | \boldsymbol{\beta}, \mathbf{c})$ , where  $\pi(\boldsymbol{\beta}, \mathbf{c})$  is the prior distribution.

Simulation of the posterior distribution is again feasible by the tactical introduction of latent variables. For the link function  $F_{t,v}$ , we introduce the latent variables  $\mathbf{z} = (z_1, \dots, z_n)$ , where  $z_i | \boldsymbol{\beta}, \lambda_i \sim \mathcal{N}(\mathbf{x}_i \boldsymbol{\beta}, \lambda_i^{-1})$ . A priori, we observe  $y_i = j$  if the latent variable  $z_i$  falls in the interval  $[c_{j-1}, c_j)$ . Define  $\delta_{ij} = 1$  if  $y_i = j$ , and 0 otherwise, then the posterior density of interest is

$$\begin{aligned} \pi(\mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\beta}, \mathbf{c} | \mathbf{y}) \propto & \pi(\boldsymbol{\beta}, \mathbf{c}) \prod_{i=1}^N \mathcal{N}(z_i | \mathbf{x}'_i \boldsymbol{\beta}, \lambda_i^{-1}) \mathcal{G}\left(\lambda_i \left| \frac{v}{2}, \frac{v}{2} \right.\right) \\ & \times \left\{ \sum_{l=0}^J I(c_{l-1} < z_i < c_l) \delta_{il} \right\}, \end{aligned}$$

where  $c_{-1} = -\infty$ . Now the basic Albert and Chib MCMC scheme draws cut-points, the latent data and the regression parameters in sequence. Albert and Chib (2001) simplified the latter step by transforming the cut-points so as to remove the ordering constraint. The transformation is defined by the one-to-one map

$$a_1 = \log c_1; \quad a_j = \log(c_j - c_{j-1}), \quad 2 \leq j \leq J-1. \quad (3.3)$$

The advantage of working with  $\mathbf{a}$  instead of  $\mathbf{c}$  is that the parameters of the tailored proposal density in the M-H step for  $\mathbf{a}$  can be obtained by an unconstrained optimization and the prior  $\pi(\mathbf{a})$  on  $\mathbf{a}$  can be an unrestricted multivariate normal. Next, given  $y_i = j$  and the cut-points the sampling of the latent data  $z_i$  is from  $\mathcal{T}_v(\mathbf{x}'_i \boldsymbol{\beta}, 1) I(c_{j-1} < z_i < c_j)$ , marginalized over  $\{\lambda_i\}$ . The sampling of the parameters  $\boldsymbol{\beta}$  is as in Algorithm 2.

ALGORITHM 3 (*Ordinal student link model*).

- 1 M-H
  - (a) Sample

$$\pi(\mathbf{a}|\mathbf{y}, \boldsymbol{\beta}) \propto \pi(\mathbf{a}) \underbrace{\prod_{j=0}^J \prod_{i: y_i=j} [F_{t,v}(c_j - \mathbf{x}'_i \boldsymbol{\beta}) - F_{t,v}(c_{j-1} - \mathbf{x}'_i \boldsymbol{\beta})]}_{p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{a})}$$

with the M-H algorithm by calculating

$$\mathbf{m} = \arg \max_{\mathbf{a}} \log p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{a})$$

and  $\mathbf{V} = \{-\partial \log f(\mathbf{y}|\boldsymbol{\beta}, \mathbf{a})/\partial \mathbf{a} \partial \mathbf{a}'\}^{-1}$  the negative inverse of the hessian at  $\mathbf{m}$ , proposing

$$\mathbf{a}' \sim f_T(\mathbf{a}|\mathbf{m}, \mathbf{V}, \xi),$$

for some choice of  $\xi$ , calculating

$$\alpha(\mathbf{a}, \mathbf{a}'|\mathbf{y}, \boldsymbol{\beta}) = \min \left\{ \frac{\pi(\mathbf{a}') p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{a}')}{\pi(\mathbf{a}) p(\mathbf{y}|\boldsymbol{\beta}, \mathbf{a})} \frac{\mathcal{T}_v(\mathbf{a}|\mathbf{m}, \mathbf{V}, \xi)}{\mathcal{T}_v(\mathbf{a}'|\mathbf{m}, \mathbf{V}, \xi)}, 1 \right\}$$

and moving to  $\mathbf{a}'$  with probability  $\alpha(\mathbf{a}, \mathbf{a}'|\mathbf{y}, \boldsymbol{\beta})$ , then transforming the new  $\mathbf{a}$  to  $\mathbf{c}$  via the inverse map  $c_j = \sum_{i=1}^j \exp(a_i)$ ,  $1 \leq j \leq J - 1$ .

(b) Sample

$$z_i | y_i, \boldsymbol{\beta}, \mathbf{c} \propto \mathcal{T}_v(z_i | \mathbf{x}'_i \boldsymbol{\beta}, 1) \left\{ \sum_{l=0}^J I(c_{l-1} < z_i < c_l)^{\delta_{il}} \right\}, \quad i \leq n$$

(c) Sample

$$\lambda_i | y_i, z_i, \boldsymbol{\beta}, \mathbf{c} \sim \mathcal{G} \left( \frac{v+1}{2}, \frac{v + (z_i - \mathbf{x}'_i \boldsymbol{\beta})^2}{2} \right), \quad i \leq n$$

2 Sample

$$\boldsymbol{\beta} | \mathbf{z}, \boldsymbol{\lambda} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \lambda_i \mathbf{x}_i z_i \right), \mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \lambda_i \mathbf{x}_i \mathbf{x}'_i \right)^{-1} \right)$$

3 Goto 1

### 3.1. Marginal likelihood of the student-t ordinal model

The marginal likelihood of the student-t ordinal model can be calculated by the method of Chib (1995) and Chib and Jeliazkov (2001). Let  $\boldsymbol{\beta}^* = E(\boldsymbol{\beta}|\mathbf{y})$  and  $\mathbf{c}^* = E(\mathbf{c}^*|\mathbf{y})$  and let  $\mathbf{a}^*$  be the transformed  $\mathbf{c}^*$ . Then from the expression of the marginal likelihood in (1.6) we can write

$$\ln m(\mathbf{y}) = \ln \pi(\boldsymbol{\beta}^*, \mathbf{a}^*) + \ln p(\mathbf{y}|\boldsymbol{\beta}^*, \mathbf{c}^*) - \ln \pi(\mathbf{a}^*|\mathbf{y}) - \ln \pi(\boldsymbol{\beta}^*|\mathbf{y}, \mathbf{a}^*).$$

The marginal posterior ordinate  $\pi(\mathbf{a}^*|\mathbf{y})$  from (1.11) can be written as

$$\pi(\mathbf{a}^*|\mathbf{y}) = \frac{E_1\{\alpha(\mathbf{a}, \mathbf{a}^*|\mathbf{y}, \boldsymbol{\beta})\mathcal{T}_v(\mathbf{a}'|\mathbf{m}, \mathbf{V}, \xi)\}}{E_2\{\alpha(\mathbf{a}^*, \mathbf{a}|\mathbf{y}, \boldsymbol{\beta})\}},$$

where  $E_1$  is the expectation with respect to  $\pi(\boldsymbol{\beta}|\mathbf{y})$  and  $E_2$  is the expectation with respect to  $\pi(\boldsymbol{\beta}|\mathbf{y}, \mathbf{a}^*)\mathcal{T}_v(\mathbf{a}|\mathbf{m}, \mathbf{V}, \xi)$ . The  $E_1$  expectation can be estimated by draws from the main run and the  $E_2$  expectation from a reduced run in which  $\mathbf{a}$  is set to  $\mathbf{a}^*$ . For each value of  $\boldsymbol{\beta}$  in this reduced run,  $\mathbf{a}$  is also sampled from  $\mathcal{T}_v(\mathbf{a}|\mathbf{m}, \mathbf{V}, \xi)$ ; these two draws (namely the draws of  $\boldsymbol{\beta}$  and  $\mathbf{a}$ ) are used to average  $\alpha(\mathbf{a}^*, \mathbf{a}|\mathbf{y}, \boldsymbol{\beta})$ . The final ordinate  $\pi(\boldsymbol{\beta}^*|\mathbf{y}, \mathbf{a}^*)$  is estimated by averaging the normal conditional density in Step 2 of Algorithm 3 over the values of  $\{z_i\}$  and  $\{\lambda_i\}$  from this same reduced run.

#### 4. Sequential ordinal model

The ordinal model presented above is appropriate when one can assume that a single unobservable continuous variable underlies the ordinal response. For example, this representation can be used for modeling ordinal letter grades in a mathematics class by assuming that the grade is an indicator of the student's unobservable continuous-valued intelligence. In other circumstances, however, a different latent variable structure may be required to model the ordinal response. For example, McCullagh (1980), and Farhmeir and Tutz (1994) consider a data set where the relative tonsil size of children is classified into the three states: "present but not enlarged", "enlarged", and "greatly enlarged". The objective is to explain the ordinal response as a function of whether the child is a carrier of the *Streptococcus pyogenes*. In this setting, it may not be appropriate to model the ordinal tonsil size in terms of a single latent variable. Rather, it may be preferable to imagine that each response is determined by two continuous latent variables where the first latent variable measures the propensity of the tonsil to grow abnormally and pass from the state "present but not enlarged" to the state "enlarged". The second latent variable measures the propensity of the tonsil to grow from the "enlarged" to the "greatly enlarged" states. Thus, the final "greatly enlarged" state is realized only when the tonsil has passed through the earlier levels. This latent variable representation leads to a *sequential* model because the levels of the response are achieved in a sequential manner.

Suppose that one observes independent observations  $y_1, \dots, y_n$ , where each  $y_i$  is an ordinal categorical response variable with  $J$  possible values  $\{1, \dots, J\}$ . Associated with the  $i$ th response  $y_i$ , let  $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})$  denote a set of  $k$  covariates. In the sequential ordinal model, the variable  $y_i$  can take the value  $j$  only after the levels  $1, \dots, j-1$  are reached. In other words, to get to the outcome  $j$ , one must pass through levels  $1, 2, \dots, j-1$  and stop (or fail) in level  $j$ . The probability of stopping in level  $j$  ( $1 \leq j \leq J-1$ ), conditional on the event that the  $j$ th level is reached, is given by

$$\Pr(y_i = j | y_i \geq j, \boldsymbol{\delta}, \mathbf{c}) = F(c_j - \mathbf{x}'_i \boldsymbol{\delta}), \quad (4.1)$$

where  $\mathbf{c} = (c_1, \dots, c_{J-1})$  are unordered cutpoints and  $\mathbf{x}'_i \boldsymbol{\delta}$  represents the effect of the covariates. This probability function is referred to as the discrete time hazard function

(Tutz, 1990; Farhmeir and Tutz, 1994). It follows that the probability of stopping at level  $j$  is given by

$$\begin{aligned} \Pr(y_i = j | \boldsymbol{\delta}, \mathbf{c}) &= \Pr(y_i = j | y_i \geq j, \boldsymbol{\delta}, \mathbf{c}) \Pr(y_i \geq j) \\ &= F(c_j - \mathbf{x}'_i \boldsymbol{\delta}) \prod_{k=1}^{j-1} \{1 - F(c_k - \mathbf{x}'_i \boldsymbol{\delta})\}, \quad j \leq J - 1, \end{aligned} \quad (4.2)$$

whereas the probability that the final level  $J$  is reached is

$$\Pr(y_i = J | \boldsymbol{\delta}, \mathbf{c}) = \prod_{k=1}^{J-1} \{1 - F(c_k - \mathbf{x}'_i \boldsymbol{\delta})\}, \quad (4.3)$$

since the event  $y_i = J$  occurs only if all previous  $J - 1$  levels are passed.

The sequential model is formally equivalent to the continuation-ratio ordinal models, discussed by Agresti (1990, Chapter 9), Cox (1972) and Ten Have and Uttal (2000). The sequential model is useful in the analysis of discrete-time survival data (Farhmeir and Tutz (1994), Chapter 9; Kalbfleish and Prentice, 1980). More generally, the sequential model can be used to model nonproportional and nonmonotone hazard functions and to incorporate the effect of time dependent covariates. Suppose that one observes the time to failure for subjects in the sample, and let the time interval be subdivided (or grouped) into the  $J$  intervals  $[a_0, a_1), [a_1, a_2), \dots, [a_{J-1}, \infty)$ . Then one defines the event  $y_i = j$  if a failure is observed in the interval  $[a_{j-1}, a_j)$ . The discrete hazard function (4.1) now represents the probability of failure in the time interval  $[a_{j-1}, a_j)$  given survival until time  $a_{j-1}$ . The vector of cutpoint parameters  $\mathbf{c}$  represents the baseline hazard of the process.

Albert and Chib (2001) develop a full Bayesian analysis of this model based on the framework of Albert and Chib (1993). Suppose that  $F$  is the distribution function of the standard normal distribution, leading to the sequential ordinal model. Corresponding to the  $i$ th observation, define latent variables  $\{w_{ij}\}$ , where  $w_{ij} = \mathbf{x}'_i \boldsymbol{\delta} + e_{ij}$  and the  $e_{ij}$  are independently distributed from  $N(0, 1)$ . We observe  $y_i = 1$  if  $w_{i1} \leq c_1$ , and observe  $y_i = 2$  if the first latent variable  $w_{i1} > c_1$  and the second variable  $w_{i2} \leq c_2$ . In general, we observe  $y_i = j$  ( $1 \leq j \leq J - 1$ ) if the first  $j - 1$  latent variables exceed their corresponding cutoffs, and the  $j$ th variable does not:  $y_i = j$  if  $w_{i1} > c_1, \dots, w_{ij-1} > c_{j-1}, w_{ij} \leq c_j$ . In this model, the latent variable  $w_{ij}$  represents one's propensity to continue to the  $(j + 1)$ st level in the sequence, given that the individual has already attained level  $j$ .

This latent variable representation can be simplified by incorporating the cutpoints  $\{\gamma_j\}$  into the mean function. Define the new latent variable  $z_{ij} = w_{ij} - c_j$ . Then it follows that  $z_{ij} | \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{x}'_{ij} \boldsymbol{\beta}, 1)$ , where  $\mathbf{x}'_{ij} = (0, 0, \dots, -1, 0, 0, \mathbf{x}'_i)$  with  $-1$  in the  $j$ th column, and  $\boldsymbol{\beta} = (c_1, \dots, c_{J-1}, \boldsymbol{\delta})'$ . The observed data are then generated

according to

$$y_i = \begin{cases} 1 & \text{if } z_{i1} \leq 0, \\ 2 & \text{if } z_{i1} > 0, z_{i2} \leq 0, \\ \vdots & \vdots \\ J - 1 & \text{if } z_{i1} > 0, \dots, z_{iJ-2} > 0, z_{iJ-1} \leq 0, \\ J & \text{if } z_{i1} > 0, \dots, z_{iJ-1} > 0. \end{cases} \tag{4.4}$$

Following the [Albert and Chib \(1993\)](#) approach, this model can be fit by MCMC methods by simulating the joint posterior distribution of  $(\{z_{ij}\}, \boldsymbol{\beta})$ . The latent data  $\{z_{ij}\}$ , conditional on  $(y, \boldsymbol{\beta})$ , are independently distributed as truncated normal. Specifically, if  $y_i = j$ , then the latent data corresponding to this observation are represented as  $\mathbf{z}_i = (z_{i1}, \dots, z_{ij_i})$  where  $j_i = \min\{j, J - 1\}$  and the simulation of  $\mathbf{z}_i$  is from a sequence of truncated normal distributions. The posterior distribution of the parameter vector  $\boldsymbol{\beta}$ , conditional on the latent data  $\mathbf{z} = (z_1, \dots, z_n)$ , has a simple form. Let  $\mathbf{X}_i$  denote the covariate matrix corresponding to the  $i$ th subject consisting of the  $j_i$  rows  $\mathbf{x}'_{i1}, \dots, \mathbf{x}'_{ij_i}$ . If  $\boldsymbol{\beta}$  is assigned a multivariate normal prior with mean vector  $\boldsymbol{\beta}_0$  and covariance matrix  $\mathbf{B}_0$ , then the posterior distribution of  $\boldsymbol{\beta}$ , conditional on  $\mathbf{z}$ , is multivariate normal.

ALGORITHM 4 (*Sequential ordinal probit model*).

1 Sample  $z_i | y_i = j, \boldsymbol{\beta}$

(a) If  $y_i = 1$ , sample

$$z_{i1} | y_i, \boldsymbol{\beta} \propto \mathcal{N}(\mathbf{x}'_{i1}\boldsymbol{\beta}, 1)I(z_{i1} < 0)$$

(b) If  $y_i = j$  ( $2 \leq j \leq J - 1$ ), sample

$$z_{ik} | y_i, \boldsymbol{\beta} \propto \mathcal{N}(\mathbf{x}'_{ik}\boldsymbol{\beta}, 1)I(z_{ik} > 0), \quad k = 1, \dots, j - 1$$

$$z_{ij} | y_i, \boldsymbol{\beta} \propto \mathcal{N}(\mathbf{x}'_{ij}\boldsymbol{\beta}, 1)I(z_{ij} < 0)$$

(c) If  $y_i = J$ , sample

$$z_{ik} | y_i, \boldsymbol{\beta} \propto \mathcal{N}(\mathbf{x}'_{ik}\boldsymbol{\beta}, 1)I(z_{ik} > 0), \quad k = 1, \dots, J - 1$$

2 Sample

$$\boldsymbol{\beta} | \mathbf{z} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \mathbf{X}'_i \mathbf{z}_i \right), \mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \mathbf{X}'_i \mathbf{X}_i \right)^{-1} \right)$$

3 Goto 1

[Albert and Chib \(2001\)](#) present generalizations of the sequential ordinal model and also discuss the computation of the marginal likelihood by the method of [Chib \(1995\)](#).

### 5. Multivariate responses

One of the considerable achievements of the recent Bayesian literature is in the development of methods for multivariate categorical data. A leading case concerns that



of multivariate binary data. From the frequentist perspective, this problem is tackled by Carey et al. (1993) and Glonek and McCullagh (1995) in the context of the multivariate logit model, and from a Bayesian perspective by Chib and Greenberg (1998) in the context of the multivariate probit (MVP) model (an analysis of the multivariate logit model from a Bayesian perspective is supplied by O'Brien and Dunson, 2004). In many respects, the multivariate probit model, along with its variants, is the canonical modeling approach for correlated binary data. This model was introduced more than twenty five years ago by Ashford and Sowden (1970) in the context of bivariate binary responses and was analyzed under simplifying assumptions on the correlation structure by Amemiya (1972), Ochi and Prentice (1984) and Lesaffre and Kaufmann (1992). The general version of the model was considered to be intractable. Chib and Greenberg (1998), building on the framework of Albert and Chib (1993), resolve the difficulties.

The MVP model provides a relatively straightforward way of modeling multivariate binary data. In this model, the marginal probability of each response is given by a probit function that depends on covariates and response specific parameters. Associations between the binary variables are incorporated by assuming that the vector of binary outcomes are a function of correlated Gaussian variables, taking the value one if the corresponding Gaussian component is positive and the value zero otherwise. This connection with a latent Gaussian random vector means that the regression coefficients can be interpreted independently of the correlation parameters (unlike the case of log-linear models). The link with Gaussian data is also helpful in estimation. By contrast, models based on marginal odds ratios Connolly and Liang (1988) tend to proliferate nuisance parameters as the number of variables increase and they become difficult to interpret and estimate. Finally, the Gaussian connection enables generalizations of the model to ordinal outcomes (see Chen and Dey, 2000), mixed outcomes (Kuhnert and Do, 2003), and spatial data (Fahrmeir and Lang, 2001; Banerjee et al., 2004).

### 5.1. Multivariate probit model

Let  $y_{ij}$  denote a binary response on the  $i$ th observation unit and  $j$ th variable, and let  $y_i = (y_{i1}, \dots, y_{iJ})'$ ,  $1 \leq i \leq n$ , denote the collection of responses on all  $J$  variables. Also let  $x_{ij}$  denote the set of covariates for the  $j$ th response and  $\beta_j \in R^{k_j}$  the conformable vector of covariate coefficients. Let  $\beta' = (\beta'_1, \dots, \beta'_J) \in R^k$ ,  $k = \sum k_j$  denote the complete set of coefficients and let  $\Sigma = \{\sigma_{jk}\}$  denote a  $J \times J$  correlation matrix. Finally let

$$X_i = \begin{pmatrix} x'_{i1} & \mathbf{0}' & \dots & \mathbf{0}' \\ \mathbf{0}' & x'_{i2} & \mathbf{0}' & \mathbf{0}' \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \dots & x'_{iJ} \end{pmatrix}$$

denote the  $J \times k$  covariate matrix on the  $i$ th subject. Then, according to the multivariate probit model, the marginal probability that  $y_{ij} = 1$  is given by the probit form

$$\Pr(y_{ij} = 1|\beta) = \Phi(x'_{ij}\beta_j)$$

and the joint probability of a given outcome  $y_i$ , conditioned on parameters  $\beta$ ,  $\Sigma$ , and covariates  $x_{ij}$ , is

$$\Pr(y_i|\beta, \Sigma) = \int_{A_{iJ}} \dots \int_{A_{i1}} \mathcal{N}_J(t|0, \Sigma) dt, \tag{5.1}$$

where  $\mathcal{N}_J(t|0, \Sigma)$  is the density of a  $J$ -variate normal distribution with mean vector 0 and correlation matrix  $\Sigma$  and  $A_{ij}$  is the interval

$$A_{ij} = \begin{cases} (-\infty, x'_{ij}\beta_j) & \text{if } y_{ij} = 1, \\ [x'_{ij}\beta_j, \infty) & \text{if } y_{ij} = 0. \end{cases} \tag{5.2}$$

Note that each outcome is determined by its own set of  $k_j$  covariates  $x_{ij}$  and covariate effects  $\beta_j$ .

The multivariate discrete mass function presented in (5.1) can be specified in terms of latent Gaussian random variables. This alternative formulation also forms the basis of the computational scheme that is described below. Let  $z_i = (z_{i1}, \dots, z_{iJ})$  denote a  $J$ -variate normal vector and let

$$z_i \sim \mathcal{N}_J(X_i\beta, \Sigma). \tag{5.3}$$

Now let  $y_{ij}$  be 1 or 0 according to the sign of  $z_{ij}$ :

$$y_{ij} = I(z_{ij} > 0), \quad j = 1, \dots, J. \tag{5.4}$$

Then, the probability in (5.1) may be expressed as

$$\int_{B_{iJ}} \dots \int_{B_{i1}} \mathcal{N}_J(z_i|X_i\beta, \Sigma) dz_i, \tag{5.5}$$

where  $B_{ij}$  is the interval  $(0, \infty)$  if  $y_{ij} = 1$  and the interval  $(-\infty, 0]$  if  $y_{ij} = 0$ . It is easy to confirm that this integral reduces to the form given above. It should also be noted that due to the threshold specification in (5.4), the scale of  $z_{ij}$  cannot be identified. As a consequence, the matrix  $\Sigma$  must be in correlation form (with units on the main diagonal).

### 5.2. Dependence structures

One basic question in the analysis of correlated binary data is the following: How should correlation between binary outcomes be defined and measured? The point of view behind the MVP model is that the correlation is modeled at the level of the latent data which then induces correlation amongst the binary outcomes. This modeling perspective is both flexible and general. In contrast, attempts to model correlation directly (as in the classical literature using marginal odds ratios) invariably lead to difficulties, partly because it is difficult to specify pair-wise correlations in general, and partly because the binary data scale is not natural for thinking about dependence.

Within the context of the MVP model, alternative dependence structures are easily specified and conceived, due to the connection with Gaussian latent data. Some of the possibilities are enumerated below.

- Unrestricted form. Here  $\Sigma$  is fully unrestricted except for the unit constraints on the diagonal. The unrestricted  $\Sigma$  matrix has  $J^* = J(J - 1)/2$  unknown correlation parameters that must be estimated.
- Equicorrelated form. In this case, the correlations are all equal and described by a single parameter  $\rho$ . This form can be a starting point for the analysis when one is dealing with outcomes where all the pair-wise correlations are believed to have the same sign. The equicorrelated model may be also be seen as arising from a random effect formulation.
- Toeplitz form. Under this case, the correlations depend on a single parameter  $\rho$  but under the restriction that  $\text{Corr}(z_{ik}, z_{il}) = \rho^{|k-l|}$ . This version can be useful when the binary outcomes are collected from a longitudinal study where it is plausible that the correlation between outcomes at different dates will diminish with the lag. In fact, in the context of longitudinal data, the  $\Sigma$  matrix can be specified in many other forms with analogy with the correlation structures that arise in standard time series ARMA modeling.

### 5.3. Student-*t* specification

Now suppose that the distribution on the latent  $z_i$  is multivariate-*t* with specified degrees of freedom  $\nu$ . This gives rise to a model that may be called the multivariate-*t* link model. Under the multivariate-*t* assumption,  $z_i | \beta, \Sigma \sim \mathcal{T}_\nu(X_i \beta, \Sigma)$  with density

$$\mathcal{T}_\nu(z_i | \beta, \Sigma) \propto |\Sigma|^{-1/2} \left\{ 1 + \frac{1}{\nu} (z_i - X_i \beta)' \Sigma^{-1} (z_i - X_i \beta) \right\}^{-(\nu+J)/2}. \tag{5.6}$$

As before, the matrix  $\Sigma$  is in correlation form and the observed outcomes are defined by (5.4). The model for the latent  $z_i$  may be expressed as a scale mixture of normals by introducing a random variable  $\lambda_i \sim \mathcal{G}(\frac{\nu}{2}, \frac{\nu}{2})$  and letting

$$z_i | \beta, \Sigma, \lambda_i \sim \mathcal{N}_J(X_i \beta, \lambda_i^{-1} \Sigma). \tag{5.7}$$

Conditionally on  $\lambda_i$ , this model is equivalent to the MVP model.

### 5.4. Estimation of the MVP model

Consider now the question of inference in the MVP model. We are given a set of data on  $n$  subjects with outcomes  $\mathbf{y} = \{y_i\}_{i=1}^n$  and interest centers on the parameters of the model  $\theta = (\beta, \Sigma)$  and the posterior distribution  $\pi(\beta, \Sigma | \mathbf{y})$ , given some prior distribution  $\pi(\beta, \Sigma)$  on the parameters.

To simplify the MCMC implementation for this model Chib and Greenberg (1998) follow the general approach of Albert and Chib (1993) and employ the latent variables  $z_i \sim \mathcal{N}_J(X_i \beta, \Sigma)$ , with the observed data given by  $y_{ij} = I(z_{ij} > 0)$ ,  $j = 1, \dots, J$ . Let  $\sigma = (\sigma_{12}, \sigma_{31}, \sigma_{32}, \dots, \sigma_{JJ})$  denote the  $J(J - 1)/2$  distinct elements of  $\Sigma$ . It can be shown that the admissible values of  $\sigma$  (that lead to a positive definite  $\Sigma$  matrix) form a convex solid body in the hypercube  $[-1, 1]^P$ . Denote this set by  $C$ . Now let  $\mathbf{z} = (z_1, \dots, z_n)$  denote the latent values corresponding to the observed data  $\mathbf{y} = \{y_i\}_{i=1}^n$ .

Then, the algorithm proceeds with the sampling of the posterior density

$$\begin{aligned} \pi(\boldsymbol{\beta}, \boldsymbol{\sigma}, \mathbf{z}|\mathbf{y}) &\propto \pi(\boldsymbol{\beta}, \boldsymbol{\sigma}) f(\mathbf{z}|\boldsymbol{\beta}, \boldsymbol{\Sigma}) \Pr(\mathbf{y}|\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) \\ &\propto \pi(\boldsymbol{\beta}, \boldsymbol{\sigma}) \prod_{i=1}^n (\phi_J(\mathbf{z}_i|\boldsymbol{\beta}, \boldsymbol{\Sigma}) \Pr(y_i|\mathbf{z}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma})), \quad \boldsymbol{\beta} \in \mathfrak{R}^k, \boldsymbol{\sigma} \in C, \end{aligned}$$

where, akin to the result in the univariate binary case,

$$\Pr(y_i|\mathbf{z}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}) = \prod_{j=1}^J \{I(z_{ij} \leq 0)^{1-y_{ij}} + I(z_{ij} > 0)^{y_{ij}}\}.$$

Conditioned on  $\{\mathbf{z}_i\}$  and  $\boldsymbol{\Sigma}$ , the update for  $\boldsymbol{\beta}$  is straightforward, while conditioned on  $(\boldsymbol{\beta}, \boldsymbol{\Sigma})$ ,  $z_{ij}$  can be sampled one at a time conditioned on the other latent values from truncated normal distributions, where the region of truncation is either  $(0, \infty)$  or  $(-\infty, 0)$  depending on whether the corresponding  $y_{ij}$  is one or zero. The key step in the algorithm is the sampling of  $\boldsymbol{\sigma}$ , the unrestricted elements of  $\boldsymbol{\Sigma}$ , from the full conditional density  $\pi(\boldsymbol{\sigma}|\mathbf{z}, \boldsymbol{\beta}) \propto p(\boldsymbol{\sigma}) \prod_{i=1}^n \mathcal{N}_J(\mathbf{z}_i|X_i\boldsymbol{\beta}, \boldsymbol{\Sigma})$ . This density, which is truncated to the complicated region  $C$ , is sampled by a M–H step with tailored proposal density  $q(\boldsymbol{\sigma}|\mathbf{z}, \boldsymbol{\beta}) = \mathcal{T}_v(\boldsymbol{\sigma}|\mathbf{m}, \mathbf{V}, \xi)$  where

$$\begin{aligned} \mathbf{m} &= \arg \max_{\boldsymbol{\sigma} \in C} \sum_{i=1}^n \ln \mathcal{N}_J(\mathbf{z}_i|X_i\boldsymbol{\beta}, \boldsymbol{\Sigma}), \\ \mathbf{V} &= - \left\{ \frac{\partial^2 \sum_{i=1}^n \ln \mathcal{N}_J(\mathbf{z}_i|X_i\boldsymbol{\beta}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\sigma} \partial \boldsymbol{\sigma}'} \right\}_{\boldsymbol{\sigma}=\mathbf{m}} \end{aligned}$$

are the mode and curvature of the target distribution, given the current values of the conditioning variables.

ALGORITHM 5 (*Multivariate probit*).

1 Sample for  $i \leq n, j \leq J$

$$\begin{aligned} z_{ij}|\mathbf{z}_{i(-j)}, \boldsymbol{\beta}, \boldsymbol{\Sigma} &\propto \mathcal{N}(\mu_{ij}, v_{ij}) \{I(z_{ij} \leq 0)^{1-y_{ij}} + I(z_{ij} > 0)^{y_{ij}}\} \\ \mu_{ij} &= E(z_{ij}|\mathbf{z}_{i(-j)}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) \\ v_{ij} &= \text{Var}(z_{ij}|\mathbf{z}_{i(-j)}, \boldsymbol{\beta}, \boldsymbol{\Sigma}) \end{aligned}$$

2 Sample

$$\boldsymbol{\beta}|\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\Sigma} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n X_i' \boldsymbol{\Sigma}^{-1} \mathbf{z}_i \right), \mathbf{B}_n \right),$$

where

$$\mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n X_i' \boldsymbol{\Sigma}^{-1} X_i \right)^{-1}$$

3 Sample

$$\pi(\sigma | \mathbf{z}, \boldsymbol{\beta}) \propto p(\sigma) \prod_{i=1}^n \mathcal{N}_J(\mathbf{z}_i | \mathbf{X}_i \boldsymbol{\beta}, \boldsymbol{\Sigma})$$

by the M-H algorithm, calculating the parameters  $(\mathbf{m}, \mathbf{V})$ , proposing  $\sigma' \sim \mathcal{T}_v(\sigma | \mathbf{m}, \mathbf{V}, \xi)$ , calculating

$$\alpha(\sigma, \sigma' | \mathbf{y}, \boldsymbol{\beta}, \{\mathbf{z}_i\}) = \min \left\{ \frac{\pi(\sigma') \prod_{i=1}^n \mathcal{N}_J(\mathbf{z}_i | \mathbf{X}_i \boldsymbol{\beta}, \boldsymbol{\Sigma}') I[\sigma' \in C] \mathcal{T}_v(\sigma | \mathbf{m}, \mathbf{V}, \xi)}{\pi(\sigma) \prod_{i=1}^n \mathcal{N}_J(\mathbf{z}_i | \mathbf{X}_i \boldsymbol{\beta}, \boldsymbol{\Sigma}) \mathcal{T}_v(\sigma' | \mathbf{m}, \mathbf{V}, \xi)}, 1 \right\}$$

and moving to  $\sigma'$  with probability  $\alpha(\sigma, \sigma' | \mathbf{y}, \boldsymbol{\beta}, \{\mathbf{z}_i\})$

4 Goto 1

EXAMPLE. As an application of this algorithm consider a data set in which the multivariate binary responses are generated by a panel structure. The data is concerned with the health effects of pollution on 537 children in Stuebenville, Ohio, each observed at ages 7, 8, 9 and 10 years, and the response variable is an indicator of wheezing status. Suppose that the marginal probability of wheeze status of the  $i$ th child at the  $j$ th time point is specified as

$$\Pr(y_{ij} = 1 | \boldsymbol{\beta}) = \Phi(\beta_0 + \beta_1 x_{1ij} + \beta_2 x_{2ij} + \beta_3 x_{3ij}), \quad i \leq 537, j \leq 4,$$

where  $\boldsymbol{\beta}$  is constant across categories,  $x_1$  is the age of the child centered at nine years,  $x_2$  is a binary indicator variable representing the mother’s smoking habit during the first year of the study, and  $x_3 = x_1 x_2$ . Suppose that the Gaussian prior on  $\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, \beta_4)$  is centered at zero with a variance of  $10\mathbf{I}_k$  and let  $\pi(\sigma)$  be the density of a normal distribution, with mean zero and variance  $\mathbf{I}_6$ , restricted to region that leads to a positive-definite correlation matrix, where  $(\sigma_{21}, \sigma_{31}, \sigma_{32}, \sigma_{41}, \sigma_{42}, \sigma_{43})$ . From 10,000 cycles of Algorithm 5 one obtains the following covariate effects and posterior distributions of the correlations.

Table 3  
Covariate effects in the Ohio wheeze data: MVP model with unrestricted correlations. In the table, NSE denotes the numerical standard error, lower is the 2.5th percentile and upper is the 97.5th percentile of the simulated draws. The results are based on 10000 draws from Algorithm 5

|           | Prior |         | Posterior |       |         |        |        |
|-----------|-------|---------|-----------|-------|---------|--------|--------|
|           | Mean  | Std dev | Mean      | NSE   | Std dev | Lower  | Upper  |
| $\beta_1$ | 0.000 | 3.162   | -1.108    | 0.001 | 0.062   | -1.231 | -0.985 |
| $\beta_2$ | 0.000 | 3.162   | -0.077    | 0.001 | 0.030   | -0.136 | -0.017 |
| $\beta_3$ | 0.000 | 3.162   | 0.155     | 0.002 | 0.101   | -0.043 | 0.352  |
| $\beta_4$ | 0.000 | 3.162   | 0.036     | 0.001 | 0.049   | -0.058 | 0.131  |

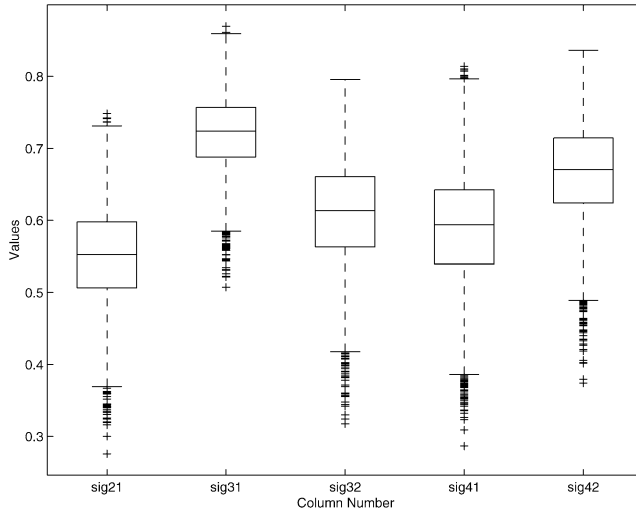


Fig. 2. Posterior boxplots of the correlations in the Ohio wheeze data: MVP model.

Notice that the summary tabular output contains not only the posterior means and standard deviations of the parameters but also the 95% credibility intervals, all computed from the sampled draws. It may be seen from Figure 2 that the posterior distributions of the correlations are similar suggesting that an equicorrelated correlation structure might be appropriate for these data.

5.5. Marginal likelihood of the MVP model

The calculation of the marginal likelihood of MVP model is considered by both Chib and Greenberg (1998) and Chib and Jeliazkov (2001). The main issue is the calculation of the posterior ordinate which may be estimated from

$$\pi(\sigma^*, \beta^* | y) = \pi(\sigma^* | y) \pi(\beta^* | y, \Sigma^*),$$

where the first ordinate (that of the correlations) is not in tractable form. To estimate the marginal ordinate one can apply (1.11) leading to

$$\pi(\sigma^* | y) = \frac{E_1\{\alpha(\sigma, \sigma' | y, \beta, \{z_i\}) \mathcal{T}_V(\sigma' | m, V, \xi)\}}{E_2\{\alpha(\sigma^*, \sigma | y, \beta, \{z_i\})\}},$$

where  $\alpha(\sigma, \sigma' | y, \beta, \{z_i\})$  is the probability of move defined in Algorithm 5. The numerator expectation can be calculated from the draws  $\{\beta^{(g)}, \{z_i^{(g)}\}, \sigma^{(g)}\}$  from the output of Algorithm 5. The denominator expectation is calculated from a reduced run consisting of the densities

$$\pi(\beta | y, \{z_i\}, \Sigma^*); \quad \pi(\{z_i\} | y, \beta, \Sigma^*),$$

after  $\Sigma$  is fixed at  $\Sigma^*$ . Then for each draw

$$\beta^{(j)}, z^{(j)} \sim \pi(\beta, z | y, \Sigma^*)$$

in this reduced run, we also draw  $\sigma^{(j)} \sim \mathcal{T}_v(\sigma|\mathbf{m}, \mathbf{V}, \xi)$ . These draws are used to average  $\alpha(\sigma^*, \sigma|\mathbf{y}, \boldsymbol{\beta}, \{z_i\})$ . In addition, the sampled variates  $\{\boldsymbol{\beta}^{(j)}, z^{(j)}\}$  from this same reduced run are also used to estimate  $\pi(\boldsymbol{\beta}^*|\mathbf{y}, \boldsymbol{\Sigma}^*)$  by averaging the normal density in Step 2 of Algorithm 5.

5.6. Fitting of the multivariate t-link model

Algorithm 5 is easily modified for the fitting of the multivariate-t link version of the MVP model. One simple possibility is to include the  $\{\lambda_i\}$  into the sampling since conditioned on the value of  $\lambda_i$ , the t-link binary model reduces to the MVP model. With this augmentation, one implements Steps 1–3 conditioned on the value of  $\{\lambda_i\}$ . The sampling is completed with an additional step involving the simulation of  $\{\lambda_i\}$ . A straightforward calculation shows that the updated full conditional distribution of  $\lambda_i$  is

$$\lambda_i|z_i, \boldsymbol{\beta}, \boldsymbol{\Sigma} \sim \mathcal{G}\left(\frac{v + J}{2}, \frac{v + (z_i - \mathbf{X}_i\boldsymbol{\beta})'\boldsymbol{\Sigma}^{-1}(z_i - \mathbf{X}_i\boldsymbol{\beta})}{2}\right), \quad i \leq n,$$

which is easily sampled. The modified sampler thus requires little extra coding.

5.7. Binary outcome with a confounded binary treatment

In many problems, one central objective of model building is to isolate the effect of a binary covariate on a binary (or categorical) outcome. Isolating this effect becomes difficult if the effect of the binary covariate on the outcome is confounded with those of missing or unobserved covariates. Such complications are the norm when dealing with observational data. As an example, suppose that  $y_i$  is a zero-one variable that represents post-surgical outcomes for hip-fracture patients (one if the surgery is successful, and zero if not) and  $x_i$  is an indicator representing the extent of delay before surgery (one if the delay exceeds say two days, and zero otherwise). In this case, it is quite likely that delay and post-surgical outcome are both affected by patient factors, e.g., patient frailty, that are usually either imperfectly measured and recorded or absent from the data. In this case, we can model the outcome and the treatment assignment jointly to control for the unobservable confounders. The resulting model is similar to the MVP model and is discussed in detail by Chib (2003).

Briefly, suppose that the outcome model is given by

$$y_i = I(z_i'\boldsymbol{\gamma}_1 + x_i\boldsymbol{\beta} + \varepsilon_i), \quad i \leq n, \tag{5.8}$$

and suppose that the effect of the binary treatment  $x_i$  is confounded with that of unobservables that are correlated with  $\varepsilon_i$ , even after conditioning on  $z_{1i}$ . To model this confounding, let the treatment assignment mechanism be given by

$$x_i = I(z_i'\boldsymbol{\gamma} + u_i),$$

where the confounder  $u_i$  is correlated with  $\varepsilon_i$  and the covariate vector  $z_i = (z_{1i}, z_{2i})$  contains the variables  $z_{2i}$  that are part of the intake model but not present in the outcome model. These covariates are called instrumental variables. They must be independent of  $(\varepsilon_i, u_i)$  given  $z_{1i}$  but correlated with the intake. Now suppose for simplicity that the joint

distribution of the unobservables  $(\varepsilon_i, u_i)$  is Gaussian with mean zero and covariance matrix  $\Omega$  that is in correlation form. It is easy to see, starting with the joint distribution of  $(\varepsilon_i, u_i)$  followed by a change of variable, that this set-up is a special case of the MVP model. It can be estimated along the lines of Algorithm 5.

Chib (2003) utilizes the posterior draws to calculate the average treatment effect (ATE) which is defined as  $E(y_{i1}|z_{1i}, \boldsymbol{\gamma}_1, \beta) - E(y_{i0}|z_{1i}, \boldsymbol{\gamma}_1, \beta)$ , where  $y_{ij}$  is the outcome when  $x_i = j$ . These are the so-called potential outcomes. From (5.8) it can be seen that the two potential outcomes are  $y_{i0} = I(z'_{1i}\boldsymbol{\gamma}_1 + \varepsilon_i)$  and  $y_{i1} = I(z'_{1i}\boldsymbol{\gamma}_1 + \beta + \varepsilon_i)$  and, therefore, under the normality assumption, the difference in expected values of the potential outcomes conditioned on  $z_{1i}$  and  $\boldsymbol{\psi} = (\boldsymbol{\gamma}_1, \beta)$  is

$$E(y_{i1}|z_{1i}, \boldsymbol{\psi}) - E(y_{i0}|z_{1i}, \boldsymbol{\psi}) = \Phi(z'_{1i}\boldsymbol{\gamma}_1 + \beta) - \Phi(z'_{1i}\boldsymbol{\gamma}_1).$$

Thus, the ATE is

$$\text{ATE} = \int \{\Phi(z'_1\boldsymbol{\gamma}_1 + \beta) - \Phi(z'_1\boldsymbol{\gamma}_1)\} \pi(z_1) dz_1, \tag{5.9}$$

where  $\pi(z_1)$  is the density of  $z_1$  (independent of  $\boldsymbol{\psi}$  by assumption). Calculation of the ATE, therefore, entails an integration with respect to  $z_1$ . In practice, a simple idea is to perform the marginalization using the empirical distribution of  $z_1$  from the observed sample data. If we are interested in more than the difference in expected value of the potential outcomes, then the posterior distribution of the ATE can be obtained by evaluating the integral as a Monte Carlo average, for each sampled value of  $(\boldsymbol{\gamma}_1, \beta)$  from the MCMC algorithm.

### 6. Longitudinal binary responses

Consider now the situation in which one is interested in modeling longitudinal (univariate) binary data. Suppose that for the  $i$ th individual at time  $t$ , the probability of observing the outcome  $y_{it} = 1$ , conditioned on parameters  $\boldsymbol{\beta}_1$  and random effects  $\boldsymbol{\beta}_{2i}$ , is given by

$$\Pr(y_{it} = 1|\boldsymbol{\beta}_i) = \Phi(\mathbf{x}_{1it}\boldsymbol{\beta}_1 + \mathbf{w}_{it}\boldsymbol{\beta}_{2i}),$$

where  $\Phi$  is the cdf of the standard normal distribution,  $\mathbf{x}_{1it}$  is a  $k_1$  vector of covariates whose effect is assumed to be constant across subjects (clusters) and  $\mathbf{w}_{it}$  is an additional and distinct  $q$  vector of covariates whose effects are cluster-specific. The objective is to learn about the parameters and the random effects given  $n$  subjects each with a set of measurements  $\mathbf{y}_i = (y_{i1}, \dots, y_{in_i})'$  observed at  $n_i$  points in time. The presence of the random effects ensures that the binary responses are correlated. Modeling of the data and the estimation of the model is aided once again by the latent variable framework of Albert and Chib (1993) and the hierarchical prior modeling of Lindley and Smith (1972).

For the  $i$ th cluster, we define the vector of latent variable

$$\mathbf{z}_i = \mathbf{X}_{1i}\boldsymbol{\beta}_1 + \mathbf{W}_i\boldsymbol{\beta}_{2i} + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \mathbf{I}_{n_i}),$$



and let

$$y_{it} = I[z_{it} > 0],$$

where  $\mathbf{z}_i = (z_{i1}, \dots, z_{in_i})'$ ,  $\mathbf{X}_{1i}$  is a  $n_i \times k_1$ , and  $\mathbf{W}_i$  is  $n_i \times q$ .

In this setting, the effects  $\beta_1$  and  $\beta_{2i}$  are estimable (without any further assumptions) under the *sequential exogeneity* assumption wherein  $\varepsilon_{it}$  is uncorrelated with  $(\mathbf{x}_{1it}, \mathbf{w}_{it})$  given past values of  $(\mathbf{x}_{1it}, \mathbf{w}_{it})$  and  $\beta_i$ . In practice, even when the assumption of sequential exogeneity of the covariates  $(\mathbf{x}_{1it}, \mathbf{w}_{it})$  holds, it is quite possible that there exist covariates  $\mathbf{a}_i: r \times 1$  (with an intercept included) that are correlated with the random-coefficients  $\beta_i$ . These subject-specific covariates may be measurements on the subject at baseline (time  $t = 0$ ) or other time-invariant covariates. In the Bayesian hierarchical approach this dependence on subject-specific covariates is modeled by a hierarchical prior. One quite general way to proceed is to assume that

$$\underbrace{\begin{pmatrix} \beta_{21i} \\ \beta_{22i} \\ \vdots \\ \beta_{2qi} \end{pmatrix}}_{\beta_{2i}} = \underbrace{\begin{pmatrix} \mathbf{a}'_i & \mathbf{0}' & \dots & \dots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{a}'_i & \dots & \dots & \mathbf{0}' \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \dots & \dots & \mathbf{a}'_i \end{pmatrix}}_{\mathbf{A}_i = \mathbf{I} \otimes \mathbf{a}'_i} \underbrace{\begin{pmatrix} \beta_{21} \\ \beta_{22} \\ \vdots \\ \beta_{2q} \end{pmatrix}}_{\beta_2} + \underbrace{\begin{pmatrix} b_{i1} \\ b_{i2} \\ \vdots \\ b_{iq} \end{pmatrix}}_{\mathbf{b}_i}$$

or in vector–matrix form

$$\beta_{2i} = \mathbf{A}_i \beta_2 + \mathbf{b}_i,$$

where  $\mathbf{A}_i$  is a  $q \times k_2$  matrix,  $k = r \times q$ ,  $\beta_2 = (\beta_{21}, \beta_{22}, \dots, \beta_{2q})$  is a  $(k_2 \times 1)$ -dimensional vector, and  $\mathbf{b}_i$  is the mean zero random effects vector (uncorrelated with  $\mathbf{A}_i$  and  $\varepsilon_i$ ) that is distributed according to the distribution (say)

$$\mathbf{b}_i | \mathbf{D} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{D}).$$

This is the second stage of the model. It may be noted that the matrix  $\mathbf{A}_i$  can be the identity matrix of order  $q$  or the zero matrix of order  $q$ . In this hierarchical model, if  $\mathbf{A}_i$  is not the zero matrix then identifiability requires that the matrices  $\mathbf{X}_{1i}$  and  $\mathbf{W}_i$  have no covariates in common. For example, if the first column of  $\mathbf{W}_i$  is a vector of ones, then  $\mathbf{X}_{1i}$  cannot include an intercept. If  $\mathbf{A}_i$  is the zero matrix, however,  $\mathbf{W}_i$  is typically a subset of  $\mathbf{X}_{1i}$ . Thus, the effect of  $\mathbf{a}_i$  on  $\beta_{21i}$  (the intercept) is measured by  $\beta_{21}$ , that on  $\beta_{21i}$  is measured by  $\beta_{22}$  and that on  $\beta_{2qi}$  by  $\beta_{2q}$ .

Inserting the model of the cluster-specific random coefficients into the first stage yields

$$y_i = \mathbf{X}_i \beta + \mathbf{W}_i \mathbf{b}_i + \varepsilon_i,$$

where

$$\mathbf{X}_i = (\mathbf{X}_{1i} \ \mathbf{W}_i \mathbf{A}_i) \quad \text{with } \beta = (\beta_1 \ \beta_2),$$

as is readily checked. We can complete the model by making suitable assumptions about the distribution of  $\boldsymbol{\varepsilon}_i$ . One possibility is to assume that

$$\boldsymbol{\varepsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \mathbf{I}_{n_i})$$

which leads to the probit link model. The prior distribution on the parameters can be specified as

$$\boldsymbol{\beta} \sim \mathcal{N}_k(\boldsymbol{\beta}_0, B_0), \quad \mathbf{D}^{-1} \sim \mathcal{W}_q(\rho_0, \mathbf{R}),$$

where the latter denotes a Wishart density with degrees of freedom  $\rho_0$  and scale matrix  $\mathbf{R}$ .

The MCMC implementation in this set-up proceeds by including the  $\mathbf{z} = (z_1, \dots, z_n)$  in the sampling (Chib and Carlin, 1999). Given  $\mathbf{z}$  the sampling resembles the steps of the continuous response longitudinal model. To improve the efficiency of the sampling procedure, the sampling of the latent data is done marginalized over  $\{\mathbf{b}_i\}$  from the conditional distribution of  $z_{it} | z_{i(-t)}, y_{it}, \boldsymbol{\beta}, \mathbf{D}$ , where  $z_{i(-t)}$  is the vector  $\mathbf{z}_i$  excluding  $z_{it}$ .

ALGORITHM 6 (*Gaussian–Gaussian panel probit*).

1 Sample

(a)

$$z_{it} | z_{i(-t)}, y_{it}, \boldsymbol{\beta}, \mathbf{D} \propto \mathcal{N}(\mu_{it}, v_{it}) \{ I(z_{it} \leq 0)^{1-y_{it}} + I(z_{it} > 0)^{y_{it}} \}$$

$$\mu_{it} = E(z_{it} | z_{i(-t)}, \boldsymbol{\beta}, \mathbf{D})$$

$$v_{it} = \text{Var}(z_{it} | z_{i(-t)}, \boldsymbol{\beta}, \mathbf{D})$$

(b)

$$\boldsymbol{\beta} | \mathbf{z}, \mathbf{D} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \mathbf{X}'_i \mathbf{V}_i^{-1} \mathbf{z}_i \right), \mathbf{B} \right)$$

$$\mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \mathbf{X}'_i \mathbf{V}_i^{-1} \mathbf{X}_i \right)^{-1}; \quad \mathbf{V}_i = \mathbf{I}_{n_i} + \mathbf{W}_i \mathbf{D} \mathbf{W}'_i$$

(c)

$$\mathbf{b}_i | \mathbf{y}, \boldsymbol{\beta}, \mathbf{D} \sim \mathcal{N}_q(\mathbf{D}_i \mathbf{W}'_i (\mathbf{z}_i - \mathbf{X}_i \boldsymbol{\beta}), \mathbf{D}_i), \quad i \leq N$$

$$\mathbf{D}_i = (\mathbf{D}^{-1} + \mathbf{W}'_i \mathbf{W}_i)^{-1}$$

2 Sample

$$\mathbf{D}^{-1} | \mathbf{y}, \boldsymbol{\beta}, \{\mathbf{b}_i\} \sim \mathcal{W}_q\{\rho_0 + n, \mathbf{R}_n\}$$

where

$$\mathbf{R}_n = \left( \mathbf{R}_0^{-1} + \sum_{i=1}^n \mathbf{b}_i \mathbf{b}'_i \right)^{-1}$$

3 Goto 1

The model with errors distributed as student- $t$

$$\begin{aligned} \boldsymbol{\varepsilon}_i | \lambda_i &\sim \mathcal{N}_{n_i}(\mathbf{0}, \lambda_i \mathbf{I}_{n_i}), \\ \eta_i &\sim G\left(\frac{\nu}{2}, \frac{\nu}{2}\right) \end{aligned}$$

and random effects distributed as student- $t$

$$\begin{aligned} \mathbf{b}_i | \mathbf{D}, \boldsymbol{\eta}_i &\sim \mathcal{N}_q(\mathbf{0}, \eta_i^{-1} \mathbf{D}), \\ \eta_i &\sim G\left(\frac{\nu_b}{2}, \frac{\nu_b}{2}\right) \end{aligned}$$

can also be tackled in ways that parallel the developments in the previous section. We present the algorithm for the student-student binary response panel model without comment.

ALGORITHM 7 (*Student-Student binary panel*).

1 Sample

(a)

$$\begin{aligned} z_{it} | \mathbf{z}_{i(-t)}, y_{it}, \boldsymbol{\beta}, \mathbf{D} &\propto \mathcal{N}(\mu_{it}, v_{it}) \{ I(z_{it} \leq 0)^{1-y_{it}} + I(z_{it} > 0)^{y_{it}} \} \\ \mu_{it} &= E(z_{it} | \mathbf{z}_{i(-t)}, \boldsymbol{\beta}, \mathbf{D}, \lambda_i, \eta_i) \\ v_{it} &= \text{Var}(z_{it} | \mathbf{z}_{i(-t)}, \boldsymbol{\beta}, \mathbf{D}, \lambda_i, \eta_i) \end{aligned}$$

(b)

$$\boldsymbol{\beta} | \mathbf{z}, \mathbf{D} \{ \lambda_i \}, \{ \eta_i \} \sim \mathcal{N}_k \left( \mathbf{B}_n \left( \mathbf{B}_0^{-1} \boldsymbol{\beta}_0 + \sum_{i=1}^n \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{z}_i \right), \mathbf{B}_n \right)$$

where

$$\mathbf{B}_n = \left( \mathbf{B}_0^{-1} + \sum_{i=1}^n \mathbf{X}_i' \mathbf{V}_i^{-1} \mathbf{X}_i \right)^{-1}; \quad \mathbf{V}_i = \lambda_i^{-1} \mathbf{I}_{n_i} + \eta_i^{-1} \mathbf{W}_i \mathbf{D} \mathbf{W}_i'$$

(c)

$$\mathbf{b}_i | \mathbf{z}_i, \boldsymbol{\beta}, \mathbf{D}, \lambda_i, \eta_i \sim \mathcal{N}_q(\mathbf{D}_i \mathbf{W}_i' \lambda_i (\mathbf{z}_i - \mathbf{X}_i \boldsymbol{\beta}), \mathbf{D}_i)$$

where

$$\mathbf{D}_i = (\eta_i \mathbf{D}^{-1} + \lambda_i \mathbf{W}_i' \mathbf{W}_i)^{-1}$$

2 Sample

(a)

$$\lambda_i | \mathbf{y}, \boldsymbol{\beta}, \{ \mathbf{b}_i \} \sim \mathcal{G} \left( \frac{\nu + n_i}{2}, \frac{\nu + \mathbf{e}_i' \mathbf{e}_i}{2} \right), \quad i \leq n$$

where  $\mathbf{e}_i = (\mathbf{z}_i - \mathbf{X}_i \boldsymbol{\beta} - \mathbf{W}_i \boldsymbol{\beta}_i)$

(b)

$$\eta_i | \mathbf{b}_i, \mathbf{D} \sim \mathcal{G}\left(\frac{\nu_b + q}{2}, \frac{\nu_b + \mathbf{b}'_i \mathbf{D}^{-1} \mathbf{b}_i}{2}\right), \quad i \leq n$$

3 Sample

$$\mathbf{D}^{-1} | \mathbf{z}, \boldsymbol{\beta}, \{\mathbf{b}_i\}, \{\lambda_i\}, \{\eta_i\} \sim \mathcal{W}_q\{\rho_0 + n, \mathbf{R}_n\}$$

where

$$\mathbf{R}_n = \left( \mathbf{R}_0^{-1} + \sum_{i=1}^n \eta_i \mathbf{b}_i \mathbf{b}'_i \right)^{-1}$$

4 Goto 1

Related modeling and analysis techniques can be developed for ordinal longitudinal data. In the estimation of such models the one change is a M–H step for the sampling of the cut-points, similar to the corresponding step in [Algorithm 3](#). It is also possible in this context to let the cut-points be subject-specific; an example is provided by [Johnson \(2003\)](#).

### 6.1. Marginal likelihood of the panel binary models

The calculation of the marginal likelihood is again straightforward. For example, in the student–student binary panel model, the posterior ordinate is expressed as

$$\pi(\mathbf{D}^{-1*}, \boldsymbol{\beta}^* | \mathbf{y}) = \pi(\mathbf{D}^{-1*} | \mathbf{y}) \pi(\boldsymbol{\beta}^* | \mathbf{y}, \mathbf{D}^*),$$

where the first term is obtained by averaging the Wishart density in [Algorithm 7](#) over draws on  $\{\mathbf{b}_i\}$ ,  $\{\lambda_i\}$  and  $\mathbf{z}$  from the full run. To estimate the second ordinate, which is conditioned on  $\mathbf{D}^*$ , we run a reduced MCMC simulation with the full conditional densities

$$\begin{aligned} &\pi(\mathbf{z} | \mathbf{y}, \boldsymbol{\beta}, \mathbf{D}^*, \{\lambda_i\}); \pi(\boldsymbol{\beta} | \mathbf{z}, \mathbf{D}^*, \{\lambda_i\}, \{\eta_i\}); \\ &\pi(\{\mathbf{b}_i\} | \mathbf{z}, \boldsymbol{\beta}, \mathbf{D}^*, \{\lambda_i\}, \{\eta_i\}); \pi(\{\lambda_i\}, \{\eta_i\} | \mathbf{z}, \boldsymbol{\beta}, \mathbf{D}^*), \end{aligned}$$

where each conditional utilizes the fixed value of  $\mathbf{D}$ . The second ordinate is now estimated by averaging the Gaussian density of  $\boldsymbol{\beta}$  in [Algorithm 7](#) at  $\boldsymbol{\beta}^*$  over the draws on  $(\mathbf{z}, \{\mathbf{b}_i\}, \{\lambda_i\}, \{\eta_i\})$  from this reduced run.

## 7. Longitudinal multivariate responses

In some circumstances it is necessary to model a set of multivariate categorical outcomes in a longitudinal setting. For example, one may be interested in the factors that influence purchase into a set of product categories (for example, egg, milk, cola, etc.) when a typical consumer purchases goods at the grocery store. We may have a longitudinal sample of such consumers each observed over many different shopping occasions.

On each occasion, the consumer has a multivariate outcome vector representing the (binary) incidence into each of the product categories. The available data also includes information on various category specific characteristics, for example, the average price, display and feature values of the major brands in each category. The goal is to model the category incidence vector as a function of available covariates, controlling for the heterogeneity in covariates effects across subjects in the panel. Another example of multidimensional longitudinal data appears in Dunson (2003) motivated by studies of an item responses battery that is used to measure traits of an individual repeatedly over time. The model discussed in Dunson (2003) allows for mixtures of count, categorical, and continuous response variables.

To illustrate a canonical situation of multivariate longitudinal categorical responses, we combine the MVP model with the longitudinal models in the previous section. Let  $y_{it}$  represent a  $J$  vector of binary responses and suppose that we have for each subject  $i$  at each time  $t$ , a covariate matrix  $X_{1it}$  in the form

$$X_{1it} = \begin{pmatrix} \mathbf{x}'_{11it} & \mathbf{0}' & \cdots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{x}'_{12it} & \mathbf{0}' & \mathbf{0}' \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \cdots & \mathbf{x}'_{1Jit} \end{pmatrix},$$

where each  $\mathbf{x}_{1jit}$  has  $k_{1j}$  elements so that the dimension of  $X_{1it}$  is  $J \times k_1$ , where  $k_1 = \sum_{j=1}^J k_{1j}$ . Now suppose that there is an additional (distinct) set of covariates  $W_{it}$  whose effect is assumed to be both response and subject specific and arranged in the form

$$W_{it} = \begin{pmatrix} \mathbf{w}'_{1it} & \mathbf{0}' & \cdots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{w}'_{2it} & \mathbf{0}' & \mathbf{0}' \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \cdots & \mathbf{w}'_{Jit} \end{pmatrix},$$

where each  $w_{jit}$  has  $q_j$  elements so that the dimension of  $W_{it}$  is  $J \times q$ , where  $q = \sum_{j=1}^J q_j$ . Also suppose that in terms of the latent variables  $z_{it}$ :  $J \times 1$ , the model generating the data is given by

$$z_{it} = X_{1it}\beta_1 + W_{it}\beta_{2i} + \varepsilon_{it}, \quad \varepsilon_{it} \sim \mathcal{N}_J(\mathbf{0}, \Sigma),$$

such that

$$y_{jit} = I[z_{jit} > 0], \quad j \leq J,$$

and where

$$\beta_1 = \begin{pmatrix} \beta_{11} \\ \beta_{12} \\ \vdots \\ \beta_{1J} \end{pmatrix} : k_1 \times 1 \quad \text{and} \quad \beta_{2i} = \begin{pmatrix} \beta_{21i} \\ \beta_{22i} \\ \vdots \\ \beta_{2Ji} \end{pmatrix} : q \times 1.$$

To model the heterogeneity in effects across subjects we can assume that

$$\underbrace{\begin{pmatrix} \beta_{21i} \\ \beta_{22i} \\ \vdots \\ \beta_{2Ji} \end{pmatrix}}_{\beta_{2i}: q \times 1} = \underbrace{\begin{pmatrix} A_{1i} & \mathbf{0} & \cdots & \cdots & \mathbf{0} \\ \mathbf{0} & A_{2i} & \cdots & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \cdots & A_{Ji} \end{pmatrix}}_{A_i: q \times k_2} \underbrace{\begin{pmatrix} \beta_{21} \\ \beta_{22} \\ \vdots \\ \beta_{2J} \end{pmatrix}}_{\beta_2: k_2 \times 1} + \underbrace{\begin{pmatrix} b_{i1} \\ b_{i2} \\ \vdots \\ b_{iJ} \end{pmatrix}}_{b_i: q \times 1},$$

where each matrix  $A_{ji}$  is of the form  $\mathbf{I}_{q_j} \otimes \mathbf{a}'_{ji}$  for some  $r_j$  category-specific covariates  $\mathbf{a}_{ji}$ . The size of the vector  $\beta_2$  is then  $k_2 = \sum_{j=1}^J q_j \times r_j$ . Substituting this model of the random-coefficients into the first stage yields

$$z_{it} = X_{it}\beta + W_{it}b_i + \varepsilon_{it}, \quad \varepsilon_{it} \sim \mathcal{N}_J(\mathbf{0}, \Sigma),$$

where

$$X_{it} = (X_{1it}, W_{it}A_i)$$

and  $\beta = (\beta_1, \beta_2): k \times 1$ . Further assembling all  $n_i$  observations on the  $i$ th subject for a total of  $n_i^* = J \times n_i$  observations, we have

$$z_i = X_i\beta + W_i b_i + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}_{n_i^*}(\mathbf{0}, \Omega_i = \mathbf{I}_{n_i} \otimes \Sigma),$$

where  $z_i = (z_{i1}, \dots, z_{in_i})': n_i^* \times 1$ ,  $X_i = (X'_{i1}, \dots, X'_{in_i})': n_i^* \times k$ , and  $W_i = (W'_{i1}, \dots, W'_{in_i})': n_i^* \times q$ . Apart from the increase in the dimension of the problem, the model is now much like the longitudinal model with univariate outcomes. The fitting, consequently, parallels that in [Algorithms 6 and 7](#).

ALGORITHM 8 (*Longitudinal MVP*).

1 Sample

(a)

$$z_{it}|z_{i(-t)}, y_{it}, \beta, \Sigma, D \propto \mathcal{N}(\mu_{it}, v_{it}) \{ I(z_{it} \leq 0)^{1-y_{it}} + I(z_{it} > 0)^{y_{it}} \}$$

$$\mu_{it} = E(z_{it}|z_{i(-t)}, \beta, \Sigma, D)$$

$$v_{it} = \text{Var}(z_{it}|z_{i(-t)}, \beta, \Sigma, D)$$

(b)

$$\beta|z, \Sigma, D \sim \mathcal{N}_k \left( B_n \left( B_0^{-1} \beta_0 + \sum_{i=1}^n X'_i V_i^{-1} z_i \right), B_n \right)$$

where

$$B_n = \left( B_0^{-1} + \sum_{i=1}^n X'_i V_i^{-1} X_i \right)^{-1}; \quad V_i = \Omega_i + W_i D W'_i$$

(c)

$$\mathbf{b}_i | \mathbf{z}_i, \boldsymbol{\beta}, \boldsymbol{\Sigma}, \mathbf{D} \sim \mathcal{N}_q(\mathbf{D}_i \mathbf{W}'_i \boldsymbol{\Omega}_i^{-1} (\mathbf{z}_i - \mathbf{X}_i \boldsymbol{\beta}), \mathbf{D}_i)$$

where

$$\mathbf{D}_i = (\mathbf{D}^{-1} + \mathbf{W}'_i \boldsymbol{\Omega}_i^{-1} \mathbf{W}_i)^{-1}$$

2 Sample

$$\mathbf{D}^{-1} | \mathbf{z}, \boldsymbol{\beta}, \{\mathbf{b}_i\} \sim \mathcal{W}_q\{\rho_0 + n, \mathbf{R}_n\}$$

where

$$\mathbf{R}_n = \left( \mathbf{R}_0^{-1} + \sum_{i=1}^n \mathbf{b}_i \mathbf{b}'_i \right)^{-1}$$

3 Goto 1

## 8. Conclusion

This chapter has summarized Bayesian modeling and fitting techniques for categorical responses. The discussion has dealt with cross-section models for binary and ordinal data and presented extensions of those models to multivariate and longitudinal responses. In each case, the Bayesian MCMC fitting technique is simple and easy to implement.

The discussion in this chapter did not explore residual diagnostics and model fit issues. Relevant ideas are contained in [Albert and Chib \(1995, 1997\)](#) and [Chen and Dey \(2000\)](#). We also did not take up the multinomial probit model as, for example, considered in [McCulloch et al. \(2001\)](#), or the class of semiparametric binary response models in which the covariate effects are modeled by regression splines or other related means and discussed by, for example, [Wood et al. \(2002\)](#), [Holmes and Mallick \(2003\)](#) and [Chib and Greenberg \(2004\)](#). Fitting of all these models relies on the framework of [Albert and Chib \(1993\)](#).

In the Bayesian context it is relatively straightforward to compare alternative models. In this chapter we restricted our attention to the approach based on marginal likelihoods and Bayes factors and showed how the marginal likelihood can be computed for the various models from the output of the MCMC simulations. Taken together, the models and methods presented in this chapter are a testament to the versatility of Bayesian ideas for dealing with categorical response data.

## References

- Agresti, A. (1990). *Categorical Data Analysis*. Wiley, New York.
- Albert, J., Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *J. Amer. Statist. Assoc.* **88**, 669–679.

- Albert, J., Chib, S. (1995). Bayesian residual analysis for binary response models. *Biometrika* **82**, 747–759.
- Albert, J., Chib, S. (1997). Bayesian tests and model diagnostics in conditionally independent hierarchical models. *J. Amer. Statist. Assoc.* **92**, 916–925.
- Albert, J., Chib, S. (2001). Sequential ordinal modeling with applications to survival data. *Biometrics* **57**, 829–836.
- Allenby, G.M., Lenk, P.J. (1994). Modeling household purchase behavior with logistic normal regression. *J. Amer. Statist. Assoc.* **89**, 1218–1231.
- Amemiya, T. (1972). Bivariate probit analysis: Minimum chi-square methods. *J. Amer. Statist. Assoc.* **69**, 940–944.
- Ashford, J.R., Sowden, R.R. (1970). Multivariate probit analysis. *Biometrics* **26**, 535–546.
- Banerjee, S., Carlin, B., Gelfand, A.E. (2004). *Hierarchical Modeling and Analysis for Spatial Data*. Chapman and Hall/CRC, Boca Raton, FL.
- Basu, S., Chib, S. (2003). Marginal likelihood and Bayes factors for Dirichlet process mixture models. *J. Amer. Statist. Assoc.* **98**, 224–235.
- Basu, S., Mukhopadhyay, S. (2000). Binary response regression with normal scale mixture links. In: Dey, D.K., Ghosh, S.K., Mallick, B.K. (Eds.), *Generalized Linear Models: A Bayesian Perspective*. Marcel Dekker, New York, pp. 231–242.
- Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems (with discussion). *J. Roy. Statist. Soc., Ser. B* **36**, 192–236.
- Carey, V., Zeger, S.L., Diggle, P. (1993). Modelling multivariate binary data with alternating logistic regressions. *Biometrika* **80**, 517–526.
- Carlin, B.P., Chib, S. (1995). Bayesian model choice via Markov chain Monte Carlo methods. *J. Roy. Statist. Soc., Ser. B* **57**, 473–484.
- Chen, M.-H., Dey, D. (2000). A unified Bayesian analysis for correlated ordinal data models. *Brazilian J. Probab. Statist.* **14**, 87–111.
- Chen, M.-H., Shao, Q.-M. (1998). On Monte Carlo methods for estimating ratios of normalizing constants. *Ann. Statist.* **25**, 1563–1594.
- Chen, M.-H., Dey, D., Shao, Q.-M. (1999). A new skewed link model for dichotomous quantal response data. *J. Amer. Statist. Assoc.* **94**, 1172–1186.
- Chib, S. (1995). Marginal likelihood from the Gibbs output. *J. Amer. Statist. Assoc.* **90**, 1313–1321.
- Chib, S. (2003). On inferring effects of binary treatments with unobserved confounders (with discussion). In: Bernardo, J.M., Bayarri, M.J., Berger, J.O., Dawid, A.P., Heckerman, D., Smith, A.F.M., West, M. (Eds.), *Bayesian Statistics, vol. 7*. Oxford University Press, London.
- Chib, S., Carlin, B.P. (1999). On MCMC sampling in hierarchical longitudinal models. *Statist. Comput.* **9**, 17–26.
- Chib, S., Greenberg, E. (1995). Understanding the Metropolis–Hastings algorithm. *Amer. Statist.* **49**, 327–335.
- Chib, S., Greenberg, E. (1998). Analysis of multivariate probit models. *Biometrika* **85**, 347–361.
- Chib, S., Greenberg, E. (2004). Analysis of additive instrumental variable models. Technical Report, Washington University in Saint Louis.
- Chib, S., Jeliazkov, I. (2001). Marginal likelihood from the Metropolis–Hastings output. *J. Amer. Statist. Assoc.* **96**, 270–281.
- Connolly, M.A., Liang, K.-Y. (1988). Conditional logistic regression models for correlated binary data. *Biometrika* **75**, 501–506.
- Cox, D.R. (1972). Regression models and life tables. *J. Roy. Statist. Soc., Ser. B* **34**, 187–220.
- DiCiccio, T.J., Kass, R.E., Raftery, A.E., Wasserman, L. (1997). Computing Bayes factors by combining simulation and asymptotic approximations. *J. Amer. Statist. Assoc.* **92**, 903–915.
- Dunson, D.B. (2003). Dynamic latent trait models for multidimensional longitudinal data. *J. Amer. Statist. Assoc.* **98**, 555–563.
- Farhmeir, L., Tutz, G. (1994). *Multivariate Statistical Modelling Based on Generalized Linear Models*. Springer-Verlag, New York.
- Fahrmeir, L., Lang, S. (2001). Bayesian semiparametric regression analysis of multicategorical time-space data. *Ann. Inst. Statist. Math.* **53**, 11–30.
- Gelfand, A.E., Smith, A.F.M. (1990). Sampling-based approaches to calculating marginal densities. *J. Amer. Statist. Assoc.* **85**, 398–409.



- Geman, S., Geman, D. (1984). Stochastic relaxation, Gibbs distributions and the Bayesian restoration of images. *IEEE Trans. Pattern Analysis and Machine Intelligence* **12**, 609–628.
- Glonek, G.F.V., McCullagh, P. (1995). Multivariate logistic models. *J. Roy. Statist. Soc., Ser. B* **57**, 533–546.
- Green, P.E. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**, 711–732.
- Hastings, W.K. (1970). Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**, 97–109.
- Holmes, C.C., Mallick, B.K. (2003). Generalized nonlinear modeling with multivariate free-knot regression splines. *J. Amer. Statist. Assoc.* **98**, 352–368.
- Johnson, T.R. (2003). On the use of heterogeneous thresholds ordinal regression models to account for individual differences in response style. *Psychometrika* **68**, 563–583.
- Kalbfleish, J., Prentice, R. (1980). *The Statistical Analysis of Failure Time Data*. Wiley, New York.
- Kuhnert, P.M., Do, K.A. (2003). Fitting genetic models to twin data with binary and ordered categorical responses: A comparison of structural equation modelling and Bayesian hierarchical models. *Behavior Genetics* **33**, 441–454.
- Lesaffre, E., Kaufmann, H.K. (1992). Existence and uniqueness of the maximum likelihood estimator for a multivariate probit model. *J. Amer. Statist. Assoc.* **87**, 805–811.
- Lindley, D.V., Smith, A.F.M. (1972). Bayes estimates for the linear model. *J. Roy. Statist. Soc., Ser. B* **34**, 1–41.
- McCullagh, P. (1980). Regression models for ordinal data. *J. Roy. Statist. Soc., Ser. B* **42**, 109–127.
- McCulloch, R.E., Polson, N.G., Rossi, P.E. (2001). A Bayesian analysis of the multinomial probit model with fully identified parameters. *J. Econometrics* **99**, 173–193.
- Metropolis, N., Rosenbluth, A.W., Rosenbluth, M.N., Teller, A.H., Teller, E. (1953). Equations of state calculations by fast computing machines. *J. Chem. Phys.* **21**, 1087–1092.
- O'Brien, S.M., Dunson, D.B. (2004). Bayesian multivariate logistic regression. *Biometrics* **60**, 739–746.
- Ochi, Y., Prentice, R.L. (1984). Likelihood inference in a correlated probit regression model. *Biometrika* **71**, 531–543.
- Ripley, B. (1987). *Stochastic Simulation*. Wiley, New York.
- Roberts, C.P. (2001). *The Bayesian Choice*. Springer-Verlag, New York.
- Smith, A.F.M., Roberts, G.O. (1993). Bayesian computation via the Gibbs sampler and related Markov chain Monte Carlo methods. *J. Roy. Statist. Soc., Ser. B* **55**, 3–24.
- Tanner, M.A., Wong, W.H. (1987). The calculation of posterior distributions by data augmentation. *J. Amer. Statist. Assoc.* **82**, 528–549.
- Ten Have, T.R., Uttal, D.H. (2000). Subject-specific and population-averaged continuation ratio logit models for multiple discrete time survival profiles. *Appl. Statist.* **43**, 371–384.
- Tierney, L. (1994). Markov chains for exploring posterior distributions (with discussion). *Ann. Statist.* **22**, 1701–1762.
- Tutz, G. (1990). Sequential item response models with an ordered response. *British J. Math. Statist. Psychol.* **43**, 39–55.
- Wood, S., Kohn, R., Shively, T., Jiang, W.X. (2002). Model selection in spline nonparametric regression. *J. Roy. Statist. Soc., Ser. B* **64**, 119–139.