

# Sampling Schemes for Generalized Linear Dirichlet Process Random Effects Models

Minjung Kyung\*      Jeff Gill<sup>†</sup>      George Casella<sup>‡</sup>

April 18, 2011

## Abstract

We evaluate MCMC sampling schemes for a variety of link functions in generalized linear models with Dirichlet process random effects. First, we find that there is a large amount of variability in the performance of MCMC algorithms, with the slice sampler typically being less desirable than either a Kolmogorov-Smirnov mixture representation or a Metropolis-Hastings algorithm. Second, in fitting the Dirichlet process, dealing with the precision parameter has troubled model specifications in the past. Here we find that incorporating this parameter into the MCMC sampling scheme is not only computationally feasible, but also results in a more robust set of estimates, in that they are marginalized-over rather than conditioned-upon. Applications are provided with social science problems in areas where the data can be difficult to model, and we find that the nonparametric nature of the Dirichlet process priors for the random effects lead to improved analyses with more reasonable inferences.

*AMS 2000 subject classifications:* Primary 62F99; secondary 62P25; secondary 62G99

*Keywords and phrases:* linear mixed models, generalized linear mixed models, hierarchical models, Gibbs sampling, Metropolis-Hastings Algorithm, Slice Sampling

---

\*Assistant Professor, Center for Applied Statistics, Washington University, One Brookings Dr., Seigle Hall, St. Louis, MO. Supported by National Science Foundation Grants DMS-0631632 and SES-0631588. Email: mkyung@artsci.wustl.edu.

<sup>†</sup>Professor, Center for Applied Statistics, Washington University, One Brookings Dr., Seigle Hall, St. Louis, MO. Supported by National Science Foundation Grants DMS-0631632 and SES-0631588. Email: jgill@wustl.edu.

<sup>‡</sup>Distinguished Professor, Department of Statistics, University of Florida, Gainesville, FL 32611. Supported by National Science Foundation Grants DMS-04-05543, DMS-0631632 and SES-0631588. Email: casella@stat.ufl.edu.

# 1 Introduction

Generalized linear models (GLMs) have enjoyed considerable attention over the years, providing a flexible framework for modeling discrete responses using a variety of error structures. If we have observations that are discrete or categorical,  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ , such data can often be assumed to be independent and from a distribution in the exponential family. The classic book by McCullagh and Nelder (1989) describes these models in detail; see also the more recent developments in Dey, Ghosh, and Mallick (2000) or Fahrmeir and Tutz (2001).

A generalized linear *mixed* model (GLMM) is an extension of a GLM that allows random effects, and can give us flexibility in developing a more suitable model when the observations are correlated, or where there may be other underlying phenomena that contribute to the resulting variability. Thus, the GLMM can be specified to accommodate outcome variables conditional on mixtures of possibly correlated random and fixed effects (Breslow and Clayton 1993, Buonaccorsi 1996, Wang, *et al.* 1998, Wolfinger and O'Connell, 1993). Details of such models, covering both statistical inferences and computational methods, can be found in the texts by McCulloch and Searle (2001) and Jiang (2007).

## 1.1 Sampling Schemes for GLMMs

There have been Markov chain Monte Carlo (MCMC) methods developed for the analysis of the GLMMs with random effects modeled with a normal distribution. Although the posteriors of parameters and the random effects are typically numerically intractable, especially when the dimension of the random effects is greater than one, there has been much progress in the development of sampling schemes. For example, Damien *et al.* (1999) proposed a Gibbs sampler using auxiliary variables for sampling non-conjugate and hierarchical models. Their methods are *slice sampling* methods derived from the full conditional posterior distribution. They mention that the assessment of convergence remains a major problem with the algorithm. However, Neal (2003) provided convergence properties of the posterior for slice sampling. Another sampling scheme was used by Chib *et al.* (1998) and Chib and Winkelmann (2001), who provided the Metropolis-Hastings (M-H) algorithms for various kinds of GLMMs. They proposed a multivariate-*t* distribution as a candidate density in an M-H implementation, taking the mean equal to the posterior mode, and variance equal to the inverse of the Hessian evaluated at the posterior mode.

To be precise about language, we discuss three types of MCMC algorithms in this work. When we refer to the *slice sampler* we mean a Gibbs sampler on an enlarged state space (augmented by auxiliary variables). When we refer to a *Gibbs sampler*, it is a sampler based on producing automatically accepted candidate values from full conditional distributions that is not the special case of the slice sampler. When *Metropolis-Hastings* algorithms are discussed, these are not the

special cases of Gibbs or slice sampling, but instead the more general process of producing candidate values from a separate distribution and deciding to accept them or not using the conventional Metropolis step.

## 1.2 Sampling Schemes for GLMDMs

Another variation of a GLMM was used by Dorazio, *et al.* (2007) and Gill and Casella (2009), where the random effects are modeled with a Dirichlet process, resulting in a Generalized Linear Mixed Dirichlet Process Model (GLMDM). Dorazio, *et al.* (2007) used a GLMDM with a log link for spatial heterogeneity in animal abundance. They proposed an empirical Bayesian approach with the Dirichlet process, instead of the regular assumption of normally distributed random effects, because they argued that for some species, the sources of heterogeneity in abundance is poorly understood or unobservable. They noted that the Dirichlet process prior is robust to errors in model specification and allows spatial heterogeneity in abundance to be specified in a data-adaptive way. Gill and Casella (2009) suggested a GLMDM with an ordered probit link to model political science data, specifically modeling the stress, from public service, of Senate-confirmed political appointees as a reason for their short tenure. For the analysis, a semi-parametric Bayesian approach was adopted, using the Dirichlet process for the random effect.

Dirichlet process mixture models were introduced by Ferguson (1973) and Antoniak (1974), with further important developments in Blackwell and MacQueen (1973), Korwar and Hollander (1973), and Sethuraman (1994). For estimation, Lo (1984) derived the analytic form of a Bayesian density estimator, and Liu (1996) derived an identity for the profile likelihood estimator of the Dirichlet precision parameter. Kyung, *et al.* (2010) looked at the properties of this MLE and found that the likelihood function can be ill-behaved. They noted that incorporating a gamma prior, and using posterior mode estimation, results in a more stable solution. McAuliffe, Blei and Jordan (2006) used a similar strategy, using a posterior mean for the estimation of the Dirichlet process precision parameter (the term  $m$ , which we describe in Section 2).

Models with Dirichlet process priors are treated as hierarchical models in a Bayesian framework, and the implementation of these models through Bayesian computation and efficient algorithms has had much attention. Escobar and West (1995) provided a Gibbs sampling algorithm for the estimation of posterior distribution for all model parameters, MacEachern and Muller (1998) presented a Gibbs sampler with non-conjugate priors by using auxiliary parameters, and Neal (2000) provided an extended and more efficient Gibbs sampler to handle general Dirichlet process mixture models. Teh *et al.* (2006) also extended the auxiliary variable method of Escobar and West (1995) for posterior sampling of the precision parameter with a gamma prior. They

developed hierarchical Dirichlet processes, with a Dirichlet prior for the base measure.

Kyung, *et al.* (2010) developed algorithms for estimation of the precision parameter and new MCMC algorithms for a linear mixed Dirichlet process random effects models that had not previously existed. In addition, they showed how to extend the developed framework to a generalized Dirichlet process mixed model with a probit link function. They derived, for the first time, a simultaneous Gibbs sampler for all of the model parameters and the subclusters of the Dirichlet process, and used a new parameterization of the hierarchical model to derive a Gibbs sampler that more fully exploits the structure of the model that mixes very well. Finally they were also able to establish a proof that the proposed sampler is an improvement, in terms of operator norm and efficiency, over other commonly used algorithms.

### 1.3 Summary

In this paper we look at MCMC sampling schemes for the generalized Dirichlet process mixture models, concentrating on logistic and log linear models. For these models, we examine a Gibbs sampling method using auxiliary parameters, based on Damien *et al.* (1999), and a Metropolis-Hastings sampler where the candidate generating distribution is a Gaussian density from log-transformed count data from a log-linear model (thus producing a form on the correct support). We incorporate the Dirichlet process precision parameter,  $m$ , into the Gibbs sampler, through the use of a gamma candidate distribution using a Laplace approximation for the calculation of the mean and variance of  $m$ , and use that in the gamma candidate. In the examples analyzed here, we find that, we find that the alternative slice sampler typically has higher autocorrelation in logistic regression and loglinear models than the proposed M-H algorithm.

Using the GLMDM with a general link function, Section 2 describes the generalized Dirichlet process mixture model. In Section 3 we estimate model parameters using a variety of algorithms, and Section 4 describes the estimation of the Dirichlet parameters. Section 5 looks at the performance of these algorithms in a variety of simulations, while Section 6 analyzes two social science data sets, further illustrating the advantage of the Dirichlet process random effects model. Section 7 summarizes these contributions and adds some perspective, and there is an Appendix with some technical details.

## 2 A Generalized Linear Mixed Dirichlet Process Model

Let  $\mathbf{X}_i$  be covariates associated with the  $i^{\text{th}}$  observation,  $\boldsymbol{\beta}$  be the coefficient vector, and  $\psi_i$  be a random effect accounting for subject-specific deviation from the underlying model. Assume that the  $Y_i|\boldsymbol{\psi}$  are conditionally independent, each with a density from the exponential family, where

$\boldsymbol{\psi} = (\psi_1, \dots, \psi_n)$ . Then, based on the notation on McCulloch and Searle (2001), the GLMDM can be expressed as follows. Start with the generalized linear model,

$$\begin{aligned} Y_i | \gamma &\stackrel{\text{ind}}{\sim} f_{Y_i | \gamma}(y_i | \gamma), \quad i = 1, \dots, n \\ f_{Y_i | \gamma}(y_i | \gamma) &= \exp [\{y_i \gamma_i - b(\gamma_i)\} / \xi^2 - c(y_i, \xi)]. \end{aligned} \quad (1)$$

where  $y_i$  is discrete valued. Here, we know that  $E[Y_i | \gamma] = \mu_i = \partial b(\gamma_i) / \partial \gamma_i$ . Using a link function  $g(\cdot)$ , we can express the transformed mean of  $Y_i$ ,  $E[Y_i | \gamma]$ , as a linear function, and we add a random effect to create the mixed model:

$$g(\mu_i) = \mathbf{X}_i \boldsymbol{\beta} + \psi_i. \quad (2)$$

Here, for the Dirichlet process mixture models, we assume that

$$\begin{aligned} \psi_i &\sim G \\ G &\sim \mathcal{DP}(mG_0), \end{aligned} \quad (3)$$

where  $\mathcal{DP}$  is the Dirichlet process with base measure  $G_0$  and precision parameter  $m$ . Blackwell and MacQueen (1973) proved that for  $\psi_1, \dots, \psi_n$  iid from  $G \sim \mathcal{DP}$ , the joint distribution of  $\boldsymbol{\psi}$  is a product of successive conditional distributions of the form:

$$\psi_i | \psi_1, \dots, \psi_{i-1}, m \sim \frac{m}{i-1+m} g_0(\psi_i) + \frac{1}{i-1+m} \sum_{l=1}^{i-1} \delta(\psi_l = \psi_i) \quad (4)$$

where  $\delta(\cdot)$  denotes the Dirac delta function and  $g_0(\cdot)$  is the density function of the base measure.

We define a *partition*  $C$  to be a clustering of the sample of size  $n$  into  $k$  groups,  $k = 1, \dots, n$ , and we call these subclusters since the grouping is done nonparametrically rather than on substantive criteria. That is, the partition assigns different distributional parameters across groups and the same parameters within groups; cases are iid only if they are assigned to the same subcluster.

Applying Lo (1984) Lemma 2 and Liu (1996) Theorem 1 to (4), we can calculate the likelihood function, which by definition is integrated over the random effects, as

$$L(\boldsymbol{\beta} | \mathbf{y}) = \frac{\Gamma(m)}{\Gamma(m+n)} \sum_{k=1}^n m^k \sum_{C:|C|=k} \prod_{j=1}^k \Gamma(n_j) \int f(\mathbf{y}_{(j)} | \boldsymbol{\beta}, \psi_j) dG_0(\psi_j),$$

where  $C$  defines the partition of subclusters of size  $n_j$ ,  $|C|$  indicates occupied subclusters,  $\mathbf{y}_{(j)}$  is the vector of  $y_i$ s that are in subcluster  $j$ , and  $\psi_j$  is the common parameter for that subcluster. There are  $\mathcal{S}_{n,k}$  different partitions  $C$ , the Stirling Number of the Second Kind (Abramowitz and Stegun 1972, 824-825).

Here, we consider an  $n \times k$  matrix  $A$  defined by

$$A = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

where each  $a_i$  is a  $1 \times k$  vector of all zeros except for a 1 in the position indicating which group the observation is from. Thus,  $A$  represents a partition of the sample of size  $n$  into  $k$  groups, with the column sums giving the subcluster sizes. Note that both the dimension  $k$ , and the placement of the 1s, are random, representing the subclustering process.

If the partition  $C$  has subclusters  $\{S_1, \dots, S_k\}$ , then if  $i \in S_j$ ,  $\psi_i = \eta_j$  and the random effect can be rewritten as

$$\boldsymbol{\psi} = A\boldsymbol{\eta}, \quad (5)$$

where  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_k)$  and  $\eta_j \stackrel{iid}{\sim} G_0$  for  $j = 1, \dots, k$ . This is the same representation of the Dirichlet process that was used in Kyung *et al.* (2010), building on the representation in McCullagh and Yang (2006).

In this paper, we consider models for the binary responses with probit and logit link function, and for count data with a log link function. First, for the binary responses,

$$Y_i \sim \text{Bernoulli}(p_i), \quad i = 1, \dots, n$$

where  $y_i$  is 1 or 0, and  $p_i = E(Y_i)$  is the probability of a success for the  $i^{\text{th}}$  observation. Using a general link function (2) leads to a sampling distribution for the observed outcome variable  $\mathbf{y}$ :

$$f(\mathbf{y}|A) = \int \prod_{i=1}^n [g^{-1}(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]^{y_i} [1 - g^{-1}(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]^{1-y_i} dG_0(\boldsymbol{\eta}),$$

which typically can only be evaluated numerically. Examples of general link functions for binary outcomes are

$$\begin{aligned} p_i &= g_1^{-1}(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i) = \Phi(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i) && \text{Probit} \\ p_i &= g_2^{-1}(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i) = (1 + \exp(-\mathbf{X}_i\boldsymbol{\beta} - (\mathbf{A}\boldsymbol{\eta})_i))^{-1} && \text{Logistic} \\ p_i &= g_3^{-1}(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i) = 1 - \exp(-\exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)) && \text{Cloglog} \end{aligned}$$

where  $\Phi()$  is the cumulative distribution function of a standard normal distribution.

For counting process data,

$$Y_i \sim \text{Poisson}(\lambda_i), \quad i = 1, \dots, n$$

where  $y_i$  is  $0, 1, \dots$ ,  $\lambda_i = E(Y_i)$  is the expected number of events for the  $i^{\text{th}}$  observation. Here, using a log link function

$$\log(\lambda_i) = \mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i,$$

the sampling distribution of  $\mathbf{y}$  is

$$f(\mathbf{y}|A) = \prod_{i=1}^n \frac{1}{y_i!} \int \prod_{i=1}^n \exp\{-\exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)\} [\exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]^{y_i} G_0(\boldsymbol{\eta}) d\boldsymbol{\eta}.$$

For the base measure of the Dirichlet process, we assume a normal distribution with mean 0 and variance  $\tau^2$ ,  $N(0, \tau^2)$ . In our experience, the model is not sensitive to this distributional assumption and others, such as the students- $t$ , could be used.

### 3 Sampling Schemes for the Model Parameters

An overview of the general sampling scheme is as follows. We have three groups of parameters:

- (i)  $m$ , the precision parameter of the Dirichlet process,
- (ii)  $\mathbf{A}$ , the indicator matrix of the partition defining the subclusters, and
- (iii)  $(\boldsymbol{\eta}, \boldsymbol{\beta}, \tau^2)$ , the model parameters.

We iterate between these three groups until convergence:

1. Conditional on  $m$  and  $A$ , generate  $(\boldsymbol{\eta}, \boldsymbol{\beta}, \tau^2)|\mathbf{A}, m$ ;
2. Conditional on  $(\boldsymbol{\eta}, \boldsymbol{\beta}, \tau^2)$  and  $m$ , generate  $A$ , a new partition matrix.
3. Conditional on  $(\boldsymbol{\eta}, \boldsymbol{\beta}, \tau^2)$  and  $A$ , generate  $m$ , the new precision parameter.

For the model parameters we add the priors

$$\begin{aligned} \boldsymbol{\beta}|\sigma^2 &\sim N(\mathbf{0}, d^* \sigma^2 I) \\ \tau^2 &\sim \text{Inverted Gamma}(a, b), \end{aligned} \tag{6}$$

where  $d^* > 1$  and  $(a, b)$  are fixed such that the inverse gamma is diffuse ( $a = 1$ ,  $b$  very small). Thus the partitioning in the algorithm assigns different normal parameters across groups and the same normal parameters within groups. For the Dirichlet process we need the previously stated priors

$$\boldsymbol{\eta} = (\eta_1, \dots, \eta_k) \text{ and } \eta_j \stackrel{iid}{\sim} G_0 \text{ for } j = 1, \dots, k. \tag{7}$$

We can either fix  $\sigma^2$  or put a prior on it and estimate it in the hierarchical model with priors; here we will fix a value for  $\sigma^2$ .

In the following sections we consider a number of sampling schemes for the estimation of the model parameters of a GLMDM. We will then turn to generation of the subclusters and the precision parameter.

### 3.1 Probit Models

Albert and Chib (1993) showed how truncated normal sampling could be used to implement the Gibbs sampler for a probit model for binary responses. They use a latent variable  $V_i$  such that

$$V_i = X_i\boldsymbol{\beta} + \psi_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad (8)$$

and

$$y_i = 1 \quad \text{if } V_i > 0 \quad \text{and} \quad y_i = 0 \quad \text{if } V_i \leq 0$$

for  $i = 1, \dots, n$ . It can be shown that  $Y_i$  are independent Bernoulli random variables with the probability of success,  $p_i = \Phi((X_i\boldsymbol{\beta} - \psi_i)/\sigma)$ , and without loss of generality, we fix  $\sigma = 1$ .

Details of implementing the Dirichlet process random effect probit model are given in Kyung *et al.* (2010) and will not be repeated here. We will use this model for comparison, but our main interest is in logistic and loglinear models.

### 3.2 Logistic Models

We look at two samplers for the logistic model. The first is based on the slice sampler of Damien *et al.* (1999), while the second exploits a mixture representation of the logistic distribution; see Andrews and Mallows (1974) or West (1987).

#### 3.2.1 Slice Sampling

The idea behind the slice sampler is the following. Suppose that the density  $f(\theta) \propto L(\theta)\pi(\theta)$ , where  $L(\theta)$  is the likelihood and  $\pi(\theta)$  is the prior, and it is not possible to sample directly from  $f(\theta)$ . Using a latent variable  $U$ , define the joint density of  $\theta$  and  $U$  by

$$f(\theta, u) \propto I\{u < L(\theta)\} \pi(\theta).$$

Then,  $U|\theta$  is uniform  $\mathcal{U}\{0, L(\theta)\}$ , and  $\theta|U = u$  is  $\pi$  restricted to the set  $A_u = \{\theta : L(\theta) > u\}$ .



The likelihood function of binary responses with logit link function can be written as

$$L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}|A, \mathbf{y}) = \prod_{i=1}^n \left[ \frac{1}{1 + \exp(-\mathbf{X}_i\boldsymbol{\beta} - (\mathbf{A}\boldsymbol{\eta})_i)} \right]^{y_i} \left[ \frac{1}{1 + \exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)} \right]^{1-y_i} \\ \times \prod_{j=1}^k \left( \frac{1}{2\pi\tau^2} \right)^{1/2} \exp \left( -\frac{1}{2\tau^2}\eta_j^2 \right), \quad (9)$$

and if we introduce latent variables  $\mathbf{U} = (U_1, \dots, U_n)$  and  $\mathbf{V} = (V_1, \dots, V_n)$ , we have the likelihood of the model parameters and the latent variables to be

$$L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}|A, \mathbf{y}) \quad (10) \\ = \prod_{i=1}^n I \left[ u_i < \left\{ \frac{1}{1 + \exp(-\mathbf{X}_i\boldsymbol{\beta} - (\mathbf{A}\boldsymbol{\eta})_i)} \right\}^{y_i}, v_i < \left\{ \frac{1}{1 + \exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)} \right\}^{1-y_i} \right] \\ \times \prod_{j=1}^k \left( \frac{1}{2\pi\tau^2} \right)^{1/2} \exp \left( -\frac{1}{2\tau^2}\eta_j^2 \right)$$

Thus, with priors that are given above, the joint posterior distribution of  $(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V})$  can be expressed as

$$\pi_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}|A, \mathbf{y}) \propto L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}|A, \mathbf{y}) \quad (11) \\ \times \left( \frac{1}{\tau^2} \right)^{a+1} \exp \left( -\frac{b}{\tau^2} \right) \exp \left( -\frac{|\boldsymbol{\beta}|^2}{2d^*\sigma^2} \right).$$

Then for fixed  $m$  and  $A$ , we can implement a Gibbs sampler using the full conditionals. Details are discussed in Appendix A.1.

### 3.2.2 A Mixture Representation

Next we consider a Gibbs sampler using truncated normal variables in a manner that is similar to the Gibbs sampler of the probit models, which arise from a mixture representation of the logistic distribution. Andrews and Mallows (1974) discussed necessary and sufficient conditions under which a random variable  $Y$  may be generated as the ratio  $Z/V$  where  $Z$  and  $V$  are independent and  $Z$  has a standard normal distribution, and establish that when  $V/2$  has the asymptotic distribution of the Kolmogorov distance statistic,  $Y$  is logistic. West (1987) generalized this result to the exponential power family of distributions, showing these distributional forms to be a subset of the class of scale mixtures of normals. The corresponding mixing distribution is explicitly obtained, identifying a close relationship between the exponential power family and a further class of normal scale mixtures, the stable distributions.

Based on Andrews and Mallows (1974), and West (1987), the logistic distribution is a scale mixture of a normal distribution with a Kolmogorov-Smirnov distribution. From Devroye (1986), the Kolmogorov-Smirnov (K-S) density function is given by

$$f_X(x) = 8 \sum_{\alpha=1}^{\infty} (-1)^{\alpha+1} \alpha^2 x e^{-2\alpha^2 x^2} \quad x \geq 0, \quad (12)$$

and we define the joint distribution

$$f_{Y,X}(y, x) = (2\pi)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left( \frac{y}{2x} \right)^2 \right\} f_X(x) \frac{1}{2x}. \quad (13)$$

From the identities in Andrews and Mallows (1974) (see also Theorem 10.2.1 in Balakrishnan 1992), the marginal distribution of  $Y$  is then given by

$$f_Y(y) = \int_0^{\infty} f_{Y,X}(y, x) dx = \sum_{\alpha=1}^{\infty} (-1)^{\alpha+1} \alpha \exp(-\alpha|y|) = \frac{e^{-y}}{(1 + e^{-y})^2}, \quad (14)$$

the density function of logistic distribution with mean 0 and variance  $\frac{\pi^2}{3}$ . Therefore,  $Y \sim \Lambda \left( 0, \frac{\pi^2}{3} \right)$ , where  $\Lambda(\cdot)$  is the logistic distribution.

Now, using the likelihood function of binary responses with logit link function (9), consider the latent variable  $W_i$  such that

$$W_i = \mathbf{X}_i \boldsymbol{\beta} + \psi_i + \boldsymbol{\eta}_i, \quad \boldsymbol{\eta}_i \sim \Lambda \left( 0, \frac{\pi^2}{3} \sigma^2 \right), \quad (15)$$

with  $y_i = 1$  if  $W_i > 0$  and  $y_i = 0$  if  $W_i \leq 0$ , for  $i = 1, \dots, n$ . It can be shown that  $Y_i$  are independent Bernoulli random variables with  $p_i = [1 + \exp(-\mathbf{X}_i \boldsymbol{\beta} - (\mathbf{A} \boldsymbol{\eta})_i)]^{-1}$ , the probability of success, and without loss of generality we fix  $\sigma = 1$ .

For given  $A$ , the likelihood function of model parameters and the latent variable is given by

$$\begin{aligned} L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U} | A, \mathbf{y}, \sigma^2) &= \prod_{i=1}^n \{ I(U_i > 0) I(y_i = 1) + I(U_i \leq 0) I(y_i = 0) \} \\ &\quad \times \int_0^{\infty} \left( \frac{1}{2\pi\sigma^2(2\xi)^2} \right)^{n/2} e^{-\frac{1}{2\sigma^2(2\xi)^2} |\mathbf{U} - \mathbf{X}\boldsymbol{\beta} - A\boldsymbol{\eta}|^2} \\ &\quad \times 8 \sum_{\alpha=1}^{\infty} (-1)^{\alpha+1} \alpha^2 \xi e^{-2\alpha^2 \xi^2} d\xi \left( \frac{1}{2\pi\tau^2} \right)^{k/2} e^{-\frac{1}{2\tau^2} |\boldsymbol{\eta}|^2}, \end{aligned}$$

where  $\mathbf{U} = (U_1, \dots, U_n)$ , and  $U_i$  is the truncated normal variable which is described in (8).

Let  $m$  and  $A$  be considered fixed for the moment. Thus, with priors given in (6) and (7), the joint posterior distribution of  $(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U})$  given the outcome  $\mathbf{y}$  is

$$\pi_k^L \propto \int_0^{\infty} L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U} | A, \mathbf{y}, \sigma^2) e^{-\frac{1}{2d^* \sigma^2} |\boldsymbol{\beta}|^2} \left( \frac{1}{\tau^2} \right)^{a+1} e^{-\frac{b}{\tau^2}} d\boldsymbol{\eta}.$$

This representation avoids the problem of generating samples from the truncated logistic distribution, which is not easy to implement. As we now have the logistic distribution expressed as a normal mixture with the K-S distribution, we now only need to generate samples from the truncated normal distribution and the K-S distribution, and we can get a Gibbs sampler for the model parameters. The details are left to Appendix A.1.2.

### 3.3 Log Linear Models

Similar to Section 3.2, we look at two samplers for the loglinear model. The first is again based on the slice sampler of Damien *et al.* (1999), while the second is an M-H algorithm based on using a Gaussian density from log-transformed data as a candidate.

#### 3.3.1 Slice Sampling

The likelihood function of the counting process data with log link function can be written as

$$L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta} | A, \mathbf{y}) = \prod_{i=1}^n \frac{1}{y_i!} e^{-\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)} [\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)]^{y_i} \quad (16)$$

$$\times \prod_{j=1}^k \left( \frac{1}{2\pi\tau^2} \right)^{1/2} \exp \left( -\frac{1}{2\tau^2} \eta_j^2 \right),$$

and the joint posterior distribution of  $(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta})$  can be obtained by appending the priors for  $\tau^2$  and  $\boldsymbol{\beta}$ . As in Section 3.2.1 we introduce latent variables  $\mathbf{U} = (U_1, \dots, U_n)$  and  $\mathbf{V} = (V_1, \dots, V_n)$ , yielding a likelihood of the model parameters and the latent variables,  $L_k(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V} | A, \mathbf{y})$ , similar to (10). Setting up the Gibbs sampler is now straightforward, with details in Appendix A.2.1.

#### 3.3.2 Metropolis-Hastings

The primary challenge in setting up an efficient Metropolis-Hastings algorithm is specifying practical candidate generating functions for each of the unknown parameters in the sampler. This involves both stipulating a close distributional form to the target *and* variances that provide a reasonable acceptance rate. Starting with the likelihood and priors described at (16), for the candidate distribution of  $\boldsymbol{\beta}$  and  $\boldsymbol{\eta}$ , we consider the model:

$$\log(Y_i) = \mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma^2).$$

which is a linear mixed Dirichlet process model (LMDPM). Sampling these model parameters is straightforward, and this enables us to have high-quality candidate values for the accept/reject

stage of the Metropolis-Hastings algorithm for the log linear setup here. Using a similar model with the same parameter support but different link function as a way to generate M-H candidate values is a standard trick in the MCMC literature (Robert and Casella 2004). Details about this process are provided in Appendix A.2.2.

### 3.3.3 Comparing Slice Sampling to Metropolis-Hastings

In a special case it is possible to directly compare slice sampling and independent Metropolis-Hastings. If we have a Metropolis-Hastings algorithm with target density  $\pi$  and candidate  $h$ , we can compare it to the slice sampler

$$U|X = x \sim \text{Uniform}\{u : 0 < u < \pi(x)/h(x)\}, \quad X|U = u \sim h(x)\{x : 0 < u < \pi(x)/h(x)\}.$$

In this setup Mira and Tierney (2002) show that the slice sampler dominates the Metropolis-Hastings algorithm in the efficiency ordering, meaning that all asymptotic variances are smaller, as well as first-order covariances.

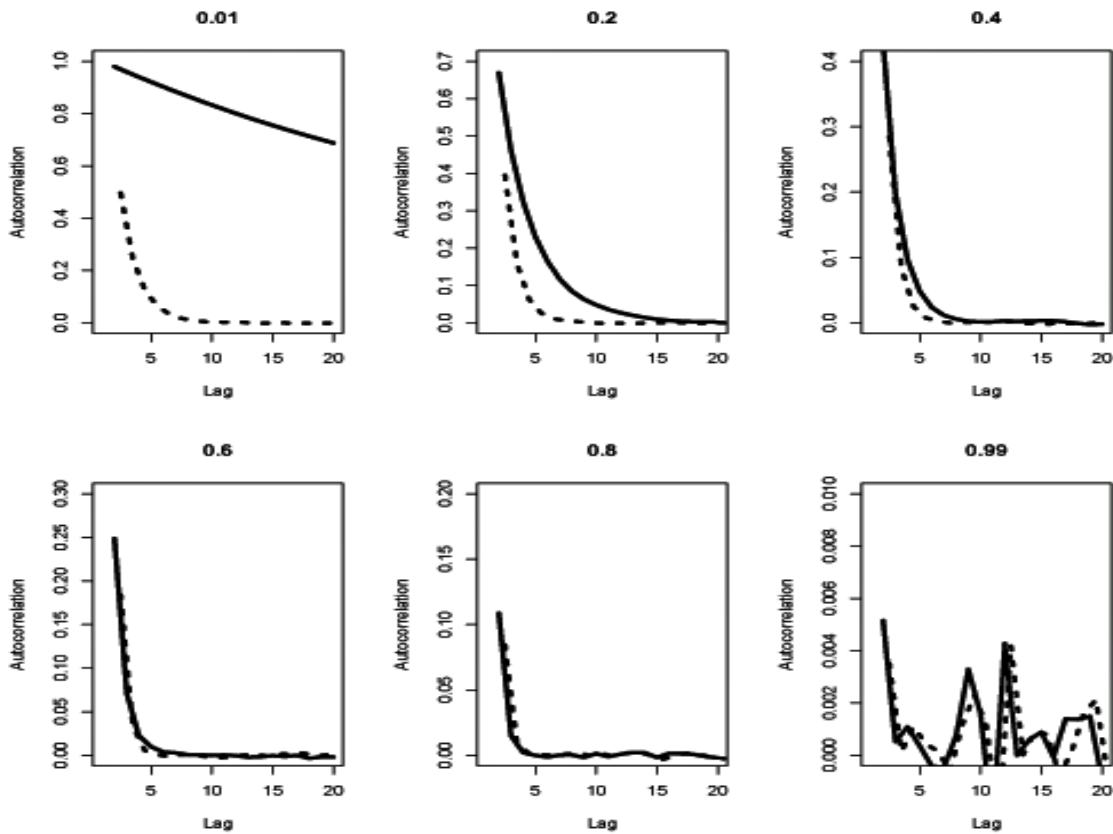
At first look this result seems to be in opposition with what we will see in Section 5; we find that Metropolis-Hastings outperforms slice sampling with respect to autocorrelations. The resolution of this discrepancy is simple; the Mira-Tierney result applies when slice sampling and Metropolis-Hastings have the relationship described above - the candidate densities must be the same. In practice, and in the examples that we will see, the candidates are chosen in each case based on ease of computation, and in the case of the Metropolis-Hastings algorithm, to try to mimic the target. Under the demanding circumstances required of our Metropolis-Hastings algorithm for the real-world data and varied link functions used, it would be a very difficult task to produce candidate generating distributions that might match a slice sampler.

As an illustration of where can actually match candidate generating distributions, consider the parameterization of Mira and Tierney (2002), where

$$\pi(x) = e^{-x} \text{ and } h(x) = qe^{-qx}, \quad 0 < q < 1. \tag{17}$$

If both slice and Metropolis-Hastings use the same value of  $q$ , then the slice sampler dominates. But if the samplers use different values of  $q$ , it can be the case that Metropolis-Hastings dominates the slice sampler. This is illustrated in Figure 1, where we show the autocorrelations for both the slice sampler and the Metropolis-Hastings algorithm, for different values of  $q$ . Compare Metropolis-Hastings with large values of  $q$ , where the candidate gets closer to the target, with a slice sampler having a smaller value of  $q$ . (Note that the different plots have different scales.) We see that in these cases the Metropolis-Hastings algorithm can dominate the slice sampler.

Figure 1: Autocorrelations for both the slice sampler (dashed) and the Metropolis-Hastings algorithm (solid), for different values of  $q$ , for the model in (17). Note that the panels have different scales.



## 4 Sampling Schemes for the Dirichlet Process Parameters

### 4.1 Generating the Partitions

We use a Metropolis-Hastings algorithm with a candidate taken from a multinomial/Dirichlet. This produces a Gibbs sampler that converges faster than the popular “stick-breaking” algorithm of Ishwaran and James (2001) See Kyung *et al.* 2010 for details on comparing stick-breaking versus “restaurant” algorithms.

For  $t = 1, \dots, T$ , at iteration  $t$

1. Starting from  $(\boldsymbol{\theta}^{(t)}, \mathbf{A}^{(t)})$ ,

$$\boldsymbol{\theta}^{(t+1)} \sim \pi(\boldsymbol{\theta} \mid \mathbf{A}^{(t)}, \mathbf{y}),$$

where  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta})$  and the updating methods are discussed above.

2. If  $\mathbf{q} = (q_1, \dots, q_n) \sim \text{Dirichlet}(r_1, \dots, r_n)$ , then for any  $k$  and  $k + 1 \leq n$

$$\mathbf{q}^{(t+1)} = (q_1^{(t+1)}, \dots, q_n^{(t+1)}) \sim \text{Dirichlet}(n_1^{(t)} + r_1, \dots, n_k^{(t)} + r_k, r_{k+1}, \dots, r_n) \quad (18)$$

3. Given  $\boldsymbol{\theta}^{(t+1)}$ ,

$$\mathbf{A}^{(t+1)} \sim P(\mathbf{A})f(\mathbf{y}|\boldsymbol{\theta}^{(t+1)}, \mathbf{A}) \binom{n}{n_1 \dots n_k} \prod_{j=1}^k [q_j^{(t+1)}]^{n_j} \quad (19)$$

where  $\mathbf{A}$  is  $n \times k$  with column sums  $n_j > 0$ ,  $n_1 + \dots + n_k = n$ .

Based on the value of the  $q_j^{(t+1)}$  in (18) we generate a candidate  $\mathbf{A}'$  that is an  $n \times n$  matrix where each row is a multinomial, and the effective dimension of the matrix, the size of the partition,  $k$ , are the non-zero column sums. Deleting the columns with column sum zero is a marginalization of the multinomial distribution. The probability of the candidate is given by

$$P(\mathbf{A}^{(t+1)}) = \frac{\Gamma(\sum_{j=1}^n r_j)}{\prod_{j=1}^{k^{(t+1)}-1} \Gamma(r_j) \Gamma(\sum_{j=k^{(t+1)}}^n r_j)} \frac{\prod_{j=1}^{k^{(t+1)}-1} \Gamma(n_j^{(t+1)} + r_j) \Gamma(n_{k^{(t+1)}}^{(t+1)} + \sum_{j=k^{(t+1)}}^n r_j)}{\Gamma(n + \sum_{j=1}^n r_j)}$$

and a Metropolis-Hastings step is then done.

## 4.2 Gibbs Sampling the Precision Parameter

To estimate the precision parameter of the Dirichlet process,  $m$ , we start with the profile likelihood,

$$L(m | \boldsymbol{\theta}, \mathbf{A}, \mathbf{y}) = \frac{\Gamma(m)}{\Gamma(m+n)} m^k \prod_{j=1}^k \Gamma(n_j) f(\mathbf{y} | \boldsymbol{\theta}, \mathbf{A}). \quad (20)$$

Rather than estimating  $m$ , a better strategy is to include  $m$  directly in the Gibbs sampler, as the maximum likelihood estimate from (20) can be very unstable (Kyung, *et al.* 2010). Using the prior  $g(m)$  we get the posterior density

$$\pi(m | \boldsymbol{\theta}, \mathbf{A}, \mathbf{y}) = \frac{\frac{\Gamma(m)}{\Gamma(m+n)} g(m) m^k}{\int_0^\infty \frac{\Gamma(m)}{\Gamma(m+n)} g(m) m^k dm}, \quad (21)$$

where  $\int \pi(m | \boldsymbol{\theta}, \mathbf{A}, \mathbf{y}) dm < \infty$  must be finite for this to be proper. Note also how far removed  $m$  is from the data, as the posterior only depends on the number of groups  $k$ . We consider a gamma distribution as a prior,  $g(m) = m^{a-1} e^{-m/b} / \Gamma(a) b^a$ , and generate  $m$  using an M-H algorithm with another gamma density as a candidate.

We choose the gamma candidate by using an approximate mean and variance of  $\pi(m)$  to set the parameters of the candidate. To get the approximate mean and variance, we will use the Laplace approximation of Tierney and Kadane (1986). Applying their results and using the log-likelihood,  $\ell()$  in place of the likelihood,  $L()$ , we have:

$$\frac{\int m^\nu \frac{\Gamma(m)}{\Gamma(m+n)} g(m) m^k dm}{\int \frac{\Gamma(m)}{\Gamma(m+n)} g(m) m^k dm} \approx \sqrt{\frac{\ell''(\hat{m})}{\ell''(\hat{m}_\nu)}} \exp \{n [\ell_\nu(\hat{m}_\nu) - \ell(\hat{m})]\}, \quad (22)$$

where

$$\begin{aligned} \ell &= \log \frac{m^{a-1} e^{-m/b}}{\Gamma(a) b^a} + \frac{1}{n} \left\{ \log \frac{\Gamma(m)}{\Gamma(m+n)} + k \log m \right\} \\ \ell_\nu &= \ell + \nu \log m \\ \ell' &= \frac{\partial}{\partial m} \ell = \frac{1}{bm} \left[ b \left( \frac{k}{n} + a - 1 \right) - m - \frac{bm}{n} \sum_{i=1}^n \frac{1}{m+i-1} \right] \\ \ell''(\hat{m}) &= \frac{\partial^2}{\partial m^2} \ell \Big|_{m=\hat{m}} = \frac{1}{\hat{m}} \left[ -\frac{1}{\hat{m}} \left( \frac{k}{n} + a - 1 \right) + \frac{\hat{m}}{n} \sum_{i=1}^n \frac{1}{(\hat{m}+i-1)^2} \right] \\ \ell'_\nu &= \ell' + \frac{\nu}{m}, \quad \ell''_\nu(\hat{m}_\nu) = \frac{\partial^2}{\partial m^2} \ell_\nu \Big|_{m=\hat{m}_\nu} = \ell''(\hat{m}_\nu) - \frac{\nu}{\hat{m}_\nu^2} \end{aligned}$$

where we get a simplification because the second derivative is evaluated at the zero of the first derivative. We use these approximations as the first and second moments of the candidate gamma distribution. Note that if  $\hat{m} \approx \hat{m}_\nu$ , then a crude approximation, which should be enough for Metropolis, is  $\text{Em}^\nu \approx (\hat{m})^\nu$ .

## 5 Simulation Study

We evaluate our sampler through a number of simulation studies. We need to generate outcomes from Bernoulli or Poisson distributions with random effects that follow the Dirichlet process. To do this we fix  $K$ , the true number of clusters (which is unknown in actual circumstances), then we set the parameter  $m$  according to the relation

$$K = \sum_{i=1}^n \frac{m}{m+i-1}, \quad (23)$$

where we note that even if  $\hat{m}$  is quite variable, there is less variability in  $\hat{K} = \sum_{i=1}^n \frac{\hat{m}}{\hat{m}+i-1}$ . When we integrate over the Dirichlet process (as done algorithmically according to Blackwell

and McQueen [1973]), the right-hand-side of (23) is the expected number of clusters, given the prior distribution on  $m$ . Neal (2000, p.252) shows this as the probability in the limit of a unique table seating, conditional on the previous table seatings, which makes intuitive sense since this expectation depends on individuals sitting at unique tables to start a new (sub)cluster in the algorithm.

## 5.1 Logistic Models

Using the GLMDM with the logistic link function of Section 3.2, we set the parameters:  $n = 100$ ,  $K = 40$ ,  $\tau^2 = 1$ , and  $\beta = (1, 2, 3)$ . Our Dirichlet process for the random effect has precision parameter  $m$  and base distribution  $G_0 = N(0, \tau^2)$ . Setting  $K = 40$ , yields  $m = 24.21$ . We then generated  $X_1$  and  $X_2$  independently from  $N(0, 1)$ , and used the fixed design matrix to generate the binary outcome  $Y$ . Then the Gibbs sampler was iterated 200 times to get values of  $m$ ,  $A$ ,  $\beta$ ,  $\tau^2$ ,  $\eta$ . This procedure was repeated 1000 times saving the last 500 draws as simulations from the posterior.

We compare the slice sampler (**Slice**) to the Gibbs sampler with the K-S distribution normal scale mixture (**K-S Mixture**) with the prior distribution of  $\beta$  from  $\beta|\sigma^2 \sim N(\mu\mathbf{1}, d^*\sigma^2I)$  and  $\mu \sim \pi(\mu) \propto c$ , a flat prior for  $\mu$ . For the estimation of  $K$ , we use the posterior mean of  $m$ ,  $\hat{m}$  and calculate  $\hat{K}$  by using equation (23). The starting points of  $\beta$  come from the maximum likelihood (ML) estimates using iteratively reweighted least squares. All summaries in the tables are posterior means and standard deviations calculated from the empirical draws of the chain in

Figure 2: ACF Plots of  $\beta$  for the GLMDM with logistic link. The left panel are the plots for  $(\beta_0, \beta_1, \beta_2)$  from the slice sampler, and the right panel are the plots for  $(\beta_0, \beta_1, \beta_2)$  from the K-S/normal mixture sampler.

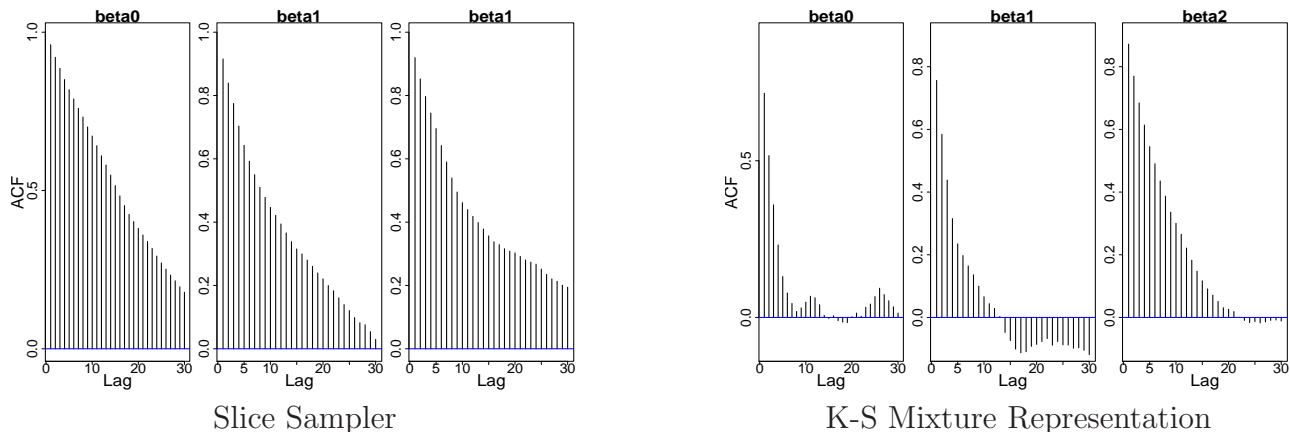




Table 1: Estimation of the coefficients of the GLMDM with logistic link function and the estimate of  $K$ , with true values  $K = 40$  and  $\beta = (1, 2, 3)$ . Standard errors are in parentheses.

Estimation Method	$\beta_0$	$\beta_1$	$\beta_2$	$K$
<b>Slice</b>	2.2796(0.4628)	3.2709(0.5558)	4.7529(0.7208)	43.0423(4.2670)
<b>K-S Mixture</b>	0.4900(0.2024)	1.0494(0.2468)	1.7787(0.2491)	43.4646(4.0844)

its apparent converged (stationary) distribution.

The numerical summary of this process is given in Table 1. The estimates of  $K$  were 43.0423 with standard error 4.2670 from **Slice** and 43.4646 with standard error 4.0844 from **K-S Mixture**. Obviously these turned out to be good estimates of the true  $K = 40$ . The estimate of  $\beta$  with **K-S Mixture** is closer to the true value than those with **Slice**, with smaller standard deviation. To evaluate the convergence of  $\beta$ , we consider the autocorrelation function (ACF) plots that are given in Figure 2. The Gibbs sampler of  $\beta$  from **Slice** exhibits strong autocorrelation, implying poor mixing.

## 5.2 Log Linear Models

We now look at the GLMDM with the log link function of Section 3.3. The setting for the data generation is the same as the procedure that we discussed in the previous section except that we take  $\beta = (3, 0.5, 1)$ . With  $K = 40$ , the solution of  $m$  from equation (23) is 24.21. As before, we generated  $X_1$  and  $X_2$  independently from  $N(0, 1)$ , and used the fixed design matrix to generate count data  $Y$ . The Gibbs sampler was iterated 200 times to produce draws of  $m, A, \beta, \tau^2, \eta$ . This procedure was repeated 1000 times, saving the last 500 values as draws from the posterior.

In this section, we compare the Gibbs sampler with the auxiliary variables (**Slice**) and the M-H sampler with a candidate density from the log-linear model (**M-H Sampler**). We use the posterior mean of  $m$ ,  $\hat{m}$ , and calculate  $\hat{K}$  by using (23) for the estimation of  $K$ . The starting points of  $\beta$  are set to the maximum likelihood (ML) estimates by using the iterative reweighted least squares. The numerical summary is given in Table 2 and the ACF plots of  $\beta$  are given in Figure 3. The resulting estimates for  $K$  are 43.5188(4.1398) from **Slice** and 43.516(4.1274) from the **M-H Sampler**, which are fairly close to the true  $K = 40$ . The estimated  $\beta$ s from the **M-H Sampler**, while not right on target, are much better than that of the slice sampler which, by standard diagnostics, has not yet converged. Once again, the consecutive draws of  $\beta$  of **Slice** from the Gibbs sampler are strongly autocorrelated. The convergence of  $\beta$  of **Slice** and **M-H Sampler** can be assessed by viewing the ACF plots in Figure 3. The M-H chain with candidate

Table 2: Estimation of the coefficients of the GLMDM with log link function and the estimate of  $K$ , with true values  $K = 40$  and  $\beta = (3, 0.5, 1)$ . Standard errors are in parentheses.

Estimation Method	$\beta_0$	$\beta_1$	$\beta_2$	$K$
<b>Slice</b>	2.7984(0.0099)	0.0907(0.0196)	0.8350(0.0184)	43.5188(4.1398)
<b>M-H Sampler</b>	2.3107(0.1407)	0.8493(1.1309)	0.9492(1.0637)	43.5161 (4.1274)

densities from log-linear models mixes better, giving additional confidence about convergence.

### 5.3 Probit Models

For completeness, we also generated data, similar to that described in Section 3.2, for a probit link. In Figure 4 we only show the ACF plot from a latent variable Gibbs sampler as described in Section 3.1. where we see that the autocorrelations are not as good as the M-H algorithm, but better than those of the slice sampler.

## 6 Data Analysis

In this section we provide two real data examples that highlight the workings of generalized linear Dirichlet process random effects models, using both logit and probit link functions. Both examples are drawn from important questions in social science research: voting behavior and

Figure 3: ACF Plots of  $\beta$  for the GLMDM with log link. The left panel are the plots for  $(\beta_0, \beta_1, \beta_2)$  from the slice sampler, and the right panel are the plots for  $(\beta_0, \beta_1, \beta_2)$  from the M-H sampler.

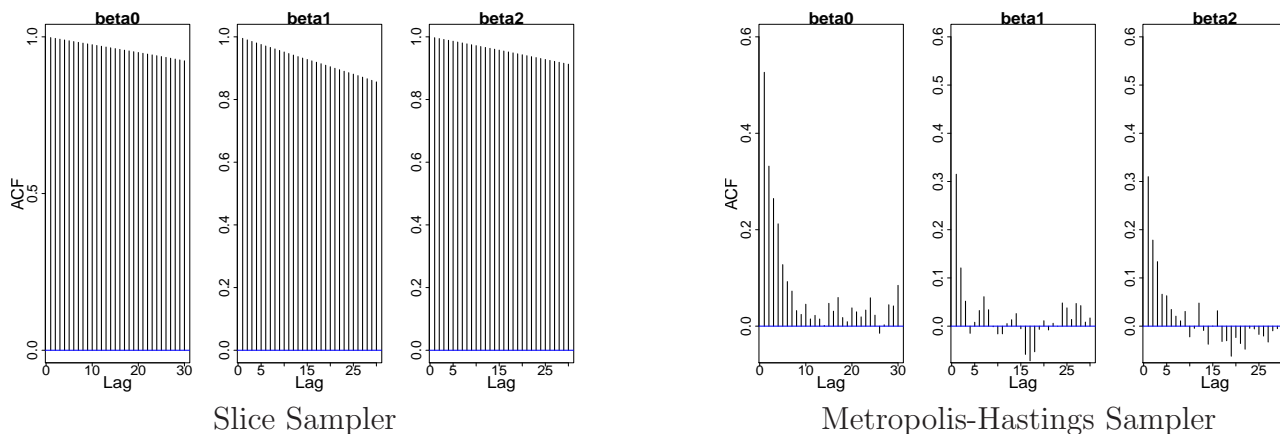
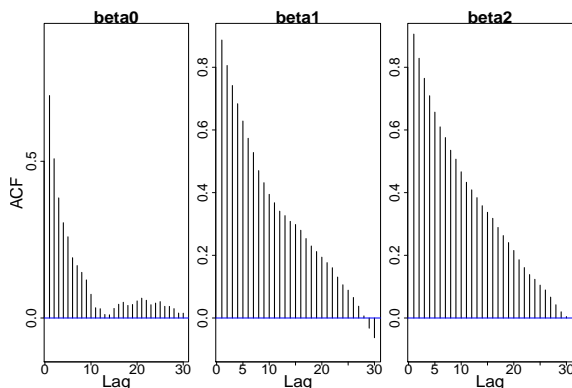


Figure 4: ACF Plots for  $(\beta_0, \beta_1, \beta_2)$  for the GLMDM with probit link, using the simulated data of Section 5.1.



terrorism studies. The voting behavior study, of social attitudes in Scotland, is fit using a logit link, while the terrorism data is fit with a probit link.

## 6.1 Social Attitudes in Scotland

The data for this example come from the Scottish Social Attitudes Survey, 2006 (UK Data Archive Study Number 5840). This study is based on face-to-face interviews conducted using computer assisted personal interviewing and a paper-based self-completion questionnaire, providing 1594 data points and 669 covariates. However, to highlight the challenge in identifying consistent attitudes with small data sizes we restrict the sample analyzed to females 18-25 years-old, giving 44 cases. This is a politically interesting group in terms of their interaction with the government, particularly with regard to healthcare and Scotland’s voice in UK public affairs. The general focus was on attitudes towards government at the UK and national level, feelings about racial groups including discrimination, views on youth and youth crime, as well as exploring the Scottish sense of national identity.

Respondents were asked whether they favored full independence for Scotland with or without membership in the European Union versus remaining in the UK under varying circumstances. This was used as a dichotomous outcome variable to explore the factors that contribute to advocating secession for Scotland. The explanatory variables used are: `househld` measuring the number of people living in the respondent’s household, `relgsums` indicating identification with the Church of Scotland versus another or no religion, `ptyallgs` measuring party allegiance with the ordering of parties given from more conservative to more liberal, `idlosem` a dichotomous variable equal to one if the respondent agreed with the statement that increased numbers of Muslims in Scotland would erode the national identity, `marrmus` another dichotomous variable

equal to one if the respondent would be unhappy or very unhappy if a family member married a Muslim, `ukintnat` for agreement that the UK government works in Scotland’s long-term interests, `natinnat` for agreement that the Scottish Executive works in Scotland’s long-term interests, `voiceuk3` indicating that the respondent believes that the Scottish Parliament gives Scotland a greater voice in the UK, `nhssat` indicating satisfaction (1) or dissatisfaction (0) with the National Health Service, `hincdif2`, a seven-point Likert scale showing the degree to which the respondent is living comfortably on current income or not (better in the positive direction), `unionsa` indicating union membership at work, `whrbrn` a dichotomous variable indicating birth in Scotland or not, and `hedqual2` the respondent’s education level. We retain the variable names from the original study for ease of replication by others. All coding decisions (along with code for the models and simulations) are documented on the webpage <http://jgill.wustl.edu/replication.html>.

We ran the Markov chain for 10,000 iterations saving the last 5,000 for analysis. All indications point towards convergence using empirical diagnostics (Geweke, Heidelberger & Welsh, graphics, etc.). The results in Table 3 are interesting in a surprising way. Notice that there are very similar results for the standard Bayesian logit model with flat priors (estimated in JAGS, see <http://www-fis.iarc.fr/~martyn/software/jags/>) and the GLMDM logit model, save for one coefficient (discussed below). This indicates that the nonparametric component does not affect all of the marginal posterior distributions and the recovered information is confined to specific aspects of the data. Figure 5 graphically displays the credible intervals, and makes it easier to see the agreement of the analyses in this case.

Several of the coefficients point towards interesting findings from these results. There is reliable evidence from the Dirichlet process results that women under 25 believe that increased numbers of Muslims in Scotland would erode the Scottish national identity. This is surprising since anecdotally and journalistically one would expect this group to be among the most welcoming in the country. There is modest evidence (the two models differ slightly here) that this group is dissatisfied by the service provided by the National Health Service. In addition, these young Scottish women report not living comfortably on their current income. It is also interesting here that the prior information provided by the GLMDM model is overwhelmed by the data as evidenced by the similarity between the two models. In line with Kyung (2010), most of the credible intervals of the GLMDM model are slightly shorter.

## 6.2 Terrorism Targeting

In this example we look at terrorist activity in 22 Asian democracies over 8 years (1990-1997) with data subsetting from Koch and Cranmer (2007). Data problems (a persistent issue in the empirical study of terrorism) reduce the number of cases to 162 and make fitting any standard model

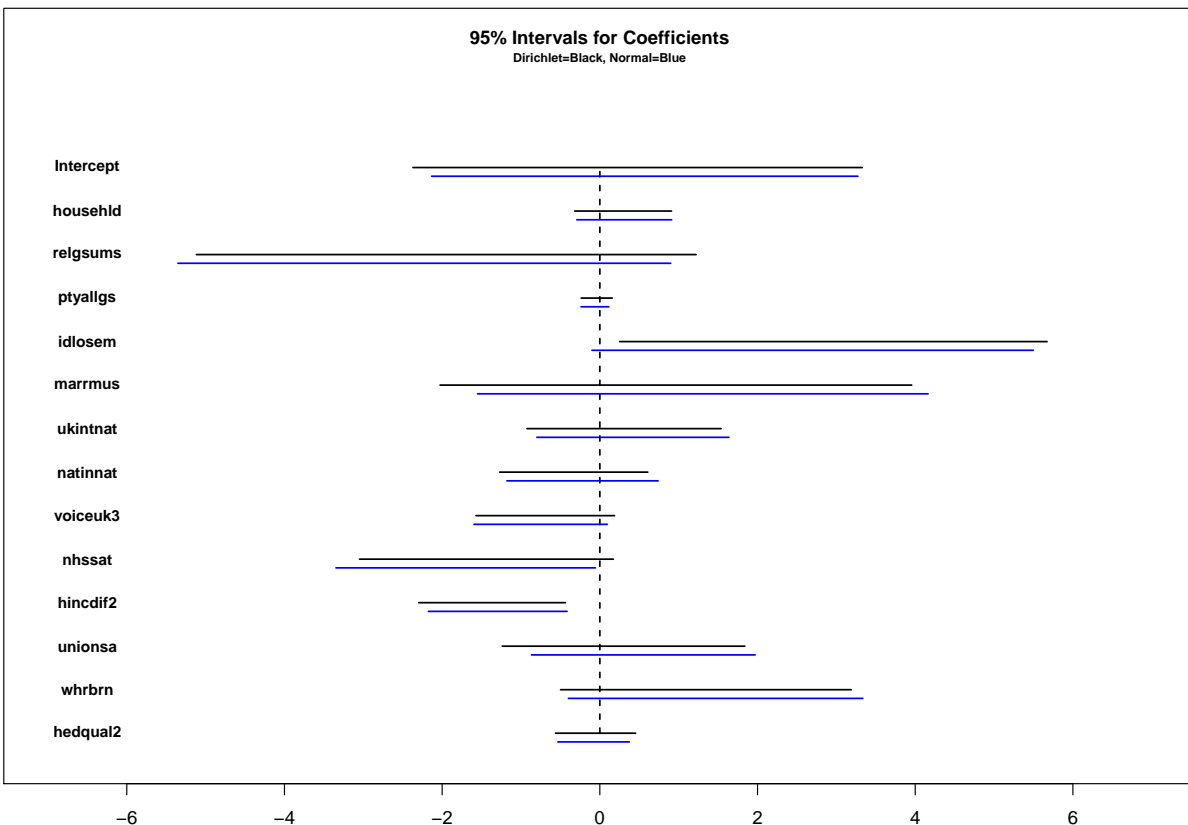
Table 3: LOGIT MODELS FOR ATTITUDES OF FEMALES 18-25 YEARS IN SCOTLAND

Coefficient	Standard Logit				GLMDM Logit			
	COEF	SE	95% CI		COEF	SE	95% CI	
Intercept	0.563	1.358	-2.133	3.274	0.351	1.396	-2.321	3.075
househld	0.281	0.303	-0.293	0.912	0.239	0.299	-0.342	0.830
relgsums	-2.006	1.604	-5.352	0.899	-1.840	1.614	-5.175	1.114
ptyallgs	-0.066	0.089	-0.239	0.114	-0.035	0.091	-0.207	0.150
idlosem	2.381	1.432	-0.101	5.498	2.663	1.343	0.219	5.487
marrmus	1.281	1.469	-1.552	4.164	1.089	1.528	-1.818	4.151
ukintnat	0.403	0.616	-0.799	1.638	0.347	0.582	-0.752	1.553
natinnat	-0.194	0.487	-1.179	0.739	-0.304	0.446	-1.174	0.575
voiceuk3	-0.708	0.433	-1.597	0.095	-0.637	0.443	-1.573	0.159
nhssat	-1.677	0.841	-3.347	-0.056	-1.405	0.812	-3.018	0.152
hincdif2	-1.219	0.446	-2.175	-0.415	-1.205	0.448	-2.114	-0.387
unionsa	0.521	0.723	-0.867	1.970	0.247	0.718	-1.117	1.692
whrbrn	1.494	0.944	-0.398	3.336	1.229	0.861	-0.461	2.924
hedqual2	-0.082	0.233	-0.532	0.374	-0.036	0.235	-0.493	0.434

difficult due to the generally poor level of measurement. The outcome of interest is dichotomous, indicating whether or not there was at least one violent terrorist act in a country/year pair. In order to control for the *level* of democracy (DEM) in these countries we use the Polity IV 21-point democracy scale ranging from -10 indicating a hereditary monarchy to +10 indicating a fully consolidated democracy (Gurr, Marshall, and Jagers 2003). The variable FED is assigned zero if sub-national governments do not have substantial taxing, spending, and regulatory authority, and one otherwise. We look at three rough classes of government structure with the variable SYS coded as: (0) for direct presidential elections, (1) for strong president elected by assembly, and (2) dominant parliamentary government. Finally, AUT is a dichotomous variable indicating whether or not there are autonomous regions not directly controlled by central government. The key substantive question evaluated here is whether specific structures of government and sub-governments lead to more or less terrorism.

We ran the Markov chain for 50,000 iterations disposing of the first half. There is no evidence of non-convergence in these runs using standard diagnostic tools. Table 4 again provides results from two approaches: a standard Bayesian probit model with flat priors, and a Dirichlet process random effects model. Notice first that while there are no changes in sign or statistical reliability for the estimated coefficients, the magnitudes of the effects are uniformly smaller

Figure 5: LENGTHS AND PLACEMENT OF CREDIBLE INTERVALS FOR THE COEFFICIENTS OF THE LOGIT MODEL FIT FOR THE SCOTTISH SOCIAL ATTITUDES SURVEY ON FEMALES 18-25 YEARS USING DIRICHLET PROCESS RANDOM EFFECTS (BLACK) AND NORMAL RANDOM EFFECTS (BLUE).



with the enhanced model: four of the estimates are roughly twice as large and the last one is about three times as large as in the standard model. This is clearly seen in Figure 6, which is a graphical display of Table 4. We feel that this indicates that there is extra variability in the data detected by the Dirichlet process random effect that tends to dampen the size of the effect of these explanatory variables on explaining incidences of terrorist attacks. Specifically, running the standard probit model would find an *exaggerated* relationship between these explanatory variables and the outcome.

The results are also interesting substantively. The more democratic a country is, the more terrorist attacks they can expect. This is consistent with the literature in that autocratic nations tend to have more security resources per capita and fewer civil rights to worry about. Secondly, the more the legislature holds central power, the fewer expected terrorist attacks. This also

Table 4: Probit Models for Terrorism Incidents

Coefficient	Standard Probit				GLMDM Probit			
	COEF	SE	95% CI		COEF	SE	95% CI	
Intercept	0.249	0.337	-0.412	0.911	0.123	0.187	-0.244	0.490
DEM	0.109	0.035	0.041	0.177	0.058	0.019	0.020	0.095
FED	0.649	0.469	-0.270	1.567	0.253	0.254	-0.245	0.750
SYS	-0.817	0.252	-1.312	-0.323	-0.418	0.136	-0.685	-0.151
AUT	1.619	0.871	-0.088	3.327	0.444	0.369	-0.279	1.167

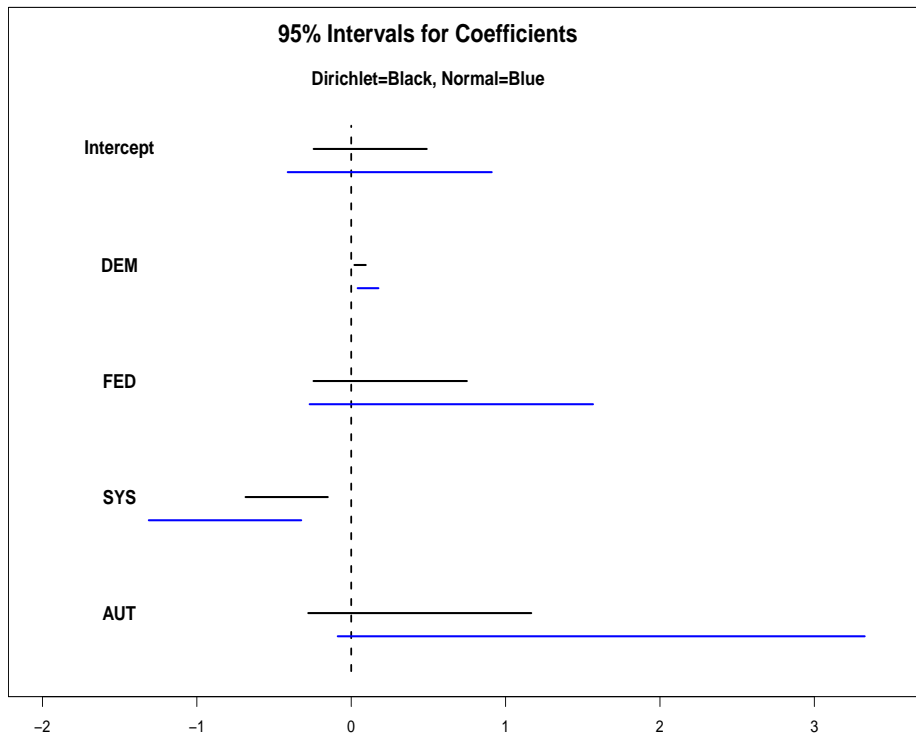
makes sense, given what is known; disparate groups in society tend to have a greater voice in government when the legislature dominates the executive. Two results are puzzling and are therefore worth further investigation. Strong sub-governments and the presence of autonomous regions both lead to more expected terrorism. This may result from strong separatist movements and typical governmental responses, an observed endogenous and cycling effect that often leads to prolonged struggles and intractable relations. We further investigate the use of Dirichlet process priors for understanding latent information in terrorism data in Kyung *et al.* (2011) with the goal of sorting out such effects.

## 7 Discussion

In this paper we demonstrate how to set up and run sampling schemes for the generalized linear mixed Dirichlet process model with a variety of link functions. We focus on the mixed effects model with a Dirichlet process prior for the random effects instead of the normal assumption, as in standard approaches. We are able to estimate model parameters as well as the Dirichlet process parameters using convenient MCMC algorithms, and to draw latent information from the data. Simulation studies and empirical studies demonstrate the effectiveness of this approach.

The major methodological contributions here are the derivation and evaluation of strategies of estimation for model parameters in Section 3 and the inclusion of the precision parameter directly into the Gibbs sampler for estimation in Section 4.2. In the latter case, including the precision parameter in the Gibbs sampler means that we are marginalizing over the parameter rather than conditioning on it leading to a more robust set of estimates. Moreover, we have seen a large amount of variability in the performance of MCMC algorithms, with the slice sampler typically being less optimal than either a K-S mixture representation or a Metropolis-Hastings algorithm.

Figure 6: Lengths and placement of credible intervals for the coefficients of the probit model fit for the terrorist activity data using Dirichlet process random effects (black) and normal random effects (blue).



The relationship of credible intervals that is quite evident in Figure 6, and less so in Figure 5, that the Dirichlet intervals tend to be shorter than those based on normal random effects, persists in other data that we have analyzed. We have found that this is not a data anomaly, but has an explanation in that the Dirichlet process random effects model results in posterior variances that are smaller than that of the normal. Kyung *et al.* (2009) are able to prove this first in a special case of the linear model (when  $\mathbf{X} = \mathbf{I}$ ), and then for almost all data vectors. The intuition follows the logic of multilevel (hierarchical) models whereby some variability at the individual-level is moved to the heterogeneous group-level thus producing a better model fit. Here, the group-level is represented by the nonparametric assignment to latent categories through the process of the Gibbs sampler.

Finally, we observed that the additional effort needed to include a Dirichlet process prior for the random effects in two empirical examples with social science data, which tends to be more messy and interrelated than that in other fields, added significant value to the data analysis. We found that the GLMDM model can detect additional variability in the data which affects parameter estimates. In particular, in the case of social attitudes in Scotland the GLMDM model



corrected a clearly illogical finding in the usual probit analysis. For the second example, we found that the GLMDM specification dampened-down over enthusiastic findings from a conventional model. In both cases either non-Bayesian or Bayesian models with flat priors would have reported results that had substantively misleading findings.

## 8 References

- Abramowitz, M. and Stegun, I. A. (1972). "Stirling Numbers of the Second Kind." Section 24.1.4 in *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, 9th printing. New York: Dover, 824-825.
- Albert, J. H. and S. Chib (1993). "Bayesian Analysis of Binary and Polychotomous Response Data." *Journal of the American Statistical Association* **88**, 669-679.
- Andrews D. F. and Mallows, C. L. (1974). "Scale Mixtures of Normal Distributions." *Journal of the Royal Statistical Society, Series B* **36**, 99-102.
- Antoniak, Charles E. (1974). "Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems." *Annals of Statistics* **2**, 1152-1174.
- Balakrishnan, N. (1992). *Handbook of the Logistic Distribution*. CRC Press.
- Blackwell, D. and MacQueen, J. B. (1973). "Discreteness of Ferguson Selections." *Annals of Statistics* **1**, 365-358.
- Breslow, N. E., and D. G. Clayton. (1993). "Approximate Inference in Generalized Linear Mixed Models." *Journal of the American Statistical Association* **88**, 9-25.
- Buonaccorsi, J. P. (1996). "Measurement Error in the Response in the General Linear Model." *Journal of the American Statistical Association* **91**, 633-642.
- Chib, S., Greenberg, E., and Chen, Y. (1998). "MCMC methods for fitting and comparing multinomial response models." Technical Report, Economics Working Paper Archive, Washington University at St. Louis, <http://129.3.20.41/econ-wp/em/papers/9802/9802001.pdf>.
- Chib, S. and Winkelmann, R. (2001). "Markov Chain Monte Carlo Analysis of Correlated Count Data." *Journal of Business and Economic Statistics* **19**, 428-435.
- Damien, P., Wakefield, J., and Walker, S. (1999) "Gibbs Sampling for Bayesian Non-Conjugate and Hierarchical Models by Using Auxiliary Variables." *Journal of the Royal Statistical Society, Series B* **61**, 331-344.

- Devroye, L. (1986). *Non-Uniform Random Variate Generation*. Springer-Verlag New York Inc.
- Dey, D. K., Ghosh, S. K., and Mallick, B. K. (2000). *Generalized Linear Models: A Bayesian Perspective*. New York: Marcel Dekker.
- Dorazio, R. M., Mukherjee, B., Zhang, L., Ghosh, M., Jelks, H. L. and Jordan, F. (2007). “Modelling Unobserved Sources of Heterogeneity in Animal Abundance Using a Dirichlet Process Prior.” *Biometrics*, Online Publication August 3, 2007.
- Escobar, M. D. and West, M. (1995). “Bayesian Density Estimation and Inference Using Mixtures.” *Journal of the American Statistical Association* **90**, 577-588.
- Fahrmeir, L. and Tutz, G. (2001). *Multivariate Statistical Modelling Based on Generalized Linear Models*. Second Edition. New York: Springer.
- Ferguson, T. S. (1973). “A Bayesian analysis of Some Nonparametric Problems.” *Annals of Statistics* **1**, 209-230.
- Gill, J. and Casella, G. (2009). “Nonparametric Priors for Ordinal Bayesian Social Science Models: Specification and Estimation.” *Journal of the American Statistical Association* **104** 453-464.
- Gurr, T. R., Marshall, M. G., and Jaggers, K. (2003). *PolityIV*, <http://www.cidcm.umd.edu/inscr/polity/>.
- Ishwaran, Hemant and James, Lancelot F. (2001). “Gibbs Sampling Methods for Stick-Breaking Priors.” *Journal of the American Statistical Association* **96**, 161-173.
- Jiang, J. (2007). *Linear and Generalized Linear Mixed Models and Their Applications*, New York: Springer-Verlag.
- Koch, M. T. Cranmer, S. (2007). “Terrorism than Governments of the Right? Testing the ‘Dick Cheney’ Hypothesis: Do Governments of the Left Attract More than Governments of the Right?” *Conflict Management and Peace Science* **24**, 311-326.
- Korwar, R. M. and Hollander, M. (1973). “Contributions to the Theory of Dirichlet Processes.” *Annals of Probability* **1**, 705-711.
- Kyung, M., Gill, J. and Casella, G. (2009). “Characterizing the variance improvement in linear Dirichlet random effects models.” *Statistics and Probability Letters* **79** 2343-2350.

- Kyung, M., Gill, J. and Casella G. (2010). “Estimation in Dirichlet Random Effects Models.” *Annals of Statistics*, 38, 979-1009.
- Kyung, M., Gill, J. and Casella G. (2011). “New Findings from Terrorism Data: Dirichlet Process Random Effects Models for Latent Groups.” *Journal of the Royal Statistical Society, Series C*, Forthcoming.
- Liu, J. S. (1996). “Nonparametric Hierarchical Bayes Via Sequential Imputations.” *Annals of Statistics* **24**, 911-930.
- Lo, A. Y. (1984). “On A Class of Bayesian Nonparametric Estimates: I. Density Estimates.” *Annals of Statistics* **12**, 351-357.
- MacEachern, S. N. and Müller, P. (1998). “Estimating Mixture of Dirichlet Process Model.” *Journal of Computational and Graphical Statistics* **7**, 223-238.
- McAuliffe, J. D., Blei, D. M. and Jordan, M. I. (2006). “Nonparametric Empirical Bayes for the Dirichlet Process Mixture Model.” *Statistics and Computing* **16**, 5-14.
- McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*, 2nd Edition. New York: Chapman & Hall.
- McCullagh, P. and Yang, J. (2006). “Stochastic Classification Models.” *International Congress of Mathematicians III* 669-686.
- McCulloch, C. E. and Searle, S. R. (2001). *Generalized, Linear, and Mixed Models*. New York: John Wiley & Sons.
- Mira, A. and Tierney, L. (2002). Efficiency and Convergence Properties of Slice Samplers. *Scandinavian Journal of Statistics* **29** 1-12.
- Neal, R. M. (2000). “Markov Chain Sampling Methods for Dirichlet Process Mixture Models.” *Journal of Computational and Graphical Statistics* **9**, 249-265.
- Neal, R. M. (2003). “Slice Sampling.” *The Annals of Statistics* **31**, 705-741.
- Robert, C. and Casella, G. (2004). *Monte Carlo Statistical Methods*. New York: Springer.
- Sethuraman, J. (1994). “A Constructive Definition of Dirichlet Priors.” *Statistica Sinica* **4**, 639-650.
- Teh, Y. W., Jordan, M. I., Beal, M. J. and Blei, D. M. (2006). “Hierarchical Dirichlet Processes.” *Journal of the American Statistical Association* **101**, 1566-1581.

Tierney, L. and Kadane, J. B. (1986). “Accurate Approximations for Posterior Moments and Marginal Densities.” *Journal of the American Statistical Association* **81**, 82-86.

Wang, N., X. Lin, R. G. Gutierrez, and R. J. Carroll. (1998). “Bias Analysis and SIMEX Approach in Generalized Linear Mixed Measurement Error Models.” *Journal of the American Statistical Association* **93**, 249-261.

West, M. (1987). “On Scale Mixtures of Normal Distributions”, *Biometrika* **74**, 646-648.

Wolfinger, R. and M. O’Connell. (1993). “Generalized Linear Mixed Models: A Pseudolikelihood Approach.” *Journal of Statistical Computation and Simulation* **48**, 233-243.

## A Generating the Model Parameters

### A.1 A Logistic Model

#### A.1.1 Slice Sampling

For fixed  $m$  and  $A$ , a Gibbs sampler of  $(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V})$  is

- for  $d = 1, \dots, p$ ,

$$\beta_d | \boldsymbol{\beta}_{-d}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \begin{cases} N(0, d^* \sigma^2) & \text{if } \beta_d \in \left[ \left\{ \max \left( \max_{X_{id} > 0} \left( \frac{\alpha_{id}}{X_{id}} \right), \left( \max_{X_{id} \leq 0} \left( \frac{\gamma_{id}}{X_{id}} \right) \right) \right\}, \right. \\ \left. \left\{ \min \left( \min_{X_{id} \leq 0} \left( \frac{\alpha_{id}}{X_{id}} \right), \left( \min_{X_{id} > 0} \left( \frac{\gamma_{id}}{X_{id}} \right) \right) \right\} \right] \right. \\ 0 & \text{otherwise} \end{cases}$$

where

$$\begin{aligned} \alpha_{id} &= -\log \left( u_i^{-\frac{1}{y_i}} - 1 \right) - \sum_{l \neq d} X_{il} \beta_l - (\mathbf{A}\boldsymbol{\eta})_i & \text{for } i \in S \\ \gamma_{id} &= \log \left( v_i^{\frac{1}{y_i-1}} - 1 \right) - \sum_{l \neq d} X_{il} \beta_l - (\mathbf{A}\boldsymbol{\eta})_i & \text{for } i \in F. \end{aligned}$$

Here,  $S = \{i : y_i = 1\}$  and  $F = \{i : y_i = 0\}$ .

- $\tau^2 | \boldsymbol{\beta}, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \text{Inverted Gamma} \left( \frac{k}{2} + a, \frac{1}{2} |\boldsymbol{\eta}|^2 + b \right)$
- for  $j = 1, \dots, k$ ,

$$\eta_j | \boldsymbol{\beta}, \tau^2, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \begin{cases} N(0, \tau^2) & \text{if } \eta_j \in \left( \max_{i \in S_j} \{\alpha_i^*\}, \min_{i \in S_j} \{\gamma_i^*\} \right) \\ 0 & \text{otherwise} \end{cases},$$

where

$$\alpha_i^* = -\log(u_i^{-1} - 1) - \mathbf{X}_i\boldsymbol{\beta} \quad \text{for } i \in S$$

$$\gamma_i^* = \log(v_i^{-1} - 1) - \mathbf{X}_i\boldsymbol{\beta} \quad \text{for } i \in F$$

- for  $i = 1, \dots, n$ ,

$$\pi_k(U_i|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{V}, A, \mathbf{y}) \propto I\left[u_i < \left\{\frac{1}{1 + \exp(-\mathbf{X}_i\boldsymbol{\beta} - \eta_j)}\right\}^{y_i}\right] \quad \text{for } i \in S$$

$$\pi_k(V_i|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y}) \propto I\left[v_i < \left\{\frac{1}{1 + \exp(\mathbf{X}_i\boldsymbol{\beta} + \eta_j)}\right\}^{1-y_i}\right] \quad \text{for } i \in F.$$

### A.1.2 K-S Mixture

Given  $\xi$ , for fixed  $m$  and  $A$ , a Gibbs sampler of  $(\mu, \boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U})$  is

$$\boldsymbol{\eta}|\mu, \boldsymbol{\beta}, \tau^2, \mathbf{U}, A, \mathbf{y}, \sigma^2 \sim N_k\left(\frac{1}{\sigma^2(2\xi)^2}\left(\frac{1}{\tau^2}I + \frac{1}{\sigma^2(2\xi)^2}A'A\right)^{-1}A'(\mathbf{U} - X\boldsymbol{\beta}), \left(\frac{1}{\tau^2}I + \frac{1}{\sigma^2(2\xi)^2}A'A\right)^{-1}\right)$$

$$\mu|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y}, \sigma^2 \sim N\left(\frac{1}{p}\mathbf{1}'_p\boldsymbol{\beta}, \frac{d^*}{p}\sigma^2\right)$$

$$\boldsymbol{\beta}|\mu, \tau^2, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y}, \sigma^2 \sim N_p\left(\left(\frac{1}{d^*}I + \frac{1}{(2\xi)^2}X'X\right)^{-1}\left(\frac{1}{d^*}\mu\mathbf{1}_p + \frac{1}{(2\xi)^2}X'(\mathbf{U} - A\boldsymbol{\eta})\right), \sigma^2\left(\frac{1}{d^*}I + \frac{1}{(2\xi)^2}X'X\right)^{-1}\right)$$

$$\tau^2|\mu, \boldsymbol{\beta}, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y}, \sigma^2 \sim \text{Inverted Gamma}\left(\frac{k}{2} + a, \frac{1}{2}|\boldsymbol{\eta}|^2 + b\right)$$

$$U_i|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, A, y_i, \sigma^2 \sim \begin{cases} N(X_i\boldsymbol{\beta} + (A\boldsymbol{\eta})_i, \sigma^2(2\xi)^2) I(U_i > 0) & \text{if } y_i = 1 \\ N(X_i\boldsymbol{\beta} + (A\boldsymbol{\eta})_i, \sigma^2(2\xi)^2) I(U_i \leq 0) & \text{if } y_i = 0 \end{cases}$$

Then we update  $\xi$  from

$$\xi|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y} \sim \left(\frac{1}{(2\xi)^2}\right)^{n/2} e^{-\frac{1}{2\sigma^2(2\xi)^2}|\mathbf{U} - X\boldsymbol{\beta} - A\boldsymbol{\eta}|^2} 8 \sum_{\alpha=1}^{\infty} (-1)^{\alpha+1} \alpha^2 \xi e^{-2\alpha^2\xi^2}.$$

The conditional posterior density of  $\xi$  is the product of a inverted gamma with parameters  $\frac{\alpha}{2} - 1$  and  $-\frac{1}{8\sigma^2}|\mathbf{U} - X\boldsymbol{\beta} - A\boldsymbol{\eta}|^2$ , and the infinite sum of the sequence  $(-1)^{\alpha+1}\alpha^2\xi e^{-2\alpha^2\xi^2}$ . To

generate samples from this target density, we consider the alternating series method that is proposed by Devroye (1986). Based on his notation, we take

$$\begin{aligned} ch(\xi) &= 8 \left( \frac{1}{\xi^2} \right)^{n/2} e^{-\frac{1}{8\sigma^2\xi^2}|\mathbf{U}-\mathbf{X}\boldsymbol{\beta}-\mathbf{A}\boldsymbol{\eta}|^2} \xi e^{-2\xi^2} \\ a_n(\xi) &= (\alpha + 1)^2 e^{-2\xi^2\{(\alpha+1)^2-1\}} \end{aligned}$$

Here, we need to generate sample from  $h(\xi)$ , and we use accept-reject sampling with candidate  $g(\xi^*) = 2e^{-2\xi^*}$ , the exponential distribution with  $\lambda = 2$ , where  $\xi^* = \xi^2$ . Then we follow Devroye's method.

## A.2 A Log Link Model

### A.2.1 Slice Sampling

Starting from the likelihood  $L(\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V})$ , and the priors on  $(\boldsymbol{\beta}, \tau^2)$ , we have the following Gibbs sampler of the model parameters.

- The conditional posterior distribution of  $\boldsymbol{\beta}$ :

$$\begin{aligned} \pi_K(\boldsymbol{\beta}|\tau^2, \boldsymbol{\eta}, A, \mathbf{y}, \mathbf{U}, \mathbf{V}) &\propto e^{-\frac{1}{2d^*\sigma^2}|\boldsymbol{\beta}|^2} \\ &\times \prod_{i=1}^n I[u_i < \exp\{y_i(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)\}, v_i > \exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]. \end{aligned}$$

For  $d = 1, \dots, p$ ,

$$\begin{aligned} \pi_K(\beta_d|\boldsymbol{\beta}_{-d}\tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) &\propto e^{-\frac{1}{2d^*\sigma^2}\beta_d^2} \\ &\times \prod_{i=1}^n I[u_i < \exp\{y_i(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)\}, v_i > \exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)], \end{aligned}$$

which can be expressed as:

$$\begin{aligned} \pi_K(\beta_d|\boldsymbol{\beta}_{-d}\tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) &\propto e^{-\frac{1}{2d^*\sigma^2}\beta_d^2} \\ &\times \prod_{i=1}^n I \left[ X_{id}\beta_d < \frac{1}{y_i} \log(u_i) - \sum_{l \neq j} X_{il}\beta_l - (\mathbf{A}\boldsymbol{\eta})_i, X_{id}\beta_d < \log(v_i) - \sum_{l \neq j} X_{il}\beta_l - (\mathbf{A}\boldsymbol{\eta})_i \right], \end{aligned}$$

where  $\boldsymbol{\beta}_{-d} = (\beta_1, \dots, \beta_{d-1}, \beta_{d+1}, \dots, \beta_p)$ . The full conditional posterior of  $\beta_d$  for  $d = 1, \dots, p$  is

$$\begin{aligned} \pi_k(\beta_d|\boldsymbol{\beta}_{-j}\tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) &\propto e^{-\frac{1}{2d^*\sigma^2}\beta_d^2} \beta_d \in \left[ \left\{ \max \left( \max_{X_{id}>0} \left( \frac{\alpha_{id}^*}{X_{id}} \right) \right), \left( \max_{X_{id}\leq 0} \left( \frac{\gamma_{id}^*}{X_{id}} \right) \right) \right\}, \right. \\ &\left. \left\{ \min \left( \min_{X_{id}\leq 0} \left( \frac{\alpha_{id}^*}{X_{id}} \right) \right), \left( \min_{X_{id}>0} \left( \frac{\gamma_{id}^*}{X_{id}} \right) \right) \right\} \right], \end{aligned}$$

where

$$\begin{aligned}\alpha_{id}^* &= \frac{1}{y_i} \log(u_i) - \sum_{l \neq d} X_{il} \beta_l - (\mathbf{A}\boldsymbol{\eta})_i \quad \text{for } i \in S \\ \gamma_{id}^* &= \log(v_i) - \sum_{l \neq d} X_{il} \beta_l - (\mathbf{A}\boldsymbol{\eta})_i \quad \text{for } i \in F\end{aligned}$$

Thus, for  $d = 1, \dots, p$ ,

$$\beta_d | \boldsymbol{\beta}_{-d}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \begin{cases} N(0, d^* \sigma^2) & \text{if } \beta_d \in \left[ \left\{ \max \left( \max_{X_{id} > 0} \left( \frac{\alpha_{id}^*}{X_{id}} \right), \max_{X_{id} \leq 0} \left( \frac{\gamma_{id}^*}{X_{id}} \right) \right) \right\}, \right. \\ \left. \left\{ \min \left( \min_{X_{id} \leq 0} \left( \frac{\alpha_{id}^*}{X_{id}} \right), \min_{X_{id} > 0} \left( \frac{\gamma_{id}^*}{X_{id}} \right) \right) \right\} \right] \\ 0 & \text{otherwise} \end{cases}$$

- The conditional posterior distribution of  $\tau^2$ :

$$\pi_k(\tau^2 | \boldsymbol{\beta}, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) \propto \left( \frac{1}{\tau^2} \right)^{k/2+a+1} e^{-\frac{1}{\tau^2}(\frac{1}{2}|\boldsymbol{\eta}|^2+b)}.$$

Thus,

$$\tau^2 | \boldsymbol{\beta}, \boldsymbol{\eta}, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \text{Inverted Gamma} \left( \frac{k}{2} + a, \frac{1}{2}|\boldsymbol{\eta}|^2 + b \right).$$

- The conditional posterior distribution of  $\boldsymbol{\eta}$ :

$$\pi_k(\boldsymbol{\eta} | \boldsymbol{\beta}, \tau^2, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) \propto \prod_{j=1}^k e^{-\frac{1}{2\tau^2}\eta_j^2} \prod_{i \in S_j} I[u_i < \exp\{y_i(\mathbf{X}_i\boldsymbol{\beta} + \eta_j)\}, v_i > \exp(\mathbf{X}_i\boldsymbol{\beta} + \eta_j)].$$

For  $j = 1, \dots, k$ ,

$$\begin{aligned}\pi_k(\eta_j | \boldsymbol{\beta}, \tau^2, \mathbf{U}, \mathbf{V}, A, \mathbf{y}) &\propto e^{-\frac{1}{2\tau^2}\eta_j^2} \prod_{i \in S_k} I[u_i < \exp\{y_i(\mathbf{X}_i\boldsymbol{\beta} + \eta_j)\}, v_i > \exp(\mathbf{X}_i\boldsymbol{\beta} + \eta_j)] \\ &\propto e^{-\frac{1}{2\tau^2}\eta_j^2} I \left[ \eta_j \in \left( \max_{i \in S_k} \{\gamma_i^*\}, \min_{i \in S_k} \{\alpha_i^*\} \right) \right],\end{aligned}$$

where

$$\begin{aligned}\alpha_i^* &= \frac{1}{y_i} \log(u_i) - \mathbf{X}_i\boldsymbol{\beta} \\ \gamma_i^* &= \log(v_i) - \mathbf{X}_i\boldsymbol{\beta}\end{aligned}$$

Thus, for  $j = 1, \dots, k$ ,

$$\eta_j | \boldsymbol{\beta}, \tau^2, \mathbf{U}, \mathbf{V}, A, \mathbf{y} \sim \begin{cases} N(0, \tau^2) & \text{if } \eta_j \in (\max_{i \in S_k} \{\gamma_i^*\}, \min_{i \in S_k} \{\alpha_i^*\}) \\ 0 & \text{otherwise} \end{cases}$$

- The conditional posterior distribution of  $\mathbf{U}$  and  $\mathbf{V}$ :

$$\begin{aligned}\pi_k(U_i|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{V}, A, \mathbf{y}) &\propto I[u_i < \exp\{y_i(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)\}] \\ \pi_k(V_i|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, \mathbf{U}, A, \mathbf{y}) &\propto e^{-v_i} I[v_i > \exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]\end{aligned}$$

### A.2.2 Metropolis-Hastings

Let  $Z_i \equiv \log(Y_i)$ , then  $Z_i$  is a linear mixed Dirichlet model (LMDM). For this model,

- the conditional posterior distribution of  $\boldsymbol{\beta}$  in the LMDM:

$$\boldsymbol{\beta}|\mu, \tau^2, \boldsymbol{\eta}, A, \mathbf{Z}, \sigma^2 \sim N_p \left( \left( \frac{1}{d^*}I + X'X \right)^{-1} \left( \frac{1}{d^*}\mu\mathbf{1}_p + X'(\mathbf{Z} - A\boldsymbol{\eta}) \right), \sigma^2 \left( \frac{1}{d^*}I + X'X \right)^{-1} \right). \quad (24)$$

- the conditional posterior distribution of  $\boldsymbol{\eta}$  in the LMDM:

$$\boldsymbol{\eta}|\boldsymbol{\beta}, \mu, \tau^2, \mathbf{Z}, A, \mathbf{y}, \sigma^2 \sim N_k \left( \frac{1}{\sigma^2} \left( \frac{1}{\tau^2}I + \frac{1}{\sigma^2}A'A \right)^{-1} A'(\mathbf{Z} - X\boldsymbol{\beta}), \left( \frac{1}{\tau^2}I + \frac{1}{\sigma^2}A'A \right)^{-1} \right). \quad (25)$$

Therefore, (24) is considered as a candidate density of  $\boldsymbol{\beta}$  and (25) for  $\boldsymbol{\eta}$ .

The Metropolis-Hastings sampler of  $(\boldsymbol{\beta}, \mu, \tau^2, \boldsymbol{\eta})$  follows.

- The conditional posterior distribution of  $\boldsymbol{\beta}$  in the log linear model:

$$\pi_k(\boldsymbol{\beta}|\mu, \tau^2, \boldsymbol{\eta}, A, \mathbf{y}, \sigma^2) \propto e^{-\frac{1}{2d^*\sigma^2}|\boldsymbol{\beta}-\mu\mathbf{1}_p|^2} \prod_{i=1}^n e^{-\exp(\mathbf{X}_i\boldsymbol{\beta}+(\mathbf{A}\boldsymbol{\eta})_i)} [\exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]^{y_i}.$$

Let

$$\pi_k^+(\boldsymbol{\beta}) \equiv e^{-\frac{1}{2d^*\sigma^2}|\boldsymbol{\beta}-\mu\mathbf{1}_p|^2} \prod_{i=1}^n e^{-\exp(\mathbf{X}_i\boldsymbol{\beta}+(\mathbf{A}\boldsymbol{\eta})_i)} [\exp(\mathbf{X}_i\boldsymbol{\beta} + (\mathbf{A}\boldsymbol{\eta})_i)]^{y_i}.$$

For given  $\boldsymbol{\beta}^{(t)}$ ,

1. Generate  $\boldsymbol{\beta}^* \sim N_p \left( \left( \frac{1}{d^*}I + X'X \right)^{-1} \left( \frac{1}{d^*}\mu\mathbf{1}_p + X'(\mathbf{Z} - A\boldsymbol{\eta}) \right), \sigma^2 \left( \frac{1}{d^*}I + X'X \right)^{-1} \right)$ .
2. Take

$$\boldsymbol{\beta}^{(t+1)} = \begin{cases} \boldsymbol{\beta}^* & \text{with probability } \min \left\{ \left( \frac{\pi_k^+(\boldsymbol{\beta}^*)}{\pi_k^+(\boldsymbol{\beta}^{(t)})} \frac{q(\boldsymbol{\beta}^{(t)})}{q(\boldsymbol{\beta}^*)} \right), 1 \right\} \\ \boldsymbol{\beta}^{(t)} & \text{otherwise} \end{cases},$$

where  $q(\cdot)$  is a density of  $N_p$  distribution in (24), and recall that  $\pi^+(\theta) = l(\theta)\pi(\theta)$ .



- The conditional posterior distribution of  $\mu$  in the log linear model:

$$\pi_k(\mu|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, A, \mathbf{y}, \sigma^2) \propto \exp \left\{ -\frac{p}{2d^*\sigma^2} \left( \mu - \frac{1}{p} \mathbf{1}'_p \boldsymbol{\beta} \right)^2 \right\}.$$

Thus,

$$\mu|\boldsymbol{\beta}, \tau^2, \boldsymbol{\eta}, A, \mathbf{y}, \sigma^2 \sim N \left( \frac{1}{p} \mathbf{1}'_p \boldsymbol{\beta}, \frac{d^*}{p} \sigma^2 \right).$$

- The conditional posterior distribution of  $\tau^2$  in the log linear model:

$$\pi_k(\tau^2|\boldsymbol{\beta}, \mu, \boldsymbol{\eta}, A, \mathbf{y}, \sigma^2) \propto \left( \frac{1}{\tau^2} \right)^{k/2+a+1} e^{-\frac{1}{\tau^2} (\frac{1}{2} |\boldsymbol{\eta}|^2 + b)}.$$

Thus,

$$\tau^2|\boldsymbol{\beta}, \mu, \boldsymbol{\eta}, A, \mathbf{y}, \sigma^2 \sim \text{Inverted Gamma} \left( \frac{k}{2} + a, \frac{1}{2} |\boldsymbol{\eta}|^2 + b \right).$$

- The conditional posterior distribution of  $\boldsymbol{\eta}$  in the log linear model:

$$\pi_k(\boldsymbol{\eta}|\boldsymbol{\beta}, \mu, \tau^2, A, \mathbf{y}, \sigma^2) \propto \prod_{k=1}^K e^{-\frac{1}{2\tau^2} \eta_k^2} \prod_{i \in S_k} e^{-\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)} [\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)]^{y_i}.$$

For  $j = 1, \dots, k$ , let

$$\begin{aligned} \pi_k^+(\eta_j) &\equiv e^{-\frac{1}{2\tau^2} \eta_j^2} \prod_{i \in S_j} e^{-\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)} [\exp(\mathbf{X}_i \boldsymbol{\beta} + (\mathbf{A} \boldsymbol{\eta})_i)]^{y_i} \\ &= e^{-\frac{1}{2\tau^2} \eta_j^2} \exp \left[ \eta_j \sum_{i \in S_j} y_i - e^{\eta_j} \sum_{i \in S_j} e^{\mathbf{X}_i \boldsymbol{\beta}} \right]. \end{aligned}$$

For given  $\boldsymbol{\eta}^{(t)}$ ,

1. Generate  $\boldsymbol{\eta}^* \sim N_k \left( \frac{1}{\sigma^2} \left( \frac{1}{\tau^2} I + \frac{1}{\sigma^2} A' A \right)^{-1} A' (\mathbf{Z} - X \boldsymbol{\beta}), \left( \frac{1}{\tau^2} I + \frac{1}{\sigma^2} A' A \right)^{-1} \right)$ .
2. Take

$$\boldsymbol{\eta}^{(t+1)} = \begin{cases} \boldsymbol{\eta}^* & \text{with probability } \min \left\{ \left( \frac{\pi_k^+(\boldsymbol{\eta}^*) q^*(\boldsymbol{\eta}^{(t)})}{\pi_k^+(\boldsymbol{\eta}^{(t)}) q^*(\boldsymbol{\eta}^*)} \right), 1 \right\} \\ \boldsymbol{\eta}^{(t)} & \text{otherwise} \end{cases},$$

where  $q^*(\cdot)$  is a density of  $N_k$  distribution in (25).