

Bernoulli 16(3), 2010, 679–704

DOI: 10.3150/09-BEJ233

Bayesian nonparametric estimation and consistency of mixed multinomial logit choice models

PIERPAOLO DE BLASI¹, LANCELOT F. JAMES² and JOHN W. LAU³

¹*Dipartimento di Statistica e Matematica Applicata and Collegio Carlo Alberto, Università degli Studi di Torino, corso Unione Sovietica 218/bis, 10134 Torino, Italy. E-mail: pierpaolo.deblasi@unito.it*

²*Department of Information Systems, Business Statistics and Operations Management, The Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong. E-mail: lancelot@ust.hk*

³*School of Mathematics and Statistics, University of Western Australia, 35 Stirling Highway, Crawley 6009, Western Australia, Australia. E-mail: john@maths.uwa.edu.au*

This paper develops nonparametric estimation for discrete choice models based on the mixed multinomial logit (MMNL) model. It has been shown that MMNL models encompass all discrete choice models derived under the assumption of random utility maximization, subject to the identification of an unknown distribution G . Noting the mixture model description of the MMNL, we employ a Bayesian nonparametric approach, using nonparametric priors on the unknown mixing distribution G , to estimate choice probabilities. We provide an important theoretical support for the use of the proposed methodology by investigating consistency of the posterior distribution for a general nonparametric prior on the mixing distribution. Consistency is defined according to an L_1 -type distance on the space of choice probabilities and is achieved by extending to a regression model framework a recent approach to strong consistency based on the summability of square roots of prior probabilities. Moving to estimation, slightly different techniques for non-panel and panel data models are discussed. For practical implementation, we describe efficient and relatively easy-to-use blocked Gibbs sampling procedures. These procedures are based on approximations of the random probability measure by classes of finite stick-breaking processes. A simulation study is also performed to investigate the performance of the proposed methods.

Keywords: Bayesian consistency; blocked Gibbs sampler; discrete choice models; mixed multinomial logit; random probability measures; stick-breaking priors

1. Introduction

Discrete choice models arise naturally in many fields of application, including marketing and transportation science. Such choice models are based on the neoclassical economic theory of random utility maximization (RUM). Given a finite set of choices $\mathbf{C} = \{1, \dots, J\}$, it is assumed that each individual has a utility function

$$U_j = \mathbf{x}'_j \boldsymbol{\beta} + \varepsilon_j \quad \text{for } j \in \mathbf{C}.$$

The values $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_J)$ are observed covariates, where $\mathbf{x}_j \in \mathbb{R}^d$ denote the covariates associated with each choice $\{j\} \in \mathbf{C}$, the coefficient $\boldsymbol{\beta}$ is an unknown (preference) vector in \mathbb{R}^d and

$(\varepsilon_1, \dots, \varepsilon_J)$ are random terms. Suppose that all U_j are distinct and that the individual makes a choice $\{j\}$ if and only if $U_j > U_l \forall l \neq j$. The introduction of the random error terms ε_j represents a departure from classical economic utility models. The random errors account for the discrepancy between the actual utility, which is known by the chooser, and that which is deduced by the experimenter who observes \mathbf{x} and the choice made by the individual. Hence, the deterministic statement of choice $\{j\}$ is replaced by the probability of choosing $\{j\}$, that is, $P\{U_j > U_l \forall l \neq j\}$. The analysis of such a model depends on the specifications of the errors. McFadden (1974) shows that the specification of independent Gumbel error terms leads to the tractable multinomial logit (MNL) model. This representation is written as

$$P(\{j\}|\boldsymbol{\beta}, \mathbf{x}) = \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}\}} \quad \text{for } j \in \mathbf{C}.$$

The MNL possesses the property of independence from irrelevant alternatives (IIA), which makes it inappropriate in many situations. The probit and the generalized extreme value models, which do not exhibit the IIA property and are models derived from dependent error structures, have been proposed as alternatives to the MNL. A drawback of the aforementioned procedures is that they are not robust against model misspecification.

The mixed multinomial logit (MMNL) model, first introduced by Cardell and Dunbar (1980), emerges as potentially the most attractive model. The book by Train (2003) includes a detailed discussion of this model. The general MMNL choice probabilities are defined by mixing an MNL model over a mixing distribution G . For a set of covariates \mathbf{x} , the MMNL model is written as

$$P(\{j\}|G, \mathbf{x}) = \int_{\mathbb{R}^d} \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}\}} G(d\boldsymbol{\beta}) \quad \text{for } j \in \mathbf{C}. \quad (1)$$

McFadden and Train (2000) establish the important result that, in theory, all RUM models can be captured by correct specification of G . Thus, a robust approach amounts to being able to employ statistical estimation methods based on a nonparametric assumption on G . However, statistical techniques have only been developed for the case where G is given a parametric form. The most popular model is when G is specified to be multivariate normal with unknown mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\tau}$:

$$P(\{j\}|\boldsymbol{\mu}, \boldsymbol{\tau}, \mathbf{x}) = \int_{\mathbb{R}^d} \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}\}} \phi(\boldsymbol{\beta}|\boldsymbol{\mu}, \boldsymbol{\tau}) d\boldsymbol{\beta} \quad \text{for } j \in \mathbf{C}, \quad (2)$$

where $\phi(\boldsymbol{\beta}|\boldsymbol{\mu}, \boldsymbol{\tau})$ represents a multivariate normal density with parameters $\boldsymbol{\mu}$ and $\boldsymbol{\tau}$. We shall refer to this as a Gaussian mixed logit (GML) model. Here, based on a sample of size n , one estimates the choice probabilities by estimating $\boldsymbol{\mu}$ and $\boldsymbol{\tau}$. Applications and discussions are found in, among others, Bhat (1998), Brownstone and Train (1999), Erdem (1996), Srinivasan and Mahmassani (2005) and Walker, Ben-Akiva and Bolduc (2007). Additionally, Dubé *et al.* (2002) provide a discussion focused on applications to marketing. The GML model is popular since it is flexible and relatively easy to estimate via simulated maximum likelihood techniques or via Bayesian MCMC procedures. Other choices for G include the lognormal and uniform distribu-

tions. Train (2003) discusses the merits and possible drawbacks of Bayesian MCMC procedures versus simulated maximum likelihood procedures for various choices of G . However, despite the attractive features of the GML, it does not encompass all RUM models, hence, it is not robust against misspecification.

In this article, we develop a nonparametric Bayesian method for the estimation of the choice probabilities and we prove consistency of the posterior distribution. The idea is to model the mixing distribution G via a random probability measure in order to fully exploit the flexibility of the MMNL model. Many nonparametric priors are currently available for modeling G , such as stick-breaking priors, normalized random measures with independent increments and Dirichlet process mixtures. We establish consistency of the posterior distribution of G under neat sufficient conditions which are readily verifiable for all of these nonparametric priors. Consistency is defined according to an L_1 -type distance on the space of choice probabilities by exploiting the square root approach to strong consistency of Walker (2003a, 2004). We essentially show that the Bayesian MMNL model is consistent if the prior on G has the true mixing distribution in its weak support and satisfies a mild condition on the tails of the prior predictive distribution. We then move to estimation and divide our discussion into methods for non-panel and panel data. Specifically, for non-panel data models, we use, as a prior for G , a mixture of Dirichlet processes. Methods for panel data instead involve a Dirichlet mixture of normal densities. For practical implementation, we describe efficient and relatively easy-to-use blocked Gibbs sampling procedures, developed in Ishwaran and Zarepour (2000) and Ishwaran and James (2001).

The rest of the paper is organized as follows. In Section 2, we describe the Bayesian nonparametric approach by placing a nonparametric prior on the mixing distribution and present the consistency result for the posterior distribution of G . In Section 3, we show how to implement a blocked Gibbs sampling for drawing inference when a discrete nonparametric prior is used. Section 4 deals with panel data with similar Bayesian nonparametric methods, where we define a class of priors for G that preserves the distinct nature of individual preferences and specialize the blocked Gibbs sampler to this setting. In Section 5, we provide an illustrative simulation study which shows the flexibility and good performance of our procedures. Finally, in Section 6, we provide a detailed proof of consistency.

2. Bayesian MMNL models

A Bayesian nonparametric MMNL model is specified by placing a nonparametric prior on the mixing distribution G in (1):

$$P(\{j\}|\tilde{G}, \mathbf{x}) = \int_{\mathbb{R}^d} \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}\}} \tilde{G}(d\boldsymbol{\beta}) \quad \text{for } j \in \mathbf{C}. \tag{3}$$

Here, \tilde{G} denotes a random probability measure which takes values over the space \mathbb{P} of probability measures on \mathbb{R}^d , the former endowed with the weak topology. The nonparametric distribution of

\tilde{G} is denoted by \mathcal{P} . Model (3) can be equivalently expressed in hierarchical form as

$$\begin{aligned}
 Y_i | \beta_i &\stackrel{\text{ind}}{\sim} \frac{\exp\{\mathbf{x}'_{iY_i} \beta_i\}}{\sum_{j \in \mathbf{C}} \exp\{\mathbf{x}'_{ij} \beta_i\}} && \text{for } i = 1, \dots, n \text{ and } Y_i \in \mathbf{C}, \\
 \beta_i | \tilde{G} &\stackrel{\text{iid}}{\sim} \tilde{G} && \text{for } i = 1, \dots, n, \\
 \tilde{G} &\sim \mathcal{P}
 \end{aligned}
 \tag{4}$$

with $\mathbf{x}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{ij})$ the covariates and Y_i the choice observed for individual i .

One can choose \tilde{G} to be a Dirichlet process (Ferguson (1973)), although there currently exist other nonparametric priors that can be used, like stick-breaking priors (Ishwaran and James (2001)) and normalized random measure with independent increments (NRMI) (Regazzini, Lijoi and Prünster (2003)). All of these priors select discrete distributions almost surely (a.s.), whereas random probability measures whose support contains continuous distributions can be obtained by using a Dirichlet process mixture of densities, in the spirit of Lo (1984). An important role in the sequel will be played by the prior predictive distribution of \tilde{G} , denoted by H , which is an element of \mathbb{P} and is defined by

$$H(B) := E[\tilde{G}(B)]
 \tag{5}$$

for all Borel sets B of \mathbb{R}^d , where $E(\cdot)$ denotes expectation. In the next section, we show that an essential condition for consistency of the posterior distribution is expressed in terms of H . This yields an easy-to-use criterion for the choice of the prior for \tilde{G} as H is readily obtained for all of the nonparametric priors listed above. Furthermore, one can embed a parametric model, such as the GML, within the nonparametric framework via a suitable specification of the distribution H .

2.1. Posterior consistency

Bayesian consistency deals with the asymptotic behavior of posterior distributions with respect to repeated sampling. The problem can be set in general terms as follows: suppose the existence of a “true” unknown distribution P_0 that generates the data, then check whether the posterior accumulates in suitably-defined neighborhoods of P_0 . There exist two main approaches to the study of strong consistency, that is, consistency when the neighborhood of P_0 is defined according to the Hellinger metric on the space of density functions. One is based on the metric entropy of the parameter space and was set forth in Barron, Schervish and Wasserman (1999) and Ghosal, Ghosh and Ramamoorthi (1999). The second approach was introduced by Walker (2003a, 2004) and has more of a Bayesian flavor, in the sense that it relies on the summability of square roots of prior probabilities. For discussion, the reader is referred to Wasserman (1998), Walker, Lijoi and Prünster (2005) and Choudhuri, Ghosal and Roy (2005). Strong consistency in mixture models for density estimation is addressed by Ghosal, Ghosh and Ramamoorthi (1999) and Lijoi, Prünster and Walker (2005), by using the metric entropy approach and the square root approach, respectively. As for the non-identically distributed case, we mention Choi and Schervish (2007) and Ghosal and Roy (2006), both of which follow the metric entropy approach. The square root

approach is adopted by Walker (2003b) for nonparametric regression models and by Ghosal and Tang (2006) for the estimation of transition densities in the context of Markov processes.

We face the issue of consistency for the MMNL model (3) by exploiting the square root approach of Walker and its variation proposed in Lijoi, Prünster and Walker (2005) which makes use of metric entropy in an instrumental way. We assume the existence of a $G_0 \in \mathbb{P}$ such that the true distribution of Y given $\mathbf{X} = \mathbf{x}$ is given by

$$P_0(\{j\}|\mathbf{x}) = \int_{\mathbb{R}^d} \frac{\exp(\mathbf{x}'_j \boldsymbol{\beta})}{\sum_{l \in \mathbf{C}} \exp(\mathbf{x}'_l \boldsymbol{\beta})} G_0(d\boldsymbol{\beta}).$$

The variables \mathbf{X}_i are taken as independent draws from a common distribution $M(d\mathbf{x})$ which is supported on $\mathcal{X} \subset \mathbb{R}^{Jd}$. The distribution of an infinite sequence $(Y_i, \mathbf{X}_i)_{i \geq 1}$ will be then denoted by $P_{(G_0, M)}^\infty$. Finally, let \mathcal{P}_n denote the posterior distribution of \tilde{G} given $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$; see also equation (19) in Section 6. In the sequel, we take the covariate distribution M to be a fixed quantity so that the posterior distribution does not depend on the specific form of M . Note, however, that the posterior evaluation is also not affected when M is considered as a parameter with an independent prior since it is reasonable to assume that the choice probabilities are unrelated to M .

We give conditions on G_0 and the prior predictive distribution of \tilde{G} such that the posterior distribution \mathcal{P}_n concentrates all probability mass in neighborhoods of G_0 defined according to strong consistency of choice probabilities. To this end, we look at the vector of choice probabilities as a vector-valued function $\mathbf{q}: \mathcal{X} \rightarrow \Delta$, where Δ is the J -dimensional probability simplex. We define

$$\mathbf{q}(\mathbf{x}; G) = [P(\{1\}|G, \mathbf{x}), \dots, P(\{J\}|G, \mathbf{x})] \tag{6}$$

for any $G \in \mathbb{P}$. On the space $\mathcal{Q} = \{\mathbf{q}(\cdot; G) : G \in \mathbb{P}\}$, we define the L_1 -type distance

$$d(\mathbf{q}_1, \mathbf{q}_2) = \int_{\mathcal{X}} |\mathbf{q}_1(\mathbf{x}) - \mathbf{q}_2(\mathbf{x})| M(d\mathbf{x}), \tag{7}$$

where $|\cdot|$ denotes the Euclidean norm in Δ .

Definition 1. \mathcal{P} is consistent at G_0 if, for any $\epsilon > 0$,

$$\mathcal{P}_n\{G : d(\mathbf{q}(\cdot; G), \mathbf{q}(\cdot; G_0)) > \epsilon\} \rightarrow 0, \quad P_{(G_0, M)}^\infty\text{-a.s.}$$

The main result is stated in the following theorem.

Theorem 1. Let \mathcal{P} be a prior on \mathbb{P} with predictive distribution H and G_0 be in the weak support of \mathcal{P} . Suppose that \mathcal{X} is a compact subset of \mathbb{R}^{Jd} . If

- (i) $P_0(\{j\}|\mathbf{x}) > 0$ for any $j \in \mathbf{C}$ and $\mathbf{x} \in \mathcal{X}$;
- (ii) $\int_{\mathbb{R}^d} |\boldsymbol{\beta}| H(d\boldsymbol{\beta}) < +\infty$,

then \mathcal{P} is consistent at G_0 .

The compactness of the covariate space is a standard assumption in nonparametric regression problems. Condition (i) is fairly reasonable since it is guaranteed by a correct specification of the RUM model: one can always redefine the set of choices or the covariate space to fulfill this requirement. Moreover, because of the compactness of \mathcal{X} , condition (i) implies that G_0 is a proper distribution on \mathbb{R}^d , that is, with no masses escaping at infinity. The verification that G_0 belongs to the weak support of \mathcal{P} is then an easy task: in general, it is sufficient that the prior predictive distribution H has full support on \mathbb{R}^d . Condition (ii) is a mild condition on the tails of H : it is satisfied by any distribution with tails lighter than the Cauchy distribution.

2.2. Illustration

It is worth considering condition (ii) in more detail for a variety of Bayesian MMNL models, obtained from different specifications of \mathcal{P} . If \tilde{G} is taken to be a Dirichlet process with base measure $\alpha = aF$, where $a > 0$ is a constant and $F \in \mathbb{P}$, then F coincides with H in (5). A larger class of Bayesian MMNL models arise when \tilde{G} is chosen to be a stick-breaking prior:

$$\tilde{G}(\cdot) = \sum_{k \geq 1} p_k \delta_{Z_k}(\cdot), \tag{8}$$

where the p_k are positive random probabilities chosen to be independent of Z_k and such that $\sum_{k \geq 1} p_k = 1$ a.s. The Z_k are random locations taken as independent draws from some nonatomic distribution F in \mathbb{P} . What characterizes a stick-breaking prior is that the random weights are expressible as $p_k = V_k \prod_{i=1}^{k-1} (1 - V_i)$, where the V_k are independent beta-distributed random variables of parameters $a_k, b_k > 0$; we write $V_k \sim \text{beta}(a_k, b_k)$. Examples of random probability measures in this class are given in Ishwaran and James (2001); see also Pitman and Yor (1997) and Ishwaran and Zarepour (2000). They represent extensions of the Dirichlet process, which has $a_k = 1$ and $b_k = a \forall k$, and they all have in common that the prior predictive distribution H coincides with F .

The class of NRMI is another valid choice for \mathcal{P} . Specifically, one can take $\tilde{G}(\cdot) = \tilde{\mu}(\cdot) / \tilde{\mu}(\mathbb{R}^d)$, where $\tilde{\mu}$ is a completely random measure with Poisson intensity measure $\nu(dv, dz) = \rho(dv|z)\alpha(dz)$ on $(0, +\infty) \times \mathbb{R}^d$. Here, $\rho(\cdot|z)$ is a Lévy density on $(0, +\infty)$ for any z and α is a finite measure on \mathbb{R}^d such that $\psi(u) := \int_{\mathbb{R}^d \times \mathbb{R}^+} (1 - e^{-uv})\rho(dv|z)\alpha(dz) < \infty$, which is needed to guarantee that $\tilde{\mu}(\mathbb{R}^d) < \infty$ a.s. It can be shown that $H(B) = \int_B \int_0^{+\infty} e^{-\psi(u)} \{ \int_0^{+\infty} e^{-uv} \times v \rho(dv|z) \} du \alpha(dz)$ for any Borel set B of \mathbb{R}^d ; see also James, Lijoi and Prünster (2009). When $\rho(dv|z) = \rho(dv)$ for each z (homogeneous case), the prior predictive distribution reduces to

$$H(B) = \frac{\alpha(B)}{\alpha(\mathbb{R}^d)} \quad \text{for any Borel } B \subset \mathbb{R}^d. \tag{9}$$

The homogeneous NRMI includes, as a special case, the Dirichlet process and belongs, together with the stick-breaking priors, to the class of *species sampling models*, for which (9) holds for some finite measure α . Note that all of the nonparametric priors belonging to this class allow an easy verification of condition (ii).

The specification of the nonparametric prior in terms of a base measure α , as in (9), allows more flexibility to be introduced via an additional level in the hierarchical structure (4). If we let the base measure be indexed by a parameter θ , say α_θ , and θ be random with probability density $\pi(\theta)$ on some Euclidean space Θ , then we obtain a mixture of Dirichlet process in the spirit of Antoniak (1974). Condition (ii) must then be verified for the convolution

$$H(B) = \int_{\Theta} \int_B H_\theta(dz) \pi(\theta) d\theta, \quad \text{where } H_\theta(dz) = \frac{\alpha_\theta(dz)}{\alpha_\theta(\mathbb{R}^d)}. \tag{10}$$

It is quite straightforward to check that condition (ii) holds for the mixture of Dirichlet processes implemented in the analysis of non-panel data in Section 3.

Finally, consider the case of Dirichlet process mixture models of Lo (1984), where \tilde{G} is absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^d with random density function specified as $\int_{\Theta} K(\boldsymbol{\beta}, \theta) \tilde{\Pi}(d\theta)$. Here, $K(\boldsymbol{\beta}, \theta)$ is a non-negative kernel defined on $\mathbb{R}^d \times \Theta$ such that, for each $\theta \in \Theta$, $\int_{\mathbb{R}^d} K(z, \theta) dz = 1$, while $\tilde{\Pi}$ is a Dirichlet process prior with base measure aF and F a probability measure on Θ . The distribution H is then absolutely continuous and is given by

$$H(B) = \int_B \int_{\Theta} K(z, \theta) F(d\theta) dz.$$

As in (10), verifying condition (ii) requires a study of the tail properties of a convolution, this time of $K(z, \theta)$ with respect to $F(d\theta)$. In the analysis of panel data (see Section 4), we adopt a Dirichlet mixture model as continuous nonparametric prior for \tilde{G} where the verification of condition (ii) can be readily established.

3. Implementation for non-panel data

Assume that we have a single observation for each individual and that we want to account for the possibility of ties among different individuals' preferences. Therefore, we use a discrete non-parametric prior for the mixing distribution. Take \tilde{G} to be a Dirichlet process with base measure aF and denote its law by $\mathcal{P}(dG|aF)$ (although the treatment can be easily extended to any other stick-breaking prior). Representation (8) then holds with random probabilities p_1, p_2, \dots at locations Z_1, Z_2, \dots , which are i.i.d. draws from F . This translates into a Bayesian model for the MMNL as

$$P(\{j\}|\tilde{G}, \mathbf{x}) = \sum_{k \geq 1} p_k \frac{\exp\{\mathbf{x}'_j Z_k\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l Z_k\}} \quad \text{for } j \in \mathbf{C}. \tag{11}$$

One can then center \tilde{G} on a parametric model like the GML in (2) by taking F to have normal density $\phi(\boldsymbol{\beta}|\boldsymbol{\mu}, \boldsymbol{\tau})$. In a parametric Bayesian framework, by placing priors on $\boldsymbol{\mu}, \boldsymbol{\tau}$, one is able to get posterior estimates of $\boldsymbol{\mu}, \boldsymbol{\tau}$, but inference is restricted to the assumption of the GML model. The flexibility of the Bayesian nonparametric approach allows one to choose F based on convenience and ease of use and to utilize, for instance, the attractive features of GML models while still maintaining the robustness of a nonparametric approach.

In the case of the Dirichlet process, the parameters associated with F , for instance, $\boldsymbol{\mu}$ and $\boldsymbol{\tau}$, are considered fixed. As observed in Section 2, one can introduce more flexibility in the model by treating such parameters as random. Specifying $\theta = (\boldsymbol{\mu}, \boldsymbol{\tau})$, $F_\theta(d\boldsymbol{\beta})$ to have density $\phi(\boldsymbol{\beta}|\theta) d\boldsymbol{\beta}$ and $\pi(\theta)$ to be the density function for θ , the law of \tilde{G} is given by the mixture $\int_{\Theta} \mathcal{P}(dG|aF_\theta)\pi(d\theta)$. Equivalently, using (8), a mixture of Dirichlet processes is defined by specifying each $Z_k|\theta$ to be i.i.d. F_θ . Note that, conditional on θ , a prior guess for the choice probabilities is

$$E[\mathbb{P}(\{j\}|\tilde{G}, \mathbf{x})|\theta] = \int_{\mathbb{R}^d} \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}\}} F_\theta(d\boldsymbol{\beta}) \quad \text{for } j \in \mathbf{C}. \tag{12}$$

By the properties of the Dirichlet process, the prediction rule for the choice probabilities given $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n$ is given by

$$\begin{aligned} E[\mathbb{P}(\{j\}|\tilde{G}, \mathbf{x})|\theta, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n] \\ = \frac{a}{a+n} \mathbb{P}(\{j\}|F_\theta, \mathbf{x}) + \sum_{i=1}^n \frac{1}{a+n} \frac{\exp\{\mathbf{x}'_j \boldsymbol{\beta}_i\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l \boldsymbol{\beta}_i\}}, \end{aligned} \tag{13}$$

where $\mathbb{P}(\{j\}|F_\theta, \mathbf{x}) := E[\mathbb{P}(\{j\}|\tilde{G}, \mathbf{x})|\theta]$ is given in (12) with a notation consistent with (1). However, the variables $\boldsymbol{\beta}_i$ are not observable and hence one needs to implement computational procedures to draw from their posterior distribution.

In this framework, a reasonable algorithm to use is the blocked Gibbs sampler developed in Ishwaran and Zarepour (2000) and Ishwaran and James (2001). Indeed, since the multinomial logistic kernel does not form a conjugate pair for $\boldsymbol{\beta}$, marginal algorithms suffer from slow convergence, although strategies for overcoming this problem can be found in MacEachern and Muller (1998).

3.1. Blocked Gibbs algorithm

In this section, we discuss how to implement a blocked Gibbs sampling algorithm for drawing inference on a nonparametric hierarchical model with the structure

$$\begin{aligned} Y_i | \boldsymbol{\beta}_i &\stackrel{\text{iid}}{\sim} L(Y_i, \boldsymbol{\beta}_i) \quad \text{for } i = 1, \dots, n \text{ and } Y_i \in \mathbf{C}, \\ \boldsymbol{\beta}_i | \tilde{G} &\stackrel{\text{iid}}{\sim} \tilde{G} \quad \text{for } i = 1, \dots, n, \\ \tilde{G} | \theta &\sim \mathcal{P}(dG|aF_\theta), \\ \theta &\sim \pi(d\theta), \end{aligned} \tag{14}$$

where $L(Y_i, \boldsymbol{\beta}) = \exp\{\mathbf{x}'_{Y_i} \boldsymbol{\beta}\} / \sum_{j \in \mathbf{C}} \exp\{\mathbf{x}'_j \boldsymbol{\beta}\}$ is the probability for Y_i conditional on $\boldsymbol{\beta}_i$. The blocked Gibbs sampler utilizes the fact that a truncated Dirichlet process, discussed in Ishwaran and Zarepour (2000) and Ishwaran and James (2001), serves as a good approximation to the

random probability measure $\tilde{G}|\theta$ in (14). We replace the conditional law $\mathcal{P}(dG|aF_\theta)$ with the law of the random probability measure

$$\tilde{G}(\cdot) = \sum_{k=1}^N p_k \delta_{Z_k}(\cdot), \quad 1 \leq N < \infty, \tag{15}$$

where $Z_k|\theta$ are i.i.d. F_θ and the random probabilities p_1, \dots, p_N are defined by the stick-breaking construction

$$p_1 = V_1 \quad \text{and} \quad p_k = (1 - V_1) \cdots (1 - V_{k-1}) V_k, \quad k = 2, \dots, N, \tag{16}$$

with V_1, V_2, \dots, V_{N-1} i.i.d. beta(1, a) and $V_N = 1$, which ensures that $\sum_{k=1}^N p_k = 1$. The law of $\tilde{G}|\theta$ in (15) is referred to as a *truncated Dirichlet process* and will be denoted $\mathcal{P}^N(dG|aF_\theta)$. Moreover, the limit as $N \rightarrow \infty$ will converge to a random probability measure with law $\mathcal{P}(dG|aF_\theta)$. Indeed, the method yields an accurate approximation of the Dirichlet process for N moderately large since the truncation is exponentially accurate. Theorem 2 in Ishwaran and James (2001) provides an L_1 -error bound for the approximation of conditional density of $\mathbf{Y} = (Y_1, \dots, Y_n)$ given θ . Let

$$\mu^N(\mathbf{Y}|\theta) = \int \left[\prod_{i=1}^n \int_{\mathbb{R}^d} L(Y_i, \boldsymbol{\beta}_i) G(d\boldsymbol{\beta}_i) \right] \mathcal{P}^N(dG|aF_\theta)$$

and $\mu(\mathbf{Y}|\theta)$ be its limit under the prior $\mathcal{P}(dG|aF_\theta)$. One then has

$$\|\mu^N - \mu\|_1 := \int |\mu^N(\mathbf{Y}|\theta) - \mu(\mathbf{Y}|\theta)| d\mathbf{Y} \sim 4ne^{-(N-1)/a},$$

where the integral above is considered over the counting measure on the n -fold product space \mathbf{C}^n . Moreover, Corollary 1 in Ishwaran and James (2002) can be used to show that the truncated Dirichlet process also leads to asymptotic approximations to the posterior that are exponentially accurate.

The key to working with random probability measures like (15) is that it allows blocked updates to be performed for $\mathbf{p} = (p_1, \dots, p_n)$ and $\mathbf{Z} = (Z_1, \dots, Z_n)$ by recasting the hierarchical model (14) completely in terms of random variables. To this aim, define the classification variables $\mathbf{K} = \{K_1, \dots, K_n\}$ such that, conditional on \mathbf{p} , each K_i is independent with distribution

$$P\{K_i \in \cdot | \mathbf{p}\} = \sum_{k=1}^N p_k \delta_k(\cdot).$$

That is, $P\{K_i = k | \mathbf{p}\} = p_k$ for $k = 1, \dots, N$ so that K_i identifies the Z_k associated with each $\boldsymbol{\beta}_i$: $\boldsymbol{\beta}_i = Z_{K_i}$. In this setting, a sample $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n$ from (15) produces $n_0 \leq \min(n, N)$ distinct values. The blocked Gibbs algorithm is based on sampling $\mathbf{K}, \mathbf{p}, \mathbf{Z}, \theta$ from the distribution proportional to

$$\left[\prod_{i=1}^n L(Y_i, \boldsymbol{\beta}_i) \right] \left[\prod_{i=1}^n \sum_{k=1}^N p_k \delta_{Z_k}(d\boldsymbol{\beta}_i) \right] \pi(\mathbf{p}) \left[\prod_{k=1}^N F_\theta(dZ_k) \right] \pi(d\theta),$$

where $\pi(\mathbf{p})$ denotes the distribution of \mathbf{p} defined in (16). This augmented likelihood is an expression of the augmented density when $\mathcal{P}(dG|aF_\theta)$ is replaced by $\mathcal{P}^N(dG|aF_\theta)$.

Before describing the algorithm, we specify choices for F_θ and θ which agree with the GML model. Set $\theta = (\boldsymbol{\mu}, \boldsymbol{\tau})$ and specify the density of F_θ to be $\phi(\boldsymbol{\beta}|\boldsymbol{\mu}, \boldsymbol{\tau})$. Let λ denote a positive scalar. We choose a multivariate normal inverse Wishart distribution for $\boldsymbol{\mu}, \boldsymbol{\tau}$, where, specifically, $\boldsymbol{\mu}|\boldsymbol{\tau}$ is a multivariate normal vector with mean parameter \mathbf{m} and scaled covariance matrix $\lambda^{-1}\boldsymbol{\tau}$ and $\boldsymbol{\tau}$ is drawn from an inverse Wishart distribution with degrees of freedom ν_0 and scale matrix \mathbf{S}_0 . We denote this distribution for $\boldsymbol{\mu}, \boldsymbol{\tau}$ as N-IW($\mathbf{m}, \lambda^{-1}\boldsymbol{\tau}, \nu_0, \mathbf{S}_0$). Our specification is similar to that used in Train (2003), Chapter 12, for a parametric GML model for panel data.

Algorithm 1.

1. *Conditional draw for \mathbf{K} .* Independently sample K_i according to $P\{K_i \in \cdot | \mathbf{p}, \mathbf{Z}, \mathbf{Y}\} = \sum_{k=1}^N p_{k,i} \delta_k(\cdot)$ for $i = 1, \dots, n$, where

$$(p_{1,i}, \dots, p_{N,i}) \propto (p_1 L(Y_i, Z_1), \dots, p_N L(Y_i, Z_N)).$$

2. *Conditional draw for \mathbf{p} .* $p_1 = V_1^*, p_k = (1 - V_1^*) \dots (1 - V_{k-1}^*) V_k^*, k = 2, \dots, N - 1$ and $V_N^* = 1$, where, if e_k records the number of K_i values which equal k ,

$$V_k^* \stackrel{\text{ind}}{\sim} \text{beta}\left(1 + e_k, a + \sum_{l=k+1}^N e_l\right), \quad k = 1, \dots, N - 1.$$

3. *Conditional draw for \mathbf{Z} .* Let $\{K_1^*, \dots, K_{n_0}^*\}$ denote the unique set of K_i values. For each $k \notin \{K_1^*, \dots, K_{n_0}^*\}$, draw $Z_k|\boldsymbol{\mu}, \boldsymbol{\tau}$ from the prior multivariate normal density $\phi(Z|\boldsymbol{\mu}, \boldsymbol{\tau})$. For $j = 1, \dots, n_0$, draw $Z_{K_j^*} := \boldsymbol{\beta}_j^*$ from the density proportional to $\phi(\boldsymbol{\beta}_j^*|\boldsymbol{\mu}, \boldsymbol{\tau}) \prod_{\{i:K_i=K_j^*\}} L(Y_i, \boldsymbol{\beta}_j^*)$ by using, for example, a standard Metropolis–Hastings procedure.
4. *Conditional draw for $\theta = (\boldsymbol{\mu}, \boldsymbol{\tau})$.* Conditional on $\boldsymbol{\tau}, \mathbf{K}, \mathbf{Z}, \mathbf{Y}$, draw $\boldsymbol{\mu}$ from a multivariate normal distribution with parameters

$$\frac{\lambda \mathbf{m} + n_0 \bar{\boldsymbol{\beta}}_{n_0}}{\lambda + n_0} \quad \text{and} \quad \frac{\boldsymbol{\tau}}{\lambda + n_0},$$

where $\bar{\boldsymbol{\beta}}_{n_0} = n_0^{-1} \sum_{j=1}^{n_0} \boldsymbol{\beta}_j^*$. Conditional on $\mathbf{K}, \mathbf{Z}, \mathbf{Y}$, draw $\boldsymbol{\tau}$ from an inverse Wishart distribution with parameters

$$\nu_0 + n_0 \quad \text{and} \quad \frac{\nu_0 \mathbf{S}_0 + n_0 \mathbf{S}_{n_0} + R(\bar{\boldsymbol{\beta}}_{n_0}, \mathbf{m})}{\nu_0 + n_0},$$

where

$$\mathbf{S}_{n_0} = \frac{1}{n_0} \sum_{j=1}^{n_0} (\boldsymbol{\beta}_j^* - \bar{\boldsymbol{\beta}}_{n_0})(\boldsymbol{\beta}_j^* - \bar{\boldsymbol{\beta}}_{n_0})' \quad \text{and} \quad R(\bar{\boldsymbol{\beta}}_{n_0}, \mathbf{m}) = \frac{\lambda n_0}{\lambda + n_0} (\bar{\boldsymbol{\beta}}_{n_0} - \mathbf{m})(\bar{\boldsymbol{\beta}}_{n_0} - \mathbf{m})'.$$

Notice that, when $n_0 = 1$, Steps 3 and 4 reduce to the MCMC steps for a parametric Bayesian model. Iterating the steps above produces a draw from the distribution $\mathbf{Z}, \mathbf{K}, \mathbf{p}, \theta | \mathbf{Y}$. Thus, each iteration m defines a probability measure $G^{(m)}(\cdot) = \sum_{k=1}^N p_k^{(m)} \delta_{Z_k^{(m)}}(\cdot)$, which eventually approximates draws from the posterior distribution of $\tilde{G} | \mathbf{Y}$. Consequently, one can approximate the posterior distributional properties of the choice probabilities $P(\{j\} | \tilde{G}, \mathbf{x})$ by constructing (iteratively)

$$P(\{j\} | G^{(m)}, \mathbf{x}) = \sum_{k=1}^N p_k^{(m)} \frac{\exp\{\mathbf{x}'_j Z_k^{(m)}\}}{\sum_{l \in \mathbf{C}} \exp\{\mathbf{x}'_l Z_k^{(m)}\}};$$

see (11). For instance, an histogram of the $P(\{j\} | G^{(m)}, \mathbf{x})$, for $m = 1, \dots, M$, approximates the posterior distribution. An approximation to the posterior mean $E[P(\{j\} | \tilde{G}, \mathbf{x}) | \mathbf{Y}]$ is obtained by $M^{-1} \sum_{m=1}^M P(\{j\} | G^{(m)}, \mathbf{x})$ or, alternatively, by

$$\hat{P}(\{j\} | \mathbf{x}) := \frac{1}{M} \sum_{m=1}^M E[P(\{j\} | \tilde{G}, \mathbf{x}) | \theta^{(m)}, \boldsymbol{\beta}_1^{(m)}, \dots, \boldsymbol{\beta}_n^{(m)}], \tag{17}$$

where $E[P(\{j\} | \tilde{G}, \mathbf{x}) | \theta, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n]$ is given in (13) and $\boldsymbol{\beta}_i^{(m)} = Z_{K_i^{(m)}}^{(m)}$.

4. Bayesian modeling for panel data

The MMNL framework may also be used to model choice probabilities based on panel data. In the panel data setting, each individual i is observed to make a sequence of choices at different time points. The random utility for choosing j for individual i in choice situation t is given by

$$U_{ijt} = \mathbf{x}'_{ijt} \boldsymbol{\beta}_i + \varepsilon_{ijt}, \quad j \in \mathbf{C},$$

for times $t = 1, \dots, T_i$. The MMNL model can be described as follows [see Train (2003), Section 6.7]: given $\boldsymbol{\beta}_i$, the probability that a person makes the sequence of choices $\mathbf{Y}_i = \{Y_{i1}, \dots, Y_{iT_i}\}$ is the product of logit formulae

$$L(\mathbf{Y}_i, \boldsymbol{\beta}_i) = \prod_{t=1}^{T_i} \frac{\exp\{\mathbf{x}'_{iY_{it}} \boldsymbol{\beta}_i\}}{\sum_{j \in \mathbf{C}} \exp\{\mathbf{x}'_{ijt} \boldsymbol{\beta}_i\}}.$$

The MMNL model is completed by taking the $\boldsymbol{\beta}_i$ to be from a distribution G so that the unconditional choice probability is specified by

$$P(\mathbf{Y}_i | G, \mathbf{x}_i) = \int_{\mathbb{R}^d} \prod_{t=1}^{T_i} \frac{\exp\{\mathbf{x}'_{iY_{it}} \boldsymbol{\beta}\}}{\sum_{j \in \mathbf{C}} \exp\{\mathbf{x}'_{ijt} \boldsymbol{\beta}\}} G(d\boldsymbol{\beta}) = \int_{\mathbb{R}^d} L(\mathbf{Y}_i, \boldsymbol{\beta}) G(d\boldsymbol{\beta}),$$

where $\mathbf{x}_i = \{\mathbf{x}_{ijt}, j \in \mathbf{C}, t = 1, \dots, T_i\}$ denotes the array of covariates associated with the sequence of choices of individual i . Similarly to the non-panel data setting, we wish to model G as a random probability measure in a Bayesian framework. While it is possible to choose \tilde{G} to follow a Dirichlet process, this would result in possible ties among the individual's preferences β_i . In order to preserve the distinct nature of each individual's preference, we assume that, given \tilde{G} , the β_i are i.i.d. with distribution \tilde{G} , where \tilde{G} is a mixture of multivariate normal distributions with random mixing distribution $\tilde{\Pi}$. That is, \tilde{G} has random density $\int_{\Theta} \phi(\beta|\mu, \tau)\tilde{\Pi}(d\mu, d\tau)$, where $\Theta = \mathbb{R}^d \times \mathcal{S}$ with \mathcal{S} the space of covariance matrices. Specifically, we take $\tilde{\Pi}$ to be a Dirichlet process with shape aF , F a probability measure on Θ . Hence, the Bayesian MMNL model for individual i is expressible as

$$P(\mathbf{Y}_i|\tilde{G}, \mathbf{x}_i) = \int_{\mathbb{R}^d} L(\mathbf{Y}_i, \beta)\tilde{G}(d\beta) = \int_{\mathbb{R}^d} \int_{\Theta} L(\mathbf{Y}_i, \beta)\phi(\beta|\mu, \tau)\tilde{\Pi}(d\mu, d\tau) d\beta.$$

While one may use any choice for F , we take $F(d\mu, d\tau)$ to be the multivariate normal inverse Wishart distribution $N\text{-IW}(\mathbf{m}, \lambda^{-1}\tau, \mathbf{S}_0, \nu_0)$ described in Section 3.

4.1. Blocked Gibbs algorithm for panel data

The explicit posterior analysis for the panel data case is quite similar to the non-panel case. The main difference is that the $(\mu_i, \tau_i), i = 1, \dots, n$, rather than β_1, \dots, β_n , are drawn from the Dirichlet process. Here, we will briefly focus on the relevant data structure and then proceed to a description of how to implement the blocked Gibbs sampler. The joint distribution of the augmented data can be expressed using a hierarchical model as follows:

$$\begin{aligned} \mathbf{Y}_i|\beta_i &\stackrel{\text{ind}}{\sim} L(\mathbf{Y}_i, \beta_i) && \text{for } i = 1, \dots, n \text{ and } Y_{it} \in \mathbf{C}, \\ \beta_i|\mu_i, \tau_i &\stackrel{\text{ind}}{\sim} \phi(\beta_i|\mu_i, \tau_i) && \text{for } i = 1, \dots, n, \\ \mu_i, \tau_i|\tilde{\Pi} &\stackrel{\text{iid}}{\sim} \tilde{\Pi} && \text{for } i = 1, \dots, n, \\ \tilde{\Pi} &\sim \mathcal{P}(d\Pi|aF). \end{aligned} \tag{18}$$

Similar to the non-panel case, the blocked Gibbs sampler works by using the $\mathcal{P}^N(d\Pi|aF)$ in place of the law of the Dirichlet process $\mathcal{P}(d\Pi|aF)$. We now sample $(\mathbf{K}, \mathbf{p}, \mathbf{Z}, \beta_1, \dots, \beta_n)$ from the distribution proportional to

$$\left[\prod_{i=1}^n L(\mathbf{Y}_i, \beta_i)\phi(\beta_i|\mu_i, \tau_i) \right] \left[\prod_{i=1}^n \sum_{k=1}^N p_k \delta_{Z_k}(d\mu_i, d\tau_i) \right] \pi(\mathbf{p}) \prod_{k=1}^N F(dZ_k).$$

Here, we use the fact that $(\mu_i, \tau_i) = Z_{K_i}$ for $i = 1, \dots, n$. To approximate the posterior law of various functionals, we cycle through the following steps.

Algorithm 2.

1. *Conditional draw for \mathbf{K} .* Independently sample K_i according to

$$P\{K_i \in \cdot | \mathbf{p}, \mathbf{Z}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n, \mathbf{Y}\} = \sum_{k=1}^N p_{k,i} \delta_k(\cdot) \quad \text{for } i = 1, \dots, n,$$

where $(p_{1,i}, \dots, p_{N,i}) \propto (p_1 \phi(\boldsymbol{\beta}_i | Z_1), \dots, p_N \phi(\boldsymbol{\beta}_i | Z_N))$.

2. *Conditional draw for \mathbf{p} .* $p_1 = V_1^*$, $p_k = (1 - V_1^*) \cdots (1 - V_{k-1}^*) V_k^*$, $k = 2, \dots, N - 1$ and $V_N^* = 1$, where, if e_k records the number of K_i values which equal k ,

$$V_k^* \stackrel{\text{ind}}{\sim} \text{beta} \left(1 + e_k, a + \sum_{l=k+1}^N e_l \right), \quad k = 1, \dots, N - 1.$$

3. *Conditional draw for \mathbf{Z} .* Let $\{K_1^*, \dots, K_{n_0}^*\}$ denote the unique set of K_i values. For each $k \notin \{K_1^*, \dots, K_{n_0}^*\}$, draw $Z_k = (\boldsymbol{\mu}_k, \boldsymbol{\tau}_k)$ from the prior N-IW($\mathbf{m}, \lambda^{-1} \boldsymbol{\tau}, \mathbf{S}_0, \nu_0$). For $j = 1, \dots, n_0$, draw $Z_{K_j^*} := (\boldsymbol{\mu}_j^*, \boldsymbol{\tau}_j^*)$ as follows: (a) conditional on $\boldsymbol{\tau}_j^*, \mathbf{K}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n, \mathbf{Y}$, draw $\boldsymbol{\mu}_j^*$ from a multivariate normal distribution with parameters

$$\frac{\lambda \mathbf{m} + e_{K_j^*} \bar{\boldsymbol{\beta}}_j^*}{\lambda + e_{K_j^*}} \quad \text{and} \quad \frac{\boldsymbol{\tau}_j^*}{\lambda + e_{K_j^*}},$$

where $\bar{\boldsymbol{\beta}}_j^* = (e_{K_j^*})^{-1} \sum_{\{i: K_i = K_j^*\}} \boldsymbol{\beta}_i$; (b) conditional on $\mathbf{K}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n, \mathbf{Y}$, draw $\boldsymbol{\tau}_j^*$ from an inverse Wishart distribution with parameters

$$\nu_0 + e_{K_j^*} \quad \text{and} \quad \frac{\nu_0 \mathbf{S}_0 + e_{K_j^*} \mathbf{S}_j + R(\bar{\boldsymbol{\beta}}_j^*, \mathbf{m})}{\nu_0 + e_{K_j^*}},$$

where

$$\mathbf{S}_j = \frac{1}{e_{K_j^*}} \sum_{\{i: K_i = K_j^*\}} (\boldsymbol{\beta}_i - \bar{\boldsymbol{\beta}}_j^*)(\boldsymbol{\beta}_i - \bar{\boldsymbol{\beta}}_j^*)' \quad \text{and} \quad R(\bar{\boldsymbol{\beta}}_j^*, \mathbf{m}) = \frac{\lambda e_{K_j^*}}{\lambda + e_{K_j^*}} (\bar{\boldsymbol{\beta}}_j^* - \mathbf{m})(\bar{\boldsymbol{\beta}}_j^* - \mathbf{m})'.$$

4. *Conditional draw for $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_n$.* For each $j = 1, \dots, n_0$, independently draw $\boldsymbol{\beta}_i$, $i \in \{l: K_l = K_j^*\}$, from the density proportional to $L(\mathbf{Y}_i, \boldsymbol{\beta}_i) \phi(\boldsymbol{\beta}_i | \boldsymbol{\mu}_j^*, \boldsymbol{\tau}_j^*)$ by using, for example, a standard Metropolis–Hastings procedure.

When $n_0 = 1$, Steps 3 and 4 equate with a parametric MCMC procedure for panel data models similar to the algorithm described in Train (2003), Section 12.

5. Simulation study

In this section, we present some empirical evidence that shows how the MMNL procedures perform overall and relative to GML models and finite mixture (FM) of MNL models. We proceed to the estimation of the choice probabilities based on simulated data. Two different artificial data sets are generated for the simulation study: data set 1 is produced for studying non-panel data models, while data set 2 is designed to study models with panel data. In both cases, we consider a RUM model with three possible responses ($J = 3$) relative to the utilities U_1, U_2 and U_3 ,

$$\begin{cases} U_1 = x_{11}\beta_1 + x_{12}\beta_2 + \varepsilon_1, \\ U_2 = x_{21}\beta_1 + x_{22}\beta_2 + \varepsilon_2, \\ U_3 = x_{31}\beta_1 + x_{32}\beta_2 + \varepsilon_3. \end{cases}$$

As for data set 1, we choose $\varepsilon_1, \varepsilon_2, \varepsilon_3 \stackrel{iid}{\sim}$ standard Gumbel and $\beta = (\beta_1, \beta_2)' \stackrel{iid}{\sim} 0.5 \times \delta_{(-5,5)} + 0.5 \times \delta_{(5,-5)}$. For individual i , we randomly generate (componentwise) the covariate vector $\mathbf{x}_i = (x_{11}, x_{12}, x_{21}, x_{22}, x_{31}, x_{32})$, independently from a Uniform(-2, 2) distribution. Set $Y_i = j$ if $U_{ij} > U_{il}, l \neq j$, for $j = 1, 2, 3$. Repeat this procedure n times independently to obtain a data set with (Y_i, \mathbf{x}_i) for $i = 1, \dots, n$. As for data set 2, we assume that there are n individuals, each making $T_i = 10$ choices for $i = 1, \dots, n$. We then simulate data using the same model used to generate data set 1. The only change is that β is drawn from the two-component mixture of bivariate normal distributions, $\beta \stackrel{iid}{\sim} 0.5 \times N((-5, 5)', 2\mathbf{I}) + 0.5 \times N((5, -5)', 2\mathbf{I})$, where \mathbf{I} is the identity matrix.

We start by applying our procedures to the estimation of choice probabilities $P(\{j\}|\mathbf{x})$, for $j = 1, 2, 3$, based on the set of covariates $\mathbf{x} = (1.0, -0.9, 1.0, 0.2, 1.0, 0.9)$. The prior parameters for the specifications of the Bayesian MMNL models for panel and non-panel data (pertaining to the explicit models in Sections 3 and 4) are set to be $a = 1, \nu_0 = 2, \mathbf{m} = (0, 0)', \mathbf{S}_0 = \mathbf{I}$ and $\lambda = 1$. Additionally, we use $N = 100$ for the truncation level in the blocked Gibbs Algorithms 1 and 2 given in Sections 3 and 4, respectively. A Bayesian GML model is also estimated for comparison with the same specifications for $\nu_0, \mathbf{m}, \mathbf{S}_0$ and λ . In all cases, we use the estimator (17) based on an initial burn-in of 10,000 cycles and an additional 10,000 Gibbs cycles ($M = 10,000$) for the estimation. In addition, to measure how good of our estimates are, we define a measure, root mean square (RMS) value, as

$$\text{RMS} = \sqrt{\frac{1}{J} \sum_{j \in \mathbf{C}} \frac{1}{M} \sum_{m=1}^M (P(\{j\}|G^{(m)}, \mathbf{x}) - P_0(\{j\}|\mathbf{x}))^2},$$

where $P_0(\{j\}|\mathbf{x})$ is the choice probability resulting from the data generating process.

Simulation results using data set 1 ($n = 500$) and data set 2 ($n = 100, T_i = 10$) are summarized in Table 1, together with RMS values, for both the GML and the MMNL models. They show that the performance of the nonparametric MMNL model is better than that of the parametric GML model in the non-panel case, as indicated by a smaller RMS value and more accurate estimates of choice probabilities, while the GML and MMNL models display similar performances in the panel case. As expected, the GML model suffers from misspecification in the non-panel case,

Table 1. Simulation results for data set 1 (columns 3–4) and for data set 2 (columns 5–6) with $\mathbf{x} = (1.0, -0.9, 1.0, 0.2, 1.0, 0.9)$ – the estimates (Est.), the credible intervals (C.I.) and the root mean square (RMS) values are presented; GML = Gaussian mixed logit, MMNL = mixed multinomial logit

		Data set 1 (non-panel case)				Data set 2 (panel case)			
		$n = 500$				$n = 100, T_i = 10$			
		True	Est.	(95% C.I.)	RMS	True	Est.	(95% C.I.)	RMS
GML	$P(\{1\} \mathbf{x})$	0.4980	0.3203	(0.2907, 0.3501)	0.2258	0.4939	0.4585	(0.4476, 0.4685)	0.0266
	$P(\{2\} \mathbf{x})$	0.0167	0.3348	(0.3308, 0.3377)		0.0279	0.0521	(0.0378, 0.0675)	
	$P(\{3\} \mathbf{x})$	0.4853	0.3449	(0.3191, 0.3715)		0.4782	0.4894	(0.4717, 0.5061)	
MMNL	$P(\{1\} \mathbf{x})$	0.4980	0.4856	(0.4748, 0.4945)	0.0137	0.4939	0.4586	(0.4495, 0.4670)	0.0265
	$P(\{2\} \mathbf{x})$	0.0167	0.0257	(0.0069, 0.0551)		0.0279	0.0494	(0.0329, 0.0679)	
	$P(\{3\} \mathbf{x})$	0.4853	0.4886	(0.4615, 0.5057)		0.4782	0.4920	(0.4705, 0.5107)	

while the two-component mixture of bivariate normals used for generating data set 2 is correctly accounted for by the GML because of the hyperprior on the parameter $(\boldsymbol{\mu}, \boldsymbol{\tau})$ we are using. We then get confirmation that the fit of the MMNL model is as good as that of the GML model. We also performed estimation of the MMNL model for different choices of the scale parameter λ (not reported here) which show two different behaviors for the non-panel and the panel case. As for the non-panel case, RMS values and the estimates remain stable, whereas, in the panel case, the estimates are more accurate when we decrease λ with slightly smaller RMS values. An interpretation of an increase of accuracy is as follows: a smaller λ corresponds to a more diffuse H , the prior predictive distribution of \tilde{G} . Since H is different from the distribution used to simulate the $\boldsymbol{\beta}$'s in the data generating process, we obtain evidence that a diffuse H helps in capturing the true form of the mixing distribution G . Also, note that a smaller λ yields a smaller RMS, the latter being a measure of the combination of the accuracy and the variability of the posterior variates of $P(\{j\}|\mathbf{x})$. An examination of their autocorrelation functions along the chain shows that a smaller λ causes a slower mixing of the blocked Gibbs sampler, which increases the component of variability in the RMS; see Figure 1. The decrease in RMS then shows that such precision loss is more than balanced by a higher accuracy of the estimate, although one should also control the convergence properties of the sampler by avoiding taking λ too small.

We investigated the sensitivity of the results to the prior parameter ν_0 , where a larger ν_0 corresponds to a more concentrated inverse Wishart distribution on \mathbf{S}_0 . However, we did not observe substantial differences in the estimation by varying ν_0 and we decided to set $\nu_0 = 2$ and $\mathbf{S}_0 = \mathbf{I}$ as a default non-informative choice for these parameters; see Train (2003), Section 12. The non-parametric prior on \tilde{G} is also dependent on the total mass a , which is positively related to the number of components in the mixture distribution of the $\boldsymbol{\beta}$'s. Generally, $a = 1$ is considered a default choice for a finite mixture model with a fixed, but uncertain, number of components. We performed estimation for larger a , obtaining almost identical results: $a = 1$ was, in fact, sufficient for detecting the two-component mixture we used in generating the data. Although we have not

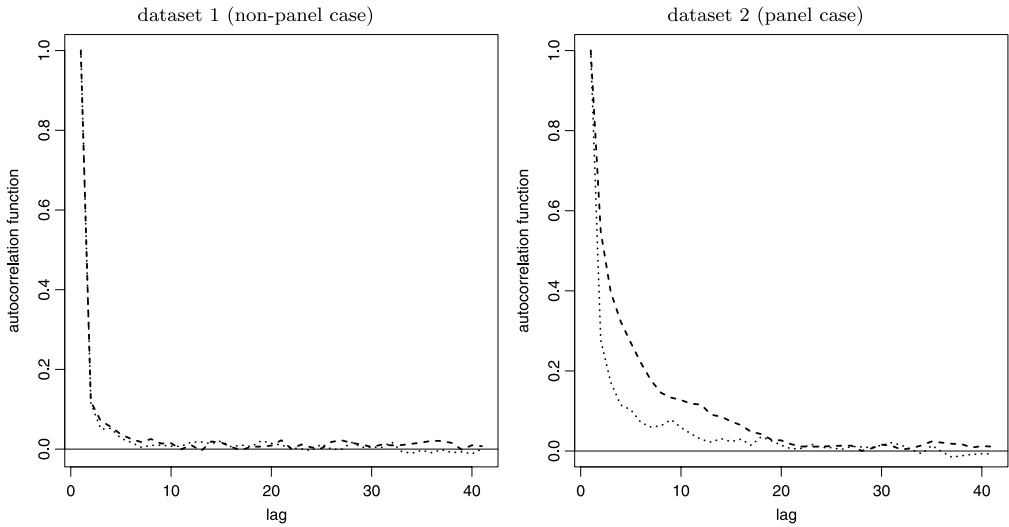


Figure 1. MMNL model: Autocorrelation functions for the choice probability $P(\{1\}|\mathbf{x})$ for data set 1 (left) and data set 2 (right), obtained from the posterior sample of the β 's for the MMNL model with prior hyperparameter $\lambda = 0.01$ (dashed) and $\lambda = 1$ (dotted).

done so, the blocked Gibbs procedures described in Sections 3 and 4 can be easily extended to place an additional prior on a . Furthermore, the truncation level of $N = 100$ in (15) is sufficiently large as we observed almost identical estimation results from runs of the blocked Gibbs sampler with larger values of N .

The second simulation study aims at the verification of the consistency result of Section 2 by estimating the MMNL model for increasing sample sizes for both data set 1 and data set 2. We also sample β variates from their posterior distribution, thus obtaining approximated evaluation of the mixing distribution G . The prior parameters are set as $a = 1$, $v_0 = 2$, $\mathbf{m} = (0, 0)'$, $\mathbf{S}_0 = \mathbf{I}$, $N = 100$ and $\lambda = 1$. Table 2 reports the results by showing, as expected, a noticeable decrease

Table 2. MMNL model: estimates and the root mean square (RMS) for data set 1 and for data set 2 with $\mathbf{x} = (1.0, -0.9, 1.0, 0.2, 1.0, 0.9)$ and different sample sizes

	Data set 1 (non-panel case)				Data set 2 (panel case)			
	True	$n = 50$	$n = 100$	$n = 500$	True	$n = 10$ $T_i = 10$	$n = 50$ $T_i = 10$	$n = 100$ $T_i = 10$
$P(\{1\} \mathbf{x})$	0.4980	0.4927	0.5145	0.4856	0.4939	0.5956	0.4176	0.4586
$P(\{2\} \mathbf{x})$	0.0167	0.1046	0.0489	0.0257	0.0279	0.0527	0.0562	0.0494
$P(\{3\} \mathbf{x})$	0.4853	0.4027	0.4366	0.4886	0.4782	0.3517	0.5261	0.4920
RMS		0.0867	0.0440	0.0137		0.0977	0.0556	0.0265

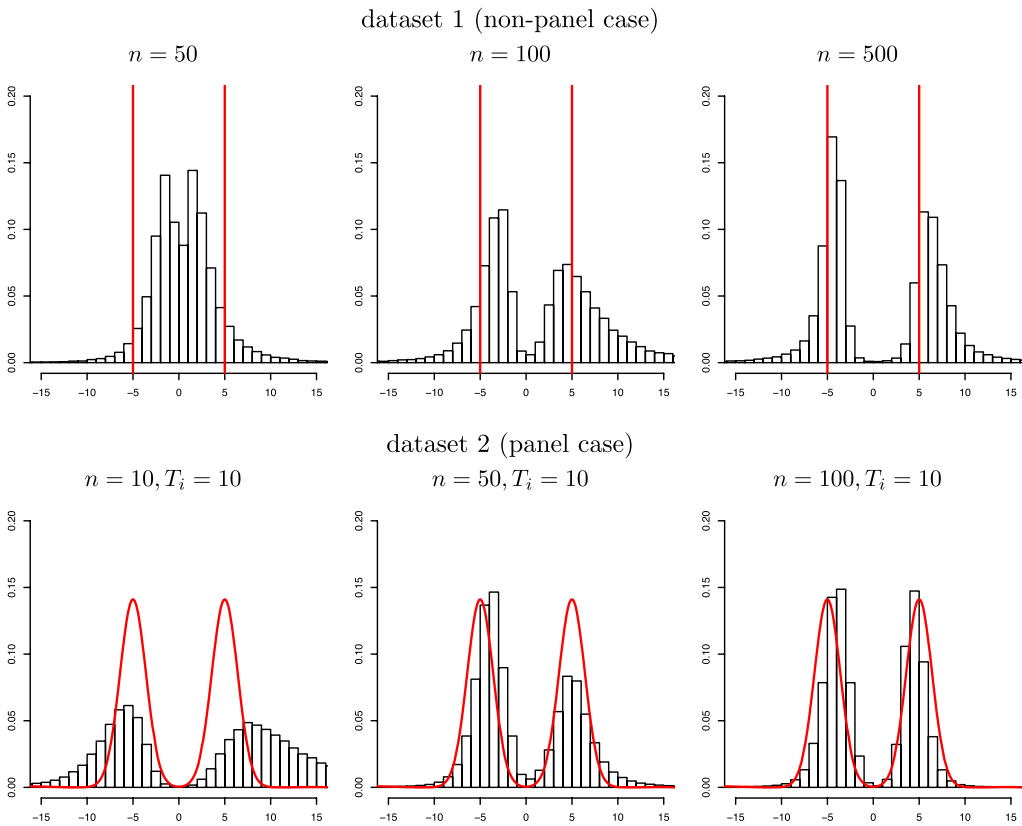


Figure 2. MMNL model: histogram estimate of the posterior marginal density of β_1 's for data set 1 (top) and for data set 2 (bottom) and different sample sizes. The solid lines represent the true mixing distribution.

of RMS for both non-panel and panel data as the number of observations increases. In addition, Figure 2 reports the histograms of samples for β_1 from its marginal posterior distribution against the mixing distribution used in the data generating process: it shows how the approximation of the true mixing distribution G improves as more and more data become available.

Finally, we evaluate the performance of the Bayesian MMNL model via a comparison with the finite mixture (FM) MNL model estimated via the EM algorithm described in Train (2008), Section 4. The FM MNL model can be considered nonparametric in the sense that the locations and weights of the mixing distribution G are both assumed to be parameters. The selection of the number of points in the mixing is based on the BIC information criterion. We consider 500 Monte Carlo replicates of each of the following 6 situations: data set 1 with sample sizes $n = 50, 100$ and 500 ; data set 2 with $(n = 10, T_i = 10)$, $(n = 50, T_i = 10)$ and $(n = 100, T_i = 10)$. For a given sample, the posterior estimate of $P(\{j\}|\mathbf{x})$ in equation (17) is computed, based on 6000 Gibbs cycles after a burn-in period of 4000 for $j = 1, 2, 3$ and for \mathbf{x} in a 6-dimensional grid of the hypercube $(-2, 2)^6$ of 5^6 equally-spaced points. At the same time, we compute the FM MNL

Table 3. Average L_1 -error from 500 Monte Carlo replicates – FM MNL = finite mixture of multinomial logit; MMNL = mixed multinomial logit

	Data set 1 (non-panel case)			Data set 2 (panel case)		
	$n = 50$	$n = 100$	$n = 500$	$n = 10$ $T_i = 10$	$n = 50$ $T_i = 10$	$n = 100$ $T_i = 10$
FM MNL	0.0521	0.0295	0.0107	0.0891	0.0505	0.0297
MMNL	0.0577	0.0316	0.011	0.0827	0.0467	0.0268

estimate of $P(\{j\}|\mathbf{x})$ for $j = 1, 2, 3$, evaluated on the same grid of \mathbf{x} -points. We call $\hat{\mathbf{q}}(\mathbf{x})$ and $\mathbf{q}_0(\mathbf{x})$ the estimated vector and the true vector of choice probabilities evaluated at \mathbf{x} , respectively. We measure the overall error of estimation with the L_1 -distance $\int_{\mathcal{X}} |\hat{\mathbf{q}}(\mathbf{x}) - \mathbf{q}_0(\mathbf{x})| d\mathbf{x}$, which corresponds to the (rescaled) distance $d(\hat{\mathbf{q}}, \mathbf{q}_0)$ in equation (7), with $M(d\mathbf{x})$ being the uniform distribution on the hypercube $(-2, 2)^6$. We compute the L_1 -error for the Bayesian MMNL estimator and the FM MNL estimator, then average over the 500 Monte Carlo replicates. The results are reported in Table 3 and show that the MMNL estimators outperform the FM MNL estimators in the panel case for all sample sizes, while in the non-panel case, the situation is reversed, with a similar performance for $n = 500$. Note, however, that data set 1 is generated exactly from a finite mixture model so that the FM MNL model is expected to perform well. Overall, the decrease in the average error for larger sample sizes is a further confirmation of the consistency result of Section 2.

6. Proof of Theorem 1

Throughout this section, we work with the family of multinomial logistic kernels

$$k_j(\mathbf{x}, \boldsymbol{\beta}) = \frac{\exp(x'_j \boldsymbol{\beta})}{\sum_{l \in \mathbf{C}} \exp(x'_l \boldsymbol{\beta})}, \quad j = 1, \dots, J.$$

With $q_j(\mathbf{x}; G)$ denoting the j th element of the vector $\mathbf{q}(\mathbf{x}; G)$, we have that $q_j(\mathbf{x}; G) = \int_{\mathbb{R}^d} k_j(\mathbf{x}, \boldsymbol{\beta}) G(d\boldsymbol{\beta})$. Note that $q_Y(\mathbf{x}; G_0)$ is the joint density of (Y, \mathbf{X}) with respect to the counting measure on the integer set \mathbf{C} and the measure $M(d\mathbf{x})$ on \mathcal{X} .

For the proof of Theorem 1, the following lemma is essential, stating that on the space \mathbb{P} , the weak topology and the topology induced by the L_1 -distance d defined in (7) are equivalent.

Lemma 1. *Let d_w be any distance that metrizes the weak topology on \mathbb{P} and $(G_n)_{n \geq 1}$ be a sequence in \mathbb{P} . Then $d_w(G_n, G_0) \rightarrow 0$ if and only if $d(\mathbf{q}(\cdot; G_n), \mathbf{q}(\cdot; G_0)) \rightarrow 0$.*

Proof. For the “if” part, it is sufficient that $d_w(G_n, G_0) \rightarrow 0$ implies that $\int_{\mathcal{X}} |q_j(\mathbf{x}; G_n) - q_j(\mathbf{x}; G_0)| M(d\mathbf{x}) \rightarrow 0$ for an arbitrary $j \in \mathbf{C}$. The latter is a consequence of the definition of weak convergence and an application of Scheffé’s theorem since $k_j(\mathbf{x}, \boldsymbol{\beta})$ is bounded and continuous in $\boldsymbol{\beta}$ for each $\mathbf{x} \in \mathcal{X}$. To show the converse, we prove that G being distant from G_0 in the

weak topology implies that $\mathbf{q}(\cdot; G)$ is distant from $\mathbf{q}(\cdot; G_0)$ in the L_1 -distance d . Define a weak neighborhood of G_0 as

$$V = \left\{ G: \left| \int_{\mathbb{R}^d} \int_{\mathcal{X}} k_j(\mathbf{x}, \boldsymbol{\beta}) M(d\mathbf{x}) G(d\boldsymbol{\beta}) - \int_{\mathbb{R}^d} \int_{\mathcal{X}} k_j(\mathbf{x}, \boldsymbol{\beta}) M(d\mathbf{x}) G_0(d\boldsymbol{\beta}) \right| < \delta, j \in \mathbf{C} \right\}.$$

Since $\int_{\mathcal{X}} k_j(\mathbf{x}, \boldsymbol{\beta}) M(d\mathbf{x})$ is a bounded continuous function on \mathbb{R}^d for each j , $G \in V^c$ implies that $d_w(G, G_0) > \delta$. Based on the inequalities

$$\begin{aligned} d(\mathbf{q}(\cdot; G), \mathbf{q}(\cdot; G_0)) &\geq \max_{j \in \mathbf{C}} \int_{\mathcal{X}} |q_j(\mathbf{x}; G_n) - q_j(\mathbf{x}; G_0)| M(d\mathbf{x}) \\ &\geq \max_{j \in \mathbf{C}} \left| \int_{\mathcal{X}} \int_{\mathbb{R}^d} k_j(\mathbf{x}, \boldsymbol{\beta}) G(d\boldsymbol{\beta}) M(d\mathbf{x}) - \int_{\mathcal{X}} \int_{\mathbb{R}^d} k_j(\mathbf{x}, \boldsymbol{\beta}) G_0(d\boldsymbol{\beta}) M(d\mathbf{x}) \right| \end{aligned}$$

and an application of Fubini's theorem, it follows that, for any $\epsilon < \delta$ and any $G \in V^c$, $d(\mathbf{q}(\cdot; G), \mathbf{q}(\cdot; G_0)) > \epsilon$. The proof is then complete. \square

Remark 1. Lemma 1 has two important consequences: (a) both \mathcal{Q} and \mathbb{P} are separable spaces under the metric d ; (b) the statement of Theorem 1 is equivalent to saying that \mathcal{P}_n accumulates all probability mass in a weak neighborhood of G_0 .

Define $\Lambda_n(G) = \prod_{i=1}^n q_{Y_i}(\mathbf{X}_i; G) / q_{Y_i}(\mathbf{X}_i; G_0)$ so that the posterior distribution of \tilde{G} can be written as

$$\mathcal{P}_n(A) = \frac{\int_A \Lambda_n(G) \mathcal{P}(dG)}{\int_{\mathbb{P}} \Lambda_n(G) \mathcal{P}(dG)}. \tag{19}$$

We now take $A = \{G: d(\mathbf{q}(\cdot; G), \mathbf{q}(\cdot; G_0)) > \epsilon\}$ and will, as is usual in the Bayesian consistency literature, separately consider the numerator and the denominator of (19). To this end, define $I_n = \int_{\mathbb{P}} \Lambda_n(G) \mathcal{P}(dG)$. Relying on the separability of \mathbb{P} under the topology induced by d (see Remark 1), for any $\eta > 0$, we can cover A with a countable union of disjoint sets A_j such that

$$A_j \subseteq A_j^* = \{G: d(\mathbf{q}(\cdot; G), \mathbf{q}(\cdot; G_j)) < \eta\} \tag{20}$$

and $\{G_j\}_{j \geq 1}$ is a countable set in \mathbb{P} such that $d(\mathbf{q}(\cdot; G_j), \mathbf{q}(\cdot; G_0)) > \epsilon$ for any j . Consider the fact that

$$\mathcal{P}_n(A) = \sum_{j \geq 1} \mathcal{P}_n(A_j) \leq \sum_{j \geq 1} \sqrt{\mathcal{P}_n(A_j)} = \sum_{j \geq 1} \sqrt{I_n^{-1} \int_{A_j} \Lambda_n(G) \mathcal{P}(dG)}.$$

Hence, Theorem 1 holds if we prove that, for all large n ,

$$\forall c > 0, \quad I_n > \exp(-nc) \quad \text{a.s.} \tag{21}$$

$$\exists b > 0: \quad \sum_{j \geq 1} \sqrt{\int_{A_j} \Lambda_n(G) \mathcal{P}(dG)} < \exp(-nb) \quad \text{a.s.} \tag{22}$$

As for (21), consider the Kullback–Leibler (KL) support condition of \mathcal{P} defined by

$$\mathcal{P} \left\{ G: \int_{\mathcal{X}} K(G_0, G|\mathbf{x})M(d\mathbf{x}) < \epsilon \right\} > 0 \quad \forall \epsilon > 0, \tag{23}$$

where $K(G_0, G|\mathbf{x}) = \sum_{j \in \mathbf{C}} q_j(\mathbf{x}; G_0) \log[q_j(\mathbf{x}; G_0)/q_j(\mathbf{x}; G)]$. If \mathcal{P} satisfies condition (23), then (21) holds. To see this, it is sufficient to note that the KL divergence of $q_Y(\mathbf{X}; G)$ from $q_Y(\mathbf{X}; G_0)$ with respect to the measure $M(d\mathbf{x})$ on \mathcal{X} and the counting measure on \mathbf{C} is given by $\int K(G, G_0|\mathbf{x})M(d\mathbf{x})$. By the compactness of \mathcal{X} , the law of large numbers then leads to

$$\frac{1}{n} \sum_{i=1}^n \log \frac{q_{Y_i}(\mathbf{X}_i; G_0)}{q_{Y_i}(\mathbf{X}_i; G)} \rightarrow \int_{\mathcal{X}} K(G_0, G|\mathbf{x})M(d\mathbf{x}) \quad \text{a.s.}$$

The result in (21) then follows from standard arguments, see, for example, Wasserman (1998). Lemma 2 below states that (23) is satisfied under the hypotheses of Theorem 1.

Lemma 2. *If G_0 lies in the weak support of \mathcal{P} and condition (i) of Theorem 1 holds, then G_0 is in the KL support of \mathcal{P} , according to (23).*

Proof. It is sufficient to show that for any $j \in \mathbf{C}$ and any $\eta < 1$, there exists a δ such that $|q_j(\mathbf{x}; G)/q_j(\mathbf{x}; G_0) - 1| \leq \eta$ whenever G is in W_δ , a δ -weak neighborhood of G_0 . In fact, this implies that

$$\begin{aligned} \int_{\mathcal{X}} q_j(\mathbf{x}; G_0) \log \left[\frac{q_j(\mathbf{x}; G_0)}{q_j(\mathbf{x}; G)} \right] M(d\mathbf{x}) &\leq \int_{\mathcal{X}} q_j(\mathbf{x}; G_0) \left| \frac{q_j(\mathbf{x}; G_0)}{q_j(\mathbf{x}; G)} - 1 \right| M(d\mathbf{x}) \\ &\leq \int_{\mathcal{X}} q_j(\mathbf{x}; G_0) \left(\frac{\eta}{1 - \eta} \right) M(d\mathbf{x}) \\ &\leq \frac{\eta}{1 - \eta}, \end{aligned}$$

which, in turn, leads to the thesis, by the arbitrary nature of j .

Let $c = \inf_{\mathbf{x} \in \mathcal{X}} q_j(\mathbf{x}; G_0)$, which is positive by condition (i) of Theorem 1, and assume that $G \in W_\delta$ for a δ that will be determined later. Note that, for any $\rho > 0$, one can set $M_\rho > 0$ such that $G_0\{\boldsymbol{\beta}: |\boldsymbol{\beta}| > M_\rho - \delta\} < \rho$. Then, using the Prokhorov metric, $G \in W_\delta$ implies that $G\{\boldsymbol{\beta}: |\boldsymbol{\beta}| > M_\rho\} < \rho + \delta$. Also, note that the family of functions $\{k_j(\mathbf{x}, \boldsymbol{\beta}), \mathbf{x} \in \mathcal{X}\}$, as $\boldsymbol{\beta}$ varies in the compact set $\{|\boldsymbol{\beta}| \leq M_\rho\}$, is uniformly equicontinuous. By an application of the Arzelà–Ascoli theorem, we know that, given a $\gamma > 0$, there exist finitely many points $\mathbf{x}_1, \dots, \mathbf{x}_m$ such that, for any $\mathbf{x} \in \mathcal{X}$, there is an index i such that

$$\sup_{|\boldsymbol{\beta}| \leq M_\rho} |k_j(\mathbf{x}, \boldsymbol{\beta}) - k_j(\mathbf{x}_i, \boldsymbol{\beta})| < \gamma. \tag{24}$$

For an arbitrary $\mathbf{x} \in \mathcal{X}$, choose the appropriate \mathbf{x}_i such that (24) holds, so that

$$\begin{aligned} \left| \frac{q_j(\mathbf{x}; G)}{q_j(\mathbf{x}; G_0)} - 1 \right| &\leq \frac{1}{c} \left(\left| \int k_j(\mathbf{x}_i, \boldsymbol{\beta}) G(d\boldsymbol{\beta}) - \int k_j(\mathbf{x}_i, \boldsymbol{\beta}) G_0(d\boldsymbol{\beta}) \right| \right. \\ &\quad \left. + \int |k_j(\mathbf{x}, \boldsymbol{\beta}) - k_j(\mathbf{x}_i, \boldsymbol{\beta})| G(d\boldsymbol{\beta}) + \int |k_j(\mathbf{x}, \boldsymbol{\beta}) - k_j(\mathbf{x}_i, \boldsymbol{\beta})| G_0(d\boldsymbol{\beta}) \right) \\ &:= \frac{I_1 + I_2 + I_3}{c}. \end{aligned}$$

We have that $G \in W_\delta$ implies $I_1 \leq \delta$. As for I_2 , we have

$$\begin{aligned} I_2 &= \int_{|\boldsymbol{\beta}| \leq M_\rho} |k_j(\mathbf{x}, \boldsymbol{\beta}) - k_j(\mathbf{x}_i, \boldsymbol{\beta})| G(d\boldsymbol{\beta}) + \int_{|\boldsymbol{\beta}| > M_\rho} |k_j(\mathbf{x}, \boldsymbol{\beta}) - k_j(\mathbf{x}_i, \boldsymbol{\beta})| G(d\boldsymbol{\beta}) \\ &\leq \gamma + 2G\{\boldsymbol{\beta}: |\boldsymbol{\beta}| > M_\rho\} \leq \gamma + 2(\rho + \delta). \end{aligned}$$

Similar arguments lead to $I_3 \leq \gamma + 2\rho$. Finally, we get

$$\left| \frac{q_j(\mathbf{x}; G)}{q_j(\mathbf{x}; G_0)} - 1 \right| \leq \frac{3\delta + 2\gamma + 4\rho}{c},$$

so that, for given $\eta < 1$, it is always possible to choose δ, ρ (by tightness of G_0) and γ (by the Arzelà–Ascoli theorem) small enough such that the right-hand side in the last inequality is smaller than η . The proof is then complete. \square

We now aim to show that (22) holds under the hypotheses of Theorem 1, by extending the method set forth by Walker (2004) for strong consistency. In order to simplify the notation, let $\Lambda_{nj} = \int_{A_j} \Lambda_n(G) \mathcal{P}(dG)$, where $(A_j)_{j \geq 1}$ is the covering of A in (20). The following identity is the key:

$$\Lambda_{n+1j} / \Lambda_{nj} = q_{Y_{n+1}}^{nA_j}(\mathbf{X}_{n+1}) / q_{Y_{n+1}}(\mathbf{X}_{n+1}; G_0), \tag{25}$$

where $q_l^{nA_j}(\mathbf{X}_{n+1}) = \int_{\mathbb{P}} q_l(\mathbf{X}_{n+1}; G) \mathcal{P}_{nA_j}(dG)$, $l \in \mathbf{C}$ and \mathcal{P}_{nA_j} is the posterior distribution restricted, and normalized, to the set A_j . Note that (25) includes the case of $n = 0$ and $\Lambda_{0j} = \mathcal{P}(A_j)$. By using conditional expectation, we have that

$$\begin{aligned} E[\Lambda_{n+1j}^{1/2} | (Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n), \mathbf{X}_{n+1}] &= \Lambda_{nj}^{1/2} \sum_{l \in \mathbf{C}} \sqrt{q_l^{nA_j}(\mathbf{X}_{n+1}) q_l(\mathbf{X}_{n+1}; G_0)} \\ &= \Lambda_{nj}^{1/2} (1 - h[\mathbf{q}^{nA_j}(\mathbf{X}_{n+1}), \mathbf{q}(\mathbf{X}_{n+1}; G_0)]), \end{aligned}$$

where $\mathbf{q}^{nA_j}(\mathbf{X}_{n+1}) = [q_1^{nA_j}(\mathbf{X}_{n+1}), \dots, q_J^{nA_j}(\mathbf{X}_{n+1})]$ and, for $\mathbf{q}_1, \mathbf{q}_2 \in \Delta$,

$$h(\mathbf{q}_1, \mathbf{q}_2) = 1 - \sum_{j \in \mathbf{C}} \sqrt{q_{1j} q_{2j}}.$$

Note that $h(\mathbf{q}_1, \mathbf{q}_2)$ is a variation of the Hellinger distance $\sqrt{\sum_{j \in \mathbf{C}} (q_{1j}^{1/2} - q_{2j}^{1/2})^2}$ on Δ and that $h(\mathbf{q}_1, \mathbf{q}_2) \leq 1$. By taking the conditional expectation with respect to $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ only, we get the following identity:

$$E\{\Lambda_{n+1j}^{1/2} | (Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)\} = \Lambda_{nj}^{1/2} \left(1 - \int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x}) \right). \tag{26}$$

Since the Hellinger distance and the Euclidean distance are equivalent metrics in Δ , it can be proven that, for $(\mathbf{q}_n)_{n \geq 1} \in \mathcal{Q}$ and $\mathbf{q}_0 \in \mathcal{Q}$,

$$\int_{\mathcal{X}} h[\mathbf{q}_n(\mathbf{x}), \mathbf{q}_0(\mathbf{x})] M(d\mathbf{x}) \rightarrow 0 \quad \text{if and only if} \quad d(\mathbf{q}_n, \mathbf{q}_0) \rightarrow 0. \tag{27}$$

The equivalence in (27) can be used to show that $\int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x})$ is bounded away from zero. In fact, take G_j defined in (20) and note that, by the triangle inequality,

$$\begin{aligned} \int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x}) &\geq \int_{\mathcal{X}} h[\mathbf{q}(\mathbf{x}; G_j), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x}) \\ &\quad - \int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_j)] M(d\mathbf{x}). \end{aligned}$$

Since $d(\mathbf{q}(\cdot; G_j), \mathbf{q}(\cdot; G_0)) > \epsilon$, (27) ensures the existence of a positive constant, say ϵ_2 , such that $\int_{\mathcal{X}} h[\mathbf{q}(\mathbf{x}; G_j), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x}) > \epsilon_2$. Now, choose η in (20) such that, for each $G \in A_j$, $\int_{\mathcal{X}} h[\mathbf{q}(\mathbf{x}; G), \mathbf{q}(\mathbf{x}; G_j)] M(d\mathbf{x}) < \epsilon_2$, where we have again used (27). Since $\mathbf{q}^{nA_j}(\mathbf{x})$ does not correspond exactly to a particular $G \in A_j$, we use the convexity of the distance $h[\mathbf{q}(\mathbf{x}; G), \mathbf{q}(\mathbf{x}; G_j)]$ in its first argument to show that $\int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_j)] M(d\mathbf{x}) < \epsilon_2$. Note that, in fact, by Jensen’s inequality,

$$\begin{aligned} \int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_j)] M(d\mathbf{x}) &= \int_{\mathcal{X}} \left(1 - \sum_{l \in \mathbf{C}} \sqrt{\int_{\mathbb{P}} q_l(\mathbf{X}_{n+1}; G) \mathcal{P}_{nA_j}(dG) q_l(\mathbf{x}; G_j)} \right) M(d\mathbf{x}) \\ &\leq \int_{\mathbb{P}} \int_{\mathcal{X}} h[\mathbf{q}(\mathbf{x}; G), \mathbf{q}(\mathbf{x}; G_j)] M(d\mathbf{x}) \mathcal{P}_{nA_j}(dG) < \epsilon_2. \end{aligned}$$

Hence, there exists a $\epsilon_3 > 0$ such that $\int_{\mathcal{X}} h[\mathbf{q}^{nA_j}(\mathbf{x}), \mathbf{q}(\mathbf{x}; G_0)] M(d\mathbf{x}) > \epsilon_3$.

From (26), it now follows that

$$E(\Lambda_{n+1j}^{1/2}) < (1 - \epsilon_3)^n \sqrt{\mathcal{P}(A_j)}$$

and an application of Markov’s inequality leads to

$$P\left\{ \sum_{j \geq 1} \Lambda_{nj}^{1/2} > \exp(-nb) \right\} < \exp(nb) (1 - \epsilon_3)^n \sum_{j \geq 1} \sqrt{\mathcal{P}(A_j)}.$$

Therefore, (21) holds for any $b < -\log(1 - \epsilon_3)$ from an application of the Borel–Cantelli lemma, provided that the following summability condition is satisfied:

$$\sum_{j \geq 1} \sqrt{\mathcal{P}(A_j)} < +\infty. \tag{28}$$

Lemma 3 below shows that \mathcal{P} satisfies condition (28) under the stated hypotheses and, in turn, completes the proof of Theorem 1.

Lemma 3. *Let $H \in \mathbb{P}$ be the prior predictive distribution of \mathcal{P} and assume that condition (ii) of Theorem 1 holds. Then (28) is verified.*

Proof. The proof follows along the lines of arguments used by Lijoi, Prünster and Walker (2005). Take δ to be any positive number in $(0, 1)$ and $(a_n)_{n \geq 1}$ any increasing sequence of positive numbers such that $a_n \rightarrow +\infty$. Also, let $a_0 = 0$. Define $C_n = \{\boldsymbol{\beta}: |\boldsymbol{\beta}| \leq a_n\}$ and consider the family of subsets of \mathbb{P} defined by

$$\mathbb{B}_{a_n, \delta} = \{G: G(C_n) \geq 1 - \delta, G(C_{n-1}) < 1 - \delta\} \tag{29}$$

for each $n \geq 1$. These sets are pairwise disjoint and $\bigcup_n \mathbb{B}_{a_n, \delta} = \mathbb{P}$. For the moment, let us assume that the metric entropy of $\mathbb{B}_{a_n, \delta}$ with respect to the distance d is uniformly bounded in n , that is, the number of η -balls in the distance d that covers $\mathbb{B}_{a_n, \delta}$ is finite for any n . Summability in (28) is then implied by

$$\sum_{n \geq 1} \sqrt{\mathcal{P}(\mathbb{B}_{a_n, \delta})} < +\infty. \tag{30}$$

In order to prove (30), note that $\mathbb{B}_{a_n, \delta} \subset \{G: G(C_{n-1}^c) > \delta'\}$ for some $\delta' > \delta$. An application of Markov’s inequality leads to $\mathcal{P}(\mathbb{B}_{a_n, \delta}) \leq (1/\delta')H(C_{n-1}^c)$, hence (30) is implied by $\sum_{n \geq 1} \sqrt{H(C_{n-1}^c)} < +\infty$. Next, we have that

$$\int_{\mathbb{R}^d} |\boldsymbol{\beta}| H(d\boldsymbol{\beta}) = \sum_{n \geq 1} \int_{C_{n-1}^c / C_n^c} |\boldsymbol{\beta}| H(d\boldsymbol{\beta}) \geq \sum_{n \geq 1} a_{n-1} [H(C_{n-1}^c) - H(C_n^c)],$$

by a second application of Markov’s inequality, so that condition (ii) of Theorem 1 ensures that $\sum_{n \geq 1} a_{n-1} [H(C_{n-1}^c) - H(C_n^c)] < +\infty$. If we now take $a_n \sim n^2$, it is easy to see that $H(C_n^c) = o(n^{-(2+r)})$ for some $r > 0$. For example,

$$\sum_{n \geq 1} (n - 1)^2 [H(C_{n-1}^c) - H(C_n^c)] = \sum_{n \geq 1} (2n - 1) H(C_n^c).$$

This, in turn, ensures the convergence of $\sum_{n \geq 1} H(C_{n-1}^c)^\alpha$ for any α such that $(2+r)^{-1} < \alpha < 1$, which includes the case $\alpha = 1/2$. Condition (30) is then verified.

In order to complete the proof, it remains to show that the metric entropy of $\mathbb{B}_{a_n, \delta}$ with respect to the distance d is uniformly bounded in n . It is actually sufficient to reason in terms of the distance over \mathbb{P} induced by

$$d_j(\mathbf{q}_1, \mathbf{q}_2) = \int_{\mathcal{X}} |q_{1j}(\mathbf{x}) - q_{2j}(\mathbf{x})| M(d\mathbf{x})$$

for an arbitrary $j \in \mathbf{C}$ since $\max_j d_j(\mathbf{q}_1, \mathbf{q}_2) \leq d(\mathbf{q}_1, \mathbf{q}_2) \leq J \max_j d_j(\mathbf{q}_1, \mathbf{q}_2)$. Let \mathcal{G} be a set in \mathcal{Q} and, for $\delta > 0$, denote by $J(\delta, \mathcal{G})$ the metric entropy of \mathcal{G} with respect to d_j , that is, the logarithm of the minimum of all k such that there exists $\mathbf{q}_1, \dots, \mathbf{q}_k \in \mathcal{Q}$ with the property that $\forall \mathbf{q} \in \mathcal{Q}$, there exists an i such that $d_j(\mathbf{q}, \mathbf{q}_i) < \delta$. The result is then stated as follows: for $\mathcal{G}_{a_n, \delta} = \{\mathbf{q}(\mathbf{x}; G) : G \in \mathbb{B}_{a_n, \delta}\}$, there exists an $M_\delta < +\infty$ depending only on δ such that, for any n ,

$$J(\delta, \mathcal{G}_{a_n, \delta}) < M_\delta. \tag{31}$$

The proof of (31) consists of a sequence of three steps.

Step (1). Define $C_a = \{\boldsymbol{\beta} : |\boldsymbol{\beta}| \leq a\}$ and $\mathcal{F}_a = \{\mathbf{q}(\mathbf{x}; G) : G(C_a) = 1\}$. Then

$$J(2\delta, \mathcal{F}_a) \leq \left(\frac{2aK}{\delta} + 1\right)^d \left(1 + \log \frac{1 + \delta}{\delta}\right), \tag{32}$$

where K is a constant that depends on the total volume of the space \mathcal{X} . It is easy to show that, for any $j \in \mathbf{C}$, the kernel $k_j(\mathbf{x}, \boldsymbol{\beta})$ is a Lipschitz function in $\boldsymbol{\beta}$ with Lipschitz constant $K_{\mathbf{x}} = \max_{i \leq J} \{|\mathbf{x}_j - \mathbf{x}_i|\}$. Hence,

$$\int_{\mathcal{X}} |k_j(\mathbf{x}, \boldsymbol{\beta}_1) - k_j(\mathbf{x}, \boldsymbol{\beta}_2)| M(d\mathbf{x}) \leq K |\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2|,$$

where $K = \sup_{\mathbf{x} \in \mathcal{X}} K_{\mathbf{x}} < +\infty$. Given δ , let N be the smallest integer greater than $4aK/\delta$ and cover C_a with a set of balls E_i of radius $2a/N$ so that, for any $\boldsymbol{\beta}_1, \boldsymbol{\beta}_2 \in E_i$, $|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2| < 4a/N$. This leads to $\int_{\mathcal{X}} |k_j(\mathbf{x}, \boldsymbol{\beta}_1) - k_j(\mathbf{x}, \boldsymbol{\beta}_2)| M(d\mathbf{x}) \leq \delta$. The number of balls necessary to cover C_a is then smaller than N^d . Using arguments similar to those used in Ghosal, Ghosh and Ramamoorthi (1999), Lemma 1, it can be shown that $J(2\delta, \mathcal{F}_a) \leq N^d(1 + \log[(1 + \delta)/\delta])$, from which (32) follows.

Step (2). Define $\mathcal{F}_{a, \delta} = \{\mathbf{q}(\mathbf{x}; G) : G(C_a) \geq 1 - \delta\}$. Then

$$J(\delta, \mathcal{F}_{a, \delta}) \leq K_\delta a^d \tag{33}$$

for a constant K_δ depending on δ . To see this, take $\mathbf{q}(\mathbf{x}; G) \in \mathcal{F}_{a, \delta}$ and denote by G^* the probability measure in \mathbb{P} defined by $G^*(A) = G(A \cap C_a)/G(C_a)$ so that $\mathbf{q}(\mathbf{x}; G^*)$ belongs to \mathcal{F}_a . It is easy to verify that $d_j(\mathbf{q}(\cdot; G^*), \mathbf{q}(\cdot; G)) < 2\delta$. It follows that $J(3\delta, \mathcal{F}_{a, \delta}) \leq J(\delta, \mathcal{F}_a)$, from which (33) follows.

Step (3). We follow here a technique used by Lijoi, Prünster and Walker (2005), Section 3.2. For the sequence $(a_n)_{n \geq 1}$ introduced before, define

$$\mathcal{F}_{a_n, \delta}^U = \{\mathbf{q}(\mathbf{x}; G) : G(C_n) \geq 1 - \delta\} \quad \text{and} \quad \mathcal{F}_{a_n, \delta}^L = \{\mathbf{q}(\mathbf{x}; G) : G(C_n) < 1 - \delta\}.$$

By construction, $\mathcal{G}_{a_n, \delta} \subset \mathcal{F}_{a_n, \delta}^U$ and $\mathcal{G}_{a_n, \delta} \subset \mathcal{F}_{a_{n-1}, \delta}^L$. Moreover, $\mathcal{F}_{a_{n-1}, \delta}^L \downarrow \emptyset$ as n increases to $+\infty$, thus, for any $\eta > 0$, there exists an integer n_0 such that, for any $n \geq n_0$, $J(\eta, \mathcal{F}_{a_n, \delta}^L) \leq J(\eta, \mathcal{F}_{a_{n_0}, \delta}^U)$. By (33), it follows that

$$J(\eta, \mathcal{G}_{a_n, \delta}) \leq K_\delta a_{n_0}^d \tag{34}$$

for any $n \geq n_0$, but, since $\mathcal{G}_{a_n, \delta} \subset \mathcal{F}_{a_n, \delta}^U$ and $\mathcal{F}_{a_n, \delta}^U \uparrow \mathcal{Q}$, (34) is also true for any $n < n_0$. Result (31) is then verified by setting $M_\delta = K_\delta a_{n_0}^d$. \square

Acknowledgments

The authors are grateful to an anonymous referee for his valuable comments and suggestions which led to a substantial improvement of the paper. Special thanks are also due to A. Lijoi and I. Prünster for some useful discussions. P. De Blasi was partially supported by Regione Piemonte. J.W. Lau’s research was partly supported by Hong Kong RGC Grant #601707. L.F. James was supported in part by grants HIA05/06.BM03, RGC-HKUST 6159/02P, DAG04/05.BM56 and RGC-HKUST 600907 of the HKSAR.

References

Antoniak, C.E. (1974). Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *Ann. Statist.* **2** 1152–1174. [MR0365969](#)

Barron, A., Schervish, M.J. and Wasserman, L. (1999). The consistency of posterior distributions in nonparametric problems. *Ann. Statist.* **27** 536–561. [MR1714718](#)

Bhat, C. (1998). Accommodating variations in responsiveness to level-of-service variables in travel mode choice models. *Transpn. Res. A* **32** 495–507.

Brownstone, D. and Train, K.E. (1999). Forecasting new product penetration with flexible substitution patterns. *J. Econometrics* **89** 109–129.

Cardell, N. and Dunbar, F. (1980). Measuring the societal impacts of automobile downsizing. *Transpn. Res. A* **14** 423–434.

Choi, T. and Schervish, M.J. (2007). On posterior consistency in nonparametric regression problems. *J. Multivariate Anal.* **98** 1969–1987. [MR2396949](#)

Choudhuri, N., Ghosal, S. and Roy, A. (2005). Bayesian methods for function estimation. In *Handbook of Statistics* (D. Dey, ed.) **25** 377–418. Amsterdam: Elsevier.

Dubé, J.P., Chintagunta, P., Bronnenberg, B., Goettler, R., Petrin, A., Seetharaman, P.B., Sudhir, K., Thomadsen, R. and Zhao, Y. (2002). Structural applications of the discrete choice model. *Marketing Letters* **13** 207–220.

Erdem, T. (1996). A dynamic analysis of market structure based on panel data. *Marketing Science* **15** 359–378.

Ferguson, T.S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1** 209–230. [MR0350949](#)

Ghosal, S., Ghosh, J.K. and Ramamoorthi, R.V. (1999). Posterior consistency of Dirichlet mixtures in density estimation. *Ann. Statist.* **27** 143–158. [MR1701105](#)

- Ghosal, S. and Roy, A. (2006). Posterior consistency of Gaussian process prior for nonparametric binary regression. *Ann. Statist.* **34** 2413–2429. [MR2291505](#)
- Ghosal, S. and Tang, Y. (2006). Bayesian consistency for Markov processes. *Sankhyā* **68** 227–239. [MR2303082](#)
- Ishwaran, H. and James, L.F. (2001). Gibbs sampling methods for stick-breaking priors. *J. Amer. Statist. Assoc.* **96** 161–173. [MR1952729](#)
- Ishwaran, H. and James, L.F. (2002). Approximate Dirichlet process computing in finite normal mixtures: Smoothing and prior information. *J. Comp. Graph. Statist.* **11** 508–532. [MR1938445](#)
- Ishwaran, H. and Zarepour, M. (2000). Markov chain Monte Carlo in approximate Dirichlet and beta two-parameter process hierarchical models. *Biometrika* **87** 371–390. [MR1782485](#)
- James, L.F., Lijoi, A. and Prünster, I. (2009). Posterior analysis for normalized random measures with independent increments. *Scand. J. Statist.* **36** 76–97. [MR2508332](#)
- Lijoi, A., Prünster, I. and Walker, S.G. (2005). On consistency of nonparametric normal mixtures for Bayesian density estimation. *J. Amer. Statist. Assoc.* **100** 1292–1296. [MR2236442](#)
- Lo, A.Y. (1984). On a class of Bayesian nonparametric estimates: I. Density estimates. *Ann. Statist.* **12** 351–257. [MR0733519](#)
- MacEachern, S.N. and Muller, P. (1998). Estimating mixture of Dirichlet process models. *J. Comput. Graph. Statist.* **7** 223–238.
- McFadden, D. (1974). Conditional logit analysis of qualitative choice behavior. In *Frontiers of Econometrics* (P. Zarembka, ed.) 105–142. New York: Academic Press.
- McFadden, D. and Train, K.E. (2000). Mixed MNL models for discrete response. *J. Appl. Econometrics* **15** 447–470.
- Pitman, J. and Yor, M. (1997). The two-parameter Poisson–Dirichlet distribution derived from a stable subordinator. *Ann. Probab.* **25** 855–900. [MR1434129](#)
- Regazzini, E., Lijoi, A. and Prünster, I. (2003). Distributional results for means of random measures with independent increments. *Ann. Statist.* **31** 560–585. [MR1983542](#)
- Srinivasan, K. and Mahmassani, H. (2005). A dynamic kernel logit model for the analysis of longitude discrete choice data: Properties and computational assessment. *Transportation Science* **39** 160–181.
- Train, K.E. (2003). *Discrete Choice Methods with Simulation*. Cambridge Univ. Press. [MR2003007](#)
- Train, K.E. (2008). EM algorithms for nonparametric estimation of mixing distributions. *Journal of Choice Modelling* **1** 40–69.
- Walker, J., Ben-Akiva, M. and Bolduc, D. (2007). Identification of parameters in normal error component logit-mixture (NECLM) models. *J. Appl. Econometrics* **22** 1095–1125. [MR2408974](#)
- Walker, S.G. (2003a). On sufficient conditions for Bayesian consistency. *Biometrika* **90** 482–488. [MR1986664](#)
- Walker, S.G. (2003b). Bayesian consistency for a class of regression problems. *South African Statist. J.* **37** 151–169. [MR2042627](#)
- Walker, S.G. (2004). New approaches to Bayesian consistency. *Ann. Statist.* **32** 2028–2043. [MR2102501](#)
- Walker, S.G., Lijoi, A. and Prünster, I. (2005). Data tracking and the understanding of Bayesian consistency. *Biometrika* **92** 765–778. [MR2234184](#)
- Wasserman, L. (1998). Asymptotic properties of nonparametric Bayesian procedures. In *Practical Nonparametric and Semiparametric Bayesian Statistics* (D. Dey, P. Muller and D. Sinha, eds.) 293–304. New York: Springer-Verlag. [MR1630088](#)

Received July 2008 and revised April 2009