# Retrospective Markov chain Monte Carlo methods for Dirichlet process hierarchical model

By OMIROS PAPASPILIOPOULOS AND GARETH O. ROBERTS

*Department of Statistics, Warwick University, Coventry, CV4 7AL, U.K.*

O.Papaspiliopoulos@warwick.ac.uk,        Gareth.O.Roberts@warwick.ac.uk

SUMMARY

Inference for Dirichlet process hierarchical models is typically performed using Markov chain Monte Carlo methods, which can be roughly categorised into marginal and conditional methods. The former integrate out analytically the infinite-dimensional component of the hierarchical model and sample from the marginal distribution of the remaining variables using the Gibbs sampler. Conditional methods impute the Dirichlet process and update it as a component of the Gibbs sampler. Since this requires imputation of an infinite-dimensional process, implementation of the conditional method has relied on finite approximations. In this paper we show how to avoid such approximations by designing two novel Markov chain Monte Carlo algorithms which sample from the exact posterior distribution of quantities of interest. The approximations are avoided by the new technique of retrospective sampling. We also show how the algorithms can obtain samples from functionals of the Dirichlet process. The marginal and the conditional methods are compared and a careful simulation study is included, which involves a non-conjugate model, different datasets and prior specifications.

*Some keywords:* Exact simulation; Mixture models; Label switching; Retrospective sampling; Stick-breaking prior

## 1. INTRODUCTION

Dirichlet process hierarchical models, also known as Dirichlet mixture models, are now standard in semiparametric inference; applications include density estimation (Müller et al.,

1996), survival analysis (Gelfand & Kottas, 2003), semi-parametric analysis of variance (Müller et al., 2005), cluster analysis and partition modelling (Quintana & Iglesias, 2003), meta-analysis (Burr & Doss, 2005) and machine learning (Teh et al., 2006). The hierarchical structure is as follows. Let $f(y \mid z, \lambda)$ be a parametric density with parameters $z$ and $\lambda$, let $H_\Theta$ be a distribution indexed by some parameters $\Theta$ and let $\mathrm{Be}(\alpha, \beta)$ denote the beta distribution with parameters $\alpha, \beta$. Then we have

$$
\begin{aligned}
Y_i \mid (Z, K, \lambda) &\sim f(Y_i \mid Z_{K_i}, \lambda), \; i = 1, \ldots, n \\
K_i \mid p &\sim \sum_{j=1}^{\infty} p_j \delta_j(\cdot) \\
Z_j \mid \Theta &\sim H_\Theta, \; j = 1, 2, \ldots \\
p_1 = V_1, \qquad p_j &= (1 - V_1)(1 - V_2) \cdots (1 - V_{j-1}) V_j, \; j \geq 2 \\
V_j &\sim \mathrm{Be}(1, \alpha).
\end{aligned}
\tag{1}
$$

Here $V = (V_1, V_2, \ldots)$ and $Z = (Z_1, Z_2, \ldots)$ are vectors of independent variables and are independent of each other, the $K_i$s are independent given $p = (p_1, p_2, \ldots)$, and the $Y_i$s are independent given $Z$ and $K = (K_1, K_2, \ldots, K_n)$. Therefore, (1) defines a mixture model in which $f(y \mid z, \lambda)$ is mixed with respect to the discrete random probability measure $P(dz)$, where

$$
P(\cdot) = \sum_{j=1}^{\infty} p_j \delta_{Z_j}(\cdot),
\tag{2}
$$

and $\delta_x(\cdot)$ denotes the Dirac delta measure centred at $x$. Note that, the discreteness of $P$ implies a distributional clustering of the data. The last three lines in the hierarchy identify $P$ with the Dirichlet process with base measure $H_\Theta$ and concentration parameter $\alpha$; this representation is due to Sethuraman (1994). For the construction of the Dirichlet process prior see Ferguson (1973,1974), for properties of Dirichlet mixture models see for example Lo (1984), Ferguson (1983), and Antoniak (1974), and for a recent article on modelling with Dirichlet processes see Green & Richardson (2001). Using more general beta distributions in the last line gives rise to the rich class of stick-breaking random measures and the general hierarchical framework introduced by Ishwaran & James

(2001). We have consciously been vague about the spaces on which $Y_i$ and $Z_j$ live; it should be obvious from the construction that a great deal of flexibility is allowed.

We refer to $Z_j$ and $p_j$ as the parameters and the weight respectively of the $j$th component in the mixture, and to $K_i$ as the allocation variable, which determines to which component the $i$th data point is allocated. Note that the prior on the component weights $p_j$ imposes a weak identifiability on the mixture components, in the sense that $E(p_j) \geq E(p_l)$ for any $j \geq l$. Throughout the paper we use the nomenclature 'cluster' to refer to a mixture component with at least one datum allocated to it. The hyperparameters $\alpha, \Theta$ and $\lambda$ will be consider fixed until §3·7. We will concentrate on inference for the allocation variables, the component parameters and weights, and $P$ itself.

Inference for Dirichlet mixture models has been made feasible by Gibbs sampling techniques which have been developed since the seminal work in an unpublished 1988 Yale University Ph. D. thesis by M. Escobar. Alternative Monte Carlo schemes do exist and include the sequential samplers of Liu (1996) and Ishwaran & James (2003), and referred to as the blocked Gibbs sampler in those papers, the particle filter of Fearnhead (2004) and the reversible jump method of Green & Richardson (2001) and Jain & Neal (2004). We concentrate on Gibbs sampling and related componentwise updating algorithms.

Broadly speaking, there are two possible Gibbs sampling strategies, corresponding to two different data augmentation schemes. The marginal method exploits convenient mathematical properties of the Dirichlet process and integrates out analytically the random probabilities $p$ from the hierarchy. Then, using a Gibbs sampler, it obtains samples from the posterior distribution of the $K_i$'s and the parameters and the weights of the clusters. Integrating out $p$ induces prior dependence among the $K_i$'s, and makes the labels of the clusters unidentifiable. This approach is easily carried out for conjugate models, where further analytic integrations are possible; (1) is said to be conjugate when $H_\Theta(dz)$ and $f(y \mid z, \lambda)$ form a conjugate pair for $z$. Implementation of the marginal method for non-conjugate models is considerably more complicated. MacEachern & Müller (1998)

3

and Neal (2000) are the state-of-the-art implementations in this context.

The conditional method works with the augmentation scheme described in (1). It consists of the imputation of $p$ and $Z$, and subsequent Gibbs sampling from the joint posterior distribution of $(K, p, Z)$. The conditional independence structure created with the imputation of $(p, Z)$ assists the simultaneous updating of large subsets of the variables. The conditional method was introduced in Ishwaran & Zarepour (2000), it was further developed in Ishwaran & James (2001, 2003) and it has two considerable advantages over the marginal method. First, it does not rely on being able to integrate out analytically components of the hierarchical model, and therefore it is much more flexible for current and future elaborations of the basic model. Such extensions include more general stick-breaking random measures, and modelling dependence of the data on covariates, as for example in Dunson & Park (2006). A second advantage is that in principle it allows inference for the latent random measure $P$.

The implementation of the conditional method poses interesting methodological challenges, since it requires the imputation of the infinite-dimensional vectors $p$ and $Z$. The modus operandi advocated in Ishwaran & Zarepour (2000) is to approximate the vectors, and thus the corresponding Dirichlet process prior, using some kind of truncation, such as $p^{(N)} = (p_1, \ldots, p_N)$ and $Z^{(N)} = (Z_1, \ldots, Z_N)$, where $N$ determines the degree of approximation. Although in some cases it is possible to control the error produced by such truncations (Ishwaran & Zarepour, 2000, 2002; Ishwaran & James, 2001) it would be desirable to avoid approximations altogether.

This paper introduces two implementations of the conditional method which avoid any approximation. The proposed Markov chain Monte Carlo algorithms are very easily implemented and can readily be extended to more general stick-breaking models. The implementation of the conditional method is achieved by retrospective sampling, which is introduced in this paper. This is a novel technique which facilitates exact simulation in finite time in problems which involve infinite-dimensional processes. We also show how to use retrospective sampling in conjunction with our Markov chain Monte Carlo

4

algorithms in order to sample from the posterior distribution of functionals of $P$.

We identify a computational problem with the conditional method. As a result of the weak identifiability imposed on the labels of the mixture components, the posterior distribution of the random measure $(p, Z)$ is multimodal. Therefore, the Gibbs sampler has to visit all these different modes. This problem is not eliminated with large datasets, since although the secondary modes become smaller, the energy gap between the modes becomes bigger, and thus the sampler can get trapped in low-probability areas of the state-space of $(p, Z)$. We design tailored label-switching moves which improve significantly the performance of our Markov chain Monte Carlo algorithm.

We also contrast our retrospective Markov chain Monte Carlo algorithm with state-of-the-art implementations of the marginal method for non-conjugate models in terms of their Monte Carlo efficiency. In particular we consider the no-gaps algorithm of MacEachern & Müller (1998) and Algorithms 7 and 8 of Neal (2000). This comparison sheds light on the relative merits of the marginal and conditional method in models where they can both be applied. We find that the marginal methods slightly outperform the conditional.

Since the original submission of this paper, the retrospective sampling ideas we introduce here have been found very useful in extensions of the Dirichlet process hierarchical model (Dunson & Park, 2006; Griffin, 2006), and in the exact simulation of diffusion processes; see for example Beskos et al. (2006b).

## 2. RETROSPECTIVE SAMPLING FROM THE DIRICHLET PROCESS PRIOR

Consider the task of simulating a sample $X = (X_1, \ldots, X_n)$ from the Dirichlet process prior (2). Such a sample has the property that $X_i \mid P \sim P$, $X_i$ is independent of $X_l$ conditionally on $P$ for each $l \neq i$, $i, l = 1, \ldots, n$, and $P$ is the Dirichlet process with concentration parameter $\alpha$ and base measure $H_\Theta$. This section introduces in a simple context the technique of retrospective sampling and the two different computational approaches, marginal and conditional, to inference for Dirichlet processes. There are

essentially two different ways of obtaining $X$.

The sample $X$ can be simulated directly from its marginal distribution. When $P$ is marginalised the joint distribution of the $X_i$s is known and it is described by a Pólya urn scheme (Blackwell & MacQueen, 1973; Ferguson, 1973). In particular, let $K_i$ be an allocation variable associated with $X_i$, where $K_i = j$ if and only if $X_i = Z_j$. Therefore, the $K_i$'s decide which component of the infinite series (2) $X_i$ is associated with. Conditionally on $P$ the $K_i$'s are independent, with

$$\mathrm{pr}(K_i = j \mid p) = p_j, \quad \text{for all } i = 1, \ldots, n, \ j = 1, 2, \ldots. \tag{3}$$

The marginal prior of $K = (K_1, \ldots, K_n)$ is obtained by the following Pólya urn scheme:

$$
\begin{aligned}
\mathrm{pr}(K_i = j \mid K_1, \ldots, K_{i-1}) &= \frac{n_{i,j}}{(\alpha + i - 1)}, \quad \text{if } K_l = j \text{ for some } l < i\,, \\
\mathrm{pr}(K_i \neq K_l \text{ for all } l < i \mid K_1, \ldots, K_{i-1}) &= \frac{\alpha}{(\alpha + i - 1)}\,,
\end{aligned}
\tag{4}
$$

where $n_{i,j}$ denotes the size of the set $\{l < i : K_l = j\}$. Thus, the probability that the $i$th sample is associated with the $j$th component is proportional to the number of samples already associated with $j$, whereas with probability $\alpha/(\alpha + i - 1)$ the $i$th sampled value is associated with a new component. Note that, whereas in (3) the labels $j$ are identifiable, in (4) the labels of the clusters are totally arbitrary. Simulation of $X$ from its marginal distribution proceeds in following way. Set $K_1 = 1$, although any other label could be chosen; simulate $\phi_1 \sim H_\Theta$; and set $X_1 = \phi_1$. For $i > 1$, let $c$ denote the number of existing clusters; simulate $K_i$ conditionally on $K_1, \ldots, K_{i-1}$ according to the probabilities (4), where $j = 1, \ldots, c$; if $K_i = j$ then set $X_i = \phi_j$, and otherwise set $c = c + 1$, simulate $\phi_c \sim H_\Theta$ independently of any previously drawn values, and set $X_i = \phi_c$. At the end of this procedure $\phi_F = (\phi_1, \ldots, \phi_c)$ are the parameters of the mixture components which are associated with at least one data point. The indexing $1, \ldots, k$ is arbitrary and it is not feasible to map these $\phi_j$'s to the $Z_j$'s in the definition of the Dirichlet process in (2).

Alternatively, $X$ can be simulated following a two-step hierarchical procedure. Initially, a realisation of $P$ is simulated, and then we simulate independently $n$ allocation

variables according to (3) and we set $X_i = Z_{K_i}$. However, simulation of $P$ entails the generation of the infinite vectors $p$ and $Z$, which is infeasible. Nevertheless, this problem can be avoided with the following retrospective simulation. The standard method for simulating from the discrete distribution defined in (3) is first to simulate $U_i$ from a uniform distribution on $(0, 1)$, and then to set $K_i = j$ if and only if

$$\sum_{l=0}^{j-1} p_l < U_i \leq \sum_{l=1}^{j} p_l \ , \tag{5}$$

where we define $p_0 = 0$; this is the inverse cumulative distribution function method for discrete random variables; see for example Ripley (1987, §3.3). Retrospective sampling simply exchanges the order of simulation between $U_i$ and the pairs $(p_j, Z_j)$. Rather than simulating first $(p, Z)$ and then $U_i$ in order to check (5), we first simulate the decision variable $U_i$ and then pairs $(p_j, Z_j)$. If for a given $U_i$ we need more $p_j$'s than we currently have in order to check (5), then we go back and simulate pairs $(p_j, Z_j)$ 'retrospectively', until (5) is satisfied. The algorithm proceeds as follows.

ALGORITHM 1. *Retrospective sampling from the Dirichlet process prior*

*Step 1. Simulate $p_1$ and $Z_1$, and set $N^* = 1$, $i = 1$ and $p_0 = 0$.*

*Step 2. Repeat the following until $i > n$.*

*Step 2.1. Simulate $U_i \sim Un[0, 1]$*

*Step 2.2. If (5) is true for some $k \leq N^*$ then set $K_i = k$, $X_i = Z_k$, for $i = i + 1$, and go to Step 2.*

*Step 2.3. If (5) is not true for any $k \leq N^*$ then set $N^* = N^* + 1$, $j = N^*$, simulate $p_j$ and $Z_j$, and go to Step 2.2.*

In this notation, $N^*$ keeps track of how far into the infinite sequence $\{(p_j, Z_j), \ j = 1, 2, \ldots\}$ we have visited during the simulation of the $X_i$'s. Note that the retrospective sampling can be easily implemented because of the Markovian structure of the $p_j$'s and the independence of the $Z_j$'s. A similar scheme was advocated in Doss (1994).

The previous simulation scheme illustrates the main principle behind retrospective

sampling: although it is impossible to simulate an infinite-dimensional random object we might still be able to take decisions which depend on such objects exactly avoiding any approximations. The success of the approach will depend on whether or not the decision problem can be formulated in a way that involves finite-dimensional summaries, possibly of random dimension, of the random object. In the previous toy example we formulated the problem of simulating draws from the Dirichlet process prior as the problem of comparing a uniform random variable with partial sums of the Dirichlet random measure. This facilitated the simulation of $X$ in finite time avoiding approximation errors. In general, the retrospective simulation scheme will require at the first stage simulation of both the decision variable and certain finite-dimensional summaries of the infinite-dimensional random object. Thus, at the second stage we will need to simulate retrospectively from the distribution of the random object conditionally on these summaries. This conditional simulation will typically be much more elaborate than the illustration we gave here.

We shall see in §3 that these ideas extend to posterior simulation in a computationally feasible way. However, the details of the method for posterior simulation are far more complicated than the simple method given above.

## 3. The retrospective conditional method for posterior simulation

### 3·1 *Posterior inference*

When we fit (1) to data $Y = (Y_1, \ldots, Y_n)$, there are a number of quantities about which we may want to make posterior inference. These include the classification variables $K_i$, which can be used to classify the data into clusters, the number of clusters in the population, the cluster parameters $\{Z_j : K_i = j$ for at least one $i\}$ and the cluster probabilities $\{p_j : K_i = j$ for at least one $i\}$, the hyperparameters $\alpha, \Theta$ and $\lambda$, the predictive distribution of future data, and the random measure $P$ itself. None of the existing methods can provide samples from the posterior distribution of $P$ without resorting to some kind of approximation; see for example Gelfand & Kottas (2002) for some suggestions. However, exact posterior simulation of finite-dimensional distributions and functionals of $P$ might

be feasible; see §3·6.

Markov chain Monte Carlo techniques have been developed for sampling-based exploration of the posterior distributions outlined above. The conditional method (Ishwaran & Zarepour, 2000) is based on an augmentation scheme in which the random probabilities $p = (p_1, p_2, \ldots)$ and the component parameters $Z = (Z_1, Z_2, \ldots)$ are imputed and the Gibbs sampler is used to sample from the joint posterior distribution of $(K, p, Z)$ according to the full-conditional distributions of the variables. However, the fact that $p$ and $Z$ have countably infinite elements poses a major challenge. Ishwaran & Zarepour (2000) advocated an approximation of these vectors, and therefore an approximation of the corresponding Dirichlet process prior. For instance a truncation of the Sethuraman representation (2) could be adopted, e.g. $p^{(N)} = (p_1, \ldots, p_N)$ and $Z^{(N)} = (Z_1, \ldots, Z_N)$, where $N$ determines the degree of approximation.

In this section we show how to avoid any such approximation. The first step of our solution is to parameterise in terms of $(K, V, Z)$, where $V = (V_1, V_2, \ldots)$ are the beta random variables used in the stick-breaking construction of $p$. We will construct algorithms which effectively can return samples from the posterior distribution of these vectors, by simulating iteratively from the full conditional posterior distributions of $K$, $V$ and $Z$. Note that we can replace direct simulation from the conditional distributions with any updating mechanism, such as a Metropolis-Hastings step, which is invariant with respect to these conditional distribution. The stationary distribution of this more general componentwise-updating sampler is still the joint posterior distribution of $(K, V, Z)$.

We now introduce some notation. For a given configuration of the classification variables $K = (K_1, \ldots, K_n)$, we define

$$m_j = \sum_{i=1}^{n} 1_{\{K_i = j\}}, \quad j = 1, 2, \ldots,$$

to be the number of data points allocated to the $j$th component of the mixture. Moreover,

again for a given configuration of $K$ let

$$
\begin{aligned}
I &= \{1, 2, \ldots\} \\
I^{(\mathrm{al})} &= \{j \in I : m_j > 0\} \\
I^{(\mathrm{d})} &= \{j \in I : m_j = 0\} = I - I^{(\mathrm{al})}.
\end{aligned}
$$

Therefore, $I$ represents all components in the infinite mixture, $I^{\mathrm{al}}$ is the set of all 'alive' components and $I^{()}$ the set of 'dead' components, where we call a component 'alive' if some data have been allocated to it. The corresponding partition of $Z$ and $V$ will be denoted by $Z^{(\mathrm{al})}, Z^{()}, V^{(\mathrm{al})}$ and $V^{()}$.

The implementation of our Markov chain Monte Carlo algorithms relies on the results contained in the following Proposition, which describes the full conditional posterior distributions of $Z, V$ and $K$.

**Proposition 1.** *Let $\pi(z \mid \Theta)$ be the density of $H_\Theta(dz)$. Conditionally on $(Y, K, \Theta, \lambda)$, $Z$ is independent of $(V, \alpha)$ and it consists of conditionally independent elements with*

$$
Z_j \mid Y, K, \Theta, \lambda \sim
\begin{cases}
H_\Theta, \text{ for all } j \in I^{()} \\
\\
\prod_{i:K_i=j} f(Y_i \mid Z_j, \lambda) \pi(Z_j \mid \Theta) \text{ for all } j \in I^{(\mathrm{al})},
\end{cases}
$$

*Conditionally on $(K, \alpha)$, $V$ is independent of $(Y, Z, \Theta, \lambda)$ and it consists of conditionally independent elements with*

$$
V_j \mid K, \alpha \sim \mathrm{Be}\left(m_j + 1, n - \sum_{l=1}^{j} m_l + \alpha\right) \text{ for all } j = 1, 2, \ldots.
$$

*Conditionally on $(Y, V, Z, \lambda)$, $K$ is independent of $(\Theta, \alpha)$ and it consist of conditionally independent elements with*

$$
\mathrm{pr}\{K_i = j \mid Y, V, Z, \lambda\} \propto p_j f(Y_i \mid Z_j, \lambda), \; j = 1, 2, \ldots.
$$

The proof of the Proposition follows directly from the conditional independence structure in the model and the stick-breaking representation of the Dirichlet process; see Ishwaran & James (2003) for a proof in a more general context.

The conditional independence structure in the model effects the simultaneous sampling of any finite subset of pairs $(V_j, Z_j)$ according to their full conditional distributions. Simulation of $(Z_j, V_j)$ when $j \in I^{()}$ is trivial. When dealing with non-conjugate models, the distribution of $Z_j$ for $j \in I^{(\mathrm{al})}$ will typically not belong to a known family and a Metropolis-Hastings step might be used instead to carry out this simulation. Note that samples from the full conditional posterior distribution of $p_j$'s are obtained using samples of $(V_h, h \leq j)$ and the stick-breaking representation in (1).

The $K_i$'s are conditionally independent given $Z$ and $V$. However, the normalising constant of the full conditional probability mass function of each $K_i$ is intractable:

$$c_i = \sum_{j=1}^{\infty} p_j f(Y_i \mid Z_j, \lambda).$$

The intractability stems from the fact that an infinite sum of random terms needs to be computed. At this stage one could resort to a finite approximation of the sums, but we wish to avoid this. The unavailability of the normalising constants renders simulation of the $K_i$'s highly nontrivial. Therefore, in order to construct a conditional Markov chain Monte Carlo algorithm we need to find ways of sampling from the conditional distribution of the $K_i$'s.

### 3·2  *A retrospective quasi-independence Metropolis-Hastings sampler for the allocation variables*

One simple and effective way of avoiding the computation of the normalising constants $c_i$ is to replace direct simulation of the $K_i$s with a Metropolis-Hastings step. Let $k = (k_1, \ldots, k_n)$ denote a configuration of $K = (K_1, \ldots, K_n)$, and let

$$\max\{k\} = \max_i k_i$$

be the maximal element of the vector $k$. We assume that we have already obtained samples from the conditional posterior distribution of $\{(V_j, Z_j) : j \leq \max\{k\}\}$ given $K = k$. Note that the distribution of $(V_j, Z_j)$ conditionally on $Y$ and $K = k$ is simply the prior, $V_j \sim \mathrm{Be}(\alpha, 1), Z_j \sim H_\Theta$, for any $j > \max\{k\}$.

We will describe an updating scheme $k \mapsto k^*$ which is invariant with respect to the full conditional distribution of $K \mid Y, Z, V, \lambda$. The scheme is a composition of $n$ Metropolis-Hastings steps which update each of the $K_i$s in turn. Let

$$k(i, j) = (k_1, \ldots, k_{i-1}, j, k_{i+1}, \ldots, k_n)$$

be the vector produced from $k$ by substituting the $i$th element by $j$. When updating $K_i$, the sampler proposes to move from $k$ to $k(i, j)$, where the proposed $j$ is generated from the probability mass function

$$q_i(k, j) \propto \begin{cases} p_j f(Y_i \mid Z_j, \lambda), & \text{for } j \leq \max\{k\} \\ \\ M_i(k) p_j, & \text{for } j > \max\{k\} . \end{cases} \tag{6}$$

The normalising constant of (6) is

$$\tilde{c}_i(k) = \sum_{j=1}^{\max\{k\}} p_j f(Y_i \mid Z_j, \lambda) + M_i(k) \left( 1 - \sum_{j=1}^{\max\{k\}} p_j \right),$$

which can be easily computed given $\{(p_j, Z_j) : j \leq \max\{k\}\}$. Note that, for $j \leq \max\{k\}$, $q_i(k, j) \propto \mathrm{pr}(K_i = j \mid Y, V, Z, \lambda)$, while, for $j > \max\{k\}$, $q_i(k, j) \propto \mathrm{pr}(K_i = j \mid V)$. Here $M_i(k)$ is a user-specified parameter which controls the probability of proposing $j > \max\{k\}$, and its choice will be discussed in §3·3.

According to this proposal distribution the $i$th data point is proposed to be reallocated to one of the alive clusters $j \in I^{(\mathrm{al})}$ with probability proportional to the conditional posterior probability $\mathrm{pr}(K_i = j \mid Y, V, Z, \lambda)$. Allocation of the $i$th data point to a new component can be accomplished in two ways: by proposing $j \in I^{()}$ for $j \leq \max\{k\}$ according to the conditional posterior probability mass function and by proposing $j \in I^{()}$ for $j > \max\{k\}$ according to the prior probability mass function. Therefore, a careful calculation yields that the Metropolis-Hastings acceptance probability of the transition

from $K = k$ to $K = k(i, j)$ is

$$\alpha_i\{k, k(i,j)\} = \begin{cases} 1, & \text{if } j \leq \max\{k\} \text{ and } \max\{k(i,j)\} = \max\{k\} \\[3mm] \min\left\{1, \frac{\tilde{c}_i(k)M\{k(i,j)\}}{\tilde{c}_i\{k(i,j)\}f(Y_i|Z_{k_i},\lambda)}\right\}, & \text{if } j \leq \max\{k\} \text{ and } \max\{k(i,j)\} < \max\{k\} \\[3mm] \min\left\{1, \frac{\tilde{c}_i(k)f(Y_i|Z_j,\lambda)}{\tilde{c}_i\{k(i,j)\}M(k)}\right\}, & \text{if } j > \max\{k\}. \end{cases}$$

Note that proposed re-allocations to a component $j \leq \max\{k\}$ are accepted with probability 1 provided that $\max\{k(i,j)\} = \max\{k\}$. If the proposed move is accepted, we set $k = k(i,j)$ and proceed to the updating of $K_{i+1}$. The composition of all these steps yields the updating mechanism for $K$.

Simulation from the proposal distribution is achieved by retrospective sampling. For each $i = 1, \ldots, n$, we simulate $U_i \sim \text{Un}[0,1]$ and propose to set $K_i = j$, where $j$ satisfies

$$\sum_{l=0}^{j-1} q_i(k, l) \; < U_i \; \leq \sum_{l=1}^{j} q_i(k, l), \tag{7}$$

with $q_i(k,0) \equiv 0$. If (7) is not satisfied for any $j \leq \max\{k\}$, then we start checking the condition for $j > \max\{k\}$ until it is satisfied. This will require the values of $(V_l, Z_l), l > \max\{k\}$. If these values have not already been simulated in the previous steps of the algorithm, they are simulated retrospectively from their prior distribution when they become needed.

We therefore have a retrospective Markov chain Monte Carlo algorithm for sampling from the joint posterior distribution of $(K, V, Z)$, which is summarised below.

ALGORITHM 2. *Retrospective Markov chain Monte Carlo*

*Give an initial allocation $k = (k_1, \ldots, k_n)$, and set $N^* = \max\{k\}$*

*Step 1. Simulate $Z_j$ from its conditional posterior, $j \leq \max\{k\}$.*

*Step 2.1 Simulate $V_j$ from its conditional posterior, $j \leq \max\{k\}$.*

*Step 2.2 Calculate $p_j = (1 - V_1) \cdots (1 - V_{j-1})V_j, \; j \leq \max\{k\}$.*

*Step 3.1. Repeat the following until $i > n$.*

13

*Step* 3.2. *Simulate* $U_i \sim Un[0, 1]$

*Step* 3.3.1 *If (7) is true for some $j \leq N^*$ then set $K_i = j$ with probability $\alpha_i\{k, k(i, j)\}$,*
*otherwise leave it unchanged. Set $i = i + 1$ and go to Step* 3.1.

*Step* 3.3.2 *If (7) is not true for any $j \leq N^*$, set $N^* = N^* + 1$, $j = N^*$. Simulate $(V_j, Z_j)$*
*from the prior, set $p_j = (1 - V_1) \cdots (1 - V_{j-1}) V_j$ and go to Step* 3.3.1

*Step* 3.4. *Set $N^* = \max\{k\}$ and go to Step* 1.

Note that both $\max\{k\}$ and $N^*$ change during the updating of the $K_i$'s, with $N^* \geq \max\{k\}$. Thus, the proposal distribution is adapting itself to improve approximation of the target distribution. At the early stages of the algorithm $N^*$ will be large, but as the cluster structure starts being identified by the data then $N^*$ will take much smaller values. Nevertheless, the adaptation of the algorithm does not violate the Markov property. It is recommended to update the $K_i$'s in a random order, to avoid using systematically less efficient proposals for some of the variables.

The algorithm we have constructed updates all the allocation variables $K$ but only a random subset of $(V, Z)$, to be precise $\{(V_j, Z_j), j \leq N^*\}$. However, pairs of the complementary set $\{(V_j, Z_j), j > N^*\}$ can be simulated when and if they are needed retrospectively from the prior distribution. This simulation can be performed off-line after the completion of the algorithm. In this sense, our algorithm is capable of exploring the joint distribution of $(K, V, Z)$.

### 3·3    *Accelerations of the main algorithm*

There are several simple modifications of the retrospective Markov chain Monte Carlo algorithm which can improve significantly its Monte Carlo efficiency.

We first discuss the choice of the user-specified parameter $M_i(k)$ in (6). This parameter relates to the probability, $\rho$ say, of proposing $j > \max\{k\}$, where

$$\rho = \frac{\left(1 - \sum_{j=1}^{\max\{k\}} p_j\right) M_i(k)}{c_i(\max\{k\})}.$$

If it is desired to propose components $j > \max\{k\}$ a specific proportion of the time

then the equation above can be solved for $M_i(k)$. For example, $\rho = \alpha/(n-1)$ is the probability of proposing new clusters in Algorithm 7 of Neal (2000). We recommend a choice of $M_i(k)$ which guarantees that the probability of proposing $j > \max\{k\}$ is greater than the prior probability assigned to the set $\{j : j > \max\{k\}\}$. Thus $M_i(k)$ should satisfy

$$\sum_{j=1}^{\max\{k\}} p_j f(Y_i \mid Z_j, \phi) + M_i(k) \left(1 - \sum_{j=1}^{\max\{k\}} p_j\right) \leq M_i(k).$$

This inequality is satisfied by setting

$$M_i(k) = \max\left\{f(Y_i \mid Z_j, \phi), j \leq \max\{k\}\right\}. \tag{8}$$

Since (6) resembles an independence sampler for $K_i$, (8) is advisable from a theoretical perspective. Mengersen & Tweedie (1996) have shown that an independence sampler is geometrically ergodic if and only if the tails of the proposal distribution are heavier than the tails of the target distribution. The choice of $M_i(k)$ according to (8) ensures that the tails of the proposal $q_i(k, j)$ are heavier than the tails of the target probability $\text{pr}\{K_i = j \mid Y, p, Z, \lambda\}$. Note that when $M_i(k)$ is chosen according to (8) then $\rho$ is random. In simulation studies we have discovered that the distribution of $\rho$ is very skewed and the choice according to (8) leads to a faster mixing algorithm than alternative schemes with fixed $\rho$.

Another interesting possibility is to update $Z$ and $V$ after each update of the allocation variables. Before Step 3.2 of Algorithm 2 we simulate $Z^{()}$ from the prior and leave $Z^{(\text{al})}$ unchanged. Moreover, we can synchronise $N^*$ and $\max\{k\}$. Theoretically, this is achieved by pretending to update $\{V_j, j > \max\{k\}\}$ from the prior before Step 3.2. In practice, the only adjustment to the existing algorithm is to set $N^* = \max\{k\}$ before Step 3.2. These extra updates have the computational advantage of storing only $Z^{(\text{al})}$ and $\{V_j, j \leq \max\{k\}\}$. In addition, simulations have verified that these extra updates improve significantly the mixing of the algorithm. Morover, one can allow more components $j \in I^{()}$ to be proposed according to the posterior probability mass function by changing $\max\{k\}$ in (6) to $\max\{k\} + l$, for some fixed integer $l$. In that case the

15

acceptance probability (3·2) needs to be slightly modified, but we have not found gains from such adjustment.

### 3·4   *Multimodality and label-switching moves*

The most significant and crucial modification of the algorithm we have introduced is the addition of label-switching moves. These moves will have to be included in any implementation of the conditional method which updates the allocation variables one at a time. The augmentation of $p$ in the conditional method makes the components in the infinite mixture (1) weakly identifiable, in the sense that $E\{p_j\} \geq E\{p_l\}$ for any $j \geq l$, but there is nonnegligible prior probability that $p_l > p_j$, in particular when $|l - j|$ is small. As a result the posterior distribution of $(p, Z)$ exhibits multiple modes. In order to highlight this phenomenon we consider below a simplified scenario, but we refer to §4 for an illustration in the context of posterior inference for a specific non-conjugate Dirichlet process hierarchical model.

In particular, assume that we have actually observed a sample of size $n$ from the Dirichlet process $X = (X_1, \ldots, X_n)$ using the notation of §2. This is the limiting case of (1) where the observation density $f(y \mid z, \lambda)$ is a point mass at $z$. In this case we directly observe the allocation of the $X_i$'s into $c$, say, clusters, each of size $n_l$, say, where $\sum_{l=1}^{c} n_l = n$. The common value of the $X_i$'s in each cluster $l$ gives precisely the parameters $\phi_l$ of the corresponding cluster $l = 1, \ldots, c$. However, there is still uncertainty regarding the component probabilities $p_j$ of the Dirichlet process which generated the sample, and the index of the component for which each cluster has been generated. Let $K_l, l = 1, \ldots, c$ denote these indices for each of the clusters. We can construct a Gibbs sampler for exploring the posterior distribution of $p$ and $(K_1, \ldots, K_c)$: this is a nontrivial simulation, and a variation of the retrospective Markov chain Monte Carlo algorithm has to be used. Figure 1 shows the posterior densities of $(p_1, p_2, p_3)$ when $n = 10$ and $c = 3$ with $n_1 = 5, n_2 = 4, n_3 = 1$, and when $n = 100$ and $c = 3$ with $n_1 = 50, n_2 = 40, n_3 = 10$. In both cases we took $\alpha = 1$.

Note that the posterior distributions of the $p_j$s exhibit multiple modes, because the labels in the mixture are only weakly identifiable. Note that the secondary modes become less prominent for larger samples, but the energy gap between the modes increases. In contrast, the probability of the largest component, $\max\{p_j\}$, has a unimodal posterior density. This is shown in the right panel of Figure 1(c); simulation from this density has been achieved by retrospective sampling, see §3·6.

In general, in the conditional method the Markov chain Monte Carlo algorithm has to explore multimodal posterior distributions as those shown in Figure 1. Therefore, we need to add label-switching moves which assist the algorithm to jump across modes. This is particularly important for large datasets, where the modes are separated by areas of negligible probability. A careful inspection of the problem suggests two types of move. The first proposes to change the labels $j$ and $l$ of two randomly chosen components $j, l \in I^{(\mathrm{al})}$. The probability of such a change is accepted with probability $\min\{1, (p_j/p_l)^{m_l - m_j}\}$. This proposal has high acceptance probability if the two components have similar weights; the probability is indeed 1 if $p_j = p_l$. On the other hand, note that the probability of acceptance is small when $|m_l - m_j|$ is close to 0. The second move proposes to change the labels $j$ and $j+1$ of two neighbouring components but at the same time to exchange $V_j$ with $V_{j+1}$. This change is accepted with probability $\min\{1, (1 - V_{j+1})^{m_j}/(1 - V_j)^{m_{j+1}}\}$. This proposal is very effective for swapping the labels of very unequal clusters. For example, it will always accept when $m_j = 0$. On the other hand, if $|p_j - p_{j+1}|$ is small, then the proposal will be rejected with high probability. For example, if $V_j \simeq 1/2$, $V_{j+1} \simeq 1$, and $m_j \simeq m_{j+1}$, then the proposal attempts to allocate $m_j$ is allocated to a component of negligible weight, $\prod_{l=1}^{j-1}(1 - V_l)(1 - V_{j+1})V_j$. Thus, the two moves we have suggested are complementary to each other.

Extensive simulation experimentation has revealed that such moves are crucial. The integrated autocorrelation time, see §4 for definition, of the updated variables can be reduced by as much as 50-90%.

The problem of multimodality has also been addressed in a recent paper by Porteous et al.

17

(2006).

### 3·5 *Exact retrospective sampling of the allocation variables*

It is possible to avoid the Metropolis-Hastings step and simulate directly from the conditional posterior distribution of the allocation variables. This is facilitated by a general retrospective simulation scheme for sampling discrete random variables whose probability mass function depends on the Dirichlet random measure.

We can formulate the problem as one of simulating a discrete random variable $J$, according to the distribution specified by the probabilities

$$r_j := \frac{p_j f_j}{\sum_{j=1}^{\infty} p_j f_j}, \ j = 1, 2, \dots. \tag{9}$$

The $f_j$'s are assumed to be positive and independent random variables, also independent of the $p_j$'s, which are given by the stick-breaking rule in (1), although more general schemes can be incorporated in our method. The retrospective sampling scheme of §2 cannot be applied here, since the normalising constant of (9), $c = \sum p_j f_j$, is unknown. In the specific application we have in mind, $f_j = f(Y_i \mid Z_j, \lambda)$ for each allocation variable $K_i$ we wish to update.

Nevertheless, a retrospective scheme can be applied if we can construct two bounding sequences $\{c_l(k)\} \uparrow c$ and $\{c_u(k)\} \downarrow c$, where $c_l(k)$ and $c_u(k)$ can be calculated simply on the basis of $(p_1, Z_1), \dots, (p_k, Z_k)$. Let $r_{u,j}(k) = r_j/c_l(k)$ and $r_{l,j}(k) = r_j/c_u(k)$. Then we first simulate a uniform $U$ and then we set $J = j$ when

$$\sum_{m=1}^{j-1} r_{u,m}(k) \leq U \leq \sum_{m=1}^{j} r_{l,m}(k). \tag{10}$$

In this algorithm $k$ should be increased, and the $p_j$'s and $f_j$'s simulated retrospectively, until (10) can be verified for some $j \leq k$.

Therefore, simulation from unnormalised discrete probability distributions is feasible provided the bounding sequences can be appropriately constructed. It turns out that this construction is generally very challenging. In this paper we will only tackle the simple case in which there exists a constant $M < \infty$ such that $f_j < M$ for all $j$ almost

surely. In our application this corresponds to $f(y \mid z, \lambda)$ being bounded in $z$. In that case one can simply take

$$
\begin{aligned}
c_l(k) &= \sum_{j=1}^{k} p_j f_j \\
c_u(k) &= c_l(k) + M\left(1 - \sum_{j=1}^{k} p_j\right).
\end{aligned}
$$

When $\limsup f_j = \infty$ almost surely, an alternative construction has to be devised, which involves an appropriate coupling of the $f_j$'s. This construction is elaborate and mathematically intricate, and will be reported elsewhere.

### 3·6    *Posterior inference for functionals of Dirichlet processes*

Suppose we are interested in estimating the posterior distribution of

$$
\mathcal{I} = \int_{\mathcal{X}} g(x) P(dx),
$$

for some real-valued function $g$, where $\mathcal{X}$ denotes the space on which the component parameters $Z_j$ are defined and $P$ is the Dirichlet random measure (2), given data $Y$ from the hierarchical model (1). Our implementation of the conditional method shows that simulation from the posterior distribution of $\mathcal{I}$ is as difficult as the simulation from its prior. Let $\{(V_j, Z_j), j \leq \max\{k\}\}$ be a sample obtained using the retrospective Markov chain Monte Carlo algorithm. We assume that the sample is taken when the chain is 'in stationarity', i.e. after a sufficient number of initial iterations. Given the $V_j$'s we can compute the corresponding $p_j, j \leq \max\{k\}$. Recall that $Z_j \sim H_\Theta$ and $V_j \sim \mathrm{Be}(1, \alpha)$ for all $j > \max\{k\}$. Then we have the following representation for the posterior distribution of $\mathcal{I}$: a draw from $\mathcal{I} \mid Y$ can be represented as

$$
\sum_{j=1}^{\max\{k\}} f(Z_j) p_j + \sum_{j=\max\{k\}+1}^{\infty} f(Z_j) p_j = \sum_{j=1}^{\max\{k\}} f(Z_j) p_j + \frac{I}{\prod_{l=1}^{\max\{k\}} (1 - V_l)}.
$$

This is equality in distribution, where $I$ is a draw from from the prior distribution of $\mathcal{I}$. Inference for linear functionals of the Dirichlet process was initiated in Cifarelli & Regazzini

(1990) and simulation aspects have been considered, for example in Guglielmi & Tweedie (2001) and Guglielmi et al. (2002).

Our algorithm can also be used for the simulation of non-linear functionals under the posterior distribution. As an illustrative example consider posterior simulation of the predominant species corresponding to $J$ such that $p_J \geq p_l$ for all $l = 1, 2, \ldots$. This can be achieved as follows given a sample $\{(V_1, Z_1), j \leq \max\{k\}\}$ obtained with any of our conditional Markov chain Monte Carlo algorithms. Let $J = \arg\max_{\{j \leq N^*\}} p_j$, where $N^* = \max\{k\}$. Then, it can be seen that if $1 - \sum_{j=1}^{N^*} p_j < p_J$ then $J$ is the predominant specie. Therefore, if $1 - \Pi_{N^*} > p_J$, we repeat the following procedure until the condition is satisfied: increase $N^*$, simulate pairs $(Z_{N^*}, V_{N^*})$ from the prior, and compute $J$. This procedure was used to obtain the results in Figure 1(c).

### 3·7  *Inference for hyperparameters*

A simple modification of the retrospective algorithms described above provides the computational machinery needed for Bayesian inference for the hyperparameters $(\Theta, \alpha, \lambda)$. If we assume that appropriate prior distributions have been chosen, the aim is to simulate the hyperparameters according to their full conditional distributions, thus adding one more step in the retrospective Markov chain Monte Carlo algorithm.

Sampling of $\lambda$ according to its full conditional distribution is standard. At first glance, sampling of $(\Theta, \alpha)$ poses a major complication. Note that, because of the hierarchical structure, $(\Theta, \alpha)$ are independent of $(Y, K)$ conditionally upon $(V, Z)$. However, $(V, Z)$ contains an infinite amount of information about $(\alpha, \Theta)$. Therefore, an algorithm which updated successively $(V, Z, K)$ and $(\Theta, \alpha)$ according to their full conditional distributions would be reducible. This type of convergence problem is not uncommon when Markov chain Monte Carlo is used to infer about hierarchical models with hidden stochastic processes; see Papaspiliopoulos et al. (2003, 2006) for reviews.

However, the conditional independence structure in $Z$ and $V$ can be used to circumvent the problem. In effect, instead of updating $(\Theta, \alpha)$ conditionally upon $(Z, V)$

we can jointly update $(\Theta, \alpha, Z^{()}, V^{()})$ conditionally upon $(Z^{(\mathrm{al})}, V^{(\mathrm{al})})$. In practice, we only need to simulate $(\Theta, \alpha)$ conditionally on $(Z^{(\mathrm{al})}, V^{(\mathrm{al})})$. That poses no complication since $(Z^{(\mathrm{al})}, V^{(\mathrm{al})})$ contains only finite amount of information about the hyperparameters. It is worth mentioning that a similar technique for updating the hyperparameters was recommended by MacEachern & Müller (1998) for the implementation of their no-gaps algorithm.

## 4. Comparison between marginal and conditional methods

In this section we attempt a comparison in terms of Monte Carlo efficiency between the retrospective Markov chain Monte Carlo algorithm and state-of-the-art implementations of the marginal approach. To this end we have carried out a large simulation study part of which is presented later in this section. The methods are tested on the non-conjugate model where $f(y \mid z)$ is a Gaussian density with $z = (\mu, \sigma^2)$ and $H_\Theta$ is the product measure $N(\mu, \sigma_z^2) \times \mathrm{IG}(\gamma, \beta)$; there are no further parameters $\lambda$ indexing $f(y \mid z)$, $\Theta = (\mu, \sigma_z^2, \gamma, \beta)$ and $\mathrm{IG}(\gamma, \beta)$ denotes the inverse Gamma distribution with density proportional to $x^{-(\gamma+1)}e^{-\beta x}$. Note that, in a density-estimation context, this model allows for local smoothing; see for example Müller et al. (1996) and Green & Richardson (2001).

An excellent overview of different implementations of the marginal approach can be found in Neal (2000). The most successful implementations for non-conjugate models are the so-called no-gaps algorithm of MacEachern & Müller (1998) and Algorithms 7 and 8 of Neal (2000), which were introduced to mitigate against certain computational inefficiencies of the no-gaps algorithm.

Received Markov chain Monte Carlo  wisdom suggests that marginal samplers ought to be preferred to conditional ones. Some limited theory supports this view; see in particular Liu (1994). However, it is common for marginalisation to destroy conditional independence structure which usually assists the conditional sampler, since conditionally independent components are effectively updated in one block. Thus, it is difficult a priori

to decide which approach is preferable.

In our comparison we have considered different simulated datasets and different prior specifications. We have simulated 4 datasets from two models. The 'lepto 100' and the 'lepto 1000' datasets consist respectively of 100 and 1000 draws from the unimodal leptokurtic mixture, $0.67N(0,1) + 0.33N(0.3, 0.25^2)$. The 'bimod 100' ('bimod 1000') dataset consists of 100 (1000) draws from the bimodal mixture, $0.5N(-1, 0.5^2) + 0.5N(1, 0.5^2)$; we have chosen these datasets following Green & Richardson (2001). In our simulation we have taken the datasets of size 100 to be subsets of those of size 1000. We fix $\Theta$ in a data-driven way as suggested by Green & Richardson (2001): if $R$ denotes the range of the data, then we set $\mu = R/2, \sigma_z = R, \gamma = 2$ and $\beta = 0.02R^2$. Data-dependent choice of hyperparameters is commonly made in mixture models, see for example Richardson & Green (1997). We consider three different values of $\alpha$, $0.2, 1$ and $5$. We use a Gibbs move to update $Z_j = (Z_j^{(1)}, Z_j^{(2)})$ for every $j \in I^{(\mathrm{al})}$: we update $Z_j^{(2)}$ given $Z_j^{(1)}$ and the rest, and then $Z_j^{(1)}$ given the new value of $Z_j^{(2)}$ and the rest. The same scheme is used to update the corresponding cluster parameters in the marginal algorithms.

We monitor the convergence of four functionals of the updated variables: the number of clusters, $M$, the deviance $D$ of the estimated density, and $Z_{K_i}^{(1)}$, for $i = 1, 2$ in 'lepto' and for $i = 2, 3$ in 'bimod'. These functionals have been used in the comparison studies in Neal (2000) and Green & Richardson (2001) to monitor algorithmic performance. In both cases, the components monitored were chosen to be ones whose allocations were particularly well- and badly-identified by the data. The deviance $D$ is calculated as follows

$$D = -2 \sum_{i=1}^{n} \log \Big\{ \sum_{j \in I^{(\mathrm{al})}} \frac{m_j}{n} f(Y_i \mid Z_j) \Big\};$$

see Green & Richardson (2001) for details. Although we have given the expression in terms of the output of the conditional algorithm, a similar expression exists given the output of the marginal algorithms. The deviance is chosen as a meaningful function of several parameters of interest.

22

The efficiency of the sampler is summarised by reporting for each of the monitored variables the estimated integrated autocorrelation time, $\tau = 1 + 2\sum_{j=1}^{\infty} \rho_j$, where $\rho_j$ is the lag-$j$ autocorrelation of the monitored chain. This is a standard way of measuring the speed of convergence of square-integrable functions of an ergodic Markov chain (Roberts, 1996; Sokal, 1997) which has also been used by Neal (2000) and Green & Richardson (2001) in their simulation studies. Recall that the integrated autocorrelation time is proportional to the asymptotic variance of the ergodic average. In particular, if $\tau_1/\tau_2 = b > 1$, where $\tau_i$ is the integrated autocorrelation time of algorithm $i$ for a specific functional, then Algorithm 1 requires roughly $b$ times as many iterations to achieve the same Monte Carlo error as Algorithm 2, for the estimation of the specific functional. Estimation of $\tau$ is a notoriously difficult problem. We have followed the guidelines in §3 of Sokal (1997). We estimate $\tau$ by summing estimated autocorrelations up to a fixed lag $L$, where $\tau << L << N$, and $N$ is the Monte Carlo sample size. Approximate standard errors of the estimate can be obtained; see formula (3.19) of Sokal (1997). For the datasets and prior specifications we have considered we have found that $N = 2 \times 10^6$ suffices in order to assess the relative performance of the competing algorithms.

The results of our comparison are reported in Tables 1 - 3. We contrast our retrospective Markov chain Monte Carlo algorithm with three marginal algorithms: an improved version of the no-gaps algorithm of MacEachern & Müller (1998), where we updated dead-cluster parameters after each allocation variable update, Algorithm 7 of Neal (2000), and Algorithm 8 of Neal (2000), where we use three auxiliary states. The results show that Algorithms 7 and 8 perform better than the competitors, although the difference among the algorithms is moderate. In this comparison we have not taken into account the computing times of the different methods. Our implementation, which however did not aim at optimizing computational time, in FORTRAN 77 suggests that no-gaps, Algorithm 7 and the retrospective Markov chain Monte Carlo algorithm all have roughly similar computing times when $\alpha = 1$. Algorithm 8 is more intensive than Algorithm 7. The computing time of the retrospective algorithm increases with the value

of $\alpha$.

Careful inspection of the output of the algorithms has suggested a possible reason why the conditional approach is outperformed by the marginal approaches. This is because it has to explore multiple modes in the posterior distribution of the random measure $(p, Z)$; see for example Figure 2 for results concerning the 'bimod' dataset. On the other hand, the ambiguity in the cluster labelling is not important in the marginal approaches, which work with the unidentifiable allocation structure and do not need to explore a multimodal distribution. This indicates that the marginal approaches achieve generally smaller integrated autocorrelation times compared to the conditional approach. Nevertheless, the label-switching moves we have included have substantially improved the performance of our algorithm.

## 5. Discussion

The appeal of the conditional approach lies in its potential for inferring for the latent random measure, which we have illustrated, and in its flexibility to be extended to more general stick-breaking random measures than the Dirichlet process. With respect to the latter, we have not explicitly shown how to extend our methods to more general models, but it should be obvious that such extensions are direct. In particular, Proposition 1 will have to be adapted accordingly, but all the crucial conditional independence structure which allows retrospective sampling will be present in the more general contexts.

In extending this work, we have already discussed in §3·3 an exact Gibbs sampler, i.e. one in which the allocation variables are simulated directly from their conditional posterior distributions. If the likelihood function is unbounded, this has to be carried out by an intricate coupling of the Dirichlet process which permits tight bounds on the normalising constants $c_i$ and also allows retrospective simulation of all related variables. Although implementation of the resulting algorithm is simple to implement, the mathematical construction behind this method is very cimplicated and will be reported elsewhere.

24

Retrospective sampling is a methodology with great potential for other problems involving simulation and inference for stochastic processes. One major application which has emerged since the completion of this research, is the exact simulation and estimation of diffusion processes (Beskos et al., 2006a, 2005, 2006b).

## REFERENCES

ANTONIAK, C. E. (1974). Mixtures of Dirichlet processes with applications to bayesian nonparametric problems. *Ann. Statist.* **2**, 1152–74.

BESKOS, A., PAPASPILIOPOULOS, O. & ROBERTS, G. O. (2005). A new factorisation of diffusion measure and finite sample path constructions. To appear in Methodology and Computing in Applied Probability.

BESKOS, A., PAPASPILIOPOULOS, O. & ROBERTS, G. O. (2006a). Retrospective exact simulation of diffusion sample paths with applications. *Bernoulli* **12**, 1077–98.

BESKOS, A., PAPASPILIOPOULOS, O., ROBERTS, G. O. & FEARNHEAD, P. (2006b). Exact and efficient likelihood–based inference for discretely observed diffusions (with Discussion). *J. Roy. Statist. Soc.* **B 68**, 333–82.

BLACKWELL, D. & MACQUEEN, J. B. (1973). Ferguson distributions via Pólya urn schemes. *Ann. Statist.* **1**, 353–5.

BURR, D. & DOSS, H. (2005). A Bayesian semiparametric model for random-effects meta-analysis. *J. Am. Statist. Assoc.* **100**, 242–51.

CIFARELLI, D. M. & REGAZZINI, E. (1990). Distribution functions of means of a Dirichlet process. *Ann. Statist.* **18**, 429–42.

DOSS, H. (1994). Bayesian nonparametric estimation for incomplete data via successive substitution sampling. *Ann. Statist.* **22**, 1763–86.

DUNSON, D. & PARK, J. (2006). Kernel stick-breaking processes. *submitted, available from http://ftp.stat.duke.edu/WorkingPapers/06-22.pdf.*

FEARNHEAD, P. (2004). Particle filters for mixture models with an unknown number of components. *Statist. Comp.* **14**, 11–21.

FERGUSON, T. S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209–30.

FERGUSON, T. S. (1974). Prior distributions on spaces of probability measures. *Ann. Statist.* **2**, 615–29.

FERGUSON, T. S. (1983). Bayesian density estimation by mixtures of normal distributions. In H. Chernoff, M. Haseeb Rizvi, J. Rustagi & D. Siegmund, eds., *Recent Advances in Statistics: Papers in Honor of Herman Chernoff on His Sixtieth Birthday.* New York: Academic Press, pp. 287–302.

GELFAND, A. & KOTTAS, A. (2003). Bayesian semiparametric regression for median residual life. *Scand. J. Statist.* **30**, 651–65.

GELFAND, A. E. & KOTTAS, A. (2002). A computational approach for full nonparametric Bayesian inference under Dirichlet process mixture models. *J. Comp. Graph. Statist.* **11**, 289–305.

GREEN, P. & RICHARDSON, S. (2001). Modelling heterogeneity with and without the Dirichlet process. *Scand. J. Statist.* 28 355–75.

GRIFFIN, J. (2006). On the Bayesian analysis of species sampling mixture models for density estimation. *submitted, available from http://www2.warwick.ac.uk/fac/sci/statistics/staff/academic/griffin/personal/densityestimation.pdf.*

GUGLIELMI, A., HOLMES, C. C. & WALKER, S. G. (2002). Perfect simulation involving functionals of a Dirichlet process. *J. Comp. Graph. Statist.* **11**, 306–10.

GUGLIELMI, A. & TWEEDIE, R. L. (2001). Markov chain Monte Carlo estimation of the law of the mean of a Dirichlet process. *Bernoulli* **7**, 573–92.

ISHWARAN, H. & JAMES, L. (2001). Gibbs sampling methods for stick-breaking priors. *J. Am. Statist. Assoc.* **96**, 161–73.

ISHWARAN, H. & JAMES, L. F. (2003). Some further developments for stick-breaking priors: finite and infinite clustering and classification. *Sankhyā, A,* **65**, 577–92.

ISHWARAN, H. & ZAREPOUR, M. (2000). Markov chain Monte Carlo in approximate Dirichlet and beta two-parameter process hierarchical models. *Biometrika* **87**, 371–90.

ISHWARAN, H. & ZAREPOUR, M. (2002). Exact and approximate sum-representations for the dirichlet process. *Can. J. Statist.* **30**, 269–83.

JAIN, S. & NEAL, R. M. (2004). A split-merge Markov chain Monte Carlo procedure for the Dirichlet process mixture model. *J. Comp. Graph. Statist.* **13**, 158–82.

LIU, J. S. (1994). The collapsed Gibbs sampler in Bayesian computations with applications to a gene regulation problem. *J. Am. Statist. Assoc* **89**, 958–66.

LIU, J. S. (1996). Nonparametric hierarchical Bayes via sequential imputations. *Ann. Statist.* **24**, 911–30.

LO, A. Y. (1984). On a class of Bayesian nonparametric estimates. I. Density estimates. *Ann. Statist.* **12**, 351–7.

MACEACHERN, S. & MÜLLER, P. (1998). Estimating mixture of Dirichlet process models. *J. Comp. Graph. Statist.* **7**, 223–38.

Mengersen, K. L. & Tweedie, R. L. (1996). Rates of convergence of the Hastings and Metropolis algorithms. *Ann. Statist.* 24 101–21.

Müller, P., Erkanli, A. & West, M. (1996). Bayesian curve fitting using multivariate normal mixtures. *Biometrika* **83**, 67–79.

Müller, P., Rosner, G. L., De Iorio, M. & MacEachern, S. (2005). A nonparametric Bayesian model for inference in related longitudinal studies. *Appl. Statist.* **54**, 611–26.

Neal, R. (2000). Markov chain sampling: Methods for Dirichlet process mixture models. *J. Comp. Graph. Statist.* **9**, 283–97.

Papaspiliopoulos, O., Roberts, G. O. & Sköld, M. (2003). Non-centered parameterisations for hierarchical models and data augmentation. In J. Bernardo, M. Bayarri, J. Berger, A. Dawid, D. Heckerman, A. Smith & M. West, eds., *Bayesian Statistics 7*. Oxford: Oxford University Press, pp. 307–27.

Papaspiliopoulos, O., Roberts, G. O. & Sköld, M. (2006). A general framework for parametrisation of hierarchical models. *to appear in Statist. Sci.* .

Porteous, I., Ihter, A., Smyth, P. & Welling, M. (2006). Gibbs sampling for (coupled) infinite mixture models in the stick breaking representation. In *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence (UAI-06)*. Arlington, Virginia: AUAI Press.

Quintana, F. & Iglesias, P. (2003). Bayesian clustering and product partition models. *J. Roy. Statist. Soc. B* **65**, 557–574.

Richardson, S. & Green, P. J. (1997). On Bayesian analysis of mixtures with an unknown number of components (with Discussion). *J. R. Statist. Soc.* **B 59**, 731–92.

Ripley, B. D. (1987). *Stochastic Simulation.* Chichester: Wiley.

ROBERTS, G. O. (1996). Markov chain concepts related to sampling algortihms. In W. Gilks, S. Richardson & D. Spiegelhalter, eds., *MCMC in Practice.* London: Chapman and Hall, pp. 45–57.

SETHURAMAN, J. (1994). A constructive definition of Dirichlet priors. *Statist. Sinica* **4**, 639–50.

SOKAL, A. (1997). Monte carlo methods in statistical mechanics: foundations and new algorithms. In *Functional Integration (Cargèse, 1996)*, vol. **361**, of *NATO Adv. Sci. Inst. Ser. B Phys.* New York: Plenum, pp. 131–92.

TEH, Y., JORDAN, M., BEAL, M. & BLEI, D. (2006). Hierarchical dirichlet processes. *to appear in J. Amer. Statist. Assoc., available from http://www.cs.princeton.edu/ blei/papers/TehJordanBealBlei2006.pdf.*

(a)                                  (b)                                  (c)
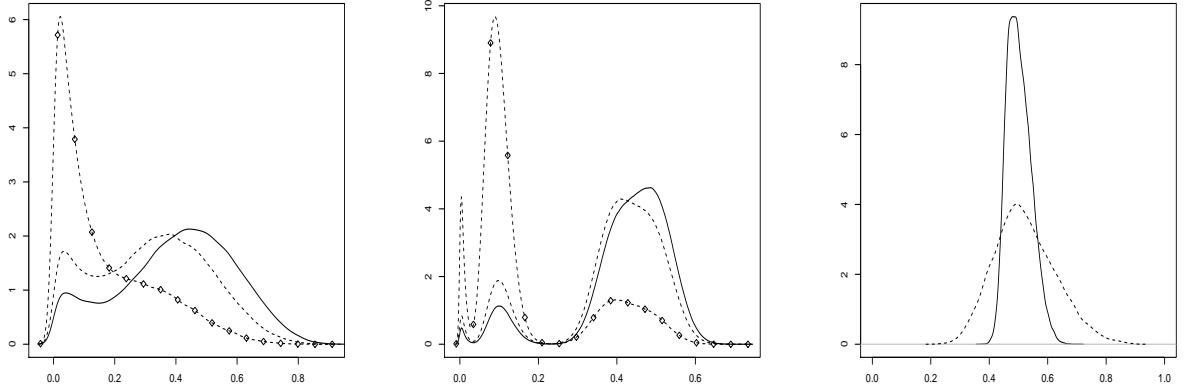


Figure 1: Posterior densities of $p_1$ (solid), $p_2$ (dashed) and $p_3$ (dashed with diamonds), corresponding to (a) a dataset with $n = 10$ separated into three clusters of sizes $n_1 = 5, n_2 = 4, n_1$ and (b) a dataset with $n = 100, n_1 = 50, n_2 = 40$ and $n_3 = 10$. (c) shows the posterior density of $\max_j\{p_j\}$ for $n = 10$ (dashed) and $n = 100$ (solid).

|  | $M$ | $D$ | $Z_{K_3}$ | $Z_{K_2}$ |
|---|---|---|---|---|
| Method | | $\alpha = 1$ | | |
| Retrospective | 41.42 (2.6) | 3.28 (0.21) | 3.9 (0.25) | 3.44 (0.22) |
| No-gaps | 45.94 (1.52) | 3.84 (0.13) | 3.66 (0.12) | 2.82 (0.09) |
| Algorithm 7 | 21.85 (0.77) | 2.48 (0.09) | 3.47 (0.12) | 2.91 (0.10) |
| Algorithm 8 | 18.21 (0.66) | 2.94 (0.11) | 3.44 (0.13) | 3.1 (0.12) |
| | | $\alpha = 0.2$ | | |
| Retrospective | 67.0 (4.24) | 6.8 (0.43) | 8.01 (0.51) | 3.44 (0.22) |
| No-gaps | 39.44 (1.31) | 3.8 (0.13) | 5.93 (0.2) | 3.1 (0.10) |
| Algorithm 7 | 24.99 (0.88) | 2.87 (0.10) | 6.32 (0.22) | 2.95 (0.10) |
| Algorithm 8 | 22.10 (0.85) | 5.30 (0.20) | 6.8 (0.26) | 3.03 (0.12) |
| | | $\alpha = 5$ | | |
| Retrospective | 21.86 (1.38) | 2.82 (0.18) | 2.01 (0.13) | 2.5 (0.16) |
| No-gaps | 57.09 (1.81) | 2.99 (0.09) | 1.67 (0.05) | 2.01 (0.06) |
| Algorithm 7 | 12.55 (0.4) | 1.77 (0.06) | 1.64 (0.05) | 1.97 (0.06) |
| Algorithm 8 | 8.2 (0.26) | 1.77 (0.06) | 1.64 (0.05) | 1.97 (0.06) |

Table 1: Estimated integrated autocorrelation times for the number of clusters $M$, the deviance $D$, $Z_{K_3}$ and $Z_{K_2}$, for the 'bimod 100' dataset. Estimates of the standard error in parenthesis. The initial state of all chains was all data allocated to the same cluster with parameters drawn from the prior.

| Method | $M$ | $D$ | $Z_{K_1}$ | $Z_{K_2}$ |
|---|---|---|---|---|
| | | $\alpha = 1$ | | |
| Retrospective | 40.71 (2.58) | 31.99 (2.01) | 46.58 (2.95) | 3.04 (0.19) |
| No-gaps | 46.08 (1.46) | 23.93 (0.76) | 33.19 (1.05) | 2.37 (0.07) |
| Algorithm 7 | 22.98 (0.73) | 20.17 (0.64) | 28.28 (0.89) | 2.33 (0.07) |
| Algorithm 8 | 18.02 (0.57) | 18.91 (0.6) | 26.71 (0.85) | 2.06 (0.07) |
| | | $\alpha = 0.2$ | | |
| Retrospective | 239.07 (15.12) | 286.49 (18.12) | 157.85 (9.99) | 14.87 (0.94) |
| No-gaps | 127.08 (6.96) | 151.90 (8.32) | 97.73 (5.35) | 7.46 (0.41) |
| Algorithm 7 | 109.37 (5.99) | 171.98 (9.42) | 86.26 (4.73) | 7.95 (0.44) |
| Algorithm 8 | 99.06 (5.43) | 142.93 (7.83) | 82.38 (4.51) | 6.98 (0.38) |
| | | $\alpha = 5$ | | |
| Retrospective | 13.69 (0.87) | 7.38 (0.47) | 5.9 (0.37) | 1.61 (0.1) |
| No-gaps | 44.25 (1.4) | 5.72 (0.18) | 4.14 (0.13) | 1.36 (0.04) |
| Algorithm 7 | 10.57 (0.33) | 5.55 (0.18) | 3.52 (0.11) | 1.33 (0.04) |
| Algorithm 8 | 6.32 (0.2) | 5.31 (0.17) | 3.23 (0.10) | 1.29 (0.04) |

Table 2: Estimated integrated autocorrelation times for the number of clusters $M$, the deviance $D$, $Z_{K_1}$ and $Z_{K_2}$, for the 'lepto 100' dataset. Estimates of the standard errors in parenthesis. The initial state of all chains was all data allocated to the same cluster with parameters drawn from the prior.

| | 'bimod 1000', $\alpha = 1$ | 'lepto 1000', $\alpha = 1$ | |
|---|---|---|---|
| | $M$ | $D$ | $M$ |
| Retrospective | 149 (7) | 254 (25) | 205 (21) |
| No-gaps | 91 (4) | 133 (6) | 102 (5) |
| Algorithm 7 | 60 (3) | 87 (4) | 99 (4) |
| Algorithm 8 | 58 (3) | 112 (5) | 104 (5) |

Table 3: Estimated integrated autocorrelation times for the 'bimod 1000' and 'lepto 1000' datasets. Ther results for $D$ in the 'bimod 1000' data set, $Z_{K_1}, Z_{K_2}$ and $Z_{K_3}$, for both datasets were not markedly different across the algorithms so are omitted.
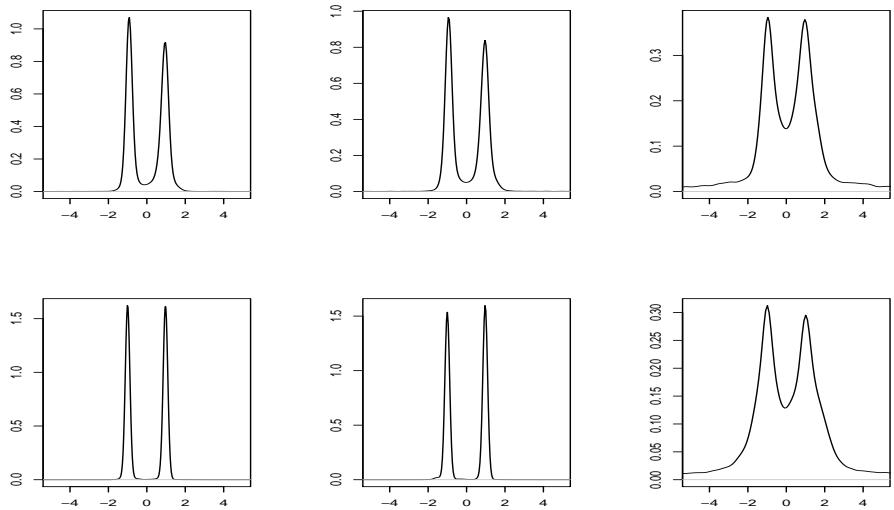
Figure 2: Posterior densities of $Z_1$ (left), $Z_2$ (middle) and $Z_3$ (right) for the 'bimod 100' (top) and the 'bimod 1000' (bottom) datasets. All results have been obtained for $\alpha = 1$.