

# Hierarchical Dirichlet Processes

Yee Whye Teh `tehyw@comp.nus.edu.sg`  
Department of Computer Science, National University of Singapore,  
Singapore 117543

Michael I. Jordan `jordan@eecs.berkeley.edu`  
Computer Science Division and Department of Statistics,  
University of California at Berkeley, Berkeley CA 94720-1776, USA

Matthew J. Beal `mbeal@cse.buffalo.edu`  
Department of Computer Science & Engineering,  
State University of New York at Buffalo, Buffalo NY 14260-2000, USA

David M. Blei `blei@eecs.berkeley.edu`  
Department of Computer Science, Princeton University,  
Princeton, NJ 08544, USA

November 15, 2005

## Abstract

We consider problems involving groups of data, where each observation within a group is a draw from a mixture model, and where it is desirable to share mixture components between groups. We assume that the number of mixture components is unknown a priori and is to be inferred from the data. In this setting it is natural to consider sets of Dirichlet processes, one for each group, where the well-known clustering property of the Dirichlet process provides a nonparametric prior for the number of mixture components within each group. Given our desire to tie the mixture models in the various groups, we consider a hierarchical model, specifically one in which the base measure for the child Dirichlet processes is itself distributed according to a Dirichlet process. Such a base measure being discrete, the child Dirichlet processes necessarily share atoms. Thus, as desired, the mixture models in the different groups necessarily share mixture components. We discuss representations of hierarchical Dirichlet processes in terms of a stick-breaking process, and a generalization of the Chinese restaurant process that we refer to as the “Chinese restaurant franchise.” We present Markov chain Monte Carlo algorithms for posterior inference in hierarchical Dirichlet process mixtures, and describe applications to problems in information retrieval and text modelling.

**Keywords:** clustering, mixture models, nonparametric Bayesian statistics, hierarchical models, Markov chain Monte Carlo

# 1 INTRODUCTION

A recurring theme in statistics is the need to separate observations into groups, and yet allow the groups to remain linked—to “share statistical strength.” In the Bayesian formalism such sharing is achieved naturally via hierarchical modeling; parameters are shared among groups, and the randomness of the parameters induces dependencies among the groups. Estimates based on the posterior distribution exhibit “shrinkage.”

In the current paper we explore a hierarchical approach to the problem of model-based clustering of grouped data. We assume that the data are subdivided into a set of groups, and that within each group we wish to find clusters that capture latent structure in the data assigned to that group. The number of clusters within each group is unknown and is to be inferred. Moreover, in a sense that we make precise, we wish to allow clusters to be shared among the groups.

An example of the kind of problem that motivates us can be found in genetics. Consider a set of  $k$  binary markers (e.g., single nucleotide polymorphisms or “SNPs”) in a localized region of the human genome. While an individual human could exhibit any of  $2^k$  different patterns of markers on a single chromosome, in real populations only a small subset of such patterns—*haplotypes*—are actually observed (Gabriel et al. 2002). Given a meiotic model for the combination of a pair of haplotypes into a *genotype* during mating, and given a set of observed genotypes in a sample from a human population, it is of great interest to identify the underlying haplotypes (Stephens et al. 2001). Now consider an extension of this problem in which the population is divided into a set of groups; e.g., African, Asian and European subpopulations. We may not only want to discover the sets of haplotypes within each subpopulation, but we may also wish to discover which haplotypes are shared between subpopulations. The identification of such haplotypes would have significant implications for the understanding of the migration patterns of ancestral populations of humans.

As a second example, consider the problem from the field of information retrieval (IR) of modeling of relationships among sets of documents. In IR, documents are generally modeled under an exchangeability assumption, the “bag of words” assumption, in which the order of words in a document is ignored (Salton and McGill 1983). It is also common to view the words in a document as arising from a number of latent clusters or “topics,” where a topic is generally modeled as a multinomial probability distribution on words from some basic vocabulary (Blei et al. 2003). Thus, in a document concerned with university funding the words in the document might be drawn from the topics “education” and “finance.” Considering a collection of such documents, we may wish to allow topics to be shared among the documents in the corpus. For example, if the corpus also contains a document concerned with university football, the topics may be “education” and “sports,” and we would want the former topic to be related to that discovered in the analysis of the document on university funding.

Moreover, we may want to extend the model to allow for multiple corpora. For example, documents in scientific journals are often grouped into themes (e.g., “empirical process theory,” “multivariate statistics,” “survival analysis”), and it would be of interest to discover to what extent the latent topics that are shared among documents are also shared across these groupings. Thus in general we wish to consider the sharing of clusters across multiple, nested groupings of data.

Our approach to the problem of sharing clusters among multiple, related groups is a nonparametric Bayesian approach, reposing on the *Dirichlet process* (Ferguson 1973). The Dirichlet process  $DP(\alpha_0, G_0)$  is a measure on measures. It has two parameters, a *scaling parameter*  $\alpha_0 > 0$  and a *base probability measure*  $G_0$ . An explicit representation of a draw from a Dirichlet process (DP)

was given by Sethuraman (1994), who showed that if  $G \sim \text{DP}(\alpha_0, G_0)$ , then with probability one:

$$G = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}, \quad (1)$$

where the  $\phi_k$  are independent random variables distributed according to  $G_0$ , where  $\delta_{\phi_k}$  is an atom at  $\phi_k$ , and where the “stick-breaking weights”  $\beta_k$  are also random and depend on the parameter  $\alpha_0$  (the definition of the  $\beta_k$  is provided in Section 3.1).

The representation in (1) shows that draws from a DP are discrete (with probability one). The discrete nature of the DP makes it unsuitable for general applications in Bayesian nonparametrics, but it is well suited for the problem of placing priors on mixture components in mixture modeling. The idea is basically to associate a mixture component with each atom in  $G$ . Introducing indicator variables to associate data points with mixture components, the posterior distribution yields a probability distribution on partitions of the data. A number of authors have studied such *Dirichlet process mixture models* (Antoniak 1974; Escobar and West 1995; MacEachern and Müller 1998). These models provide an alternative to methods that attempt to select a particular number of mixture components, or methods that place an explicit parametric prior on the number of components.

Let us now consider the setting in which the data are subdivided into a number of groups. Given our goal of solving a clustering problem within each group, we consider a set of random measures  $G_j$ , one for each group  $j$ , where  $G_j$  is distributed according to a group-specific Dirichlet process  $\text{DP}(\alpha_{0j}, G_{0j})$ . To link these clustering problems, we link the group-specific DPs. Many authors have considered ways to induce dependencies among multiple DPs via links among the parameters  $G_{0j}$  and/or  $\alpha_{0j}$  (Cifarelli and Regazzini 1978; MacEachern 1999; Tomlinson 1998; Müller et al. 2004; De Iorio et al. 2004; Kleinman and Ibrahim 1998; Mallick and Walker 1997; Ishwaran and James 2004). Focusing on the  $G_{0j}$ , one natural proposal is a hierarchy in which the measures  $G_j$  are conditionally independent draws from a single underlying Dirichlet process  $\text{DP}(\alpha_0, G_0(\tau))$ , where  $G_0(\tau)$  is a parametric distribution with random parameter  $\tau$  (Carota and Parmigiani 2002; Fong et al. 2002; Muliere and Petrone 1993). Integrating over  $\tau$  induces dependencies among the DPs.

That this simple hierarchical approach will not solve our problem can be observed by considering the case in which  $G_0(\tau)$  is absolutely continuous with respect to Lebesgue measure for almost all  $\tau$  (e.g.,  $G_0$  is Gaussian with mean  $\tau$ ). In this case, given that the draws  $G_j$  arise as conditionally independent draws from  $G_0(\tau)$ , they necessarily have no atoms in common (with probability one). Thus, although clusters arise *within* each group via the discreteness of draws from a DP, the atoms associated with the different groups are different and there is no sharing of clusters *between* groups. This problem can be skirted by assuming that  $G_0$  lies in a discrete parametric family, but such an assumption would be overly restrictive.

Our proposed solution to the problem is straightforward: to force  $G_0$  to be discrete and yet have broad support, we consider a nonparametric hierarchical model in which  $G_0$  is itself a draw from a Dirichlet process  $\text{DP}(\gamma, H)$ . This restores flexibility in that the modeler can choose  $H$  to be continuous or discrete. In either case, with probability one,  $G_0$  is discrete and has a stick-breaking representation as in (1). The atoms  $\phi_k$  are shared among the multiple DPs, yielding the desired sharing of atoms among groups. In summary, we consider the hierarchical specification:

$$\begin{aligned} G_0 \mid \gamma, H &\sim \text{DP}(\gamma, H) \\ G_j \mid \alpha_0, G_0 &\sim \text{DP}(\alpha_0, G_0) \quad \text{for each } j, \end{aligned} \quad (2)$$

which we refer to as a *hierarchical Dirichlet process*. The immediate extension to *hierarchical Dirichlet process mixture models* yields our proposed formalism for sharing clusters among related clustering problems.

Related nonparametric approaches to linking multiple DPs have been discussed by a number of authors. Our approach is a special case of a general framework for “dependent Dirichlet processes” due to MacEachern (1999) and MacEachern et al. (2001). In this framework the random variables  $\beta_k$  and  $\phi_k$  in (1) are general stochastic processes (i.e., indexed collections of random variables); this allows very general forms of dependency among DPs. Our hierarchical approach fits into this framework; we endow the stick-breaking weights  $\beta_k$  in (1) with a second subscript indexing the groups  $j$ , and view the weights  $\beta_{jk}$  as dependent for each fixed value of  $k$ . Indeed, as we show in Section 4, the definition in (2) yields a specific, canonical form of dependence among the weights  $\beta_{jk}$ .

Our approach is also a special case of a framework referred to as *analysis of densities* (AnDe) by Tomlinson (1998) and Tomlinson and Escobar (2003). The AnDe model is a hierarchical model for multiple DPs in which the common base measure  $G_0$  is random, but rather than treating  $G_0$  as a draw from a DP, as in our case, it is treated as a draw from a mixture of DPs. The resulting  $G_0$  is continuous in general (Antoniak 1974), which, as we have discussed, is ruinous for our problem of sharing clusters. It is an appropriate choice, however, for the problem addressed by Tomlinson (1998), which is that of sharing statistical strength among multiple sets of density estimation problems. Thus, while the AnDe framework and our hierarchical DP framework are closely related formally, the inferential goal is rather different. Moreover, as we will see, our restriction to discrete  $G_0$  has important implications for the design of efficient MCMC inference algorithms.

The terminology of “hierarchical Dirichlet process” has also been used by Müller et al. (2004) to describe a different notion of hierarchy than the one discussed here. These authors consider a model in which a coupled set of random measures  $G_j$  are defined as  $G_j = \epsilon F_0 + (1 - \epsilon)F_j$ , where  $F_0$  and the  $F_j$  are draws from DPs. This model provides an alternative approach to sharing clusters, one in which the shared clusters are given the same stick-breaking weights (those associated with  $F_0$ ) in each of the groups. By contrast, in our hierarchical model, the draws  $G_j$  are based on the same underlying base measure  $G_0$ , but each draw assigns different stick-breaking weights to the shared atoms associated with  $G_0$ . Thus, atoms can be partially shared.

Finally, the terminology of “hierarchical Dirichlet process” has been used in yet a third way by Beal et al. (2002) in the context of a model known as the *infinite hidden Markov model*, a hidden Markov model with a countably infinite state space. The “hierarchical Dirichlet process” of Beal et al. (2002) is, however, not a hierarchy in the Bayesian sense; rather, it is an algorithmic description of a coupled set of urn models. We discuss this model in more detail in Section 7, where we show that the notion of hierarchical DP presented here yields an elegant treatment of the infinite hidden Markov model.

In summary, the notion of hierarchical Dirichlet process that we explore is a specific example of a dependency model for multiple Dirichlet processes, one specifically aimed at the problem of sharing clusters among related groups of data. It involves a simple Bayesian hierarchy where the base measure for a set of Dirichlet processes is itself distributed according to a Dirichlet process. While there are many ways to couple Dirichlet processes, we view this simple, canonical Bayesian hierarchy as particularly worthy of study. Note in particular the appealing recursiveness of the definition; a hierarchical Dirichlet process can be readily extended to multiple hierarchical levels. This is natural in applications. For example, in our application to document modeling, one level of hierarchy is needed to share clusters among multiple documents within a corpus, and second level of hierarchy is needed to share clusters among multiple corpora. Similarly, in the genetics example, it is of interest to consider nested subdivisions of populations according to various criteria (geographic, cultural, economic), and to consider the flow of haplotypes on the resulting tree.

As is the case with other nonparametric Bayesian methods, a significant component of the chal-

length in working with the hierarchical Dirichlet process is computational. To provide a general framework for designing procedures for posterior inference in the hierarchical Dirichlet process that parallel those available for the Dirichlet process, it is necessary to develop analogs for the hierarchical Dirichlet process of some of the representations that have proved useful in the Dirichlet process setting. We provide these analogs in Section 4 where we discuss a stick-breaking representation of the hierarchical Dirichlet process, an analog of the Pólya urn model that we refer to as the “Chinese restaurant franchise,” and a representation of the hierarchical Dirichlet process in terms of an infinite limit of finite mixture models. With these representations as background, we present MCMC algorithms for posterior inference under hierarchical Dirichlet process mixtures in Section 5. We present experimental results in Section 6 and present our conclusions in Section 8.

## 2 SETTING

We are interested in problems where the observations are organized into *groups*, and assumed exchangeable both within each group and across groups. To be precise, letting  $j$  index the groups and  $i$  index the observations within each group, we assume that  $x_{j1}, x_{j2}, \dots$  are exchangeable within each group  $j$ . We also assume that the observations are exchangeable at the group level, that is, if  $\mathbf{x}_j = (x_{j1}, x_{j2}, \dots)$  denote all observations in group  $j$ , then  $\mathbf{x}_1, \mathbf{x}_2, \dots$  are exchangeable.

Assuming each observation is drawn independently from a mixture model, there is a mixture component associated with each observation. Let  $\theta_{ji}$  denote a parameter specifying the mixture component associated with the observation  $x_{ji}$ . We will refer to the variables  $\theta_{ji}$  as *factors*. Note that these variables are not generally distinct; we will develop a different notation for the distinct values of factors. Let  $F(\theta_{ji})$  denote the distribution of  $x_{ji}$  given the factor  $\theta_{ji}$ . Let  $G_j$  denote a prior distribution for the factors  $\boldsymbol{\theta}_j = (\theta_{j1}, \theta_{j2}, \dots)$  associated with group  $j$ . We assume that the factors are conditionally independent given  $G_j$ . Thus we have the following probability model:

$$\begin{aligned} \theta_{ji} \mid G_j &\sim G_j && \text{for each } j \text{ and } i, \\ x_{ji} \mid \theta_{ji} &\sim F(\theta_{ji}) && \text{for each } j \text{ and } i, \end{aligned} \tag{3}$$

to augment the specification given in (2).

## 3 DIRICHLET PROCESSES

In this section, we provide a brief overview of Dirichlet processes. After a discussion of basic definitions, we present three different perspectives on the Dirichlet process: one based on the stick-breaking construction, one based on a Pólya urn model, and one based on a limit of finite mixture models. Each of these perspectives has an analog in the hierarchical Dirichlet process, which is described in Section 4.

Let  $(\Theta, \mathcal{B})$  be a measurable space, with  $G_0$  a probability measure on the space. Let  $\alpha_0$  be a positive real number. A *Dirichlet process*  $DP(\alpha_0, G_0)$  is defined to be the distribution of a random probability measure  $G$  over  $(\Theta, \mathcal{B})$  such that, for any finite measurable partition  $(A_1, A_2, \dots, A_r)$  of  $\Theta$ , the random vector  $(G(A_1), \dots, G(A_r))$  is distributed as a finite-dimensional Dirichlet distribution with parameters  $(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r))$ :

$$(G(A_1), \dots, G(A_r)) \sim \text{Dir}(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r)). \tag{4}$$

We write  $G \sim DP(\alpha_0, G_0)$  if  $G$  is a random probability measure with distribution given by the Dirichlet process. The existence of the Dirichlet process was established by Ferguson (1973).

### 3.1 The stick-breaking construction

Measures drawn from a Dirichlet process are discrete with probability one (Ferguson 1973). This property is made explicit in the *stick-breaking construction* due to Sethuraman (1994). The stick-breaking construction is based on independent sequences of i.i.d. random variables  $(\pi'_k)_{k=1}^\infty$  and  $(\phi_k)_{k=1}^\infty$ :

$$\pi'_k \mid \alpha_0, G_0 \sim \text{Beta}(1, \alpha_0) \quad \phi_k \mid \alpha_0, G_0 \sim G_0 . \quad (5)$$

Now define a random measure  $G$  as

$$\pi_k = \pi'_k \prod_{l=1}^{k-1} (1 - \pi'_l) \quad G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k} , \quad (6)$$

where  $\delta_\phi$  is a probability measure concentrated at  $\phi$ . Sethuraman (1994) showed that  $G$  as defined in this way is a random probability measure distributed according to  $\text{DP}(\alpha_0, G_0)$ .

It is important to note that the sequence  $\pi = (\pi_k)_{k=1}^\infty$  constructed by (5) and (6) satisfies  $\sum_{k=1}^\infty \pi_k = 1$  with probability one. Thus we may interpret  $\pi$  as a random probability measure on the positive integers. For convenience, we shall write  $\pi \sim \text{GEM}(\alpha_0)$  if  $\pi$  is a random probability measure defined by (5) and (6) (GEM stands for Griffiths, Engen and McCloskey; e.g. see Pitman 2002b).

### 3.2 The Chinese restaurant process

A second perspective on the Dirichlet process is provided by the *Pólya urn scheme* (Blackwell and MacQueen 1973). The Pólya urn scheme shows that draws from the Dirichlet process are both discrete and exhibit a clustering property.

The Pólya urn scheme does not refer to  $G$  directly; it refers to draws from  $G$ . Thus, let  $\theta_1, \theta_2, \dots$  be a sequence of i.i.d. random variables distributed according to  $G$ . That is, the variables  $\theta_1, \theta_2, \dots$  are conditionally independent given  $G$ , and hence exchangeable. Let us consider the successive conditional distributions of  $\theta_i$  given  $\theta_1, \dots, \theta_{i-1}$ , where  $G$  has been integrated out. Blackwell and MacQueen (1973) showed that these conditional distributions have the following form:

$$\theta_i \mid \theta_1, \dots, \theta_{i-1}, \alpha_0, G_0 \sim \sum_{\ell=1}^{i-1} \frac{1}{i-1 + \alpha_0} \delta_{\theta_\ell} + \frac{\alpha_0}{i-1 + \alpha_0} G_0 . \quad (7)$$

We can interpret the conditional distributions in terms of a simple urn model in which a ball of a distinct color is associated with each atom. The balls are drawn equiprobably; when a ball is drawn it is placed back in the urn together with another ball of the same color. In addition, with probability proportional to  $\alpha_0$  a new atom is created by drawing from  $G_0$  and a ball of a new color is added to the urn.

Expression (7) shows that  $\theta_i$  has positive probability of being equal to one of the previous draws. Moreover, there is a positive reinforcement effect; the more often a point is drawn, the more likely it is to be drawn in the future. To make the clustering property explicit, it is helpful to introduce a new set of variables that represent distinct values of the atoms. Define  $\phi_1, \dots, \phi_K$  to be the distinct values taken on by  $\theta_1, \dots, \theta_{i-1}$ , and let  $m_k$  be the number of values  $\theta_{i'}$  that are equal to  $\phi_k$  for  $1 \leq i' < i$ . We can re-express (7) as

$$\theta_i \mid \theta_1, \dots, \theta_{i-1}, \alpha_0, G_0 \sim \sum_{k=1}^K \frac{m_k}{i-1 + \alpha_0} \delta_{\phi_k} + \frac{\alpha_0}{i-1 + \alpha_0} G_0 . \quad (8)$$

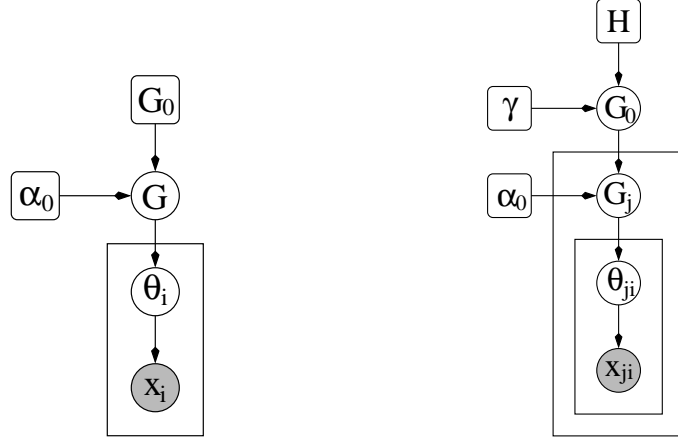


Figure 1: (Left) A representation of a Dirichlet process mixture model as a graphical model. (Right) A hierarchical Dirichlet process mixture model. In the graphical model formalism, each node in the graph is associated with a random variable, where shading denotes an observed variable. Rectangles denote replication of the model within the rectangle. Sometimes the number of replicates is given in the bottom right corner of the rectangle.

Using a somewhat different metaphor, the Pólya urn scheme is closely related to a distribution on partitions known as the *Chinese restaurant process* (Aldous 1985). This metaphor has turned out to be useful in considering various generalizations of the Dirichlet process (Pitman 2002a), and it will be useful in this paper. The metaphor is as follows. Consider a Chinese restaurant with an unbounded number of tables. Each  $\theta_i$  corresponds to a customer who enters the restaurant, while the distinct values  $\phi_k$  correspond to the tables at which the customers sit. The  $i^{\text{th}}$  customer sits at the table indexed by  $\phi_k$ , with probability proportional to the number of customers  $m_k$  already seated there (in which case we set  $\theta_i = \phi_k$ ), and sits at a new table with probability proportional to  $\alpha_0$  (increment  $K$ , draw  $\phi_K \sim G_0$  and set  $\theta_i = \phi_K$ ).

### 3.3 Dirichlet process mixture models

One of the most important applications of the Dirichlet process is as a nonparametric prior on the parameters of a mixture model. In particular, suppose that observations  $x_i$  arise as follows:

$$\begin{aligned} \theta_i &| G \sim G \\ x_i &| \theta_i \sim F(\theta_i), \end{aligned} \tag{9}$$

where  $F(\theta_i)$  denotes the distribution of the observation  $x_i$  given  $\theta_i$ . The factors  $\theta_i$  are conditionally independent given  $G$ , and the observation  $x_i$  is conditionally independent of the other observations given the factor  $\theta_i$ . When  $G$  is distributed according to a Dirichlet process, this model is referred to as a *Dirichlet process mixture model*. A graphical model representation of a Dirichlet process mixture model is shown in Figure 1 (Left).

Since  $G$  can be represented using a stick-breaking construction (6), the factors  $\theta_i$  take on values  $\phi_k$  with probability  $\pi_k$ . We may denote this using an indicator variable  $z_i$  which takes on positive integral values and is distributed according to  $\pi$  (interpreting  $\pi$  as a random probability measure on

the positive integers). Hence an equivalent representation of a Dirichlet process mixture is given by the following conditional distributions:

$$\begin{aligned} \boldsymbol{\pi} \mid \alpha_0 &\sim \text{GEM}(\alpha_0) & z_i \mid \boldsymbol{\pi} &\sim \boldsymbol{\pi} \\ \phi_k \mid G_0 &\sim G_0 & x_i \mid z_i, (\phi_k)_{k=1}^\infty &\sim F(\phi_{z_i}). \end{aligned} \quad (10)$$

Moreover,  $G = \sum_{k=1}^\infty \pi_k \delta_{\phi_k}$  and  $\theta_i = \phi_{z_i}$ .

### 3.4 The infinite limit of finite mixture models

A Dirichlet process mixture model can be derived as the limit of a sequence of finite mixture models, where the number of mixture components is taken to infinity (Neal 1992; Rasmussen 2000; Green and Richardson 2001; Ishwaran and Zarepour 2002). This limiting process provides a third perspective on the Dirichlet process.

Suppose we have  $L$  mixture components. Let  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_L)$  denote the mixing proportions. Note that we previously used the symbol  $\boldsymbol{\pi}$  to denote the weights associated with the atoms in  $G$ . We have deliberately overloaded the definition of  $\boldsymbol{\pi}$  here; as we shall see later, they are closely related. In fact, in the limit  $L \rightarrow \infty$  these vectors are equivalent up to a random *size-biased permutation* of their entries (Pitman 1996).

We place a Dirichlet prior on  $\boldsymbol{\pi}$  with symmetric parameters  $(\alpha_0/L, \dots, \alpha_0/L)$ . Let  $\phi_k$  denote the parameter vector associated with mixture component  $k$ , and let  $\phi_k$  have prior distribution  $G_0$ . Drawing an observation  $x_i$  from the mixture model involves picking a specific mixture component with probability given by the mixing proportions; let  $z_i$  denote that component. We thus have the following model:

$$\begin{aligned} \boldsymbol{\pi} \mid \alpha_0 &\sim \text{Dir}(\alpha_0/L, \dots, \alpha_0/L) & z_i \mid \boldsymbol{\pi} &\sim \boldsymbol{\pi} \\ \phi_k \mid G_0 &\sim G_0 & x_i \mid z_i, (\phi_k)_{k=1}^L &\sim F(\phi_{z_i}). \end{aligned} \quad (11)$$

Let  $G^L = \sum_{k=1}^L \pi_k \delta_{\phi_k}$ . Ishwaran and Zarepour (2002) show that for every measurable function  $f$  integrable with respect to  $G_0$ , we have, as  $L \rightarrow \infty$ :

$$\int f(\theta) dG^L(\theta) \xrightarrow{\mathcal{D}} \int f(\theta) dG(\theta). \quad (12)$$

A consequence of this is that the marginal distribution induced on the observations  $x_1, \dots, x_n$  approaches that of a Dirichlet process mixture model.

## 4 HIERARCHICAL DIRICHLET PROCESSES

We propose a nonparametric Bayesian approach to the modeling of grouped data, where each group is associated with a mixture model, and where we wish to link these mixture models. By analogy with Dirichlet process mixture models, we first define the appropriate nonparametric prior, which we refer to as the *hierarchical Dirichlet process*. We then show how this prior can be used in the grouped mixture model setting. We present analogs of the three perspectives presented earlier for the Dirichlet process—a stick-breaking construction, a Chinese restaurant process representation, and a representation in terms of a limit of finite mixture models.

A hierarchical Dirichlet process is a distribution over a set of random probability measures over  $(\Theta, \mathcal{B})$ . The process defines a set of random probability measures  $G_j$ , one for each group, and a

global random probability measure  $G_0$ . The global measure  $G_0$  is distributed as a Dirichlet process with concentration parameter  $\gamma$  and base probability measure  $H$ :

$$G_0 \mid \gamma, H \sim \text{DP}(\gamma, H), \quad (13)$$

and the random measures  $G_j$  are conditionally independent given  $G_0$ , with distributions given by a Dirichlet process with base probability measure  $G_0$ :

$$G_j \mid \alpha_0, G_0 \sim \text{DP}(\alpha_0, G_0). \quad (14)$$

The hyperparameters of the hierarchical Dirichlet process consist of the baseline probability measure  $H$ , and the concentration parameters  $\gamma$  and  $\alpha_0$ . The baseline  $H$  provides the prior distribution for the factors  $\theta_{ji}$ . The distribution  $G_0$  varies around the prior  $H$ , with the amount of variability governed by  $\gamma$ . The actual distribution  $G_j$  over the factors in the  $j^{\text{th}}$  group deviates from  $G_0$ , with the amount of variability governed by  $\alpha_0$ . If we expect the variability in different groups to be different, we can use a separate concentration parameter  $\alpha_j$  for each group  $j$ . In this paper, following Escobar and West (1995), we put vague gamma priors on  $\gamma$  and  $\alpha_0$ .

A hierarchical Dirichlet process can be used as the prior distribution over the factors for grouped data. For each  $j$  let  $\theta_{j1}, \theta_{j2}, \dots$  be i.i.d. random variables distributed as  $G_j$ . Each  $\theta_{ji}$  is a factor corresponding to a single observation  $x_{ji}$ . The likelihood is given by:

$$\begin{aligned} \theta_{ji} \mid G_j &\sim G_j \\ x_{ji} \mid \theta_{ji} &\sim F(\theta_{ji}). \end{aligned} \quad (15)$$

This completes the definition of a *hierarchical Dirichlet process mixture model*. The corresponding graphical model is shown in Figure 1 (Right).

The hierarchical Dirichlet process can readily be extended to more than two levels. That is, the base measure  $H$  can itself be a draw from a DP, and the hierarchy can be extended for as many levels as are deemed useful. In general, we obtain a tree in which a DP is associated with each node, in which the children of a given node are conditionally independent given their parent, and in which the draw from the DP at a given node serves as a base measure for its children. The atoms in the stick-breaking representation at a given node are thus shared among all descendant nodes, providing a notion of shared clusters at multiple levels of resolution.

#### 4.1 The stick-breaking construction

Given that the global measure  $G_0$  is distributed as a Dirichlet process, it can be expressed using a stick-breaking representation:

$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}, \quad (16)$$

where  $\phi_k \sim H$  independently and  $\beta = (\beta_k)_{k=1}^{\infty} \sim \text{GEM}(\gamma)$  are mutually independent. Since  $G_0$  has support at the points  $\phi = (\phi_k)_{k=1}^{\infty}$ , each  $G_j$  necessarily has support at these points as well, and can thus be written as:

$$G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\phi_k}. \quad (17)$$

Let  $\pi_j = (\pi_{jk})_{k=1}^\infty$ . Note that the weights  $\pi_j$  are independent given  $\beta$  (since the  $G_j$  are independent given  $G_0$ ). We now describe how the weights  $\pi_j$  are related to the global weights  $\beta$ .

Let  $(A_1, \dots, A_r)$  be a measurable partition of  $\Theta$  and let  $K_l = \{k : \phi_k \in A_l\}$  for  $l = 1, \dots, r$ . Note that  $(K_1, \dots, K_r)$  is a finite partition of the positive integers. Further, assuming that  $H$  is non-atomic, the  $\phi_k$ 's are distinct with probability one, so any partition of the positive integers corresponds to some partition of  $\Theta$ . Thus, for each  $j$  we have:

$$\begin{aligned} (G_j(A_1), \dots, G_j(A_r)) &\sim \text{Dir}(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r)) \\ \Rightarrow \left( \sum_{k \in K_1} \pi_{jk}, \dots, \sum_{k \in K_r} \pi_{jk} \right) &\sim \text{Dir} \left( \alpha_0 \sum_{k \in K_1} \beta_k, \dots, \alpha_0 \sum_{k \in K_r} \beta_k \right), \end{aligned} \quad (18)$$

for every finite partition of the positive integers. Hence each  $\pi_j$  is independently distributed according to  $\text{DP}(\alpha_0, \beta)$ , where we interpret  $\beta$  and  $\pi_j$  as probability measures on the positive integers. If  $H$  is non-atomic then a weaker result still holds: if  $\pi_j \sim \text{DP}(\alpha_0, \beta)$  then  $G_j$  as given in (17) is still  $\text{DP}(\alpha_0, G_0)$  distributed.

As in the Dirichlet process mixture model, since each factor  $\theta_{ji}$  is distributed according to  $G_j$ , it takes on the value  $\phi_k$  with probability  $\pi_{jk}$ . Again let  $z_{ji}$  be an indicator variable such that  $\theta_{ji} = \phi_{z_{ji}}$ . Given  $z_{ji}$  we have  $x_{ji} \sim F(\phi_{z_{ji}})$ . Thus we obtain an equivalent representation of the hierarchical Dirichlet process mixture via the following conditional distributions:

$$\begin{aligned} \beta \mid \gamma &\sim \text{GEM}(\gamma) \\ \pi_j \mid \alpha_0, \beta &\sim \text{DP}(\alpha_0, \beta) & z_{ji} \mid \pi_j &\sim \pi_j \\ \phi_k \mid H &\sim H & x_{ji} \mid z_{ji}, (\phi_k)_{k=1}^\infty &\sim F(\phi_{z_{ji}}). \end{aligned} \quad (19)$$

We now derive an explicit relationship between the elements of  $\beta$  and  $\pi_j$ . Recall that the stick-breaking construction for Dirichlet processes defines the variables  $\beta_k$  in (16) as follows:

$$\beta'_k \sim \text{Beta}(1, \gamma) \quad \beta_k = \beta'_k \prod_{l=1}^{k-1} (1 - \beta'_l). \quad (20)$$

Using (18), we show that the following stick-breaking construction produces a random probability measure  $\pi_j \sim \text{DP}(\alpha_0, \beta)$ :

$$\pi'_{jk} \sim \text{Beta} \left( \alpha_0 \beta_k, \alpha_0 \left( 1 - \sum_{l=1}^k \beta_l \right) \right) \quad \pi_{jk} = \pi'_{jk} \prod_{l=1}^{k-1} (1 - \pi'_{jl}). \quad (21)$$

To derive (21), first notice that for a partition  $(\{1, \dots, k-1\}, \{k\}, \{k+1, k+2, \dots\})$ , (18) gives:

$$\left( \sum_{l=1}^{k-1} \pi_{jl}, \pi_{jk}, \sum_{l=k+1}^{\infty} \pi_{jl} \right) \sim \text{Dir} \left( \alpha_0 \sum_{l=1}^{k-1} \beta_l, \alpha_0 \beta_k, \alpha_0 \sum_{l=k+1}^{\infty} \beta_l \right). \quad (22)$$

Removing the first element, and using standard properties of the Dirichlet distribution, we have:

$$\frac{1}{1 - \sum_{l=1}^{k-1} \pi_{jl}} \left( \pi_{jk}, \sum_{l=k+1}^{\infty} \pi_{jl} \right) \sim \text{Dir} \left( \alpha_0 \beta_k, \alpha_0 \sum_{l=k+1}^{\infty} \beta_l \right). \quad (23)$$

Finally, define  $\pi'_{jk} = \frac{\pi_{jk}}{1 - \sum_{l=1}^{k-1} \pi_{jl}}$  and observe that  $1 - \sum_{l=1}^k \beta_l = \sum_{l=k+1}^{\infty} \beta_l$  to obtain (21). Together with (20), (16) and (17), this completes the description of the stick-breaking construction for hierarchical Dirichlet processes.

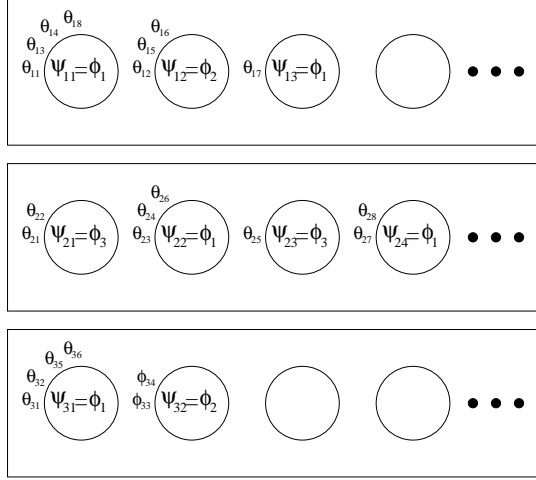


Figure 2: A depiction of a Chinese restaurant franchise. Each restaurant is represented by a rectangle. Customers ( $\theta_{ji}$ 's) are seated at tables (circles) in the restaurants. At each table a dish is served. The dish is served from a global menu ( $\phi_k$ ), whereas the parameter  $\psi_{jt}$  is a table-specific indicator that serves to index items on the global menu. The customer  $\theta_{ji}$  sits at the table to which it has been assigned in (24).

## 4.2 The Chinese restaurant franchise

In this section we describe an analog of the Chinese restaurant process for hierarchical Dirichlet processes that we refer to as the *Chinese restaurant franchise*. In the Chinese restaurant franchise, the metaphor of the Chinese restaurant process is extended to allow multiple restaurants which share a set of dishes.

The metaphor is as follows (see Figure 2). We have a restaurant franchise with a shared menu across the restaurants. At each table of each restaurant one dish is ordered from the menu by the first customer who sits there, and it is shared among all customers who sit at that table. Multiple tables in multiple restaurants can serve the same dish.

In this setup, the restaurants correspond to groups and the customers correspond to the factors  $\theta_{ji}$ . We also let  $\phi_1, \dots, \phi_K$  denote  $K$  i.i.d. random variables distributed according to  $H$ ; this is the global menu of dishes. We also introduce variables  $\psi_{jt}$  which represent the table-specific choice of dishes; in particular,  $\psi_{jt}$  is the dish served at table  $t$  in restaurant  $j$ .

Note that each  $\theta_{ji}$  is associated with one  $\psi_{jt}$ , while each  $\psi_{jt}$  is associated with one  $\phi_k$ . We introduce indicators to denote these associations. In particular, let  $t_{ji}$  be the index of the  $\psi_{jt}$  associated with  $\theta_{ji}$ , and let  $k_{jt}$  be the index of  $\phi_k$  associated with  $\psi_{jt}$ . In the Chinese restaurant franchise metaphor, customer  $i$  in restaurant  $j$  sat at table  $t_{ji}$  while table  $t$  in restaurant  $j$  serves dish  $k_{jt}$ .

We also need a notation for counts. In particular, we need to maintain counts of customers and counts of tables. We use the notation  $n_{jtk}$  to denote the number of customers in restaurant  $j$  at table  $t$  eating dish  $k$ . Marginal counts are represented with dots. Thus,  $n_{jt}$  represents the number of customers in restaurant  $j$  at table  $t$  and  $n_{j\cdot k}$  represents the number of customers in restaurant  $j$  eating dish  $k$ . The notation  $m_{jk}$  denotes the number of tables in restaurant  $j$  serving dish  $k$ . Thus,  $m_j$  represents the number of tables in restaurant  $j$ ,  $m_{\cdot k}$  represents the number of tables serving dish  $k$ , and  $m_{\cdot}$  the total number of tables occupied.

Let us now compute marginals under a hierarchical Dirichlet process when  $G_0$  and  $G_j$  are

integrated out. First consider the conditional distribution for  $\theta_{ji}$  given  $\theta_{j1}, \dots, \theta_{j,i-1}$  and  $G_0$ , where  $G_j$  is integrated out. From (8):

$$\theta_{ji} \mid \theta_{j1}, \dots, \theta_{j,i-1}, \alpha_0, G_0 \sim \sum_{t=1}^{m_j} \frac{n_{jt}}{i-1+\alpha_0} \delta_{\psi_{jt}} + \frac{\alpha_0}{i-1+\alpha_0} G_0, \quad (24)$$

This is a mixture, and a draw from this mixture can be obtained by drawing from the terms on the right-hand side with probabilities given by the corresponding mixing proportions. If a term in the first summation is chosen then we set  $\theta_{ji} = \psi_{jt}$  and let  $t_{ji} = t$  for the chosen  $t$ . If the second term is chosen then we increment  $m_j$  by one, draw  $\psi_{jm_j} \sim G_0$  and set  $\theta_{ji} = \psi_{jm_j}$  and  $t_{ji} = m_j$ .

Now we proceed to integrate out  $G_0$ . Notice that  $G_0$  appears only in its role as the distribution of the variables  $\psi_{jt}$ . Since  $G_0$  is distributed according to a Dirichlet process, we can integrate it out by using (8) again and write the conditional distribution of  $\psi_{jt}$  as:

$$\psi_{jt} \mid \psi_{11}, \psi_{12}, \dots, \psi_{21}, \dots, \psi_{jt-1}, \gamma, H \sim \sum_{k=1}^K \frac{m_{..k}}{m_{..} + \gamma} \delta_{\phi_k} + \frac{\gamma}{m_{..} + \gamma} H. \quad (25)$$

If we draw  $\psi_{jt}$  via choosing a term in the summation on the right-hand side of this equation, we set  $\psi_{jt} = \phi_k$  and let  $k_{jt} = k$  for the chosen  $k$ . If the second term is chosen then we increment  $K$  by one, draw  $\phi_K \sim H$  and set  $\psi_{jt} = \phi_K$  and  $k_{jt} = K$ .

This completes the description of the conditional distributions of the  $\theta_{ji}$  variables. To use these equations to obtain samples of  $\theta_{ji}$ , we proceed as follows. For each  $j$  and  $i$ , first sample  $\theta_{ji}$  using (24). If a new sample from  $G_0$  is needed, we use (25) to obtain a new sample  $\psi_{jt}$  and set  $\theta_{ji} = \psi_{jt}$ .

Note that in the hierarchical Dirichlet process the values of the factors are shared between the groups, as well as within the groups. This is a key property of hierarchical Dirichlet processes.

### 4.3 The infinite limit of finite mixture models

As in the case of a Dirichlet process mixture model, the hierarchical Dirichlet process mixture model can be derived as the infinite limit of finite mixtures. In this section, we present two apparently different finite models that both yield the hierarchical Dirichlet process mixture in the infinite limit, each emphasizing a different aspect of the model.

Consider the following collection of finite mixture models, where  $\beta$  is a global vector of mixing proportions and  $\pi_j$  is a group-specific vector of mixing proportions:

$$\begin{aligned} \beta \mid \gamma &\sim \text{Dir}(\gamma/L, \dots, \gamma/L) \\ \pi_j \mid \alpha_0, \beta &\sim \text{Dir}(\alpha_0 \beta) & z_{ji} \mid \pi_j &\sim \pi_j \\ \phi_k \mid H &\sim H & x_{ji} \mid z_{ji}, (\phi_k)_{k=1}^L &\sim F(\phi_{z_{ji}}). \end{aligned} \quad (26)$$

The parametric hierarchical prior for  $\beta$  and  $\pi$  in (26) has been discussed by MacKay and Peto (1994) as a model for natural languages. We will show that the limit of this model as  $L \rightarrow \infty$  is the hierarchical Dirichlet process. Let us consider the random probability measures  $G_0^L = \sum_{k=1}^L \beta_k \delta_{\phi_k}$  and  $G_j^L = \sum_{k=1}^L \pi_{jk} \delta_{\phi_k}$ . As in Section 3.4, for every measurable function  $f$  integrable with respect to  $H$  we have

$$\int f(\theta) dG_0^L(\theta) \xrightarrow{\mathcal{D}} \int f(\theta) dG_0(\theta), \quad (27)$$

as  $L \rightarrow \infty$ . Further, using standard properties of the Dirichlet distribution, we see that (18) still holds for the finite case for partitions of  $\{1, \dots, L\}$ ; hence we have:

$$G_j^L \sim \text{DP}(\alpha_0, G_0^L). \quad (28)$$

It is now clear that as  $L \rightarrow \infty$  the marginal distribution this finite model induces on  $\mathbf{x}$  approaches the hierarchical Dirichlet process mixture model.

There is an alternative finite model whose limit is also the hierarchical Dirichlet process mixture model. Instead of introducing dependencies between the groups by placing a prior on  $\beta$  (as in the first finite model), each group can instead choose a subset of  $T$  mixture components from a model-wide set of  $L$  mixture components. In particular consider the following model:

$$\begin{aligned} \beta &| \gamma \sim \text{Dir}(\gamma/L, \dots, \gamma/L) & k_{jt} &| \beta \sim \beta \\ \pi_j &| \alpha_0 \sim \text{Dir}(\alpha_0/T, \dots, \alpha_0/T) & t_{ji} &| \pi_j \sim \pi_j \\ \phi_k &| H \sim H & x_{ji} &| t_{ji}, (k_{jt})_{t=1}^T, (\phi_k)_{k=1}^L \sim F(\phi_{k_{jt_{ji}}}). \end{aligned} \quad (29)$$

As  $T \rightarrow \infty$  and  $L \rightarrow \infty$ , the limit of this model is the Chinese restaurant franchise process; hence the infinite limit of this model is also the hierarchical Dirichlet process mixture model.

## 5 INFERENCE

In this section we describe three related Markov chain Monte Carlo sampling schemes for the hierarchical Dirichlet process mixture model. The first is a straightforward Gibbs sampler based on the Chinese restaurant franchise, the second is based upon an augmented representation involving both the Chinese restaurant franchise and the posterior for  $G_0$ , while the third is a variation on the second sampling scheme with streamlined bookkeeping. To simplify the discussion we assume that the base distribution  $H$  is conjugate to the data distribution  $F$ ; this allows us to focus on the issues specific to the hierarchical Dirichlet process. The nonconjugate case can be approached by adapting to the hierarchical Dirichlet process techniques developed for nonconjugate DP mixtures (Neal 2000). Moreover, in this section we assume fixed values for the concentration parameters  $\alpha_0$  and  $\gamma$ ; we present a sampler for these parameters in the appendix.

We recall the random variables of interest. The variables  $x_{ji}$  are the observed data. Each  $x_{ji}$  is assumed to arise as a draw from a distribution  $F(\theta_{ji})$ . Let the factor  $\theta_{ji}$  be associated with the table  $t_{ji}$  in the restaurant representation; i.e., let  $\theta_{ji} = \psi_{jt_{ji}}$ . The random variable  $\psi_{jt}$  is an instance of mixture component  $k_{jt}$ ; i.e.,  $\psi_{jt} = \phi_{k_{jt}}$ . The prior over the parameters  $\phi_k$  is  $H$ . Let  $z_{ji} = k_{jt_{ji}}$  denote the mixture component associated with the observation  $x_{ji}$ . We use the notation  $n_{jtk}$  to denote the number of customers in restaurant  $j$  at table  $t$  eating dish  $k$ , while  $m_{jk}$  denotes the number of tables in restaurant  $j$  serving dish  $k$ . Marginal counts are represented with dots.

Let  $\mathbf{x} = (x_{ji} : \text{all } j, i)$ ,  $\mathbf{x}_{jt} = (x_{ji} : \text{all } i \text{ with } t_{ji} = t)$ ,  $\mathbf{t} = (t_{ji} : \text{all } j, i)$ ,  $\mathbf{k} = (k_{jt} : \text{all } j, t)$ ,  $\mathbf{z} = (z_{ji} : \text{all } j, i)$ ,  $\mathbf{m} = (m_{jk} : \text{all } j, k)$  and  $\phi = (\phi_1, \dots, \phi_K)$ . When a superscript is attached to a set of variables or a count, e.g.,  $x^{-ji}$ ,  $\mathbf{k}^{-jt}$  or  $n_{jt}^{-ji}$ , this means that the variable corresponding to the superscripted index is removed from the set or from the calculation of the count. In the examples,  $x^{-ji} = \mathbf{x} \setminus x_{ji}$ ,  $\mathbf{k}^{-jt} = \mathbf{k} \setminus k_{jt}$  and  $n_{jt}^{-ji}$  is the number of observations in group  $j$  whose factor is associated with  $\psi_{jt}$ , leaving out item  $x_{ji}$ .

Let  $F(\theta)$  have density  $f(\cdot|\theta)$  and  $H$  have density  $h(\cdot)$ . Since  $H$  is conjugate to  $F$  we integrate out the mixture component parameters  $\phi$  in the sampling schemes. Denote the conditional density

of  $x_{ji}$  under mixture component  $k$  given all data items except  $x_{ji}$  as

$$f_k^{-x_{ji}}(x_{ji}) = \frac{\int f(x_{ji}|\phi_k) \prod_{j'i' \neq ji, z_{j'i'}=k} f(x_{j'i'}|\phi_k) h(\phi_k) d\phi_k}{\int \prod_{j'i' \neq ji, z_{j'i'}=k} f(x_{j'i'}|\phi_k) h(\phi_k) d\phi_k}. \quad (30)$$

Similarly denote  $f_k^{-x_{jt}}(\mathbf{x}_{jt})$  as the conditional density of  $x_{jt}$  given all data items associated with mixture component  $k$  leaving out  $\mathbf{x}_{jt}$ .

Finally, we will suppress references to all variables except those being sampled in the conditional distributions to follow, in particular we omit references to  $\mathbf{x}$ ,  $\alpha_0$  and  $\gamma$ .

## 5.1 Posterior sampling in the Chinese restaurant franchise

The Chinese restaurant franchise presented in Section 4.2 can be used to produce samples from the prior distribution over the  $\theta_{ji}$ , as well as intermediary information related to the tables and mixture components. This framework can be adapted to yield a Gibbs sampling scheme for posterior sampling given observations  $\mathbf{x}$ .

Rather than dealing with the  $\theta_{ji}$ 's and  $\psi_{jt}$ 's directly, we shall sample their index variables  $t_{ji}$  and  $k_{jt}$  instead. The  $\theta_{ji}$ 's and  $\psi_{jt}$ 's can be reconstructed from these index variables and the  $\phi_k$ 's. This representation makes the Markov chain Monte Carlo sampling scheme more efficient (cf. Neal 2000). Notice that the  $t_{ji}$  and the  $k_{jt}$  inherit the exchangeability properties of the  $\theta_{ji}$  and the  $\psi_{jt}$ —the conditional distributions in (24) and (25) can be adapted to be expressed in terms of  $t_{ji}$  and  $k_{jt}$ . The state space consists of values of  $\mathbf{t}$  and  $\mathbf{k}$ . Notice that the number of  $k_{jt}$  variables represented explicitly by the algorithm is not fixed. We can think of the actual state space as consisting of an infinite number of  $k_{jt}$ 's; only finitely many are actually associated to data and represented explicitly.

**Sampling  $t$ .** To compute the conditional distribution of  $t_{ji}$  given the remainder of the variables, we make use of exchangeability and treat  $t_{ji}$  as the last variable being sampled in the last group in (24) and (25). We obtain the conditional posterior for  $t_{ji}$  by combining the conditional prior distribution for  $t_{ji}$  with the likelihood of generating  $x_{ji}$ .

Using (24), the prior probability that  $t_{ji}$  takes on a particular previously used value  $t$  is proportional to  $n_{jt}^{-j_i}$ , whereas the probability that it takes on a new value (say  $t^{\text{new}} = m_j + 1$ ) is proportional to  $\alpha_0$ . The likelihood due to  $x_{ji}$  given  $t_{ji} = t$  for some previously used  $t$  is  $f_k^{-x_{ji}}(x_{ji})$ . The likelihood for  $t_{ji} = t^{\text{new}}$  can be calculated by integrating out the possible values of  $k_{jt^{\text{new}}}$  using (25):

$$p(x_{ji} | \mathbf{t}^{-j_i}, t_{ji} = t^{\text{new}}, \mathbf{k}) = \sum_{k=1}^K \frac{m_{.k}}{m_{..} + \gamma} f_k^{-x_{ji}}(x_{ji}) + \frac{\gamma}{m_{..} + \gamma} f_{k^{\text{new}}}^{-x_{ji}}(x_{ji}), \quad (31)$$

where  $f_{k^{\text{new}}}^{-x_{ji}}(x_{ji}) = \int f(x_{ji}|\phi)h(\phi)d\phi$  is simply the prior density of  $x_{ji}$ . The conditional distribution of  $t_{ji}$  is then

$$p(t_{ji} = t | \mathbf{t}^{-j_i}, \mathbf{k}) \propto \begin{cases} n_{jt}^{-j_i} f_{k_{jt}}^{-x_{ji}}(x_{ji}) & \text{if } t \text{ previously used,} \\ \alpha_0 p(x_{ji} | \mathbf{t}^{-j_i}, t_{ji} = t^{\text{new}}, \mathbf{k}) & \text{if } t = t^{\text{new}}. \end{cases} \quad (32)$$

If the sampled value of  $t_{ji}$  is  $t^{\text{new}}$ , we obtain a sample of  $k_{jt^{\text{new}}}$  by sampling from (31):

$$p(k_{jt^{\text{new}}} = k | \mathbf{t}, \mathbf{k}^{-j_t^{\text{new}}}) \propto \begin{cases} m_{.k} f_k^{-x_{ji}}(x_{ji}) & \text{if } k \text{ previously used,} \\ \gamma f_{k^{\text{new}}}^{-x_{ji}}(x_{ji}) & \text{if } k = k^{\text{new}}. \end{cases} \quad (33)$$

If as a result of updating  $t_{ji}$  some table  $t$  becomes unoccupied, i.e.,  $n_{jt} = 0$ , then the probability that this table will be reoccupied in the future will be zero, since this is always proportional to  $n_{jt}$ . As a result, we may delete the corresponding  $k_{jt}$  from the data structure. If as a result of deleting  $k_{jt}$  some mixture component  $k$  becomes unallocated, we delete this mixture component as well.

**Sampling  $k$ .** Since changing  $k_{jt}$  actually changes the component membership of all data items in table  $t$ , the likelihood obtained by setting  $k_{jt} = k$  is given by  $f_k^{-\mathbf{x}_{jt}}(\mathbf{x}_{jt})$ , so that the conditional probability of  $k_{jt}$  is

$$p(k_{jt} = k \mid \mathbf{t}, \mathbf{k}^{-jt}) \propto \begin{cases} m_{\cdot k}^{-jt} f_k^{-\mathbf{x}_{jt}}(\mathbf{x}_{jt}) & \text{if } k \text{ is previously used,} \\ \gamma f_{k^{\text{new}}}^{-\mathbf{x}_{jt}}(\mathbf{x}_{jt}) & \text{if } k = k^{\text{new}}. \end{cases} \quad (34)$$

## 5.2 Posterior sampling with an augmented representation

In the Chinese restaurant franchise sampling scheme, the sampling for all groups is coupled since  $G_0$  is integrated out. This complicates matters in more elaborate models (e.g., in the case of the hidden Markov model considered in Section 7). In this section we describe an alternative sampling scheme where in addition to the Chinese restaurant franchise representation,  $G_0$  is instantiated and sampled from so that the posterior conditioned on  $G_0$  factorizes across groups.

Given a posterior sample  $(\mathbf{t}, \mathbf{k})$  from the Chinese restaurant franchise representation, we can obtain a draw from the posterior of  $G_0$  by noting that  $G_0 \sim \text{DP}(\gamma, H)$  and  $\psi_{jt}$  for each table  $t$  is a draw from  $G_0$ . Conditioning on the  $\psi_{jt}$ 's,  $G_0$  is now distributed as  $\text{DP}(\gamma + m_{\cdot}, \frac{\gamma H + \sum_{k=1}^K m_{\cdot k} \delta_{\phi_k}}{\gamma + m_{\cdot}})$ . An explicit construction for  $G_0$  is now given as

$$\begin{aligned} \boldsymbol{\beta} = (\beta_1, \dots, \beta_K, \beta_u) &\sim \text{Dir}(m_{\cdot 1}, \dots, m_{\cdot K}, \gamma) & G_u &\sim \text{DP}(\gamma, H) \\ p(\phi_k \mid \mathbf{t}, \mathbf{k}) &\propto h(\phi_k) \prod_{ji: k_{jt_{ji}}=k} f(x_{ji} \mid \phi_k) & G_0 &= \sum_{k=1}^K \beta_k \delta_{\phi_k} + \beta_u G_u \end{aligned} \quad (35)$$

Given a sample of  $G_0$  the posterior for each group is factorized and sampling in each group can be performed separately. The variables of interest in this scheme are  $\mathbf{t}$  and  $\mathbf{k}$  as in the Chinese restaurant franchise sampling scheme and  $\boldsymbol{\beta}$  above, while both  $\phi$  and  $G_u$  are integrated out (this introduces couplings into the sampling for each group but is easily handled).

**Sampling for  $\mathbf{t}$  and  $\mathbf{k}$**  is almost identical to the Chinese restaurant franchise sampling scheme. The only novelty is that we replace  $m_{\cdot k}$  by  $\beta_k$  and  $\gamma$  by  $\beta_u$  in (31), (32), (33) and (34), and when a new component  $k^{\text{new}}$  is instantiated we draw  $b \sim \text{Beta}(1, \gamma)$  and set  $\beta_{k^{\text{new}}} = b\beta_u$  and  $\beta_u^{\text{new}} = (1 - b)\beta_u$ . We can understand  $b$  as follows: when a new component is instantiated, it is instantiated from  $G_u$  by choosing an atom in  $G_u$  with probability given by its weight  $b$ . Using the fact that the sequence of stick-breaking weights is a size-biased permutation of the weights in a draw from a Dirichlet process (Pitman 1996), the weight  $b$  corresponding to the chosen atom in  $G_u$  will have the same distribution as the first stick-breaking weight, i.e.,  $\text{Beta}(1, \gamma)$ .

**Sampling for  $\boldsymbol{\beta}$**  has already been described in (35):

$$(\beta_1, \dots, \beta_K, \beta_u) \mid \mathbf{t}, \mathbf{k} \sim \text{Dir}(m_{\cdot 1}, \dots, m_{\cdot K}, \gamma). \quad (36)$$

## 5.3 Posterior sampling by direct assignment

In both the Chinese restaurant franchise and augmented representation sampling schemes, data items are first assigned to some table  $t_{ji}$ , and the tables are then assigned to some mixture component  $k_{jt}$ .

This indirect association to mixture components can make the bookkeeping somewhat involved. In this section we describe a variation on the augmented representation sampling scheme that directly assigns data items to mixture components via a variable  $z_{ji}$  which is equivalent to  $k_{jt_{ji}}$  in the earlier sampling schemes. The tables are only represented in terms of the numbers of tables  $m_{jk}$ .

**Sampling  $z$**  can be realized by grouping together terms associated with each  $k$  in (31) and (32):

$$p(z_{ji} = k \mid \mathbf{z}^{-ji}, \mathbf{m}, \boldsymbol{\beta}) = \begin{cases} (n_{j \cdot k}^{-ji} + \alpha_0 \beta_k) f_k^{-x_{ji}}(x_{ji}) & \text{if } k \text{ previously used,} \\ \alpha_0 \beta_u f_{k^{\text{new}}}^{-x_{ji}}(x_{ji}) & \text{if } k = k^{\text{new}}. \end{cases} \quad (37)$$

where we have replaced  $m_{\cdot k}$  with  $\beta_k$  and  $\gamma$  with  $\beta_u$ .

**Sampling  $\mathbf{m}$ .** In the augmented representation sampling scheme, conditioned on the assignment of data items to mixture components  $\mathbf{z}$ , the only effect of  $\mathbf{t}$  and  $\mathbf{k}$  on other variables is via  $\mathbf{m}$  in the conditional distribution of  $\boldsymbol{\beta}$  in (36). As a result it is sufficient to sample  $\mathbf{m}$  in place of  $\mathbf{t}$  and  $\mathbf{k}$ . To obtain the distribution of  $m_{jk}$  conditioned on other variables, consider the distribution of  $t_{ji}$  assuming that  $k_{jt_{ji}} = z_{ji}$ . The probability that data item  $x_{ji}$  is assigned to some table  $t$  such that  $k_{jt} = k$  is

$$p(t_{ji} = t \mid k_{jt} = k, \mathbf{t}^{-ji}, \mathbf{k}, \boldsymbol{\beta}) \propto n_{jt}^{-ji}, \quad (38)$$

while the probability that it is assigned a new table under component  $k$  is

$$p(t_{ji} = t^{\text{new}} \mid k_{jt^{\text{new}}} = k, \mathbf{t}^{-ji}, \mathbf{k}, \boldsymbol{\beta}) \propto \alpha_0 \beta_k. \quad (39)$$

These equations form the conditional distributions of a Gibbs sampler whose equilibrium distribution is the prior distribution over the assignment of  $n_{j \cdot k}$  observations to components in an ordinary Dirichlet process with concentration parameter  $\alpha_0 \beta_k$ . The corresponding distribution over the number of components is then the desired conditional distribution of  $m_{jk}$ . Antoniak (1974) has shown that this is:

$$p(m_{jk} = m \mid \mathbf{z}, \mathbf{m}^{-jk}, \boldsymbol{\beta}) = \frac{\Gamma(\alpha_0 \beta_k)}{\Gamma(\alpha_0 \beta_k + n_{j \cdot k})} s(n_{j \cdot k}, m) (\alpha_0 \beta_k)^m, \quad (40)$$

where  $s(n, m)$  are unsigned Stirling numbers of the first kind. We have by definition that  $s(0, 0) = s(1, 1) = 1$ ,  $s(n, 0) = 0$  for  $n > 0$  and  $s(n, m) = 0$  for  $m > n$ . Other entries can be computed as  $s(n + 1, m) = s(n, m - 1) + n s(n, m)$ .

**Sampling for  $\boldsymbol{\beta}$**  is the same as in the augmented sampling scheme and is given by (36).

## 5.4 Comparison of Sampling Schemes

Let us now consider the relative merits of these three sampling schemes. In terms of ease of implementation, the direct assignment scheme is preferred because its bookkeeping is straightforward. The two schemes based on the Chinese restaurant franchise involve more substantial effort. In addition, both the augmented and direct assignment schemes sample rather than integrate out  $G_0$ , and as a result the sampling of the groups is decoupled given  $G_0$ . This simplifies the sampling schemes and makes them applicable in elaborate models such as the hidden Markov model in Section 7.

In terms of convergence speed, the direct assignment scheme changes the component membership of data items one at a time, while in both schemes using the Chinese restaurant franchise changing the component membership of one table will change the membership of multiple data items at the same time, leading to potentially improved performance. This is akin to split-and-merge

techniques in Dirichlet process mixture modeling (Jain and Neal 2000). This analogy is, however, somewhat misleading in that unlike split-and-merge methods, the assignment of data items to tables is a consequence of the *prior* clustering effect of a Dirichlet process with  $n_{j,k}$  samples. As a result, we expect that the probability of obtaining a successful reassignment of a table to another previously used component will often be small, and we do not necessarily expect the Chinese restaurant franchise schemes to dominate the direct assignment scheme.

The inference methods presented here should be viewed as first steps in the development of inference procedures for hierarchical Dirichlet process mixtures. More sophisticated methods—such as split-and-merge methods (Jain and Neal 2000) and variational methods (Blei and Jordan 2005)—have shown promise for Dirichlet processes and we expect that they will prove useful for hierarchical Dirichlet processes as well.

## 6 EXPERIMENTS

We describe two experiments in this section to highlight the two aspects of the hierarchical Dirichlet process: its nonparametric nature and its hierarchical nature. In the next section we present a third experiment highlighting the ease with which we can extend the framework to more complex models, specifically a hidden Markov model with a countably infinite state space.

### 6.1 Document modeling

Recall the problem of document modeling discussed in Section 1. Following standard methodology in the information retrieval literature (Salton and McGill 1983), we view a document as a “bag of words”; that is, we make an exchangeability assumption for the words in the document. Moreover, we model the words in a document as arising from a mixture model, in which a mixture component—a “topic”—is a multinomial distribution over words from some finite and known vocabulary. The goal is to model a corpus of documents in such a way as to allow the topics to be shared among the documents in a corpus.

A parametric approach to this problem is provided by the *latent Dirichlet allocation* (LDA) model of Blei et al. (2003). This model involves a finite mixture model in which the mixing proportions are drawn on a document-specific basis from a Dirichlet distribution. Moreover, given these mixing proportions, each word in the document is an independent draw from the mixture model. That is, to generate a word, a mixture component (i.e., a topic) is selected, and then a word is generated from that topic.

Note that the assumption that each word is associated with a possibly different topic differs from a model in which a mixture component is selected once per document, and then words are generated i.i.d. from the selected topic. Moreover, it is interesting to note that the same distinction arises in population genetics, where multiple words in a document are analogous to multiple markers along a chromosome. Indeed, Pritchard et al. (2000) have developed a model in which marker probabilities are selected once per marker; their model is essentially identical to LDA.

As in simpler finite mixture models, it is natural to try to extend LDA and related models by using Dirichlet processes to capture uncertainty regarding the number of mixture components. This is somewhat more difficult than in the case of a simple mixture model, however, because in the LDA model the documents have document-specific mixing proportions. We thus require multiple DPs, one for each document. This then poses the problem of sharing mixture components across multiple DPs, precisely the problem that the hierarchical DP is designed to solve.

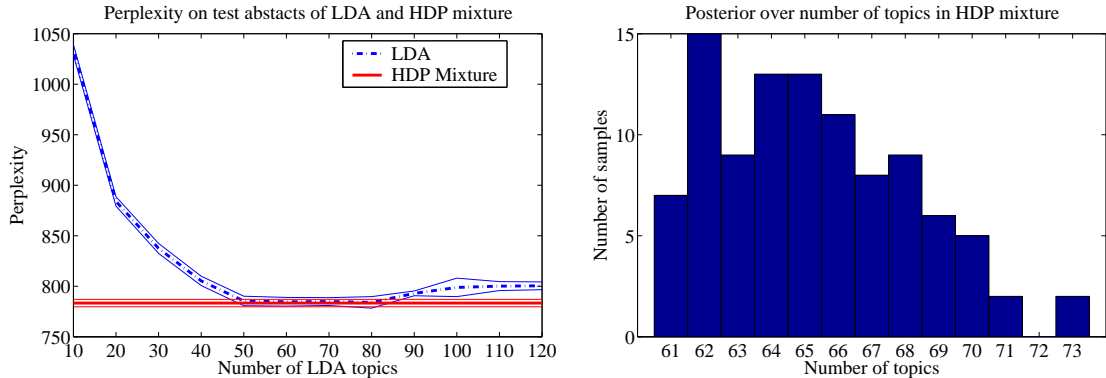


Figure 3: (Left) Comparison of latent Dirichlet allocation and the hierarchical Dirichlet process mixture. Results are averaged over 10 runs; the error bars are one standard error. (Right) Histogram of the number of topics for the hierarchical Dirichlet process mixture over 100 posterior samples.

The hierarchical DP extension of LDA thus takes the following form. Given an underlying measure  $H$  on multinomial probability vectors, we select a random measure  $G_0$  which provides a countably infinite collection of multinomial probability vectors; these can be viewed as the set of all topics that can be used in a given corpus. For the  $j$ th document in the corpus we sample  $G_j$  using  $G_0$  as a base measure; this selects specific subsets of topics to be used in document  $j$ . From  $G_j$  we then generate a document by repeatedly sampling specific multinomial probability vectors  $\theta_{ji}$  from  $G_j$  and sampling words  $x_{ji}$  with probabilities  $\theta_{ji}$ . The overlap among the random measures  $G_j$  implements the sharing of topics among documents.

We fit both the standard parametric LDA model and its hierarchical DP extension to a corpus of nematode biology abstracts (see <http://elegans.swmed.edu/wli/cgcbib>). There are 5838 abstracts in total. After removing standard stop words and words appearing fewer than 10 times, we are left with 476441 words in total. Following standard information retrieval methodology, the vocabulary is defined as the set of distinct words left in all abstracts; this has size 5699.

Both models were as similar as possible beyond the distinction that LDA assumes a fixed finite number of topics while the hierarchical Dirichlet process does not. Both models used a symmetric Dirichlet distribution with parameters of 0.5 for the prior  $H$  over topic distributions. The concentration parameters were given vague gamma priors,  $\gamma \sim \text{Gamma}(1, .1)$  and  $\alpha_0 \sim \text{Gamma}(1, 1)$ . The distribution over topics in LDA is assumed to be symmetric Dirichlet with parameters  $\alpha_0/L$  with  $L$  being the number of topics;  $\gamma$  is not used in LDA. Posterior samples were obtained using the Chinese restaurant franchise sampling scheme, while the concentration parameters were sampled using the auxiliary variable sampling scheme presented in the appendix.

We evaluated the models via 10-fold cross-validation. The evaluation metric was the *perplexity*, a standard metric in the information retrieval literature. The perplexity of a held-out abstract consisting of words  $w_1, \dots, w_I$  is defined to be:

$$\exp\left(-\frac{1}{I} \log p(w_1, \dots, w_I | \text{Training corpus})\right) \quad (41)$$

where  $p(\cdot)$  is the probability mass function for a given model.

The results are shown in Figure 3. For LDA we evaluated the perplexity for mixture component cardinalities ranging between 10 and 120. As seen in Figure 3 (Left), the hierarchical DP mixture approach—which integrates over the mixture component cardinalities—performs as well as the

best LDA model, doing so without any form of model selection procedure. Moreover, as shown in Figure 3 (Right), the posterior over the number of topics obtained under the hierarchical DP mixture model is consistent with this range of the best-fitting LDA models.

## 6.2 Multiple corpora

We now consider the problem of sharing clusters among the documents in multiple corpora. We approach this problem by extending the hierarchical Dirichlet process to a third level. A draw from a top-level DP yields the base measure for each of a set of corpus-level DPs. Draws from each of these corpus-level DPs yield the base measures for DPs associated with the documents within a corpus. Finally, draws from the document-level DPs provide a representation of each document as a probability distribution across topics (which are distributions across words). The model allows topics to be shared both within each corpus and between corpora.

The documents that we used for these experiments consist of articles from the proceedings of the *Neural Information Processing Systems* (NIPS) conference for the years 1988-1999. The original articles are available at <http://books.nips.cc>; we use a preprocessed version available at <http://www.cs.utoronto.ca/~roweis/nips>. The NIPS conference deals with a range of topics covering both human and machine intelligence. Articles are separated into nine sections: algorithms and architectures (AA), applications (AP), cognitive science (CS), control and navigation (CN), implementations (IM), learning theory (LT), neuroscience (NS), signal processing (SP), vision sciences (VS). (These are the sections used in the years 1995-1999. The sectioning in earlier years differed slightly; we manually relabeled sections from the earlier years to match those used in 1995-1999.) We treat these sections as “corpora,” and are interested in the pattern of sharing of topics among these corpora.

There were 1447 articles in total. Each article was modeled as a bag-of-words. We culled standard stop words as well as words occurring more than 4000 or fewer than 50 times in the whole corpus. This left us with on average slightly more than 1000 words per article.

We considered the following experimental setup. Given a set of articles from a single NIPS section that we wish to model (the VS section in the experiments that we report below), we wish to know whether it is of value (in terms of prediction performance) to include articles from other NIPS sections. This can be done in one of two ways: we can lump all of the articles together without regard for the division into sections, or we can use the hierarchical DP approach to link the sections. Thus we consider three models (see Figure 4 for graphical representations of these models):

- **M1:** This model ignores articles from the other sections and simply uses a hierarchical DP mixture of the kind presented in Section 6.1 to model the VS articles. This model serves as a baseline. We used  $\gamma \sim \text{Gamma}(5, 0.1)$  and  $\alpha_0 \sim \text{Gamma}(0.1, 0.1)$  as prior distributions for the concentration parameters.
- **M2:** This model incorporates articles from other sections, but ignores the distinction into sections, using a single hierarchical DP mixture model to model all of the articles. Priors of  $\gamma \sim \text{Gamma}(5, 0.1)$  and  $\alpha_0 \sim \text{Gamma}(0.1, 0.1)$  were used.
- **M3:** This model takes a full hierarchical approach and models the NIPS sections as multiple corpora, linked via the hierarchical DP mixture formalism. The model is a tree, in which the root is a draw from a single DP for all articles, the first level is a set of draws from DPs for the NIPS sections, and the second level is set of draws from DPs for the articles within sections. Priors of  $\gamma \sim \text{Gamma}(5, 0.1)$ ,  $\alpha_0 \sim \text{Gamma}(5, 0.1)$ , and  $\alpha_1 \sim \text{Gamma}(0.1, 0.1)$  were used.

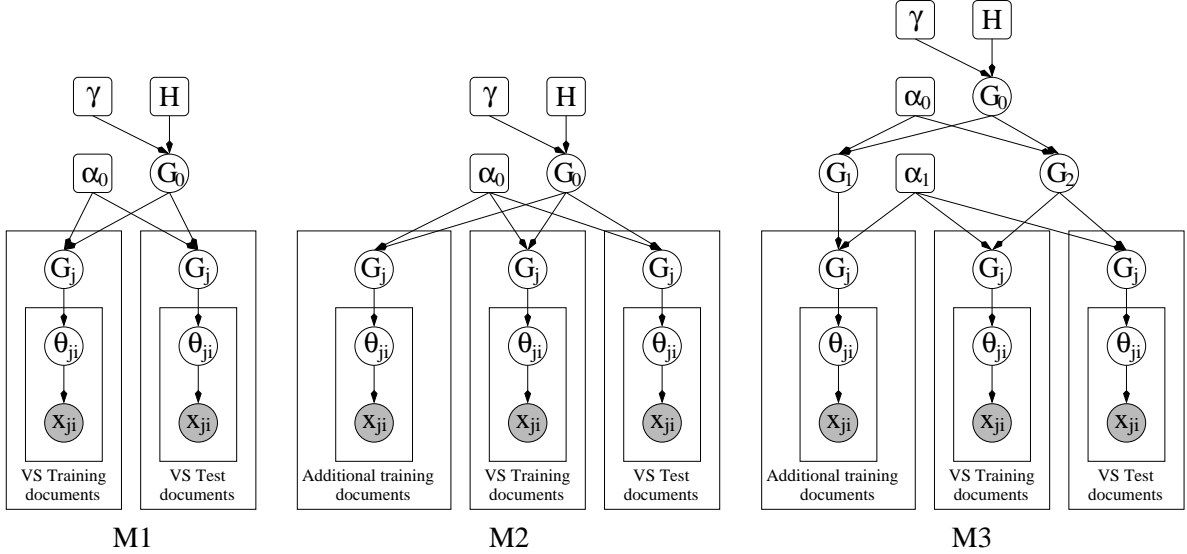


Figure 4: Three models for the NIPS data. From left to right: M1, M2 and M3.

In all models a finite and known vocabulary is assumed and the base measure  $H$  used is a symmetric Dirichlet distribution with parameters of 0.5.

We conducted experiments in which a set of 80 articles were chosen uniformly at random from one of the sections other than VS (this was done to balance the impact of different sections, which are of different sizes). A training set of 80 articles were also chosen uniformly at random from the VS section, as were an additional set of 47 test articles distinct from the training articles. Results report predictive performance on VS test articles based on a training set consisting of the 80 articles in the additional section and  $N$  VS training articles where  $N$  varies between 0 and 80. The direct assignment sampling scheme is used, while concentration parameters are sampled using the auxiliary variable sampling scheme in the appendix.

Figure 5 (Left) presents the average predictive performance for all three models over 5 runs as the number  $N$  of VS training articles ranged from 0 to 80. The performance is measured in terms of the perplexity of single words in the test articles given the training articles, averaged over the choice of which additional section was used. As seen in the figure, the fully hierarchical model M3 performs best, with perplexity decreasing rapidly with modest values of  $N$ . For small values of  $N$ , the performance of M1 is quite poor, but the performance approaches that of M3 when more than 20 articles are included in the VS training set. The performance of the partially-hierarchical M2 was poorer than the fully-hierarchical M3 throughout the range of  $N$ . M2 dominated M1 for small  $N$ , but yielded poorer performance than M1 for  $N$  greater than 14. Our interpretation is that the sharing of strength based on other articles is useful when little other information is available (small  $N$ ), but that eventually (medium to large  $N$ ) there is crosstalk between the sections and it is preferable to model them separately and share strength via the hierarchy.

While the results in Figure 5 (Left) are an average over the sections, it is also of interest to see which sections are the most beneficial in terms of enhancing the prediction of the articles in VS. Figure 5 (Right) plots the predictive performance for model M3 when given data from each of three particular sections: LT, AA and AP. While articles in the LT section are concerned mostly with theoretical properties of learning algorithms, those in AA are mostly concerned with models and methodology, and those in AP are mostly concerned with applications of learning algorithms to

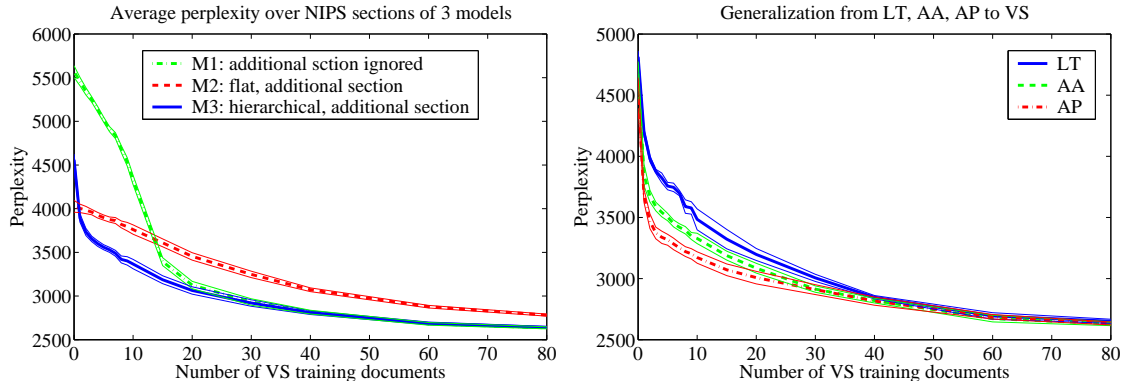


Figure 5: (Left) Perplexity of single words in test VS articles given training articles from VS and another section for 3 different models. Curves shown are averaged over the other sections and 5 runs. (Right) Perplexity of test VS articles given LT, AA and AP articles respectively, using M3, averaged over 5 runs. In both plots, the error bars represent one standard error.

various problems. As seen in the figure, we see that predictive performance is enhanced the most by prior exposure to articles from AP, less by articles from AA, and still less by articles from LT. Given that articles in VS tend to be concerned with the practical application of learning algorithms to problems in computer vision, this pattern of transfer seems reasonable.

Finally, it is of interest to investigate the subject matter content of the topics discovered by the hierarchical DP model. We did so in the following experimental setup. For a given section other than VS (e.g., AA), we fit a model based on articles from that section. We then introduced articles from the VS section and continued to fit the model, while holding the topics found from the earlier fit fixed, and recording which topics from the earlier section were allocated to words in the VS section. Table 1 displays representations of the two most frequently occurring topics in this setup (a topic is represented by the set of words which have highest probability under that topic). These topics provide qualitative confirmation of our expectations regarding the overlap between VS and other sections.

## 7 HIDDEN MARKOV MODELS

The simplicity of the hierarchical DP specification—the base measure for a DP is distributed as a DP—makes it straightforward to exploit the hierarchical DP as a building block in more complex models. In this section we demonstrate this in the case of the hidden Markov model.

Recall that a hidden Markov model (HMM) is a doubly stochastic Markov chain in which a sequence of multinomial “state” variables  $(v_1, v_2, \dots, v_T)$  are linked via a state transition matrix, and each element  $y_t$  in a sequence of “observations”  $(y_1, y_2, \dots, y_T)$  is drawn independently of the other observations conditional on  $v_t$  (Rabiner 1989). This is essentially a dynamic variant of a finite mixture model, in which there is one mixture component corresponding to each value of the multinomial state. As with classical finite mixtures, it is interesting to consider replacing the finite mixture underlying the HMM with a Dirichlet process.

Note that the HMM involves not a single mixture model, but rather a set of mixture models—one for each value of the current state. That is, the “current state”  $v_t$  indexes a specific row of the transition matrix, with the probabilities in this row serving as the mixing proportions for the choice

Table 1: Topics shared between VS and the other NIPS sections. These topics are the most frequently occurring in the VS fit, under the constraint that they are associated with a significant number of words (greater than 2500) from the other section.

CS	task representation pattern processing trained representations three process unit patterns examples concept similarity bayesian hypotheses generalization numbers positive classes hypothesis
NS	cells cell activity response neuron visual patterns pattern single fig visual cells cortical orientation receptive contrast spatial cortex stimulus tuning
LT	signal layer gaussian cells fig nonlinearity nonlinear rate eq cell large examples form point see parameter consider random small optimal
AA	algorithms test approach methods based point problems form large paper distance tangent image images transformation transformations pattern vectors convolution simard
IM	processing pattern approach architecture single shows simple based large control motion visual velocity flow target chip eye smooth direction optical
SP	visual images video language image pixel acoustic delta lowpass flow signals separation signal sources source matrix blind mixing gradient eq
AP	approach based trained test layer features table classification rate paper image images face similarity pixel visual database matching facial examples
CN	ii tree pomdp observable strategy class stochastic history strategies density policy optimal reinforcement control action states actions step problems goal

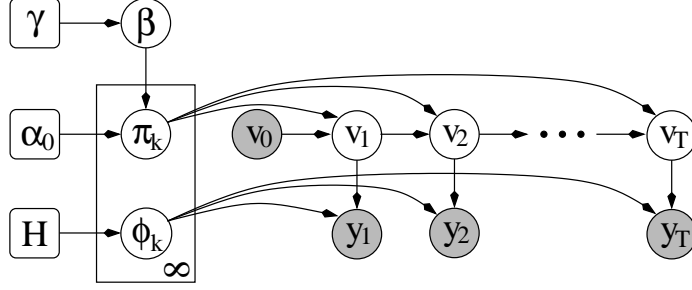


Figure 6: A graphical representation of a hierarchical Dirichlet process hidden Markov model.

of the “next state”  $v_{t+1}$ . Given the next state  $v_{t+1}$ , the observation  $y_{t+1}$  is drawn from the mixture component indexed by  $v_{t+1}$ . Thus, to consider a nonparametric variant of the HMM which allows an unbounded set of states, we must consider a set of DPs, one for each value of the current state. Moreover, these DPs must be linked, because we want the same set of “next states” to be reachable from each of the “current states.” This amounts to the requirement that the atoms associated with the state-conditional DPs should be shared—exactly the framework of the hierarchical DP.

Thus, we can define a nonparametric hidden Markov model by simply replacing the set of conditional finite mixture models underlying the classical HMM with a hierarchical Dirichlet process mixture model. We refer to the resulting model as a *hierarchical Dirichlet process hidden Markov model* (HDP-HMM). The HDP-HMM provides an alternative to methods that place an explicit parametric prior on the number of states or make use of model selection methods to select a fixed number of states (Stolcke and Omohundro 1993).

In work that served as an inspiration for the HDP-HMM, Beal et al. (2002) discussed a model known as the *infinite hidden Markov model*, in which the number of hidden states of a hidden Markov model is allowed to be countably infinite. Indeed, Beal et al. (2002) defined a notion of “hierarchical Dirichlet process” for this model, but their “hierarchical Dirichlet process” is not hierarchical in the Bayesian sense—involving a distribution on the parameters of a Dirichlet process—but is instead a description of a coupled set of urn models. We briefly review this construction, and relate it to our formulation.

Beal et al. (2002) considered the following two-level procedure for determining the transition probabilities of a Markov chain with an unbounded number of states. At the first level, the probability of transitioning from a state  $u$  to a state  $v$  is proportional to the number of times the same transition is observed at other time steps, while with probability proportional to  $\alpha_0$  an “oracle” process is invoked. At this second level, the probability of transitioning to state  $v$  is proportional to the number of times state  $v$  has been chosen by the oracle (regardless of the previous state), while the probability of transitioning to a novel state is proportional to  $\gamma$ . The intended role of the oracle is to tie together the transition models so that they have destination states in common, in much the same way that the baseline distribution  $G_0$  ties together the group-specific mixture components in the hierarchical Dirichlet process.

To relate this two-level urn model to the hierarchical DP framework, let us describe a representation of the HDP-HMM using the stick-breaking formalism. In particular, consider the hierarchical Dirichlet process representation shown in Figure 6. The parameters in this representation have the following distributions:

$$\beta \mid \gamma \sim \text{GEM}(\gamma) \quad \pi_k \mid \alpha_0, \beta \sim \text{DP}(\alpha_0, \beta) \quad \phi_k \mid H \sim H, \quad (42)$$

for each  $k = 1, 2, \dots$ , while for time steps  $t = 1, \dots, T$  the state and observation distributions are:

$$v_t \mid v_{t-1}, (\boldsymbol{\pi}_k)_{k=1}^{\infty} \sim \boldsymbol{\pi}_{v_{t-1}} \quad y_t \mid v_t, (\phi_k)_{k=1}^{\infty} \sim F(\phi_{v_t}), \quad (43)$$

where we assume for simplicity that there is a distinguished initial state  $v_0$ . If we now consider the Chinese restaurant franchise representation of this model as discussed in Section 5, it turns out that the result is equivalent to the coupled urn model of Beal et al. (2002), hence the infinite hidden Markov model is an HDP-HMM.

Unfortunately, posterior inference using the Chinese restaurant franchise representation is awkward for this model, involving substantial bookkeeping. Indeed, Beal et al. (2002) did not present an MCMC inference algorithm for the infinite hidden Markov model, proposing instead a heuristic approximation to Gibbs sampling. On the other hand, both the augmented representation and direct assignment representations lead directly to MCMC sampling schemes that are straightforward to implement. In the experiments reported in the following section we used the direct assignment representation.

Practical applications of hidden Markov models often consider sets of sequences, and treat these sequences as exchangeable at the level of sequences. Thus, in applications to speech recognition, a hidden Markov model for a given word in the vocabulary is generally trained via replicates of that word being spoken. This setup is readily accommodated within the hierarchical DP framework by simply considering an additional level of the Bayesian hierarchy, letting a master Dirichlet process couple each of the HDP-HMMs, each of which is a set of Dirichlet processes.

## 7.1 Alice in Wonderland

In this section we report experimental results for the problem of predicting strings of letters in sentences taken from Lewis Carroll’s *Alice’s Adventures in Wonderland*, comparing the HDP-HMM to other HMM-related approaches.

Each sentence is treated as a sequence of letters and spaces (rather than as a sequence of words). There are 27 distinct symbols (26 letters and space); cases and punctuation marks are ignored. There are 20 training sentences with average length of 51 symbols, and there are 40 test sentences with an average length of 100. The base distribution  $H$  is a symmetric Dirichlet distribution over 27 symbols with parameters 0.1. The concentration parameters  $\gamma$  and  $\alpha_0$  are given Gamma(1, 1) priors.

Using the direct assignment sampling method for posterior predictive inference, we compared the HDD-HMM to a variety of other methods for prediction using hidden Markov models: (1) a classical HMM using maximum likelihood (ML) parameters obtained via the Baum-Welch algorithm (Rabiner 1989), (2) a classical HMM using maximum a posteriori (MAP) parameters, taking the priors to be independent, symmetric Dirichlet distributions for both the transition and emission probabilities, and (3) a classical HMM trained using an approximation to a full Bayesian analysis—in particular, a variational Bayesian (VB) method due to MacKay (1997) and described in detail in Beal (2003). For each of these classical HMMs, we conducted experiments for each value of the state cardinality ranging from 1 to 60.

We present the perplexity on test sentences in Figure 7 (Left). For VB, the predictive probability is intractable to compute, so the modal setting of parameters was used. Both MAP and VB models were given optimal settings of the hyperparameters found using the HDP-HMM. We see that the HDP-HMM has a lower perplexity than all of the models tested for ML, MAP, and VB. Figure 7 (Right) shows posterior samples of the number of states used by the HDP-HMM.

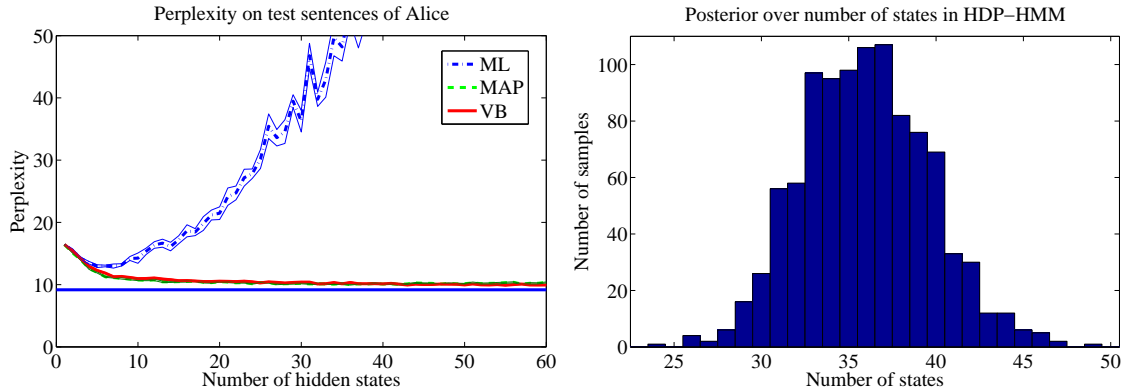


Figure 7: (Left) Comparing the HDP-HMM (solid horizontal line) with ML, MAP and VB trained hidden Markov models. The error bars represent one standard error (those for the HDP-HMM are too small to see). (Right) Histogram for the number of states in the HDP-HMM over 1000 posterior samples.

## 8 DISCUSSION

We have described a nonparametric approach to the modeling of groups of data, where each group is characterized by a mixture model and we allow mixture components to be shared between groups. We have proposed a hierarchical Bayesian solution to this problem, in which a set of Dirichlet processes are coupled via their base measure, which is itself distributed according to a Dirichlet process.

We have described three different representations that capture aspects of the hierarchical Dirichlet process. In particular, we described a stick-breaking representation that describes the random measures explicitly, a representation of marginals in terms of an urn model that we referred to as the “Chinese restaurant franchise,” and a representation of the process in terms of an infinite limit of finite mixture models.

These representations led to the formulation of three Markov chain Monte Carlo sampling schemes for posterior inference under hierarchical Dirichlet process mixtures. The first scheme is based directly on the Chinese restaurant franchise representation, the second scheme represents the posterior using both a Chinese restaurant franchise and a sample from the global measure, while the third uses a direct assignment of data items to mixture components.

Clustering is an important activity in many large-scale data analysis problems in engineering and science, reflecting the heterogeneity that is often present when data are collected on a large scale. Clustering problems can be approached within a probabilistic framework via finite mixture models (Fraley and Raftery 2002; Green and Richardson 2001), and recent years have seen numerous examples of applications of finite mixtures and their dynamical cousins the HMM in areas such as bioinformatics (Durbin et al. 1998), speech recognition (Huang et al. 2001), information retrieval (Blei et al. 2003) and computational vision (Forsyth and Ponce 2002). These areas also provide numerous instances of data analyses which involve multiple, linked sets of clustering problems, for which classical clustering methods (model-based or non-model-based) provide little in the way of leverage. In bioinformatics we have already alluded to the problem of finding haplotype structure in subpopulations. Other examples in bioinformatics include the use of HMMs for amino acid sequences, where a hierarchical DP version of the HMM would allow motifs to be discovered and shared among different families of proteins. In speech recognition multiple HMMs are

already widely used, in the form of word-specific and speaker-specific models, and adhoc methods are generally used to share statistical strength among models. We have discussed examples of grouped data in information retrieval; other examples include problems in which groups are indexed by author or by language. Finally, computational vision and robotics problems often involve sets of descriptors or objects that are arranged in a taxonomy. Examples such as these, in which there is substantial uncertainty regarding appropriate numbers of clusters, and in which the sharing of statistical strength among groups is natural and desirable, suggest that the hierarchical nonparametric Bayesian approach to clustering presented here may provide a generally useful extension of model-based clustering.

## A Posterior sampling for concentration parameters

MCMC samples from the posterior distributions for the concentration parameters  $\gamma$  and  $\alpha_0$  of the hierarchical Dirichlet process can be obtained using straightforward extensions of analogous techniques for Dirichlet processes. Let the number of observed groups be equal to  $J$ , with  $n_{j\cdot}$  observations in the  $j^{\text{th}}$  group. Consider the Chinese restaurant franchise representation. The concentration parameter  $\alpha_0$  governs the distribution of the number of  $\psi_{jt}$ 's in each mixture. As noted in Section 5.3 this is given by:

$$p(m_1, \dots, m_J | \alpha_0, n_{1\cdot}, \dots, n_{J\cdot}) = \prod_{j=1}^J s(n_{j\cdot}, m_{j\cdot}) \alpha_0^{m_{j\cdot}} \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_0 + n_{j\cdot})}. \quad (44)$$

Further,  $\alpha_0$  does not govern other aspects of the joint distribution, hence (44) along with the prior for  $\alpha_0$  is sufficient to derive MCMC updates for  $\alpha_0$  given all other variables.

In the case of a single mixture model ( $J = 1$ ), Escobar and West (1995) proposed a gamma prior and derived an auxiliary variable update for  $\alpha_0$ , while Rasmussen (2000) observed that (44) is log-concave in  $\log(\alpha_0)$  and proposed using adaptive rejection sampling instead. The adaptive rejection sampler of Rasmussen (2000) can be directly applied to the case  $J > 1$  since the conditional distribution of  $\log(\alpha_0)$  is still log-concave. The auxiliary variable method of Escobar and West (1995) requires a slight modification for the case  $J > 1$ . Assume that the prior for  $\alpha_0$  is a gamma distribution with parameters  $a$  and  $b$ . For each  $j$  we can write

$$\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_0 + n_{j\cdot})} = \frac{1}{\Gamma(n_{j\cdot})} \int_0^1 w_j^{\alpha_0} (1 - w_j)^{n_{j\cdot} - 1} \left(1 + \frac{n_{j\cdot}}{\alpha_0}\right) dw_j. \quad (45)$$

We define auxiliary variables  $\mathbf{w} = (w_j)_{j=1}^J$  and  $\mathbf{s} = (s_j)_{j=1}^J$  where each  $w_j$  is a variable taking on values in  $[0, 1]$ , and each  $s_j$  is a binary  $\{0, 1\}$  variable, and define the following distribution:

$$q(\alpha_0, \mathbf{w}, \mathbf{s}) \propto \alpha_0^{a-1+m_{\cdot\cdot}} e^{-\alpha_0 b} \prod_{j=1}^J w_j^{\alpha_0} (1 - w_j)^{n_{j\cdot} - 1} \left(\frac{n_{j\cdot}}{\alpha_0}\right)^{s_j}. \quad (46)$$

Now marginalizing  $q$  to  $\alpha_0$  gives the desired conditional distribution for  $\alpha_0$ . Hence  $q$  defines an auxiliary variable sampling scheme for  $\alpha_0$ . Given  $\mathbf{w}$  and  $\mathbf{s}$  we have:

$$q(\alpha_0 | \mathbf{w}, \mathbf{s}) \propto \alpha_0^{a-1+m_{\cdot\cdot} - \sum_{j=1}^J s_j} e^{-\alpha_0 (b - \sum_{j=1}^J \log w_j)}, \quad (47)$$

which is a gamma distribution with parameters  $a + m_{..} - \sum_{j=1}^J s_j$  and  $b - \sum_{j=1}^J \log w_j$ . Given  $\alpha_0$ , the  $w_j$  and  $s_j$  are conditionally independent, with distributions:

$$q(w_j|\alpha_0) \propto w_j^{\alpha_0} (1 - w_j)^{n_{j..}-1} \quad (48)$$

$$q(s_j|\alpha_0) \propto \left(\frac{n_{j..}}{\alpha_0}\right)^{s_j}, \quad (49)$$

which are beta and binomial distributions respectively. This completes the auxiliary variable sampling scheme for  $\alpha_0$ . We prefer the auxiliary variable sampling scheme as it is easier to implement and typically mixes quickly (within 20 iterations).

Given the total number  $m_{..}$  of  $\psi_{jt}$ 's, the concentration parameter  $\gamma$  governs the distribution over the number of components  $K$ :

$$p(K|\gamma, m_{..}) = s(m_{..}, K) \gamma^K \frac{\Gamma(\gamma)}{\Gamma(\gamma + m_{..})}. \quad (50)$$

Again other variables are independent of  $\gamma$  given  $m_{..}$  and  $K$ , hence we may apply the techniques of Escobar and West (1995) or Rasmussen (2000) directly to sampling  $\gamma$ .

## Acknowledgments

We wish to acknowledge helpful discussions with Lancelot James. We also wish to acknowledge support from Intel Corporation, Microsoft Research, and a grant from Darpa in support of the CALO program.

## References

- Aldous, D. (1985), "Exchangeability and Related Topics," in *École d'Été de Probabilités de Saint-Flour XIII–1983*, Springer, Berlin, pp. 1–198.
- Antoniak, C. (1974), "Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems," *Annals of Statistics*, 2(6), pp. 1152–1174.
- Beal, M. (2003), "Variational Algorithms for Approximate Bayesian Inference," Ph.D. thesis, Gatsby Computational Neuroscience Unit, University College London.
- Beal, M. J., Ghahramani, Z., and Rasmussen, C. (2002), "The Infinite Hidden Markov Model," in T. G. Dietterich, S. Becker, and Z. Ghahramani (eds.) *Advances in Neural Information Processing Systems*, Cambridge, MA: MIT Press, vol. 14, pp. 577–584.
- Blackwell, D. and MacQueen, J. (1973), "Ferguson Distributions via Pólya Urn Schemes," *Annals of Statistics*, 1, pp. 353–355.
- Blei, D., Jordan, M., and Ng, A. (2003), "Hierarchical Bayesian Models for Applications in Information Retrieval," in *Bayesian Statistics*, vol. 7, pp. 25–44.
- Blei, D. M. and Jordan, M. I. (2005), "Variational methods for Dirichlet process mixtures," *Bayesian Analysis*, 1, pp. 121–144.

- Carota, C. and Parmigiani, G. (2002), “Semiparametric Regression for Count Data,” *Biometrika*, 89(2), pp. 265–281.
- Cifarelli, D. and Regazzini, E. (1978), “Problemi Statistici Non Parametrici in Condizioni di Scambiabilità Parziale e Impiego di Medie Associate,” Tech. rep., Quaderni Istituto Matematica Finanziaria dell’Università di Torino.
- De Iorio, M., Müller, P., and Rosner, G. (2004), “An ANOVA Model for Dependent Random Measures,” *Journal of the American Statistical Association*, 99(465), pp. 205–215.
- Durbin, R., Eddy, S., Krogh, A., and Mitchison, G. (1998), *Biological Sequence Analysis*, Cambridge, UK: Cambridge University Press.
- Escobar, M. and West, M. (1995), “Bayesian Density Estimation and Inference Using Mixtures,” *Journal of the American Statistical Association*, 90, pp. 577–588.
- Ferguson, T. (1973), “A Bayesian Analysis of Some Nonparametric Problems,” *Annals of Statistics*, 1(2), pp. 209–230.
- Fong, D., Pammer, S., Arnold, S., and Bolton, G. (2002), “Reanalyzing Ultimatum Bargaining—Comparing Nondecreasing Curves Without Shape Constraints,” *Journal of Business and Economic Statistics*, 20, pp. 423–440.
- Forsyth, D. A. and Ponce, J. (2002), *Computer Vision—A Modern Approach*, Upper Saddle River, NJ: Prentice-Hall.
- Fraley, C. and Raftery, A. E. (2002), “Model-Based Clustering, Discriminant Analysis, and Density Estimation,” *Journal of the American Statistical Association*, 97, pp. 611–631.
- Gabriel, S. et al. (2002), “The Structure of Haplotype Blocks in the Human Genome,” *Science*, 296, pp. 2225–2229.
- Green, P. and Richardson, S. (2001), “Modelling Heterogeneity with and without the Dirichlet Process,” *Scandinavian Journal of Statistics*, 28, pp. 355–377.
- Huang, X., Acero, A., and Hon, H.-W. (2001), *Spoken Language Processing*, Upper Saddle River, NJ: Prentice-Hall.
- Ishwaran, H. and James, L. (2004), “Computational Methods for Multiplicative Intensity Models using Weighted Gamma Processes: Proportional Hazards, Marked Point Processes and Panel Count Data,” *Journal of the American Statistical Association*, 99, pp. 175–190.
- Ishwaran, H. and Zarepour, M. (2002), “Exact and Approximate Sum-Representations for the Dirichlet Process,” *Canadian Journal of Statistics*, 30, pp. 269–283.
- Jain, S. and Neal, R. (2000), “A Split-Merge Markov Chain Monte Carlo Procedure for the Dirichlet Process Mixture Model,” Tech. Rep. 2003, Department of Statistics, University of Toronto.
- Kleinman, K. and Ibrahim, J. (1998), “A Semi-parametric Bayesian Approach to Generalized Linear Mixed Models,” *Statistics in Medicine*, 17, pp. 2579–2596.
- MacEachern, S. (1999), “Dependent Nonparametric Processes,” in *Proceedings of the Section on Bayesian Statistical Science*, American Statistical Association.

- MacEachern, S., Kottas, A., and Gelfand, A. (2001), “Spatial Nonparametric Bayesian Models,” Tech. Rep. 01-10, Institute of Statistics and Decision Sciences, Duke University.
- MacEachern, S. and Müller, P. (1998), “Estimating Mixture of Dirichlet Process Models,” *Journal of Computational and Graphical Statistics*, 7, pp. 223–238.
- MacKay, D. and Peto, L. (1994), “A Hierarchical Dirichlet Language Model,” *Natural Language Engineering*.
- MacKay, D. J. C. (1997), “Ensemble Learning for Hidden Markov Models,” Tech. rep., Cavendish Laboratory, University of Cambridge.
- Mallick, B. and Walker, S. (1997), “Combining Information from Several Experiments with Nonparametric Priors,” *Biometrika*, 84, pp. 697–706.
- Muliere, P. and Petrone, S. (1993), “A Bayesian Predictive Approach to Sequential Search for an Optimal Dose: Parametric and Nonparametric Models,” *Journal of the Italian Statistical Society*, 2, pp. 349–364.
- Müller, P., Quintana, F., and Rosner, G. (2004), “A Method for Combining Inference Across Related Nonparametric Bayesian Models,” *Journal of the Royal Statistical Society*, 66, pp. 735–749.
- Neal, R. (1992), “Bayesian Mixture Modeling,” in *Proceedings of the Workshop on Maximum Entropy and Bayesian Methods of Statistical Analysis*, vol. 11, pp. 197–211.
- (2000), “Markov Chain Sampling Methods for Dirichlet Process Mixture Models,” *Journal of Computational and Graphical Statistics*, 9, pp. 249–265.
- Pitman, J. (1996), “Random discrete distributions invariant under size-biased permutation,” *Advances in Applied Probability*, 28, pp. 525–539.
- (2002a), “Combinatorial Stochastic Processes,” Tech. Rep. 621, Department of Statistics, University of California at Berkeley, lecture notes for St. Flour Summer School.
- (2002b), “Poisson-Dirichlet and GEM invariant distributions for split-and-merge transformations of an interval partition,” *Combinatorics, Probability and Computing*, 11, pp. 501–514.
- Pritchard, J., Stephens, M., and Donnelly, P. (2000), “Inference of Population Structure using Multilocus Genotype Data,” *Genetics*, 155, pp. 945–959.
- Rabiner, L. (1989), “A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition,” *Proceedings of the IEEE*, 77, pp. 257–285.
- Rasmussen, C. (2000), “The Infinite Gaussian Mixture Model,” in S. Solla, T. Leen, and K.-R. Müller (eds.) *Advances in Neural Information Processing Systems*, Cambridge, MA: MIT Press, vol. 12.
- Salton, G. and McGill, M. (1983), *An Introduction to Modern Information Retrieval*, New York: McGraw-Hill.
- Sethuraman, J. (1994), “A Constructive Definition of Dirichlet Priors,” *Statistica Sinica*, 4, pp. 639–650.

- Stephens, M., Smith, N., and Donnelly, P. (2001), "A New Statistical Method for Haplotype Reconstruction from Population Data," *American Journal of Human Genetics*, 68, pp. 978–989.
- Stolcke, A. and Omohundro, S. (1993), "Hidden Markov Model Induction by Bayesian Model Merging," in C. Giles, S. Hanson, and J. Cowan (eds.) *Advances in Neural Information Processing Systems*, San Mateo CA: Morgan Kaufmann, vol. 5, pp. 11–18.
- Tomlinson, G. (1998), "Analysis of Densities," Ph.D. thesis, Department of Public Health Sciences, University of Toronto.
- Tomlinson, G. and Escobar, M. (2003), "Analysis of Densities," Tech. rep., Department of Public Health Sciences, University of Toronto.