



Hierarchical Dirichlet Processes

Yee Whye Teh, Michael I Jordan, Matthew J Beal & David M Blei

To cite this article: Yee Whye Teh, Michael I Jordan, Matthew J Beal & David M Blei (2006) Hierarchical Dirichlet Processes, Journal of the American Statistical Association, 101:476, 1566-1581, DOI: [10.1198/016214506000000302](https://doi.org/10.1198/016214506000000302)

To link to this article: <https://doi.org/10.1198/016214506000000302>



Published online: 01 Jan 2012.



Submit your article to this journal [↗](#)



Article views: 7365



View related articles [↗](#)



Citing articles: 191 View citing articles [↗](#)

Hierarchical Dirichlet Processes

Yee Whye TEH, Michael I. JORDAN, Matthew J. BEAL, and David M. BLEI

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 modeling.

KEY WORDS: Clustering; Hierarchical model; Markov chain Monte Carlo; Mixture model; Nonparametric Bayesian statistics.

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 through 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 this article 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 sharing of clusters 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 (SNPs)] in a localized region of the human genome. Although 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*—is 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, Smith, and Donnelly 2001). Now consider an extension of this problem in which the population is divided into a set of groups, such as, African, Asian, and European subpopulations. We not only may want to discover the sets of haplotypes within each subpopulation, but also may 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, Jordan, and Ng 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, then 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 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 $G \sim DP(\alpha_0, G_0)$, then, with probability 1,

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

Yee Whye Teh is Lee Kuan Yew Postdoctoral Fellow, Department of Computer Science, National University of Singapore, Singapore (E-mail: tehyw@comp.nus.edu.sg). Michael I. Jordan is Professor of Electrical Engineering and Computer Science and Professor of Statistics, University of California, Berkeley, CA 94720 (E-mail: jordan@eecs.berkeley.edu). Matthew J. Beal is Assistant Professor of Computer Science and Engineering, SUNY Buffalo, Buffalo, NY 14260 (E-mail: mbeal@cse.buffalo.edu). David M. Blei is Assistant Professor of Computer Science, Princeton University, Princeton, NJ (E-mail: blei@eecs.berkeley.edu). Correspondence should be directed to Michael I. Jordan. This work was supported in part by Intel Corporation, Microsoft Research, and a grant from Darpa under contract number NBCHD030010. The authors wish to acknowledge helpful discussions with Lancelot James and Jim Pitman, and to thank the referees for useful comments.

where the ϕ_k are independent random variables distributed according to G_0 , where δ_{ϕ_k} is an atom at ϕ_k and the “stick-breaking weights,” β_k , are also random and depend on the parameter α_0 . (The definition of the β_k is provided in Sec. 3.1.)

The representation in (1) shows that draws from a DP are discrete (with probability 1). 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 *DP 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 DP, $\text{DP}(\alpha_{0j}, G_{0j})$. To link these clustering problems, we link the group-specific DPs. Many authors have considered ways to induce dependencies among multiple DPs through links among the parameters G_{0j} and/or α_{0j} (Cifarelli and Regazzini 1978; MacEachern 1999; Tomlinson 1998; Müller, Quintana, and Rosner 2004; De Iorio, Müller, and Rosner 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 DP, $\text{DP}(\alpha_0, G_0(\tau))$, where $G_0(\tau)$ is a parametric distribution with random parameter τ (Carota and Parmigiani 2002; Fong, Pammer, Arnold, and Bolton 2002; Muliere and Petrone 1993). Integrating over τ 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 τ (e.g., G_0 is Gaussian with mean τ). 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 1). Thus, although clusters arise *within* each group through 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 DP, $\text{DP}(\gamma, H)$. This restores flexibility in that the modeler can choose H to be continuous or discrete. In either case, with probability 1, G_0 is discrete and has a stick-breaking representation as in (1). The atoms ϕ_k are shared among the multiple DPs, yielding the desired sharing of atoms among groups. In summary, we consider the hierarchical specification

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

which we refer to as a *hierarchical DP*. The immediate extension to *hierarchical DP 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 DPs” due to MacEachern (1999) and MacEachern, Kottas, and Gelfand (2001). In this framework the random variables β_k and ϕ_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 β_k in (1) with a second subscript indexing the groups j and view the weights β_{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 β_{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, we treat it 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. But it is an appropriate choice for the problem addressed by Tomlinson (1998), that of sharing statistical strength among multiple sets of density estimation problems. Thus, whereas the AnDe framework and our hierarchical DP framework are closely related formally, the inferential goal is rather different. Moreover, as we show later, our restriction to discrete G_0 has important implications for the design of efficient Markov chain Monte Carlo (MCMC) inference algorithms.

The terminology of “hierarchical DP” has also been used by Müller et al. (2004) to describe a different notion of hierarchy than that discussed here. These authors considered 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, in which the shared clusters are given the same stick-breaking weights (those associated with F_0) in each of the groups. In 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 term “hierarchical DP” has been used in yet a third way by Beal, Ghahramani, and Rasmussen (2002) in the context of a model known as the *infinite hidden Markov model*, a hidden Markov model with a countably infinite state space. But the “hierarchical DP” of Beal et al. (2002) is 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 DP that we explore is a specific example of a dependency model for multiple DPs, 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 DPs is itself distributed according to a DP. Although there are many ways to couple DPs, we view this simple, canonical Bayesian hierarchy as particularly worthy of study. Note in particular the appealing recursiveness of the definition; a hierarchical DP 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 consider the flow of haplotypes on the resulting tree.

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

2. SETTING

We are interested in problems in which the observations are organized into *groups* and assumed to be 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 that each observation is drawn independently from a mixture model, there is a mixture component associated with each observation. Let θ_{ji} denote a parameter specifying the mixture component associated with the observation x_{ji} . We refer to the variables θ_{ji} as *factors*. Note that these variables are not generally distinct; we develop a different notation for the distinct values of factors. Let $F(\theta_{ji})$ denote the distribution of x_{ji} given the factor θ_{ji} . Let G_j denote a prior distribution for the factors $\theta_j = (\theta_{j1}, \theta_{j2}, \dots)$ associated with group j . We assume that the factors are conditionally independent given G_j . Thus we have the following probability model:

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

to augment the specification given in (2).

3. DIRICHLET PROCESSES

In this section we provide a brief overview of DPs. After a discussion of basic definitions, we present three different perspectives on the DP: 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 DP, as described in Section 4.

Let (Θ, \mathcal{B}) be a measurable space, with G_0 a probability measure on the space. Let α_0 be a positive real number. A *Dirichlet process*, $\text{DP}(\alpha_0, G_0)$, is defined as the distribution of a random probability measure G over (Θ, \mathcal{B}) such that, for any finite measurable partition (A_1, A_2, \dots, A_r) of Θ , the random vector $(G(A_1), \dots, G(A_r))$ is distributed as a finite-dimensional Dirichlet distribution with parameters $(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r))$,

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

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

3.1 The Stick-Breaking Construction

Measures drawn from a DP are discrete with probability 1 (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 iid random variables $(\pi'_k)_{k=1}^\infty$ and $(\phi_k)_{k=1}^\infty$,

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

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

3.2 The Chinese Restaurant Process

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

The Pólya urn scheme does not refer to G directly; rather, it refers to draws from G . Thus let $\theta_1, \theta_2, \dots$ be a sequence of iid random variables distributed according to G . That is, the variables $\theta_1, \theta_2, \dots$ are conditionally independent given G , and hence are exchangeable. Let us consider the successive conditional distributions of θ_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:

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

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 α_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 θ_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 ϕ_1, \dots, ϕ_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 ϕ_k for $1 \leq i' < i$. We can re-express (7) as

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

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 DP (Pitman 2002a), and it is useful in this article as well. The metaphor is as follows. Consider a Chinese restaurant with an unbounded number of tables. Each θ_i corresponds to a customer who enters the restaurant, whereas the distinct values ϕ_k correspond to the tables at which the customers sit. The i th customer sits at the table indexed by ϕ_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 α_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 DP is as a non-parametric prior on the parameters of a mixture model. In particular, suppose that observations x_i arise as

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

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

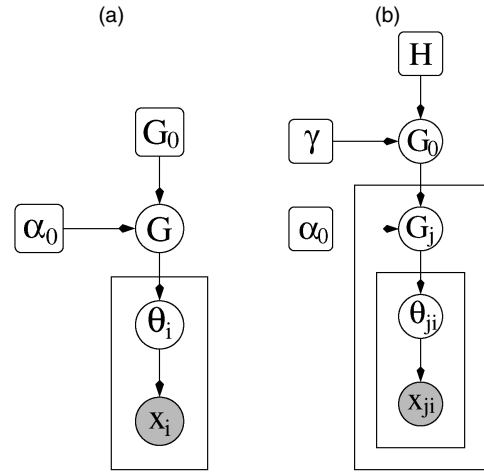


Figure 1. Graphical Model Representation of a DP Mixture Model (a) and a Hierarchical DP Mixture Model (b). In the graphical model formalism, each node in the graph is associated with a random variable, where shading denotes an observed variable. Rectangles denote replication of the model within the rectangle. Sometimes the number of replicates is given in the bottom right corner of the rectangle.

Because G can be represented using a stick-breaking construction (6), the factors θ_i take on values ϕ_k with probability π_k . We may denote this using an indicator variable, z_i , that takes on positive integral values and is distributed according to π (interpreting π as a random probability measure on the positive integers). Hence an equivalent representation of a DP mixture is given by the following conditional distributions:

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

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 DP 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 DP.

Suppose that we have L mixture components. Let $\pi = (\pi_1, \dots, \pi_L)$ denote the mixing proportions. Note that we previously used the symbol π to denote the weights associated with the atoms in G . We have deliberately overloaded the definition of π here; as we show 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 π with symmetric parameters $(\alpha_0/L, \dots, \alpha_0/L)$. Let ϕ_k denote the parameter vector associated with mixture component k , and let ϕ_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:

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

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

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

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

4. HIERARCHICAL DIRICHLET PROCESSES

We propose a nonparametric Bayesian approach to the modeling of grouped data, in which each group is associated with a mixture model and we wish to link these mixture models. By analogy with DP mixture models, we first define the appropriate nonparametric prior, which we call the *hierarchical DP*. 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 DP: a stick-breaking construction, a Chinese restaurant process representation, and a representation in terms of a limit of finite mixture models.**

A hierarchical DP is a distribution over a set of random probability measures over (Θ, \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 DP with concentration parameter γ and base probability measure H ,

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

and the random measures G_j are conditionally independent given G_0 , with distributions given by a DP with base probability measure G_0 ,

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

The hyperparameters of the hierarchical DP consist of the baseline probability measure H , and the concentration parameters γ and α_0 . The baseline H provides the prior distribution for the factors θ_{ji} . The distribution G_0 varies around the prior H , with the amount of variability governed by γ . The actual distribution G_j over the factors in the j th group deviates from G_0 , with the amount of variability governed by α_0 . If we expect the variability in different groups to be different, then we can use a separate concentration parameter α_j for each group j . In this article, following Escobar and West (1995), we put vague gamma priors on γ and α_0 .

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

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

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

The hierarchical DP 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 DP, 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. Because 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 thus can 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 π_j are independent given β , because the G_j 's are independent given G_0 . We now describe how the weights π_j are related to the global weights β .

Let (A_1, \dots, A_r) be a measurable partition of Θ 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 nonatomic, the ϕ_k 's are distinct with probability 1, and so any partition of the positive integers corresponds to some partition of Θ . 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 π_j is independently distributed according to $\text{DP}(\alpha_0, \beta)$, where we interpret β and π_j as probability measures on the positive integers. If H is nonatomic, 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 DP mixture model, because each factor θ_{ji} is distributed according to G_j , it takes on the value ϕ_k with probability π_{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 DP mixture through the following conditional distributions:

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

We now derive an explicit relationship between the elements of β and π_j . Recall that the stick-breaking construction for DPs defines the variables β_k in (16) as

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

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

$$\pi'_{jk} \sim \text{beta}\left(\alpha_0 \beta_k, \alpha_0 \left(1 - \sum_{l=1}^k \beta_l\right)\right),$$

$$\pi_{jk} = \pi'_{jk} \prod_{l=1}^{k-1} (1 - \pi'_{jl}). \quad (21)$$

To derive (21), first note 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 $\pi'_{jk} = \pi_{jk} / (1 - \sum_{l=1}^{k-1} \pi_{jl})$ 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 DPs.

4.2 The Chinese Restaurant Franchise

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

The metaphor is as follows (see Fig. 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 this dish is shared among all of the 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 θ_{ji} . We also let ϕ_1, \dots, ϕ_K denote K iid random variables distributed according to H ; this is the global menu of dishes. We also introduce variables, ψ_{jt} , that represent the table-specific choice of dishes; in particular, ψ_{jt} is the dish served at table t in restaurant j .

Note that each θ_{ji} is associated with one ψ_{jt} , whereas each ψ_{jt} is associated with one ϕ_k . We introduce indicators to denote these associations. In particular, let t_{ji} be the index of the ψ_{jt} associated with θ_{ji} , and let k_{jt} be the index of ϕ_k associated with ψ_{jt} . In the Chinese restaurant franchise metaphor, customer i in restaurant j sits at table t_{ji} whereas table t in restaurant j serves dish k_{jt} .

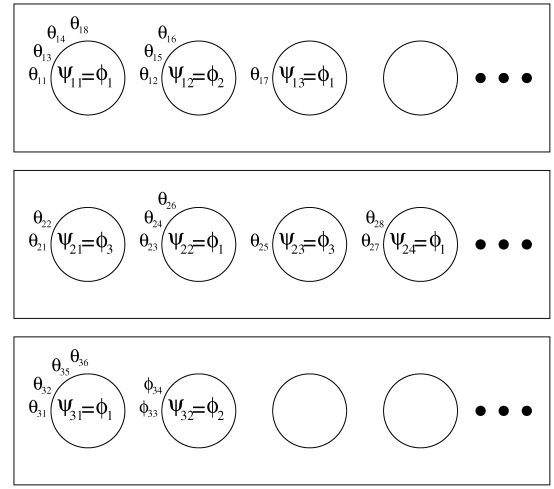


Figure 2. A Depiction of a Chinese Restaurant Franchise. Each restaurant is represented by a rectangle. Customers (θ_{ji} 's) are seated at tables (circles) in the restaurants. At each table a dish is served. The dish is served from a global menu (ϕ_k), whereas the parameter ψ_{jt} is a table-specific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

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 $m_{\cdot\cdot}$ represents the total number of tables occupied.

Let us now compute marginals under a hierarchical DP when G_0 and G_j are integrated out. First, consider the conditional distribution for θ_{ji} given $\theta_{j1}, \dots, \theta_{j,i-1}$ and G_0 , where G_j is integrated out. From (8),

$$\theta_{ji} | \theta_{j1}, \dots, \theta_{j,i-1}, \alpha_0, G_0$$

$$\sim \sum_{t=1}^{m_j} \frac{n_{jt}}{i-1+\alpha_0} \delta_{\psi_{jt}} + \frac{\alpha_0}{i-1+\alpha_0} G_0. \quad (24)$$

This is a mixture, a draw from which can be obtained by drawing from the terms on the right side with probabilities given by the corresponding mixing proportions. If a term in the first summation is chosen, then we increment n_{jt} , 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 . Note that G_0 appears only in its role as the distribution of the variables ψ_{jt} . Because G_0 is distributed according to a DP, we can integrate it out by using (8) again and write the conditional distribution of ψ_{jt} as

$$\psi_{jt} | \psi_{11}, \psi_{12}, \dots, \psi_{21}, \dots, \psi_{j,t-1}, \gamma, H$$

$$\sim \sum_{k=1}^K \frac{m_{\cdot k}}{m_{\cdot\cdot} + \gamma} \delta_{\phi_k} + \frac{\gamma}{m_{\cdot\cdot} + \gamma} H. \quad (25)$$

If we draw ψ_{jt} by choosing a term in the summation on the right side of this equation, then we set $\psi_{jt} = \phi_k$ and let $k_{jt} = k$ for the chosen k . If we choose the second term, 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 θ_{ji} variables. To use these equations to obtain samples of θ_{ji} , we proceed as follows. For each j and i , we first sample θ_{ji} using (24). If a new sample from G_0 is needed, then we use (25) to obtain a new sample ψ_{jt} and set $\theta_{ji} = \psi_{jt}$.

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

4.3 The Infinite Limit of Finite Mixture Models

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

Consider the following collection of finite mixture models, where β is a global vector of mixing proportions and π_j is a group-specific vector of mixing proportions:

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

The parametric hierarchical prior for β and π in (26) has been discussed by MacKay and Peto (1994) as a model for natural languages. We show that the limit of this model as $L \rightarrow \infty$ is the hierarchical DP. 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

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

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 that this finite model induces on \mathbf{x} approaches the hierarchical DP mixture model.

There is an alternative finite model whose limit is also the hierarchical DP mixture model. Instead of introducing dependencies between the groups by placing a prior on β (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:

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

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 DP mixture model.

5. INFERENCE

In this section we describe three related MCMC sampling schemes for the hierarchical DP mixture model. The first scheme is a straightforward Gibbs sampler based on the Chinese restaurant franchise; the second is based on an augmented representation involving both the Chinese restaurant franchise and the posterior for G_0 ; and 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 DP. The nonconjugate case can be approached by adapting to the hierarchical DP techniques developed for nonconjugate DP mixtures (Neal 2000). Moreover, in this section we assume fixed values for the concentration parameters α_0 and γ ; 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 θ_{ji} be associated with the table t_{ji} in the restaurant representation, that is, let $\theta_{ji} = \psi_{jt_{ji}}$. The random variable ψ_{jt} is an instance of mixture component k_{jt} , that is, $\psi_{jt} = \phi_{k_{jt}}$. The prior over the parameters ϕ_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 , m_{jk} to denote the number of tables in restaurant j serving dish k , and K to denote the number of dishes being served throughout the franchise. 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 ψ_{jt} , leaving out item x_{ji} .

Let $F(\theta)$ have density $f(\cdot | \theta)$ and let H have density $h(\cdot)$. Because H is conjugate to F , we integrate out the mixture component parameters ϕ 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 let $f_k^{-\mathbf{x}_{jt}}(\mathbf{x}_{jt})$ denote the conditional density of \mathbf{x}_{jt} given all data items associated with mixture component k leaving out \mathbf{x}_{jt} .

Finally, we suppress references to all variables except those being sampled in the conditional distributions to follow. In particular, we omit references to \mathbf{x} , α_0 , and γ .

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 θ_{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 θ_{ji} 's and ψ_{jt} 's directly, we instead sample their index variables t_{ji} and k_{jt} . The θ_{ji} 's and ψ_{jt} 's can be reconstructed from these index variables and the ϕ_k 's. This representation makes the MCMC sampling scheme more efficient (cf. Neal 2000). Note that the t_{ji} and the k_{jt} inherit the exchangeability properties of the θ_{ji} and the ψ_{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} . Note 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 a finitely number of which are actually associated with the data and represented explicitly.

Sampling \mathbf{t} . To compute the conditional distribution of t_{ji} given the rest 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}^{-ji} , whereas the probability that it takes on a new value (say, $t^{\text{new}} = m_j + 1$) is proportional to α_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}^{new} using (25),

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

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

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

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

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

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 0, because this is always proportional to n_{jt} . Consequently, 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, then we delete this mixture component as well.

Sampling \mathbf{k} . Because 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 $f_k^{-x_{jt}}(\mathbf{x}_{jt})$, so that the conditional probability of k_{jt} is

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

5.2 Posterior Sampling With an Augmented Representation

In the Chinese restaurant franchise sampling scheme, the sampling for all groups is coupled because G_0 is integrated out. This complicates matters in more elaborate models (e.g., in the case of the hidden Markov model considered in Sec. 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 that ψ_{jt} for each table t is a draw from G_0 . Conditioning on the ψ_{jt} 's, G_0 is now distributed as $\text{DP}(\gamma + m_{..}, (\gamma H + \sum_{k=1}^K m_k \delta_{\phi_k})/(\gamma + m_{..}))$. An explicit construction for G_0 is now given as

$$\beta = (\beta_1, \dots, \beta_K, \beta_u) \sim \text{Dir}(m_{.1}, \dots, m_{.K}, \gamma),$$

$$G_u \sim \text{DP}(\gamma, H),$$

$$p(\phi_k|\mathbf{t}, \mathbf{k}) \propto h(\phi_k) \prod_{ji: k_{jt_{ji}}=k} f(x_{ji}|\phi_k), \quad (35)$$

$$G_0 = \sum_{k=1}^K \beta_k \delta_{\phi_k} + \beta_u G_u.$$

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 β earlier, whereas both ϕ and G_u are integrated out (which introduces couplings into the sampling for each group but is easily handled).

Sampling \mathbf{t} and \mathbf{k} . This is almost identical to the Chinese restaurant franchise sampling scheme, with the only novelty being that we replace m_k by β_k and γ by β_u in (31), (32), (33), and (34), and when a new component k^{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 DP (Pitman 1996), the weight b corresponding to the chosen atom in G_u will have the same distribution as the first stick-breaking weight, that is, $\text{beta}(1, \gamma)$.

Sampling β . This has already been described in (35):

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

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 with 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 through a variable, z_{ji} , which is equivalent to k_{jti} in the earlier sampling schemes. The tables are represented only in terms of the numbers of tables m_{jk} .

Sampling \mathbf{z} . This can be realized by grouping together terms associated with each k in (31) and (32),

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

where we have replaced $m_{j,k}$ with β_k and γ with β_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 through \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 $k_{jti} = z_{ji}$. The probability that data item x_{ji} is assigned to some table t such that $k_{jt} = k$ is

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

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

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

These equations form the conditional distributions of a Gibbs sampler whose equilibrium distribution is the prior distribution over the assignment of $n_{j,k}$ observations to components in an ordinary DP 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

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

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 $s(n+1, m) = s(n, m-1) + ns(n, m)$.

Sampling $\boldsymbol{\beta}$. This 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 thus 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, whereas 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 DP 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 DP 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 DP 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 DPs, and we expect that they will prove useful for hierarchical DPs as well.

6. EXPERIMENTS

In this section we describe two experiments to highlight the two aspects of the hierarchical DP: 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.

The software that we used for these experiments is available at <http://www.cs.berkeley.edu/~jordan/hdp>. The software implements a hierarchy of DPs of arbitrary depth.

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 iid 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, Stephens, and Donnelly (2000) 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 DPs 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.

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 θ_{ji} from G_j and sampling words x_{ji} with probabilities θ_{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 5,838 abstracts in total. After removing standard stop words and words appearing fewer than 10 times, we are left with a total of 476,441 words. Following standard information retrieval methodology, the vocabulary is defined as the set of distinct words left in all abstracts; this has size 5,699.

Both models were as similar as possible beyond the distinction that LDA assumes a fixed finite number of topics, whereas the hierarchical DP does not. Both models used a symmetric Dirichlet distribution with parameters of .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 was assumed to be symmetric Dirichlet with parameters α_0/L , with L being the number of topics; γ was not used in LDA. Posterior samples were obtained using the Chinese restaurant franchise sampling scheme, whereas the concentration parameters were sampled using the auxiliary variable sampling scheme presented in the Appendix.

We evaluated the models through 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 as

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

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 from 10 to 120. As shown in Figure 3(a), 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(b), 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 DP to a third level. A draw from a

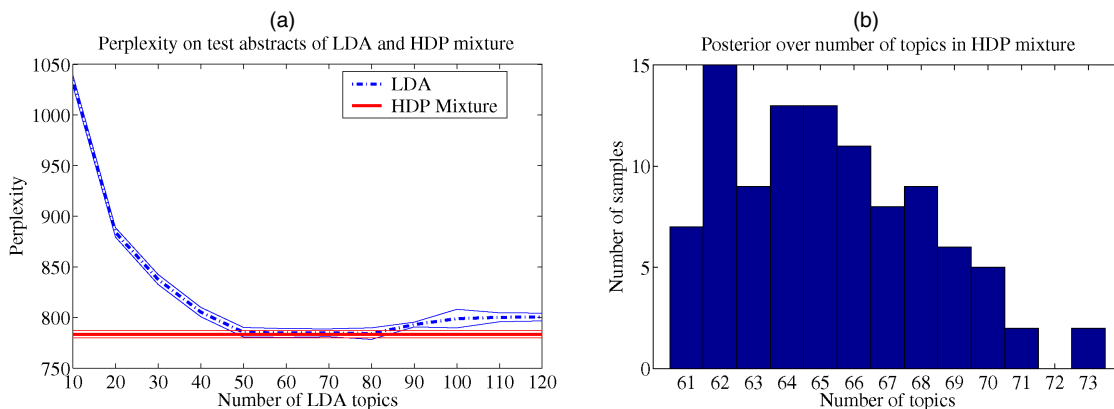


Figure 3. Results for Document Topic Modeling. (a) Comparison of LDA (---) and the HDP (—) Mixtures, With Results Averaged Over 10 Runs (error bars are one standard error); and (b) Histogram of the Number of Topics for the Hierarchical Dirichlet Process Mixture Over 100 Posterior Samples.

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 sharing of topics 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), and 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 1,447 articles in total. Each article was modeled as a bag of words. We culled standard stop words as well as words occurring more than 4,000 times or fewer than 50 times in the whole corpus. This left us with an average of slightly more than 1,000 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 later), 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 Fig. 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 $\gamma \sim \text{gamma}(5, .1)$ and $\alpha_0 \sim \text{gamma}(.1, .1)$ 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. We used priors of $\gamma \sim \text{gamma}(5, .1)$ and $\alpha_0 \sim \text{gamma}(.1, .1)$.
- M3. This model takes a full hierarchical approach and models the NIPS sections as multiple corpora, linked through 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. We used priors of $\gamma \sim \text{gamma}(5, .1)$, $\alpha_0 \sim \text{gamma}(5, .1)$, and $\alpha_1 \sim \text{gamma}(.1, .1)$.

In all models a finite and known vocabulary is assumed, and the base measure H used is a symmetric Dirichlet distribution with parameters of .5.

We conducted experiments in which a set of 80 articles was 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 was also chosen uniformly at random from the VS section, as was an additional set of 47 test articles distinct from the training articles. Our 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 with N varying between 0 and 80. The direct assignment sampling scheme is

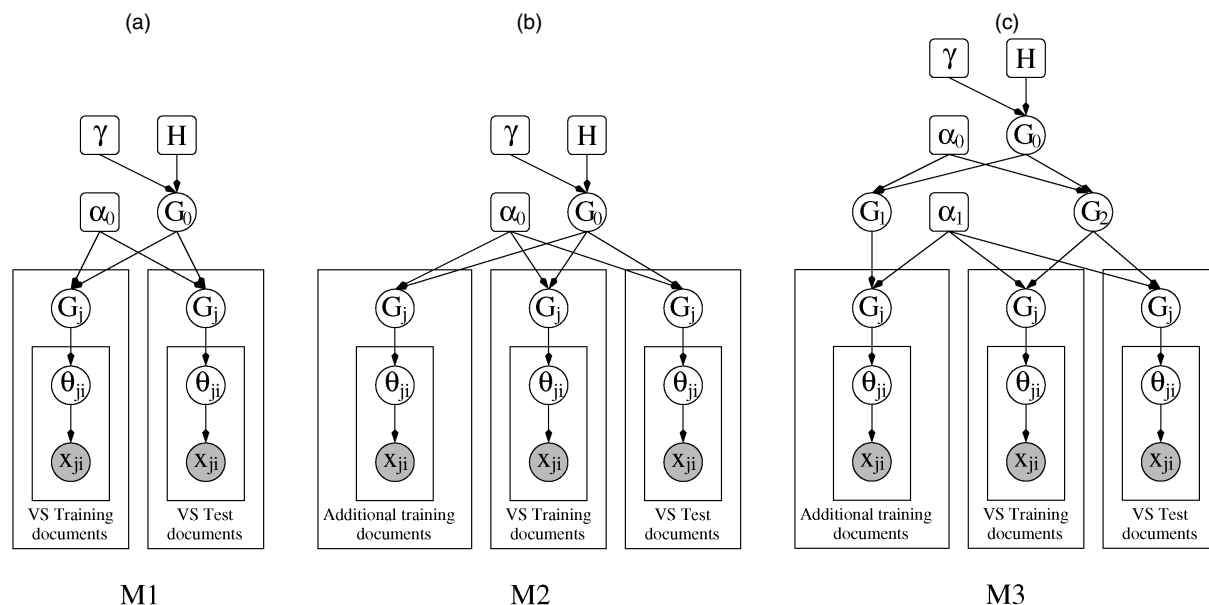


Figure 4. Three Models for the NIPS Data: (a) M1, (b) M2, and (c) M3.

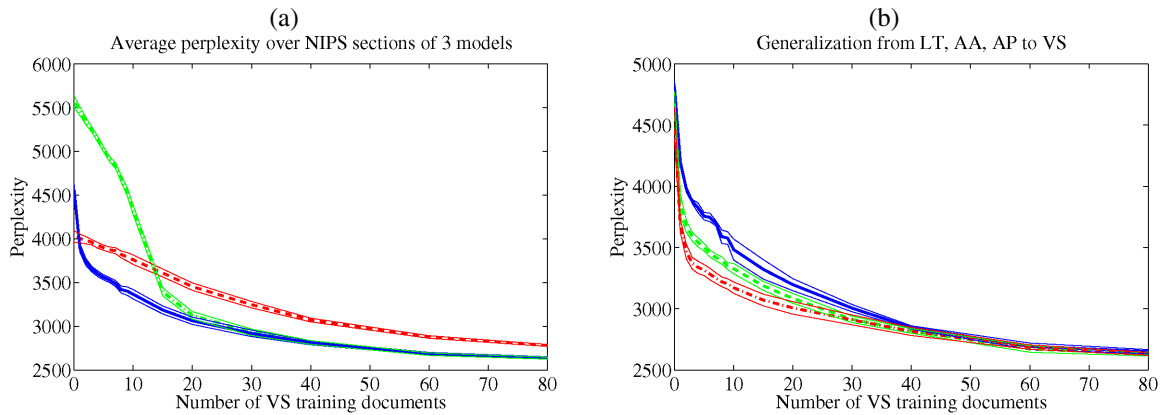


Figure 5. Results for Multi-Corpora Document Topic Modeling. (a) Perplexity of single words in test VS articles given training articles from VS and another section for three different models. Curves shown are averaged over the other sections and five runs (— M1: additional section ignored; - - - M2: flat, additional section; — M3: hierarchical, additional section). (b) Perplexity of test VS articles given LT, AA, and AP articles, using M3, averaged over five runs (— LT; - - - AA; - - - AP). In both plots, the error bars represent one standard error.

used, and concentration parameters are sampled using the auxiliary variable sampling scheme given in the Appendix.

Figure 5(a) presents the average predictive performance for all three models over five 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 the figure shows, 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 that of the fully hierarchical M3 throughout the range of N . M2 dominated M1 for small N , but yielded poorer performance than M1 for $N > 14$. 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 through the hierarchy.

Although the results in Figure 5(a) 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(b) plots the predictive performance for model M3 when given data from each of three particular sections: LT, AA, and AP. Whereas articles in the LT section are concerned mostly with theoretical properties of learning algorithms, those in AA are concerned mostly with models and methodology, and those in AP are concerned mostly with applications of learning algorithms to various problems. As the figure shows, predictive performance is enhanced the most by previous 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 that 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 (HMM).

Recall that an HMM is a doubly stochastic Markov chain in which a sequence of multinomial “state” variables (v_1, v_2, \dots, v_T) are linked through 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 one mixture component corresponds 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 DP.

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 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 non-parametric variant of the HMM that 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.

Table 1. Topics Shared Between VS and the Other NIPS Sections

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

NOTE: These topics are the most frequently occurring in the VS fit, under the constraint that they are associated with a significant number of words (>2,500) from the other section.

Thus, we can define a nonparametric HMM by simply replacing the set of conditional finite mixture models underlying the classical HMM with a hierarchical DP 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 use 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 HMM*, 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 DP” for this model, but their “hierarchical DP” was not hierarchical in the Bayesian sense—involving a distribution on the parameters of a DP—but was 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 that the same transition is observed at other time steps, whereas with probability proportional to α_0 , an “oracle” process is invoked. At this second level, the probability of transitioning to state v is proportional to the number of times that state v has been chosen by the oracle (regardless of the previous state), whereas the probability of transitioning to a novel state is proportional to γ . 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 DP.

To relate this two-level urn model to the hierarchical DP framework, we describe a representation of the HDP-HMM using the stick-breaking formalism. In particular, consider the hierarchical DP representation shown in Figure 6. The parameters

in this representation have the following distributions:

$$\begin{aligned}\beta|\gamma &\sim \text{GEM}(\gamma), \\ \pi_k|\alpha_0, \beta &\sim \text{DP}(\alpha_0, \beta), \\ \phi_k|H &\sim H,\end{aligned}\quad (42)$$

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

$$\begin{aligned}v_t|v_{t-1}, (\pi_k)_{k=1}^\infty &\sim \pi_{v_{t-1}} \quad \text{and} \\ y_t|v_t, (\phi_k)_{k=1}^\infty &\sim F(\phi_{v_t}),\end{aligned}\quad (43)$$

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, then it turns out that the result is equivalent to the coupled urn model of Beal et al. (2002), and hence the infinite HMM 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 HMM, proposing instead a heuristic approximation to Gibbs sampling. On the other hand, both the augmented representation and direct assignment representation lead directly to MCMC sampling schemes that can be implemented straightforwardly. In

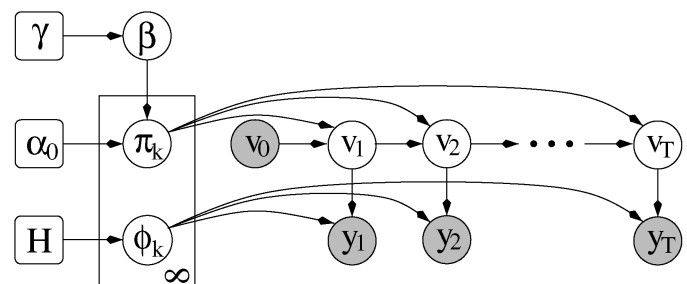


Figure 6. A Graphical Representation of an HDP-HMM.

the experiments reported in the following section, we used the direct assignment representation.

Practical applications of HMMs often consider sets of sequences and treat these sequences as exchangeable at the level of sequences. Thus, in applications to speech recognition, an HMM for a given word in the vocabulary is generally trained through 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 DP couple each of the HDP-HMMs, each of which is a set of DPs.

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 with 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, along with 40 test sentences with an average length of 100. The base distribution H is a symmetric Dirichlet distribution over 27 symbols with parameters .1. The concentration parameters γ and α_0 are given gamma(1, 1) priors.

Using the direct assignment sampling method for posterior predictive inference, we compared the HDP-HMM with various other methods for prediction using HMMs: (1) a classical HMM using maximum likelihood (ML) parameters obtained through 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 by 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(a). For VB, computing the predictive probability is intractable, so

we used the modal setting of parameters. Both the 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(b) shows posterior samples of the number of states used by the HDP-HMM.

8. DISCUSSION

In this article we have described a nonparametric approach to modeling groups of data in which each group is characterized by a mixture model and we allow sharing of mixture components between groups. We have proposed a hierarchical Bayesian solution to this problem, in which a set of DPs is coupled through their base measure, which is itself distributed according to a DP.

We have described three different representations that capture aspects of the hierarchical DP: a stick-breaking representation that describes the random measures explicitly, a representation of marginals in terms of an urn model that we call 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 MCMC sampling schemes for posterior inference under hierarchical DP 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, and the third scheme 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 often present when data are collected on a large scale. Clustering problems can be approached within a probabilistic framework through 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 HMMs in such areas as bioinformatics (Durbin, Eddy, Krogh, and Mitchison 1998), speech recognition (Huang, Acero, and Hon 2001), information retrieval (Blei et al. 2003), and computational vision (Forsyth and Ponce 2002). These areas also provide numerous instances of data analyses that involve multiple linked sets of clustering

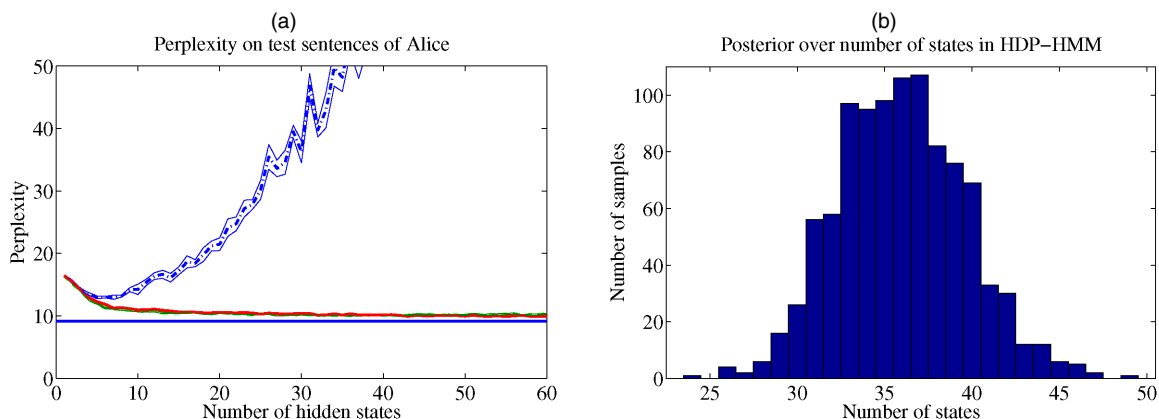


Figure 7. Results for HMMs. (a) Comparing the HDP-HMM (solid horizontal line) with ML (---), MAP (---), and VB (—) trained hidden Markov models. The error bars represent one standard error (those for the HDP-HMM are too small to see). (b) Histogram for the number of states in the HDP-HMM over 1,000 posterior samples.

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 the discovery of and sharing of motifs 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 ad hoc 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.

APPENDIX: POSTERIOR SAMPLING FOR CONCENTRATION PARAMETERS

MCMC samples from the posterior distributions for the concentration parameters γ and α_0 of the hierarchical DP can be obtained using straightforward extensions of analogous techniques for DP. Let the number of observed groups be equal to J , with $n_{j..}$ observations in the j th group. Consider the Chinese restaurant franchise representation. The concentration parameter α_0 governs the distribution of the number of ψ_{jt} 's in each mixture. As noted in Section 5.3, this is given by

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

Further, α_0 does not govern other aspects of the joint distribution; hence (A.1) along with the prior for α_0 is sufficient to derive MCMC updates for α_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 α_0 , and Rasmussen (2000) observed that (A.1) 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 where $J > 1$, because 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 where $J > 1$. Assume that the prior for α_0 is a gamma distribution with parameters a and b . For each j , we can write

$$\frac{\Gamma(\alpha_0)}{\Gamma(\alpha_0 + n_{j..})} = \frac{1}{\Gamma(n_{j..})} \int_0^1 w_j^{\alpha_0} (1 - w_j)^{n_{j..}-1} \left(1 + \frac{n_{j..}}{\alpha_0}\right) dw_j. \quad (\text{A.2})$$

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:

$$q(\alpha_0, \mathbf{w}, \mathbf{s}) \propto \alpha_0^{a-1+m_{..}} e^{-\alpha_0 b} \prod_{j=1}^J w_j^{\alpha_0} (1 - w_j)^{n_{j..}-1} \left(\frac{n_{j..}}{\alpha_0}\right)^{s_j}. \quad (\text{A.3})$$

Now marginalizing q to α_0 gives the desired conditional distribution for α_0 . Hence q defines an auxiliary variable sampling scheme for α_0 . Given \mathbf{w} and \mathbf{s} , we have

$$q(\alpha_0 | \mathbf{w}, \mathbf{s}) \propto \alpha_0^{a-1+m_{..}-\sum_{j=1}^J s_j} e^{-\alpha_0(b-\sum_{j=1}^J \log w_j)}, \quad (\text{A.4})$$

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 α_0 , the w_j and s_j are conditionally independent, with distributions

$$q(w_j | \alpha_0) \propto w_j^{\alpha_0} (1 - w_j)^{n_{j..}-1} \quad (\text{A.5})$$

and

$$q(s_j | \alpha_0) \propto \left(\frac{n_{j..}}{\alpha_0}\right)^{s_j}, \quad (\text{A.6})$$

which are beta and Bernoulli distributions. This completes the auxiliary variable sampling scheme for α_0 . We used the auxiliary variable sampling scheme in our simulations, because it is easier to implement and typically mixes quickly (within 20 iterations).

Given the total number, $m_{..}$, of the ψ_{jt} 's, the concentration parameter γ governs the distribution over the number of components K ,

$$p(K | \gamma, m_{..}) = s(m_{..}, K) \gamma^K \frac{\Gamma(\gamma)}{\Gamma(\gamma + m_{..})}. \quad (\text{A.7})$$

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

[Received October 2004. Revised December 2005.]

REFERENCES

- Aldous, D. (1985), "Exchangeability and Related Topics," in *École d'Été de Probabilités de Saint-Flour XIII-1983*, Berlin: Springer-Verlag, pp. 1-198.
- Antoniak, C. E. (1974), "Mixtures of Dirichlet Processes With Applications to Bayesian Nonparametric Problems," *The Annals of Statistics*, 2, 1152-1174.
- Beal, M. J. (2003), "Variational Algorithms for Approximate Bayesian Inference," unpublished doctoral thesis, University College London, Gatsby Computational Neuroscience Unit.
- Beal, M. J., Ghahramani, Z., and Rasmussen, C. (2002), "The Infinite Hidden Markov Model," in *Advances in Neural Information Processing Systems*, eds. T. G. Dietterich, S. Becker, and Z. Ghahramani, Cambridge, MA: MIT Press, pp. 577-584.
- Blackwell, D., and MacQueen, J. B. (1973), "Ferguson Distributions via Pólya Urn Schemes," *The Annals of Statistics*, 1, 353-355.
- Blei, D. M., and Jordan, M. I. (2005), "Variational Methods for Dirichlet Process Mixtures," *Bayesian Analysis*, 1, 121-144.
- Blei, D. M., Jordan, M. I., and Ng, A. Y. (2003), "Hierarchical Bayesian Models for Applications in Information Retrieval," *Bayesian Statistics*, 7, 25-44.
- Carota, C., and Parmigiani, G. (2002), "Semiparametric Regression for Count Data," *Biometrika*, 89, 265-281.
- Cifarelli, D., and Regazzini, E. (1978), "Problemi Statistici Non Parametrici in Condizioni di Scambiabilità Parziale e Impiego di Medie Associate," technical report, Quaderni Istituto Matematica Finanziaria dell'Università di Torino.
- De Iorio, M., Müller, P., and Rosner, G. L. (2004), "An ANOVA Model for Dependent Random Measures," *Journal of the American Statistical Association*, 99, 205-215.
- Durbin, R., Eddy, S., Krogh, A., and Mitchison, G. (1998), *Biological Sequence Analysis*, Cambridge, U.K.: Cambridge University Press.
- Escobar, M. D., and West, M. (1995), "Bayesian Density Estimation and Inference Using Mixtures," *Journal of the American Statistical Association*, 90, 577-588.
- Ferguson, T. S. (1973), "A Bayesian Analysis of Some Nonparametric Problems," *The Annals of Statistics*, 1, 209-230.
- Fong, D. K. H., Pammer, S. E., Arnold, S. F., and Bolton, G. E. (2002), "Re-analyzing Ultimatum Bargaining: Comparing Nondecreasing Curves Without Shape Constraints," *Journal of Business & Economic Statistics*, 20, 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, 611-631.

- Gabriel, S. B., Schaffner, S. F., Nguyen, H., Moore, J. M., Roy, J., Blumenstiel, B., Higgins, J., DeFelice, M., Lochner, A., Faggart, M., Liu-Cordero, S. N., Rotimi, C., Adeyemo, A., Cooper, R., Ward, R., Lander, E. S., Daly, M. J., and Altshuler, D. (2002), "The Structure of Haplotype Blocks in the Human Genome," *Science*, 296, 2225–2229.
- Green, P., and Richardson, S. (2001), "Modelling Heterogeneity With and Without the Dirichlet Process," *Scandinavian Journal of Statistics*, 28, 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. F. (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, 175–190.
- Ishwaran, H., and Zarepour, M. (2002), "Exact and Approximate Sum-Representations for the Dirichlet Process," *Canadian Journal of Statistics*, 30, 269–283.
- Jain, S., and Neal, R. M. (2000), "A Split-Merge Markov Chain Monte Carlo Procedure for the Dirichlet Process Mixture Model," Technical Report 2003, University of Toronto, Dept. of Statistics.
- Kleinman, K. P., and Ibrahim, J. G. (1998), "A Semi-Parametric Bayesian Approach to Generalized Linear Mixed Models," *Statistics in Medicine*, 17, 2579–2596.
- MacEachern, S. N. (1999), "Dependent Nonparametric Processes," in *Proceedings of the Bayesian Statistical Science Section*, American Statistical Association.
- MacEachern, S. N., Kottas, A., and Gelfand, A. E. (2001), "Spatial Nonparametric Bayesian Models," Technical Report 01-10, Duke University, Institute of Statistics and Decision Sciences.
- MacEachern, S. N., and Müller, P. (1998), "Estimating Mixture of Dirichlet Process Models," *Journal of Computational and Graphical Statistics*, 7, 223–238.
- MacKay, D. J. C. (1997), "Ensemble Learning for Hidden Markov Models," technical report, Cambridge University, Cavendish Laboratory.
- MacKay, D. J. C., and Peto, L. C. B. (1994), "A Hierarchical Dirichlet Language Model," *Natural Language Engineering*, 1, 289–307.
- Mallick, B. K., and Walker, S. G. (1997), "Combining Information From Several Experiments With Nonparametric Priors," *Biometrika*, 84, 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, 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, Ser. B*, 66, 735–749.
- Neal, R. M. (1992), "Bayesian Mixture Modeling," *Proceedings of the Workshop on Maximum Entropy and Bayesian Methods of Statistical Analysis*, 11, 197–211.
- (2000), "Markov Chain Sampling Methods for Dirichlet Process Mixture Models," *Journal of Computational and Graphical Statistics*, 9, 249–265.
- Pitman, J. (1996), "Random Discrete Distributions Invariant Under Size-Biased Permutation," *Advances in Applied Probability*, 28, 525–539.
- (2002a), "Combinatorial Stochastic Processes," Technical Report 621, University of California at Berkeley, Dept. of Statistics, 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, 501–514.
- Pritchard, J. K., Stephens, M., and Donnelly, P. (2000), "Inference of Population Structure Using Multilocus Genotype Data," *Genetics*, 155, 945–959.
- Rabiner, L. (1989), "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition," *Proceedings of the IEEE*, 77, 257–285.
- Rasmussen, C. E. (2000), "The Infinite Gaussian Mixture Model," in *Advances in Neural Information Processing Systems*, eds. S. Solla, T. Leen, and K.-R. Müller, Cambridge, MA: MIT Press.
- Salton, G., and McGill, M. J. (1983), *An Introduction to Modern Information Retrieval*, New York: McGraw-Hill.
- Sethuraman, J. (1994), "A Constructive Definition of Dirichlet Priors," *Statistica Sinica*, 4, 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, 978–989.
- Stolcke, A., and Omohundro, S. (1993), "Hidden Markov Model Induction by Bayesian Model Merging," in *Advances in Neural Information Processing Systems*, Vol. 5, eds. C. Giles, S. Hanson, and J. Cowan, San Mateo CA: Morgan Kaufmann, pp. 11–18.
- Tomlinson, G. A. (1998), "Analysis of Densities," unpublished doctoral thesis, University of Toronto, Dept. of Public Health Sciences.
- Tomlinson, G. A., and Escobar, M. (2003), "Analysis of Densities," technical report, University of Toronto, Dept. of Public Health Sciences.