

# Exact slice sampler for Hierarchical Dirichlet Processes

Arash A. Amini, Marina Paez, Lizhen Lin and Zahra S. Razaee

March 22, 2019

## Abstract

We propose an exact slice sampler for Hierarchical Dirichlet process (HDP) and its associated mixture models (Teh et al., 2006). Although there are existing MCMC algorithms for sampling from the HDP, a slice sampler has been missing from the literature. Slice sampling is well-known for its desirable properties including its fast mixing and its natural potential for parallelization. On the other hand, the hierarchical nature of HDPs poses challenges to adopting a full-fledged slice sampler that automatically truncates all the infinite measures involved without ad-hoc modifications. In this work, we adopt the powerful idea of Bayesian variable augmentation to address this challenge. By introducing new latent variables, we obtain a full factorization of the joint distribution that is suitable for slice sampling. Our algorithm has several appealing features such as (1) fast mixing; (2) remaining exact while allowing natural truncation of the underlying infinite-dimensional measures, as in (Kalli et al., 2011), resulting in updates of only a finite number of necessary atoms and weights in each iteration; and (3) being naturally suited to parallel implementations. The underlying principle for joint factorization of the full likelihood is simple and can be applied to many other settings, such as designing sampling algorithms for general dependent Dirichlet process (DDP) models.

## 1 Introduction

Hierarchical Dirichlet process (HDP) [Teh+06] is one of the popular Bayesian nonparametric models for modeling the hierarchy of groups of data. It has been widely applied in various learning tasks in statistics and machine learning, such as topic modeling [Teh+06; New+09] and information retrieval [Teh+06]. It is often used as the mixing measure in a mixture model for modeling the cluster structure for groups of data [SX09]. Other applications include using HDP for modeling the transition probabilities between hidden states in a hidden Markov model [Fox+11]. Being able to sample efficiently from a HDP is crucial for making inference in HDP-related models, and to do so, both approximating algorithms and sampling-based algorithms have been proposed in the literature. Approximating algorithms for HDP are mostly based on variational approaches [KWT07; SWA09; New+09]. These algorithms can often scale to large datasets but suffer from some obvious drawbacks: the variational posterior tends to underestimate the variability of the true posterior.

In this work, we will focus on the sampling-based methods. One of the prominent methods is the Chinese Restaurant Franchise (CRF)-based Gibbs sampler [Teh+06]. This algorithm, however, is known to mix slowly, encouraging the search for more efficient sampling approaches. Attempting that, [Fox+11] adopted a truncated approximation to the full posterior distribution for sampling HDP in a hidden Markov model, and [WB12] proposed a split/merge MCMC

algorithm, which, however, according to [CFI14] only shows a marginal improvement over the sampler of [Teh+06]. More recently, [CFI14] proposed an algorithm that also considers split/merge steps but has the advantage of allowing parallel sampling, and was reported to exhibit significantly improved convergence compared to [WB12]. In [Kim+16], online algorithms based on mini-batch ideas are proposed.

An efficient approach to sample from DP mixtures is through slice sampling. This approach has the advantage of allowing natural truncation of the infinite-dimensional Dirichlet measures, as shown in the seminal work of [KGW11]. Slice samplers are also known to have great mixing properties. To the best of our knowledge, however, no slice sampler has been proposed for HDPs. We fill this gap with the development of an *exact slice sampler* for HDP, extending the ideas of the slice sampling for DPs to the hierarchical setup.

One difficulty in performing slice sampling for HDPs arises from the hierarchical nature of the HDP model. We refer to Section 2.1 for more details on the difficulty facing the hierarchical stick-breaking representation and the marginalization approach. To circumvent this, we adopt the powerful idea of Bayesian variable augmentation which allows a full factorization of the joint distribution. The slice sampler that we derive for HDPs retains the same advantage of being exact while truncating the underlying infinite-dimensional measures ( $\beta$  and  $\gamma_j$  in our notation). In other words, the sampler only updates the necessary atoms and weights in each iteration, whose number is guaranteed to be finite.

The slice sampler proposed here enjoys several key features: it is very simple and its updates are fairly intuitive. It is also naturally parallel. Moreover, as a by-product, we derive a complete factorization of the joint density of HDP that could be of independent interest. The expression for the joint density allows one to easily verify the validity of the slice sampler updates; i.e., not much knowledge of Dirichlet processes and their intricacies is required, and the derivations are accessible to most practitioners with a basic understanding of Bayesian statistics.

Our motivation for this new algorithm came from a network point of view, and the need to propose inference models for multiplex networks that can take into account potential dependency across different layers, particularly when the aim is community detection. In this work (that is currently in progress), we specify HDP as a natural random partition prior for the partitions across different layers in the multiplex network. Despite our original motivation, it is important to point out that our algorithm can be easily generalized to other models if one replaces the mixture part with a general likelihood, therefore immediately getting a nicer sampling scheme.

## 2 Representation of HDPs

We start by deriving a representation of the HDPs which is suitable for slice sampling. We assume familiarity with the setup of a HDP as in [Teh+06] and the associated metaphor of a Chinese restaurant franchise (CRF). We first focus on the label generation part of the HDP, and then discuss how one can add the bottom mixture layer.

## 2.1 Label-only HDP

Consider the HDP as defined in Equ. (19) of [Teh+06]. Using mostly the same notation and the CRF metaphor, we write

$$\begin{aligned}\boldsymbol{\beta} &| \gamma_0 \sim \text{GEM}(\gamma_0) \\ \boldsymbol{\pi}_j &| \alpha_o, \boldsymbol{\beta} \sim \text{DP}(\alpha_o, \boldsymbol{\beta}), \\ z_{ji} &| \boldsymbol{\pi}_j \sim \boldsymbol{\pi}_j\end{aligned}$$

where  $i = 1, \dots, n_j$  is the customer index and  $j = 1, \dots, J$  is the restaurant index. We are ignoring the downstream mixture model for the moment, since this is the main part of the sampling problem. Thus, we assume that we observe the labels  $\{z_{ji}\}$  where  $z_{ji}$  is the dish of customer  $i$  in restaurant  $j$ , and we want to estimate  $\boldsymbol{\beta}$  and  $\{\boldsymbol{\pi}_j\}$ .

The stick-breaking representation of  $\boldsymbol{\pi}_j$ —see Equ. (21) in [Teh+06]—is not suitable for slice (or Gibbs) sampling due to the complicated dependence on  $\boldsymbol{\beta}$ . Another idea is to marginalize  $\boldsymbol{\pi}_j$ , but that would lead to distributions with ratios of Gamma functions as densities which are not easy to sample from. Instead, we just use the fact that  $\boldsymbol{\pi}_j$  is itself a Dirichlet measure. That is, we can write the model as

$$\begin{aligned}\boldsymbol{\beta} &| \gamma_0 \sim \text{GEM}(\gamma_0) \\ \boldsymbol{\gamma}_j &| \alpha_0 \sim \text{GEM}(\alpha_0), \quad \boldsymbol{\gamma}_j = (\gamma_{jt}) \\ k_{jt} &| \boldsymbol{\beta} \sim \boldsymbol{\beta}, \quad t \in \mathbb{N} \\ \boldsymbol{\pi}_j &= \sum_{t=1}^{\infty} \gamma_{jt} \delta_{k_{jt}}, \quad z_{ji} | \boldsymbol{\pi}_j \sim \boldsymbol{\pi}_j.\end{aligned}$$

Note that  $\boldsymbol{\gamma}_j$  and  $(k_{jt})$  are independently drawn. Here, index “ $t$ ” is interpreted as indexing the tables;  $k_{jt}$  is the dish (type) of table  $t$  in restaurant  $j$ , while  $\gamma_{jt}$  represents the fraction of the customers in restaurant  $j$  that would sit at table  $t$  (eventually).

This representation is still not suitable for sampling. Instead of sampling directly from  $\boldsymbol{\pi}_j$ , we sample from the weights  $\boldsymbol{\gamma}_j$  first, i.e., we pick the table of customer  $i$  and then assign them the dish of the table. This gives us access to the last missing piece which is  $t_{ji}$ , the table of customer  $i$  in restaurant  $j$ . We can write the model equivalently as

$$\begin{aligned}\boldsymbol{\beta} &| \gamma_0 \sim \text{GEM}(\gamma_0), \quad \boldsymbol{\beta} = (\beta_k) \\ \boldsymbol{\gamma}_j &| \alpha_0 \sim \text{GEM}(\alpha_0), \quad \boldsymbol{\gamma}_j = (\gamma_{jt}) \\ k_{jt} &| \boldsymbol{\beta} \sim \boldsymbol{\beta}, \quad t \in \mathbb{N}, \\ t_{ji} &| \boldsymbol{\gamma}_j \sim \boldsymbol{\gamma}_j, \quad i = 1, \dots, n_j \\ z_{ji} &| \boldsymbol{t}_j, \boldsymbol{k}_j = k_{j,t_{ji}}\end{aligned}\tag{1}$$

where  $\boldsymbol{k}_j = (k_{jt}, t \in \mathbb{N})$  is the collection of all the dishes at restaurant  $j$ . Note that we sample  $t_{ji}$  (which table to sit customer  $i$  in restaurant  $j$ ) from the eventual distributions of customers among tables in restaurant  $j$ , i.e.,  $\boldsymbol{\gamma}_j$ . Equation  $z_{ji} = k_{j,t_{ji}}$  means that the dish of customer  $i$  in restaurant  $j$ , i.e.  $z_{ji}$ , is completely determined by looking at which table they are sitting at,  $t_{ji}$ , and what dish is presented at that table  $k_{j,t_{ji}}$ .

Since given everything else,  $z_{ji}$  is deterministic, we only need to worry about sampling  $\boldsymbol{\beta}$ ,  $\boldsymbol{\gamma}_j$ ,  $\mathbf{k}_j$  and  $\mathbf{t}_j = (t_{ji})$ . Let us define  $F : [0, 1]^{\mathbb{N}} \rightarrow [0, 1]^{\mathbb{N}}$  by

$$[F(\mathbf{x})]_1 := x_1, \quad [F(\mathbf{x})]_j := x_j \prod_{\ell=1}^{j-1} (1 - x_\ell) \quad (2)$$

where  $\mathbf{x} = (x_j, j \in \mathbb{N})$ . Both  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}_j$  have stick-breaking representations [Set94; IJ01]:

$$\begin{aligned} \gamma'_{jt} &\sim \text{Beta}(1, \alpha_0), & \beta'_k &\sim \text{Beta}(1, \gamma_0), \\ \boldsymbol{\gamma}_j &= F(\boldsymbol{\gamma}'_j), & \boldsymbol{\beta} &= F(\boldsymbol{\beta}'), \end{aligned}$$

where  $\boldsymbol{\gamma}'_j = (\gamma'_{jt})$  and  $\boldsymbol{\beta}' = (\beta'_k)$ . Let us write  $x \mapsto b_{\alpha_0}(x)$  for the density of  $\text{Beta}(1, \alpha_0)$ , that is,  $b_{\alpha_0}(x) \propto (1 - x)^{\alpha_0 - 1}$ .

Note that  $\mathbb{P}(t_{ji} = t) = \gamma_{jt}$ ,  $t \in \mathbb{N}$ . Thus, we can write down the joint density as

$$\begin{aligned} p(\mathbf{t}, \mathbf{k}, \boldsymbol{\gamma}', \boldsymbol{\beta}') &= \prod_{j=1}^J \left[ p(\mathbf{t}_j | \boldsymbol{\gamma}_j) p(\boldsymbol{\gamma}'_j) p(\mathbf{k}_j | \boldsymbol{\beta}') \right] p(\boldsymbol{\beta}') \\ &= \prod_{j=1}^J \left( \prod_{i=1}^{n_j} \gamma_{j,t_{ji}} \prod_{t=1}^{\infty} b_{\alpha_0}(\gamma'_{jt}) \prod_{t=1}^{\infty} \beta_{k_{jt}} \right) \prod_{k=1}^{\infty} b_{\gamma_0}(\beta'_k). \end{aligned} \quad (3)$$

Interestingly, this decomposition works for any other stick-breaking distributions on  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}_j$ . Using  $\boldsymbol{\gamma}_j = F(\boldsymbol{\gamma}'_j)$  and  $\boldsymbol{\beta} = F(\boldsymbol{\beta}')$ , a more explicit formula is

$$p(\mathbf{t}, \mathbf{k}, \boldsymbol{\gamma}', \boldsymbol{\beta}') = \prod_{j=1}^J \left( \prod_{i=1}^{n_j} [F(\boldsymbol{\gamma}'_j)]_{t_{ji}} \prod_{t=1}^{\infty} b_{\alpha_0}(\gamma'_{jt}) \prod_{t=1}^{\infty} [F(\boldsymbol{\beta}')]_{k_{jt}} \right) \prod_{k=1}^{\infty} b_{\gamma_0}(\beta'_k) \quad (4)$$

which gives the complete joint density of HDP in (1).

## 2.2 Mixture part

Finally, we can add in the mixture component as

$$\begin{aligned} f_k | \mathcal{F} &\sim \mathcal{F}, & \mathbf{f} &= (f_k) \\ y_{ji} | z_{ji}, \mathbf{f} &\sim f_{z_{ji}}, \end{aligned} \quad (5)$$

where  $\mathbf{f}$  is an infinite collection of possible mixture components, where each coordinate  $f_k$  is a density drawn from a distribution  $\mathcal{F}$  on densities. We can assume  $f_k = K(\cdot, \phi_k)$  for some kernel  $K$ , where  $\phi_k | H \sim H$ , to get back the more common parametric mixture model; the more general setup however is easier to work with conceptually.

Since  $p(\mathbf{y} | \mathbf{t}, \mathbf{k}, \mathbf{f}) = \prod_{j=1}^J \prod_{i=1}^{n_j} f_{z_{ji}}(y_{ji})$ , the overall joint density is

$$\begin{aligned} p(\mathbf{y}, \mathbf{f}, \mathbf{t}, \mathbf{k}, \boldsymbol{\gamma}', \boldsymbol{\beta}') &= p(\mathbf{y} | \mathbf{t}, \mathbf{k}, \mathbf{f}) p(\mathbf{t}, \mathbf{k}, \boldsymbol{\gamma}', \boldsymbol{\beta}') p(\mathbf{f}) \\ &= \prod_{j=1}^J \left( \prod_{i=1}^{n_j} [f_{k_{j,t_{ji}}}(y_{ji}) \gamma_{j,t_{ji}}] \prod_{t=1}^{\infty} b_{\alpha_0}(\gamma'_{jt}) \prod_{t=1}^{\infty} \beta_{k_{jt}} \right) \prod_{k=1}^{\infty} [b_{\gamma_0}(\beta'_k) \mathcal{F}(f_k)]. \end{aligned} \quad (6)$$

We note that this joint density is completely factorized over all its variables. The diagram of the model is shown in Figure 1.

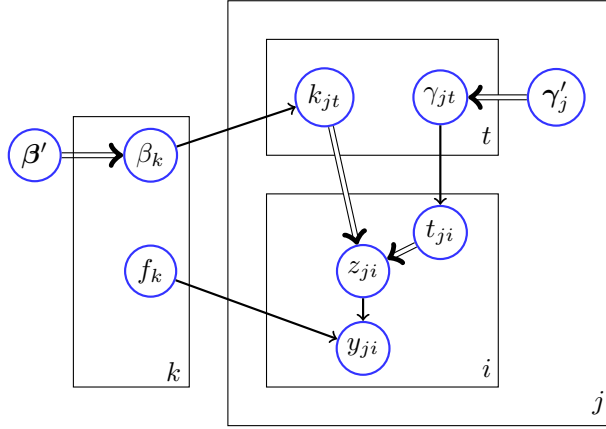


Figure 1: Schematic diagram of the HDP mixture with latent variables introduced for sampling. Double arrows show deterministic relations.

### 3 Sampling

For the most part, when sampling, we can ignore the mixture part. That is, for the most part it is enough to sample from (4). Only in sampling  $\mathbf{t}$  and  $\mathbf{k}$  the mixture part comes in. The factorized form of the density in (6) allows us to easily derive the Gibbs updates.

First, we state a key lemma. Let  $x \mapsto b(x; \alpha, \beta)$  be the density of Beta( $\alpha, \beta$ ). The derivations in this paper can be extended to any stick-breaking prior for which a conjugacy relation similar to the one described in the lemma holds:

**Lemma 1.** *Assume that the joint density of  $\mathbf{x} = (x_1, x_2, \dots) \in [0, 1]^{\mathbb{N}}$  is proportional to*

$$\prod_{i=1}^n [F(\mathbf{x})]_{z_i} \prod_{j=1}^{\infty} b(x_j; \alpha, \beta),$$

where  $\mathbf{z} = (z_1, \dots, z_n) \in \mathbb{N}^n$  and  $F: [0, 1]^{\mathbb{N}} \mapsto [0, 1]^{\mathbb{N}}$  is defined as in (2). Then

$$x_j \mid \mathbf{x}_{-j} \sim \text{Beta}(n_j(\mathbf{z}) + \alpha, n_{>j}(\mathbf{z}) + \beta)$$

where  $n_j(\mathbf{z}) = |\{i : z_i = j\}|$  and  $n_{>j}(\mathbf{z}) = |\{i : z_i > j\}|$ .

The proof is given in Appendix A.

#### 3.1 Usual block Gibbs sampling

**Sampling  $\gamma' \mid \mathbf{t}, \mathbf{k}, \beta'$ .** This posterior factorizes over  $\gamma'_j$ , and the posterior of  $\gamma'_j$  given the rest is proportional to  $\prod_i [F(\gamma'_j)]_{t_{ji}} \prod_t b_{\alpha_0}(\gamma'_{jt})$ . Using Lemma 1, we have

$$\gamma'_{jt} \mid \gamma'_{-jt}, \mathbf{t}, \mathbf{k}, \beta' \sim \text{Beta}(n_t(\mathbf{t}_j) + 1, n_{>t}(\mathbf{t}_j) + \alpha_0). \quad (7)$$

Note that, for a fixed  $j$ , we are applying Lemma 1 to the factorization indexed by  $t$ . Here,  $n_t(\mathbf{t}_j) = |\{i : t_{ji} = t\}|$  and  $n_{>t}(\mathbf{t}_j) = |\{i : t_{ji} > t\}|$ .

**Sampling  $\beta' \mid \mathbf{t}, \mathbf{k}, \gamma'$ .** This posterior is proportional to  $\prod_j \prod_t [F(\beta')]_{k_{jt}} \prod_k b_{\gamma_0}(\beta'_k)$ . Applying Lemma 1 to the factorization over  $k$ , we obtain

$$\beta'_k \mid \beta'_{-k}, \gamma', \mathbf{t}, \mathbf{k} \sim \text{Beta}(n_k(\mathbf{k}) + 1, n_{>k}(\mathbf{k}) + \gamma_0) \quad (8)$$

where  $n_k(\mathbf{k}) = |\{(j, t) : k_{j,t} = k\}|$  and similarly for  $n_{>k}(\mathbf{k})$ .

**Sampling  $t \mid \mathbf{k}, \gamma', \beta', \mathbf{y}, \mathbf{f}$ .** This posterior factorizes over  $j$  and  $i$  and is given by

$$t_{ji} \mid \gamma'_j, \mathbf{k}_j, \mathbf{y}_j, \mathbf{f} \sim (f_{k_{jt}}(y_{ji}) \gamma_{jt})_{t \in \mathbb{N}}.$$

**Sampling  $\mathbf{k} \mid \mathbf{t}, \gamma', \beta', \mathbf{y}, \mathbf{f}$ .** This posterior factorizes over  $j$ . Writing

$$f_{k_{j,t_{ji}}}(y_{ji}) = \prod_{t=1}^{\infty} [f_{k_{jt}}(y_{ji})]^{1\{t_{ji}=t\}},$$

the posterior for  $\mathbf{k}_j$ , given the rest, is proportional to

$$\prod_{i=1}^{n_j} [f_{k_{j,t_{ji}}}(y_{ji})] \prod_{t=1}^{\infty} \beta_{k_{jt}} = \prod_{t=1}^{\infty} \left( \beta_{k_{jt}} \prod_{i=1}^{n_j} [f_{k_{jt}}(y_{ji})]^{1\{t_{ji}=t\}} \right) \quad (9)$$

which also factorizes over  $t$ . Thus, it is enough to sample (independently over  $j$  and  $t$ ),

$$k_{jt} \mid \beta', \mathbf{t}_j, \mathbf{y}_j, \mathbf{f} \sim \left( \beta_k \prod_{i: t_{ij}=t} f_k(y_{ji}) \right)_{k \in \mathbb{N}}.$$

**Sampling  $\mathbf{f} \mid \dots$ .** Recalling  $z_{ji} = k_{j,t_{ji}}$ , we sample independently over  $k$ ,

$$p(f_k \mid \dots) \propto \mathcal{F}(f_k) \prod_{(i,j): z_{ji}=k} f_k(y_{ji}). \quad (10)$$

### 3.2 Slice sampling

We recall the basic idea of slice sampling, which itself is a form of variable augmentation: In order to sample from density  $f(x)$ , we introduce the nonnegative variable  $u$ , and look at the joint density  $g(x, u) = 1\{u \leq f(x)\}$  whose marginal over  $x$  is  $f(x)$ . Then, we perform Gibbs sampling on the joint  $g$ . In the end, we only keep samples of  $x$  and discard those of  $u$ . This idea has been successfully employed in [KGW11] to sample from the classical DP mixture. We now extend the ideas in [KGW11] to sample from HDP mixtures.

To carry the idea over to the HDPs, we augment the model by adding variables  $\mathbf{u}_j = (u_{ji})$  and  $\mathbf{v}_j = (v_{jt})$  and consider the joint density

$$p(\mathbf{y}, \mathbf{f}, \mathbf{t}, \mathbf{k}, \gamma', \mathbf{u}, \beta', \mathbf{v}) = \prod_{j=1}^J \left( \prod_{i=1}^{n_j} f_{k_{j,t_{ji}}}(y_{ji}) 1\{u_{ji} \leq \gamma_{j,t_{ji}}\} \prod_{t=1}^{\infty} b_{\alpha_0}(\gamma'_{jt}) \prod_{t=1}^{\infty} 1\{v_{jt} \leq \beta_{k_{jt}}\} \right) \prod_{k=1}^{\infty} [b_{\gamma_0}(\beta'_k) \mathcal{F}(f_k)]. \quad (11)$$

Note that by integrating out the variables  $(u_{ji})$  and  $(v_{jt})$ , we get back original joint density (6). The idea is that we sample  $(\gamma', \mathbf{u})$  jointly given the rest of variables, and similarly for  $(\beta', \mathbf{v})$ .

**Sampling**  $(\gamma', \mathbf{u}) \mid \mathbf{t}, \mathbf{k}, \beta', \mathbf{v}$ . First we sample  $(\mathbf{u} \mid \gamma', \mathbf{t}, \mathbf{k}, \beta', \mathbf{v})$  which factorizes and the coordinate posteriors are  $p(u_{ji} \mid \gamma', \mathbf{t} \dots) \propto 1\{u_{ji} \leq \gamma_{j,t_{ji}}\}$ , that is

$$u_{ji} \mid \gamma', \mathbf{t}, \mathbf{k}, \beta', \mathbf{v} \sim \text{Unif}(0, \gamma_{j,t_{ji}}).$$

Next we sample from  $(\gamma' \mid \mathbf{t}, \mathbf{k}, \beta', \mathbf{v})$ . This would be the same as (7).

**Sampling**  $(\beta', \mathbf{v}) \mid \mathbf{k}, \mathbf{t}, \gamma', \mathbf{u}$ . First we sample  $(\mathbf{v} \mid \beta', \mathbf{k}, \mathbf{t}, \gamma', \mathbf{u})$  which factorizes and the coordinate posteriors are  $p(v_{jt} \mid \beta', \mathbf{k} \dots) \propto 1\{v_{jt} \leq \beta_{k_{jt}}\}$ , that is

$$v_{jt} \mid \beta', \mathbf{k}, \mathbf{t}, \gamma', \mathbf{u} \sim \text{Unif}(0, \beta_{k_{jt}}).$$

Next, we sample from  $(\beta' \mid \mathbf{k}, \mathbf{t}, \gamma', \mathbf{u})$ . This would be the same as (8).

**Sampling**  $\mathbf{t} \mid \dots$ . This posterior also factorizes over  $i$  and  $j$ . From (11), we have

$$\mathbb{P}(t_{ji} = t \mid \mathbf{t}_{-ji}, \mathbf{k}, \gamma', \mathbf{u}, \beta', \mathbf{v}) \propto f_{k_{jt}}(y_{ji}) 1\{u_{ji} \leq \gamma_{jt}\}. \quad (12)$$

Let  $T_{ji} := T(\gamma_j; u_{ji}) := \sup\{t : u_{ji} \leq \gamma_{jt}\}$ . According to the above,  $t_{ji}$  given everything else will be distributed as

$$t_{ji} \mid \dots \sim \left( f_{k_{jt}}(y_{ji}) \right)_{t \in [T_{ji}]}.$$

In the CRF metaphor,  $T_{ji}$  is the maximum table index ( $t$ ) that customer  $i$  in restaurant  $j$  can hop to at current iteration. Note that different customers are allowed different ranges of tables for their wandering. The update for  $\mathbf{t}$  is an instance of how the slice sampler truncates an infinite measure. Due to the presence of the indicator in (12), only values of  $t$  for which  $\gamma_{jt} \geq u_{ji}$  lead to a nonzero probability. In other words, the support of distribution (12) is contained in  $[T_{ji}]$ .

**Sampling**  $\mathbf{k} \mid \dots$ . This posterior also factorizes over  $j$  and  $t$ . From (11), the posterior for  $\mathbf{k}_j$  given the rest is proportional to the same expression (9) but with  $\beta_{k_{jt}}$  replaced with  $1\{v_{jt} \leq \beta_{k_{jt}}\}$ . Thus, we have

$$\mathbb{P}(k_{jt} = k \mid \dots) \propto 1\{v_{jt} \leq \beta_k\} \prod_{i: t_{ji}=t} f_k(y_{ji}).$$

Let  $K_{jt} := K(\beta; v_{jt}) := \sup\{k : v_{jt} \leq \beta_k\}$ . According to the above,  $k_{jt}$  given everything else will be distributed as

$$k_{jt} \mid \dots \sim \left( \prod_{i: t_{ji}=t} f_k(y_{ji}) \right)_{k \in [K_{jt}]}.$$

In CRF metaphor,  $K_{jt}$  is the maximum dish index ( $k$ ) available for substitution at table  $t$  in restaurant  $j$  at current iteration. Again different tables in the same restaurant have potentially different options for dish exchange. The update of  $\mathbf{k}$  is another instance where the slice sampler is truncating the infinite measures involved.

**Sampling**  $\mathbf{f} \mid \dots$ . This will be the same as (10).

---

**Algorithm 1** Slice sampler for HDP mixture

---

- 1: Initialize  $T_j^{\text{cap}}$  and  $K^{\text{cap}}$  to pre-specified values (say 10).
- 2: Initialize  $\mathbf{t}_j$  and  $\mathbf{k}_j$  to all-ones vectors.
- 3: Initialize  $(u_{ji})$  and  $(v_{jt})$  to independent uniform variables.
- 4: **while** not CONVERGED, nor maximum iterations reached **do**
- 5:   **for**  $j = 1, \dots, J$  **do**
- 6:     Sample  $\gamma'_{jt} \sim \text{Beta}(n_t(\mathbf{t}_j) + 1, n_{>t}(\mathbf{t}_j) + \alpha_0)$  for all  $t \in [T_j^{\text{cap}}]$ .
- 7:     Let  $[\gamma_j]_{1:T_j^{\text{cap}}} \leftarrow [F(\gamma'_j)]_{1:T_j^{\text{cap}}}$ .
- 8:     Let  $T_{ji} \leftarrow \max\{t : u_{ji} \leq \gamma_{jt}\}, \forall i \in [n_j]$  and  $T_j \leftarrow \max_{i=1, \dots, n_j} T_{ji}$ .
- 9:     doubbling( $T_j, T_j^{\text{cap}}$ )
- 10:   **end for**
- 11:   Sample  $\beta'_k \sim \text{Beta}(n_k(\mathbf{k}) + 1, n_{>k}(\mathbf{k}) + \gamma_0)$  for all  $k \in [K^{\text{cap}}]$ .
- 12:   Let  $[\beta]_{1:K^{\text{cap}}} \leftarrow [F(\beta')]_{1:K^{\text{cap}}}$ .
- 13:   Let  $K_{jt} \leftarrow \max\{k : v_{jt} \leq \beta_k\}, \forall t \in [T_j^{\text{cap}}], j \in [J]$ .
- 14:   Let  $K_j \leftarrow \max_t K_{jt}$ , and  $K \leftarrow \max_j K_j$ .
- 15:   doubbling( $K, K^{\text{cap}}$ )
- 16:   Sample  $f_k$  from density  $p(f | \dots) \propto \mathcal{F}(f) \prod_{(i,j): z_{ji}=k} f(y_{ji})$  for all  $k \in [K^{\text{cap}}]$ .
- 17:   **for**  $j = 1, \dots, J$  **do**
- 18:     Sample  $k_{jt} \sim \left( \prod_{i: t_{ji}=t} f_k(y_{ji}) \right)_{k \in [K_{jt}]}$  for all  $t \in [T_j^{\text{cap}}]$ . Set  $\mathbf{k}_j \leftarrow (k_{jt})$ .
- 19:     Sample  $v_{jt} \sim \text{Unif}(0, \beta_{k_{jt}})$  for all  $t \in [T_j^{\text{cap}}]$ .
- 20:     Sample  $t_{ji} \sim (f_{k_{jt}}(y_{ji}))_{t \in [T_{ji}]}$  for all  $i \in [n_j]$ . Set  $\mathbf{t}_j \leftarrow (t_{ji})$ .
- 21:     Sample  $u_{ji} \sim \text{Unif}(0, \gamma_{j,t_{ji}})$  for all  $i \in [n_j]$ .
- 22:     Set  $z_{ji} \leftarrow k_{j,t_{ji}}$ .
- 23:   **end for**
- 24: **end while**
- 25: **macro** doubling( $K, K^{\text{cap}}$ )
- 26:   **if**  $K < K^{\text{cap}}$ , **then** continue **else**  $K^{\text{cap}} \leftarrow 1.5K^{\text{cap}}$  and go to the previous iteration.

---

### 3.2.1 How many atoms to keep?

Let us define

$$T_j := \max_{i=1, \dots, n_j} T_{ji}, \quad K_j := \max_{t=1, \dots, T_j} K_{jt}, \quad K := \max_j K_j$$

so that  $T_j$  determines the maximum table index “ $t$ ” which we need to keep track of for restaurant  $j$ . Given  $T_j$ , one can compute  $K_j$  which is the maximum number of dish index “ $k$ ” we need to keep track of in restaurant  $j$ . Note that quantities  $T_j$  and  $K_j$  will be finite and random; they depend on  $\gamma, \beta, \mathbf{u}$  and  $\mathbf{v}$  and get updated in each iteration.



This completes the description of the slice sampler which is summarized in Algorithm 1. The implementation, however, uses a few other ideas besides the update equations derived earlier. The difficulty is that the count parameters  $T_{ji}$ ,  $T_j$ ,  $K_{jt}$ ,  $K_j$  and  $K$  are interrelated among themselves and with other latent parameters of the model. Updating some of the parameters while keeping others fixed would create a chain of dependencies which is hard to track.

The easiest way to assure that we always have sufficiently enough atoms, from all the infinite measures, is to put caps on their numbers, i.e.,  $T_j^{\text{cap}}$  and  $K^{\text{cap}}$  in Algorithm 1, and increase the cap whenever we hit it, and repeat the previous iteration (for which we need to keep record of the state of the chain one-step into the past). This is reflected in the `doubling` macro in Algorithm 1. This procedure creates a little bit of redundancy, but after a few steps, the chain will remain within the caps, and one avoids resampling for all but a few early updates. The algorithm guarantees that we always have  $T_j < T_j^{\text{cap}}$  and  $K < K^{\text{cap}}$  and so the chain is sampling exactly.

An advantage of the slice sampler is that all the updates in each step can be done in parallel over the underlying coordinates. Even updates at multiple steps involving disjoint sets of parameters can be performed in parallel.

**Explicit mixture densities.** Using the more common notation  $f_k(y) = K(y; \phi_k)$ , with  $\phi_k | H \sim H$ , Step 16 can be written as follows: Sample  $\phi_k$  from

$$p(\phi | \dots) \propto H(\phi) \prod_{(i,j): z_{ji}=k} K(y_{ji}; \phi) \quad \text{for all } k \in [K^{\text{cap}}] \quad (13)$$

and set  $f_k = K(\cdot; \phi_k)$ . Note that in Step 16, if the set  $\{(i, j) : z_{ji} = k\}$  is empty for some  $k \in [K]$ —which could happen since  $z_{ji}$  has not yet been updated from the previous iteration while  $K$  has just been updated—then the product evaluates to 1, and we draw  $f_k$  from the prior  $\mathcal{F}(f)$  itself.

We also note that in updating  $k_{jt}$  in Step 18, if the set  $\{i : t_{ji} = t\}$  is empty, it means that customers are no longer sitting at table  $t$ . As before, we interpret products over empty sets as evaluating to 1, hence the dish of the vacant table is updated uniformly at random.

### 3.3 Examples

Let us consider a few examples of the mixture densities  $f_k = K(\cdot; \phi_k)$ . As a first example, consider a *hierarchical Gaussian mixture*: We assume that

$$\begin{aligned} \phi_k | H &\sim H = N(0, I_d/\tau_\phi^2) \\ y_{ji} | z_{ji} &\sim N(\phi_{z_{ji}}, I_d/\tau_y^2) \end{aligned}$$

where  $\tau_\phi^2$  and  $\tau_y^2$  are the precision parameters for the prior and the likelihood, respectively. For this model, it is not hard to see that the posterior update in Step 16 is equivalent to drawing the atoms from a Gaussian distribution  $\phi_k | \dots \sim N(\mu_k, \tau_k^2 I_d)$  where

$$\mu_k = \frac{\tau_y^2}{\tau_\phi^2 + n_k(\mathbf{z})\tau_y^2} \sum_{(j,i): z_{ji}=k} y_{ji}, \quad \text{and} \quad \tau_k^2 = \tau_\phi^2 + n_k(\mathbf{z})\tau_y^2. \quad (14)$$

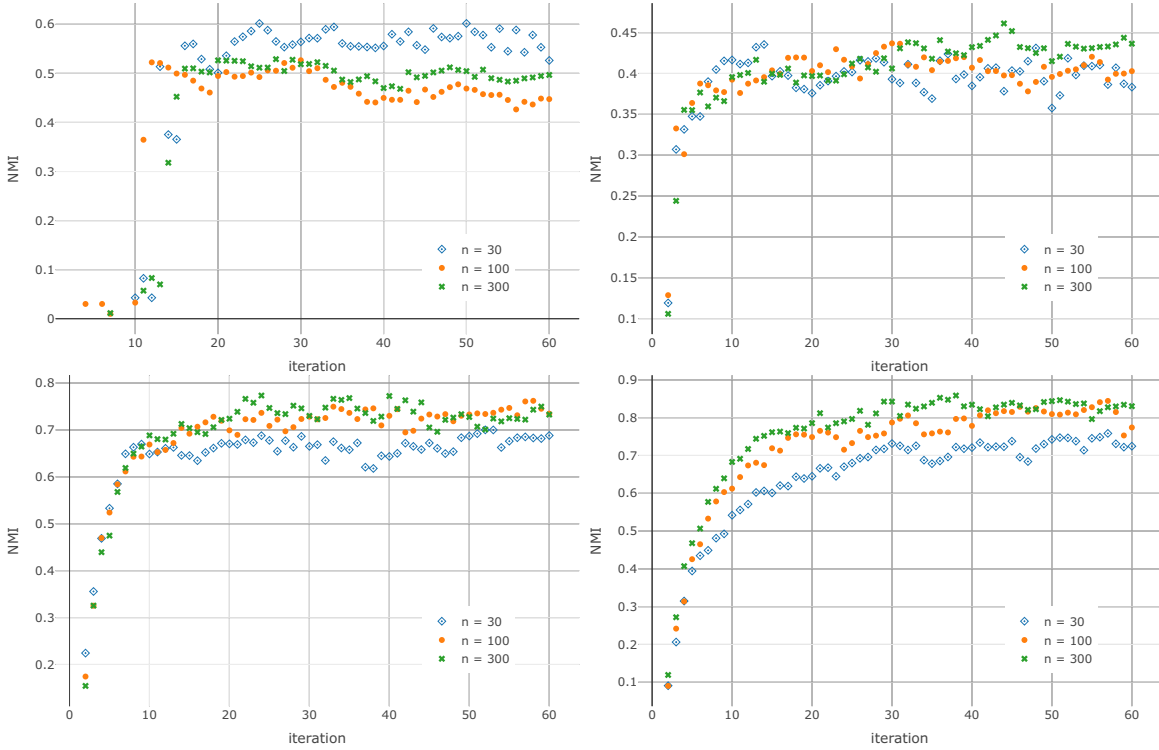


Figure 2: Typical Mixing behavior of the slice sampler for multinomial HDP-mixtures for various sample sizes (Section 3.4). From top-left clockwise:  $J = W = 10, 20, 50, 200$ . Each plot shows the aggregate normalized mutual information (NMI) at each iteration. The NMI is computed for the labels from the posterior against the “true” labels.

As another example, consider a *topic modeling* setup, where  $y_{ji}$  represents the  $i$ th word in document  $j$ . We assume  $y_{ji} \in [W]$  where  $[W] := \{1, 2, \dots, W\}$  is a vocabulary of  $W$  words, identified with their index in a dictionary. Each atom  $\phi \in [0, 1]^W$  in this case represents a probability distribution over words in vocabulary  $[W]$ . A natural prior on  $\phi$  is  $\text{Dir}((\alpha_w))$ , i.e.,

$$H(\phi) \propto \prod_{w=1}^W \phi_w^{\alpha_w - 1},$$

and the likelihood is  $y_{ji} \mid z_{ji} \sim \text{Categorical}(\phi_{z_{ji}})$  corresponding to the kernel

$$K(y; \phi) = \phi_y = \prod_{w=1}^W \phi_w^{1\{y=w\}}.$$

The posterior update in Step 16 (cf. 13) will be  $\phi_k \mid \dots \sim \text{Dir}(\alpha'_k)$  where  $\alpha'_k$  has coordinates

$$\alpha'_{kw} = \alpha_w + \sum_{j,i} 1\{y_{ji} = w, z_{ji} = k\}.$$

We also note that updating  $\mathbf{k}$ —Step (18)—simplifies to

$$k_{jt} \sim \left( \prod_w \phi_{kw}^{\nu'_{jtw}} \right)_{k \in [K_{jt}]} \quad \text{where} \quad \nu'_{jtw} = \sum_i 1\{y_{ji} = w, t_{ji} = t\}.$$

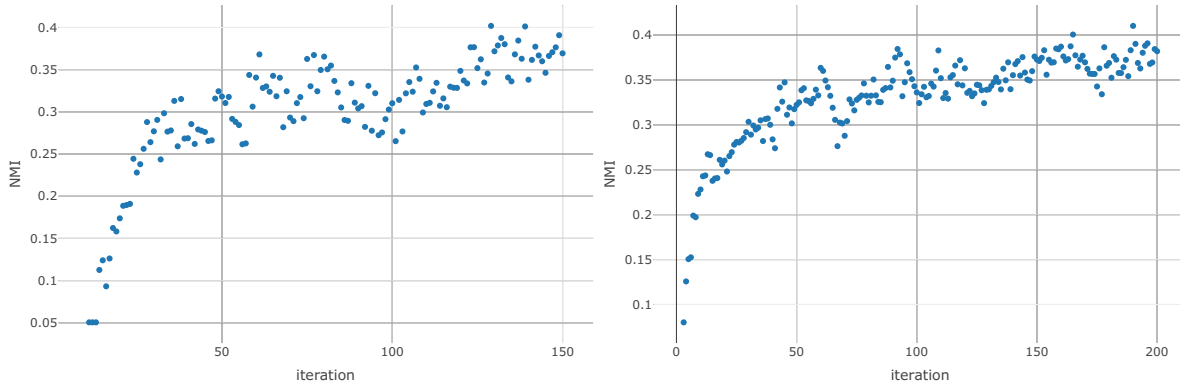


Figure 3: Results for a real data experiment. The NMI (relative to true labels) versus iteration for the real-world paper-title network. (left) performance on a random subset of size  $J = 100$  of the data (right) on the whole dataset  $J = 894$ .

### 3.4 Experiments

We now present some numerical experiments to illustrate the performance of the slice sampler. Since HDP mixtures are very popular and their performance well-known, we will mostly focus on studying the mixing time of the sampler. Figure 2 illustrates the mixing behavior for the multinomial HDP mixture discussed in Section 3.3. We have also experimented with the Gaussian mixtures but we omit them here due to similarity. The code for these experiments is available on GitHub, repository [aaamini/hdpslicer](https://github.com/aaamini/hdpslicer).

In each case we simulated from HDP with concentration parameters  $\gamma_0 = 3$  and  $\alpha_0 = 1$  and have run the slice sampler on a single sample. The multinomial parameter  $W$  is varied and we set  $\alpha_w = 1/W$ . In each case, we have  $n_j = n$  for all  $j$  and three values  $n = 30, 100, 300$  are considered. For simplicity, we have set the number of restaurants to  $J = W = 10, 20, 50, 200$ . Figure 2 illustrates single typical runs of the algorithm without burn-in or thinning; there is also no averaging over multiple runs and the labels are all initialized to 1 as in Algorithm 1.

We have calculated the normalized mutual information (NMI) between estimated  $(z_{ji})$  and true labels  $(z_{ji}^*)$ , aggregated over all  $(i, j)$ . NMI measures the matching between two clusterings, its value being in  $[0, 1]$  with a value of 1 corresponding to a perfect match. Figure 2 shows the quality of recovered labels relative to the true data-generating labels, over the iterations of the sampler. The plots clearly indicate a fast mixing time, somewhere between 10 to 20 iterations. We note the decrease in the variance of the posterior as  $n$  increases which is expected.

We have also applied the algorithm to a real world example where the documents are papers and bag-of-words are made from the words in their titles. A vocabulary of a total of  $W = 189$  was used after running standard text mining procedures for removing the stopwords, stemming, and so on. The information on a total of  $J = 894$  papers was collected from the DBLP website. The papers were published in 2017 in three CS topics: machine learning, multimedia and security. We treated DBLP subject classification as the true cluster of each paper. The HDP mixture is run on the dataset which recovers a clustering for every word in each document. We then assign an estimated cluster to each paper by majority voting (among the estimated clusters for their words) and compare with the true labels. Figure 3 illustrates the resulting NMIs versus iteration. Both a random subset of the papers (with  $J = 100$ ) and the whole set is considered. Again, we observe that the algorithm is mixing very fast, and by

about 100 iterations we already have pretty good quality labels ( $\text{NMI} \in [0.35, 0.4]$ ).

## References

- [1] J. Chang and J. W. Fisher III. Parallel sampling of HDPs using sub-cluster splits. *Advances in Neural Information Processing Systems*. 2014, pp. 235–243.
- [2] E. B. Fox et al. A sticky HDP-HMM with application to speaker diarization. *Ann. Appl. Stat.* 5.2A (June 2011), pp. 1020–1056.
- [3] H. Ishwaran and L. James. Gibbs Sampling Methods for Stick-Breaking Priors. English. *Journal of the American Statistical Association* 96.453 (Mar. 2001), pp. 161–173.
- [4] M. Kalli, J. E. Griffin, and S. G. Walker. Slice sampling mixture models. English. *Statistics and Computing* 21.1 (2011), pp. 93–105.
- [5] Y. Kim et al. An Online Gibbs Sampler Algorithm for Hierarchical Dirichlet Processes Prior. *Machine Learning and Knowledge Discovery in Databases*. Ed. by P. Frasconi et al. Cham: Springer International Publishing, 2016, pp. 509–523.
- [6] K. Kurihara, M. Welling, and Y. W. Teh. Collapsed Variational Dirichlet Process Mixture Models. *Proceedings of the 20th International Joint Conference on Artificial Intelligence*. IJCAI’07. Hyderabad, India: Morgan Kaufmann Publishers Inc., 2007, pp. 2796–2801.
- [7] D. Newman et al. Distributed Algorithms for Topic Models. *J. Mach. Learn. Res.* 10 (Dec. 2009), pp. 1801–1828.
- [8] J. Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica* 4 (1994), pp. 639–650.
- [9] P. Smyth, M. Welling, and A. U. Asuncion. Asynchronous Distributed Learning of Topic Models. *Advances in Neural Information Processing Systems 21*. Ed. by D. Koller et al. Curran Associates, Inc., 2009, pp. 81–88.
- [10] K.-A. Sohn and E. P. Xing. A hierarchical Dirichlet process mixture model for haplotype reconstruction from multi-population data. *Ann. Appl. Stat.* 3.2 (June 2009), pp. 791–821.
- [11] Y. W. Teh et al. Hierarchical Dirichlet Processes. *Journal of the American Statistical Association* 101.476 (2006), pp. 1566–1581.
- [12] C. Wang and D. M. Blei. A Split-Merge MCMC Algorithm for the Hierarchical Dirichlet Process. *CoRR* abs/1201.1657 (2012).

## A Proof of Lemma 1

Since  $[F(\mathbf{x})]_j$  only depends on  $x_1, \dots, x_{j-1}$ , we have

$$\begin{aligned}
 p(x_j \mid \mathbf{x}_{-j}) &\propto b(x_j; \alpha, \beta) \prod_{i: z_i \geq j} [F(\mathbf{x})]_{z_i} \\
 &= b(x_j; \alpha, \beta) \prod_{i: z_i = j} [F(\mathbf{x})]_j \prod_{i: z_i > j} [F(\mathbf{x})]_{z_i} \\
 &\propto b(x_j; \alpha, \beta) \prod_{i: z_i = j} x_j \prod_{i: z_i > j} (1 - x_j) \\
 &= b(x_j; \alpha, \beta) x_j^{n_j(\mathbf{z})} (1 - x_j)^{n_{>j}(\mathbf{z})}
 \end{aligned}$$

which gives the desired result.

### A.1 Remarks

Due to the sick-breaking interpretation of  $F(\mathbf{x})$ , it is not hard to see that

$$\sum_{\ell=1}^j [F(\mathbf{x})]_{\ell} + [P(\mathbf{x})]_j = 1, \quad \forall j \in \mathbb{N}.$$

where  $P : [0, 1]^{\mathbb{N}} \rightarrow [0, 1]^{\mathbb{N}}$  is defined by  $[P(\mathbf{x})]_j := \prod_{\ell \leq j} (1 - x_{\ell})$ . That is,  $[F(\mathbf{x})]_j = x_j [P(\mathbf{x})]_{j-1}$  hence  $[P(\mathbf{x})]_{j-1} < \tau$  implies  $[F(\mathbf{x})]_j < \tau$ . In other words,

$$\{j : [F(\mathbf{x})]_j \geq \tau\} \subseteq \{j : [P(\mathbf{x})]_{j-1} \geq \tau\}$$

and we can use the latter set to guarantee that we have enough atoms when truncating  $F(\mathbf{x})$  at level  $\tau$ . This is due to the fact the  $j \mapsto [P(\mathbf{x})]_j$  is nonincreasing in  $j$  as opposed to  $j \mapsto [F(\mathbf{x})]_j$  which is not necessarily monotone. Note that  $[P(\mathbf{x})]_{j-1}$  is easy to keep track of since it is the complement to the cumulative distribution associated with  $F(\mathbf{x})$  up to index  $j - 1$ . We take  $[P(\mathbf{x})]_0 = 1$ .