

---

# Supplementary Material: Ordered Stick-Breaking Prior for Sequential MCMC Inference of Bayesian Nonparametric Models

---

Mrinal Das<sup>†</sup>  
 Trapit Bansal<sup>†</sup>  
 Chiranjib Bhattacharyya<sup>†</sup>

MRINAL@CSA.IISC.ERNET.IN  
 TRAPIT@CSA.IISC.ERNET.IN  
 CHIRU@CSA.IISC.ERNET.IN

<sup>†</sup>Department of Computer Science and Automation,  
 Indian Institute of Science, Bangalore, India

## S.1. Introduction

We discuss some relevant mathematical background first in Section S.2, those are directly used in the paper or proofs. We include some examples and properties related to OSBP and PPFs of OSBP in Section S.3. There are two lemmas and four theorems in the paper. We prove them here in Sections S.4, S.5, S.6 related to OSBP, PPF of OSBP and SUMO respectively. We additionally provide one theorem (Theorem A) and a lemma (Lemma 1) which are strongly related to OSBP but could not be included in the paper due to space constraint. Then we provide construction of dependency over mini-batches using OSBP on PYP, SBP, HDP in Section S.7. Finally we give inference details for DPMM for text datasets in Section S.8.

## S.2. Mathematical background

This is not a comprehensive review, we cover only those definitions and properties that will be referred later in this material.

### S.2.1. Gamma distribution and Gamma process

**Definition 1.** (*Gamma Distribution*). A non-negative real-valued random variable  $X$  is said to have a Gamma distribution with shape parameter  $\alpha$  and scale parameter  $\beta$ , denoted by  $X \sim \text{Gamma}(\alpha; \beta)$ , if its probability density function is given by

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} e^{-x/\beta}}{\beta^\alpha \Gamma(\alpha)} \quad (1)$$

**Proposition 1.** Let  $X_1, X_2, \dots$  be a countable collection of independent Gamma distributed variables as  $X_k \sim \text{Gamma}(\alpha_k; \beta)$ . Then

$$\sum_{k=1}^{\infty} X_k \sim \text{Gamma}\left(\sum_{k=1}^{\infty} \alpha_k, \beta\right) \quad (2)$$

**Definition 2.** (*Gamma process*) A random measure  $G$  on  $\Omega$

is called a Gamma process with base measure  $H$  and scale parameter  $\alpha$ , denoted by  $G \sim \Gamma P(\alpha H)$ , if it satisfies

- for each measurable subset  $A \in \mathcal{B}$ ,  $G(A)$  has a Gamma distribution as  $G(A) \sim \text{Gamma}(\alpha H(A))$ , and
- $G$  is completely random

**Proposition 2.** If  $G_j \sim \Gamma P(\alpha_j H_j)$  for  $j = 1, \dots, k$ , then  $\sum_{j=1}^k G_j \sim \Gamma P(\sum_{j=1}^k \alpha_j H_j)$ .

### S.2.2. Dirichlet distribution and Dirichlet process

Let  $S_d$  denote the probability simplex in the  $d$ -dimensional real vector space  $\mathbf{R}_d$ , as

$$S_d = \{(x_1, \dots, x_d) \in \mathbf{R}_d : x_i \geq 0, \forall i; \sum_{i=1}^d x_i = 1\} \quad (3)$$

**Definition 3.** (*Dirichlet distribution*) An  $S_d$ -valued random variable  $X$  is said to have a Dirichlet distribution, denoted by  $X \sim \text{Dir}(\alpha_1, \dots, \alpha_d)$  with  $\alpha_1, \dots, \alpha_d > 0$ , if it has a probability density function given by

$$f(x_1, \dots, x_d; \alpha_1, \dots, \alpha_d) = \frac{\Gamma(\sum_{i=1}^d \alpha_i)}{\prod_{i=1}^d \Gamma(\alpha_i)} x_i^{\alpha_i-1} \quad (4)$$

**Proposition 3.** Let  $X_1, X_2, \dots, X_k$  be  $k$  independent Gamma distributed variables as  $X_j \sim \text{Gamma}(\alpha_j; \beta)$ . Then for  $Y_j = \frac{X_j}{\sum_{j=1}^k X_j}$ ,  $(Y_1, \dots, Y_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ .

*Proof.* This can be seen using the procedure of transformation of random variables.  $\square$

**Proposition 4.** If  $(x_1, \dots, x_{k-1}, x_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_{k-1}, \alpha_k)$ , then  $(y_1, \dots, y_{k-1}) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_{k-1})$ , where  $y_j = \frac{x_j}{\sum_{l=1}^{k-1} x_l}$  for  $j = 1, \dots, k-1$ .

*Proof.* Let  $Z_1, \dots, Z_k$  be  $k$  independent variables such that  $Z_j \sim \text{Gamma}(\alpha_j, \beta)$  for  $j = 1, \dots, k$ . Then  $(x_1, \dots, x_k) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$ , for  $x_j = \frac{Z_j}{\sum_{l=1}^k Z_l}$  using Proposition 3.

We can write,

$$y_j = \frac{x_j}{\sum_{l=1}^{k-1} x_l} = \frac{\frac{Z_j}{\sum_{r=1}^k Z_r}}{\sum_{l=1}^{k-1} \frac{Z_l}{\sum_{r=1}^k Z_r}} = \frac{Z_j}{\sum_{l=1}^{k-1} Z_l} \quad (5)$$

Thus, by Proposition 3,  $(y_1, \dots, y_{k-1}) \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_{k-1})$ .  $\square$

**Definition 4.** (*Dirichlet Process*). Let,  $\mathbb{H}$  is a probability measure over a measurable space  $(\Omega, \mathcal{B})$ , and  $\gamma$  is a positive real number. A random measure  $G$  on  $\Omega$  is called a Dirichlet process with base measure  $\mathbb{H}$ , denoted by  $G \sim DP(\gamma, \mathbb{H})$  if for any finite measurable partition  $(B_1, B_2, \dots, B_k)$  of  $\Omega$ ,

$$(G(B_1), \dots, G(B_k)) \sim \text{Dirichlet}(\gamma \mathbb{H}(B_1), \dots, \gamma \mathbb{H}(B_k)) \quad (6)$$

**Stick-breaking representation of DP.** (Sethuraman, 1994) proposed a stick-breaking construction of DP such that if  $G \sim DP(\gamma, \mathbb{H})$ , then

$$G = \sum_{j=1}^{\infty} \theta_j \delta_{\beta_j}, \quad \beta_j \sim \mathbb{H} \quad (7)$$

$\delta_{\beta_j}$  denotes an atomic distribution where the entire probability mass is concentrated at  $\beta_j$ .  $\{\theta_j\}$  are constructed as follows.

$$\theta_1 = v_1, \theta_j = v_j \prod_{l=1}^{j-1} (1 - v_l), \quad v_j \sim \text{Beta}(1, \gamma) \quad (8)$$

The above construction can be understood as breaking a unit length stick using stick fractions  $v_j$ . (Sethuraman, 1994) showed that  $\sum_{j=1}^{\infty} \theta_j = 1$  when  $\{\theta_j\}$  are constructed as above. Often  $\theta$  is said to be distributed as  $GEM(\gamma)$ .

### S.2.3. Generalized Dirichlet distribution

**Definition 5.** (*Generalized Dirichlet distribution*) An  $S_k$ -valued random variable  $X$  is said to have a generalized Dirichlet distribution, denoted by

$$X \sim GDD(a_1, b_1, \dots, a_{k-1}, b_{k-1}) \quad (9)$$

with  $a_j, b_j > 0, \forall j$  if it has a probability density function given by

$$f(x_1, \dots, x_k; a_1, b_1, \dots, a_{k-1}, b_{k-1}) = \left( \prod_{j=1}^{k-1} B(a_j, b_j) \right)^{-1} x_k^{b_{k-1}-1} \prod_{j=1}^{k-1} \left( x_j^{a_j-1} \left( \sum_{i=j}^k x_i \right)^{b_{j-1}-a_j-b_j} \right) \quad (10)$$

where  $x_k = 1 - \sum_{j=1}^k x_j$ .  $B(a_j, b_j) = \frac{\Gamma(a_j)\Gamma(b_j)}{\Gamma(a_j+b_j)}$ .

*Example.* Let  $k = 4$ , then the density function of  $(x_1, x_2, x_3, x_4)$  is

$$\left( \prod_{j=1}^3 B(a_j, b_j) \right)^{-1} x_1^{a_1-1} x_2^{a_2-1} x_3^{a_3-1} x_4^{b_3-1} (x_2 + x_3 + x_4)^{b_1-a_2-b_2} (x_3 + x_4)^{b_2-a_3-b_3} \quad (11)$$

**Proposition 5.** By setting  $b_{j-1} = a_j + b_j, 2 \leq j \leq k-1$  ( $b_0$  is arbitrary),  $X \sim GDD(a_1, b_1, \dots, a_{k-1}, b_{k-1})$  is equivalently  $X \sim \text{Dirichlet}(a_1, a_2, \dots, a_{k-1}, b_{k-1})$ .

*Proof.* This follows directly from Eq. (10) and Eq. (4).  $\square$

### S.2.4. Stick-breaking process

We have defined SBP in the paper, however we re-iterate the discussion to show one useful result in Lemma A regarding SBP.

Any almost sure (a.s.) discrete probability measure  $G$  is a stick-breaking process (SBP) (Ishwaran & James, 2001) if it can be represented as

$$G = \sum_{j=1}^{\infty} \theta_j \delta_{\beta_j}, \theta_1 = v_1, \theta_j = v_j \prod_{l=1}^{j-1} (1 - v_l) \\ a_j, b_j > 0, v_j \sim \text{Beta}(a_j, b_j), \beta_j \sim \mathbb{H} \quad (12)$$

$\mathbb{H}$  is a diffuse measure over a measurable space  $(\Omega, \mathcal{B})$  and  $\{a_j, b_j\}$  are set of parameters.

The following lemma gives a condition over  $\{a_j, b_j\}$  so that  $\sum_{j=1}^{\infty} \theta_j = 1$  a.s.

**Lemma A.** (Ishwaran & James, 2001). For the random weights in an SBP,  $\sum_{j=1}^{\infty} \theta_j = 1$  a.s. iff  $\sum_{j=1}^{\infty} \mathbb{E}[\log(1 - v_j)] = -\infty$ . Alternatively, it is sufficient to check that  $\sum_{j=1}^{\infty} \log(1 + \frac{a_j}{b_j}) = +\infty$ .

*Proof.* See appendix by Ishwaran & James (2001).  $\square$

**Important special cases.** SBP subsumes many well known BNP priors. When  $a_j = 1$  and  $b_j = \gamma$  for all  $j$ , SBP becomes  $DP(\gamma, \mathbb{H})$  following the constructive definition of Dirichlet process by Sethuraman (1994). Another popular BNP prior, the two parameter Poisson-Dirichlet process or Pitman-Yor process (PYP) (Pitman & Yor, 1997) can also be obtained as a special case when  $a_j = 1 - \lambda$  and  $b_j = \gamma + j\lambda$  for all  $j$ . There are many other existing priors which are special cases of SBP, see (Ishwaran & James, 2001) for a detailed discussion.

### S.3. Appearance in order and OSBP

In this section, we first give an example of appearance in order phenomenon, and then we recall the definition of OSBP, followed by one essential information about OSBP.

### S.3.1. Example of appearance in order

Here, we give an example of the appearance in order phenomenon defined in Section 2.1.

Let,  $t = 9$  and  $(Y_i)$  is  $(a, a, b, a, c, a, b, a, a)$ . Notice that,  $k_9 = 3$  with  $\{a, b, c\}$  as  $\{\bar{Y}_1, \bar{Y}_2, \bar{Y}_3\}$ . Now we have

$$B_1 = \{1, 2, 4, 6, 8, 9\}, B_2 = \{3, 7\}, \text{ and } B_3 = \{5\}$$

Then we say that it is appearing in order as

$$\begin{aligned} [9] - B_1 &= \{3, 7, 5\} \Rightarrow 3 \in B_2 \\ [9] - (B_1 \cup B_2) &= \{5\} \Rightarrow 5 \in B_3 \end{aligned}$$

whereas if  $\bar{Y} = \{a, c, b\}$  then

$$B_1 = \{1, 2, 4, 6, 8, 9\}, B_2 = \{5\} \text{ and } B_3 = \{3, 7\}$$

is *not* appearing in order as

$$[9] - B_1 = \{3, 7, 5\} \text{ but } 3 \notin B_2$$

### S.3.2. Definition of OSBP

As we will be referring to OSBP in later parts of this material, for the sake of easy reading we present them here again.

Let  $\Gamma$  be a *diffuse* probability measure over random measures, and  $\mu, \nu$  denote the set of scalar hyper-parameters  $\{\mu_j\}$  and  $\{\nu_j\}$  respectively such that  $0 < \mu_j < 1$ ,  $\nu_j > 0$ ,  $\forall j$ .  $(G_1, G_2, \dots)$  is an *appearing in order* sequence of random measures.  $(Q_1, \dots, Q_{k_{t-1}})$  is the set of  $k_{t-1}$  unique values among  $G_{1:t-1}$ . We define,  $G_1, G_2, \dots \sim \text{OSBP}(\mu, \nu, \Gamma)$  if  $G_1 \sim \Gamma$  and for any  $t \geq 2$ , the following holds:

$$\begin{aligned} G_t | G_{1:t-1}, (\rho_j), \Gamma &\sim \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} \Gamma \\ \rho_1 &= v_1, \quad \forall j > 1, \rho_j = v_j \prod_{l=1}^{j-1} (1 - v_l) \\ v_j | \mu_j, \nu_j &\sim \text{Beta}(\mu_j \nu_j, (1 - \mu_j) \nu_j) \\ \alpha_{k_{t-1}} &= 1 - \sum_{j=1}^{k_{t-1}} \rho_j \end{aligned} \quad (13)$$

### S.3.3. Diffuse base measure of OSBP ensures appearance in order

The need of the base measure  $\Gamma$  to be a diffuse measure is explained with the following Theorem.

**Theorem A.** *The samples from OSBP,  $(G_1, G_2, \dots) \sim \text{OSBP}(\mu, \nu, \Gamma)$  will be appearing in order almost surely iff the base measure of OSBP,  $\Gamma$  in Eq. (13) is a diffuse probability measure.*

*Proof.* When  $\Gamma$  is diffuse, for any two samples  $Q_j$  and  $Q_l$  sampled from  $\Gamma$  will be almost sure distinct iff  $j \neq l$ . By definition of OSBP, if for any  $t$ ,  $G_t \sim \Gamma$  then  $k_t = k_{t-1} + 1$  and  $Q_{k_t} = G_t$ . As  $\Gamma$  is diffuse measure,  $Q_{k_t}$  is a.s. distinct

from all  $Q_j$ ,  $j < k_t$ . Thus  $G_t \neq G_l$  for all  $l < t$ . Hence,  $[t] \setminus \bigcup_{l=1}^{k_t-1} B_l = t$  and  $B_{k_t} = [t]$ .

We show the sufficient condition by contradiction. Suppose,  $\Gamma$  is atomic. Let  $t = 4$ ,  $k_4 = 2$ , and  $Q_2 \neq Q_1$ . There are two partitions  $B_1$  and  $B_2$ . Now when  $G_5$  is sampled suppose it is sampled from  $\Gamma$ , then  $Q_3 = G_5$ . Then by definition of appearance in order  $[5] \setminus (B_1 \cup B_2)$  should be in  $B_3$  which is  $[5]$ . As  $\Gamma$  is atomic let  $Q_3 = Q_1$ . Then  $G_5$  becomes equal to  $G_1$  and so  $5 \in B_1$  and  $[5] \setminus (B_1 \cup B_2) = \emptyset$ . Contradiction. So, whenever a  $Q_j$  is sampled from  $\Gamma$ ,  $k_t$  must increase.  $k_t$  will increase iff  $Q_j$  is different from all  $Q_l$ ,  $l < j$ . Hence  $\Gamma$  has to be a diffuse probability measure.  $\square$

This is a slightly strict condition on the base measure than that for DP and PYP which also points out one key difference with the common BNP priors such as DP, PYP.

### S.3.4. Comparison with DP and PYP on PPF

DP and Pitman-Yor process (PYP) (Pitman & Yor, 1997) are the only other two existing SBP class of priors possessing PPFs. It is worth to note the difference of OSBP from DP, PYP in terms of PPFs due to modeling *appearance in order*. Recall that, PPF  $(\pi_j, j \in [k_{t-1}]$  and  $\sigma_{k_{t-1}})$  are defined by Pitman (1996) as

$$\begin{aligned} \pi_j &= p(z_t = j | z_{1:t-1}, \Theta), \quad j \in [k_{t-1}], \\ \sigma_{k_{t-1}} &= p(z_t = k_{t-1} + 1 | z_{1:t-1}, \Theta) \end{aligned} \quad (14)$$

where  $\Theta$  denotes the set of hyper-parameters. The PPFs corresponding to  $DP(\gamma, \mathbb{H})$ , also popularly referred as Chinese restaurant process (CRP) are

$$\begin{aligned} \pi_j &= \frac{g_j}{\gamma + t - 1}, \quad j \in [k_{t-1}], \\ \sigma_{k_{t-1}} &= \frac{\gamma}{\gamma + t - 1} \end{aligned} \quad (15)$$

where  $g_j = |B_j|$ , and  $B_j = \{i | z_i = j\}$ . Thus,  $\sum_{j=1}^{k_{t-1}} g_j = t - 1$  and  $\sum_{j=1}^{k_{t-1}} \pi_j + \alpha_{k_{t-1}} = 1$ . Similarly, PPFs of  $PYP(a, b, \mathbb{H})$  ( $0 \leq a < 1$  and  $b > -a$ ) are

$$\begin{aligned} \pi_j &= \frac{g_j - a}{b + t - 1}, \quad j \in [k_{t-1}], \\ \sigma_{k_{t-1}} &= \frac{b + a k_{t-1}}{b + t - 1} \end{aligned} \quad (16)$$

Note that,  $\sum_{j=1}^{k_{t-1}} g_j - a = t - 1 - a k_{t-1}$  and hence  $\sum_{j=1}^{k_{t-1}} \pi_j + \alpha_{k_{t-1}} = 1$ . By using  $a = 0$  and  $b = \gamma$ , PYP becomes equivalent to DP.

Notice from Theorem 3 that,  $\pi_j$  for OSBP can be written as  $a_j \prod_{l=1}^{j-1} (1 - a_l)$ , where  $a_j = \frac{\mu_j \nu_j + m_j - 1}{\nu_j + m_j + r_j - 1}$ . From Lemma 2,  $a_j$  is the posterior expectation of  $v_j$  conditioned on  $G_{1:t-1}$ . Thus in OSBP, the probability of joining partition  $B_j$  directly depends on *not* joining the partitions  $\{B_1, B_2, \dots, B_{j-1}\}$ . Whereas, in case of DP and PYP the probability of joining partition  $B_j$  depends only on the size

of  $B_j$  and the probabilities of joining partitions only loosely depend because of summing up to one.

Moreover, Corollary 1 of (Lee et al., 2013) shows that the only PPF which lead to an exchangeable sequence are those for which  $\pi_j$  is a function of partition  $B_j$  only. This is true for DP and PYP but not for OSBP. In Case of OSBP,  $\pi_j$  is a function of all the existing partitions ( $B_j$ ).

An important implication of this interpretation is that even though the partitions (and hence atoms) can be assumed to appear in order for DP and PYP, the effect of ordering of the partitions is lost due to the symmetric nature of this function which leads to an exchangeable partition probability function (EPPF by (Pitman, 1996)), or equivalently an exchangeable sequence illustrating another view why DP and PYP are not suitable to model appearance in order.

### S.3.5. Example related to Theorem 2

*Example.* Let  $\mu_j > 1/2$ ,  $\forall j$ , and  $\epsilon = 0.01$ . For  $k = 14$ ,  $\alpha_k \leq 0.01$  with probability more than 0.99.

## S.4. Proofs related to OSBP

### S.4.1. Proof of Theorem 1

**Theorem 1.** *If  $P_1 = \Gamma$ ,  $P_t = \sum_{j=1}^{k_t-1} \rho_j \delta_{Q_j} + \alpha_{k_t-1} \Gamma$  for  $t > 1$  and  $P^* = \sum_{j=1}^{\infty} \rho_j \delta_{Q_j}$  such that  $\sum_{j=1}^{\infty} \rho_j = 1$ , where  $(\rho_j)$ ,  $(Q_j)$ ,  $\alpha_{k_t}$  and  $\Gamma$  as defined in Eq. (13) with parameter  $\mu, \nu$ , then  $\lim_{t \rightarrow \infty} P_t = P^*$  a.s.*

*Proof.* By definition,  $k_t$  is the cardinality of the set  $(Q_1, Q_2, \dots, Q_{k_t})$ . So for any  $t > 0$ ,  $k_t = k_{t-1}$  if no new atom is sampled, and  $k_t = k_{t-1} + 1$  if a new atom is sampled from the base measure  $\Gamma$ . From Eq. (13), the probability of  $k_t = k_{t-1} + 1$  is  $\alpha_{k_{t-1}}$  and probability of  $k_t = k_{t-1}$  is  $\sum_{j=1}^{k_{t-1}} \rho_j$  which by definition is  $1 - \alpha_{k_{t-1}}$ . Hence, we get

$$k_{t-1} \leq k_t \quad a.s. \quad (17)$$

As,  $k_{t+1} \geq k_t$  a.s., and  $\alpha_{k_t} = 1 - \sum_{j=1}^{k_t} \rho_j$  by definition, with  $\rho_j > 0$  a.s. for all  $j$ , we get

$$\alpha_{k_{t-1}} \geq \alpha_{k_t} \quad a.s. \quad (18)$$

$k_t \geq k_{t-1}$  and not bounded above. For any  $K > 0$ , there is a  $t'$  such that  $k_{t'} > K$ , otherwise  $K$  is the upperbound of  $k_t$ . So we can say

$$\lim_{t \rightarrow \infty} k_t = \infty \quad a.s. \quad (19)$$

On the other hand,  $\alpha_{k_t} \leq \alpha_{k_{t-1}}$  and bounded below by zero. For any  $\epsilon > 0$  there is a  $t'$  such that  $\alpha_{k_{t'}} < \epsilon$ , otherwise  $\epsilon$  is the lower bound of  $\alpha_{k_t}$ . Hence,

$$\lim_{t \rightarrow \infty} \alpha_{k_t} = 0 \quad a.s. \quad (20)$$

Thus, we can write  $\lim_{t \rightarrow \infty} P_t = \lim_{t \rightarrow \infty} \sum_{j=1}^{k_t} \rho_j \delta_{Q_j} + \lim_{t \rightarrow \infty} \alpha_{k_t} \Gamma = \lim_{k_t \rightarrow \infty} \sum_{j=1}^{k_t} \rho_j \delta_{Q_j} + \lim_{t \rightarrow \infty} \alpha_{k_t} \Gamma = \sum_{j=1}^{\infty} \rho_j \delta_{Q_j} = P^*$ . That proves the Theorem.  $\square$

**Corollary 1.** *For  $t \in \mathbb{N}$  and  $\alpha_{k_t}$  as defined in OSBP,  $\lim_{k_t \rightarrow \infty} \alpha_{k_t} = 0$  a.s.*

*Proof.* This corollary is immediate from the above result. However we give one alternative proof here. Note that  $(1 + \frac{\mu_j \nu_j}{(1-\mu_j)\nu_j}) > 1$ , hence  $\sum_{j=1}^{\infty} \log(1 + \frac{\mu_j \nu_j}{(1-\mu_j)\nu_j}) = +\infty$ . By Lemma A it follows that  $\sum_{j=1}^{\infty} \rho_j = 1$  a.s. Therefore,

$$\lim_{k_t \rightarrow \infty} \alpha_{k_t} = 0 \quad a.s. \quad (21)$$

$\square$

### S.4.2. Proof of Lemma 1

**Lemma 1.** *For any  $t \in \mathbb{N}$ ,  $R_t = (\rho_1, \rho_2, \dots, \rho_{k_{t-1}}, \alpha_{k_{t-1}})$  as defined in Eq. (13) is distributed as generalized Dirichlet distribution (Connor & Mosimann, 1969). Furthermore, if  $(1 - \mu_{j-1})\nu_{j-1} = \nu_j$  for  $j$ ,  $2 \leq j \leq k_{t-1}$ , then  $R_t \sim \text{Dirichlet}(\mu_1 \nu_1, \mu_2 \nu_2, \dots, \mu_{k_{t-1}} \nu_{k_{t-1}}, (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}})$ .*

*Proof.* From Eq. (13), notice that  $v_j \sim \text{Beta}(\mu_j \nu_j, (1 - \mu_j) \nu_j)$  and  $(\rho_1, \rho_2, \dots, \rho_{k_{t-1}})$  is constructed by transforming  $(v_j)$ . Hence, Jacobian is  $\prod_{j=1}^{k_{t-1}} \left( \prod_{l=1}^{j-1} (1 - v_l) \right)^{-1}$ . Applying the transformation, we obtain the density function as

$$f_{R_t} = \left( \prod_{j=1}^{k_{t-1}} B(\mu_j \nu_j, (1 - \mu_j) \nu_j) \right)^{-1} \alpha_{k_t}^{(1 - \mu_{k_{t-1}}) \nu_{k_{t-1}} - 1} \prod_{j=1}^{k_{t-1}} \left( \rho_j^{\mu_j \nu_j - 1} \left( \sum_{i=j}^{k_{t-1}} \rho_i + \alpha_{k_{t-1}} \right)^{\kappa_j} \right)$$

where  $\alpha_{k_{t-1}} = 1 - \sum_{j=1}^{k_{t-1}} \rho_j$ .

$$B(\mu_j \nu_j, (1 - \mu_j) \nu_j) = \frac{\Gamma(\mu_j \nu_j) \Gamma((1 - \mu_j) \nu_j)}{\Gamma(\nu_j)}$$

and

$$\kappa_j = (1 - \mu_{j-1}) \nu_{j-1} - \mu_j \nu_j - (1 - \mu_j) \nu_j$$

Now let us write  $a_j = \mu_j \nu_j$  and  $b_j = (1 - \mu_j) \nu_j$ . Then Eq. (22) becomes equivalent to Eq. (10). Hence

$$R_t \sim GDD(\mu_1 \nu_1, (1 - \mu_1) \nu_1, \dots, \mu_{k_{t-1}} \nu_{k_{t-1}}, (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}})$$

This proves the first part.

We prove the second part as follows. When

$$(1 - \mu_{j-1}) \nu_{j-1} = \mu_j \nu_j + (1 - \mu_j) \nu_j = \nu_j$$

for  $2 \leq j \leq k_{t-1}$  by Proposition 5 we get

$$R_t \sim \text{Dirichlet}(\mu_1 \nu_1, \dots, \mu_{k_{t-1}} \nu_{k_{t-1}}, (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}})$$

### S.4.3. Proof of Theorem 2

**Theorem 2.** For  $\alpha_{k_t}$  as defined in Eq. (13) with parameters  $\mu, \nu$ , and any  $\epsilon \in (0, 1)$ , if  $\mu_j > 1/2$  for all  $j$ , then  $\alpha_k \leq \epsilon$  whenever  $k \geq \frac{2}{\log 2} \log \frac{1}{\epsilon}$  with probability more than  $1 - \epsilon$ .

*Proof.* From Eq. (13), using direct algebra one can rewrite  $\alpha_r = \prod_{j=1}^r (1 - v_j)$ , and using the independence of  $v_j$  we find that

$$\mathbb{E}[\alpha_r] = \prod_{j=1}^r \mathbb{E}[1 - v_j] \quad (22)$$

Using Markov inequality, we get

$$p(\alpha_{2r} > \frac{1}{2^r}) \leq 2^r \mathbb{E}[\alpha_{2r}] \quad (23)$$

Now using the fact that  $\mathbb{E}[1 - v_j] = 1 - \mu_j < 1/2$  and choosing a positive integer  $r$  such that  $\epsilon > \frac{1}{2^r}$  one obtains

$$p(\alpha_{2r} \geq \epsilon) \leq \epsilon \quad (24)$$

Then the proof follows by putting  $k = 2r$ .  $\square$

## S.5. Proofs related to PPF of OSBP

### S.5.1. Proof of Lemma 2

**Lemma 2.** Let,  $(v_j)$  be defined as in Eq. (13), and  $G_{1:t-1} | \mu, \nu, \Gamma \sim \text{OSBP}(\mu, \nu, \Gamma)$ . Then  $\forall j$ ,  $v_j | z_{1:t-1}, \mu_j, \nu_j \sim \text{Beta}(\mu_j \nu_j + g_j - 1, (1 - \mu_j) \nu_j + h_j)$ .

*Proof.* By definition of  $z$ , and OSBP, the following holds.

$$p(z_t = j | z_{1:t-1}, v_{1:k_{t-1}}) = v_j \prod_{l=1}^{j-1} (1 - v_l), \quad j \in [k_{t-1}]$$

$$p(z_t = k_{t-1} + 1 | z_{1:t-1}, v_{1:k_{t-1}}) = \prod_{l=1}^{k_{t-1}} (1 - v_l)$$

Now, following (Pitman, 1995), it is straight forward to see that,

$$p(z_1, \dots, z_{t-1} | v_{1:k_{t-1}}) = \prod_{j=1}^{k_{t-1}} \left( (1 - v_j)^{h_j} v_j^{g_j - 1} \right)$$

Now, we compute the posterior  $p(v_1, \dots, v_{k_{t-1}} | z_{1:t-1})$  as follows.

$$\propto p(z_1, \dots, z_{t-1} | v_{1:k_{t-1}}) p(v_1, \dots, v_{k_{t-1}})$$

$$\propto \prod_{j=1}^{k_{t-1}} \left( (1 - v_j)^{h_j} v_j^{g_j - 1} v_j^{\mu_j \nu_j - 1} (1 - v_j)^{(1 - \mu_j) \nu_j - 1} \right)$$

After marginalizing over all other  $v_l$ ,  $l \in [k_{t-1}] \setminus j$ , the lemma follows.  $\square$

### S.5.2. Proof of Theorem 3

**Theorem 3.** Let  $(\pi_j)$ ,  $\sigma_{k_{t-1}}$  be defined in Eq. (14), and  $G_{1:t-1} | \mu, \nu, \Gamma \sim \text{OSBP}(\mu, \nu, \Gamma)$ . Then, we have:

$$\pi_j = \frac{\mu_j \nu_j + g_j - 1}{\nu_j + g_j + h_j - 1} \prod_{l=1}^{j-1} \frac{(1 - \mu_l) \nu_l + h_l}{\nu_l + g_l + h_l - 1}, \quad j \in [k_{t-1}],$$

$$\sigma_{k_{t-1}} = \prod_{l=1}^{k_{t-1}} \frac{(1 - \mu_l) \nu_l + h_l}{\nu_l + g_l + h_l - 1} \quad (25)$$

*Proof.* By definition of OSBP, for  $1 \leq j \leq k_{t-1}$ ,  $p(G_t = Q_j | G_{1:t-1}, \{v_l\}) = p(z_t = j | z_{1:t-1}, \{v_l\}) = \rho_j$ . Now by definition of PPFs in Eq. (14), one can write  $\pi_j$  as

$$\mathbb{E}[\rho_j | z_{1:t-1}, \mu, \nu] = \mathbb{E} \left[ v_j \prod_{l=1}^{j-1} (1 - v_l) | z_{1:t-1}, \mu, \nu \right]$$

$$= \mathbb{E} [v_j | z_{1:t-1}, \mu_j, \nu_j] \prod_{l=1}^{j-1} \mathbb{E} [(1 - v_l) | z_{1:t-1}, \mu_l, \nu_l]$$

The second equation follows using the independence property of  $\{v_j\}$ . Following definition of OSBP, we similarly get  $\beta_{k_{t-1}}$  defined in Eq. (14),  $\beta_{k_{t-1}} = \prod_{l=1}^{k_{t-1}} \mathbb{E}[(1 - v_l) | z_{1:t-1}, \mu_l, \nu_l]$ . Theorem follows using Lemma 2.  $\square$

## S.6. Proofs related to SUMO

DPMM can be described as

$$\forall i x_i \sim f(\phi_i), \quad \phi_i | G \sim G, \quad G \sim DP(\gamma, H) \quad (26)$$

Using OSBP we propose following for DPMM

$$\forall t, \quad G_t | G_{1:t-1}, H \sim \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} \delta_{Q_{k_{t-1}+1}}$$

$$\forall i, \quad x_{ti} | \phi_{ti} \sim \text{mult}(\phi_{ti}), \quad \phi_{ti} | G_t \sim G_t \quad (27)$$

$(\rho_j)$ ,  $(Q_j)$  and  $\alpha_{k_{t-1}}$  are as defined in OSBP. Each  $G_t$  takes value from  $(Q_1, \dots, Q_{k_{t-1}}, Q_{k_{t-1}+1})$ , where  $Q_1 = G_1$  and all other  $Q_j$  are sampled from  $DP(\gamma_j, H)$ . The second line in Eq. (27) models DPMM with  $G_t$  similar to Eq. (26).

**Parameter settings.** Regarding the parameters  $(\mu_j, \nu_j, \gamma_j)$ , we set for all  $j$   $\mu_j = \mu$ , for some  $0.5 < \mu < 1$ .  $\nu_j = (1 - \mu) \nu_{j-1}$  and  $\nu_1 = \gamma$  ( $\gamma > 0$  as in Eq. (26)). Hence  $\nu_j = (1 - \mu)^{j-1} \gamma$ .

We use,  $\gamma_j = \mu \nu_j$  and hence  $\gamma_j = \mu (1 - \mu)^{j-1} \gamma$ . Thus, there are two parameters  $\mu$  and DPMM parameter  $\gamma$ .

**Equivalence with DPMM.** Before prove the Theorem 4, we need to revise one useful result as follows.

**Theorem B.** Let,  $Q_j \sim DP(\gamma_j, H_j)$  for  $j = 1, \dots, k$  and  $(c_1, \dots, c_k) \sim \text{Dirichlet}(\gamma_1, \gamma_2, \dots, \gamma_k)$  be independent of  $Q_1, \dots, Q_k$ , then

$$\sum_{j=1}^k c_j Q_j \sim DP\left(\sum_{j=1}^k \gamma_j, \frac{\sum_{j=1}^k \gamma_j H_j}{\sum_{j=1}^k \gamma_j}\right) \quad (28)$$

*Proof.*  $g_j \sim \text{Gamma}(\gamma_j, \beta)$  for  $j = 1, \dots, k$  independently. Let,  $G_j = g_j Q_j$ , then  $G_1, \dots, G_k$  are independent Gamma processes with  $G_j \sim \Gamma P(\gamma_j H_j)$ .

Let  $G' = \sum_{j=1}^k G_j$ , then by Proposition 2,  $G' \sim \Gamma P(\sum_{j=1}^k \gamma_j H_j)$ . Let  $g' = \sum_{j=1}^k g_j$ , then

$$\frac{G'}{G'(\Omega)} = \frac{G'}{\sum_{j=1}^k G_j(\Omega)} = \frac{G'}{\sum_{j=1}^k g_j Q_j(\Omega)} = \frac{G'}{g'}$$

Hence,  $\frac{G'}{g'}$  is a normalized Gamma process, hence

$$\frac{G'}{g'} \sim DP\left(\sum_{j=1}^k \gamma_j, \frac{\sum_{j=1}^k \gamma_j H_j}{\sum_{j=1}^k \gamma_j}\right)$$

Again, if  $c_j = \frac{g_j}{g'}$ , then

$$(c_1, \dots, c_k) \sim \text{Dirichlet}(a_1, \dots, a_k)$$

and we can say

$$\frac{G'}{g'} = \sum_{j=1}^k \frac{g_j}{g'} Q_j = \sum_{j=1}^k c_j Q_j$$

Thus  $\sum_{j=1}^k c_j Q_j \sim DP\left(\sum_{j=1}^k \gamma_j, \frac{\sum_{j=1}^k \gamma_j H_j}{\sum_{j=1}^k \gamma_j}\right)$ .  $\square$

### S.6.1. Proof of Theorem 4

**Theorem 4.** For any  $t \in \mathbb{N}$ , each  $x_{ti}$  sampled using model Eq. (27) has marginal distribution same as  $x_i$  sampled with DPMM in Eq. (26) with  $G \sim DP(c_t, H)$ , where  $c_t = \sum_{j=1}^{k_{t-1}} \gamma_j + (1 - \mu)^{k_{t-1}} \gamma$ . Furthermore, for any  $\epsilon > 0$  and  $t > 0$ , with probability greater than  $1 - \epsilon$ , each  $x_{ti}$  in Eq. (27) has marginal distribution same as  $x_i$  in Eq. (26) with  $G \sim DP(\sum_{j=1}^k \gamma_j, H)$ , when  $k_t \geq k \geq \frac{2}{\log 2} \log \frac{1}{\epsilon}$ . Also, for  $t \rightarrow \infty$ , each  $x_{ti}$  in Eq. (27) has marginal distribution same as  $x_i$  in Eq. (26) with  $G \sim DP(\gamma, H)$ .

*Proof.* There are three parts of the Theorem, we prove them one by one.

Let us write,

$$R_t = (\rho_1, \rho_2, \dots, \rho_{k_{t-1}}, \alpha_{k_{t-1}})$$

then by Lemma 1,  $R_t$  is GDD distributed.

Now, we use

$$\gamma_j = \mu \nu_j = \mu(1 - \mu)^{j-1} \nu$$

that satisfies

$$\nu_j = (1 - \mu_{j-1}) \nu_{j-1}$$

for  $j > 1$ . Hence, by Lemma 1,

$R_t \sim \text{Dirichlet}(\mu_1 \nu_1, \mu_2 \nu_2, \dots, \mu_{k_{t-1}} \nu_{k_{t-1}}, (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}}) DP(\sum_{j=1}^k \gamma_j, H)$  for  $k \geq \frac{2}{\log 2} \log \frac{1}{\epsilon}$ . Now, for  $t$  such that

Now we can apply Theorem B, and by that

$$G_t \sim DP(c_t, H')$$

$$H' = \frac{(\sum_{j=1}^{k_{t-1}} \gamma_j + (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}}) H}{\sum_{j=1}^{k_{t-1}} \gamma_j + (1 - \mu_{k_{t-1}}) \nu_{k_{t-1}}} = H$$

For any  $t$ , using the parameter setting we get

$$G_t \sim DP(c_t, H) \quad (29)$$

$$c_t = \left( \sum_{j=1}^{k_{t-1}} \gamma_j + (1 - \mu)^{k_{t-1}} \gamma \right)$$

Hence the marginal distribution of each  $x_{ti}$  sampled from model in Eq. (27) is always equivalent to that of DPMM. The difference will be in scale.

Then the first part of the Theorem follows immediately from Eq. (27) (second line) and Eq. (26).

For the second part, let us write

$$\begin{aligned} & \sum_{j=1}^k \rho_j \delta_{Q_j} + \alpha_k \delta_{Q_{k+1}} \\ &= \left( \sum_{l=1}^k \rho_l \right) \left( \sum_{j=1}^k \frac{\rho_j}{\sum_{l=1}^k \rho_l} \delta_{Q_j} + \frac{\alpha_k}{\sum_{l=1}^k \rho_l} \delta_{Q_{k+1}} \right) \\ &= (1 - \alpha_k) \left( \sum_{j=1}^k \frac{\rho_j}{\sum_{l=1}^k \rho_l} \delta_{Q_j} + \frac{\alpha_k}{1 - \alpha_k} \delta_{Q_{k+1}} \right) \\ &= (1 - \alpha_k) \left( \sum_{j=1}^k \frac{\rho_j}{\sum_{l=1}^k \rho_l} \delta_{Q_j} \right) + \alpha_k \delta_{Q_{k+1}} \quad (30) \end{aligned}$$

Now we can say

$$P_t = \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} \delta_{Q_{k_{t-1}+1}}$$

is a mixture of two distributions  $J_{k_{t-1}}$  and  $\delta_{Q_{k_{t-1}+1}}$  for any  $t$ , where

$$J_{k_{t-1}} = \sum_{j=1}^{k_{t-1}} \frac{\rho_j}{\sum_{l=1}^{k_{t-1}} \rho_l} \delta_{Q_j}$$

. Probability of sampling  $G_t$  from  $J_{k_{t-1}}$  is  $1 - \alpha_{k_{t-1}}$ .

Again as,

$$(\rho_1, \dots, \rho_{k_{t-1}}, \alpha_{k_{t-1}}) \sim \text{Dir}(\gamma_1, \dots, \gamma_{k_{t-1}}, (1 - \mu)^{k_{t-1}} \gamma)$$

by Proposition 4,

$$\left( \frac{\rho_1}{1 - \alpha_{k_{t-1}}}, \dots, \frac{\rho_{k_{t-1}}}{1 - \alpha_{k_{t-1}}} \right) \sim \text{Dirichlet}(\gamma_1, \dots, \gamma_{k_{t-1}})$$

Hence  $J_{k_{t-1}}$  is equivalent to

$$DP\left(\sum_{j=1}^{k_{t-1}} \gamma_j, H\right)$$

Now, from Theorem 2, for  $k_{t-1} \geq \frac{2}{\log 2} \log \frac{1}{\epsilon}$ ,  $\alpha_{k_{t-1}} < \epsilon$  for any  $\epsilon > 0$  with probability at least  $1 - \epsilon$ .

Hence, with probability at least  $1 - \epsilon$  we sample  $G_t$  from  $DP(\sum_{j=1}^k \gamma_j, H)$  for  $k \geq \frac{2}{\log 2} \log \frac{1}{\epsilon}$ . Now, for  $t$  such that

$k_t \geq k$ ,  $x_{ti}$  sampled from model in Eq. (27), is marginally equivalent to  $x_i$  sampled from model DPMM in Eq. (26) with  $G \sim DP(\sum_{j=1}^k \gamma_j, H)$ .

For the third part, from Theorem 1, as  $t \rightarrow \infty$ ,  $P_t \rightarrow P^* = \sum_{j=1}^{\infty} \rho_j \delta_{Q_j}$  such that  $\sum_{j=1}^{\infty} \rho_j = 1$  a.s. By the first part,  $P^*$  becomes equivalent to  $DP(\sum_{j=1}^{\infty} \gamma_j, H)$ .

Now we can write

$$\sum_{j=1}^{\infty} \gamma_j = \mu \gamma \sum_{j=1}^{\infty} (1 - \mu)^{j-1} = \frac{\mu \gamma}{1 - (1 - \mu)} = \gamma$$

Hence,  $P^*$  becomes equivalent to  $DP(\gamma, H)$ .

Thus, for  $t$ ,  $P_t$  going to  $P^*$ ,  $x_{ti}$  sampled from model in Eq. (27), is marginally equivalent to  $x_i$  sampled from model DPMM in Eq. (26) with  $G \sim DP(\gamma, H)$ . This proves the Theorem.  $\square$

## S.7. Application of OSBP on other BNP models

Recall that, in the mini-batch setup, we consider a streaming dataset as  $(X_t) = (X_1, X_2, \dots, X_{\bar{d}})$ , where  $X_t = \{x_i\}_{i=\bar{n}(t-1)+1}^{\bar{n}t}$ . For clarity in notation, we represent  $X_t$  as  $\{x_{ti}\}_{i=1}^{\bar{n}}$ . We will describe application of OSBP on Pitman-Yor process, stick-breaking process and higherarchical Dirichlet process to share information across mini-batches. However, unlike DPMM the equivalence relationship is not straight forward to prove in these cases, and are left for future work.

We will consider  $H$  as a probability measure over a measurable space  $(\Omega, \mathcal{B})$ , and  $f(\cdot)$  is the distribution for the data model.  $(\rho_j)$  and  $\alpha_{k_{t-1}}$  are as defined in OSBP.

The inference mechanism follows from Theorem 3 and the inference techniques for the corresponding models. It is straight forward to derive them as outlined in Section 3.2 for DPMM. We do not describe them here.

### S.7.1. OSBP on Pitman-Yor process

Pitman-Yor process (PYP) (Pitman & Yor, 1997) can be described following the stick-breaking representation as follows. For  $0 \leq a < 1$ , and  $b > -a$ , any random probability measure  $G \sim PYP(a, b, H)$  if

$$G = \sum_{j=1}^{\infty} \theta_j \delta_{\beta_j}, \quad \theta_1 = v_1, \quad \theta_j = v_j \prod_{l=1}^{j-1} (1 - v_l) \\ \forall j, \quad v_j \sim \text{Beta}(1 - a, b + ja); \quad \beta_j \sim H \quad (31)$$

PYP mixture model can be described as follows.

$$\forall i, \quad x_i \sim f(\phi_i); \quad \forall i, \quad \phi_i | G \sim G \quad (32)$$

<sup>1</sup>for simplicity we have assumed  $n = \bar{n}\bar{d}$ .

Now we apply OSBP on PYP mixture model to get

$$\forall t, \quad G_t | Q_{1:k_{t-1}}, H \sim \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} PYP(a, b, H) \\ \forall i, \quad x_{ti} | \phi_{ti} \sim f(\phi_{ti}), \quad \forall i, \quad \phi_{ti} | G_t \sim G_t \quad (33)$$

Application of OSBP on PYP is similar to that of DP. The inference follows from the Theorem 3, and PPFs of PYP.

### S.7.2. OSBP on stick-breaking process

Recall that, stick-breaking process (SBP) (Ishwaran & James, 2001) can be described following the stick-breaking representation as follows. Let,  $a_j, b_j > 0$ , and  $a = (a_1, a_2, \dots), b = (b_1, b_2, \dots)$ . Any random probability measure  $G \sim SBP(a, b, H)$  if following holds.

$$G = \sum_{j=1}^{\infty} \theta_j \delta_{\beta_j}, \quad \theta_1 = v_1, \quad \theta_j = v_j \prod_{l=1}^{j-1} (1 - v_l) \\ \forall j, \quad v_j \sim \text{Beta}(a_j, b_j); \quad \beta_j \sim H \quad (34)$$

SBP mixture model can be described as follows.

$$\forall i, \quad x_i \sim f(\phi_i); \quad \forall i, \quad \phi_i | G \sim G \quad (35)$$

Imposing OSBP on SBP mixture model yields the following model.

$$\forall t, \quad G_t | Q_{1:k_{t-1}}, H \sim \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} SBP(a, b, H) \\ \forall i, \quad x_{ti} | \phi_{ti} \sim f(\phi_{ti}), \quad \forall i, \quad \phi_{ti} | G_t \sim G_t \quad (36)$$

SBP being generalized version subsumes many BNP priors including DP and PYP. The construction of OSBP based sequential model for DP, PYP and SBP are similar. Essentially, following this structure it is easy to build such sequential models for a wide range of BNP models. The inference will follow from Theorem 3 and inference procedure of SBP. Unfortunately, SBP does not have PPFs and truncated methods are applied (Ishwaran & James, 2001).

### S.7.3. OSBP on hierarchical Dirichlet process

Hierarchical Dirichlet process (HDP) (Teh et al., 2006) is defined as below for  $\gamma, \lambda > 0$

$$G_0 \sim DP(\gamma, H) \\ \forall i, \quad G_i \sim DP(\lambda, G_0) \quad (37)$$

HDP assumes grouped data, that  $x_i$  represent a group which consists of data points  $\{x_{il}\}$ . HDP mixture model can be described as

$$G_0 \sim DP(\gamma, H) \\ \forall i, \quad G_i \sim DP(\lambda, G_0) \\ \forall l, \quad x_{il} \sim f(\phi_{il}); \quad \phi_{il} | G_i \sim G_i \quad (38)$$

Imposing OSBP on HDP mixture model by using the base measure  $\Gamma$  of OSBP as  $DP(\gamma, H)$ , we get

$$\begin{aligned} \forall t, G_t | Q_{1:k_{t-1}}, H &\sim \sum_{j=1}^{k_{t-1}} \rho_j \delta_{Q_j} + \alpha_{k_{t-1}} DP(\gamma, H) \\ \forall i, G_{ti} &\sim DP(\lambda, G_t) \\ \forall l, x_{til} | \phi_{til} &\sim f(\phi_{til}), \quad \phi_{til} | G_{ti} \sim G_{ti} \end{aligned} \quad (39)$$

The inference will follow from Theorem 3 and the prediction rule for HDP, Chinese restaurant franchise (CRF) (Teh et al., 2006).

## S.8. SUMO for DPMM on document clustering

We describe SUMO here for text datasets. Each data point  $x_i$  is a document which has multiple examples (words). A words in document  $i$  is denoted by  $x_{il}$ . The data model is  $x_{il} \sim \text{mult}(\phi_i)$ . In order to maintain conjugacy,  $\phi_i$  has Dirichlet prior.

( $\phi_i$ ) are sampled from ( $Q_j$ ) where  $Q_j \sim DP(\gamma_j, H)$  in Eq. (27). We can say  $Q_j = \sum_{r=1}^{\infty} \zeta_{jr} \delta_{\psi_{jr}}$  following Eq. (12), where ( $\zeta_{jr}$ ) form the stick-breaking weights and atoms are  $\psi_{jr}$ . Let,  $\{\beta_s\}$  is set of global components. Then each  $\psi_{jr} \in \{\beta_s\}$  ensures same components across  $t$ . We can create global components by ad-hoc merging of components across  $t$ . But we prefer a more technical approach of using a.s. discrete  $H$  by  $H \sim DP(\lambda, \text{Dirichlet}(\eta))$ . We can write  $H = \sum_{s=1}^{\infty} \theta_s \delta_{\beta_s}$ , where  $\beta_s \sim \text{Dirichlet}(\eta)$  and ( $\theta_s$ ) form the stick-breaking weights.

Given this setup, we introduce alternative variables to speed up the mixing of the Markov chain following standard approach. Recall that,  $z_t = j$  if  $G_t = Q_j$ . Let,  $a_{ti} = r$  if  $\phi_{ti} = \psi_{jr}$  and  $z_t = j$ . So  $r$  is the index of the mixture component in prior  $G_t$  assigned to document  $i$  of mini-batch  $t$ . If  $s$  is the index of global mixture component represented by  $\psi_{jr}$  in  $Q_j$ , then we define  $b_{jr} = s$  if  $\psi_{jr} = \beta_s$ . Furthermore, let  $y_{ti} = s$  if  $z_t = j$  and  $b_{jr} = s$ .  $y_{ti}$  is the index of the global mixture component assigned to document  $i$  in mini-batch  $t$ .  $\phi$  and  $\psi$  can be retrieved from  $z, a, b$  and  $\beta$ .

Due to this representation, the equivalent random quantities are  $A_{1:t} = \{\{a_{li}\}_{i=1}^{\bar{n}}\}_{l=1}^t$ ,  $B_{1:k_t} = \{b_{jr}\}_{j=1}^{k_t}$ , and  $Y_{1:t} = \{\{y_{li}\}_{i=1}^{\bar{n}}\}_{l=1}^t$ . We integrate out ( $Q_j$ ) and  $H$  following Chinese restaurant process (CRP), ( $\rho_j$ ) following Theorem 3, and  $\{\beta_s\}$  following Dirichlet multinomial conjugacy. So, we need to infer  $A_t, B$ , and  $z_t$  at time  $t$  after observing  $X_t$ . The posterior of  $\{\beta_s\}$  and other variables can be retrieved after the inference through  $a, b, z$  and ( $X_t$ ).

*Notation.* Superscript with hyphen denotes set minus, e.g.  $X_t^{-i} = X_t \setminus x_{ti}$ , and  $X_{tr}^{-r} = X_{tr} \setminus x_{tr}$ , where  $X_{tr} = \{x_{ti} | a_{ti} = r\}$ .  $X_{1:t}^{-tr} = X_{1:t} \setminus X_{tr}$ , and  $X_{1:t}^{-ti} = X_{1:t} \setminus x_{ti}$ .  $A_{1:t}^{-ti} = A_{1:t} \setminus a_{ti}$ .  $B_{z_t}^{-r} = B_{z_t} \setminus b_{z_t r}$ .  $L_s(x_{ti})$  and  $L_s(X_{tr})$  are the likelihood of  $x_{ti}$  and  $X_{tr}$  respectively for mixture component  $s$ .

**Recursive computation of likelihood.**  $L_s(x_{ti})$  is the likelihood of  $x_{ti}$  under mixture component  $s$ , that is  $L_s(x_{ti}) = p(x_{ti} | Y_{1:t}, X_{1:t-1}, X_t^{-i})$ . After observing  $X_{1:t-1}$  and  $X_t^{-i}$ ,  $L_s(x_{ti})$  can be computed by recursively applying Bayes theorem using Dirichlet multinomial conjugacy as follows.

$$\begin{aligned} p(x_{ti} | Y_{1:t}, X_{1:t-1}, X_t^{-i}) &= \\ \int \prod_f p(x_{tif} | y_{ti} = s, \beta_s) p(\beta_s | X_{t-1}^{-i}, X_{1:t-1}, Y_{1:t}) d\beta_s \\ &= \int \prod_f \beta_{sx_{tif}} \frac{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}))}{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i})} \prod_v \beta_{sv}^{\eta_v + C_{sv} + c_{sv}^{-i} - 1} d\beta_s \\ &= \frac{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}))}{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i})} \frac{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i} + c_{sv}^i)}{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}) + c_{sv}^i)} \\ \int \frac{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}) + c_{sv}^i)}{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i} + c_{sv}^i)} \prod_v \beta_{sv}^{\eta_v + C_{sv} + c_{sv}^{-i} + c_{sv}^i - 1} d\beta_s \\ &= \frac{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}))}{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i})} \frac{\prod_v \Gamma(\eta_v + C_{sv} + c_{sv}^{-i} + c_{sv}^i)}{\Gamma(\sum_v (\eta_v + C_{sv} + c_{sv}^{-i}) + c_{sv}^i)} \end{aligned} \quad (40)$$

Integration happens following the property that  $\beta_s \sim \text{Dirichlet}(\eta)$  and using Dirichlet multinomial conjugacy. Please refer to the Appendix for detailed steps. We define the sufficient statistics as below.

$$\begin{aligned} C_{sv} &= \sum_{l=1}^{t-1} \sum_{i=1}^{\bar{n}} \sum_f \mathbb{I}[y_{li} = s, x_{lif} = v] \\ c_{sv} &= \sum_{i=1}^{\bar{n}} \sum_f \mathbb{I}[y_{ti} = s, x_{tif} = v] \\ c_{sv}^{-i} &= \sum_{q=1, q \neq i}^{\bar{n}} \sum_f \mathbb{I}[y_{tq} = s, x_{lqf} = v] \\ c_{sv}^i &= \sum_f \mathbb{I}[y_{ti} = s, x_{lif} = v] \end{aligned} \quad (41)$$

Similarly, we compute  $L_s(X_{tr})$  the likelihood of  $X_{tr}$  for mixture component  $s$ ,  $p(X_{tr} | Y_{1:t}, X_{1:t-1}, X_{tr}^{-r})$  as follows.

$$\begin{aligned} \int \prod_{i=1: a_{ti}=r}^{\bar{n}} \prod_f p(x_{tif} | y_{ti} = s, \beta_s) p(\beta_s | X_{1:t-1}, Y_{1:t}) d\beta_s \\ = \int \prod_{i=1: a_{ti}=r}^{\bar{n}} \prod_f \beta_{sx_{tif}} \frac{\Gamma(\sum_v (\eta_v + c_{sv} + c_{sv}^{-r}))}{\prod_v \Gamma(\eta_v + c_{sv} + c_{sv}^{-r})} \\ \prod_v \beta_{sv}^{\eta_v + c_{sv} + c_{sv}^{-r} - 1} d\beta_s \\ = \frac{\Gamma(\sum_v (\eta_v + c_{sv} + c_{sv}^{-r}))}{\prod_v \Gamma(\eta_v + c_{sv} + c_{sv}^{-r})} \frac{\prod_v \Gamma(\eta_v + c_{sv} + c_{sv}^{-r} + c_{sv}^r)}{\Gamma(\sum_v (\eta_v + c_{sv} + c_{sv}^{-r}) + c_{sv}^r)} \end{aligned}$$

We define the required sufficient statistics as below.

$$\begin{aligned} c_{sv}^r &= \sum_{i=1}^{\bar{n}} \sum_f \mathbb{I}[a_{ti} = r, z_t = j, b_{jr} = s, x_{tif} = v] \\ c_{sv}^{-r} &= \sum_{i=1}^{\bar{n}} \sum_f \mathbb{I}[a_{ti} = q, q \neq r, b_{z_t q} = s, x_{tif} = v] \end{aligned} \quad (42)$$

**Inference of  $a$ .** We infer  $a$  as below.

$$p(a_{ti} = r | A_{1:t}^{-ti}, B_{1:k_t}, z_{1:t}, X_{1:t}) \propto \quad (43)$$

$$p(x_{ti} | a_{ti} = r, z_{1:t}, A_{1:t}^{-i}, B_{1:k_t}, X_{1:t}^{-i}) p(a_{ti} = r | A_{1:t}^{-ti}, z_t)$$

where  $p(x_{ti} | a_{ti} = r, z_{1:t}, A_{1:t}^{-i}, B_{1:k_t}, X_{1:t}^{-i})$  is  $L_{b_{z_t r}}(x_{ti})$ .  $p(a_{ti} = r | A_{1:t}^{-ti}, z_t)$  comes from CRP as

$$\propto L_{b_{z_t r}}(x_{ti}) (m_{z_t r}^{-i} + M_{z_t r}) (1 - \iota_r) + \gamma_{b_{z_t r}} L_{b_{z_t r \text{new}}}(x_{ti}) \iota_r \quad (44)$$

$$\iota_r = \mathbb{I}[r = r_{\text{new}}], \quad m_{z_t r} = \sum_{i=1}^{\bar{n}} \mathbb{I}[a_{ti} = r],$$

$$M_{j_r} = \sum_{l=1}^{t-1} \sum_{i=1}^{\bar{n}} \mathbb{I}[z_l = j, a_{li} = r] \quad (45)$$

$m_{z_t r}$  denotes the number of time component  $\psi_{j_r}$  is assigned in the current mini-batch, whereas  $M_{j_r}$  denotes how



many times  $\psi_{jr}$  is assigned across all the mini-batches seen so far excluding the current mini-batch. When a new  $r_{new}$  is sampled we obtain  $b_{z_t r_{new}}$  from  $p(b_{z_t r} = s_{new} | z_{1:t}, A_{1:t}, B_{1:k_t}, X_{1:t})$  which is shown later.

**Inference of  $z$ .** Following the dependence structure in Eq. (27),  $z_t$  is independent of  $X_t$  given  $Y_t$ . So, we can infer  $z$  from  $p(z_t = j | z_{1:t-1}, Y_t, B_{1:k_t})$  as

$$\propto \left[ \prod_{i=1}^{\bar{n}} p(y_{ti} = s | z_{1:t}, B_{1:k_t}) \right] p(z_t = j | z_{1:t-1}) \quad (46)$$

$p(z_t = j | z_{1:t-1})$  comes from Theorem 3. Recall that  $y_{ti} = b_{z_t a_{ti}}$ . So  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1})$  comes from CRP by integrating out  $G_t$  and  $H$ .

Let,  $\iota_j = \mathbb{I}[z_t = j_{new}]$ ,  $\iota_s^j = \mathbb{I}[z_t = j, s = s_{new}]$ ,  $\iota_s^0 = \prod_{l=1}^{k_t} \iota_s^l$ ,  $J_{js} = \sum_r \mathbb{I}[b_{jr} = s, z_t = j]$  and  $J_{.s} = \sum_{j=1}^{k_t-1} J_{js}$ .  $\iota_s^j$ ,  $\iota_s^0$  denote if  $\beta_s$  is present in  $Q_j$ ,  $H$  respectively.  $J_{js}$  counts number of times  $\beta_s$  is present among  $\{\psi_{jr}\}$ .

Notice that, when  $\iota_s^j = 0$ ,  $\iota_s^0$  must be 0 and that implies the situation that the global component  $\beta_s$  is present in  $Q_j$ . When,  $\iota_s^j = 1$ ,  $\iota_s^0 = 0$  signifies that  $\beta_s$  is not present in  $Q_j$ , but is present in  $H$ . Whereas  $\iota_s^0 = 1$  implies  $\iota_s^j = 1$  and  $\beta_s$  is not present in any prior. When,  $\iota_j = 1$ ,  $\iota_s^j$  must be 1, but  $\iota_s^0$  may be 1 or 0. Hence there are following scenarios.

i.  $\iota_j = 0$ ,  $\iota_s^j = 0$ : then we can say  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1}) \propto J_{js}$ .

ii.  $\iota_j = 0$ ,  $\iota_s^j = 1$ ,  $\iota_s^0 = 0$ : then we need to sample a global component from  $H$  which is proportional to  $\gamma_j J_{.s}$ .  $J_{.s} = \sum_{j=1}^{k_t-1} \sum_r \mathbb{I}[b_{jr} = s, z_t = j]$ , sum over all existing priors. So  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1}) \propto \lambda J_{.s}$ .

iii.  $\iota_j = 0$ ,  $\iota_s^j = 1$ ,  $\iota_s^0 = 1$ : then we need to sample a new global component from  $Dirichlet(\eta)$  which is proportional to  $\lambda$ . So  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1}) \propto \lambda \gamma_j$ .

iv.  $\iota_j = 1$ ,  $\iota_s^j = 1$ ,  $\iota_s^0 = 0$ : then we need to sample a new global component from  $Dirichlet(\eta)$  which is proportional to  $\lambda$ . So  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1}) \propto J_{.s}$ .  $\gamma_j$  does not appear here as there is not  $Q_j$  and no corresponding CRP.

v.  $\iota_j = 1$ ,  $\iota_s^j = 1$ ,  $\iota_s^0 = 1$ : then we need to sample a new global component from  $Dirichlet(\eta)$  which is proportional to  $\lambda$ . So  $p(y_{ti} = s | z_t = j, B_{1:t}, z_{1:t-1}) \propto \lambda$ .

Combining them together we get  $p(z_t = j | z_{1:t-1}, Y_t, B_{1:k_t})$

$$\propto \left[ \prod_{i=1}^{\bar{n}} J_{js} (1 - \iota_s^j) + \gamma_j \iota_s^j (J_{.s} (1 - \iota_s^0) + \lambda \iota_s^0) \right] \pi_j (1 - \iota_j) + \left[ J_{.s} (1 - \iota_s^0) + \lambda \iota_s^0 \right] \sigma_{k_t-1} \iota_j \quad (47)$$

**Algorithm 2** SUMO for DPMM on text datasets.

**Require:**  $(X_t)$ ,  $\mu$ ,  $\lambda$ ,  $\gamma$  and  $\eta$

```

1: for  $t = 1, 2, \dots$  do
2:   Initialize global component assignments  $Y_t$ 
3:   for  $iter = 1$  to  $I$  do
4:     Sample  $z_t$  from  $p(z_t | z_{1:t-1}, Y_t, J)$ 
5:     for  $i = 1$  to  $\bar{n}$  do
6:       Sample  $a_{ti}$  from  $p(a_{ti} = r | A_t^{-i}, z_t, X_t, M, C)$ 
7:     end for
8:     Sample  $B_{z_t}$  from  $p(B_{z_t} = s | B_{z_t}^{-r}, z_t, X_t, M, N, C)$ 
9:   end for
10:  Compute  $c, m, n$  and update  $C, M, N$ , and  $J$ 
11:  Discard local variables  $X_t, A_t$ , and  $Y_t$ 
12: end for

```

**Ensure:**  $z, A, B, C, M, N$

$\pi_j$  and  $\sigma_{k_t-1}$  are as defined in Eq. (25).

**Inference of  $b$ .** We infer  $b$  as below.

$$p(b_{z_t r} = s | z_{1:t}, A_{1:t}, B_{1:k_t}, X_{1:t}) \propto p(X_{tr} | z_{1:t}, A_t, B_{1:k_t}, X_{1:t}^{-tr}) p(b_{z_t r} = s | B_{z_t}^{-r}, z_{1:t}, A_{1:t}, B_{1:k_t}) \quad (48)$$

where  $p(X_{tr} | z_{1:t}, A_t, B_{1:k_t}, X_{1:t}^{-tr})$  is  $L_s(X_{tr})$  and  $p(b_{z_t r} = s | B_{z_t}^{-r}, z_{1:t}, A_{1:t}, B_{1:k_t})$  comes from CRP as

$$\propto L_s(X_{tr}) (n_{z_t s}^{-r} + N_s^{-z_t}) (1 - \iota_s) + \lambda L_{s_{new}}(X_{tr}) \iota_s \quad (49)$$

we define the variables as

$$\begin{aligned} \iota_s &= \mathbb{I}[s = s_{new}], \quad n_{z_t s}^{-r} = \sum_{q \neq r} \mathbb{I}[b_{z_t q} = s], \\ N_s^{-z_t} &= \sum_{l=1}^{k_t-1} \sum_q \mathbb{I}[b_{lq} = s, l \neq z_t] \end{aligned} \quad (50)$$

$n_{z_t s}^{-r}$  denotes the number of times component  $\beta_s$  has been used in the mixing distribution  $Q_{z_t}$  excluding  $\psi_{jr}$ . Whereas  $N_s^{-z_t}$  denotes how many times component  $\beta_s$  is used in unique mixing distributions ( $Q_j$ ) except  $Q_{z_t}$ .

**SUMO for DPMM on text datasets.** Using Eq. (43) in step 5, and Eq. (46), Eq. (48) in step 7 of SUMO (Algorithm 1), we obtain SUMO for text datasets presented in Algorithm 2.

Notice that from Eq. (44), Eq. (47) and Eq. (49) that by maintaining sufficient statistics  $M, J, N$  and  $L$ , we need not store the local variables  $A_{1:t-1}, Y_{1:t-1}, X_{1:t-1}$ .

## References

- Connor, R. and Mosimann, J. Concepts of independence for proportions with a generalization of the dirichlet distribution. *Journal of the American Statistical Association*, 64(325):194–206, 1969.
- Ishwaran, H. and James, L. Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 96(453):161–173, 2001.

- Lee, J., Quintana, F., Muller, P., and Trippa, L. Defining predictive probability functions for species sampling models. *Statistical Science*, 28(2):209–222, 2013.
- Pitman, J. Exchangeable and partially exchangeable random partitions. *Probability Theory and Related Fields*, 102(2):145–158, 1995.
- Pitman, J. Some developments of the Blackwell-MacQueen urn scheme. *Lecture Notes-Monograph Series*, pp. 245–267, 1996.
- Pitman, J. and Yor, M. The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *The Annals of Probability*, 25(2):855–900, 1997.
- Sethuraman, J. A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–650, 1994.
- Teh, Y., Jordan, M., Beal, M., and Blei, D. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.