

# DISTRIBUTION THEORY FOR HIERARCHICAL PROCESSES

BY FEDERICO CAMERLENGHI<sup>\*</sup>, ANTONIO LIJOI<sup>\*,†</sup> PETER ORBANZ AND IGOR PRÜNSTER<sup>†</sup>

*University of Milano–Bicocca, Bocconi University, Columbia University and Bocconi University*

Hierarchies of discrete probability measures are remarkably popular as nonparametric priors in applications, arguably due to two key properties: (i) they naturally represent multiple heterogeneous populations; (ii) they produce ties across populations, resulting in a shrinkage property often described as “sharing of information”. In this paper we establish a distribution theory for hierarchical random measures that are generated via normalization, thus encompassing both the hierarchical Dirichlet and hierarchical Pitman–Yor processes. These results provide a probabilistic characterization of the induced (partially exchangeable) partition structure, including the distribution and the asymptotics of the number of partition sets, and a complete posterior characterization. They are obtained by representing hierarchical processes in terms of completely random measures, and by applying a novel technique for deriving the associated distributions. Moreover, they also serve as building blocks for new simulation algorithms, and we derive marginal and conditional algorithms for Bayesian inference.

**1. Introduction.** The random partition structure induced by discrete nonparametric priors plays a pivotal role in a number of inferential problems related to clustering, density estimation, and prediction. It appears in applications such as species sampling, computational linguistics and topic modeling, genomics, and networks. The theory for the exchangeable case is now well understood and extensively studied. See, e.g., [26, 38, 39, 18] for probabilistic investigations and, e.g., [20, 21, 23, 10] for statistical contributions. However, in most applications data are intrinsically heterogeneous and consistent with a dependence assumption more general than exchangeability.

---

<sup>\*</sup>Also affiliated to Collegio Carlo Alberto, Moncalieri, Italy.

<sup>†</sup>Supported by the European Research Council (ERC) through StG “N-BNP” 306406  
*AMS 2000 subject classifications:* Primary 60G57, 62G05, 62F15

*Keywords and phrases:* Bayesian Nonparametrics, Distribution theory, Hierarchical processes, Partition structure, Posterior distribution, Prediction, Random measures, Species sampling models

Starting from the seminal contributions of MacEachern [34, 35], an extensive literature has been developed to address inferential issues arising with non-exchangeable observations in a Bayesian nonparametric setting. See [11, 43] for reviews. In document analysis, for example, the overall population consists of all words in a collection of documents, but each document constitutes a sub-population with its own distribution. Latent Dirichlet Allocation (LDA) was developed in [2] as a simple and effective solution; its enormous popularity is testament to the importance of the problem. The hierarchical Dirichlet process [42] is a natural nonparametric extension. Further contributions in this direction include [17, 43, 45, 37]. In these models, the induced partition structure determines the inferential outcomes but, due to the analytical complexity, its investigation and that of the associated prediction rules have been quite limited; first contributions in this direction, under different dependence assumptions, can be found in [30, 36, 46]. As far as posterior characterizations are concerned, no results are known beyond the hierarchical Dirichlet case [43]. Such characterizations are of theoretical interest, but also a prerequisite for inference algorithms, which simulate draws from (unobserved) random measures conditionally on data. See [6, 19, 31, 46] for examples, and [16] for a comprehensive list of references.

The present paper deals with a general class of hierarchical processes obtained by normalizing random measures, which encompass hierarchical Dirichlet and Pitman-Yor processes. We establish a distribution theory for this class of processes and determine the two distributional quantities essential for Bayesian inference, namely the induced partition structure and a posterior characterization. These allow to perform prediction density estimation, clustering and the assessment of distributional homogeneity across different samples. The focus on a general class of priors rather than on special cases has a two-fold motivation. On the one hand, it helps to clarify the underlying, probabilistic structure of hierarchical models and its statistical implications. On the other hand, the Dirichlet process has well-known limitations in the plain exchangeable framework, and that is similarly true in the non-exchangeable case. In the former, various extensions of the Dirichlet process have been introduced to provide more flexibility; our results provide counterparts in the latter more general framework.

1.1. *Partial exchangeability.* A random infinite sequence is exchangeable if its distribution is invariant under the group of all finitary permutations (those which permute an arbitrary but finite number of indices of the sequence). It is *partially exchangeable* if invariance holds under a subgroup of such permutations; see [24] for an extensive bibliography. In the prob-

lems considered in the following, partial exchangeability arises naturally: if a population decomposes into (conditionally independent) multiple sub-populations that are each exchangeable in their own right, the overall population is partially exchangeable.

More formally, suppose  $\mathbb{X}$  is a complete and separable metric space endowed with the Borel  $\sigma$ -field  $\mathcal{X}$ . Consider  $d$  partially exchangeable sequences  $\{(X_{i,j})_{j \geq 1} : i = 1, \dots, d\}$  defined on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and taking values in  $(\mathbb{X}, \mathcal{X})$ . By de Finetti's representation theorem this is equivalent to assuming

$$(1) \quad \mathbb{P}\left[\left\{\mathbf{X}^{(N_i)} \in A_i : i = 1, \dots, d\right\}\right] = \int_{\mathbb{P}_{\mathbb{X}}^d} \prod_{i=1}^d p_i^{(N_i)}(A_i) Q_d(dp_1, \dots, dp_d)$$

for any integer  $N_i \geq 1$  and  $A_i \in \mathcal{X}^{N_i}$ , where  $\mathbf{X}^{(N_i)} = (X_{i,1}, \dots, X_{i,N_i})$  and  $p^{(q)} = p \times \dots \times p$  is the  $q$ -fold product measure on  $\mathbb{X}^q$ , for any  $q \geq 1$ . Moreover,  $\mathbb{P}_{\mathbb{X}}$  is the space of all probability measures on  $\mathbb{X}$ , which we suppose is endowed with the topology of weak convergence and denote as  $\mathcal{P}_{\mathbb{X}}$  the corresponding Borel  $\sigma$ -algebra. The mixing or de Finetti measure  $Q_d$  is a probability measure on  $(\mathbb{P}_{\mathbb{X}}^d, \mathcal{P}_{\mathbb{X}}^d)$  that plays the role of a prior distribution. Hence, (1) amounts to assuming that, given a vector of random probability measures  $(\tilde{p}_1, \dots, \tilde{p}_d) \sim Q_d$ , the  $d$  samples are independent and the observations  $\mathbf{X}^{(N_i)}$  of the  $i$ -th sample are independent and identically distributed from  $\tilde{p}_i$ .

As in most of the current literature, here we focus on choices of  $Q_d$  that select, with probability 1, vectors of discrete probability measures. This implies that there will be ties, with positive probability, within each sample and typically also across different samples. From a modeling perspective this is a desirable feature since it allows clustering both within and across samples or, in other terms, to have models accounting for heterogeneity in a flexible way. The appearance of ties then naturally leads to look at the induced partition structure. In the exchangeable framework, the partition structure is uniquely characterized by the exchangeable partition probability function (EPPF) (see [39]), which is a key tool for studying clustering properties, deriving prediction rules and sampling schemes.

In the partially exchangeable context one can define an analogous object, which we term partially exchangeable partition probability function (pEPPF), and plays exactly the same role of the EPPF in this more general setup. In order to provide a probabilistic description of the pEPPF, let  $k$  be the number of distinct values recorded among the  $N = N_1 + \dots + N_d$  observations in  $\{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$ . Each distinct value identifies a specific cluster of the partition. Accordingly,  $\mathbf{n}_i = (n_{i,1}, \dots, n_{i,k})$  denotes the

vector of frequency counts and  $n_{i,j}$  is the number of elements of the  $i$ -th sample that coincide with the  $j$ -th distinct value. Clearly  $n_{i,j} \geq 0$  for any  $i = 1, \dots, d$  and  $j = 1, \dots, k$ , and  $\sum_{i=1}^d n_{i,j} \geq 1$  for any  $j = 1, \dots, k$ . Note that  $n_{i,j} = 0$  means that the  $j$ -th distinct value does not appear in the  $i$ -th sample. The  $j$ -th distinct value is shared by any two samples  $i$  and  $\kappa$  if and only if  $n_{i,j}n_{\kappa,j} \geq 1$ . To sum up, the pEPPF is defined as

$$(2) \quad \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) = \mathbb{E} \int_{\mathbb{X}^k} \prod_{j=1}^k \tilde{p}_1^{n_{1,j}}(dx_j) \dots \tilde{p}_d^{n_{d,j}}(dx_j)$$

with the obvious constraint  $\sum_{j=1}^k n_{i,j} = N_i$ , for each  $i = 1, \dots, d$ .

1.2. *Outline.* The main goal of the paper is to establish a distribution theory for prior distributions  $Q_d$  displaying a hierarchical structure and selecting discrete random probabilities. We focus on two key aspects. On the one hand, we investigate the random partitions induced by an array of partially exchangeable sequences as in (1), including the distribution of the number of partition sets and its asymptotics when the sample size increases. On the other hand, we provide a posterior characterization for a vector of hierarchical random probability measures  $(\tilde{p}_1, \dots, \tilde{p}_d)$ , conditional on the data. The former allows one to address two relevant issues in Bayesian non-parametric inference, namely inference on the clustering structure of the data and prediction. The latter is crucial for accurate uncertainty quantification and for devising simulation algorithms that generate trajectories of hierarchical random probability measure, from their posterior distribution.

In Section 2 we introduce some basic elements on completely random measures and provide a description of hierarchical normalized random measures. A probabilistic characterization of the induced partially exchangeable random partition is detailed in Section 3 and this forms the basis for investigating the distributional properties of the number of distinct values in  $d$  partially exchangeable samples in Section 4. The main results for establishing a posterior representation of  $(\tilde{p}_1, \dots, \tilde{p}_d)$  are, then, stated in Section 5. Finally, the computational algorithms that can be obtained from our theoretical results are described in Section 6. Proofs are deferred to the the supplementary material [3].

**2. Hierarchical normalized random measures.** In the present work we rely on random measures as the basic building blocks for the construction of discrete nonparametric priors having a hierarchical structure. Let  $\mathbb{M}_{\mathbb{X}}$  be the space of boundedly finite measures on  $(\mathbb{X}, \mathcal{X})$ , i.e.  $m(A) < \infty$

for any  $m \in \mathbf{M}_{\mathbb{X}}$  and for any bounded set  $A \in \mathcal{X}$ , equipped with the corresponding Borel  $\sigma$ -algebra  $\mathcal{M}_{\mathbb{X}}$ . See [9] for details. We consider random elements  $\tilde{\mu}$  defined on some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  and taking values in  $(\mathbf{M}_{\mathbb{X}}, \mathcal{M}_{\mathbb{X}})$ . Furthermore,  $\tilde{\mu}$  is assumed to be almost surely discrete and without fixed points of discontinuity. Hence, they can be represented as  $\tilde{\mu} = \sum_{i \geq 1} J_i \delta_{Y_i}$ . We shall henceforth focus on random probabilities obtained as suitable transformations of  $\tilde{\mu}$ . In particular, we will focus on normalization. Indeed, if  $0 < \tilde{\mu}(\mathbb{X}) < \infty$  a.s., we define

$$(3) \quad \tilde{p} = \frac{\tilde{\mu}}{\tilde{\mu}(\mathbb{X})} = \sum_{i \geq 1} \frac{J_i}{\bar{J}} \delta_{Y_i} \sim \text{NRM}(P)$$

where  $\bar{J} := \sum_{i \geq 1} J_i = \tilde{\mu}(\mathbb{X})$  and  $P = \mathbb{E}\tilde{p}$  is a probability distribution on  $(\mathbb{X}, \mathcal{X})$ . In order to obtain a hierarchical structure, one then assumes that  $(Y_i)_{i \geq 1}$  in (3) is exchangeable with  $Y_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \tilde{p}_0$ . Moreover,  $\tilde{p}_0 = \tilde{\mu}_0 / \tilde{\mu}_0(\mathbb{X})$  is obtained by normalizing a random measure  $\tilde{\mu}_0 = \sum_{i \geq 1} J_{i,0} \delta_{Y_{i,0}}$ , where  $(Y_{i,0})_{i \geq 1}$  is an i.i.d. sequence taking values in  $\mathbb{X}$  and whose common probability distribution  $P_0$  is non-atomic. Therefore, we deal with  $d$  sequences  $\{(X_{i,j})_{j \geq 1} : i = 1, \dots, d\}$  that are partially exchangeable according to (1) and the mixing measure  $Q_d$  is characterized by

$$(4) \quad \begin{aligned} \tilde{p}_i | \tilde{p}_0 &\stackrel{\text{iid}}{\sim} \text{NRM}(\tilde{p}_0) \quad i = 1, \dots, d \\ \tilde{p}_0 &\sim \text{NRM}(P_0). \end{aligned}$$

The almost sure discreteness of  $\tilde{\mu}$  is clearly inherited by the  $\tilde{p}_i$ 's and hence, as desired, we have nonparametric priors  $Q_d$  selecting discrete distributions and inducing ties within and across the samples  $\mathbf{X}^{(N_1)}, \dots, \mathbf{X}^{(N_d)}$ .

The following subsections focus on two specifications of  $(\tilde{\mu}, \tilde{\mu}_0)$ , and hence of (4), that will be thoroughly investigated in the paper.

**2.1. Hierarchical NRMIs.** A first natural choice is to set  $\tilde{\mu}$  as a *completely random measure* (CRM), i.e. a random element taking values in  $\mathbf{M}_{\mathbb{X}}$  such that for any collection of pairwise disjoint sets  $A_1, \dots, A_k$  in  $\mathcal{X}$ , and for any  $k \geq 1$ , the random variables  $\tilde{\mu}(A_1), \dots, \tilde{\mu}(A_k)$  are mutually independent. See [25]. An appealing feature of CRMs is the availability of their Laplace functional. Indeed, if it is further assumed that  $\tilde{\mu}$  does not have fixed points of discontinuity, for any measurable function  $f : \mathbb{X} \rightarrow \mathbb{R}^+$  one has

$$(5) \quad \mathbb{E} e^{-\int_{\mathbb{X}} f(x) \tilde{\mu}(dx)} = e^{-\int_{\mathbb{R}^+ \times \mathbb{X}} [1 - e^{-sf(x)}] \nu(ds, dx)}$$

where  $\nu$  is the Lévy intensity uniquely characterizing the CRM  $\tilde{\mu}$ . See [25, 27] for an exhaustive account. Though the treatment can be developed for any

CRM, for the ease of illustration henceforth we consider the case where the jumps  $J_i$ 's and the locations  $Y_i$ 's are independent and specifically that

$$(6) \quad \nu(ds, dx) = \rho(s) ds c P_0(dx)$$

for some measurable function  $\rho : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ , constant  $c > 0$  and probability measure  $P_0$  on  $(\mathbb{X}, \mathcal{X})$ . Noteworthy examples are the gamma process and the  $\sigma$ -stable process, which correspond to CRMs having  $\rho(s) = s^{-1} e^{-s}$  and  $\rho(s) = \sigma s^{-1-\sigma} / \Gamma(1 - \sigma)$ , for some  $\sigma \in (0, 1)$ . If  $\tilde{p} = \tilde{\mu} / \tilde{\mu}(\mathbb{X})$  we use the notation

$$\tilde{p} \sim \text{NRMI}(\rho, c, P_0),$$

which recalls the acronym of [41], where normalized random measures have first been introduced and studied in the exchangeable framework. The corresponding hierarchical model in (4) is thus termed *hierarchical NRMI*.

For hierarchical NRMI's one can evaluate the correlation between  $\tilde{p}_i(A)$  and  $\tilde{p}_j(A)$ , for any  $i \neq j$  and measurable subset  $A$  of  $\mathbb{X}$ , in terms of the underlying parameters  $(c, \rho, c_0, \rho_0)$ . In order to ease the statement of the result, set  $\psi(u) = \int_0^\infty [1 - e^{-us}] \rho(s) ds$  and  $\psi_0(u) = \int_0^\infty [1 - e^{-us}] \rho_0(s) ds$  as the Laplace exponents corresponding to  $\tilde{p}$  and  $\tilde{p}_0$ , respectively.

**THEOREM 1.** *Suppose that  $\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{NRMI}(\rho, c, \tilde{p}_0)$ , for  $i = 1, \dots, d$ , and  $\tilde{p}_0 \sim \text{NRMI}(\rho_0, c_0, P_0)$ . Then, for any  $A \in \mathcal{X}$  and  $i \neq j$*

$$(7) \quad \text{corr}(\tilde{p}_i(A), \tilde{p}_j(A)) = \left\{ 1 + c_0 c \frac{\int_0^\infty u e^{-c\psi(u)} \tau_2(u) du \int_0^\infty u e^{-c_0\psi_0(u)} \tau_{1,0}^2(u) du}{\int_0^\infty u e^{-c_0\psi_0(u)} \tau_{2,0}(u) du} \right\}^{-1},$$

where  $\tau_q(u) = \int_0^\infty s^q e^{-us} \rho(s) ds$  and  $\tau_{q,0}(u) = \int_0^\infty s^q e^{-us} \rho_0(s) ds$ .

It is worth stressing two important facts. The correlation coefficient between  $\tilde{p}_i(A)$  and  $\tilde{p}_j(A)$  is always positive. It does not depend on the specific set  $A$ . Moreover, by specifying  $(c, \rho, c_0, \rho_0)$  the correlation coefficient (7) becomes readily available as shown in the following examples.

**EXAMPLE 1.** If  $\rho(s) = \rho_0(s) = s^{-1} e^{-s}$ , then  $\tilde{p}_0$  is a Dirichlet process and the  $\tilde{p}_i$ 's are, conditionally on  $\tilde{p}_0$ , independent and identically distributed Dirichlet processes. Hence,  $(\tilde{p}_1, \dots, \tilde{p}_d)$  is a vector of hierarchical Dirichlet processes as in [42]. A straightforward application of Theorem 1 yields

$$\text{corr}(\tilde{p}_i(A), \tilde{p}_j(A)) = \frac{c + 1}{c + 1 + c_0}.$$

Note that the correlation is increasing in  $c$  and decreasing in  $c_0$ . As  $c_0 \uparrow \infty$  the distribution of  $\tilde{p}_0$  degenerates on  $P_0$  and the  $\tilde{p}_i$ 's are independent, which is consistent with  $\text{corr}(\tilde{p}_i(A), \tilde{p}_j(A))$  converging to 0. On the other hand, if  $c \uparrow \infty$ , then the distribution of each  $\tilde{p}_i$ , conditional on  $\tilde{p}_0$ , degenerates on  $\tilde{p}_0$  and it is, thus, not surprising that the correlation coefficient between any pair of  $\tilde{p}_i(A)$ 's converges to 1, for any  $A$  in  $\mathcal{X}$ .

EXAMPLE 2. The hierarchical stable NRMI arises by setting  $\rho(s) = \sigma s^{-1-\sigma}/\Gamma(1-\sigma)$  and  $\rho_0(s) = \sigma_0 s^{-1-\sigma_0}/\Gamma(1-\sigma_0)$ , for some  $\sigma$  and  $\sigma_0$  in  $(0, 1)$ . This implies that  $\tilde{p}_0$  is a  $\sigma_0$ -stable NRMI and, conditionally on  $\tilde{p}_0$ , the  $\tilde{p}_i$ 's are independent and identically distributed  $\sigma$ -stable NRMI's. We will say that  $(\tilde{p}_1, \dots, \tilde{p}_d)$  is a vector of hierarchical stable NRMI's. A plain application of Theorem 1 leads to

$$\text{corr}(\tilde{p}_i(A), \tilde{p}_j(A)) = \frac{1 - \sigma_0}{1 - \sigma\sigma_0},$$

which is increasing in  $\sigma$  and decreasing in  $\sigma_0$ . Due to the properties of the stable CRM, unsurprisingly the correlation coefficient does not depend on the total masses  $c_0$  and  $c$ .

2.2. *Hierarchical Pitman–Yor processes.* The second relevant construction arises when  $\tilde{\mu}$  has a distribution obtained by a suitable transformation of the distribution of a CRM. In particular, let  $\mathbb{P}_\sigma$  be the probability distribution on  $(\mathbb{M}_\mathbb{X}, \mathcal{M}_\mathbb{X})$  of a  $\sigma$ -stable CRM, with  $\sigma \in (0, 1)$ . For  $\theta > 0$  define  $\mathbb{P}_{\sigma,\theta}$  on  $(\mathbb{M}_\mathbb{X}, \mathcal{M}_\mathbb{X})$  as absolutely continuous w.r.t.  $\mathbb{P}_\sigma$  and such that its Radon–Nikodym derivative is

$$(8) \quad \frac{d\mathbb{P}_{\sigma,\theta}}{d\mathbb{P}_\sigma}(m) = \frac{\sigma \Gamma(\theta)}{\Gamma(\theta/\sigma)} m^{-\theta}(\mathbb{X})$$

The resulting random measure  $\tilde{\mu}_{\sigma,\theta}$  with distribution  $\mathbb{P}_{\sigma,\theta}$  is not completely random. Nonetheless, via normalization

$$(9) \quad \tilde{p} = \frac{\tilde{\mu}_{\sigma,\theta}}{\tilde{\mu}_{\sigma,\theta}(\mathbb{X})} \sim \text{PY}(\sigma, \theta; P)$$

one obtains a fundamental process, the *Pitman–Yor process* or *two-parameter Poisson–Dirichlet process*. A different equivalent construction, simpler but less convenient for our purposes, starts from a specific NRMI  $(\rho, c, P_0)$  and puts a gamma prior on the parameter  $c$ . See [40] for details on both derivations.

The following results provides the correlation structure for the hierarchical Pitman–Yor process and nicely describes the role of the parameters  $(\sigma, \sigma_0, \theta, \theta_0)$ .

THEOREM 2. *Suppose that  $\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{PY}(\sigma, \theta, \tilde{p}_0)$ , for  $i = 1, \dots, d$ , and  $\tilde{p}_0 \sim \text{PY}(\sigma_0, \theta_0, P_0)$ . Then, for any  $A \in \mathcal{X}$  and  $i \neq j$*

$$(10) \quad \text{corr}(\tilde{p}_i(A), \tilde{p}_j(A)) = \left\{ 1 + \frac{1 - \sigma}{1 - \sigma_0} \frac{\theta_0 + \sigma_0}{\theta + 1} \right\}^{-1}$$

Unsurprisingly, also for hierarchical Pitman–Yor processes the correlation between  $\tilde{p}_i(A)$  and  $\tilde{p}_j(A)$ , for any  $i \neq j$ , is positive and does not depend on  $A \in \mathcal{X}$ . Moreover, from (10) the impact of  $(\theta_0, \sigma_0, \theta, \sigma)$  on  $\text{corr}(\tilde{p}_i(A), \tilde{p}_j(A))$  can be easily deduced.

**3. Random partitions induced by hierarchical NRMs.** Consider an array of  $d$  partially exchangeable sequences with de Finetti measure  $Q_d$  given by hierarchies of normalized measures as in (4). As already mentioned, the discreteness of the  $\tilde{p}_i$ 's and  $\tilde{p}_0$  entails that  $\mathbb{P}[X_{\ell,i} = X_{\kappa,j}] > 0$  for any  $\ell$  and  $\kappa$ , i.e. there is a positive probability of ties both within each sample and across the different samples  $\mathbf{X}^{(N_i)} = (X_{i,1}, \dots, X_{i,N_i})$ . A random partition of the samples is, thus, induced, whereby any two elements  $X_{\ell,i}$  and  $X_{\kappa,j}$  are in the same partition group (or cluster) if and only if they take on the same value. Its probability distribution is identified by the pEPPF  $\Pi_k^{(N)}$  in (2). Here we determine a closed form expression for hierarchical NRMs and the hierarchical Pitman–Yor process.

We first focus on hierarchical NRMs. In order to gain some intuition on the structure of  $\Pi_k^{(N)}$ , it is worth recalling the so-called *Chinese restaurant franchise* metaphor described in [42] for the hierarchical Dirichlet process. According to this scheme, a franchise of  $d$  restaurants shares the same menu, which includes an infinite number of dishes and is generated by the top level base measure  $P_0$ . Each restaurant has infinitely many tables. The first customer sitting at each table of restaurant  $i$  chooses the dish and this dish is shared by all other customers who afterwards join the same table. In contrast to the well-known Chinese restaurant process, the same dish can be served at different tables within the same restaurant and across different restaurants. According to this scheme,  $X_{i,j}$  represents the dish served in the  $i$ -th restaurant to the  $j$ -th customer for  $j = 1, \dots, N_i$  and  $i = 1, \dots, d$ . Furthermore, the frequency  $n_{i,j}$  in (2) is the number of customers in restaurant  $i$  eating the  $j$ -th dish and we further let  $\ell_{i,j} \in \{1, \dots, n_{i,j}\}$  be the number of tables in restaurant  $i$  at which the  $j$ -th dish is served, if  $n_{i,j} \geq 1$ . When  $n_{i,j} = 0$  it is obvious that  $\ell_{i,j} = 0$  as well. Hence

$$\bar{\ell}_{\bullet j} = \sum_{i=1}^d \ell_{i,j}, \quad \bar{\ell}_{i\bullet} = \sum_{j=1}^k \ell_{i,j}.$$



denote, respectively, the number of tables serving dish  $j$  (across restaurants) and the number of tables in restaurant  $i$  (regardless of the served dishes). Moreover, if we further label the tables, with  $q_{i,j,t}$  we can identify the number of customers in restaurant  $i$  eating dish  $j$  at table  $t$  so that  $\sum_{t=1}^{\ell_{i,j}} q_{i,j,t} = n_{i,j}$ . This additional notation suggests we are going to consider a combinatorial structure arising from the composition of random partitions acting at different levels of the hierarchy: one yields a partition where the  $N = N_1 + \dots + N_d$  customers are allocated to  $|\ell| = \sum_{i=1}^d \sum_{j=1}^k \ell_{i,j} \geq k$  tables and these tables are, then, clustered into  $k$  groups, with each group being identified by a different distinct dish.

Before providing the pEPPF, we introduce the notation that identifies the composing random partitions. If  $\tilde{p}_0 \sim \text{NRMI}(\rho_0, c_0, P_0)$  and  $P_0$  is a diffuse probability measure on  $\mathbb{X}$ , for any  $k \in \{1, \dots, n\}$  and any vector of positive integers  $(r_1, \dots, r_k)$  such that  $\sum_{i=1}^k r_i = n$ , we set

$$(11) \quad \Phi_{k,0}^{(n)}(r_1, \dots, r_k) = \frac{c_0^k}{\Gamma(n)} \int_0^\infty u^{n-1} e^{-c_0 \psi_0(u)} \prod_{j=1}^k \tau_{r_j,0}(u) \, du.$$

Note that according to [23, Proposition 3],  $\Phi_{k,0}^{(n)}$  is the EPPF induced by an exchangeable sequence drawn from a NRMI with parameter  $(c_0, \rho_0)$ .

**THEOREM 3.** *Suppose the sequences  $\{(X_{i,j})_{j \geq 1} : i = 1, \dots, d\}$  are partially exchangeable according to (1), with  $Q_d$  characterized by*

$$\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{NRMI}(\rho, c, \tilde{p}_0) \quad (i = 1, \dots, d), \quad \tilde{p}_0 \sim \text{NRMI}(\rho_0, c_0, P_0).$$

Then

$$(12) \quad \begin{aligned} \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) &= \sum_{\ell} \sum_{\mathbf{q}} \Phi_{k,0}^{(|\ell|)}(\bar{\ell}_{\bullet,1}, \dots, \bar{\ell}_{\bullet,k}) \\ &\times \prod_{i=1}^d \prod_{j=1}^k \frac{1}{\ell_{i,j}!} \binom{n_{i,j}}{q_{i,j,1}, \dots, q_{i,j,\ell_{i,j}}} \Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}(\mathbf{q}_{i,1}, \dots, \mathbf{q}_{i,k}) \end{aligned}$$

where, if  $n_{i,j} \geq 1$ ,  $\mathbf{q}_{i,j} = (q_{i,j,1}, \dots, q_{i,j,\ell_{i,j}})$  is a vector of positive integers such that  $|\mathbf{q}_{i,j}| = n_{i,j}$ , for any  $i = 1, \dots, d$  and  $j = 1, \dots, k$ , and

$$(13) \quad \Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}(\mathbf{q}_{i,1}, \dots, \mathbf{q}_{i,k}) = \frac{c^{\bar{\ell}_{i,\bullet}}}{\Gamma(N_i)} \int_0^\infty u^{N_i-1} e^{-c\psi(u)} \prod_{j=1}^k \prod_{t=1}^{\ell_{i,j}} \tau_{q_{i,j,t}}(u) \, du.$$

Note that, if  $n_{i,j} = 0$ , then  $\mathbf{q}_{i,j} = (0, \dots, 0)$  and

$$\Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}(\mathbf{q}_{i,1}, \dots, \mathbf{q}_{i,k}) = \Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}(\mathbf{q}_{i,1}, \dots, \mathbf{q}_{i,j-1}, \mathbf{q}_{i,j+1}, \dots, \mathbf{q}_{i,k})$$

The backbone of (12) is

$$(14) \quad \Phi_{k,0}^{(|\ell|)}(\bar{\ell}_{\bullet 1}, \dots, \bar{\ell}_{\bullet k}) \prod_{i=1}^d \Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}(\mathbf{q}_{i,1}, \dots, \mathbf{q}_{i,k})$$

which displays the random partitions' composition acting at the two levels of the hierarchy: the single samples (or restaurants) and the whole collection of samples (or the franchise). The former is captured by  $\prod_{i=1}^d \Phi_{\bar{\ell}_{i,\bullet}, i}^{(N_i)}$  while the latter is identified by  $\Phi_{k,0}^{(|\ell|)}$ . The resulting expression of  $\Pi_k^{(N)}$  then follows from plain marginalization.

We now illustrate the result by considering again the hierarchical Dirichlet process and the hierarchical stable NRM.

EXAMPLE 3. Let  $(\tilde{p}_1, \dots, \tilde{p}_d)$  be a vector of hierarchical Dirichlet processes as in Example 1. Let  $(a)_n = \Gamma(a+n)/\Gamma(a)$  be the ascending factorial and  $|\mathfrak{s}(n, k)|$  the signless Stirling number of the first kind. It is then straightforward to show that

$$\begin{aligned} \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) &= \frac{c_0^k}{\prod_{i=1}^d (c)_{N_i}} \sum_{\ell} \frac{c^{|\ell|}}{(c_0)^{|\ell|}} \prod_{j=1}^k (\bar{\ell}_{\bullet j} - 1)! \prod_{i=1}^d |\mathfrak{s}(n_{i,j}, \ell_{i,j})| \\ &= c_0^k \left( \prod_{i=1}^d \frac{\prod_{j=1}^k (c)_{n_{i,j}}}{(c)_{N_i}} \right) \sum_{\ell} \frac{1}{(c_0)^{|\ell|}} \prod_{j=1}^k (\bar{\ell}_{\bullet j} - 1)! \prod_{i=1}^d \mathbb{P}[K_{n_{i,j}} = \ell_{i,j}] \end{aligned}$$

where  $K_{n_{i,j}}$  is a random variable denoting the number of distinct observations, out of  $n_{i,j}$  drawn from an exchangeable sequence whose de Finetti measure is a Dirichlet process with concentration parameter  $c$ . Alternatively, one can rely on properties of  $|\mathfrak{s}(n, k)|$  and deduce the following integral representation

$$\begin{aligned} \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) &= \\ &= \frac{c_0^k c^{\xi}}{\prod_{i=1}^d (c)_{N_i}} \int_{\Delta_k} D_k(d\mathbf{p}; \xi_1, \dots, \xi_k, c_0) \prod_{i=1}^d \prod_{\{j: n_{i,j} \geq 1\}} (c p_j + 1)^{n_{i,j}-1} \end{aligned}$$

where  $\xi_j = \sum_{i=1}^d \mathbb{1}_{\{1,2,\dots\}}(n_{i,j})$  is the number of restaurants sharing the  $j$ -th dish,  $\xi = \sum_{j=1}^k \xi_j$  and  $D_k(\cdot; a_1, \dots, a_{k+1})$  is the multivariate Dirichlet distribution on the  $k$ -dimensional simplex  $\Delta_k$ , with parameters  $(a_1, \dots, a_{k+1})$ .

EXAMPLE 4. Let  $(\tilde{p}_1, \dots, \tilde{p}_d)$  be a vector of hierarchical stable NRMIs defined as in Example 2 and  $\mathcal{C}(n, k; \sigma)$  be the generalized factorial coefficients defined by

$$(\sigma t)_n = \sum_{k=1}^n \mathcal{C}(n, k; \sigma) (t)_k.$$

As for the pEPPF, Theorem 3 and some algebra lead to

$$\begin{aligned} \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) &= \frac{\sigma_0^{k-1} \Gamma(k)}{\prod_{i=1}^d \Gamma(N_i)} \sum_{\ell} \frac{\sigma^{|\ell|-d} \prod_{i=1}^d \Gamma(\bar{\ell}_{i\bullet})}{\Gamma(|\ell|)} \prod_{j=1}^k (1 - \sigma_0)_{\bar{\ell}_{\bullet j} - 1} \\ &\quad \times \prod_{i=1}^d \prod_{j=1}^k \frac{\mathcal{C}(n_{i,j}, \ell_{i,j}; \sigma)}{\sigma^{\ell_{i,j}}} \\ &= \sigma_0^{k-1} \sigma^{\xi-d} \prod_{i=1}^d \left( \frac{\prod_{\{j: n_{i,j} \geq 1\}} \Gamma(n_{i,j})}{\Gamma(N_i)} \right) \sum_{\ell} \frac{\prod_{j=1}^k (1 - \sigma_0)_{\bar{\ell}_{\bullet j} - 1}}{\Gamma(|\ell|)} \\ &\quad \times \prod_{i=1}^d \left( \frac{\Gamma(\bar{\ell}_{i\bullet})}{\prod_{\{j: n_{i,j} \geq 1\}} \Gamma(\ell_{i,j})} \right) \prod_{i=1}^d \prod_{j=1}^k \mathbb{P}[K_{n_{i,j}} = \ell_{i,j}] \end{aligned}$$

with  $K_{n_{i,j}}$  denoting the number of distinct observations generated by  $n_{i,j}$  observations from an exchangeable sequence whose de Finetti measure is a normalized  $\sigma$ -stable process.

The combinatorial structure yielding the pEPPF in (12) is not specific to hierarchical NRMIs. Indeed, it can be established also for the Pitman–Yor process, which arises as the normalization of a measure that is not completely random.

THEOREM 4. Let  $\{(X_{i,j})_{j \geq 1} : i = 1, \dots, d\}$  be partially exchangeable as in (1), with  $Q_d$  characterized by

$$\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{PY}(\sigma, \theta; \tilde{p}_0) \quad (i = 1, \dots, d), \quad \tilde{p}_0 \sim \text{PY}(\sigma_0, \theta_0; P_0)$$

Then

$$\begin{aligned} (15) \quad \Pi_k^{(N)}(\mathbf{n}_1, \dots, \mathbf{n}_d) &= \sum_{\ell} \frac{\prod_{r=1}^{k-1} (\theta_0 + r\sigma_0)}{(\theta_0 + 1)^{|\ell|-1}} \prod_{j=1}^k (1 - \sigma_0)_{\bar{\ell}_{\bullet j} - 1} \\ &\quad \times \prod_{i=1}^d \frac{\prod_{r=1}^{\bar{\ell}_{i\bullet} - 1} (\theta + r\sigma)}{(\theta + 1)^{N_i - 1}} \prod_{j=1}^k \frac{\mathcal{C}(n_{i,j}, \ell_{i,j}; \sigma)}{\sigma^{\ell_{i,j}}} \end{aligned}$$

This result is related to the findings in [17], whose construction leads to a tree structure used as a language model. In models of this type, termed *sequence memoizer*, the observations take values in the space  $\Sigma^*$  of finite sequences of elements from a countable (typically finite) set symbols  $\Sigma$ . Each random probability measure involved in the hierarchies of the model is supported by  $\Sigma$  and it is, then, apparent that the base measure at the root of the hierarchy is atomic. Our treatment is different, in the sense that the state space coincides with any separable and complete metric space  $\mathbb{X}$  and the probability distribution at the root of the hierarchy is diffuse. The latter is crucial for obtaining the expressions of the pEPPF displayed in this paper.

**4. Distribution of the number of clusters  $K_N$ .** Having determined the pEPPF of hierarchical NRMI and hierarchical Pitman–Yor processes, a natural issue to address is the determination of the probability distribution of the number  $K_N$  of distinct values out of  $N = N_1 + \dots + N_d$  partially exchangeable observations. This can be achieved by relying on the composition of random partitions in the pEPPF representations in Theorems 3 and 4 and highlighted in (14). For the derivation of the result, it is useful to introduce a collection of sequences of latent random variables  $\{(T_{i,j})_{j \geq 1} : i = 1, \dots, d\}$ . They are such that  $T_{i,j} | \tilde{q}_i \stackrel{\text{iid}}{\sim} \tilde{q}_i$ , with  $\tilde{q}_i \stackrel{\text{iid}}{\sim} \text{NRMI}(c, \rho, G)$  for hierarchical NRMI and  $\tilde{q}_i \stackrel{\text{iid}}{\sim} \text{PY}(\sigma, \theta, G)$  for hierarchical Pitman–Yor processes, while  $G$  is some diffuse probability measure. In terms of the Chinese restaurant franchise metaphor,  $T_{i,j}$  is the label of the table where the  $j$ -th customer of the  $i$ -th restaurant is seated. In view of this, the probability distribution of  $K_N$  arises by considering:

- (i) independent random variables  $K'_{i,N_i}$  that equal, for each  $i = 1, \dots, d$ , the number of distinct values in  $\mathbf{T}^{(N_i)} = (T_{i,1}, \dots, T_{i,N_i})$ ;
- (ii)  $K_{0,t}$ , which represents the number of distinct values out of  $t$  exchangeable random elements generated from  $\tilde{p}_0$ .

According to the Chinese restaurant metaphor,  $K'_{i,N_i}$  is the number of tables where the  $N_i$  customers of restaurant  $i$  are seated, while  $K_{0,t}$  is the number of distinct dishes allocated to the  $t$  tables where the  $N$  customers of the whole franchise are seated.

**THEOREM 5.** *Suppose  $K_N$  is the number of distinct values in the  $d$  partially exchangeable samples  $\{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$  governed by a vector of hierarchical NRMI, i.e.  $\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{NRMI}(c, \rho, \tilde{p}_0)$  and  $\tilde{p}_0 \sim \text{NRMI}(c_0, \rho_0, P_0)$ ,*

with  $P_0$  being non-atomic. Then, for any  $k = 1 \dots, N$  one has

$$(16) \quad \mathbb{P}[K_N = k] = \sum_{t=k}^N \mathbb{P}[K_{0,t} = k] \mathbb{P}\left[\sum_{i=1}^d K'_{i,N_i} = t\right]$$

The probability distributions of  $K_{0,t}$  and of  $K'_{i,N_i}$  are readily derived from their EPPFs and coincide with

$$(17) \quad \mathbb{P}[K_{0,t} = k] = \frac{1}{k!} \sum_{(r_1, \dots, r_k) \in \Delta_{k,t}} \binom{t}{r_1 \dots r_k} \Phi_{k,0}^{(t)}(r_1, \dots, r_k)$$

for any  $k \in \{1, \dots, t\}$ , where  $\Delta_{j,n} = \{(r_1, \dots, r_j) : r_i \geq 1, \sum_{i=1}^j r_i = n\}$ , and

$$(18) \quad \mathbb{P}[K'_{i,N_i} = \zeta] = \frac{1}{\zeta!} \sum_{(r_1, \dots, r_\zeta) \in \Delta_{\zeta, N_i}} \binom{N_i}{r_1 \dots r_\zeta} \Phi_{\zeta,i}^{(N_i)}(r_1, \dots, r_\zeta)$$

for any  $\zeta \in \{1, \dots, N_i\}$ .

A similar result holds for the hierarchical Pitman–Yor process.

**THEOREM 6.** *Suppose  $K_N$  is the number of distinct values in the  $d$  partially exchangeable samples  $\{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$  governed by a vector of hierarchical Pitman–Yor processes, i.e.  $\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{PY}(\sigma, \theta; \tilde{p}_0)$  and  $\tilde{p}_0 \sim \text{PY}(\sigma_0, \theta_0; P_0)$ . Then*

$$(19) \quad \mathbb{P}[K_N = k] = \sum_{t=k}^N \frac{\prod_{r=1}^{k-1} (\theta_0 + r\sigma_0)}{(\theta_0 + 1)_{t-1}} \frac{\mathcal{C}(t, k; \sigma_0)}{\sigma_0^k} \\ \times \sum_{\{(\zeta_1, \dots, \zeta_d) \in \Delta_{d,t}\}} \prod_{i=1}^d \frac{\prod_{r=1}^{\zeta_i-1} (\theta + r\sigma)}{(\theta + 1)_{N_i-1}} \frac{\mathcal{C}(N_i, \zeta_i; \sigma)}{\sigma^{\zeta_i}}$$

**REMARK 1.** In the proofs of Theorems 5–6, based on the expressions of the pEPPFs, we give an alternative equivalent representation of  $K_N$ : if  $\xi(\mathbf{N}) = K'_{1,N_1} + \dots + K'_{d,N_d}$ , from (16) and (19) one deduces for both hierarchical NRMI and Pitman–Yor processes

$$K_N \stackrel{d}{=} K_{0,\xi(\mathbf{N})}.$$

The equality between  $K_N$  and  $K_{0,\xi(\mathbf{N})}$  can be strengthened, and actually holds almost surely. This fact is useful for the determination of the asymptotic behaviour of  $K_N$ .

Before establishing the asymptotic behavior of  $K_N$ , as  $N \rightarrow \infty$ , introduce two positive sequences  $(\lambda_0(n))_{n \geq 1}$  and  $(\lambda(n))_{n \geq 1}$  such that  $\lim_n \lambda_0(n) = \lim_n \lambda(n) = \infty$  and assume  $\lambda_0$  satisfies the following condition:

(H1) for any pair of positive sequences  $(b_1(n))_{n \geq 1}$  and  $(b_2(n))_{n \geq 1}$  such that  $\lim_n b_1(n) = \lim_n b_2(n) = \infty$  and  $\lim_n (b_1(n)/b_2(n)) = 1$

$$\lim_{n \rightarrow \infty} \frac{\lambda_0(b_1(n))}{\lambda_0(b_2(n))} = 1$$

We would like to stress that assumption (H1) is satisfied when  $\lambda_0$  is a *regularly varying function*.

In the sequel we agree that  $Y_n \simeq \lambda(n)$ , for  $n \rightarrow \infty$ , means that  $\lim_n Y_n/\lambda(n)$  almost surely exists and equals a finite and positive random variable, then one can state the following.

**THEOREM 7.** *Suppose  $K_N$  is the number of distinct values in the  $d$  partially exchangeable samples  $\{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$  governed by a vector of hierarchical NRMI's such that  $K_{0,N} \simeq \lambda_0(n)$  and  $K'_{i,N} \simeq \lambda(N)$  as  $N \rightarrow \infty$ , where  $(\lambda_0(n))_{n \geq 1}$  satisfies (H1). Moreover, let  $N_1 = \dots = N_d = N^* = N/d$ . Then*

$$K_N \simeq \lambda_0(\eta \lambda(N/d)) \quad \text{as } N \rightarrow \infty,$$

for some positive and finite random variable  $\eta$ .

In particular, if  $(\tilde{p}_1, \dots, \tilde{p}_d)$  is a vector of hierarchical Dirichlet processes, then

$$K_N \simeq \log \log N \quad \text{as } N \rightarrow \infty.$$

Note that the rate of increase of  $K_N$  for the hierarchical Dirichlet process has been also displayed in [43] based on a more informal argument. The corresponding result for hierarchical Pitman–Yor process is as follows.

**THEOREM 8.** *Suppose  $K_N$  is the number of distinct values in the  $d$  partially exchangeable samples  $\{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$  governed by a vector of hierarchical Pitman–Yor processes. Furthermore, let  $N_1 = \dots = N_d = N^* = N/d$ . Then*

$$K_N \simeq N^{\sigma \sigma_0} \quad \text{as } N \rightarrow \infty$$

**REMARK 2.** These results can be extended to the case where only a subset of the  $N_i$ 's diverge and the others stay finite. Indeed, if for some  $m \leq d$  one has  $N_{j_1} = \dots = N_{j_m} = N^*$ , where  $N^* \rightarrow \infty$ , and  $N_i < L < \infty$  for any other  $i \notin \{j_1, \dots, j_m\}$ , then it is possible to conclude that

$$K_N \simeq \lambda_0(\eta \lambda(N/m))$$

as  $N^* \rightarrow \infty$ , which entails  $N \rightarrow \infty$ . This leaves the rates of increase for  $K_N$  displayed in Theorems 7–8 unchanged .

REMARK 3. With some care the results can be generalized to cover the case of the  $N_i$ 's diverging at different rates. Indeed, considering the asymptotics as  $\max_{1 \leq i \leq d} N_i \rightarrow \infty$ ,  $K_N$  increases at rates similar to those displayed Theorems 7–8.

**5. Posterior characterizations.** In order to complete the description of distributional properties of hierarchical processes, it is essential to determine a posterior characterization. To the best of our knowledge, no posterior characterization is available for dependent processes in a partially exchangeable framework, whether constructed in terms of hierarchies or by different means. Hence, our following results are the very first. Despite the theoretical interest, note that, while for prediction the partition probability functions of Theorems 3–4 suffice, inference on non-linear functionals of  $(\tilde{p}_1, \dots, \tilde{p}_d)$  requires the posterior distribution of the vector of hierarchical random probabilities.

5.1. *Hierarchical NRMI posterior.* In the following let  $X_1^*, \dots, X_k^*$  denote the distinct observations featured by the whole collection of samples  $\mathbf{X} = \{\mathbf{X}^{(N_i)} : i = 1, \dots, d\}$  and assume  $U_0$  is a positive random variable whose density function, conditional on  $\mathbf{X}$  and on the latent tables' labels  $\mathbf{T} = \{\mathbf{T}^{(N_i)} : i = 1, \dots, d\}$  introduced in Section 4, equals

$$(20) \quad f_0(u|\mathbf{X}, \mathbf{T}) \propto u^{|\ell|-1} e^{-c_0 \psi_0(u)} \prod_{j=1}^k \tau_{\tilde{\ell}_{\bullet, j}, 0}(u).$$

The posterior characterization is then composed of two blocks, the first concerning the root of the hierarchy in terms of  $\tilde{\mu}_0$  and the second concerning the vector of random probabilities.

THEOREM 9. *Suppose the data  $\mathbf{X}$  are partially exchangeable and are modeled as in (4). Then*

$$(21) \quad \tilde{\mu}_0 | (\mathbf{X}, \mathbf{T}, U_0) \stackrel{d}{=} \eta_0^* + \sum_{j=1}^k I_j \delta_{X_j^*}$$

where the two summands on the right-hand side of the distributional identity are independent and

(i)  $\eta_0^*$  is a CRM with intensity

$$\nu_0(ds, dx) = e^{-U_0 s} \rho_0(s) ds c_0 P_0(dx).$$

(ii) the  $I_j$ 's are independent and non-negative jumps with density

$$f_j(s|\mathbf{X}, \mathbf{T}) \propto s^{\bar{\ell} \cdot j} e^{-s U_0} \rho_0(s)$$

It is worth noting that the posterior of  $\tilde{\mu}_0$  depends on sample information across the populations rather than population-specific, most notably the number of different dishes served across restaurants. This clearly serves the purpose of directing the dependence across populations. Theorem 9 allows us then to establish the posterior distribution of a vector  $(\tilde{\mu}_1, \dots, \tilde{\mu}_d)$  of hierarchical CRMs, conditional a vector  $\mathbf{U} = (U_1, \dots, U_d)$  whose components are conditionally independent, given  $(\mathbf{X}, \mathbf{T})$ , and with respective densities

$$(22) \quad f_i(u|\mathbf{X}, \mathbf{T}) \propto u^{N_i-1} e^{-c\psi(u)} \prod_{j=1}^k \prod_{t=1}^{\ell_{i,j}} \tau_{q_{i,j,t}}(u) \quad i = 1, \dots, d.$$

The fundamental posterior characterization, where population-specific characteristics come into play, can then be stated as follows.

**THEOREM 10.** *Suppose the data  $\mathbf{X}$  are partially exchangeable and are modeled as in (4). Then*

$$(23) \quad (\tilde{\mu}_1, \dots, \tilde{\mu}_d) | (\mathbf{X}, \mathbf{T}, \mathbf{U}, \tilde{\mu}_0) \stackrel{d}{=} (\tilde{\mu}_1^*, \dots, \tilde{\mu}_d^*) + \left( \sum_{j=1}^k \sum_{t=1}^{\ell_{1,j}} J_{1,j,t} \delta_{X_j^*}, \dots, \sum_{j=1}^k \sum_{t=1}^{\ell_{d,j}} J_{d,j,t} \delta_{X_j^*} \right),$$

where the two summands on the right-hand-side are independent,  $\sum_{t=1}^{\ell_{i,j}} J_{i,j,t} \equiv 0$  if  $n_{i,j} = 0$  and

(i)  $(\tilde{\mu}_1^*, \dots, \tilde{\mu}_d^*)$  is a vector of hierarchical CRMs and, conditional on  $\tilde{\mu}_0^* = \eta_0^* + \sum_{j=1}^k I_j \delta_{X_j^*}$  in (21), each  $\tilde{\mu}_i^*$  has intensity

$$\nu_i(ds, dx) = e^{-U_i s} \rho(s) ds c \tilde{p}_0^*(dx),$$

with  $\tilde{p}_0^* = \tilde{\mu}_0^* / \tilde{\mu}_0^*(\mathbb{X})$ ;

(ii) the jumps  $J_{i,j,t}$  are independent and non-negative random variables whose density equals

$$f_{i,j,t}(s) \propto e^{-U_i s} s^{q_{i,j,t}} \rho(s).$$

when  $n_{i,j} \geq 1$ , whereas  $J_{i,j,t} = 0$ , almost surely, if  $n_{i,j} = 0$ .



The expressions involved in the posterior characterization of Theorem 10 are somehow reminiscent of the ones provided in [23] for the exchangeable case. This is due to the fact that, once accounted for the dependence structure inherited from the hierarchical construction, one has exchangeability within each population.

We now illustrate the general results by means of two examples, related to the hierarchical Dirichlet process and the hierarchical stable NRMI.

EXAMPLE 5. Assume that  $\rho(s) = \rho_0(s) = e^{-s}/s$ , so we are considering a vector of hierarchical Dirichlet processes. Recall that  $\psi(u) = \psi_0(u) = \log(1+u)$  and  $\tau_q(u) = \tau_{q,0}(u) = \Gamma(q)/(1+u)^q$ . In this case

$$f_0(u) = \frac{\Gamma(|\ell| + c_0)}{\Gamma(|\ell|)\Gamma(c_0)} \frac{u^{|\ell|-1}}{(1+u)^{c_0+|\ell|}} \mathbb{1}_{(0,\infty)}(u)$$

implying that  $U_0/(1+U_0) \sim \text{Beta}(|\ell|, c_0)$ . In the posterior representation of  $\tilde{\mu}_0$  as stated in Theorem 9, one has

- (a)  $\eta_0^*$  is a gamma CRM with intensity  $e^{-(1+U_0)s} s^{-1} ds c_0 P_0(dx)$ ,
- (b)  $I_j \stackrel{\text{ind}}{\sim} \text{Ga}(\bar{\ell}_{\bullet j}, 1+U_0)$ , meaning that its density function is

$$\frac{(1+U_0)^{\bar{\ell}_{\bullet j}}}{\Gamma(\bar{\ell}_{\bullet j})} x^{\bar{\ell}_{\bullet j}-1} e^{-(1+U_0)x} \mathbb{1}_{(0,\infty)}(x)$$

Now, since the normalized distributions of (a) and (b) do not depend on the scale  $U_0$ , it follows that

$$\tilde{p}_0^* = \tilde{p}_0 | (\mathbf{X}, \mathbf{T}) \sim \mathcal{D}(c_0 P_0 + \sum_{j=1}^k \bar{\ell}_{\bullet j} \delta_{X_j^*}).$$

with  $\mathcal{D}$  indicating a Dirichlet process. As far as the vector of random probabilities  $(\tilde{p}_1, \dots, \tilde{p}_d)$  is concerned, by Theorem 10 one has that, conditional on  $\tilde{p}_0^*$  and on  $(\mathbf{X}, \mathbf{T}, \mathbf{U})$ , the CRMs  $\tilde{\mu}_1, \dots, \tilde{\mu}_d$  are independent, and the distribution of each  $\tilde{\mu}_i$  equals the probability distribution of the random measure  $\tilde{\mu}_i^* + \sum_{j=1}^k H_{i,j} \delta_{X_{i,j}^*}$  where

- (a')  $\tilde{\mu}_i^*$  a gamma CRM having intensity  $e^{-(1+U_i)s} s^{-1} ds c \tilde{p}_0^*(dx)$
- (b')  $H_{i,j} = \sum_{t=1}^{\ell_{i,j}} J_{i,j,t}$ , where  $J_{i,j,t} \stackrel{\text{ind}}{\sim} \text{Ga}(q_{i,j,t}, U_i + 1)$ , for  $t = 1, \dots, \ell_{i,j}$ , thus implying that  $H_{i,j} \sim \text{Ga}(n_{i,j}, U_i + 1)$  if  $n_{i,j} \geq 1$  and  $H_{i,j} = 0$  almost surely if  $n_{i,j} = 0$ , by virtue of Theorem 10(ii).

Moreover, note that  $U_i/(1+U_i) \stackrel{\text{ind}}{\sim} \text{Beta}(c, N_i)$ . Hence, by the same arguments as before, one has

$$\tilde{p}_i | (\mathbf{X}, \mathbf{T}, \tilde{p}_0^*) \sim \mathcal{D} \left( c\tilde{p}_0^* + \sum_{j=1}^k n_{i,j} \delta_{X_{i,j}^*} \right)$$

for  $i = 1, \dots, d$ . Note that the dependence on the table configuration  $\mathbf{T}$  is induced solely by  $\tilde{p}_0^*$ , arguably a quite restrictive feature.  $\square$

**EXAMPLE 6.** For a hierarchical stable NRMI one has  $\rho(s) = \sigma s^{-1-\sigma} ds / \Gamma(1-\sigma)$ , for some  $\sigma \in (0, 1)$ ,  $\psi(u) = u^\sigma$  and  $\tau_q(u) = \sigma(1-\sigma)_{q-1} u^{\sigma-q}$ . Similar expressions hold true for  $\rho_0, \tau_{q,0}$  and  $\psi_0$ , with  $\sigma_0 \in (0, 1)$  replacing  $\sigma$ . It is easily seen that  $U_0$  is such that  $U_0^{\sigma_0} \sim \text{Ga}(k, c_0)$  and note that the distribution of  $U_0$  depends on the observations only through  $k$ . Moreover

(a)  $\eta_0^*$  is a CRM with intensity

$$\frac{\sigma_0}{\Gamma(1-\sigma_0)} \frac{e^{-U_0 s}}{s^{\sigma_0+1}} ds c_0 P_0(dx),$$

which is known as generalized gamma CRM (see, e.g. [29]).

(b)  $I_j \stackrel{\text{ind}}{\sim} \text{Ga}(\bar{\ell}_{\bullet,j} - \sigma_0, U_0)$ .

Hence  $\tilde{p}_0^* = (\eta_0^* + \sum_{j=1}^k I_j \delta_{X_j^*}) / (\eta_0^*(\mathbb{X}) + \sum_{j=1}^k I_j)$ . Conditional on  $\tilde{p}_0^*$ , and on  $(\mathbf{X}, \mathbf{T}, \mathbf{U})$ , the CRMs  $\tilde{\mu}_1, \dots, \tilde{\mu}_d$  are independent and each  $\tilde{\mu}_i$  equals, in distribution,  $\tilde{\mu}_i^* + \sum_{j=1}^{k_i} H_{i,j} \delta_{X_{i,j}^*}$ , where

(a')  $\tilde{\mu}_i^*$  is a generalized gamma CRM whose intensity is

$$\frac{\sigma}{\Gamma(1-\sigma)} \frac{e^{-U_i s}}{s^{\sigma+1}} ds c \tilde{p}_0^*(dx);$$

(b')  $H_{i,j} := \sum_{t=1}^{\ell_{i,j}} J_{i,j,t}$ , where  $J_{i,j,t} \stackrel{\text{ind}}{\sim} \text{Ga}(q_{i,j,t} - \sigma, U_i)$ , for  $t = 1, \dots, \ell_{i,j}$ , thus implying that  $H_{i,j} \stackrel{\text{ind}}{\sim} \text{Ga}(n_{i,j} - \ell_{i,j}\sigma, U_i)$  if  $n_{i,j} \geq 1$ , while  $H_{i,j} = 0$  almost surely if  $n_{i,j} = 0$ .

Finally,  $U_i$  is such that  $U_i^\sigma \sim \text{Ga}(k, c)$ . This implies that the posterior distribution of  $(\tilde{p}_1, \dots, \tilde{p}_d)$ , conditional on the data and a suitable latent structure, is a vector of normalized generalized gamma CRMs with fixed points of discontinuity at the data points.  $\square$

5.2. *Hierarchical PY posterior.* Even if not obtained through the normalization of a CRM, the techniques used in Theorems 9–10 apply, with suitable modifications, to the determination of a posterior characterization of the Pitman–Yor process. Hence, assume that data  $\mathbf{X}$  are partially exchangeable as in (1) and the prior  $Q_d$  is characterized by

$$\tilde{p}_i | \tilde{p}_0 \stackrel{\text{iid}}{\sim} \text{PY}(\sigma, \theta; \tilde{p}_0) \quad (i = 1, \dots, d), \quad \tilde{p}_0 \sim \text{PY}(\sigma_0, \theta_0; P_0)$$

where  $\tilde{p}_0 = \tilde{\mu}_0 / \tilde{\mu}_0(\mathbb{X})$  and  $\tilde{p}_i = \tilde{\mu}_i / \tilde{\mu}_i(\mathbb{X})$ , for  $i = 1, \dots, d$  and, recall that, in view of (8), here the random measures  $\tilde{\mu}_0$  and  $\tilde{\mu}_i$  are not completely random. The first step is again the posterior characterization of the root of the hierarchy in terms of  $\tilde{\mu}_0$ .

**THEOREM 11.** *Let  $V_0$  be such that  $V_0^{\sigma_0} \sim \text{Ga}(k + \theta_0 / \sigma_0, 1)$ . Then  $\tilde{\mu}_0 | (\mathbf{X}, \mathbf{T}, V_0)$  equals, in distribution, the random measure  $\eta_0^* + \sum_{j=1}^k I_j \delta_{X_j^*}$ , where  $\eta_0^*$  is a generalized gamma CRM whose intensity is*

$$\frac{\sigma_0}{\Gamma(1 - \sigma_0)} \frac{e^{-V_0 s}}{s^{1 + \sigma_0}} ds P_0(dx),$$

the jumps  $\{I_j : j = 1, \dots, k\}$  and  $\eta_0^*$  are independent and  $I_j \stackrel{\text{ind}}{\sim} \text{Ga}(\bar{\ell}_{\bullet j} - \sigma_0, V_0)$ , for  $j = 1, \dots, k$ .

Given this result, one can establish the following posterior characterization of the vector of random measures  $(\tilde{\mu}_1, \dots, \tilde{\mu}_d)$  whose normalization yields a vector of hierarchical PY processes.

**THEOREM 12.** *Let  $V_i$  be such that  $V_i^\sigma \stackrel{\text{ind}}{\sim} \text{Ga}(\bar{\ell}_{i\bullet} + \theta / \sigma, 1)$ , for  $i = 1, \dots, d$ . Then  $(\tilde{\mu}_1, \dots, \tilde{\mu}_d) | (\mathbf{X}, \mathbf{T}, \mathbf{V}, \tilde{p}_0^*)$  equals, in distribution, the random measure*

$$(\tilde{\mu}_1^*, \dots, \tilde{\mu}_d^*) + \left( \sum_{j=1}^k H_{1,j} \delta_{X_j^*}, \dots, \sum_{j=1}^k H_{d,j} \delta_{X_j^*} \right)$$

where the two summands in the above expression are independent,  $\tilde{p}_0^* = (\eta_0^* + \sum_{j=1}^k I_j \delta_{X_j^*}) / (\eta_0^*(\mathbb{X}) + \sum_{j=1}^k I_j)$  and

- (i)  $\tilde{\mu}_1^*, \dots, \tilde{\mu}_d^*$  are independent and each  $\tilde{\mu}_i^*$  is a generalized gamma CRM with intensity

$$\frac{\sigma}{\Gamma(1 - \sigma)} \frac{e^{-V_i s}}{s^{1 + \sigma}} ds \tilde{p}_0^*(dx)$$

- (ii)  $H_{i,j} \stackrel{\text{ind}}{\sim} \text{Ga}(n_{i,j} - \ell_{i,j} \sigma, V_i)$  if  $n_{i,j} \geq 1$  and  $H_{i,j} = 0$ , almost surely, if  $n_{i,j} = 0$ .

From Theorems 11–12 the posterior distribution of  $\tilde{p}_0$  and of the  $\tilde{p}_i$ 's, conditional on  $\tilde{p}_0$ , immediately follow. However, given the special features of the PY process, one can further simplify such a representation and discard the dependence on the latent random elements  $V_0$  and  $\mathbf{V} = (V_1, \dots, V_d)$  leading to a simple posterior representation, which completes the picture of the posterior behaviour of hierarchical PY process. In stating the result, we set  $k_i = \text{card}\{j : n_{i,j} \geq 1\}$  and agree that the Dirichlet distribution with parameters  $(n_{i,1} - \ell_{i,1}\sigma, \dots, n_{i,k} - \ell_{i,k}\sigma, \theta + \bar{\ell}_{i\bullet}\sigma)$  is on the  $k_i$ -dimensional simplex, after removing the parameters having  $n_{i,j} = 0$ .

**THEOREM 13.** *The posterior distribution of  $\tilde{p}_0$ , conditional on  $(\mathbf{X}, \mathbf{T})$ , equals the distribution of the random probability measure*

$$(24) \quad \sum_{j=1}^k W_j \delta_{X_j^*} + W_{k+1} \tilde{p}_{0,k}$$

where  $(W_1, \dots, W_k)$  is a  $k$ -variate Dirichlet random vector with parameters  $(\bar{\ell}_{\bullet 1} - \sigma_0, \dots, \bar{\ell}_{\bullet k} - \sigma_0, \theta_0 + k\sigma_0)$ ,  $W_{k+1} = 1 - \sum_{i=1}^k W_i$  and  $\tilde{p}_{0,k} \sim \text{PY}(\sigma_0, \theta_0 + k\sigma_0; P_0)$ . Moreover, conditional on  $(\tilde{p}_0, \mathbf{X}, \mathbf{T})$ , the posterior distribution of each  $\tilde{p}_i^* = (\tilde{\mu}_i^* + \sum_{j=1}^k H_{i,j} \delta_{X_j^*}) / (\tilde{\mu}_i^*(\mathbb{X}) + \sum_{j=1}^k H_{i,j})$  equals the distribution of the random measure

$$(25) \quad \sum_{j=1}^k W_{i,j} \delta_{X_j^*} + W_{i,k+1} \tilde{p}_{i,k}$$

where  $(W_{i,1}, \dots, W_{i,k})$  is a  $k$ -variate Dirichlet random vector with parameters  $(n_{i,1} - \ell_{i,1}\sigma, \dots, n_{i,k} - \ell_{i,k}\sigma, \theta + \bar{\ell}_{i\bullet}\sigma)$ ,  $W_{i,k+1} = 1 - \sum_{j=1}^k W_{i,j}$  and  $\tilde{p}_{i,k} | \tilde{p}_0 \stackrel{\text{ind}}{\sim} \text{PY}(\sigma, \theta + \bar{\ell}_{i\bullet}\sigma; \tilde{p}_0)$ .

As previously mentioned, in (25) one has  $\mathbb{P}[W_{i,j} = 0] = 1$  whenever  $n_{i,j} = 0$  and the distribution of  $(W_{i,1}, \dots, W_{i,k})$  degenerates on a lower-dimensional simplex. Both representations (24) and (25) are reminiscent of the one given in the exchangeable case by [38]. The common thread is the so-called *quasi-conjugacy* property characteristic of the PY process. See [33].

**6. Algorithms.** The theoretical findings in Sections 3 and 5 are essential for deriving, respectively, marginal and conditional sampling schemes. Note that, based on the pEPPFs provided in Theorems 3–4, one can derive the predictive distributions associated to hierarchical normalized random measures. However, the analytical complexity inherent to the hierarchical

construction does not allow to deduce closed form expressions. Therefore, the best route for a concrete implementation is represented by the derivation of suitable sampling schemes. In Section 6.1 we state the marginal sampler arising from the pEPPF in the context of prediction problems, when  $\tilde{p}_1, \dots, \tilde{p}_d$  model directly the data and one is interested in specific features of additional samples  $(X_{i,N_i+1}, \dots, X_{i,N_i+m})$ , conditional on  $\mathbf{X}^{(N_i)} = (X_{i,1}, \dots, X_{i,N_i})$ , for  $i = 1, \dots, d$ . The algorithm can be adapted in a straightforward way to mixture models with  $\tilde{p}_1, \dots, \tilde{p}_d$  modeling latent random variables in dependent mixtures. Finally, in Section 6.2 we devise a conditional algorithm, which allows to simulate the trajectories of  $(\tilde{p}_1, \dots, \tilde{p}_d)$  from its posterior distribution. These posterior trajectories can then be immediately used for prediction and mixture modeling.

6.1. *Blackwell–MacQueen urn scheme.* The pEPPFs established in Theorems 3–4 arise upon marginalizing out the hierarchical random probability measures and naturally lend themselves to be used for addressing predictive inferential issues. To be more specific, conditional on observed data  $\mathbf{X}^{(N_i)}$ , we aim at determining the probability distribution of the  $m_i$  additional outcomes for each population  $i = 1, \dots, d$

$$(26) \quad \mathbb{P}[\cap_{i=1}^d \{\mathbf{X}^{(m_i|N_i)} \in A_i\} | \mathbf{X}^{(N_1)}, \dots, \mathbf{X}^{(N_d)}] \\ = \int_{\mathbb{P}_{\mathbb{X}}} \prod_{i=1}^d p_i^{(m_i)}(A_i) Q_d(dp_1, \dots, dp_d | \mathbf{X}^{(N_1)}, \dots, \mathbf{X}^{(N_d)})$$

where  $\mathbf{X}^{(m_i|N_i)} = (X_{i,N_i+1}, \dots, X_{i,N_i+m_i})$  and  $A_i \in \mathcal{X}^{m_i}$ . Based on (26), one can predict specific features of  $\mathbf{X}^{(m_i|N_i)}$ , for  $i = 1, \dots, d$ , such as, e.g., the number of new distinct values in the additional  $m_i$  sample data or the number of distinct values that have appeared  $r$  times in the observed sample  $\mathbf{X}^{(N_i)}$  that will be recorded in  $\mathbf{X}^{(m_i|N_i)}$ . These, and a number of related problems, have been extensively studied in the exchangeable case in view of species sampling applications where such quantities can be seen as measures of species diversity. See, e.g., [13, 28]. The results of this paper allow to cover also the more realistic partially exchangeable case for the first time.

The direct evaluation of (26) is unfeasible and one needs to resort to some simulation scheme. To this end, one may rely on the pEPPF in (12)–(15) to devise a Blackwell–MacQueen urn scheme, for any  $d \geq 2$ , that generates  $\mathbf{X}^{(m_i|N_i)}$  for any hierarchical NRMI. In order to simplify the notation and the description of the algorithm, we consider the case  $d = 2$ . The goal is to generate samples  $X_{1,N_1+1}, \dots, X_{1,N_1+m_1}$  and  $X_{2,N_2+1}, \dots, X_{2,N_2+m_2}$ , conditional on  $\mathbf{X}^{(N_1)}$  and  $\mathbf{X}^{(N_2)}$ , for any two positive integers  $m_1$  and  $m_2$ . One

needs to introduce  $N_1 + m_1 + N_2 + m_2$  latent variables  $T_{1,1}, \dots, T_{1,N_1+m_1}, T_{2,1}, \dots, T_{2,N_2+m_2}$ , which are the labels identifying the tables at which the different costumers are seated in the restaurants. The determination of the full conditionals follows immediately from Theorems 3–4 and, more specifically, (14). The sampler allows one to generate  $(T_{i,1}, \dots, T_{i,N_i})$  and  $(X_{i,N_i+r}, T_{i,N_i+r})$ , for  $r = 1, \dots, m_i$  and  $i = 1, 2$ . In order to provide details on this, the label  $-r$  is used to identify a quantity determined after removing  $r$ -th element. Hence, for each  $i = 1, 2$ , one has

- (1) At  $t = 0$ , start from an initial configuration  $X_{l,N_l+1}^{(0)}, \dots, X_{l,N_l+m_l}^{(0)}$  and  $T_{l,1}^{(0)}, \dots, T_{l,N_l+m_l}^{(0)}$ , for  $l = 1, 2$ .
  - (2) At iteration  $t \geq 1$
- (2.a) With  $X_{i,r} = X_h^*$  generate latent variables  $T_{i,r}^{(t)}$ , for  $r = 1, \dots, N_i$ , from

$$\mathbb{P}(T_{i,r} = \text{“new”} \mid \dots) = w_{h,r} \frac{\Phi_{\bar{\ell}_{i,\bullet}^{-r}+1,i}^{(N_i)}(\mathbf{q}_{i,1}^{-r}, \dots, (\mathbf{q}_{i,h}^{-r}, 1), \dots, \mathbf{q}_{i,k}^{-r})}{\Phi_{\bar{\ell}_{i,\bullet}^{-r},i}^{(N_i-1)}(\mathbf{q}_{1,1}^{-r}, \dots, \mathbf{q}_{1,h}^{-r}, \dots, \mathbf{q}_{i,k}^{-r})}$$

and, for  $\kappa = 1, \dots, \ell_{i,h}^{-r}$ ,

$$\mathbb{P}(T_{i,r} = T_{i,h,\kappa}^{*, -r} \mid \dots) = \frac{\Phi_{\bar{\ell}_{i,\bullet}^{-r},i}^{(N_i)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,h}^{-r} + \mathbf{1}_\kappa, \dots, \mathbf{q}_{i,k}^{-r})}{\Phi_{\bar{\ell}_{i,\bullet}^{-r},i}^{(N_i-1)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,h}^{-r}, \dots, \mathbf{q}_{i,k}^{-r})}$$

where

$$w_{h,r} = \frac{\Phi_{k,0}^{(|\ell^{-r}|+1)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{\bullet h}^{-r} + 1, \dots, \bar{\ell}_{\bullet k}^{-r})}{\Phi_{k,0}^{(|\ell^{-r}|)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{\bullet h}^{-r}, \dots, \bar{\ell}_{\bullet k}^{-r})} \mathbb{1}_{\{0\}^c}(\bar{\ell}_{i,h}^{-r}) + \mathbb{1}_{\{0\}}(\bar{\ell}_{i,h}^{-r})$$

and  $\mathbf{1}_\kappa$  is a vector of dimension  $\ell_{i,h}^{-r}$  with all components being zero but the  $\kappa$ -th which equals 1. Moreover,  $T_{i,h,1}^{*, -r}, \dots, T_{i,h,\ell_{i,h}^{-r}}^{*, -r}$  are the tables at the first restaurant where the  $h$ -th dish is served, after the removal of  $T_{i,r}$ .

- (2.b) For  $r = 1, \dots, m_i$ , generate  $(X_{i,N_i+r}^{(t)}, T_{i,N_i+r}^{(t)})$  from the following predictive distributions

$$\mathbb{P}(X_{i,N_i+r} = \text{“new”}, T_{i,N_i+r} = \text{“new”} \mid \dots) = \frac{\Phi_{k+j-r+1,0}^{(|\ell^{-r}|+1)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{\bullet k+j-r}^{-r}, 1)}{\Phi_{k+j-r,0}^{(|\ell^{-r}|)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{\bullet k+j-r}^{-r})} \frac{\Phi_{\bar{\ell}_{i,\bullet}^{-r}+1,i}^{(N_i+m_i)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,k}^{-r}, 1)}{\Phi_{\bar{\ell}_{i,\bullet}^{-r},i}^{(N_i+m_i-1)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,k}^{-r})}$$

while, for any  $h = 1, \dots, k + j^{-r}$  and  $\kappa = 1, \dots, \ell_{i,h}^{-r}$ ,

$$\begin{aligned} \mathbb{P}(X_{i,N_i+r} = X_h^{*, -r}, T_{i,N_i+r} = \text{“new”} | \dots) &= \\ &= \frac{\Phi_{k+j^{-r},0}^{(|\ell^{-r}|+1)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{1,h}^{-r} + 1, \dots, \bar{\ell}_{\bullet k+j^{-r}}^{-r})}{\Phi_{k+j^{-r},0}^{(|\ell^{-r}|)}(\bar{\ell}_{\bullet 1}^{-r}, \dots, \bar{\ell}_{\bullet k+j^{-r}}^{-r})} \\ &\quad \times \frac{\Phi_{\bar{\ell}_{i\bullet}^{-r}+1,i}^{(N_i+m_i)}(\mathbf{q}_{i,1}^{-r}, \dots, (\mathbf{q}_{i,h}^{-r}, 1), \dots, \mathbf{q}_{i,k}^{-r})}{\Phi_{\bar{\ell}_{i\bullet}^{-r},i}^{(N_i+m_i-1)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,k}^{-r})} \\ \mathbb{P}(X_{i,N_i+r} = X_h^{*, -r}, T_{i,N_i+r} = T_{i,h,\kappa}^{*, -r} | \dots) &= \\ &= \frac{\Phi_{\bar{\ell}_{i\bullet}^{-r},i}^{(N_i+m_i)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,h}^{-r} + \mathbf{1}_\kappa, \dots, \mathbf{q}_{i,k}^{-r}, 1)}{\Phi_{\bar{\ell}_{i\bullet}^{-r},i}^{(N_i+m_i-1)}(\mathbf{q}_{i,1}^{-r}, \dots, \mathbf{q}_{i,h}^{-r}, \dots, \mathbf{q}_{i,k}^{-r})} \mathbb{1}_{\{n_{i,h}^{-r} > 0\}} \end{aligned}$$

where  $X_h^{*, -r}$ , for  $h = 1, \dots, k + j^{-r}$  denote the distinct dishes in the whole franchise after the removal of the  $r$ -th observation, while the condition  $n_{i,h}^{-r} > 0$  entails that the  $h$ -th dish is served in the  $i$ -th restaurant.

The above algorithm holds for any hierarchical NRMI and only requires insertion of the specific  $\rho$ ,  $\rho_0$  and  $P_0$  to specialize to a particular instance of hierarchical NRMI. The sampling schemes outlined above can also be tailored, in a quite straightforward way, to the hierarchical Pitman–Yor case (see the supplementary material [3] for details and [4] for applications). Finally note that the proposed algorithm can also be adapted to yield a marginal sampling schemes for mixture models with dependent hierarchical mixing measures.

*6.2. Simulation of  $(\tilde{p}_1, \dots, \tilde{p}_d)$  from its posterior distribution.* The posterior representations derived in Theorems 10 and 13 are of great importance also from a computational standpoint as they allow to establish algorithms that generate the trajectories of  $\tilde{p}_1, \dots, \tilde{p}_d$  from their posterior distributions, conditional on  $\mathbf{T}$ . The resulting sampling scheme can be viewed as an extension of a Ferguson & Klass-type algorithm (see [15, 44] for additional details) to a partially exchangeable setting. With respect to the generalized Blackwell–MacQueen urn scheme described in Section 6.1, the possibility of generating posterior samples of hierarchical processes is a significant addition. Just to give an example, it allows to obtain estimates of non-linear

functionals, such as credible intervals, of the vector  $(\tilde{p}_1, \dots, \tilde{p}_d)$  that cannot be otherwise achieved.

For the sake of simplicity assume that  $\mathbb{X} = \mathbb{R}^+$ . Using a representation of  $X_t$  given in [15] and the notation of Theorems 9–10, one has

$$(27) \quad \eta_0^*((0, t]) = \sum_{h=1}^{\infty} J_h^{(0)} \mathbb{1}\{V_h \leq P_0((0, t])\},$$

with  $V_1, V_2, \dots \stackrel{\text{iid}}{\sim} U(0, 1)$ . The jumps  $J_h^{(0)}$  are in decreasing order and can be obtained by solving the following

$$(28) \quad S_{h,0} = c_0 \int_{J_h^{(0)}}^{\infty} e^{-U_0 s} \rho_0(s) \, ds.$$

where  $S_{1,0}, S_{2,0}, \dots$  are the points of a standard Poisson process on  $\mathbb{R}^+$ , that is to say  $S_{h,0} - S_{h-1,0}$  are i.i.d. exponential random variables having unit mean. Similarly, one has

$$(29) \quad \tilde{\mu}_i^*((0, t]) = \sum_{h=1}^{\infty} J_h^{(i)} \mathbb{1}\{V_h \leq \tilde{p}_0^*((0, t])\},$$

where the ordered jumps  $J_h^{(i)}$  are now the solution of

$$(30) \quad S_{h,i} = c \int_{J_h^{(i)}}^{\infty} e^{-U_i s} \rho(s) \, ds,$$

where  $S_{1,i}, S_{2,i} - S_{1,i}, \dots$  are i.i.d. exponential random variables having unit mean. In view of these representations, once one has sampled the latent variables  $\mathbf{T}$  through the algorithm described in Section 6.1, one can proceed as follows:

- (1) Generate  $\tilde{p}_0$  from its posterior distribution, described in Theorem 9, namely:
  - (1.a) Generate  $U_0$  from  $f_0(\cdot | \mathbf{X}, \mathbf{T})$  in (20);
  - (1.b) Generate  $I_j$  from  $f_j(\cdot | \mathbf{X}, \mathbf{T})$  in Theorem 9(ii), for any  $j = 1, \dots, k$ ;
  - (1.c) Fix  $\varepsilon > 0$  and for any  $h \geq 1$ 
    - Generate unit mean exponential random variables  $S_{h,0} - S_{h-1,0}$
    - Determine jumps  $J_h^{(0)}$  according to (28)
    - Stop at  $\bar{h} = \min\{h \geq 1 : J_h^{(0)} \leq \varepsilon\}$



- Generate i.i.d.  $V_1, \dots, V_{\bar{h}}$  from a  $U(0, 1)$
- and evaluate an approximate draw of  $\eta_0^*$  on  $(0, t]$  as

$$\eta_0^*((0, t]) \approx \sum_{h=1}^{\bar{h}} J_h^{(0)} \mathbf{1}\{V_h \leq P_0((0, t])\},$$

- (1. d) Evaluate an approximate draw of a posterior sample of  $\tilde{p}_0$  as

$$\tilde{p}_0^*((0, t]) \approx \frac{\sum_{h=1}^{\bar{h}} J_h^{(0)} \mathbf{1}\{V_h \leq P_0((0, t])\} + \sum_{j=1}^k I_j \delta_{X_j^*}((0, t])}{\sum_{h=1}^{\bar{h}} J_h^{(0)} + \sum_{j=1}^k I_j}.$$

Having drawn  $\tilde{p}_0^*$ , one can now rely on Theorem 10 in order to approximately sample  $(\tilde{p}_1, \dots, \tilde{p}_d)$  from its posterior distribution. This can be easily deduced and described as follows.

- (2) For any  $i = 1, \dots, d$ , generate  $\tilde{p}_i | (\mathbf{X}, \mathbf{T}, \tilde{p}_0^*)$  as follows
  - (2. a) Generate  $U_i$  from  $f_i(\cdot | \mathbf{X}, \mathbf{T})$  in (22);
  - (2. b) Generate  $J_{i,j,t}$  from  $f_{i,j,t}(\cdot | \mathbf{X}, \mathbf{T})$  in Theorem 10(ii)
  - (2. c) Fix  $\varepsilon > 0$  and for any  $h \geq 1$ 
    - Generate unit mean exponential random variables  $S_{h,i} - S_{h-1,i}$
    - Determine jumps  $J_h^{(i)}$  according to (30)
    - Stop at  $\bar{h}_i = \min\{h \geq 1 : J_h^{(i)} \leq \varepsilon\}$

and evaluate an approximate sample of the posterior trajectory of  $\tilde{p}_i$  as follows

$$\tilde{p}_i((0, t]) \approx \frac{\sum_{h=1}^{\bar{h}_i} J_h^{(i)} \mathbf{1}\{V_h \leq \tilde{p}_0^*((0, t])\} + \sum_{j=1}^k \sum_{t=1}^{\ell_{i,j}} J_{i,j,t} \delta_{X_j^*}((0, t])}{\sum_{h=1}^{\bar{h}_i} J_h^{(i)} + \sum_{j=1}^{k_i} \sum_{t=1}^{\ell_{i,j}} J_{i,j,t}}.$$

An important, and well-known, advantage of the procedure is the fact that it generates jumps  $J_h^{(0)}$  and  $J_{\bar{h}_i}^{(i)}$ , for  $i = 1, \dots, d$ , in decreasing order. This entails that the truncation at  $\bar{h}$  or  $\bar{h}_i$  is such that the most relevant jumps are taken into account and one is discarding a negligible random mass of the actual trajectory. Future work, of more computational nature, will aim at: (i) investigating the implementation of the algorithm to applied problems, such as density estimation with accurate uncertainty quantification, allowed by the conditional structure of the algorithm and (ii) comparing the performance of our proposal with the so-called direct assignment algorithm, which is widely used within estimation problems involving the hierarchical Dirichlet process.

## SUPPLEMENTARY MATERIAL

**Supplement A: Distribution theory for hierarchical processes: supplementary material**

(doi: [COMPLETED BY THE TYPESETTER](#); supplementary.pdf). We provide the proofs of the theoretical results and specialize the Blackwell–MacQueen urn scheme of Section 6.1 to the case of hierarchies of Pitman–Yor processes.

**Acknowledgments.** The authors are grateful to an Associate Editor and three anonymous Referees for their valuable comments and insightful suggestions, which led to a substantial improvement of the paper. F. Camerlenghi has completed the paper while being a postdoctoral fellow at the University of Bologna. He is deeply grateful to the Department of Statistical Sciences for the support.

**References.**

- [1] ADAMS, M., KELLEY, J., POLYMEROPOULOS, M., XIAO, H., MERRIL, C., WU, A., OLDE, B., MORENO, R., KERLAVAGE, A., MCCOMBE, W. and VENTER, J. (1991). Complementary DNA sequencing: Expressed Sequence Tags and human genome project. *Science* **252**, 1651–1656.
- [2] BLEI, D.M., NG, A.Y. and JORDAN, M.I. (2003). Latent Dirichlet allocation. *J. Mach. Learn. Res.* **3**, 993–1022.
- [3] CAMERLENGHI, F., LIJOI, A., ORBANZ, P. and PRÜNSTER, I. (2017). Distribution theory for hierarchical processes: supplementary material.
- [4] CAMERLENGHI, F., LIJOI, A. and PRÜNSTER, I. (2017). Bayesian prediction with multiple-sample information. *J. Multivariate Anal.*, **156**, 18–28.
- [5] CARON, F., DAVY, M., DOUCET, A. (2007). Generalized Pólya urn for time-varying Dirichlet process mixtures. *23rd Conference on Uncertainty in Artificial Intelligence (UAI'2007)*, Vancouver, Canada.
- [6] CARON, F., TEH, Y.W., MURPHY, T.B. (2014). Bayesian nonparametric Plackett-Luce models for the analysis of preferences for college degree programmes. *Ann. Appl. Statist.* **8**, 1145–1181.
- [7] CHARALAMBIDES, C.A. (2005). *Combinatorial methods in discrete distributions*. Hoboken, NJ: Wiley.
- [8] CONSTANTINES G.M., SAVITS T.H. (1996). A multivariate version of the Faa di Bruno formula. *Trans. Amer. Math. Soc.* **348**, 503–520.
- [9] DALEY, D.J. and VERE–JONES, D. (2008). *An introduction to the theory of point processes. Volume II*, Springer, New York.
- [10] DE BLASI, P., FAVARO, S., LIJOI, A., MENA, R.H., RUGGIERO, M. and PRÜNSTER, I. (2015). Are Gibbs–type priors the most natural generalization of the Dirichlet process? *IEEE Trans. Pattern Anal. Mach. Intell.* **37**, 212–229.
- [11] DUNSON, D.B. (2010). Nonparametric Bayes applications to biostatistics. In *Bayesian Nonparametrics*, pp. 223–273. Cambridge University Press, Cambridge.
- [12] EWENS, W.J. (1972). The sampling theory of selectively neutral alleles. *Theor. Popul. Biol.* **3**, 87–112.
- [13] FAVARO, S., LIJOI, A., MENA, R.H. and PRÜNSTER, I. (2009). Bayesian nonparametric inference for species variety with a two parameter Poisson–Dirichlet process prior. *J. Roy. Statist. Soc. Ser. B* **71**, 993–1008.

- [14] FERGUSON, T.S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209–230.
- [15] FERGUSON, T.S. and KLASS, M.J. (1972). A representation of independent increment processes without Gaussian components. *Ann. Math. Statist.* **43**, 1634–1643.
- [16] FOTI, N.J. and WILLIAMSON, S.A. (2015). A survey of non-exchangeable priors for Bayesian nonparametric models. *IEEE Trans. Pattern Anal. Mach. Intell.* **37**, 359–371.
- [17] GASTHAUS, J. and TEH, Y.W. (2010). Improvements to the sequence memoizer. *Advances in Neuronal Information Processing Systems* **23**.
- [18] GNEDIN, A.V. and PITMAN, J. (2005). Exchangeable Gibbs partitions and Stirling triangles. *Zap. Nauchn. Sem. S.Peterb. Otdel. Mat. Inst. Steklov (POMI)* **325**, 83–102.
- [19] GRIFFIN, J. E., KOLOSSIATIS, M. and STEEL, M. F. J. (2013) Comparing distributions by using dependent normalized random-measure mixtures. *J. R. Stat. Soc. Ser. B.* **75**, 499–529.
- [20] JAMES, L.F. (2005). Bayesian Poisson process partition calculus with an application to Bayesian Lévy moving averages. *Ann. Statist.* **33**, 1771–1799.
- [21] JAMES, L.F. (2006). Poisson calculus for spatial neutral to the right processes. *Ann. Statist.* **34**, 416–440.
- [22] JAMES, L.F., LIJOI, A. and PRÜNSTER, I. (2006). Conjugacy as a distinctive feature of the Dirichlet process. *Scand. J. Statist.* **33**, 105–120.
- [23] JAMES, L.F., LIJOI, A. and PRÜNSTER, I. (2009). Posterior analysis for normalized random measures with independent increments. *Scand. J. Statist.* **36**, 76–97.
- [24] KALLENBERG, O. (2005) *Probabilistic Symmetries and Invariance Principles*. Springer, New York.
- [25] KINGMAN, J.F.C. (1967). Completely random measures. *Pacific J. Math.* **21**, 59–78.
- [26] KINGMAN, J.F.C. (1982). The coalescent. *Stochastic Process. Appl.* **13**, 235–248.
- [27] KINGMAN, J.F.C. (1993). *Poisson processes*. Oxford University Press, Oxford.
- [28] LIJOI, A., MENA, R.H. and PRÜNSTER, I. (2007). Bayesian nonparametric estimation of the probability of discovering a new species *Biometrika.* **94**, 769–786.
- [29] LIJOI, A., MENA, R.H. and PRÜNSTER, I. (2007a). Controlling the reinforcement in Bayesian nonparametric mixture models. *J. Roy. Statist. Soc. Ser. B* **69**, 715–740.
- [30] LIJOI, A., NIPOTI, B. AND PRÜNSTER, I. (2014). Bayesian inference with dependent normalized completely random measures. *Bernoulli* **20**, 1260–1291.
- [31] LIJOI, A., NIPOTI, B. AND PRÜNSTER, I. (2014). Dependent mixture models: clustering and borrowing information. *Comput. Statist. Data Anal.* **71**, 417–433.
- [32] LIJOI, A. and PRÜNSTER, I. (2010). Models beyond the Dirichlet process. In *Bayesian Nonparametrics*, pp. 80–136. Cambridge University Press, Cambridge.
- [33] LIJOI, A., PRÜNSTER, I. and WALKER, S.G. (2008a). Bayesian nonparametric estimators derived from conditional Gibbs structures. *Ann. Appl. Probab.* **18**, 1519–1547.
- [34] MACEACHERN, S.N. (1999). Dependent nonparametric processes. In *ASA Proceedings of the SBSS*. Alexandria: American Statistical Association, 50–55.
- [35] MACEACHERN, S.N. (2000). Dependent Dirichlet processes. *Technical Report*. Department of Statistics, Ohio State University.
- [36] NGUYEN, X. (2010). Inference of global clusters from locally distributed data. *Bayesian Analysis* **5**, 817–846.
- [37] NGUYEN, X. (2016). Borrowing strength in hierarchical Bayes: convergence of the Dirichlet base measure. *Bernoulli* **22**, 1535–1571.
- [38] PITMAN, J. (1996). Some developments of the Blackwell-MacQueen urn scheme. In *Statistics, Probability and Game Theory* (T.S. Ferguson, L.S. Shapley and J.B. MacQueen, Eds.). IMS Lecture Notes Monogr. Ser., Vol. **30**. IMS, Hayward, 245–267.

- [39] PITMAN, J. (2006). *Combinatorial stochastic processes*. Ecole d'Été de Probabilités de Saint-Flour XXXII. Lecture Notes in Mathematics N. 1875, Springer, New York.
- [40] PITMAN, J. and YOR, M. (1997). The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *Ann. Probab.* **25**, 855–900.
- [41] REGAZZINI, E., LIJOI, A. and PRÜNSTER, I. (2003). Distributional results for means of normalized random measures with independent increments. *Ann. Statist.* **31**, 560–585.
- [42] TEH, Y.W., JORDAN, M.I., BEAL, M.J. and BLEI, D.M. (2006). Hierarchical Dirichlet processes. *J. Amer. Statist. Assoc.* **101**, 1566–1581.
- [43] TEH, Y.W., JORDAN, M.I. (2010). Hierarchical Bayesian nonparametric models with applications. In *Bayesian Nonparametrics*, pp. 158-207, Cambridge Univ. Press, Cambridge.
- [44] WALKER, S. and DAMIEN, P. (2000). Representations of Lévy processes without Gaussian components. *Biometrika.* **87**, 477–483.
- [45] WOOD, F., GASTHAUS, J., ARCHAMBEAU, C., JAMES, L.F. and TEH, Y.W. (2011). The sequence memoizer. *Communications ACM* **54**, 91–98.
- [46] ZHOU, M. and CARIN, L. (2015). Negative binomial process count and mixture modeling. *IEEE Trans. Pattern Anal. Mach. Intell.* **37**, 307-320.

FEDERICO CAMERLENGHI  
 DEPARTMENT OF ECONOMICS,  
 MANAGEMENT AND STATISTICS  
 UNIVERSITY OF MILANO–BICOCCA  
 PIAZZA DELL'ATENEO NUOVO 1, 20126 MILANO, ITALY  
 E-MAIL: [federico.camerlenghi@unimib.it](mailto:federico.camerlenghi@unimib.it)

ANTONIO LIJOI AND IGOR PRÜNSTER  
 DEPARTMENT OF DECISION SCIENCES AND BIDS  
 BOCCONI UNIVERSITY  
 VIA RÖNTGEN 1, 20136 MILANO, ITALY  
 E-MAIL: [antonio.ljoi@unibocconi.it](mailto:antonio.ljoi@unibocconi.it); [igor@unibocconi.it](mailto:igor@unibocconi.it)

PETER ORBANZ  
 DEPARTMENT OF STATISTICS  
 COLUMBIA UNIVERSITY  
 1255 AMSTERDAM AVENUE 10027 NEW YORK, USA  
 E-MAIL: [porbanz@stat.columbia.edu](mailto:porbanz@stat.columbia.edu)