

UNORDERED AND ORDERED SAMPLE FROM DIRICHLET DISTRIBUTION

THIERRY HUILLET

*Laboratoire de Physique Théorique et Modélisation, CNRS-UMR 8089 et Université de Cergy-Pontoise, 5 mail Gay-Lussac, 95031, Neuville sur Oise, France,
e-mail: Thierry.Huillet@ptm.u-cergy.fr*

(Received March 2, 2004; revised September 1, 2004)

Abstract. Consider the random Dirichlet partition of the interval into n fragments with parameter $\theta > 0$. Explicit results on the statistical structure of its size-biased permutation are recalled, leading to (unordered) Ewens and (ordered) Donnelly-Tavaré-Griffiths sampling formulae from finite Dirichlet partitions. We use these preliminary statistical results on frequencies distribution to address the following sampling problem: what are the intervals between new sampled categories when sampling is from Dirichlet populations? The results obtained are in accordance with the ones found in sampling theory from random proportions with $GEM(\gamma)$ distribution. These can be obtained from Dirichlet model when considering the Kingman limit $n \uparrow \infty$, $\theta \downarrow 0$ while $n\theta = \gamma > 0$.

Key words and phrases: Random discrete distribution, Dirichlet partition, size-biased permutation, GEM, Ewens and Donnelly-Tavaré-Griffiths sampling formulae, intervals between new sampled species.

1. Introduction

The joint distribution of unordered (or ordered) frequencies of a sample from random proportions with Griffiths-Engen-McCloskey, or $GEM(\gamma)$, distribution is the Donnelly-Tavaré-Griffiths formula (or the Ewens sampling formulae). The $GEM(\gamma)$ distribution will be defined later (see also Kingman (1993), Chapter 9, for example).

We first reconsider the same sampling problems and formulae when sampling is from random proportions with Dirichlet $D_n(\theta)$ distribution, hence with a finite number n of fragments in the partition. The usual Ewens and Donnelly-Tavaré-Griffiths sampling formulae can be found when passing to the Kingman limit $n \uparrow \infty$, $\theta \downarrow 0$, $n\theta = \gamma > 0$. These preliminary results which are recalled are used to study the intervals between consecutive categories and number of distinct categories when sampling is from $D_n(\theta)$. The results obtained are in accordance with particular cases of those of Yamato-Sibuya-Nomachi when sampling is from $GEM(\gamma)$, while passing to the Kingman limit.

The organization of this manuscript is the following. Consider the random Dirichlet partition of the interval into n fragments with parameter $\theta > 0$. Elementary properties of its $D_n(\theta)$ distribution are first recalled in Section 2. Explicit results on the law of its size-biased permutation are also recalled there. These are useful when considering the Donnelly-Tavaré-Griffiths sampling formula from Dirichlet partitions in subsequent Subsection 3.3, as given in Theorem 3.3. A size-biased permutation of the fragments

sizes is the one obtained in a size-biased sampling process without replacement from a Dirichlet partition. The main points which we recall are the following: in Lemma 2.1, its residual allocation model structure is recalled. In Lemma 2.2, the order in which the consecutive fragments are visited is considered. In Theorem 2.1, we recall the joint law of the size-biased permutation fragments sizes explicitly. Using this, it is recalled in Corollary 2.1 that consecutive fragments in the size-biased permuted partition are arranged in stochastic descending order.

Section 3 recalls the (unordered) Ewens and (ordered) Donnelly-Tavaré-Griffiths sampling formulae when sampling is from finite Dirichlet partitions. In more details, Subsection 3.1 is devoted to the first Ewens sampling formula when sampling is from Dirichlet partition $D_n(\theta)$. Here the order in which sequentially sampled species arise is irrelevant. Subsection 3.2 concerns the second Ewens sampling formula under the same hypothesis (as a problem of random partitioning of the integers) and Subsection 3.3 deals with the finite Dirichlet version of the Donnelly-Tavaré-Griffiths sampling formula. Here, the order of appearance of sampled species is taken into account. Main results are displayed in Theorems 3.1, 3.2 and 3.3 for each of the problems alluded to. As corollaries to these theorems, the usual well-known sampling formulae can be deduced in each case when sampling is from *GEM* distribution which is the limiting version of the size-biased permutation of Dirichlet partitions in the sense of Kingman.

The main body of our new sampling results is in Section 4. We use the previous statistical results on frequencies distribution to address the following sampling problem: what are the intervals between new sampled categories until exhaustion of the list, when sampling is from Dirichlet populations? One intuitively expects these intervals to be increasing while approaching complete exhaustion. Our main results on the statistical structure of these intervals are summarized in Theorems 4.1, 4.2 and 4.3, considering both cases of a fixed and then unlimited sample size.

The results obtained are in accordance with the Yamato-Sibuya-Nomachi ones when sampling is from random proportions with *GEM*(γ) distribution. These can naturally be obtained from Dirichlet model when considering the Kingman limit. Results in this direction are displayed in Corollaries 4.1, 4.2 and 4.3 to Theorems 4.1, 4.2 and 4.3, respectively.

2. Size-biased permutation of the Dirichlet distribution

2.1 Dirichlet partition of the interval

Consider the following random partition into n fragments of the unit interval. Let $\theta > 0$ be some parameter and assume that the random fragments' sizes $\mathbf{S}_n := (S_1, \dots, S_n)$ (with $\sum_{m=1}^n S_m = 1$) is distributed according to the (exchangeable) Dirichlet $D_n(\theta)$ density function on the simplex, that is to say

$$(2.1) \quad f_{S_1, \dots, S_n}(s_1, \dots, s_n) = \frac{\Gamma(n\theta)}{\Gamma(\theta)^n} \prod_{m=1}^n s_m^{\theta-1} \cdot \delta_{(\sum_{m=1}^n s_m - 1)}.$$

Alternatively, the law of $\mathbf{S}_n := (S_1, \dots, S_n)$ is characterized by its joint moment function

$$(2.2) \quad \mathbf{E} \left[\prod_{m=1}^n S_m^{q_m} \right] = \frac{\Gamma(n\theta)}{\Gamma(n\theta + \sum_{m=1}^n q_m)} \prod_{m=1}^n \frac{\Gamma(\theta + q_m)}{\Gamma(\theta)}.$$

We shall put $\mathbf{S}_n \stackrel{d}{\sim} D_n(\theta)$ if \mathbf{S}_n is Dirichlet distributed with parameter θ .

If this is so, $S_m \stackrel{d}{=} S_n$, $m = 1, \dots, n$, independently of m and the individual fragments sizes are all identically distributed. Their common density on the interval $(0, 1)$ is a beta($\theta, (n - 1)\theta$) density.

When $\theta = 1$, the partition model equations (2.1), (2.2) corresponds to the standard uniform partition model of the interval.

In the random division of the interval as in equation (2.1), although all fragments are identically distributed with sizes of order n^{-1} , the smallest fragment's size grows like $n^{-(\theta+1)/\theta}$ while the one of the largest is of order $\frac{1}{n\theta} \log(n \log^{\theta-1} n)$. The smaller θ is, the larger (smaller) the largest (smallest) fragments' size is: hence, the smaller θ is, the more the values of the S_m s are, with high probability, disparate. Smaller the parameter, the size of the largest fragment $S_{(1)}$ tends to dominate the other ones. On the contrary, large values of θ correspond to situations in which the range of fragments' sizes is lower: the fragments' sizes look more homogeneous and distribution equation (2.1) concentrates on its centre. For large θ , the diversity of the partition is small.

Although \mathbf{S}_n has a degenerate weak limit when $n \uparrow \infty$, $\theta \downarrow 0$ while $n\theta = \gamma > 0$, this situation is worth being considered (as first noted by Kingman (1975)). Indeed, many interesting statistical features emerge.

2.2 Size-biased permutation of Dirichlet partitions

The results on size-biased permutation of Dirichlet distributions presented in this section are not new. When $\theta = 1$, they can be found in Huillet (2003); they were generalized to all $\theta > 0$ in Barrera *et al.* (2005), to solve a problem consisting in computing the search-cost distribution arising from heaps processes. Part of them are reproduced here for the sake of completeness and to make the presentation self-contained. They will prove useful in the sequel.

Assume some observer is sampling the unit interval as follows: drop points at random onto this randomly broken interval and record the corresponding numbers of visited fragments. We shall consider the problem of determining the order in which the various fragments are discovered in such a sampling process. To avoid revisiting the same fragment many times, once it has been discovered, we need to remove it from the population as soon as it has been met in the sampling process. But to do that, the law of its size is needed. Once this is done, after renormalizing the remaining fragments' sizes, we are left with a population of $n - 1$ fragments, the sampling of which will necessarily supply a so far undiscovered fragment. The distribution of its size can itself be computed and so forth, renormalizing again, until the whole available fragments population has been visited. In this way, not only the visiting order of the different fragments can be understood but also their sizes. The purpose of this section is to describe the statistical structure of the size-biased permutation of the fragments' sizes as those obtained while avoiding the ones previously encountered in a sampling process from Dirichlet partition.

- *The RAM structure of size-biased permutation.* Let $\mathbf{S}_n := (S_1, \dots, S_n)$ be the random partition of the interval $[0, 1]$ considered here. Let L_1 be the length of the first randomly chosen fragment $M_1 := M$, so with $L_1 := S_{M_1}$ and $\mathbf{P}(M_1 = m_1 | \mathbf{S}_n) = S_{m_1}$. One may check that $L_1 \stackrel{d}{\sim} \text{beta}(1 + \theta, (n - 1)\theta)$. A standard problem is to iterate the size-biased picking procedure, by avoiding the fragments already encountered: by doing so, a size-biased permutation (SBP) of the fragments is obtained. It turns out that $\text{SBP}(\mathbf{S}_n)$ has a residual allocation model (RAM) structure.

In the first step of this size-biased picking procedure indeed, $\mathcal{S}_n =: \mathcal{S}_n^{(0)} \rightarrow (L_1, S_1, \dots, S_{M_1-1}, S_{M_1+1}, \dots, S_n)$ which may be written as $\mathcal{S}_n \rightarrow (L_1, (1 - L_1)\mathcal{S}_{n-1}^{(1)})$, with $\mathcal{S}_{n-1}^{(1)} := (S_1^{(1)}, \dots, S_{M_1-1}^{(1)}, S_{M_1+1}^{(1)}, \dots, S_n^{(1)})$ a new random partition of the unit interval into $n - 1$ random fragments.

Given $L_1 \stackrel{d}{\sim} \text{beta}(1 + \theta, (n - 1)\theta)$, the conditional joint distribution of the remaining components of \mathcal{S}_n is the same as that of $(1 - L_1)\mathcal{S}_{n-1}^{(1)}$ where the $(n - 1)$ -vector $\mathcal{S}_{n-1}^{(1)} \stackrel{d}{\sim} D_{n-1}(\theta)$ has the distribution of a Dirichlet random partition into $n - 1$ fragments (see Kingman (1993), Chapter 9). Furthermore, $\mathcal{S}_{n-1}^{(1)}$ is independent of $1 - L_1$. Pick next at random an interval in $\mathcal{S}_{n-1}^{(1)}$ and call V_2 its length, now with distribution $\text{beta}(1 + \theta, (n - 2)\theta)$, and iterate until all fragments have been exhausted.

With $V_1 := L_1$, the length of the second fragment by avoiding the first reads $L_2 = (1 - V_1)V_2$. Iterating, the final SBP of \mathcal{S}_n is $L_n := (L_1, \dots, L_n)$ and we shall put $L_n = \text{SBP}(\mathcal{S}_n)$. From this construction, we easily get

LEMMA 2.1. *Let $L_n = \text{SBP}(\mathcal{S}_n)$. If (V_1, \dots, V_{n-1}) is an independent sample with distribution $V_k \stackrel{d}{\sim} \text{beta}(1 + \theta, (n - k)\theta)$, $k = 1, \dots, n - 1$, then,*

$$(2.3) \quad L_k = \prod_{i=1}^{k-1} (1 - V_i)V_k, \quad k = 1, \dots, n - 1$$

$$(2.4) \quad L_n = 1 - \sum_{k=1}^{n-1} L_k = \prod_{k=1}^{n-1} (1 - V_k)$$

is the RAM representation of the size-biased permutation L_n of \mathcal{S}_n .

Note that $\bar{V}_i := 1 - V_i \stackrel{d}{\sim} \text{beta}((n - i)\theta, 1 + \theta)$ and that V_n should be set to 1. These are well-known construction and properties; see Kingman ((1993), Chapter 9, Section 6) and Donnelly (1986, 1991).

This RAM representation allows to compute the joint distribution of the size-biased permutation L_n of \mathcal{S}_n . We shall say in the sequel that, if $L_n = \text{SBP}(\mathcal{S}_n)$, then $L_n \stackrel{d}{\sim} \text{SBD}_n(\theta)$ assuming that $\mathcal{S}_n \stackrel{d}{\sim} D_n(\theta)$. Before addressing this problem, we shall first consider the order in which fragments are visited.

• *The visiting order of the fragments in the SBP process.* For any permutation $\{m_1, \dots, m_n\}$ of $\{1, \dots, n\}$, with M'_1, \dots, M'_k , $k = 1, \dots, n$, the first k distinct fragments numbers which have been visited in the SBP sampling process, we have

$$(2.5) \quad P(M'_1 = m_1, \dots, M'_k = m_k \mid \mathcal{S}_n) = \prod_{i=1}^{k-1} \frac{S_{m_i}}{1 - \sum_{l=1}^i S_{m_l}} S_{m_k},$$

so that

$$(2.6) \quad P(M'_k = m_k \mid \mathcal{S}_n, M'_1 = m_1, \dots, M'_{k-1} = m_{k-1}) = \frac{S_{m_k}}{1 - \sum_{l=1}^{k-1} S_{m_l}}.$$

As a result,

$$(2.7) \quad P(M'_k = m \mid \mathcal{S}_n) = S_m \sum_{(m_1 \neq \dots \neq m_{k-1}) \neq m} \prod_{i=1}^{k-1} \frac{S_{m_i}}{1 - \sum_{l=1}^i S_{m_l}}$$

is the conditional (given S_n) probability that the k -th visited fragment is fragment number m from $D_n(\theta)$.

LEMMA 2.2. *Given $M'_1 = m_1, \dots, M'_{k-1} = m_{k-1}$, M'_k is uniformly distributed on the set $m \in \{1, \dots, n\} \setminus \{m_1 \neq \dots \neq m_{k-1}\}$. The joint probability distribution of M'_1, \dots, M'_k , $k = 1, \dots, n$, is Bose-Einstein distribution*

$$(2.8) \quad P_\theta(M'_1 = m_1, \dots, M'_k = m_k) = \frac{1}{\prod_{i=0}^{k-1} (n-i)},$$

with $m_1 \neq \dots \neq m_k \in \{1, \dots, n\}$.

PROOF. Although this result is immediate by symmetry, we shall supply a short proof of it. From equation (2.6), the function $S_n \rightarrow \frac{S_{m_k}}{1 - \sum_{l=1}^{k-1} S_{m_l}} = \frac{S_{m_k}}{\sum_{m \neq \{m_1, \dots, m_{k-1}\}} S_m}$ is homogeneous with degree 0. Applying (ii) of Theorem 1 p. 471 of Huillet and Martinez (2003), with $\{T_1, \dots, T_n\}$ iid $\text{gamma}(\theta)$ distributed random variables, we get

$$\begin{aligned} P_\theta(M'_k = m_k \mid M'_1 = m_1, \dots, M'_{k-1} = m_{k-1}) &= E \left[\frac{S_{m_k}}{\sum_{m \neq \{m_1, \dots, m_{k-1}\}} S_m} \right] \\ &= E \left[\frac{T_k}{\sum_{l=k}^n T_l} \right] = \frac{1}{n-k+1}. \end{aligned}$$

The Bose-Einstein distribution of M'_1, \dots, M'_k results from Bayes formula. \square

• *The joint distribution of the size-biased permutation.* The SBP of S_n is L_n with $L_n \stackrel{d}{\sim} SBD_n(\theta)$ and $S_n \stackrel{d}{\sim} D_n(\theta)$. First, we have

$$(2.9) \quad (L_1, \dots, L_n) = (S_{M'_1}, \dots, S_{M'_n}),$$

and consequently

$$(2.10) \quad P(L_1 = S_{m_1}, \dots, L_n = S_{m_n} \mid S_n) = \prod_{k=1}^{n-1} \frac{S_{m_k}}{1 - \sum_{l=1}^k S_{m_l}} S_{m_n}.$$

Consider now the joint moment function of the random size-biased permutation $L_n = (L_1, \dots, L_n)$. In Barrera *et al.* (2005), using the RAM representation of L_n , the following result was proven.

THEOREM 2.1. *The joint moment function of the SBP $L_n = (L_1, \dots, L_n) \stackrel{d}{\sim} SBD_n(\theta)$ reads*

$$\begin{aligned} (2.11) \quad E \left[\prod_{k=1}^n L_k^{q_k} \right] &= \sum_{\{m_1 \neq \dots \neq m_n\}} E \left[\prod_{k=1}^{n-1} \frac{S_{m_k}^{q_k+1}}{1 - \sum_{l=1}^k S_{m_l}} S_{m_n}^{q_n+1} \right] \\ &= \prod_{k=1}^{n-1} \left\{ \frac{\Gamma(1 + (n-k+1)\theta)}{\Gamma(1 + \theta)\Gamma((n-k)\theta)} \right. \\ &\quad \left. \times \frac{\Gamma(1 + \theta + q_k)\Gamma((n-k)\theta + q_{k+1} + \dots + q_n)}{\Gamma(1 + (n-k+1)\theta + q_k + \dots + q_n)} \right\}. \end{aligned}$$

• *One-dimensional marginals.* From Theorem 2.1, we get the one-dimensional law of the L_k s, $k = 1, \dots, n$. Furthermore, one may check that the L_k s are arranged in stochastically decreasing order (denoted by \succeq). More precisely (see Barrera *et al.* (2005))

COROLLARY 2.1. (i) *The law of L_k , for $k = 1, \dots, n$, is characterized by*

$$\begin{aligned}
 (2.12) \quad \mathbf{E}[L_k^q] &= \prod_{i=1}^{k-1} \mathbf{E}[\overline{V}_i^q] \mathbf{E}[V_k^q] \\
 &= \prod_{i=1}^{k-1} \frac{\Gamma((n-i)\theta + q)\Gamma((n-i+1)\theta + 1)}{\Gamma((n-i)\theta)\Gamma((n-i+1)\theta + 1 + q)} \\
 &\quad \times \frac{\Gamma(1 + \theta + q)\Gamma(1 + (n-k+1)\theta)}{\Gamma(1 + \theta)\Gamma(1 + (n-k+1)\theta + q)}.
 \end{aligned}$$

(ii) *Let $B_{(n-k+1)\theta,1} \stackrel{d}{\sim} \text{beta}((n-k+1)\theta, 1)$. Then,*

$$(2.13) \quad L_k \stackrel{d}{=} B_{(n-k+1)\theta,1} \cdot L_{k-1}, \quad k = 2, \dots, n,$$

where pairs $B_{(n-k+1)\theta,1}$ and L_{k-1} are mutually independent for $k = 2, \dots, n$.

(iii) $L_1 \succeq \dots \succeq L_k \succeq \dots \succeq L_n$.

The Kingman limit

Consider the limit $n \uparrow \infty$, $\theta \downarrow 0$ while $n\theta = \gamma > 0$. Such an asymptotic was first considered by Kingman (1975); we shall “star” the results (as in $\stackrel{d}{*}$) when referring to such an asymptotic. As noted by Kingman, $\mathcal{S}_n \stackrel{d}{\sim} D_n(\theta)$ itself has no non-degenerate limit.

When $k = o(n)$, recalling $V_k \stackrel{d}{\sim} \text{beta}(1 + \theta, (n-k)\theta)$, we have $V_k \stackrel{d}{*} V_k^* \stackrel{d}{\sim} \text{beta}(1, \gamma)$ and the $SBD_n(\theta)$ distribution converges weakly from equations (2.3), (2.4) to a Griffiths-Engen-McCloskey or $GEM(\gamma)$ distribution.

Namely, $(L_1, \dots, L_n) \stackrel{d}{*} (L_1^*, \dots, L_k^*, \dots) =: \mathbf{L}^*$ where

$$(2.14) \quad L_k^* = \prod_{i=1}^{k-1} \overline{V}_i^* V_k^*, \quad k \geq 1.$$

Here $(V_k^*, k \geq 1)$ are iid with common law $V_1^* \stackrel{d}{\sim} \text{beta}(1, \gamma)$ and $\overline{V}_1^* := 1 - V_1^* \stackrel{d}{\sim} \text{beta}(\gamma, 1)$. Note that $L_1^* \succeq \dots \succeq L_k^* \succeq \dots$, and that \mathbf{L}^* is invariant under size-biased permutation. In the Kingman limit, $(S_{(m)}, m = 1, \dots, n)$ converges in law to a Poisson-Dirichlet distribution $(L_{(k)}^*, k \geq 1) \stackrel{d}{\sim} PD(\gamma)$ with $L_{(1)}^* > \dots > L_{(k)}^* > \dots$. The size-biased permutation of $(L_{(k)}^*, k \geq 1)$ is $(L_k^*, k \geq 1) \stackrel{d}{\sim} GEM(\gamma)$ (see Kingman (1993), Chapter 9).

The model (2.14) generates a random countable partition of the unit interval, with many fundamental invariance properties (for a review of these results and applications to Computer Science, Combinatorial Structures, Physics, Biology ..., see Tavaré and

Ewens (1997) and the references therein for example; this model and related ones are also fundamental in Probability Theory; see Pitman (1996, 1999, 2002) and Pitman and Yor (1997).

3. Dirichlet partitions: sampling formulae for unordered and ordered sequences

Ewens’ sampling formula (ESF) gives the distribution of alleles (different types of genes) in a sample with size k from the Poisson-Dirichlet process $PD(\gamma)$. Alternatively, it can be described in terms of sequential sampling of animals from a countable collection of distinguishable species drawn from $GEM(\gamma)$. It provides the probability of the partition of a sample of k selectively equivalent genes into a number of alleles as population size becomes indefinitely large. Depending on whether the order of appearance of sequentially sampled species matters or not, we are led to the first ESF for unordered sequences or to the Donnelly-Tavaré-Griffiths (DTG) sampling formula for ordered sequences. A third way to describe the sample is to record the number of species in the k -sample with exactly i representatives, $i = 0, \dots, k$. When doing this while assuming the species have random frequencies following $GEM(\gamma)$ distribution, we are led to a second Ewens Sampling Formula. We recall here the exact expressions of both first, second Ewens and DTG sampling formulae, when sampling is from finite Dirichlet random partitions. These sampling formulae give both ESF and DTG formulae from $GEM(\gamma)$ when passing to the Kingman limit (see Sibuya and Yamato (1995) for further results). Most of the results (and their proofs) presented therein can be found in Huillet (2005).

3.1 The first Ewens sampling formula for Dirichlet partitions

We first consider a sampling formula from Dirichlet partitions for which the order in which the consecutive fragments are being discovered in the sampling process is irrelevant.

Let \mathbf{S}_n be the above Dirichlet random partition with parameter $\theta > 0$. Let $k > 1$ and (U_1, \dots, U_k) be k iid uniform random throws on $[0, 1]$. Let then (M_1, \dots, M_k) be the (conditionally iid) corresponding fragments numbers (or animals’ species), with common conditional and unconditional distributions

$$(3.1) \quad P(M = m \mid \mathbf{S}_n) = S_m, \quad m \in \{1, \dots, n\}$$

$$(3.2) \quad P_\theta(M = m) := E[P(M = m \mid \mathbf{S}_n)] = ES_m = \frac{1}{n}.$$

Let $\mathcal{B}_{n,k}(m) = \sum_{l=1}^k \mathbf{I}(M_l = m)$ count the random number of occurrences of fragment m in the k -sample and $P_{n,k} := \sum_{m=1}^n \mathbf{I}(\mathcal{B}_{n,k}(m) > 0)$ count the number of distinct fragments which have been visited in the k -sampling process. We let $(\theta)_k := \theta(\theta + 1)(\theta + k - 1)$. In Huillet (2005) we obtained

THEOREM 3.1. (i) We have

$$(3.3) \quad P_\theta(\mathcal{B}_{n,k}(1) = k_1, \dots, \mathcal{B}_{n,k}(p) = k_p; P_{n,k} = p) \\ = \binom{n}{p} \frac{k!}{\prod_{q=1}^p k_q!} \frac{1}{(n\theta)_k} \prod_{q=1}^p (\theta)_{k_q}.$$

(ii) *With*

$$B_{k,p}(x_1, x_2, \dots) := \sum_{\substack{a_i \geq 0: \sum_{i=1}^k ia_i = k; \\ \sum_{i=1}^k a_i = p}} \frac{k!}{\prod_{i=1}^k i^{a_i} a_i!} \prod_{i=1}^k x_i^{a_i}$$

the Bell polynomials,

$$(3.4) \quad P_\theta(P_{n,k} = p) = \frac{n!}{(n-p)!} \frac{1}{(n\theta)_k} B_{k,p}((\theta)_1, (\theta)_2, \dots).$$

(iii) *We have*

$$(3.5) \quad \begin{aligned} P_\theta(\mathcal{B}_{n,k}(1) = k_1, \dots, \mathcal{B}_{n,k}(p) = k_p \mid P_{n,k} = p) \\ = \frac{k!}{p!} \frac{1}{B_{k,p}((\theta)_1, (\theta)_2, \dots)} \prod_{q=1}^p \frac{(\theta)_{k_q}}{k_q!}. \end{aligned}$$

Note in particular that if $p = k = 2$, with $k_1 = k_2 = 1$, equation (3.3) gives

$$P_\theta(\mathcal{B}_{n,2}(1) = 1, \mathcal{B}_{n,2}(2) = 1; P_{n,2} = 2) = \frac{(n-1)\theta}{n\theta + 1}.$$

This is the probability that the first 2 random throws will visit any 2 distinct fragments. The complementary probability that it does not is thus $1 - \frac{(n-1)\theta}{n\theta + 1} = \frac{\theta + 1}{n\theta + 1}$.

Remark. (the law of succession) We would like to briefly recall a related question raised in Donnelly (1986) and Ewens (1996), concerning the law of succession.

(i) Let the “ M_{k+1} is new” denote the event that M_{k+1} is none of the previously visited fragments. One can prove (see Huillet (2005))

$$(3.6) \quad \begin{aligned} P_\theta(M_{k+1} \text{ is new} \mid \mathcal{B}_{n,k}(1) = k_1, \dots, \mathcal{B}_{n,k}(p) = k_p; P_{n,k} = p) \\ = \frac{(n-p)\theta}{n\theta + k}, \end{aligned}$$

which is independent of cell occupancies k_1, \dots, k_p but depends on the number p of distinct fragments already visited by the k -sample.

(ii) Similarly, let the event “ $M_{k+1} \in$ species seen k_r times” denote the fact that the $(k + 1)$ -th sample is one from the previously encountered fragment already visited k_r times. We easily get

$$(3.7) \quad \begin{aligned} P_\theta(M_{k+1} \in \text{species seen } k_r \text{ times} \mid \mathcal{B}_{n,k}(1) = k_1, \dots, \mathcal{B}_{n,k}(p) = k_p; P_{n,k} = p) \\ = \frac{\theta + k_r}{n\theta + k}, \end{aligned}$$

which is independent of occupancy numbers $k_q, q \in \{1, \dots, p\} \setminus \{r\}$ and also of the number p of distinct observations.

Remark. (number of distinct observations) From equations (3.6) and (3.7), we also have the transition probabilities

$$P_\theta(P_{n,k+1} = p + 1 \mid P_{n,k} = p) = \frac{(n - p)\theta}{n\theta + k}$$

$$P_\theta(P_{n,k+1} = p \mid P_{n,k} = p) = \frac{\sum_{r=1}^p (\theta + k_r)}{n\theta + k} = \frac{p\theta + k}{n\theta + k}.$$

Next,

$$P_\theta(P_{n,k+1} = p) = \frac{(n - p + 1)\theta}{n\theta + k} P_\theta(P_{n,k} = p - 1) + \frac{p\theta + k}{n\theta + k} P_\theta(P_{n,k} = p).$$

Using equation (3.4), we obtain the following triangular recurrence for Bell polynomials $B_{k,p}((\theta)_1, (\theta)_2, \dots)$,

$$B_{k+1,p}((\theta)_1, (\theta)_2, \dots) = \theta B_{k,p-1}((\theta)_1, (\theta)_2, \dots) + (p\theta + k) B_{k,p}((\theta)_1, (\theta)_2, \dots).$$

These should be considered with boundary conditions

$$B_{k,0}((\theta)_1, (\theta)_2, \dots) = B_{0,p}((\theta)_1, (\theta)_2, \dots) = 0,$$

except for $B_{0,0}((\theta)_1, (\theta)_2, \dots) := 1$.

This leads in particular to $B_{k,1}((\theta)_1, (\theta)_2, \dots) = (\theta)_k$, $k \geq 1$ and to

$$P_\theta(P_{n,k} = 1) = \frac{n(\theta)_k}{(n\theta)_k}.$$

The Kingman limit

Consider the situation where $n \uparrow \infty$, $\theta \downarrow 0$ while $n\theta = \gamma > 0$ in equations (3.3), (3.4) of Theorem 3.1. Proceeding in this way, we recover the first Ewens sampling formula (1972):

COROLLARY 3.1. (i) *In the Kingman limit,*

$$P_\theta(\mathcal{B}_{n,k}(1) = k_1, \dots, \mathcal{B}_{n,k}(p) = k_p; P_{n,k} = p)$$

converges to

$$(3.8) \quad P_\gamma^*(\mathcal{B}_k(1) = k_1, \dots, \mathcal{B}_k(p) = k_p; P_k = p) = \frac{k!}{p!} \frac{\gamma^p}{(\gamma)_k \prod_{q=1}^p k_q}.$$

(ii) *With $s_{k,p}$ the absolute value of the first kind Stirling numbers, we get*

$$(3.9) \quad P_\gamma^*(P_k = p) = \frac{\gamma^p s_{k,p}}{(\gamma)_k}, \quad p = 1, \dots, k$$

and

$$(3.10) \quad P_\gamma^*(\mathcal{B}_k(1) = k_1, \dots, \mathcal{B}_k(p) = k_p \mid P_k = p) = \frac{k!}{p!} \frac{1}{s_{k,p} \prod_{q=1}^p k_q}.$$

Remark. (the law of succession) In the Kingman limit, the probabilities displayed in examples (3.6) and (3.7) converge respectively to

$$(3.11) \quad \frac{\gamma}{\gamma + k} \quad \text{and} \quad \frac{k_r}{\gamma + k}.$$

3.2 *The second Ewens formula for Dirichlet populations*

Let now $\mathcal{A}_{n,k}(i), i \in \{0, \dots, k\}$ count the number of fragments in the k -sample with i representatives, that is

$$(3.12) \quad \mathcal{A}_{n,k}(i) = \# \{m \in \{1, \dots, n\} : \mathcal{B}_{n,k}(m) = i\} = \sum_{m=1}^n \mathbf{I}(\mathcal{B}_{n,k}(m) = i).$$

Then $\sum_{i=0}^k \mathcal{A}_{n,k}(i) = n, \sum_{i=1}^k \mathcal{A}_{n,k}(i) = p$ is the number of fragments visited by the k -sample and $\mathcal{A}_{n,k}(0)$ the number of unvisited ones. Note that $\sum_{i=1}^k i \mathcal{A}_{n,k}(i) = k$ is the sample size.

The vector $(\mathcal{A}_{n,k}(1), \dots, \mathcal{A}_{n,k}(k))$ is called the fragments vector count or the species vector count in biology, see Ewens (1990). In this case, see Huillet (2005), we have

THEOREM 3.2. (i) *For any $a_i \geq 0, i = 1, \dots, k$ satisfying $\sum_{i=1}^k i a_i = k$ and $\sum_{i=1}^k a_i = p$, we have*

$$(3.13) \quad \begin{aligned} P_\theta(\mathcal{A}_{n,k}(1) = a_1, \dots, \mathcal{A}_{n,k}(k) = a_k; P_{n,k} = p) \\ = \frac{n!}{(n-p)!} \frac{k!}{\prod_{i=1}^k i^{a_i} a_i!} \frac{1}{(n\theta)_k} \prod_{i=1}^k (\theta)_i^{a_i}. \end{aligned}$$

(ii) *With $B_{k,p}(x_1, x_2, \dots)$, the Bell polynomials, we have*

$$(3.14) \quad P_\theta(P_{n,k} = p) = \frac{n!}{(n-p)!} \frac{\Gamma(n\theta)}{\Gamma(n\theta + k)} B_{k,p}((\theta)_1, (\theta)_2, \dots).$$

(iii) *We have*

$$(3.15) \quad \begin{aligned} P_\theta(\mathcal{A}_{n,k}(1) = a_1, \dots, \mathcal{A}_{n,k}(k) = a_k \mid P_{n,k} = p) \\ = \frac{k!}{B_{k,p}((\theta)_1, (\theta)_2, \dots)} \prod_{i=1}^k \frac{(\theta)_i^{a_i}}{i^{a_i} a_i!}. \end{aligned}$$

The Kingman limit

Consider the limit $n \uparrow \infty, \theta \downarrow 0$ while $n\theta = \gamma > 0$. We recover the celebrated Ewens Sampling Formula (1972).

COROLLARY 3.2. (i) *In the Kingman limit, the probability displayed in (3.13) converges to*

$$(3.16) \quad P_\gamma^*(\mathcal{A}_k(1) = a_1, \dots, \mathcal{A}_k(k) = a_k; P_k = p) = \frac{k! \gamma^p}{(\gamma)_k \prod_{i=1}^k i^{a_i} a_i!}.$$

(ii) *With $s_{k,p}$ the absolute value of the first kind Stirling numbers,*

$$(3.17) \quad P_\gamma^*(P_k = p) = \frac{\gamma^p s_{k,p}}{(\gamma)_k}, \quad p = 1, \dots, k$$

and

$$(3.18) \quad P_\gamma^*(\mathcal{A}_k(1) = a_1, \dots, \mathcal{A}_k(p) = a_k \mid P_k = p) = \frac{k!}{s_{k,p} \prod_{i=1}^k i^{a_i} a_i!}.$$

3.3 *Donnelly-Tavaré-Griffiths sampling formula for the Dirichlet partition*

We now consider sampling formulae from Dirichlet partitions for which the order in which the consecutive fragments are being discovered in the sampling process matters.

Consider a k -sample and let $m_1 \neq m_2 \neq \dots \neq m_p$ denote the ordered number of the first, second, \dots , the p -th distinct animals sampled from S_n when only $P_{n,k} = p$ distinct fragments were visited. Let $C_{n,k}(q)$, $q = 1, \dots, p$ be the number of animals of the q -th species to appear. Using Theorem 2.1, we can prove (see Huillet (2005))

THEOREM 3.3. *For any $c_q \geq 1$, $q = 1, \dots, p$ satisfying $\sum_1^p c_q = k$ and any $p = 1, \dots, n \wedge k$,*

$$\begin{aligned}
 (3.19) \quad & \mathbf{P}_\theta(\mathcal{C}_{n,k}(1) = c_1, \dots, \mathcal{C}_{n,k}(p) = c_p; P_{n,k} = p) \\
 &= \frac{(k-1)!}{\prod_{q=1}^{p-1} (k - \sum_1^q c_i)} \frac{\Gamma(1 + (n-p+1)\theta)\Gamma(\theta + c_p)}{\Gamma(1 + \theta)\Gamma((n-p+1)\theta + c_p)\Gamma(c_p)} \\
 & \times \prod_{q=1}^{p-1} \left\{ \frac{\Gamma(1 + (n-q+1)\theta)\Gamma(\theta + c_q)}{\Gamma(1 + \theta)\Gamma((n-q)\theta)\Gamma(c_q)} \frac{\Gamma((n-q)\theta + c_{q+1} + \dots + c_p)}{\Gamma((n-q+1)\theta + c_q + \dots + c_p)} \right\}.
 \end{aligned}$$

Remark. (the law of succession) (i) Consider equation (3.19) and, with $m_r \in \{m_1, \dots, m_p\}$, let

$$\mathbf{P}(M_{k+1} = m_r \mid \mathcal{C}_{n,k}(1) = c_1, \dots, \mathcal{C}_{n,k}(p) = c_p; P_{n,k} = p)$$

be the conditional probability that the $(k+1)$ -th sample is one from the previously encountered species already visited c_r times. One can show, see Huillet (2005), that, as for the ESF,

$$\begin{aligned}
 (3.20) \quad & \mathbf{P}_\theta(M_{k+1} \in \text{species seen } c_r \text{ times} \mid \mathcal{C}_{n,k}(1) = c_1, \dots, \mathcal{C}_{n,k}(p) = c_p; P_{n,k} = p) \\
 &= \frac{\theta + c_r}{n\theta + k},
 \end{aligned}$$

which is again independent of c_q , $q \in \{1, \dots, p\} \setminus \{r\}$ and also of p .

(ii) Summing over $r = 1, \dots, p$, the conditional probability that $M_{k+1} \in \{\text{any one of the species previously seen}\}$ is thus $\sum_{r=1}^p \frac{\theta + c_r}{n\theta + k} = \frac{p\theta + k}{n\theta + k}$. Taking its complement to 1, we obtain

$$\begin{aligned}
 (3.21) \quad & \mathbf{P}_\theta(M_{k+1} \text{ is new} \mid \mathcal{C}_{n,k}(1) = c_1, \dots, \mathcal{C}_{n,k}(p) = c_p; P_{n,k} = p) \\
 &= \frac{(n-p)\theta}{n\theta + k},
 \end{aligned}$$

which is independent of cell occupancies c_1, \dots, c_p but depends on the number p of distinct fragments already visited by the k -sample.

The Kingman limit

Passing to the Kingman limit in Theorem 3.3 gives the celebrated Donnelly-Tavaré-Griffiths sampling formula given in Donnelly and Tavaré (1986), p. 10 (see Huillet (2005) for details).

COROLLARY 3.3. (i) In the Kingman limit, the probability (3.19) converges to

$$(3.22) \quad \mathbf{P}_\gamma^*(C_k(1) = c_1, \dots, C_k(p) = c_p; P_k = p) = \frac{k! \gamma^p}{(\gamma)_k \prod_{q=1}^p (c_q + \dots + c_p)}.$$

(ii) With $s_{k,p}$ the absolute value of the first kind Stirling numbers,

$$(3.23) \quad \mathbf{P}_\gamma^*(P_k = p) = \frac{\gamma^p s_{k,p}}{(\gamma)_k}, \quad p = 1, \dots, k$$

and

$$(3.24) \quad \mathbf{P}_\gamma^*(C_k(1) = c_1, \dots, C_k(p) = c_p \mid P_k = p) = \frac{k!}{s_{k,p} \prod_{q=1}^p (c_q + \dots + c_p)}.$$

4. Intervals between new sampled species from Dirichlet partition

In this section, we shall consider the following sampling problem from Dirichlet partition. From the knowledge of the order in which new fragments are being discovered, what can be said about the random number of samples separating the discovery of consecutive new fragments until exhaustion of the list? We shall first consider the case of a fixed sample size k and then the case of an unlimited size.

Assume there are $P_{n,k} = p \leq n \wedge k := k_n$ distinct species visited by the k -sample drawn from $\mathcal{S}_n \stackrel{d}{\sim} D_n(\theta)$, with $k \geq 2$. Let $m_1 \neq \dots \neq m_p$ be their fragments numbers in order of their appearance. Let $\mathcal{D}_{n,k}(r)$, $r = 1, \dots, p-1$ be the sample size between the discovery of the r -th and the $(r+1)$ -th new fragments. Let $\chi_1 = 1$ and $\chi_l = \prod_{k=1}^{l-1} \mathbf{I}(M_l \neq M_k)$, $l = 2, \dots, k$. This binary sequence takes the value 1 each time a new species is visited by the k -sample; otherwise its value is 0. Let $d_r \geq 1$, $r = 1, \dots, p-1$, $d_p := k - \sum_{r=1}^{p-1} d_r \geq 1$ and $\bar{d}_r = \sum_{q=1}^r d_q$, $r = 1, \dots, p$, $\bar{d}_0 := 0$. The following two events are identical

$$“\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(p-1) = d_{p-1}; P_{n,k} = p”$$

and

$$“\chi_1 = 1, \chi_2 = \dots = \chi_{d_1} = 0, \chi_{d_1+1} = 1, \chi_{d_1+2} = \dots \\ = \chi_{\bar{d}_2} = 0, \dots, \chi_{\bar{d}_{p-1}+1} = 1, \chi_{\bar{d}_{p-1}+2} = \dots = \chi_{\bar{d}_p} = 0”.$$

Next, from the above law of succession displayed in (3.21), with $b_1, \dots, b_l \in \{0, 1\}^l$, we get

$$\mathbf{P}_\theta(\chi_{l+1} = 1 \mid \chi_1 = b_1, \dots, \chi_l = b_l) = \frac{(n - \sum_{q=1}^l b_q) \theta}{n\theta + l}$$

$$\mathbf{P}_\theta(\chi_{l+1} = 0 \mid \chi_1 = b_1, \dots, \chi_l = b_l) = \frac{l + \theta \cdot \sum_{q=1}^l b_q}{n\theta + l}.$$

Using this, we obtain

THEOREM 4.1. Let $p \leq k_n$. With $d_r \geq 1$, $r = 1, \dots, p-1$, $d_p := k - \sum_{r=1}^{p-1} d_r \geq 1$ and $\bar{d}_r = \sum_{q=1}^r d_q$, $r = 1, \dots, p$, $\bar{d}_0 := 0$, we get

(i)

$$(4.1) \quad \begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(p-1) = d_{p-1}; P_{n,k} = p) \\ = \frac{\theta^p n!}{(n-p)!(n\theta)_k} \prod_{r=1}^p (1 + r\theta + \bar{d}_{r-1})_{d_{r-1}}. \end{aligned}$$

(ii)

$$(4.2) \quad \begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(p-1) = d_{p-1} \mid P_{n,k} = p) \\ = \frac{\theta^p}{B_{k,p}((\theta)_1, (\theta)_2, \dots)} \prod_{r=1}^p (1 + r\theta + \bar{d}_{r-1})_{d_{r-1}}. \end{aligned}$$

PROOF. (i) We have:

$$\begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(p-1) = d_{p-1}; P_{n,k} = p) \\ = \prod_{r=0}^{p-1} \frac{(n-r)\theta}{\bar{d}_r + n\theta} \cdot \prod_{r=1}^p \prod_{l=\bar{d}_{r-1}+1}^{\bar{d}_r-1} \frac{l+r\theta}{l+n\theta} \\ = \frac{\prod_{r=0}^{p-1} [(n-r)\theta]}{\prod_{l=0}^{k-1} (l+n\theta)} \prod_{r=1}^p \prod_{l=\bar{d}_{r-1}+1}^{\bar{d}_r-1} (l+r\theta) \\ = \frac{n!}{(n-p)!} \frac{\theta^p}{(n\theta)_k} \prod_{r=1}^p (1 + r\theta + \bar{d}_{r-1})_{d_{r-1}}, \end{aligned}$$

(ii) follows from the expression (3.14) of $P_\theta(P_{n,k} = p)$. \square

Remark. (Waring distribution) Let us recall here some elementary facts on Waring distributions in the generalized hypergeometric of type B3 class. A discrete random variable $W \in \{0, 1, \dots\}$ has Waring distribution $W(b, a)$ with parameters $b > a > 0$ if $P(W = x) = (b-a) \frac{(a)_x}{(b)_{x+1}}$. We have $E(W) = a/(b-a-1)$ if $b-a > 1$, $= +\infty$ if not. Regrouping the events " $W \geq n$ ", a random variable $bW \in \{0, 1, \dots, n\}$ has bounded Waring distribution $bW(n, b, a)$ if $P(bW = x) = (b-a) \frac{(a)_x}{(b)_{x+1}}$, if $x \in \{0, 1, \dots, n-1\}$, $P(bW = n) = \frac{(a)_n}{(b)_n}$ if $bW = n$. We obtain $E(bW) = \frac{a}{b-a-1} (1 - \frac{(a+1)_n}{(b)_n})$ if $b \neq a+1$ and $E(bW) = (b-1)[\psi'(b+n) - \psi'(b)]$ if $b = a+1$, where $\psi'(b) := \Gamma'(b)/\Gamma(b)$.

One may check that if $W \stackrel{d}{\sim} W(b, a)$, then $W \stackrel{d}{\sim} \text{geom}(\text{beta}(a, b-a))$ where $G \stackrel{d}{\sim} \text{geom}(p)$ is a geometrically distributed random variable with law $P(G = x) = p^x(1-p)$, $x \geq 0$. Here, p is the random success probability, with $p \stackrel{d}{\sim} \text{beta}(a, b-a)$, independent of G . From this, we get

$$\begin{aligned} P(W = x) &= \frac{\Gamma(b)}{\Gamma(a)\Gamma(b-a)} \int_0^1 p^{x+a-1} (1-p)^{b-a} dp \\ &= \frac{\Gamma(b)}{\Gamma(a)\Gamma(b-a)} \frac{\Gamma(a+x)\Gamma(b-a+1)}{\Gamma(b+x+1)} = (b-a) \frac{(a)_x}{(b)_{x+1}}. \end{aligned}$$

Furthermore, $P(W \geq n) = \frac{\Gamma(b)}{\Gamma(a)\Gamma(b-a)} \int_0^1 p^{n+a-1}(1-p)^{b-a-1} dp = \frac{(a)_n}{(b)_n}$. Let bG have bounded geometric distribution such that $P(bG = x) = p^x(1-p)$, if $x \in \{0, 1, \dots, n-1\}$, $P(bG = n) = p^n$ if $bG = n$. Using $E(G) = p/(1-p)$ and $E(bG) = \frac{p}{1-p} - \frac{p^{n+1}}{1-p}$, the claims on $E(W)$ and $E(bW)$ follow from randomization of p .

Now, we are ready for presenting subsequent results. As a consequence of Theorem 4.1, we have

THEOREM 4.2. *Let $1 \leq r < k_n$.*

(i) *With $d_1, \dots, d_r \geq 1, \bar{d}_r < k$, we have*

$$(4.3) \quad \begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(r) = d_r) \\ = \frac{(n-1)!}{(n-r-1)!} \frac{\theta^r}{(1+n\theta)_{\bar{d}_r}} \prod_{q=1}^r (1+q\theta + \bar{d}_{q-1})_{d_q-1}. \end{aligned}$$

In particular, if $r = 1$

$$(4.4) \quad \begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1) &= \theta(n-1) \frac{(1+\theta)_{d_1-1}}{(1+n\theta)_{d_1}}, \quad d_1 \in \{1, \dots, k_n-1\} \\ &= \frac{(1+\theta)_{d_1-1}}{(1+n\theta)_{d_1-1}}, \quad \text{if } d_1 = k_n, \end{aligned}$$

and $\mathcal{D}_{n,k}(1) - 1$ has bounded Waring distribution $bW(k_n-1, 1+n\theta, 1+\theta)$. Consequently,

$$(4.5) \quad E_\theta[\mathcal{D}_{n,k}(1)] = 1 + \frac{1+\theta}{(n-1)\theta-1} \left[1 - \frac{(2+\theta)_{k_n-1}}{(1+n\theta)_{k_n-1}} \right].$$

(ii) *Let $1 \leq r < k_n$. With $d_1, \dots, d_r \geq 1, \bar{d}_{r-1} =: d \geq r-1, \bar{d}_r = d + d_r < k$, we have*

$$(4.6) \quad \begin{aligned} P_\theta(\mathcal{D}_{n,k}(r) = d_r \mid \mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(r-1) = d_{r-1}) \\ = \theta(n-r) \frac{(1+r\theta+d)_{d_r-1}}{(1+n\theta+d)_{d_r}}, \quad \text{if } d_r \in \{1, \dots, k_n-d-1\} \\ = \frac{(1+\theta+d)_{d_r-1}}{(1+n\theta+d)_{d_r-1}}, \quad \text{if } d_r = k_n-d, \end{aligned}$$

showing that, given $\bar{\mathcal{D}}_{n,k}(r-1) := \sum_{q=1}^{r-1} \mathcal{D}_{n,k}(q) = d, \mathcal{D}_{n,k}(r) - 1$ has bounded Waring distribution $bW(k_n-d-1, 1+n\theta+d, 1+r\theta+d)$. Note from the preceding remark that

$$(4.7) \quad E_\theta[\bar{\mathcal{D}}_{n,k}(r-1)] = 1 + \frac{1+r\theta+d}{(n-r)\theta-1} \left[1 - \frac{(2+r\theta+d)_{k_n-d-1}}{(1+n\theta+d)_{k_n-d-1}} \right].$$

(iii) *For $r = 1, \dots, k_n-1$, the cumulative random variable $\bar{\mathcal{D}}_{n,k}(r) := \sum_{q=1}^r \mathcal{D}_{n,k}(q)$ has distribution*

$$(4.8) \quad \begin{aligned} P_\theta(\bar{\mathcal{D}}_{n,k}(r) = d) &= \frac{(n-1)! B_{d,r}((\theta)_1, (\theta)_2, \dots)}{(n-r-1)!(1+n\theta)_d}, \quad d = r, \dots, k_n-1 \\ &= \frac{n!}{(n-r)!} \frac{B_{d,r}((\theta)_1, (\theta)_2, \dots)}{(n\theta)_d}, \quad d = k_n, \end{aligned}$$

where $B_{d,r}((\theta)_1, (\theta)_2, \dots)$ are Bell polynomials.

PROOF. (i) The following two events coincide

$$“\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(r) = d_r”$$

and “ $\chi_1 = 1, \chi_2 = \dots = \chi_{d_1} = 0, \chi_{d_1+1} = 1, \chi_{d_1+2} = \dots = \chi_{\bar{d}_2} = 0, \dots, \chi_{\bar{d}_{r-1}+1} = 1, \chi_{\bar{d}_{r-1}+2} = \dots = \chi_{\bar{d}_r} = 0, \chi_{\bar{d}_r+1} = 1$ ”. Proceeding as for the proof of Theorem 4.3 gives (i).

In particular,

$$\begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) = d_1) &= \frac{\theta^2 n!}{(n-2)!} \frac{(1+\theta)_{d_1-1}}{(n\theta)_{d_1+1}} \\ &= \frac{(n-1)(\theta)_{d_1}}{(1+n\theta)_{d_1}} = \theta(n-1) \frac{(1+\theta)_{d_1-1}}{(1+n\theta)_{d_1}}, \end{aligned}$$

so that, with $d = 0, \dots, k_n - 1$

$$\begin{aligned} P_\theta(\mathcal{D}_{n,k}(1) - 1 = d) &= \theta(n-1) \frac{(1+\theta)_d}{(1+n\theta)_{d+1}}, \quad d = 0, \dots, k_n - 2 \\ &= \frac{(1+\theta)_d}{(1+n\theta)_d}, \quad \text{if } d = k_n - 1 \end{aligned}$$

which is a bounded Waring distribution $bW(k_n - 1, 1 + n\theta, 1 + \theta)$. The expected value of $\mathcal{D}_{n,k}(1)$ displayed in (4.5) follows from the preceding remark.

(ii) If $d_r \in \{1, \dots, k - d - 1\}$, it follows from the law of $\mathcal{D}_{n,k}(1), \dots, \mathcal{D}_{n,k}(r)$ that

$$\begin{aligned} P_\theta(\mathcal{D}_{n,k}(r) = d_r \mid \mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(r-1) = d_{r-1}) \\ &= \theta(n-r) \frac{(n\theta)_{d+1}}{(n\theta)_{d+d_r+1}} (1+r\theta+d)_{d_r-1} \\ &= \theta(n-r) \frac{(1+r\theta+d)_{d_r-1}}{(1+n\theta+d)_{d_r}}, \end{aligned}$$

where $d := d_1 + \dots + d_{r-1}$. This conditional distribution depends on d_1, \dots, d_{r-1} only through their sum.

From this, one may check that given $\bar{\mathcal{D}}_{n,k}(r-1) = d$, $\mathcal{D}_{n,k}(r) - 1$ has bounded Waring distribution $bW(k_n - d - 1, 1 + n\theta + d, 1 + r\theta + d)$. The expected value of $\bar{\mathcal{D}}_{n,k}(r-1)$ displayed in (4.7) follows from the preceding remark.

(iii) This results from the Ewens sampling formula displayed in Theorem 3.2.

Suppose that $(\mathcal{D}_{n,k}(1), \dots, \mathcal{D}_{n,k}(p-1); P_{n,k})$ has distribution given by (4.1). Let r and \bar{d}_r be integers ≥ 1 such that $r \leq \bar{d}_r < k_n$. Let $a_1, \dots, a_{\bar{d}_r} \geq 0$ be an unordered partition of \bar{d}_r for which $\sum_1^{\bar{d}_r} a_i = r$ and $\sum_1^{\bar{d}_r} i a_i = \bar{d}_r$. Then

$$\begin{aligned} P_\theta(\mathcal{A}_{n,k}^{(r)}(1) = a_1, \dots, \mathcal{A}_{n,k}^{(r)}(\bar{d}_r) = a_{\bar{d}_r}) \\ &= \sum P_\theta(\mathcal{D}_{n,k}(1) = d_1, \dots, \mathcal{D}_{n,k}(r) = d_r) \end{aligned}$$

where summation runs over all distinct ordered partitions d_1, \dots, d_r of \bar{d}_r which give rise to the unordered partition $a_1, \dots, a_{\bar{d}_r}$ of \bar{d}_r . From (4.3), we get

$$\begin{aligned} P_\theta(\mathcal{A}_{n,k}^{(r)}(1) = a_1, \dots, \mathcal{A}_{n,k}^{(r)}(\bar{d}_r) = a_{\bar{d}_r}) &= \\ &= \frac{n!}{(n-r)!} \frac{\bar{d}_r!}{\prod_{i=1}^{\bar{d}_r} i!^{a_i} a_i!} \frac{1}{(n\theta)^{\bar{d}_r}} \prod_{i=1}^{\bar{d}_r} (\theta)_i^{a_i} \quad \text{if } \bar{d}_r = k_n \\ &= \frac{(n-1)!}{(n-r-1)!} \frac{\bar{d}_r!}{\prod_{i=1}^{\bar{d}_r} i!^{a_i} a_i!} \frac{1}{(1+n\theta)^{\bar{d}_r}} \prod_{i=1}^{\bar{d}_r} (\theta)_i^{a_i} \quad \text{if } r \leq \bar{d}_r < k_n \end{aligned}$$

with

$$\sum \frac{\bar{d}_r!}{\prod_{i=1}^{\bar{d}_r} i!^{a_i} a_i!} \prod_{i=1}^{\bar{d}_r} (\theta)_i^{a_i} = B_{\bar{d}_r,r}((\theta)_1, (\theta)_2, \dots).$$

Next, for $d = r, r + 1, \dots, k_n - 1$, we have

$$P_\theta(\bar{\mathcal{D}}_{n,k}(r) = d) = \sum P_\theta(\mathcal{A}_{n,k}^{(r)}(1) = a_1, \dots, \mathcal{A}_{n,k}^{(r)}(d) = a_d)$$

where summation runs over $a_i \geq 0$ satisfying $\sum_1^d a_i = r$ and $\sum_1^d i a_i = d$. \square

Lastly, we shall turn to the related question of computing the law of $\mathcal{D}_n(r)$, the sample size separating consecutive visits to the r -th and the $(r + 1)$ -th new species, when sampling is from $D_n(\theta)$ with unlimited sample size. For this problem, we obtain

THEOREM 4.3. *Let $1 \leq r < n - 1$.*

(i) *The law of $\mathcal{D}_n(r)$ is given by*

$$\begin{aligned} (4.9) \quad P_\theta(\mathcal{D}_n(r) > d) &= \sum_{k_1, \dots, k_r \geq 0: \sum_1^r k_q = d} \binom{d}{k_1, \dots, k_r} \\ &\times \prod_{q=1}^r \left\{ \frac{\Gamma(1 + (n - q + 1)\theta)}{\Gamma(1 + \theta)\Gamma((n - q)\theta)} \right. \\ &\quad \left. \times \frac{\Gamma(1 + \theta + k_q)\Gamma((n - q)\theta + k_{q+1} + \dots + k_r)}{\Gamma(1 + (n - q + 1)\theta + k_q + \dots + k_r)} \right\} \end{aligned}$$

and the following stochastic domination property holds

$$\mathcal{D}_n(r) \succeq \mathcal{D}_n(r - 1).$$

(ii) *With $1 \leq r < n - 1$, the mean value of $\mathcal{D}_n(r)$ reads*

$$\begin{aligned} (4.10) \quad E_\theta(\mathcal{D}_n(r)) &= \prod_{q=1}^r \frac{(n - q + 1)\theta}{(n - q)\theta - 1}, \quad \text{if } (n - r)\theta > 1 \\ &= +\infty, \quad \text{if not.} \end{aligned}$$

PROOF. (i) Given \mathbf{S}_n , we have

$$P(\mathcal{D}_n(r) > d \mid M'_1 = m_1, \dots, M'_r = m_r, \mathbf{S}_n) = \left(\sum_{q=1}^r S_{m_q} \right)^d$$

since, if $\mathcal{D}_n(r) > d$, at least d sample need to fall in the already visited fragments. Averaging over \mathbf{S}_n and summing over all realizations $m_1 \neq \dots \neq m_r$ of M'_1, \dots, M'_r , we get

$$(4.11) \quad P_\theta(\mathcal{D}_n(r) > d) = E \left[\left(\sum_{q=1}^r L_q \right)^d \right].$$

Developing with the help of the multinomial identity and using the joint moment function of L_1, \dots, L_r , as obtained from (2.11), gives the result. The stochastic domination property $\mathcal{D}_n(r) \succeq \mathcal{D}_n(r-1)$ follows from the fact that $(\sum_{q=1}^r L_q)^d \geq (\sum_{q=1}^{r-1} L_q)^d$ which is maintained when taking the expectation.

(ii) Concerning the mean value of $\mathcal{D}_n(r)$, recalling $1 - \sum_{q=1}^r L_q = \prod_{q=1}^r \bar{V}_q$ where $\bar{V}_q, q = 1, \dots, r$ are independent with law $\bar{V}_q \stackrel{d}{\sim} \text{beta}((n-q)\theta, 1+\theta)$ we have from (4.11)

$$\begin{aligned} E_\theta(\mathcal{D}_n(r)) &= \sum_{d \geq 0} P_\theta(\mathcal{D}_n(r) > d) = E_\theta \left(\frac{1}{1 - \sum_{q=1}^r L_q} \right) \\ &= E \left(\frac{1}{\prod_{q=1}^r \bar{V}_q} \right) = \prod_{q=1}^r E \left(\frac{1}{\bar{V}_q} \right). \end{aligned}$$

The expected value of $E_\theta(\frac{1}{\bar{V}_q})$ is finite and equal to $\frac{(n-q+1)\theta}{(n-q)\theta-1}$ if and only if $(n-q)\theta > 1$. We note that when $r \geq n - \theta^{-1}$ is such that $(n-r)\theta \leq 1$, the expected time separating the visit from the r -th to the $(r+1)$ -th new fragment becomes infinitely large. In particular, for values of θ such that $\theta \leq 1/(n-1)$, $E_\theta(\mathcal{D}_n(r)) = \infty$ for each $1 \leq r < n-1$. On the contrary, if $\theta \geq 1$, $E_\theta(\mathcal{D}_n(r)) < \infty$ for each $1 \leq r < n-1$.

Note that, if $(n-r)\theta > 1$, since $\theta > 0$

$$E_\theta(\mathcal{D}_n(r)) = \frac{(n-r+1)\theta}{(n-r)\theta-1} E_\theta(\mathcal{D}_n(r-1)) > E_\theta(\mathcal{D}_n(r-1)).$$

The expected visiting times sequence to consecutive new species, when these exist (when they are finite), is an increasing one, as conventional wisdom suggests. These results are consistent with the ones of Huillet (2003), obtained in the particular case $\theta = 1$, namely

$$E_{\theta=1}(\mathcal{D}_n(r)) = \prod_{q=1}^r \frac{n-q+1}{n-q-1} = \frac{n(n-1)}{(n-r)(n-r-1)},$$

for each $1 \leq r < n-1$. \square

Remark. With $\sum_{q=1}^r k_q = d$, the following term arising from (4.9)

$$\binom{d}{k_1, \dots, k_r} \prod_{q=1}^r \frac{\Gamma(1 + (n-q+1)\theta)\Gamma(1 + \theta + k_q)\Gamma((n-q)\theta + k_{q+1} + \dots + k_r)}{\Gamma(1 + \theta)\Gamma((n-q)\theta)\Gamma(1 + (n-q+1)\theta + k_q + \dots + k_r)}$$

is the joint probability $P_\theta(\mathcal{D}_n(r) > d; \mathcal{B}_{n,d}(1) = k_1, \dots, \mathcal{B}_{n,d}(r) = k_r)$ that $\mathcal{D}_n(r) > d$ and that cell occupancies $\mathcal{B}_{n,d}(q) = k_q, q = 1, \dots, r$ where $\mathcal{B}_{n,d}(q)$ counts the random number of occurrences of the q -th previously visited fragment in a d -sample from $D_n(\theta)$.

The Kingman limit

Consider the situation where $n \uparrow \infty, \theta \downarrow 0$ while $n\theta = \gamma > 0$. In this case, we recover part of the results given by Yamato *et al.* (2001) in a broader (two-parameter) model.

As a corollary to Theorem 4.1, we have

COROLLARY 4.1. *Let $p \leq k$. Assume $d_r \geq 1, r = 1, \dots, p-1, d_p := k - \sum_{r=1}^{p-1} d_r \geq 1$ and $\bar{d}_r = \sum_{q=1}^r d_q, r = 1, \dots, p, \bar{d}_0 := 0$. Then, the weak limits $\mathcal{D}_k(r) := \lim_* \mathcal{D}_{n,k}(r), r = 1, \dots, p-1$ exist and*

$$(i) \quad (4.12) \quad \begin{aligned} P_\gamma^*(\mathcal{D}_k(1) = d_1, \dots, \mathcal{D}_k(p-1) = d_{p-1}; P_k = p) \\ = \frac{\gamma^p}{(\gamma)_k} \prod_{r=1}^p (1 + \bar{d}_{r-1})_{d_{r-1}}. \end{aligned}$$

$$(ii) \quad (4.13) \quad \begin{aligned} P_\gamma^*(\mathcal{D}_k(1) = d_1, \dots, \mathcal{D}_k(p-1) = d_{p-1} \mid P_k = p) \\ = \frac{1}{s_{k,p}} \prod_{r=1}^p (1 + \bar{d}_{r-1})_{d_{r-1}}. \end{aligned}$$

PROOF. Immediate. This is consistent with the result of Yamato *et al.* (2001), Theorem 2, p. 21, putting $\alpha = 0$ in the general two-parameter Pitman class where sampling takes place. This distribution is called *DTGII* by Yamato (1997). \square

As a corollary to Theorem 4.2, we get

COROLLARY 4.2. *Let $1 \leq r < k$.*

$$(i) \quad (4.14) \quad \begin{aligned} \text{With } d_1, \dots, d_r \geq 1, \bar{d}_r < k, \\ P_\gamma^*(\mathcal{D}_k(1) = d_1, \dots, \mathcal{D}_k(r) = d_r) \\ = \frac{\gamma^{r+1}}{(\gamma)_{\bar{d}_r+1}} \prod_{q=1}^r (1 + \bar{d}_{q-1})_{d_{q-1}}. \end{aligned}$$

In particular, if $r = 1$

$$(4.15) \quad \begin{aligned} P_\gamma^*(\mathcal{D}_k(1) = d_1) &= \gamma \frac{(d_1 - 1)!}{(1 + \gamma)_{d_1}}, \quad d_1 \in \{1, \dots, k-1\} \\ &= \frac{(k-1)!}{(1 + \gamma)_{k-1}}, \quad d_1 = k, \end{aligned}$$

showing that $\mathcal{D}_k(1) - 1$ has bounded Waring distribution $bW(k-1, 1 + \gamma, 1)$, with

$$(4.16) \quad E_\gamma^*(\mathcal{D}_k(1)) = 1 + \frac{1}{\gamma - 1} \left(1 - \gamma \frac{k!}{(\gamma)_k} \right).$$

(ii) Let $1 \leq r < k$. With $d_1, \dots, d_r \geq 1, \bar{d}_{r-1} =: d \geq r - 1, \bar{d}_r = d + d_r < k$,

$$\begin{aligned}
 (4.17) \quad \mathbf{P}_\gamma^*(\mathcal{D}_k(r) = d_r \mid \mathcal{D}_k(1) = d_1, \dots, \mathcal{D}_k(r-1) = d_{r-1}) \\
 &= \gamma \frac{(1+d)_{d_r-1}}{(1+\gamma+d)_{d_r}}, \quad \text{if } d_r \in \{1, \dots, k-d-1\} \\
 &= \frac{(1+d)_{k-d-1}}{(1+\gamma+d)_{k-d-1}}, \quad \text{if } d_r = k-d,
 \end{aligned}$$

showing that, given $\bar{\mathcal{D}}_k(r-1) := \sum_{q=1}^{r-1} \mathcal{D}_k(q) = d$, $\mathcal{D}_k(r) - 1$ has bounded Waring distribution $bW(k-d-1, 1+\gamma+d, 1+d)$. In particular,

$$\begin{aligned}
 (4.18) \quad \mathbf{E}_\gamma^* \left(\mathcal{D}_k(r) \mid \sum_{q=1}^{r-1} \mathcal{D}_k(q) = d \right) \\
 = 1 + \frac{1+d}{\gamma-1} \left(1 - \frac{(2+d)_{k-d-1}}{(1+\gamma+d)_{k-d-1}} \right).
 \end{aligned}$$

(iii) For $r = 1, \dots, k-1$, the cumulative variable $\bar{\mathcal{D}}_k(r) := \sum_{q=1}^r \mathcal{D}_k(q)$ has distribution

$$\begin{aligned}
 (4.19) \quad \mathbf{P}_\gamma^*(\bar{\mathcal{D}}_k(r) = d) &= \frac{\gamma^r s_{d,r}}{(1+\gamma)_d}, \quad d = r, r+1, \dots, k-1 \\
 &= \frac{\gamma^r s_{d,r}}{(\gamma)_d}, \quad d = k,
 \end{aligned}$$

where $s_{d,r}$ are the absolute values of first kind Stirling numbers.

PROOF. Immediate. These are again consistent with results of Yamato *et al.* (2001), as from Theorem 3, Propositions 3 and 4, putting $\alpha = 0$ in their formula. \square

Finally, as a corollary to Theorem 4.3, we easily obtain

COROLLARY 4.3. Let $r \geq 1$.

(i) The law of $\mathcal{D}(r) := \lim_* \mathcal{D}_n(r)$ is given by

$$(4.20) \quad \mathbf{P}_\gamma^*(\mathcal{D}(r) > d) = \frac{\gamma^r d!}{(\gamma)_d} \sum_{k_1, \dots, k_r \geq 0: \sum_1^r k_q = d} \frac{1}{\prod_{q=1}^r \{\gamma + k_q + \dots + k_r\}}$$

and the following stochastic domination property holds

$$\mathcal{D}(r) \succeq \mathcal{D}(r-1).$$

(ii) With $1 \leq r$, the mean value of $\mathcal{D}(r)$ reads

$$\begin{aligned}
 (4.21) \quad \mathbf{E}_\gamma^*(\mathcal{D}(r)) &= \left(\frac{\gamma}{\gamma-1} \right)^r, \quad \text{if } \gamma > 1 \\
 &= +\infty, \quad \text{if } 0 < \gamma \leq 1.
 \end{aligned}$$

PROOF. Immediate from Theorem 4.3. \square

Remark. With $\sum_{q=1}^r k_q = d$, $k_q \geq 0$, the following term arising from (4.20)

$$\frac{\gamma^r d!}{(\gamma)_d \prod_{q=1}^r \{\gamma + k_q + \cdots + k_r\}}$$

is the joint probability $P_\gamma^*(\mathcal{D}(r) > d; \mathcal{B}_d(1) = k_1, \dots, \mathcal{B}_d(r) = k_r)$ that $\mathcal{D}(r) > d$ and that cell occupancies are $\mathcal{B}_d(q) = k_q$, $q = 1, \dots, r$, where $\mathcal{B}_d(q)$ counts the random number of occurrences of the q -th previously visited fragment in a d -sample from $GEM(\gamma)$.

REFERENCES

- Barrera, J., Huillet, T. and Paroissin, C. (2005). Size-biased permutation of Dirichlet partitions and search-cost distribution, *Probability in the Engineering & Informational Sciences*, **19**(1), 83–97.
- Donnelly, P. (1986). Partition structures, Pólya urns, the Ewens sampling formula and the age of alleles, *Theoretical Population Biology*, **30**, 271–288.
- Donnelly, P. (1991). The heaps process, libraries and size-biased permutation, *Journal of Applied Probability*, **28**, 321–335.
- Donnelly, P. and Tavaré, S. (1986). The age of alleles and a coalescent, *Advances in Applied Probability*, **18**, 1–19.
- Ewens, W. J. (1972). The sampling theory of selectively neutral alleles, *Theoretical Population Biology*, **3**, 87–112.
- Ewens, W. J. (1990). Population genetics theory—the past and the future, *Mathematical and Statistical Developments of Evolutionary Theory* (ed. S. Lessard), Kluwer, Dordrecht.
- Ewens, W. J. (1996). Some remarks on the law of succession, *Athens Conference on Applied Probability and Time Series Analysis (1995)*, Vol. I, Lecture Notes in Statistics, **114**, 229–244, Springer, New York.
- Huillet, T. (2003). Sampling problems for randomly broken sticks, *Journal of Physics A*, **36**(14), 3947–3960.
- Huillet, T. (2005). Sampling formulae arising from random Dirichlet populations, *Communications in Statistics: Theory and Methods* (to appear).
- Huillet, T. and Martinez, S. (2003). Sampling from finite random partitions, *Methodology and Computing in Applied Probability*, **5**(4), 467–492.
- Kingman, J. F. C. (1975). Random discrete distributions, *Journal of the Royal Statistical Society. Series B*, **37**, 1–22.
- Kingman, J. F. C. (1993). *Poisson Processes*, Clarendon Press, Oxford.
- Pitman, J. (1996). Random discrete distributions invariant under size-biased permutation, *Advances in Applied Probability*, **28**, 525–539.
- Pitman, J. (1999). Coalescents with multiple collisions, *Annals of Probability*, **27**(4), 1870–1902.
- Pitman, J. (2002). Poisson-Dirichlet and GEM invariant distributions for split-and-merge transformation of an interval partition, *Combinatorics, Probability and Computing*, **11**(5), 501–514.
- Pitman, J. and Yor, M. (1997). The two parameter Poisson-Dirichlet distribution derived from a stable subordinator, *Annals of Probability*, **25**, 855–900.
- Sibuya, M. and Yamato, H. (1995). Ordered and unordered random partitions of an integer and the GEM distribution, *Statistics & Probability Letters*, **25**(2), 177–183.
- Tavaré, S. and Ewens, W. J. (1997). Multivariate Ewens distribution, *Discrete Multivariate Distributions* (eds. N. L. Johnson, S. Kotz and N. Balakrishnan), **41**, 232–246, Wiley, New York.
- Yamato, H. (1997). On the Donnelly-Tavaré-Griffiths formula associated with the coalescent, *Communications in Statistics: Theory and Methods*, **26**(3), 589–599.
- Yamato, H., Sibuya, M. and Nomachi, T. (2001). Ordered sample from two-parameter GEM distribution, *Statistics & Probability Letters*, **55**(1), 19–27.