



HAL
open science

Occupancy distributions arising in sampling from Gibbs-Poisson abundance models

Thierry Huillet, Servet Martinez

► **To cite this version:**

Thierry Huillet, Servet Martinez. Occupancy distributions arising in sampling from Gibbs-Poisson abundance models. 2013. hal-00797149v1

HAL Id: hal-00797149

<https://hal.science/hal-00797149v1>

Preprint submitted on 5 Mar 2013 (v1), last revised 27 Sep 2013 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

OCCUPANCY DISTRIBUTIONS ARISING IN SAMPLING FROM GIBBS-POISSON ABUNDANCE MODELS

THIERRY HUILLET¹, SERVET MARTÍNEZ²

ABSTRACT. Estimating the number n of unseen species from a k -sample displaying only $p \leq k$ distinct sampled species has received attention for long. It requires a model of species abundance together with a sampling model. We start with a discrete model of iid stochastic species abundances, each with Gibbs-Poisson distribution. A k -sample drawn from the n -species abundances vector is the one obtained while conditioning it on summing to k . We discuss the sampling formulae (species occupancy distributions, frequency of frequencies) in this context. We then develop some aspects of the estimation of n problem from the size k of the sample and the observed value of $P_{n,k}$, the number of distinct sampled species.

It is shown that it always makes sense to study these occupancy problems from a Gibbs-Poisson abundance model in the context of a population with infinitely many species. From this extension, a parameter γ naturally appears, which is a measure of richness or diversity of species. We rederive the sampling formulae for a population with infinitely many species, together with the distribution of the number P_k of distinct sampled species. We investigate the estimation of γ problem from the sample size k and the observed value of P_k .

We then exhibit a large special class of Gibbs-Poisson distributions having the property that sampling from a discrete abundance model may equivalently be viewed as a sampling problem from a random partition of unity, now in the continuum. When n is finite, this partition may be built upon normalizing n infinitely divisible iid positive random variables by its partial sum. It is shown that the sampling process in the continuum should generically be biased on the total length appearing in the latter normalization. A construction with size-biased sampling from the ranked normalized jumps of a subordinator is also supplied, would the problem under study present infinitely many species. We illustrate our point of view with many examples, some of which being new ones.

Keywords: Occupancy distributions. Sampling from Gibbs-Poisson distribution. Species abundance and frequencies. Biodiversity. Combinatorial probability. Subordinators.

Running title: Gibbs-Poisson sampling and occupancies.

1. INTRODUCTION AND OUTLINE OF MAIN RESULTS

Estimating the number n of unseen species from a k -sample displaying only $p \leq k$ distinct sampled species has been a challenging problem since the mid-twentieth century, [20]. It requires a model of species abundance together with a sampling model [16], and the answer to the latter question is of course model-dependent. In this work, we start with a discrete model of independent and identically distributed

(iid) stochastic species abundances $(\xi \stackrel{d}{=} \xi_1, \dots, \xi_n)$, based on Gibbs-Poisson distributions for ξ . We discuss the sampling formulae (species occupancy distributions, frequency of frequencies) in this discrete context; Typically, a k -sample drawn from the n -species abundances vector is the one obtained while conditioning this vector on summing to k (the sample size). It has to do with random allocation of balls into boxes, [33]. Various combinatorial identities arising in this setup are discussed. A distribution for the number of distinct visited species $P_{n,k}$ in a k -sample from a population of size n with Gibbs-Poisson abundance is derived. For this class of sampling problems, a ‘temperature’ type parameter $\theta > 0$ pops in naturally. It is a measure of how similar the box occupancy numbers look like statistically, after the sampling process: the smaller the values of θ , the more likely it is that these occupancy numbers are disparate. When sampling from ξ , we then discuss some aspects of the problem of the estimation of the number of species n from the size k of the sample and the number $P_{n,k}$ of distinct sampled species, assuming θ to be known. These results are supplied in Propositions 1 and 3.

It turns out that it always makes sense to study these occupancy problems from a Gibbs-Poisson abundance model in the context of a population with infinitely many species, provided n goes to ∞ together with θ going to 0 while $n\theta \rightarrow \gamma > 0$. From this construction, γ then appears as a measure of species richness or diversity. We rederive the sampling formulae (species occupancy distributions, frequency of frequencies) for a population with infinitely many species, together with the distribution of the number P_k of distinct sampled species. We discuss the problem of the estimation of the diversity parameter γ from the size k of the sample and the number P_k .

One particular model in the Gibbs-Poisson class has been discussed at length in the literature: the sampling problem from a population with discrete negative binomial distribution abundance ξ , both when the population is made of a finite number of species and when there are infinitely many of them. For this particular model, the sampling formulae are the ones of Ewens, [18]. It is also well-known that the Ewens sampling formulae may also be viewed as sampling from a continuous random Dirichlet partition of unity when the number of species is finite or as sampling from a random Poisson-Dirichlet partition of unity when there are infinitely many classes, [23]. This property is remarkable. By sampling from a continuous partition of unity, we mean that we draw independently k uniform random variables on the unit interval, looking at the subintervals of the partition which are being hit in the process to form the occupancy distributions of classes.

In this work, we exhibit a large class of Gibbs-Poisson distributions sharing with the negative binomial distribution this property that sampling from a discrete abundance model may equivalently be viewed as a sampling problem from a random partition of unity in the continuum. When n is finite, this partition may be built upon normalizing n infinitely divisible independent and identically distributed positive random variables $(Y \stackrel{d}{=} Y_1, Y_2, \dots, Y_n)$ by its partial sum. We exhibit the one-to-one correspondence between the laws of ξ and Y , assuming ξ to be in the special class. It is however shown that the sampling process in the continuum should generically be biased on the total length appearing in the latter normalization. A construction with size-biased sampling from the ranked normalized jumps of a subordinator is also supplied, would the problem under study present infinitely many species.

With this correspondence in mind, we discuss several examples, among which the Engen extended negative binomial model [15], the Berestycki-Pitman model [4] for the enumeration of forests of trees with generalized binomial generator, the polylog and the Mittag-Leffler models. When there are some reasons to suspect that the ranked species frequencies decay algebraically with the rank number, then the Engen model is well suited. Would one think of the ranked species frequencies as decaying exponentially with the rank number, then the Ewens model seems relevant. If the ranked species frequencies are believed to decay exponentially as some power of the rank number, then one should opt for the polylog model.

We end up giving a new example of ξ sharing some common issues with the Engen's model (in particular the algebraic decay property of the ranked frequencies). For this precise model, we are able to give an exact estimator of the biodiversity parameter.

2. SAMPLING FROM DISCRETE GIBBS-POISSON DISTRIBUTIONS

The sampling problem from a negative binomial abundance model and its Dirichlet counterpart in the continuum suggest to study the following general construction (see [24], [25] and [4] for similar recent interest).

2.1. Generating and partition function. With $\phi_\bullet := (\phi_m; m \geq 1)$ a sequence of non-negative real numbers with $\phi_1 > 0$, let

$$(1) \quad \phi(x) := \sum_{m \geq 1} \frac{\phi_m}{m!} x^m$$

be a formal power series in x . Assume that $x_0 := \sup(x > 0 : \phi(x) < \infty) \in (0, +\infty)$ is its convergence radius. Then $\phi(x)$ defines a convergent series on $|x| < x_0$ and it is absolutely monotone on $(0, x_0)$ in the sense that $\phi^{(n)}(x) \geq 0$ for all $n \geq 0$ and $x \in (0, x_0)$. We call it the local generating function.

Let $\theta > 0$ and consider the 'partition' function

$$(2) \quad Z_\theta(x) = e^{\theta\phi(x)}.$$

This function also defines a convergent series on $|x| < x_0$ with $Z_\theta(0) = 1$. Further, with $\sigma_k(\theta) = k! [x^k] Z_\theta(x)$

$$Z_\theta(x) = 1 + \sum_{k \geq 1} \frac{x^k}{k!} \sigma_k(\theta)$$

where, since $\partial_x Z_\theta(x) = \theta\phi'(x) Z_\theta(x)$, the Taylor coefficients $(\sigma_k(\theta); k \geq 1)$ of $Z_\theta(x)$ satisfy the general recurrence:

$$(3) \quad \sigma_{k+1}(\theta) = \theta \sum_{l=0}^k \binom{k}{l} \phi_{k-l+1} \sigma_l(\theta), \quad k \geq 0, \quad \sigma_0(\theta) \equiv 1.$$

Similarly, since $\partial_\theta Z_\theta(x) = \phi(x) Z_\theta(x)$, the Taylor coefficients $(\sigma_k(\theta); k \geq 1)$ of $Z_\theta(x)$ also satisfy the difference-differential recursion:

$$(4) \quad \sigma'_k(\theta) = \sum_{l=0}^{k-1} \binom{k}{l} \phi_{k-l} \sigma_l(\theta), \quad k \geq 1, \quad \sigma_0(\theta) = 1.$$

Let $[x^k] f(x)$ be the x^k -coefficient in the series expansion of the function $f(x)$. Then, clearly,

$$(5) \quad \sigma_k(\theta) = \sum_{l=1}^k B_{k,l}(\phi_\bullet) \theta^l,$$

with:

$$B_{k,l}(\phi_\bullet) = \frac{k!}{l!} [x^k] \phi(x)^l = \frac{k!}{l!} \sum_{\sum_{j=1}^l m_j = k}^* \prod_{j=1}^l \frac{\phi_{m_j}}{m_j!} \geq 0.$$

In the latter star-sum, summation runs over the integers $(m_1, \dots, m_l) \geq 1$, there are $\binom{k-1}{l-1}$ terms in such sums (In the sequel, the star-sums will always take into account only indexes ≥ 1). So $\sigma_k(\theta)$ is a degree- k Bell polynomial in θ whose θ^l coefficient is $B_{k,l}(\phi_\bullet)$ which is known as the Bell exponential polynomial in the variables ϕ_\bullet (see [10]). On $\theta > 0$, the function $\sigma_k(\theta)$ is convex and log-concave, for all k . As a polynomial with non-negative coefficients of degree k , $\sigma_k(\theta)$ has no strictly positive real root and (by Descartes rule sign) at most k real negative roots (including 0), counting roots with their order of multiplicity.

Remarks (Bell polynomials and convolutions).

(i) Define $(\phi * \phi)_m := \sum_{l=1}^{m-1} \binom{m}{l} \phi_l \phi_{m-l}$, $m \geq 2$, as the binomial self-convolution sequence of ϕ_m . Define ϕ_m^{*p} as the m^{th} term, $m \geq p$, of the sequence $\phi^{*p} := \phi * \dots * \phi$, p times; then the following convolution identity is well-known to hold:

$$B_{k,p}(\phi_\bullet) = \phi_k^{*p} / p!.$$

(ii) Because $Z_{\theta+\theta'}(x) = Z_\theta(x) Z_{\theta'}(x)$, the polynomials $\sigma_k(\theta)$ satisfy

$$(6) \quad \sigma_k(\theta + \theta') = \sum_{l=0}^k \binom{k}{l} \sigma_l(\theta) \sigma_{k-l}(\theta') \text{ for all } \theta, \theta' > 0,$$

and so they form a so-called binomial convolution sequence of polynomials.

If $p \geq 1$ is an integer, with $\sigma(1)_k^{*p} := (\sigma(1)^{*p})_k$

$$\sigma_k(p) = \sigma(1)_k^{*p} = \sum_{k_1 + \dots + k_p = k} \binom{k}{k_1 \dots k_p} \prod_{q=1}^p \sigma_{k_q}(1).$$

We clearly have

$$\sigma_k(p) = \sum_{q=1}^p \binom{p}{q} \sum_{k_1 + \dots + k_q = k}^* \binom{k}{k_1 \dots k_q} \prod_{r=1}^q \sigma_{k_r}(1).$$

In other words,

$$(7) \quad \sigma_k(p) = \sum_{q=1}^k \binom{p}{q} \sum_{k_1 + \dots + k_q = k}^* \binom{k}{k_1 \dots k_q} \prod_{r=1}^q \sigma_{k_r}(1),$$

where it is tacitly understood that $\binom{p}{q} = 0$ if $q > p$. This expression extends to non-integral arguments $\theta > 0$ of $\sigma_k(\cdot)$ as

$$(8) \quad \sigma_k(\theta) =: \sigma(1)_k^{*\theta} = \sum_{q=1}^k \binom{\theta}{q} \sum_{k_1+\dots+k_q=k}^* \binom{k}{k_1\dots k_q} \prod_{r=1}^q \sigma_{k_r}(1)$$

where $\binom{\theta}{q} =: \{\theta\}_q / q!$ with $\{\theta\}_q := \Gamma(\theta+1) / \Gamma(\theta-q+1) = \theta(\theta-1)\dots(\theta-q+1)$, the usual extension of $\binom{p}{q}$ for the expansion of $(1+x)^\theta$. From (8), it is clear again that $\sigma_k(\theta)$ is a degree- k polynomial in θ with no constant term. This expression should be used instead of (5) whenever the values at $\theta = 1$ of $\sigma_k(\cdot)$ are available in the first place, instead of the ϕ_\bullet .

(iii) Putting the expression of $\sigma_k(\theta)$ in (5) into the recurrence equation (3) which $(\sigma_k(\theta); k \geq 1)$ satisfies gives

$$(9) \quad l \cdot B_{k,l}(\phi_\bullet) = \sum_{j=l-1}^{k-1} \binom{k}{j} \phi_{k-j} B_{j,l-1}(\phi_\bullet).$$

Recalling the boundary conditions

$$B_{k,0}(\phi_\bullet) = B_{0,l}(\phi_\bullet) = 0,$$

except for $B_{0,0}(\phi_\bullet) := 1$, we get in particular

$$(10) \quad B_{k,1}(\phi_\bullet) = \phi_k \text{ and } B_{k,k}(\phi_\bullet) = \phi_1^k.$$

(iv) While performing the substitution $\theta \rightarrow 1/\theta$, $\sigma_k(\theta)$ should be mapped into the new polynomial with respect to $1/\theta$

$$\sigma_k(1/\theta) = \theta^{-(k+1)} \sum_{l=1}^k B_{k,k-l+1}(\phi_\bullet) \theta^l,$$

involving the ‘reversed’ Bell sequence $B_{k,k-l+1}(\phi_\bullet)$.

2.2. Discrete Gibbs-Poisson distributions arising from $Z_\theta(x)$. Let now $\xi \in \mathbb{N}_0 := \{0, 1, 2, \dots\}$ be a discrete random variable whose probability generating (pgf) is given by:

$$\Phi(u) := \mathbf{E}[u^\xi] = \frac{Z_\theta(xu)}{Z_\theta(x)}, \quad |u| \leq 1.$$

Since

$$(11) \quad \mathbf{E}[u^\xi] = e^{-\theta\phi(x)(1-\frac{\phi(xu)}{\phi(x)})},$$

ξ is in the compound Poisson class (as a Poisson sum of independent and identically distributed, say iid, jumps), hence infinitely divisible. The jumps’ height law is given by its pgf $\mathbf{E}[u^\delta] = \frac{\phi(xu)}{\phi(x)}$, where $\delta \in \mathbb{N} := \{1, 2, \dots\}$ is one of these jumps. Note that both $\mathbf{E}[\delta] = x \frac{\phi'(x)}{\phi(x)}$ and $\mathbf{E}[\xi] = \theta\phi(x) \mathbf{E}[\delta] = \theta x \phi'(x)$ are finite when $|x| < x_0$. Clearly

$$\mathbf{P}(\delta = m) = \frac{\phi_m x^m}{\phi(x) \cdot m!}, \quad m \geq 1 \text{ and}$$

$$\mathbf{P}(\xi = k) = \frac{\sigma_k(\theta) x^k}{Z_\theta(x) \cdot k!}, \quad k \geq 0.$$

We note that $\mathbf{P}(\xi = k)$ is also a Gibbs distribution with partition function $Z_\theta(x)$. With y defined by $x =: e^{-y}$, y is indeed the Legendre conjugate of $\mu := \mathbf{E}(\xi)$. So the parameter x in (11) can serve to adjust the mean μ of ξ . We call such distributions for ξ Gibbs-Poisson (GP) distributions. The random variable ξ will be used in the sequel as the typical abundance of some species in a population with n species. Due to its compound Poisson structure, it is tacitly assumed that the number of species is modelled as a Poisson sum of iid ‘clusters’ each with random size distributed like $\delta \geq 1$.

Consider now a sequence $(\xi \stackrel{d}{=} \xi_1, \dots, \xi_n, \dots)$ of iid Gibbs-Poisson random variables on \mathbb{N}_0 . Let $\zeta_n := \sum_{m=1}^n \xi_m$ denote their partial sum. Then, because ξ is in the compound-Poisson class due to $Z_\theta(x)^n = Z_{n\theta}(x)$

$$\mathbf{P}(\zeta_n = k) = \frac{\sigma_k(n\theta) x^k}{Z_{n\theta}(x) \cdot k!}, \quad k \geq 0.$$

This is also a Gibbs-Poisson distribution with corresponding partition function $Z_{n\theta}(x)$.

Remark: One could think of starting with $\phi(x) := \phi_0 + \sum_{m \geq 1} \frac{\phi_m}{m!} x^m$ with $\phi_0 \geq 0$ but because we shall deal with GP distributions whose pgfs are given by (11), ϕ_0 plays no role in our problem.

2.3. Sampling from infinitely divisible GP distributions. Define a random allocation scheme of k distinguishable particles or balls into n distinguishable boxes by

$$\mathbf{K}_{n,k} := (K_{n,k}(1), \dots, K_{n,k}(n)) \stackrel{d}{=} (\xi \stackrel{d}{=} \xi_1, \dots, \xi_n \mid \zeta_n = k),$$

so that $K_{n,k}(m)$ counts the number of particles in box m , $m = 1, \dots, n$ in a k -sample. Defining $\mathbf{K}_{n,k}$ from n iid ξ 's conditioned on summing to k , we get the generalized allocation scheme defined by Kolchin, (see [33]). When the ξ 's are in addition GP distributed, we call this model sampling from GP distributions.

For such random allocation models, each ξ_m may be viewed as the theoretical abundance of species $m = 1, \dots, n$ (the m^{th} species size). In this context, the random allocation scheme of k balls into n boxes accounts equivalently for a k -sampling process designed to model a random pick from $\boldsymbol{\xi}_n := (\xi_1, \dots, \xi_n)$ coming out from some measurement campaign which counts the number of times each species is being encountered, proportional to species abundances.

Remark: Since $\mathbf{E}[\xi] = \Phi'(1) = \theta x \phi'(x)$, $\theta > 0$ and $x \in (0, x_0)$, we could adjust the mean μ of ξ so that $\mathbf{E}[\xi] = \mu$. Then we would have the relation $\mu/\theta = x \phi'(x)$ (Legendre conjugation of x and μ) from which, by Lagrange inversion formula, an expression of $x = x(\mu/\theta)$ would follow. However, as we shall see, the actual value of the mean μ does not really matter after the sampling process.

Taking now into account the conditioning on the sample size in the definition of $\mathbf{K}_{n,k}$'s law, with $\mathbf{k}_n := (k_1, \dots, k_n) \geq 0$ summing to k

$$(12) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \frac{\mathbf{P}(\xi_1 = k_1, \dots, \xi_n = k_n)}{\mathbf{P}(\zeta_n = k)} = \frac{1}{\sigma_k(n\theta)} \binom{k}{k_1 \dots k_n} \prod_{m=1}^n \sigma_{k_m}(\theta),$$

this (Maxwell-Boltzmann) joint law being independent of x and so of the mean μ of the ξ 's. In other words, the joint probability generating function of $\mathbf{K}_{n,k}$ reads ($|u_m| \leq 1; m = 1, \dots, n$):

$$(13) \quad \mathbf{E} \left[\prod_{m=1}^n u_m^{K_{n,k}(m)} \right] = \frac{1}{\sigma_k(n\theta)} \sum_{|\mathbf{k}_n| := k_1 + \dots + k_n = k} \binom{k}{k_1 \dots k_n} \prod_{m=1}^n \sigma_{k_m}(\theta) u_m^{k_m}.$$

From (12), $w_{k_m}(\theta) := \sigma_{k_m}(\theta)/k_m!$ is seen to be the Boltzmann weight of box m with $e_{k_m}(\theta) := -\log(\sigma_{k_m}(\theta)/k_m!)$ being the energy required to put k_m balls into box number m . More precisely, for our random allocation GP model of particles (13) and from (5), the price to pay for having the l^{th} particle, $l \in \{1, \dots, k_m\}$, in box m simply is l and this event is assigned the weight $B_{k_m,l}(\phi_\bullet)/k_m!$. From this, one may view θ as a box fugacity parameter which, under our assumptions, is here common to all boxes (or species). Due to $\sigma_{k_m}(\theta)$ being a polynomial in θ with positive coefficients, the energy $e_{k_m}(\theta)$ is a decreasing function of θ and one may as well interpret θ as some temperature of the boxes (maybe through the monotone transformation $\theta \leftrightarrow e^{-1/T}$). Note that when θ approaches 0, the energy $e_{k_m}(\theta) \sim -\log \theta$ tends to $+\infty$: because the price to pay to put any number of particles into a box is extremely high, the optimal strategy is to put them all into a single box. One therefore expects that, as θ gets very small, the vector $\mathbf{K}_{n,k}$ gets very skewed (most balls into a single box), that is, completely opposite to the balanced multinomial $(k; \frac{1}{n}, \dots, \frac{1}{n})$ situation

$$\mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \frac{k!}{\prod_{m=1}^n k_m!} n^{-k}, \quad |\mathbf{k}_n| = k,$$

which is obtained for $\theta \rightarrow \infty$, as a result of $\sigma_{k_m}(\theta) \sim (\phi_1 \theta)^{k_m}$. In the latter balanced case, the most probable occupancy state is the centre $(k/n, \dots, k/n)$. As a conclusion, smaller the values of θ , the more likely it is that the occupancy numbers $K_{n,k}(m)$ are disparate.

From (12), the random vector-count $\mathbf{K}_{n,k}$ has exchangeable distribution (invariance under any permutation of the boxes numbers). But obviously, in the ordered version $\mathbf{K}_{(n),k}$ of the box occupancies $\mathbf{K}_{n,k}$, say with $K_{(n),k}(1) \geq \dots \geq K_{(n),k}(n)$, the boxes are not equally filled.

- Let us now compute the distribution of one of its typical component, say $K_{n,k}(1)$. With $l \in \{0, \dots, k\}$, we get

$$\begin{aligned} \mathbf{P}(K_{n,k}(1) = l) &= \mathbf{P}(\xi_1 = l) \frac{[u^{k-l}] \Phi(u)^{n-1}}{[u^k] \Phi(u)^n} = \\ &= \frac{\sigma_l(\theta) x^l [u^{k-l}] Z_\theta(xu)^{n-1}}{l! [u^k] Z_\theta(xu)^n} = \binom{k}{l} \frac{\sigma_l(\theta) \sigma_{k-l}((n-1)\theta)}{\sigma_k(n\theta)}. \end{aligned}$$

- Proceeding similarly, with $l \in \{0, \dots, k\}$, we would obtain the law of the partial sums $K_{n,k}(1) + \dots + K_{n,k}(m)$, $m < n$, as

$$\mathbf{P}(K_{n,k}(1) + \dots + K_{n,k}(m) = l) = \binom{k}{l} \frac{\sigma_l(m\theta) \sigma_{k-l}((n-m)\theta)}{\sigma_k(n\theta)}.$$

As required, $\sum_{l=0}^k \mathbf{P}(K_{n,k}(1) + \dots + K_{n,k}(m) = l) = 1$, as a result of $\sigma_k(\theta)$ being a convolution sequence of polynomials, from (6).

- Finally, define $\{k\}_l := k(k-1)\dots(k-l+1)$ with $\{k\}_0 := 1$ and let us now consider the falling factorial moments of $\mathbf{K}_{n,k}$.

Fix $\mathbf{l}_n := (l_1, \dots, l_n) \geq \mathbf{0}$ summing to $l \leq k$. We have

$$\mathbf{E} \left[\prod_{m=1}^n \{K_{n,k}(m)\}_{l_m} \right] = \prod_{m=1}^n l_m! \frac{[v^k] \prod_{m=1}^n [v^{l_m}] Z_\theta(xv(v_m+1))}{[v^k] Z_{n\theta}(xv)}.$$

Since $l_m! [v^{l_m}] Z_\theta(xv(v_m+1)) = \sum_{\mathbf{k}_m \geq \mathbf{l}_m} \frac{\sigma_{\mathbf{k}_m}(\theta) \cdot (xv)^{k_m}}{(k_m - l_m)!}$, with \mathbf{k}_n summing to $|\mathbf{k}_n| = k$, we get

$$(14) \quad \mathbf{E} \left[\prod_{m=1}^n \{K_{n,k}(m)\}_{l_m} \right] = \frac{\sum_{\mathbf{k}_n \geq \mathbf{l}_n} \prod_{m=1}^n \sigma_{k_m}(\theta) / (k_m - l_m)!}{\sigma_k(n\theta) / k!}.$$

These combinatorial quantities arise in the following resampling problem:

Subsampling without replacement from $\mathbf{K}_{n,n}$. Suppose $K_{n,n}(m)$, $m = 1, \dots, n$ are the random box occupancies of some sample with size exactly equal to the number n of boxes, generated by some compound-Poisson $\boldsymbol{\xi}_n$. So there are at most n boxes filled by a singleton as a result of $\sum_{m=1}^n K_{n,n}(m) = n$. Let $p \leq k \leq n$. We are interested in the event that after a random k -subsampling without replacement from $\mathbf{K}_{n,n}$, balls are reassigned at random into boxes so as to end up in a new occupancy $\mathbf{K}'_{n,k} := (K'_{n,k}(q); q = 1, \dots, p)$ where only $\Pi_{n,k} = p$ boxes (labeled in arbitrary order) are being filled. So $\mathbf{K}'_{n,k}$ obeys $\sum_{q=1}^p K'_{n,k}(q) = k$ and $K'_{n,k}(q) \geq 1$. Then, with $(k_1, \dots, k_p) \geq \mathbf{1}$ summing to k , the sampling without replacement strategy yields:

$$\begin{aligned} \mathbf{P}(K'_{n,k}(1) = k_1, \dots, K'_{n,k}(p) = k_p; \Pi_{n,k} = p) &= \binom{n}{p} \binom{k}{k_1 \dots k_p} \frac{\mathbf{E} \left(\prod_{q=1}^p \{K_{n,n}(q)\}_{k_q} \right)}{\{n\}_k} \\ &= \frac{\binom{n}{p}}{\binom{n}{k}} \mathbf{E} \prod_{q=1}^p \binom{K_{n,n}(q)}{k_q}. \end{aligned}$$

Summing over $(k_1, \dots, k_p) \geq \mathbf{1}$

$$\mathbf{P}(\Pi_{n,k} = p) = \frac{\binom{n}{p}}{\{n\}_k} \sum_{k_1 + \dots + k_p = k}^* \binom{k}{k_1 \dots k_p} \mathbf{E} \left(\prod_{q=1}^p \{K_{n,n}(q)\}_{k_q} \right)$$

is the probability that in a k -subsampling without replacement from $\mathbf{K}_{n,n}$ exactly $p \leq k \leq n$ boxes will be filled. Using (14), with $\mathbf{k}_p = (k_1, \dots, k_p) \geq \mathbf{1}$ satisfying

$|\mathbf{k}_p| = k$, we have

$$\mathbf{E} \left[\prod_{q=1}^p \{K_{n,n}(q)\}_{k_q} \right] = \frac{\sum_{\mathbf{l}_p \geq \mathbf{0}} \prod_{q=1}^p \sigma_{k_q + l_q}(\theta) / l_q!}{\sigma_n(n\theta) / n!}$$

and the full expression of the probabilities $\mathbf{P}(\Pi_{n,k} = p)$ can be obtained in terms of the original weights $w_k(\theta) = \sigma_k(\theta) / k!$.

These questions arise in the discrete theory of compound-Poisson coalescent processes where $\mathbf{K}_{n,n}$ is the random reproduction law of some Markov branching process preserving the total number n of individuals over the subsequent generations, [28]. The (m, l) entry of the transition matrix of this Markov process on the state-space $\{0, \dots, n\}$ is

$$\mathbf{P}(K_{n,n}(1) + \dots + K_{n,n}(m) = l) = \binom{n}{l} \frac{\sigma_l(m\theta) \sigma_{n-l}((n-m)\theta)}{\sigma_n(n\theta)}, \quad m, l \in \{0, \dots, n\},$$

looking at the descent of all size- m subsample of the full population with size n . Clearly, the states $\{0, n\}$ are both absorbing.

Looking at this process backward in time, individuals are seen to merge, giving rise to the ancestral process where individuals are identified if they share a common ancestor one generation backward in time.

The quantity $\mathbf{P}(K'_{n,k}(1) = k_1, \dots, K'_{n,k}(p) = k_p; \Pi_{n,k} = p)$ is then the probability that a one-step back (k_1, \dots, k_p) to p merger for a subsample of size k occurs in the ancestral process. The lower-triangular stochastic matrix $Q_{k,p}^{(n)} := \mathbf{P}(\Pi_{n,k} = p)$ is the transition matrix of this pure death Markov process on $\{0, \dots, n\}$.

Number of filled boxes in $\mathbf{K}_{n,k}$: Let $P_{n,k} := \sum_{m=1}^n \mathbf{I}(K_{n,k}(m) > 0)$ count the number of non empty boxes in the sampling process from ξ_n . With $1 \leq p \leq n \wedge k$, the probability that there are only $P_{n,k} = p \in [n]$ visited boxes in the sampling process, the $n - p$ remaining ones remaining empty, is easily obtained as follows: Using exchangeability of $\mathbf{K}_{n,k}$, with $\mathbf{k}_p := (k_1, \dots, k_p) \geq \mathbf{1}$ summing to k , using (13),

$$(15) \quad \mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(p) = k_p; P_{n,k} = p) = \binom{n}{p} \frac{k!}{\sigma_k(n\theta)} \prod_{q=1}^p \frac{\sigma_{k_q}(\theta)}{k_q!}$$

is the joint probability that there are $p \in [n]$ non-empty boxes and that (k_1, \dots, k_p) are the respective occupancies of the p filled boxes (labeled in arbitrary order). Note that

$$P_{k,p}^{(n)} := \mathbf{P}(P_{n,k} = p) = \binom{n}{p} \frac{k!}{\sigma_k(n\theta)} \sum_{k_1 + \dots + k_p = k}^* \prod_{q=1}^p \frac{\sigma_{k_q}(\theta)}{k_q!}$$

is the probability that in a k -sample from n species with abundance ξ_n , the exact number of distinct visited species is p . In particular, $P_{k,1}^{(n)} := n \frac{\sigma_k(\theta)}{\sigma_k(n\theta)}$ is the probability that in this k -sample, only one species is discovered (whichever it is).

The expression (15) turns out to be the canonical Gibbs distribution on finite size- n partitions of k into p distinct clusters (the filled boxes), derived from the weight

sequence ϕ_\bullet . In this language, the normalizing quantity $\sigma_k(n\theta)/k!$ is called the canonical Gibbs partition function.

Now, from (15), with $\{n\}_p := n!/(n-p)!$

$$(16) \quad \mathbf{P}(P_{n,k} = p) = \frac{\{n\}_p}{\sigma_k(n\theta)} B_{k,p}(\sigma_\bullet(\theta)), \quad p \in \{1, \dots, n \wedge k\},$$

where

$$(17) \quad B_{k,p}(\sigma_\bullet(\theta)) := \frac{k!}{p!} \sum_{\sum_{q=1}^p k_q = k}^* \prod_{q=1}^p \frac{\sigma_{k_q}(\theta)}{k_q!} = \frac{k!}{p!} [x^k] (Z_\theta(x) - 1)^p$$

is now a Bell polynomial in the polynomial variables $\sigma_\bullet(\theta) := (\sigma_1(\theta), \sigma_2(\theta), \dots)$.

Conditioning the canonical Gibbs distribution on the number of filled cells being equal to p yields the corresponding micro-canonical distribution as

$$\begin{aligned} \mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(p) = k_p \mid P_{n,k} = p) \\ = \frac{k!}{p!} \frac{1}{B_{k,p}(\sigma_\bullet(\theta))} \prod_{q=1}^p \frac{\sigma_{k_q}(\theta)}{k_q!}. \end{aligned}$$

The new normalizing constant $B_{k,p}(\sigma_\bullet(\theta))/k!$ may be called the microcanonical partition function.

The microcanonical distribution is independent of n . So, for all models studied here, the map $P \rightarrow \mathbf{P}(P_{n,k} = P)$ is a sufficient statistics in the estimation of n problem from occupancy data (assuming θ known).

Let us now give some additional details on the distribution of $P_{n,k}$.

Proposition 1. (a) *Assume $k \geq n$. The probability generating function of $P_{n,k}$ is given by*

$$(18) \quad \mathbf{E}(u^{P_{n,k}}) = \sum_{p=0}^{n-1} \binom{n}{p} u^{n-p} (1-u)^p \frac{\sigma_k((n-p)\theta)}{\sigma_k(n\theta)},$$

with:

$$(19) \quad \mathbf{P}(P_{n,k} = p) = \binom{n}{p} \sum_{q=1}^p (-1)^{p-q} \binom{p}{q} \frac{\sigma_k(q\theta)}{\sigma_k(n\theta)}, \quad p \in \{1, \dots, n\}.$$

In addition,

$$\begin{aligned} \mathbf{E}(P_{n,k}) &= n \left(1 - \frac{\sigma_k((n-1)\theta)}{\sigma_k(n\theta)} \right) \\ \sigma^2(P_{n,k}) &= n \left(\frac{\sigma_k((n-1)\theta)}{\sigma_k(n\theta)} + (n-1) \frac{\sigma_k((n-2)\theta)}{\sigma_k(n\theta)} - n \left(\frac{\sigma_k((n-1)\theta)}{\sigma_k(n\theta)} \right)^2 \right) \end{aligned}$$

(b) *If $k < n$, (18) and (19) still hold, but now with a modified support for $P_{n,k}$ s law:*

$$(20) \quad \mathbf{P}(P_{n,k} = p) = \binom{n}{p} \sum_{q=1}^p (-1)^{p-q} \binom{p}{q} \frac{\sigma_k(q\theta)}{\sigma_k(n\theta)}, \quad p \in \{1, \dots, k\}.$$

Proof: (a) This follows from $B_{k,p}(\sigma_{\bullet}(\theta)) = \frac{k!}{p!} [x^k] (Z_{\theta}(x) - 1)^p$. Indeed, from (16)

$$\begin{aligned} \mathbf{E}(u^{P_{n,k}}) &= \sum_{p=0}^n u^p \{n\}_p \frac{B_{k,p}(\sigma_{\bullet}(\theta))}{\sigma_k(n\theta)} = \frac{k!}{\sigma_k(n\theta)} \sum_{p=0}^n \binom{n}{p} [x^k] (u(Z_{\theta}(x) - 1))^p \\ &= \frac{k!}{\sigma_k(n\theta)} [x^k] (1 - u + uZ_{\theta}(x))^n = \frac{k!}{\sigma_k(n\theta)} \sum_{p=0}^n \binom{n}{p} u^{n-p} (1-u)^p [x^k] Z_{\theta}(x)^{n-p} \\ &= \sum_{p=0}^{n-1} \binom{n}{p} u^{n-p} (1-u)^p \frac{\sigma_k((n-p)\theta)}{\sigma_k(n\theta)} \end{aligned}$$

The alternating sum expression of $\mathbf{P}(P_{n,k} = p)$ follows from extracting $[u^p] \mathbf{E}(u^{P_{n,k}})$ and the mean and variance of $P_{n,k}$ from the evaluations of the first and second derivatives of $\mathbf{E}(u^{P_{n,k}})$ with respect to u at $u = 1$.

(b) follows from similar considerations. Indeed, in principle, we should start with $\mathbf{E}(u^{P_{n,k}}) = \sum_{p=0}^k u^p \{n\}_p \frac{B_{k,p}(\sigma_{\bullet}(\theta))}{\sigma_k(n\theta)}$ where the p -sum now stops at $p = k = k \wedge n$. But the upper bound of this p -sum can be extended to n because $B_{k,p}(\sigma_{\bullet}(\theta)) = 0$ if $p > k$. \diamond

In (16), the new combinatorial coefficients $B_{k,p}(\sigma_{\bullet}(\theta))$ come into the game. They are given by

Corollary 2. *With $S_{l,p}$ the second kind Stirling numbers,*

$$B_{k,p}(\sigma_{\bullet}(\theta)) = \sum_{l=p}^k B_{k,l}(\phi_{\bullet}) S_{l,p} \theta^l = \theta^p \cdot \sum_{l=0}^{k-p} B_{k,p+l}(\phi_{\bullet}) S_{l+p,p} \theta^l,$$

showing that $B_{k,p}(\sigma_{\bullet}(\theta))$ is itself a polynomial in θ with larger (smaller) degree k (respectively p).

Proof: From (16) and (19), we have

$$B_{k,p}(\sigma_{\bullet}(\theta)) = \frac{1}{p!} \sum_{q=1}^p (-1)^{p-q} \binom{p}{q} \sigma_k(q\theta) \binom{1}{q}.$$

Recalling $\sigma_k(\theta) = \sum_{l=1}^k \theta^l B_{k,l}(\phi_{\bullet})$ and observing $S_{l,p} = \sum_{q=1}^p (-1)^{p-q} \binom{p}{q} q^l$ gives the result after reversing the sums. This result actually is in accordance with the Faà di Bruno formula (see [10]) giving the Taylor coefficients of the composition function g of the two analytic functions $g(x) := e_{\lambda,\theta} \circ \phi(x)$ where $e_{\lambda,\theta}(x) := e^{\lambda(e^{\theta x} - 1)}$ as

$$S_k(\lambda) = \sum_{l=1}^k e_l(\theta, \lambda) B_{k,l}(\phi_{\bullet}),$$

with $e_l(\theta, \lambda) = \theta^l \sum_{p=1}^l \lambda^p S_{l,p}$ the l^{th} Taylor coefficient of $e_{\lambda,\theta}(x)$. Clearly indeed,

$$g(x) = e^{\lambda(Z_{\theta}(x) - 1)} = 1 + \sum_{k \geq 1} \frac{x^k}{k!} S_k(\lambda) =: 1 + \sum_{k \geq 1} \frac{x^k}{k!} \left(\sum_{p=1}^k \lambda^p B_{k,p}(\sigma_{\bullet}(\theta)) \right)$$

¹This identity was derived in a different way in [43].

and the λ^p -coefficient of $S_k(\lambda)$ is exactly $\sum_{l=p}^k B_{k,l}(\phi_\bullet) S_{l,p} \theta^l$. \diamond

2.4. The estimation of n problem. Let us now to discuss the important question of estimating the unknown number of species n based on the data k and P (assuming θ is known), recalling $\mathbf{P}(P_{n,k} = P)$ is a sufficient statistics in this estimation problem. Our forthcoming statement holds for a class of ϕ which is such that the degree- k polynomial $\sigma_k(\theta) \in ZR_-$ (has only real non-positive zeroes). We recall that $\sigma_k(\theta) \in ZR_-$ iff the matrix M with entries $M_{i,j} = B_{k,i-j}(\phi_\bullet)$, $i, j = 0, \dots, k$, with $B_{k,l}(\phi_\bullet) = 0$ if $l \notin \{l : B_{k,l}(\phi_\bullet) > 0\}$ is totally positive of order k (with $l = 1, \dots, k$, each $l \times l$ minor of M has a nonnegative determinant), [39]. Therefore, there is no simple way to check whether or not $\sigma_k(\theta) \in ZR_-$.

We also recall here, [39], that if and only if the matrix $M = M_{i,j}$ would be such that all its 2×2 minors have a nonnegative determinant, then the sequence $B_{k,l}(\phi_\bullet)$, $l = 1, \dots, k$ (with no internal zeros) is l -log-concave (the l -sequence $B_{k,l}(\phi_\bullet)$ is a Pòlya frequency sequence of order 2). If this is the case, we shall say $\sigma_k(\theta) \in PF_2$.

Proposition 3. *Suppose $\sigma_k(\theta) \in ZR_-$. Then the log-likelihood $\log \mathbf{P}(P_{n,k} = p)$ attains its maximum in n at least once and at most twice in which latter case, the two values are adjacent integers. This leads to the maximum likelihood estimator \hat{n} of n characterized by:*

$$\hat{n} = \sup \left\{ n > 0 : \frac{\mathbf{P}(P_{n,k} = P)}{\mathbf{P}(P_{n-1,k} = P)} > 1 \right\},$$

where this last quantity verifies $\hat{n} = \lceil P \rceil$, the smallest integer $\geq P$, when the set of integers $\left\{ n > 0 : \frac{\mathbf{P}(P_{n,k} = P)}{\mathbf{P}(P_{n-1,k} = P)} > 1 \right\}$ is empty.

When this is not the case and for large n , an approximation of the estimator \hat{n} of n is given by the implicit equation

$$P = \hat{n} \left(1 - \frac{\sigma_k((\hat{n} - 1)\theta)}{\sigma_k(\hat{n}\theta)} \right).$$

Proof: We extend (16) to n a real variable, so we can differentiate $\log \mathbf{P}(P_{n,k} = p)$ with respect to $n > p$. In this domain, we have $\partial_n \log \{n\}_p = \sum_{q=0}^{p-1} \frac{1}{n-q}$, and so we get

$$\partial_n \log \mathbf{P}(P_{n,k} = p) = \sum_{q=0}^{p-1} \frac{1}{n-q} - \partial_n \log \sigma_k(n\theta).$$

Suppose the polynomial $\sigma_k(\theta) \in ZR_-$ has zeroes $-r_{l,k}$ where: $0 = r_{1,k} \leq \dots \leq r_{k,k}$. Then $\sigma_k(n\theta) = \prod_{l=1}^k (n\theta + r_{l,k})$ and $\partial_n \log \sigma_k(n\theta) = \sum_{l=1}^k (n + r_{l,k}/\theta)^{-1}$, together with $\partial_n^2 \log \sigma_k(n\theta) = -\sum_{l=1}^k (n + r_{l,k}/\theta)^{-2} < 0$.

If $\sum_{q=0}^{p-1} \frac{1}{n-q} - \sum_{l=1}^k (n + r_{l,k}/\theta)^{-1} \stackrel{(*)}{=} 0$, then

$$\partial_n^2 \log \mathbf{P}(P_{n,k} = p) = -\sum_{q=0}^{p-1} \frac{1}{(n-q)^2} + \sum_{l=1}^k (n + r_{l,k}/\theta)^{-2} < 0,$$

showing that the likelihood is log-concave around the critical points. Hence, if \hat{n} solves $(*)$ it is a local maximum and there is no local minimum. The maximum likelihood estimator of real n is thus unique.

Coming back to n integer, we deduce that the maximum likelihood estimator of n is the integer $\sup \left\{ n > 0 : \frac{\mathbf{P}(P_{n,k}=P)}{\mathbf{P}(P_{n-1,k}=P)} > 1 \right\}$. When n is large, it may thus be approximated by $\frac{\mathbf{P}(P_{\hat{n},k}=P)}{\mathbf{P}(P_{\hat{n}-1,k}=P)} = 1$, leading to

$$\frac{\{\hat{n}\}_P \sigma_k((\hat{n}-1)\theta)}{\{\hat{n}-1\}_P \sigma_k(\hat{n}\theta)} = 1 \text{ or } P = \hat{n} \left(1 - \frac{\sigma_k((\hat{n}-1)\theta)}{\sigma_k(\hat{n}\theta)} \right). \diamond$$

An alternative estimator. Let us now come to an alternative estimator of n (see [29] for a similar approach in the particular context of the Dirichlet model given by $\phi(x) = -\alpha \log(1-x)$). Suppose that for all $\theta > 0$ and $k \geq 1$, $B_{k,p}(\sigma_\bullet(\theta))$ is a log-concave p -sequence (equivalently, each degree- k λ -polynomial $S_k(\lambda) \in PF_2$). Then, by Darroch Theorem [12], $B_{k,p}(\sigma_\bullet(\theta))$ is p -unimodal or bimodal at two consecutive p . Because the p -sequence $\{n\}_p$ is also log-concave, $\{n\}_p B_{k,p}(\sigma_\bullet(\theta))$ is itself p -log-concave. For each n therefore, there is a unique \tilde{p} defined as $\tilde{p} = \sup \left\{ p > 0 : \frac{\mathbf{P}(P_{n,k}=p)}{\mathbf{P}(P_{n,k}=p-1)} > 1 \right\}$. Inverting the map $n \rightarrow \tilde{p}(n)$, given $p = P$, there exists a unique \tilde{n} , approximately characterized by $\frac{\mathbf{P}(P_{\tilde{n},k}=P-1)}{\mathbf{P}(P_{\tilde{n},k}=P)} = 1$, which can serve as an alternative estimator of n given the data (k, P) . From (16), it is thus given by

$$\tilde{n} = P + \frac{B_{k,P-1}(\sigma_\bullet(\theta))}{B_{k,P}(\sigma_\bullet(\theta))}.$$

If $k \geq n$, taking the expectation with respect to P , we have

$$\begin{aligned} \mathbf{E}(\tilde{n}) &= \mathbf{E}(P) + \sum_{p=1}^n \frac{B_{k,p-1}(\sigma_\bullet(\theta))}{B_{k,p}(\sigma_\bullet(\theta))} \frac{\{n\}_p}{\sigma_k(n\theta)} B_{k,p}(\sigma_\bullet(\theta)) \\ &= \mathbf{E}(P) + \sum_{p=2}^n \frac{\{n\}_p B_{k,p-1}(\phi_\bullet)}{\sigma_k(n\theta)} = \mathbf{E}(P) + \sum_{p=2}^n (n - (p-1)) \frac{\{n\}_{p-1} B_{k,p-1}(\phi_\bullet)}{\sigma_k(n\theta)} \\ &= \mathbf{E}(P) + n \left(1 - \frac{\{n\}_n B_{k,n}(\phi_\bullet)}{\sigma_k(n\theta)} \right) - \left(\mathbf{E}(P) - n \frac{\{n\}_n B_{k,n}(\phi_\bullet)}{\sigma_k(n\theta)} \right) = n. \end{aligned}$$

So, when $k \geq n$, \tilde{n} is an unbiased estimator of n . The Fisher information of n is

$$I(n) = -\mathbf{E}(\partial_n^2 \log \mathbf{P}(P_{n,k}=P)) = \mathbf{E} \left(\sum_{q=0}^{P-1} \frac{1}{(n-q)^2} \right) - \sum_{l=1}^k (n + r_{l,k}/\theta)^{-2} > 0,$$

giving the Cramér-Rao bound for the variance: $\sigma^2(\tilde{n}) \geq I(n)^{-1}$.

2.5. Frequency of frequencies. This suggests to look at the frequency of frequencies distribution problem. For $i = 0, \dots, k$, let now

$$(21) \quad A_{n,k}(i) = \sum_{m=1}^n \mathbf{I}(K_{n,k}(m) = i)$$

count the number of boxes visited i times by the k -sample, with $A_{n,k}(0) = n - P_{n,k}$, the number of empty boxes.

Let (a_1, a_2, \dots) be non-negative integers satisfying $\sum_{i \geq 1} a_i = p$ and $\sum_{i \geq 1} i a_i = k$.

It follows from (12) that

$$(22) \quad \mathbf{P}(A_{n,k}(1) = a_1, A_{n,k}(2) = a_2, \dots) = \frac{\{n\}_p \cdot k!}{\sigma_k(n\theta)} \prod_{i \geq 1} \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i} \frac{1}{a_i!} \right\}.$$

Taking $A_{n,k}(0)$ into account, let (a_0, a_1, \dots, a_k) be non-negative integers satisfying $\sum_{i=0}^k a_i = n$ and $\sum_{i=1}^k i a_i = k$. Then

$$\mathbf{P}(A_{n,k}(0) = a_0, A_{n,k}(1) = a_1, \dots, A_{n,k}(k) = a_k) = \frac{n! \cdot k!}{\sigma_k(n\theta)} \prod_{i=0}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i} \frac{1}{a_i!} \right\}.$$

Note from this that, with $\sum_{i=1}^k i a_i = k$ and $\sum_1^k a_i \leq n$, the normalization condition gives the identity

$$(23) \quad \sum_{a_1, \dots, a_k} \frac{k!}{(n - \sum_1^k a_i)!} \prod_{i=1}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i} \frac{1}{a_i!} \right\} = \frac{\sigma_k(n\theta)}{n!}.$$

From this, we get:

Proposition 4. *If $p = n - a_0$, the joint distribution of $(A_{n,k}(1), \dots, A_{n,k}(k))$ and $P_{n,k}$ reads*

$$(24) \quad \mathbf{P}(A_{n,k}(1) = a_1, \dots, A_{n,k}(k) = a_k; P_{n,k} = p) = \frac{\{n\}_p \cdot k!}{\sigma_k(n\theta)} \prod_{i=1}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i} \frac{1}{a_i!} \right\}.$$

Let us compute the falling factorial moments of $A_{n,k}(i)$, $i = 1, \dots, k$.

Proposition 5. *Let r_i , $i = 1, \dots, k$ be non-negative integers satisfying $\sum_1^k r_i = r \leq n$ and $\sum_1^k i r_i = \kappa \leq k$. We have*

$$(25) \quad \mathbf{E} \left[\prod_{i=1}^k \{A_{n,k}(i)\}_{r_i} \right] = \{n\}_r \{k\}_\kappa \frac{\sigma_{k-\kappa}((n-r)\theta)}{\sigma_k(n\theta)} \prod_{i=1}^k \left(\frac{\sigma_i(\theta)}{i!} \right)^{r_i}.$$

Proof:

$$\begin{aligned} \mathbf{E} \left[\prod_{i=1}^k \{A_{n,k}(i)\}_{r_i} \right] &= \frac{n! \cdot k!}{\sigma_k(n\theta)} \sum_{a_1, \dots, a_k} \frac{1}{(n - \sum_1^k a_i)!} \prod_{i=1}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i} \frac{1}{(a_i - r_i)!} \right\} \\ &= \frac{n! \cdot k!}{\sigma_k(n\theta)} \prod_{i=1}^k \left(\frac{\sigma_i(\theta)}{i!} \right)^{r_i} \sum_{a_1, \dots, a_k} \frac{1}{(n - \sum_1^k a_i)!} \prod_{i=1}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i - r_i} \frac{1}{(a_i - r_i)!} \right\}. \end{aligned}$$

The normalization condition (23) gives:

$$\sum_{a_1, \dots, a_k} \frac{1}{(n - \sum_1^k a_i)!} \prod_{i=1}^k \left\{ \left(\frac{\sigma_i(\theta)}{i!} \right)^{a_i - r_i} \frac{1}{(a_i - r_i)!} \right\} = \frac{\sigma_{k-\kappa}((n-r)\theta)}{(n-r)! \cdot (k-\kappa)!}.$$

Finally, we get

$$\mathbf{E} \left[\prod_{i=1}^k \{A_{n,k}(i)\}_{r_i} \right] = \{n\}_r \{k\}_\kappa \frac{\sigma_{k-\kappa}((n-r)\theta)}{\sigma_k(n\theta)} \prod_{i=1}^k \left(\frac{\sigma_i(\theta)}{i!} \right)^{r_i} \cdot \diamond$$

In particular, if all $r_i = 0$, except for one i for which $r_i = 1$ ($r = 1$, $\kappa = i$), then

$$(26) \quad \mathbf{E}[A_{n,k}(i)] = n \{k\}_i \frac{\sigma_{k-i}((n-1)\theta) \sigma_i(\theta)}{\sigma_k(n\theta) i!} = n \mathbf{P}(K_{n,k}(1) = i).$$

This shows that the expected number of cells visited i times is n times the probability that there are i visits to (say) cell one. In fact, we have the more general statement:

Corollary 6. *If $r_i = \#\{m \in \{1, \dots, n\} : k_m = i\}$, then*

$$\mathbf{E} \left[\prod_{i=1}^k \{A_{n,k}(i)\}_{r_i} \right] = n! \mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(n) = k_n),$$

so that the joint falling factorial moments of the A s can directly be obtained in terms of the joint distribution of the K s.

Proof: With the r_i as stated, using a sampling without replacement argument

$$\begin{aligned} \mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(n) = k_n \mid A_{n,k}(1), \dots, A_{n,k}(k)) = \\ \frac{1}{n!} \prod_{i=1}^k \{A_{n,k}(i)\}_{r_i}. \end{aligned}$$

Averaging over the A s gives the announced result. \diamond

2.6. The $*$ -limit of sampling distributions (the infinitely many species abundance model). Theoretical biologists work in a framework of a population with infinitely many species, with the more frequent one occurring with abundance $\xi_{(1)}$, second more frequent with abundance $\xi_{(2)}$, ... with $\xi_{(1)} \geq \xi_{(2)} \geq \dots$. Sampling from $(\xi_{(1)}, \xi_{(2)}, \dots)$ turns out to be a challenging problem. This requires the introduction of a model with infinitely many species (not only n) with ordered abundance $\xi_{(m)}$, $m \geq 1$. For such abundance models, a k -sample will represent the met individuals of various species when sampling from a population with infinitely many species, [8]. One can think of obtaining such models while considering the limit $n \rightarrow \infty$ and $\theta \rightarrow 0$ in the preceding model with n species. Indeed, as we saw, small values of the temperature $\theta > 0$ was an indication on how disparate the abundance numbers ξ_n were. But it may happen that some (necessarily few) of the $(\xi_m)_{m=1}^n$ are not so small with the hope that the ranked $\xi_{(m)}$ s would have a non-degenerate limit as $n \rightarrow \infty$, $\theta \rightarrow 0$ while $n\theta \rightarrow \gamma > 0$. We call such a limit the $*$ -limit.

It turns out that for the class of Gibbs-Poisson allocation models considered in this Section, the $*$ -limit always makes sense. This illustrates that limiting models should come down from some finitary counterpart, [21]. We first verify our claim intuitively.

Observing indeed that

$$\sigma_k(\theta) \sim_{\theta \downarrow 0} \theta B_{k,1}(\phi_\bullet) = \theta \phi_k \text{ and } B_{k,p}(\sigma_\bullet(\theta)) \sim_{\theta \downarrow 0} \theta^p B_{k,p}(\phi_\bullet)$$

and recalling $\{n\}_p \sim_{n \rightarrow \infty} n^p$, we easily get:

Proposition 7. From (15), with $(k_1, \dots, k_p) \geq 1$ summing to k and $p \leq k$

$$\mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(p) = k_p; P_{n,k} = p) \rightarrow_*$$

$$(27) \quad \mathbf{P}^*(K_k(1) = k_1, \dots, K_k(p) = k_p; P_k = p) = \frac{k!}{p!} \frac{\gamma^p}{\sigma_k(\gamma)} \prod_{q=1}^p \frac{\phi_{k_q}}{k_q!}$$

and

$$(28) \quad \mathbf{P}(P_{n,k} = p) \rightarrow_* \mathbf{P}^*(P_k = p) = \frac{\gamma^p}{\sigma_k(\gamma)} B_{k,p}(\phi_\bullet).$$

Equivalently, the limiting probability generating function of P_k also reads

$$(29) \quad \mathbf{E}^*(u^{P_k}) = \frac{\sigma_k(\gamma u)}{\sigma_k(\gamma)},$$

with mean $\mathbf{E}^*(P_k) = \gamma \frac{\sigma'_k(\gamma)}{\sigma_k(\gamma)}$. From this,

$$(30) \quad \mathbf{P}^*(K_k(1) = k_1, \dots, K_k(p) = k_p \mid P_k = p) = \frac{k!}{p!} \frac{1}{B_{k,p}(\phi_\bullet)} \prod_{q=1}^p \frac{\phi_{k_q}}{k_q!}$$

which is independent of γ .

Further, from (22), with (a_1, a_2, \dots) satisfying $\sum_{i \geq 1} i a_i = k$ and $\sum_{i \geq 1} a_i = p$

$$\mathbf{P}(A_{n,k}(1) = a_1, A_{n,k}(2) = a_2, \dots) \rightarrow_*$$

$$(31) \quad \mathbf{P}^*(A_k(1) = a_1, A_k(2) = a_2, \dots) = \frac{\gamma^p k!}{\sigma_k(\gamma)} \prod_{i=1}^k \frac{(\phi_i/i!)^{a_i}}{a_i!}.$$

Equivalently, from (24)

$$\mathbf{P}(A_{n,k}(1) = a_1, \dots, A_{n,k}(k) = a_k; P_{n,k} = p) \rightarrow_*$$

$$(32) \quad \mathbf{P}^*(A_k(1) = a_1, \dots, A_k(k) = a_k; P_k = p) = \frac{\gamma^p k!}{\sigma_k(\gamma)} \prod_{i=1}^k \frac{(\phi_i/i!)^{a_i}}{a_i!}.$$

and

$$(33) \quad \mathbf{P}^*(A_k(1) = a_1, \dots, A_k(k) = a_k \mid P_k = p) = \frac{k!}{B_{k,p}(\phi_\bullet)} \prod_{i=1}^k \frac{(\phi_i/i!)^{a_i}}{a_i!},$$

which is also independent of γ .

(27) or (32) are the canonical Gibbs distributions on partitions of k into p distinct clusters, derived from the weight sequence ϕ_\bullet . In this context, the normalizing quantity $\sigma_k(\gamma)/k!$ is called the canonical Gibbs partition polynomial⁽²⁾. Conditioning the canonical Gibbs distribution on the number of filled boxes being equal to p yields the corresponding micro-canonical distributions (30) or (33). The new normalizing constant $B_{k,p}(\phi_\bullet)/k!$ is called the microcanonical partition function.

Let us finally compute the falling factorial moments of $A_k(i)$, $i = 1, \dots, k$.

²The occupancy distribution (32) also appears in Ecology in a species abundance model occurring in the Hubbell's unified neutral theory of biodiversity. In this context, γ is the fundamental biodiversity number, [26].

Proposition 8. *Let $r_i, i = 1, \dots, k$ be non-negative integers satisfying $\sum_1^k r_i = r$ and $\sum_1^k i r_i = \kappa \leq k$. We have*

$$(34) \quad \mathbf{E}^* \left[\prod_{i=1}^k \{A_k(i)\}_{r_i} \right] = \gamma^r \{k\}_\kappa \frac{\sigma_{k-\kappa}(\gamma)}{\sigma_k(\gamma)} \prod_{i=1}^k \left(\frac{\phi_i}{i!} \right)^{r_i}.$$

Proof: This follows straightforwardly from Proposition 5 while taking the $*$ -limit and using $\sigma_i(\theta) \sim \theta \phi_i$ for small θ . This formula is a generalization of the Watterson expression [42] obtained in the special Ewens context when $\phi(x) = -\log(1-x)$, with $\phi_i = (i-1)!$ and $\sigma_k(\gamma) = \Gamma(\gamma+k)/\Gamma(\gamma) =: (\gamma)_k$; see Section 3 for a special account on this model. From (34), we easily get a closed-form expression for the mean $\mathbf{E}^*(A_k(i))$, $i \leq k$, the variance $\sigma^{*2}(A_k(i))$, for all i with $2i \leq k$ and the covariance $\text{Cov}^*(A_k(i_1), A_k(i_2))$ for all $i_1 \neq i_2, i_1 + i_2 \leq k$. \diamond

We observed that (30) or (33) were independent of γ , meaning that $P \rightarrow \mathbf{P}^*(P_k = P)$ is a sufficient statistics in the estimation of γ problem. Let us now briefly investigate this problem.

2.7. The estimation of γ problem. We wish now to discuss the question of estimating γ from the data k and P . From (28)

$$\partial_\gamma \log \mathbf{P}^*(P_k = p) = p/\gamma - \partial_\gamma \log \sigma_k(\gamma).$$

Suppose the polynomial $\sigma_k(\gamma) \in ZR_-$ with zeroes $-r_{l,k}$ where: $0 = r_{1,k} \leq \dots \leq r_{k,k}$. Then $\sigma_k(\gamma) = \prod_{l=1}^k (\gamma + r_{l,k})$ and $\partial_\gamma \log \sigma_k(\gamma) = \sum_{l=1}^k (\gamma + r_{l,k})^{-1}$, together with $\partial_\gamma^2 \log \sigma_k(\gamma) = -\sum_{l=1}^k (\gamma + r_{l,k})^{-2} < 0$ ($\gamma \rightarrow \sigma_k(\gamma)$ is log-concave).

If $p/\gamma - \sum_{l=1}^k (\gamma + r_{l,k})^{-1} \stackrel{(*)}{=} 0$, then

$$\partial_\gamma^2 \log \mathbf{P}^*(P_k = p) = -p/\gamma^2 + \sum_{l=1}^k (\gamma + r_{l,k})^{-2} < 0,$$

showing that $\hat{\gamma}$ solving $(*)$ is a local maximum and that $\log \mathbf{P}^*(P_k = p)$ has no local minima. So $\hat{\gamma}$ is the maximum likelihood estimator of γ . Even though $\sigma_k(\gamma)$ ($1/\sigma_k(\gamma)$) is a log-concave (respectively log-convex) function of γ , the log-likelihood is a log-concave function of γ leading to the existence of $\hat{\gamma}$. To summarize, there exists a maximum likelihood estimator $\hat{\gamma}$ of γ which is characterized by the implicit equation:

$$P = \hat{\gamma} \frac{\sigma'_k(\hat{\gamma})}{\sigma_k(\hat{\gamma})}.$$

Let us now come to another estimator of γ . If $\sigma_k(\gamma) \in ZR_-$, then by Newton's inequality ([22], p.52)

$$B_{k,p}(\phi_\bullet)^2 \geq B_{k,p-1}(\phi_\bullet) B_{k,p+1}(\phi_\bullet) \left(1 + \frac{1}{p}\right) \left(1 + \frac{1}{k-p}\right) > B_{k,p-1}(\phi_\bullet) B_{k,p+1}(\phi_\bullet).$$

So $B_{k,p}(\phi_\bullet)$ is p -log-concave and by Darroch Theorem, $B_{k,p}(\phi_\bullet)$ is p -unimodal or bimodal at two consecutive p , with mode (maybe up to one unit) equal to $\sigma'_k(1)/\sigma_k(1)$. Because the p -sequence γ^p is also log-concave (and log-convex),

$\gamma^p B_{k,p}(\phi_\bullet)$ is itself p -log-concave and therefore there exists a unique $\tilde{\gamma}$ such that $\frac{\mathbf{P}^*(P_k=p)}{\mathbf{P}^*(P_k=p-1)} = 1$. It is thus defined by

$$\frac{\gamma^p B_{k,p}(\phi_\bullet)}{\gamma^{p-1} B_{k,p-1}(\phi_\bullet)} = 1, \text{ or } \tilde{\gamma} = \frac{B_{k,p-1}(\phi_\bullet)}{B_{k,p}(\phi_\bullet)}.$$

This $\tilde{\gamma}$ is an alternative explicit estimator of γ based on the data k and P .

Taking the expectation with respect to P , we have

$$\begin{aligned} \mathbf{E}^*(\tilde{\gamma}) &= \sum_{p=1}^k \frac{B_{k,p-1}(\phi_\bullet)}{B_{k,p}(\phi_\bullet)} \frac{\gamma^p}{\sigma_k(\gamma)} B_{k,p}(\phi_\bullet) = \gamma \sum_{p=2}^k B_{k,p-1}(\phi_\bullet) \frac{\gamma^{p-1}}{\sigma_k(\gamma)} \\ &= \gamma \sum_{p=1}^{k-1} B_{k,p}(\phi_\bullet) \frac{\gamma^p}{\sigma_k(\gamma)} = \gamma \left(1 - \frac{(\phi_1 \gamma)^k}{\sigma_k(\gamma)} \right) < \gamma. \end{aligned}$$

This shows that $\tilde{\gamma}$ is not an unbiased estimator of γ .

Remark: The estimator $\tilde{\gamma}$ only requires that the sequence $B_{k,p}(\phi_\bullet)$ be p -log-concave and, although sufficient, it is therefore not necessary that $\sigma_k(\gamma) \in ZR_-$; the sequence $B_{k,p}(\phi_\bullet)$ only needs to be a Pòlya frequency sequence of order 2 (so $\sigma_k(\gamma) \in PF_2$) for $\tilde{\gamma}$ to be well-defined. In this spirit, we draw the attention on a result in [3], stating that if the non-null roots of $\sigma_k(\gamma)$ all lie in the angular cone $\phi \in (2\pi/3, 4\pi/3)$ of the complex plane, then $\sigma_k(\gamma)$ has p -log-concave coefficients.

3. SAMPLING FROM DIRICHLET PARTITION: A SPECIAL CASE

We now briefly investigate one particular model of species abundance ξ_n .

• Sampling from a binomial negative sample.

Assume $\phi(x) = -\log(1-x)$, with $\phi_m = (m-1)!$ and let $Z_\theta(x) = (1-x)^{-\theta}$. Thus, with $(\theta)_k := \theta(\theta+1)\dots(\theta+k-1)$ denoting the (rising factorial) Pochhammer symbol, $\sigma_k(\theta) = (\theta)_k$ and ξ is a binomial negative random variable with parameters θ and $1-x$. Note that $\sigma_k(\theta) \in ZR_-$. From (11), the jumps' height δ of ξ is seen to obey a logarithmic series distribution.

When sampling from this discrete species-abundance model $\xi_n = (\xi_1, \dots, \xi_n)$, for instance (12) takes the particular form:

$$(35) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \frac{\mathbf{P}(\xi_1 = k_1, \dots, \xi_n = k_n)}{\mathbf{P}(\zeta_n = k)} = \frac{k!}{(n\theta)_k} \prod_{m=1}^n \frac{(\theta)_{k_m}}{k_m!}.$$

Substituting $(\theta)_k$ to $\sigma_k(\theta)$ in (15) gives its particular expression.

Because $\sigma_{k+1}(\theta) = (k+\theta)\sigma_k(\theta)$, it follows from (3) and (4) that with $S_k(\lambda) = k! [x^k] e^{\lambda((1-x)^{-\theta}-1)}$, $S_{k+1}(\lambda) = (\theta\lambda+k)S_k(\lambda) + \theta\lambda S'_k(\lambda)$. Thus, the Bell coefficients $B_{k,p}(\sigma_\bullet(\theta)) = B_{k,p}((\theta)_\bullet) = [\lambda^p] S_k(\lambda)$, appearing in (16), obey a simple 3-term recurrence [14], [27]

$$B_{k+1,p}((\theta)_\bullet) = \theta B_{k,p-1}((\theta)_\bullet) + (p\theta+k) B_{k,p}((\theta)_\bullet),$$

which should be considered with the boundary conditions

$$B_{k,0}((\theta)_\bullet) = B_{0,p}((\theta)_\bullet) = 0,$$

except for $B_{0,0}((\theta)_\bullet) := 1$. This observation is important because it follows from (16), that, there exist transition probabilities

$$\begin{aligned} \mathbf{P}(P_{n,k+1} = p+1 \mid P_{n,k} = p) &= \frac{(n-p)\theta}{n\theta+k} \text{ and} \\ \mathbf{P}(P_{n,k+1} = p \mid P_{n,k} = p) &= \frac{\sum_{r=1}^p (\theta+k_r)}{n\theta+k} = \frac{p\theta+k}{n\theta+k}. \end{aligned}$$

such that,

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

The first transition probability gives the probability of the event that a new species is discovered given $p < n$ of them were discovered from a previous sample of size $k \geq p$ (the so-called law of succession, [17]) in a size- n population.

Considering the sampling formulae in the $*$ -limit, the expressions (30) and (33) with $\phi_i = (i-1)!$ and $B_{k,p}(\phi_\bullet) = s_{k,p}$ (the absolute first kind Stirling numbers) are the Ewens sampling formulae [18]. Due to $\sigma_{k+1}(\theta) = (k+\theta)\sigma_k(\theta)$, the Bell coefficients $B_{k,p}(\phi_\bullet) = B_{k,p}((\bullet-1)!)$ also obey a 3-term recurrence

$$B_{k+1,p}((\bullet-1)!) = B_{k,p-1}((\bullet-1)!) + kB_{k,p}((\bullet-1)!).$$

• Sampling from a Dirichlet partition of unity in the continuum.

It turns out that this sampling formula can be obtained while following a different path for the sampling procedure:

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

$$(36) \quad f_{S_{1,\theta}, \dots, S_{n,\theta}}(s_1, \dots, s_n) = \frac{\Gamma(n\theta)}{\Gamma(\theta)^n} \prod_{m=1}^n s_m^{\theta-1} \cdot \delta_{(\sum_{m=1}^n s_m - 1)}.$$

Alternatively, with $(\theta)_q := \Gamma(\theta+q)/\Gamma(\theta)$, the law of $\mathbf{S}_n(\theta)$ is characterized by its joint moment function

$$(37) \quad \mathbf{E} \left(\prod_{m=1}^n S_{m,\theta}^{q_m} \right) = \frac{1}{(n\theta)_{\sum_{m=1}^n q_m}} \prod_{m=1}^n (\theta)_{q_m}.$$

We shall put $\mathbf{S}_n(\theta) \stackrel{d}{\sim} D_n(\theta)$ if $\mathbf{S}_n(\theta)$ is Dirichlet distributed with parameter θ . $\mathbf{S}_n(\theta)$ can be obtained while considering $(Y_\theta \stackrel{d}{=} Y_{1,\theta}, \dots, Y_{n,\theta})$, an iid random vector with $Y_\theta \stackrel{d}{\sim} \text{gamma}(\theta)$ and letting $S_{m,\theta} = Y_{m,\theta}/(Y_{1,\theta} + \dots + Y_{n,\theta})$, $m = 1, \dots, n$ (normalizing the $Y_{m,\theta}$ s by their sum). $\mathbf{S}_n(\theta)$ accounts now for a n -species frequency (proportion) model, but now in the continuum. We now come to the sampling procedure from $\mathbf{S}_n(\theta)$.

Let (U_1, \dots, U_k) be k iid uniform throws on the unit interval partitioned according to $\mathbf{S}_n(\theta)$. Let

$$\mathbf{K}_{n,k} := (K_{n,k}(1), \dots, K_{n,k}(n)) \geq 0$$

be an integral-valued random vector which counts the number of visits to the different fragments of $\mathbf{S}_n(\theta)$ in this k -sample. Hence, if M_l is the random fragment label in which the l^{th} trial U_l falls, $K_{n,k}(m) := \sum_{l=1}^k \mathbf{I}(M_l = m)$, $m = 1, \dots, n$.

With $|\mathbf{k}_n| = k$ and $\mathbf{k}_n := (k_1, \dots, k_n) \geq 0$, $\mathbf{K}_{n,k}$ follows the conditional multinomial distribution:

$$(38) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n \mid \mathbf{S}_n(\theta)) = \frac{k!}{\prod_{m=1}^n k_m!} \prod_{m=1}^n S_{m,\theta}^{k_m}.$$

Averaging over $\mathbf{S}_n(\theta)$, we find

$$(39) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \mathbf{E}\mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n \mid \mathbf{S}_n(\theta)) = \frac{k!}{(n\theta)_k} \prod_{m=1}^n \frac{(\theta)_{k_m}}{k_m!},$$

which is the Dirichlet-multinomial distribution, with $\mathbf{E}(K_{n,k}(m)) = k/n$. We shall put $\mathbf{K}_{n,k} \stackrel{d}{\sim} D_{n,k}(\theta)$.

The sampling from $\mathbf{S}_n(\theta) \stackrel{d}{\sim} D_n(\theta)$ formula (39) coincides with the one (35) obtained while sampling from a discrete species abundance model ξ_n with negative binomial distributions. The $*$ -limit of this Dirichlet model is known to lead to the Ewens sampling formulae which are particular incarnation of (30) and (33) with $\phi_i = (i-1)!$ and $B_{k,p}(\phi_\bullet) = s_{k,p}$. See [31] and [32].

It is worthwhile to explore if this remarkable property (or maybe a weaker one) propagates to sampling from other discrete species abundance model.

4. SAMPLING PROBLEMS FROM A SPECIAL GP CLASS

We shall now exhibit a sub-class of GP models whose statistical properties are very similar to the ones developed in the latter Section for the Dirichlet model.

4.1. Sampling from a special GP class. Let us first define the class of ϕ we will be interested in.

The special class \mathcal{S} .

We first recall that a function $h(x)$ defined on some interval $x \in (-\infty, x_0)$ is absolutely monotone on some open interval $I \subseteq (-\infty, x_0)$ if it is C^∞ with $h^{(n)}(x) \geq 0$ for all $n \geq 0$ and $x \in I$.

We shall consider the following special class model

Definition 1. *Suppose that $\phi(x)$ (with $\phi_1 > 0$ and $\phi_m \geq 0$, $m \geq 2$) as from (1), is defined (finite) on the unbounded half-domain $x \in (-\infty, x_0)$ with $0 < x_0 \leq \infty$ and that $\phi'(x)$ is absolutely monotone for all $x \in (-\infty, x_0)$. If this is the case, we shall put $\phi \in \mathcal{S}$. If $\phi \in \mathcal{S}$, $Z_\theta(x) = \exp(\theta\phi(x))$ is also defined on $x \in (-\infty, x_0)$ and absolutely monotone there.*

- Examples of $\phi \in \mathcal{S}$ are x , $e^x - 1$ (Bell), $-\log(1-x)$, $(1-x)^{-\alpha} - 1$, $\alpha > 0$ and $1 - (1-x)^\alpha$, $\alpha \in (0, 1)$.
- Examples of $\phi \notin \mathcal{S}$ are polynomials with positive coefficients $\sum_{l=1}^d c_l x^l$ ($d \geq 2$), $x e^x$, $\sinh(x)$, $\cosh(x) - 1$ and $\tan(x)$. Although the latter ϕ 's can be expanded as in (1) and all have non-negative Taylor coefficients ϕ_m ($\phi_1 > 0$), the corresponding $\phi'(x)$ are not absolutely monotone on $(-\infty, x_0)$ although they are of course on $(0, x_0)$.

Remarks and properties:

- If $\phi \in \mathcal{S}$, so does clearly $\tilde{\phi}(x) := a\phi(bx)$ for all $a, b > 0$. We can check that: $B_{k,p}(\tilde{\phi}_\bullet) = a^p b^k B_{k,p}(\phi_\bullet)$.
- If $\phi^1, \phi^2 \in \mathcal{S}$, then $\phi^1 + \phi^2 \in \mathcal{S}$ and the composition $\phi^1 \circ \phi^2 \in \mathcal{S}$. This allows to produce a lot of new examples of ϕ 's in \mathcal{S} from the ones already introduced. For instance because $\phi^1 = (1-x)^{-\alpha} - 1$ and $\phi^2 = 1 - (1-x)^\alpha$ both belong to \mathcal{S} , would $\alpha \in (0, 1)$, $\phi^1 + \phi^2 = 2 \sinh(-\alpha \log(1-x))$ belongs to \mathcal{S} , together with $\phi^1 \circ \phi^2 = (1-x)^{-\alpha^2} - 1$ and $\phi^2 \circ \phi^1 = 1 - \left(2 - (1-x)^{-\alpha}\right)^\alpha$.
- If $\phi^1, \phi^2 \in \mathcal{S}$, the product $\phi := \phi^1 \cdot \phi^2 \notin$ (in the first place because $\phi_1 = 0$). The Taylor coefficients ϕ_m of ϕ are

$$\phi_m = \sum_{l=1}^{m-1} \binom{m}{l} \phi_l^1 \phi_{m-l}^2 = (\phi^1 * \phi^2)_m, \quad m \geq 2$$

and the ϕ_m do not necessarily form a log-convex sequence, even though ϕ_m^1, ϕ_m^2 , $m \geq 1$, would be log-convex themselves. This is not in contradiction with the Davenport and Pólya theorem [13] stating that the binomial convolution of two log-convex sequences is log-convex because the ϕ^1, ϕ^2 sequences here have no constant terms: $\phi_0^1 = \phi_0^2 = 0$ (resulting in $\phi_1 = 0$). The reason why, when $\phi(x) \in \mathcal{S}$, log-convexity of the sequences $(\phi_m)_{m \geq 1}$ pops in is:

Proposition 9. *When $\phi \in \mathcal{S}$, the function $h(x) := \phi'(-x)$ is completely monotone on the domain $x \in (-x_0, \infty)$, meaning it is C^∞ with $(-1)^n h^{(n)}(x) \geq 0$ for all $n \geq 0$ and $x \in (-x_0, \infty)$. So (from Bernstein theorem [5]), $h(x)$ is the Laplace-Stieltjes transform (LST) of some finite non-negative measure μ on $[0, +\infty)$: $h(x) = \int_0^\infty e^{-xt} d\mu(t)$. We have*

$$h(x) = \sum_{m \geq 0} \frac{\phi_{m+1}}{m!} (-x)^m$$

and so ϕ_{m+1} is the m^{th} moment of $d\mu$, with finite total mass ϕ_1 . By the Cauchy-Schwarz inequality, for all $m \geq 2$, $\phi_{m+1}\phi_{m-1} \geq \phi_m^2$, showing that when $\phi \in \mathcal{S}$, $(\phi_m)_{m \geq 1}$ is a log-convex sequence. Upon shifting, $(\phi_m)_{m \geq 1}$ is the moment sequence of some non-negative measure $d\pi(t) := t^{-1}d\mu(t)$.

Let us now consider $Z_\theta(-x) = e^{\theta\phi(-x)} =: e^{-\theta\psi(x)}$, with

$$\psi(x) := -\phi(-x), \quad x > -x_0.$$

Proposition 10. *When $\phi \in \mathcal{S}$, it holds that $\psi'(x) = h(x) = \int_0^\infty e^{-xt} d\mu(t)$ is completely monotone, so $Z_\theta(-x) = e^{-\theta\psi(x)}$ is the LST of some infinitely divisible random variable (or process) Y_θ on $[0, +\infty)$, whose integral moments are all finite. The function ψ is the Laplace exponent of Y_θ with $\psi(x) = \int_0^\infty (1 - e^{-xt}) d\pi(t)$ for some positive Lévy measure $d\pi(t) = t^{-1}d\mu(t)$, integrating $1 \wedge t$ [40]. Therefore, when $\phi \in \mathcal{S}$,*

$$Z_\theta(-x) = \mathbf{E}(e^{-xY_\theta}) = e^{-\theta\psi(x)} = 1 + \sum_{k \geq 1} \frac{(-x)^k}{k!} \sigma_k(\theta),$$

with $(\sigma_k(\theta), k \geq 0)$ being the Stieltjes moment sequence of Y_θ : $\sigma_k(\theta) = \mathbf{E}(Y_\theta^k)$. Thus, when $\phi \in \mathcal{S}$, for all $\theta > 0$, $(\sigma_k(\theta))_{k \geq 0}$ forms a k -log-convex sequence and for all $k \geq 1$, all $\theta > 0$: $\sigma_{k+1}(\theta) \sigma_{k-1}(\theta) \geq \sigma_k(\theta)^2$.

Since $\mathbf{E}(e^{-x\bar{Y}_{n,\theta}}) = e^{-n\theta\psi(x)}$, $\sigma_k(n\theta)$ is also the k^{th} moment of the sum $\bar{Y}_{n,\theta} := Y_{1,\theta} + \dots + Y_{n,\theta}$ of n iid terms $Y_{m,\theta}$. So, $\sigma_k(n\theta) = \mathbf{E}(\bar{Y}_{n,\theta}^k) = \mathbf{E}(Y_{n\theta}^k)$.

Note finally that taking $Z_\theta(x) = Z_\theta^1(x) Z_\theta^2(x)$ where $Z_\theta^i(x) = e^{\theta\phi_i(x)}$ for two ϕ_i in \mathcal{S} , with $\sigma_k^i(\theta)$ defined by $Z_\theta^i(x) = 1 + \sum_{k \geq 1} \frac{x^k}{k!} \sigma_k^i(\theta)$, two k -log-convex sequences, the sequence $\sigma_k(\theta)$ defined by $Z_\theta(x) = 1 + \sum_{k \geq 1} \frac{x^k}{k!} \sigma_k(\theta)$ obeys

$$\sigma_k(\theta) = \sum_{l=0}^k \binom{k}{l} \sigma_l^1(\theta) \sigma_{k-l}^2(\theta) = (\sigma^1(\theta) * \sigma^2(\theta))_k, \quad k \geq 0,$$

and is k -log-convex by Davenport and Pólya theorem, as a binomial convolution of two log-convex sequences.

Sampling from ξ_n when $\phi \in \mathcal{S}$.

Assume $\phi \in \mathcal{S}$ and consider the sampling problem from ξ_n , where ξ is constructed as in Section 2 from ϕ , but now for $\phi \in \mathcal{S}$. Note that in this case

$$\mathbf{E}(u^\xi) = e^{\theta[\phi(xu) - \phi(x)]} = e^{-\theta[\psi(-xu) - \psi(-x)]}.$$

In a general sampling problem from ξ_n , the joint probability generating function of $\mathbf{K}_{n,k}$ was given by (13). From (12) and making use of $\phi \in \mathcal{S}$, we have

$$(40) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \frac{k!}{\sigma_k(n\theta)} \prod_{m=1}^n \frac{\sigma_{k_m}(\theta)}{k_m!} = \binom{k}{k_1 \dots k_n} \frac{\prod_{m=1}^n \mathbf{E}(Y_{m,\theta}^{k_m})}{\mathbf{E}(\bar{Y}_{n,\theta}^k)},$$

Remark: Because (40) does not depend on the common mean of the $Y_{m,\theta}$'s, we can as well define the reduced random variables with mean 1 : $X_{m,\theta} := Y_{m,\theta}/(\theta\phi_1)$, $m = 1, \dots, n$ and $\bar{X}_{n,\theta} := \sum_{m=1}^n X_{m,\theta}$. Then, with $S_{m,\theta} := X_{m,\theta}/\bar{X}_{n,\theta}$, $m = 1, \dots, n$ defining a random partition $\mathbf{S}_n(\theta) = (S_{1,\theta}, \dots, S_{n,\theta})$ of unity into n id (mean $1/n$) parts

$$(41) \quad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \binom{k}{k_1 \dots k_n} \frac{\prod_{m=1}^n \mathbf{E}(X_{m,\theta}^{k_m})}{\mathbf{E}(\bar{X}_{n,\theta}^k)}$$

$$= \binom{k}{k_1 \dots k_n} \frac{\prod_{m=1}^n \mathbf{E} \left(\overline{X}_{n,\theta}^k S_{m,\theta}^{k_m} \right)}{\mathbf{E} \left(\overline{X}_{n,\theta}^k \right)}$$

as well. The latter expression is identified to an occupancy distribution arising from sampling from the random partition of unity $\mathbf{S}_n(\theta)$ but size-biased by the total length $\overline{X}_{n,\theta}$. In the occupancy distribution (41) indeed, realizations of $(X_{m,\theta})_{m=1}^n$ giving rise to large values of the sum $\overline{X}_{n,\theta}$ are favored, compared to the “neutral” multinomial one say $\mathbf{Q}(\mathbf{K}_{n,k} = \mathbf{k}_n) := \binom{k}{k_1 \dots k_n} \prod_{m=1}^n \mathbf{E} \left(S_{m,\theta}^{k_m} \right)$, based on the same $\mathbf{S}_n(\theta)$.

Whenever $\overline{X}_{n,\theta}$ would be independent of $S_{m,\theta} = X_{m,\theta}/\overline{X}_{n,\theta}$, $m = 1, \dots, n$, (the only possible way to have this is when $\mathbf{S}_n(\theta)$ has Dirichlet(θ) distribution, [23]), this expression boils down to the usual sampling one

$$\mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \binom{k}{k_1 \dots k_n} \mathbf{E} \left(\prod_{m=1}^n (X_{m,\theta}/\overline{X}_{n,\theta})^{k_m} \right) = \mathbf{Q}(\mathbf{K}_{n,k} = \mathbf{k}_n).$$

Alternatively, from (41), the joint pgf of $\mathbf{K}_{n,k}$ also reads

$$\mathbf{E} \left[\prod_{m=1}^n u_m^{K_{n,k}(m)} \right] = \frac{\mathbf{E} \left[(\sum_{m=1}^n u_m X_{m,\theta})^k \right]}{\mathbf{E} \left(\overline{X}_{n,\theta}^k \right)} = \frac{\mathbf{E} \left[\overline{X}_{n,\theta}^k (\sum_{m=1}^n u_m S_{m,\theta})^k \right]}{\mathbf{E} \left(\overline{X}_{n,\theta}^k \right)}.$$

Its computation is thus amenable to the normalized k^{th} moment of the weighted sum $\sum_1^n u_m X_{m,\theta}$ of iid mean 1 infinitely divisible random variables with LST $\mathbf{E} \left(e^{-xX_\theta} \right) = e^{\theta\phi(-x/(\theta\phi_1))} = e^{-\theta\psi(x/(\theta\phi_1))}$ and moments $\mathbf{E} \left(X_\theta^k \right) = \sigma_k(\theta) / (\theta\phi_1)^k$, $k \geq 1$.

Note also that with $\mathbf{k}_p := (k_1, \dots, k_p) \geq 1$ summing to k

$$\begin{aligned} \mathbf{P}(K_{n,k}(1) = k_1, \dots, K_{n,k}(p) = k_p; P_{n,k} = p) \\ = \binom{n}{p} \binom{k}{k_1 \dots k_p} \frac{\prod_{q=1}^p \mathbf{E} \left(\overline{X}_{n,\theta}^k S_{q,\theta}^{k_q} \right)}{\mathbf{E} \left(\overline{X}_{n,\theta}^k \right)}. \end{aligned}$$

is the joint probability that there are $p \in [n]$ non-empty boxes and that (k_1, \dots, k_p) are the respective occupancies of the p filled boxes, labeled in arbitrary order. Again

$$P_{k,p}^{(n)} := \binom{n}{p} \sum_{k_1 + \dots + k_p = k}^* \binom{k}{k_1 \dots k_p} \frac{\prod_{q=1}^p \mathbf{E} \left(\overline{X}_{n,\theta}^k S_{q,\theta}^{k_q} \right)}{\mathbf{E} \left(\overline{X}_{n,\theta}^k \right)}$$

is the probability that in a k -sample from n species with abundance ξ_n in the special class \mathcal{S} , the exact number of distinct visited species is p .

To summarize, we conclude

Proposition 11. *When $\phi \in \mathcal{S}$ and when the discrete species abundance model ξ_n is built on ϕ , its occupancy distribution (12) can alternatively be interpreted as an occupancy distribution (41) arising from sampling from the random partition of*

unity $\mathbf{S}_n(\theta)$ but size-biased by the total length $\bar{X}_{n,\theta}$ appearing in the normalization of $S_{m,\theta} := X_{m,\theta}/\bar{X}_{n,\theta}$. The positive random variable $X_\theta \stackrel{d}{=} X_{1,\theta}$ is infinitely divisible. The correspondence between ξ and (mean 1) X_θ is:

$$\mathbf{E}[u^\xi] = e^{-\theta\phi(x)\left(1 - \frac{\phi(xu)}{\phi(x)}\right)} \text{ and } \mathbf{E}(e^{-xX_\theta}) = e^{\theta\phi(-x/(\theta\phi_1))} = e^{-\theta\psi(x/(\theta\phi_1))}.$$

Note finally that $\psi(x/(\theta\phi_1))$ being the Laplace exponent of X_θ :

$$\mathbf{E}(e^{-xX_\theta}) = e^{-\theta \int_0^\infty (1 - e^{-xt}) d\pi_\theta(t)},$$

where the Lévy measure $d\pi_\theta(t)$ integrates $1 \wedge t$. The measure $d\mu_\theta(t) = td\pi_\theta(t)$ is a finite positive measure with all finite m -moments: $\int_0^\infty t^m d\mu_\theta(t) = \phi_{m+1}/(\theta\phi_1)^{m+1}$, $m \geq 0$. So $\left((\theta\phi_1)^{-m} \phi_m\right)_{m \geq 1}$ is the moment sequence of $d\pi_\theta(t)$.

With $S_{1,\theta} := X_{1,\theta}/\bar{X}_{n,\theta}$, define finally $\mu_k := \mathbf{E}\left[S_{1,\theta}^k\right]$, $k \geq 1$, the sequence of the moments of $S_{1,\theta}$; then $(\mu_k; k \geq 1)$ is a Hausdorff sequence which is completely monotonic in the sense that

$$(42) \quad (-1)^l \Delta^l \mu_k \geq 0 \text{ for each } l, k \geq 0$$

where $\Delta^l \mu_k$ is the l^{th} iterate of the difference operator $\Delta \mu_k := \mu_{k+1} - \mu_k$.

Examples. Examples of admissible $\phi \in \mathcal{S}$ were $-\log(1-x)$, $(1-x)^{-\alpha} - 1$, $\alpha > 0$ and $1 - (1-x)^\alpha$, $\alpha \in (0, 1)$.

The LST $\mathbf{E}(e^{-xX_\theta})$ of X_θ in each case is $(1+x/\theta)^{-\theta}$, $\exp\left[-\theta\left(1 - \left(1 + \frac{x}{\alpha\theta}\right)^{-\alpha}\right)\right]$ and $\exp\left[-\theta\left(\left(1 + \frac{x}{\alpha\theta}\right)^\alpha - 1\right)\right]$ corresponding respectively to a Gamma(θ, θ) distribution, a compound Poisson sum of iid gamma($\alpha, \alpha\theta$) random variables and an exponentially damped stable(θ, α). For this last case, let $\Sigma > 0$ be a stable(θ, α) random variable i.e. with LST $\mathbf{E}(e^{-x\Sigma}) := \exp[-\theta x^\alpha]$, $x \geq 0$. Let f_Σ be its density. Define a random variable Y_θ with damped density $f_{Y_\theta}(t) = \frac{1}{\mathbf{E}(e^{-\Sigma})} e^{-t} f_\Sigma(t)$, $t > 0$. Its LST is $\mathbf{E}(e^{-xY_\theta}) = \mathbf{E}(e^{-(x+1)\Sigma}) / \mathbf{E}(e^{-\Sigma}) = \exp -\theta[(1+x)^\alpha - 1]$. Upon scaling Y_θ , $X_\theta := Y_\theta/(\theta\alpha)$ is mean 1. In the sampling context, the last example was recently considered in ([15], [16], [24] and [25]). They were named the generalized inverse Gaussian or Engen models. \diamond

Remark. in the degenerate case, $\phi(x) = x$, X_θ is purely atomic with $X_\theta \stackrel{d}{\sim} \delta_1$. The LST of X_θ can be obtained from the one of the first gamma(θ, θ) example: $\mathbf{E}(e^{-xX_\theta}) = (1+x/\theta)^{-\theta}$ as $\theta \rightarrow \infty$. In this very particular (admissible) case, $\mathbf{S}_n = (1/n, \dots, 1/n)$ is the uniform deterministic partition of unity (the Maxwell-Boltzmann case). \diamond

4.2. The *-limit. We now come back to the *-limit.

Let $\phi \in \mathcal{S}$. With $\gamma > 0$, let $(Y_\gamma)_{\gamma \geq 0}$ be a subordinator with $Y_0 = 0$ and LST

$$\mathbf{E}(e^{-xY_\gamma}) = e^{-\gamma\psi(x)}, \quad \psi(x) = -\phi(-x).$$

Under assumptions on ϕ , $\mathbf{E}(Y_\gamma) = \gamma\phi_1 < \infty$. Then the Laplace exponent ψ reads

$$(43) \quad \psi(x) = \int_0^\infty (1 - e^{-xt}) d\pi(t),$$

for some positive Lévy measure π on $(0, \infty)$, integrating $1 \wedge t$, [6]. Let $\bar{\pi}(t) := \int_t^\infty d\pi(s)$ be the tail function of π and assume $\bar{\pi}(t) \rightarrow \infty$ as $t \rightarrow 0$ ⁽³⁾. Then

$$(44) \quad Y_\gamma \stackrel{d}{=} \sum_{k \geq 1} \bar{\pi}^{-1}(\Gamma_k/\gamma)$$

where $(\Gamma_k)_{k \geq 1}$ are the points of a standard Poisson Point Process (PPP) on $(0, \infty)$ with intensity 1. The random variables

$$\Delta_{(k)}(\gamma) := \bar{\pi}^{-1}(\Gamma_k/\gamma)$$

with $\Delta_{(1)}(\gamma) \geq \Delta_{(2)}(\gamma) \geq \dots$ constitute the ranked jumps' heights of the subordinator Y_γ (they are infinitely many, with 0 as a limit point). They form a PPP on the half-line with intensity $\gamma d\pi(t)$, and the law of $\Delta_{(k)}(\gamma)$ can easily be computed to be

$$(45) \quad \mathbf{P}(\Delta_{(k)}(\gamma) \in dt) = \frac{\gamma^k \bar{\pi}(t)^{k-1}}{(k-1)!} e^{-\gamma \bar{\pi}(t)} d\pi(t).$$

By Campbell formula indeed (see [34], [32]), for all measurable function g for which $\int_0^\infty (1 - e^{-xg(t)}) d\pi(t) < \infty$, we have

$$\mathbf{E} \left(\exp \left\{ -x \sum_{k \geq 1} g(\bar{\pi}^{-1}(\Gamma_k/\gamma)) \right\} \right) = \exp \left\{ -\gamma \int_0^\infty (1 - e^{-xg(t)}) d\pi(t) \right\}.$$

Putting $g(t) = t$, $\mathbf{E}(e^{-xY_\gamma}) = e^{-\gamma\psi(x)}$, as claimed.

From the above construction, we can define a random distribution on the infinite-dimensional 1-simplex by normalizing the ranked jumps' heights of Y_γ by itself. Consider again Y_γ and, with $\theta := \gamma/n$, define $Y_{m,\theta} := Y_{m\theta} - Y_{(m-1)\theta}$, $m = 1, \dots, n$ which are mutually independent. Then, $\bar{Y}_{n,\theta} := \sum_{m=1}^n Y_{m,\theta} = Y_{n\theta} = Y_\gamma$. If we rank the $Y_{m,\theta}$'s, with $Y_{(1),\theta} \geq \dots \geq Y_{(n),\theta}$ ⁽⁴⁾, then, [30], as $n \rightarrow \infty$, $\theta \rightarrow 0$, $n\theta = \gamma$

$$(46) \quad (Y_{(1),\theta}, \dots, Y_{(n),\theta}, 0, 0, \dots) \xrightarrow[*]{d} (\Delta_{(1)}(\gamma), \Delta_{(2)}(\gamma), \dots).$$

Normalizing,

$$(Y_{(1),\theta}/Y_\gamma, \dots, Y_{(n),\theta}/Y_\gamma, 0, 0, \dots) \xrightarrow[*]{d}$$

$$(47) \quad (\Delta_{(1)}(\gamma)/Y_\gamma, \Delta_{(2)}(\gamma)/Y_\gamma, \dots) =: \mathbf{S}_\infty(\gamma) := (S_{(1),\gamma}, S_{(2),\gamma}, \dots),$$

with $\mathbf{S}_\infty(\gamma)$ defining a random partition of unity with infinitely many (ordered) pieces.

³If $\bar{\pi}$ has a finite limit, the random partition of unity defined in (47) is finite with a random Poisson number of pieces (see Example (iii) below). This case deserves a special treatment.

⁴If Y_θ has a density (π has no atom), these inequalities are strict.

If $t > 0$ is some (small) cutoff or threshold value, let $N_+(t) := \sum_{k \geq 1} \mathbf{I}(\Delta_{(k)}(\gamma) > t)$ count the numbers of atoms of the partition of Y_γ exceeding t . By Campbell formula

$$\begin{aligned} \mathbf{E}(\exp\{-xN_+(t)\}) &= \exp\left\{-\gamma \int_0^\infty (1 - e^{-x\mathbf{I}(s>t)}) d\pi(s)\right\} \\ (48) \qquad \qquad \qquad &= \exp\{-\gamma\bar{\pi}(t)(1 - e^{-x})\} \end{aligned}$$

is the full LST of $N_+(t)$. This shows that $N_+(t)$ is Poisson distributed with mean $\gamma\bar{\pi}(t)$. Recalling $\bar{\pi}(t) \xrightarrow[t \rightarrow 0]{} \infty$, the law of large numbers gives

$$(49) \qquad \qquad \qquad N_+(t) / \bar{\pi}(t) \xrightarrow{a.s.} \gamma, \text{ as } t \rightarrow 0.$$

The fact that $N_+(t)$ is Poisson may be also checked as follows. We have $N_+(t) = \inf(k \geq 1 : \Delta_{(k)}(\gamma) \leq t) - 1$ and $\mathbf{P}(N_+(t) \geq k) = \mathbf{P}(\Delta_{(k)}(\gamma) > t) = \mathbf{P}(\Gamma_k \leq \gamma\bar{\pi}(t)) = e^{-\gamma\bar{\pi}(t)} \sum_{l \geq k} \frac{[\gamma\bar{\pi}(t)]^l}{l!}$. So $N_+(t)$ is Poisson with mean $\gamma\bar{\pi}(t)$.

Because also, by the strong law of large numbers, $\Gamma_k/k \rightarrow 1$ a.s. as $k \rightarrow \infty$, recalling $\Gamma_k = \gamma\bar{\pi}(Y_\gamma S_{(k),\gamma})$, we get

$$\gamma\bar{\pi}(Y_\gamma S_{(k),\gamma}) / k \rightarrow 1 \text{ a.s. as } k \rightarrow \infty.$$

From the behavior of $\bar{\pi}(t)$ near $t = 0$, the decay rate of $S_{(k),\gamma}$ to 0 as $k \rightarrow \infty$ follows.

Sampling from $S_{m,\theta} := Y_{m,\theta}/Y_\gamma$, $m = 1, \dots, n$. Define as in (40) a size-biased (SB) sampling procedure for which $(|\mathbf{k}_n| = k)$

$$(50) \qquad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \binom{k}{k_1 \dots k_n} \frac{\mathbf{E}(Y_\gamma^k \prod_{m=1}^n S_{m,\theta}^{k_m})}{\mathbf{E}(Y_\gamma^k)}.$$

Recall that this SB procedure is not the standard sampling from a k uniform throw on $S_{m,\theta}$, $m = 1, \dots, n$, obtained while counting the number of uniform hits within each $S_{m,\theta}$. Indeed, would the latter sampling model hold, instead of (50), one would rather expect the multinomial occupancy distribution

$$\mathbf{Q}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \binom{k}{k_1 \dots k_n} \mathbf{E}\left(\prod_{m=1}^n (Y_{m,\theta}/Y_\gamma)^{k_m}\right),$$

and in general, we have $\mathbf{Q}(\mathbf{K}_{n,k} = \mathbf{k}_n) \neq \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n)$

However, as from the usual size-biased sampling point of view, see [2] for example, we have

$$(51) \qquad \mathbf{P}(\mathbf{K}_{n,k} = \mathbf{k}_n) = \frac{\mathbf{E}_{\mathbf{Q}}(Y_\gamma^k \cdot I(\mathbf{K}_{n,k} = \mathbf{k}_n))}{\mathbf{E}_{\mathbf{Q}}(Y_\gamma^k)},$$

consistently with size-biasing $\mathbf{K}_{n,k} \stackrel{d}{\sim} \mathbf{Q}$ on the total length Y_γ .

According to (50), the joint pgf of $\mathbf{K}_{n,k}$ is

$$\mathbf{E}\left(\prod_{m=1}^n u_m^{K_{n,k}(m)}\right) = \frac{1}{\mathbf{E}(Y_\gamma^k)} \sum_{\substack{k_1, \dots, k_n \geq 0 \\ |\mathbf{k}_n| := k_1 + \dots + k_n = k}} \binom{k}{k_1 \dots k_n} \prod_{m=1}^n u_m^{k_m} \mathbf{E}\left(\prod_{m=1}^n Y_{m,\theta}^{k_m}\right)$$

$$(52) \quad = \frac{\mathbf{E} \left[\left(\sum_{m=1}^n u_m Y_{m,\theta} \right)^k \right]}{\mathbf{E} \left(Y_\gamma^k \right)},$$

which is akin to (40).

Finally, by symmetry or exchangeability, SB sampling from $S_{(m),\theta} := Y_{(m),\theta}/Y_\gamma$, $m = 1, \dots, n$ can similarly be defined by

$$\mathbf{E} \left(\prod_{m=1}^n u_{(m)}^{K_{n,k}((m))} \right) = \frac{\mathbf{E} \left[\left(\sum_{m=1}^n u_{(m)} Y_{(m),\theta} \right)^k \right]}{\mathbf{E} \left(Y_\gamma^k \right)},$$

where in the latter formula, $K_{n,k}((m))$, $m = 1, \dots, n$ are now the relabeled occupation numbers of the boxes with sizes $S_{(m),\theta}$ arranged in decreasing order.

Proposition 12. *The distribution of $K_{n,k}((m))$, $m = 1, \dots, n$ is exchangeable.*

Proof: With π the random permutation transforming Y_m into $Y_{(m)}$, $m = 1, \dots, n$, due to the exchangeability of the $Y_{m,\theta}$ s, for all deterministic permutation σ of $\{1, \dots, n\}$:

$$\sum_{m=1}^n u_{\pi_m} Y_{\pi_m,\theta} = \sum_{m=1}^n u_m Y_{m,\theta} \stackrel{d}{\sim} \sum_{m=1}^n u_m Y_{\sigma_m,\theta} = \sum_{m=1}^n u_{\pi_m} Y_{\sigma \circ \pi_m,\theta}$$

and so $\mathbf{E} \left(\prod_{m=1}^n u_{(m)}^{K_{n,k}((m))} \right)$ is a symmetric function of the $u_{(m)} = u_{\pi_m}$. \diamond

From these considerations, we can state the following result:

Proposition 13. *Let $\gamma = n\theta$. When $\phi \in \mathcal{S}$, with $(\sigma_k(\theta), k \geq 0)$ the Stieltjes moment sequence of some infinitely divisible subordinator Y_γ with Laplace exponent $\psi(x) = -\phi(-x)$, the occupancy distributions (12), (15) and (24) are size-biased sampling distributions from $S_{m,\theta} := Y_{m,\theta}/Y_\gamma$, $m = 1, \dots, n$ as defined by (50) or (51). They are also the joint distribution of the occupation numbers $K_{n,k}((m))$, $m = 1, \dots, n$ obtained by SB sampling from $S_{(m),\theta} = Y_{(m),\theta}/Y_\gamma$, $m = 1, \dots, n$.*

Corollary 14. *When $\phi \in \mathcal{S}$ and π has infinite mass ($\phi(x) \xrightarrow{x \rightarrow -\infty} -\infty$), the occupancy distributions (27), (31) and (32) are size-biased sampling distributions from $\mathbf{S}_\infty(\gamma) = (S_{(1),\gamma}, S_{(2),\gamma}, \dots)$ defined in (47).*

Proof: The proof follows from the previous Proposition, the fact that (27) and (32) were obtained as weak $*$ -limits of (15) and (24), from (47) and from symmetry. \diamond

For instance, from (27)

$$\begin{aligned} \mathbf{P}^* (K_k((1)) = k_1, \dots, K_k((p)) = k_p; P_k = p) &= \frac{k!}{p!} \frac{\gamma^p}{\sigma_k(\gamma)} \prod_{q=1}^p \frac{\phi_{k_q}}{k_q!} \\ &= \frac{\mathbf{E}^* \left(Y_\gamma^k \sum_{m_1 < \dots < m_p} \prod_{q=1}^p S_{(m_q),\gamma}^{k_q} \right)}{\mathbf{E} \left(Y_\gamma^k \right)} \end{aligned}$$

is the probability that there are p visited species in the k -sample, each visited k_q times, and that they were obtained after SB sampling from $S_{(m_1),\gamma} > \dots > S_{(m_p),\gamma}$ for any ordered sequence $m_1 < \dots < m_p$.

In particular, the probability that, in a size-biased sampling procedure from $\mathbf{S}_\infty(\gamma)$, all elements of the k -sample are of the same species (whichever species it can be) is thus

$$(53) \quad \frac{\mathbf{E}^* \left(Y_\gamma^k \sum_{m \geq 1} S_{(m),\gamma}^k \right)}{\mathbf{E} \left(Y_\gamma^k \right)} = \gamma \frac{\phi_k}{\sigma_k(\gamma)} = \mathbf{E}^* \left(A_k(k) \right).$$

This identity also follows from (31) with $a_1 = \dots = a_{k-1} = 0$, $a_k = 1$ and $p = 1$ (only one species visited k times).

We observe that, as $\gamma \rightarrow 0$ (or $\mathbf{E}^*(Y_\gamma) \rightarrow 0$ as well), due to $\sigma_k(\gamma) \sim \gamma \phi_k$, this probability tends to 1, showing that γ itself may be viewed as some temperature parameter for the population with infinitely many species: the smaller γ , the larger the probability is that any k -sample visits a single one species (among which the one with largest frequency $S_{(1),\gamma}$).

Similarly, the probability that all elements of the k -sample reveal only two species (whichever species they can be) is

$$\sum_{l=1}^{k-1} \frac{\mathbf{E}^* \left(Y_\gamma^k \sum_{m_1 < m_2}^* S_{(m_1),\gamma}^l S_{(m_2),\gamma}^{k-l} \right)}{\mathbf{E}^* \left(Y_\gamma^k \right)} = \frac{1}{2} \frac{\gamma^2 k!}{\sigma_k(\gamma)} \sum_{l=1}^{k-1} \frac{\phi_l}{l!} \frac{\phi_{k-l}}{(k-l)!} = \frac{\gamma^2}{\sigma_k(\gamma)} B_{k,2}(\phi_\bullet).$$

This identity follows from (31) with $a_l = 1$, $a_{k-l} = 1$, $a_j = 0$ if $j \neq \{l, k-l\}$ and $p = 2$ (only two species visited, one l times and the other one $k-l$ times), summing on $l = 1, \dots, k-1$ and from $\phi_k^{*2} = 2B_{k,2}(\phi_\bullet)$. More generally, if $p \leq k$, $\frac{\gamma^p}{\sigma_k(\gamma)} B_{k,p}(\phi_\bullet)$ is the probability that all elements of the k -sample reveal p distinct species (consistently with (28)), $\frac{(\gamma \phi_1)^k}{\sigma_k(\gamma)}$ the probability that all species in the k -sample are of distinct types. When γ is small this latter probability is polynomially small $\sim \gamma^{k-1}$.

Finally, the probability that only one species is visited by the k -sample and that it is the m^{th} more abundant one is

$$(54) \quad \frac{\mathbf{E}^* \left(Y_\gamma^k S_{(m),\gamma}^k \right)}{\mathbf{E}^* \left(Y_\gamma^k \right)} = \frac{\mathbf{E} \left(\Delta_{(m)}(\gamma)^k \right)}{\mathbf{E} \left(Y_\gamma^k \right)} = \frac{1}{(m-1)!} \frac{\int_0^\infty e^{-x} x^{m-1} \bar{\pi}^{-1}(x/\gamma)^k dx}{\sigma_k(\gamma)}.$$

$$= \frac{\gamma}{\sigma_k(\gamma)} \frac{1}{(m-1)!} \int_0^\infty t^k (\gamma \bar{\pi}(t))^{m-1} e^{-\gamma \bar{\pi}(t)} d\pi(t),$$

consistently with (45). Summing (54) over $m \geq 1$, we recover from (53), that $\phi_k = \frac{1}{\gamma} \int_0^\infty \bar{\pi}^{-1}(x/\gamma)^k dx = \int_0^\infty t^k d\pi$ is the k^{th} moment of the Lévy measure $d\pi$. In particular, the probability that only one species is visited by the k -sample and that it is the more abundant one is (compare with (53))

$$\frac{\mathbf{E}^* \left(Y_\gamma^k S_{(1),\gamma}^k \right)}{\mathbf{E}^* \left(Y_\gamma^k \right)} = \frac{\gamma}{\sigma_k(\gamma)} \int_0^\infty t^k e^{-\gamma \bar{\pi}(t)} d\pi(t)$$

$$= \frac{\gamma \phi_k}{\sigma_k(\gamma)} \left[1 - \frac{1}{\phi_k} \int_0^\infty t^k \left(1 - e^{-\gamma \bar{\pi}(t)} \right) d\pi(t) \right].$$

When γ gets very small, this probability approaches 1 from below, up to an $O(\gamma)$ residual term: again, $S_{(1),\gamma}$ dominates the other smaller $S_{(m),\gamma}$ and for small values of the biodiversity parameter γ therefore, the species frequencies $S_{(m),\gamma}$; $m \geq 1$

turn out to be very disparate.

5. EXAMPLES

Let us supply some Examples illustrating our results.

(i) Take the Fisher logarithmic series model $\phi(x) = -\log(1-x) \in \mathcal{S}$, resulting in ξ obeying a negative binomial distribution with parameters $\theta > 0$ and $1-x \in (0, 1)$, [20]. Here $\phi_\bullet = (\bullet - 1)!$. Then Y_γ is a Moran subordinator with Lévy-measure: $d\pi(t) = t^{-1}e^{-t}dt$. The Laplace exponent of Y_γ is $\psi(x) = \log(1+x)$, in accordance with $\psi(x) = -\phi(-x)$. In that particular case, $(S_{(1),\gamma}, S_{(2),\gamma}, \dots) \sim PD(0, \gamma)$, a Poisson-Dirichlet partition with parameter γ , [23]. Because, due to well-known properties of Gamma-distributed random variables, Y_γ is independent of $S_{m,\theta} = Y_{m,\theta}/Y_\gamma$, $m = 1, \dots, n$, the size-biased sampling distributions from $(S_{1,\theta}, \dots, S_{n,\theta})$ corresponds to the usual multinomial one. In this well-known model for species frequency, $\sigma_k(\theta) = (\theta)_k$. So $\sigma_k(\theta) \in ZR_-$.

Because $\bar{\pi}(t) \sim -\log t$ as $t \rightarrow 0$, $N_+(t) := \#\{k : \Delta_{(k)}(\gamma) > t\}$ grows like $-\gamma \log t$ as $t \rightarrow 0$. Besides,

$$-\log S_{(k),\gamma} \sim k/\gamma \text{ as } k \rightarrow \infty$$

and the ordered frequencies decay exponentially fast with k : species with small frequency get exponentially rare.

Assuming θ known, the Maximum Likelihood Estimator (MLE) estimator of n in the finitely many species model is given implicitly by $P = \hat{n} \left(1 - \frac{\sigma_k((\hat{n}-1)\theta)}{\sigma_k(\hat{n}\theta)}\right)$, so here

$$P = \hat{n} \left(1 - \frac{((\hat{n}-1)\theta)_k}{(\hat{n}\theta)_k}\right).$$

When $\theta = 1$, this estimator is explicitly given by

$$\hat{n} = \frac{(k-1)P}{k-P},$$

where, as conventional wisdom suggests, \hat{n} will be large when the difference between $1/P$ and $1/k$ is small (new species are being frequently discovered). The MLE estimator of γ in the infinitely many species model is given implicitly by $P = \hat{\gamma} \frac{\sigma_k(\hat{\gamma})}{\sigma_k(\hat{\gamma})}$, [41], so here

$$P = \sum_{l=0}^{k-1} \frac{\hat{\gamma}}{\hat{\gamma} + l}.$$

The estimator $\hat{\gamma}$ is biased but its bias decreases as k grows. The alternative estimator $\tilde{\gamma} = \frac{B_{k,P-1}(\phi_\bullet)}{B_{k,P}(\phi_\bullet)}$ with $B_{k,p}(\phi_\bullet) = s_{k,p}$ is also biased and can be computed using the recursion for third kind Stirling numbers

$$B_{k+1,p}((\bullet-1)!) = B_{k,p-1}((\bullet-1)!) + kB_{k,p}((\bullet-1)!).$$

(ii) The full two-parameters $PD(\alpha, \gamma)$ defined in [37] can be obtained while subordinating the damped α -stable subordinator (see (iii) below) to an independent

Moran one with parameter γ/α . And considering the normalized ranked sizes of the subordinate jumps: here, independently of this partition of unity, Y_γ again is gamma(γ) distributed. As shown in [37], $PD(\alpha, \gamma)$ has many interesting properties, [35]. This partition of unity leads to a generalized (unbiased) Ewens' sampling formula called Pitman's sampling formula, [36].

(iii) Take $\phi(x) = (1-x)^{-\alpha} - 1 \in \mathcal{S}$ where $\alpha > 0$. Here $\phi_\bullet = (\alpha)_\bullet$ resulting in ξ being a Poisson sum of negative binomial increments δ . The Lévy-measure corresponding to Y_γ is the Gamma($\alpha, 1$) probability density: $d\pi(t) = 1/\Gamma(\alpha) \cdot t^{\alpha-1} e^{-t} dt$. The Laplace exponent of Y_γ is $\psi(x) = 1 - (1+x)^{-\alpha}$, in accordance with $\psi(x) = -\phi(-x)$. Because π is integrable with mass 1, Y_γ is a subordinator in the compound Poisson class (a Poisson(γ) sum of iid positive jumps with Gamma($\alpha, 1$) density). For this reason

$$Y_\gamma \stackrel{d}{=} \left[\sum_{k=1}^{P(\gamma)} \bar{\pi}^{-1}(\Gamma_k) \right] \cdot \mathbf{I}(P(\gamma) \geq 1) + 0 \cdot \mathbf{I}(P(\gamma) = 0),$$

where $(\Gamma_k)_{k \geq 1}$ are the points of a standard PPP on $(0, \infty)$ with intensity 1, independent of $P(\gamma)$ which is Poisson(γ) distributed. The random variables

$$\Delta_{(k)}(\gamma) := \bar{\pi}^{-1}(\Gamma_k); \quad k = 1, \dots, P(\gamma)$$

with $\Delta_{(1)}(\gamma) \geq \dots \geq \Delta_{(P(\gamma))}(\gamma)$ constitute the ranked jumps' heights of the subordinator Y_γ (they are here finitely many); normalizing with Y_γ , size-bias sampling is therefore from a finite random partition of unity. Note that when π is integrable, the biodiversity parameter γ interprets directly as the expected number of species in the population.

We first recall that for $\phi_\bullet = (\alpha)_\bullet$ as in the case study

$$B_{k+1,p}(\phi_\bullet) = \alpha B_{k,p-1}(\phi_\bullet) + (k + p\alpha) B_{k,p}(\phi_\bullet).$$

When $\alpha = 1$, $B_{k,p}(\bullet!) = \binom{k-1}{p-1} \frac{k!}{p!}$ are the Lah numbers.

Recalling also $\mathbf{P}^*(P_k = p) = \frac{\gamma^p}{\sigma_k(\gamma)} B_{k,p}(\phi_\bullet)$, we get the recursion

$$\mathbf{P}^*(P_{k+1} = p) = \frac{\gamma^p}{\sigma_{k+1}(\gamma)} (\alpha B_{k,p-1}(\phi_\bullet) + (k + p\alpha) B_{k,p}(\phi_\bullet)) =$$

$$\frac{\sigma_k(\gamma)}{\sigma_{k+1}(\gamma)} (\alpha \gamma \mathbf{P}^*(P_k = p-1) + (k + p\alpha) \mathbf{P}^*(P_k = p)).$$

This shows that the event $P_{k+1} = p$ only depends on the event $P_k = p-1$ (respectively $P_k = p$), when a new species (respectively no new species) is being discovered as the sample size is increased by one unit. And not on further past events such as $P_l = p-1$ for $p-1 \leq l < k$. The transition rates are $\lambda_{p,p+1} = \alpha \gamma \frac{\sigma_k(\gamma)}{\sigma_{k+1}(\gamma)}$ (independent of p but dependent on k) and $\lambda_{p,p} = (k + p\alpha) \frac{\sigma_k(\gamma)}{\sigma_{k+1}(\gamma)}$. $\lambda_{p,p+1}$ is the rate at which a new species is being discovered given p of them were previously discovered in a size- k sample. This suggests an underlying sequential urn scheme, [7], [41].

The estimator $\tilde{\gamma} = \frac{B_{k,P-1}(\phi_\bullet)}{B_{k,P}(\phi_\bullet)}$ of γ can easily be evaluated numerically thanks to the three-term recurrence which $B_{k,p}(\phi_\bullet)$ fulfills. When $\alpha = 1$, it is

$$\tilde{\gamma} = \frac{P(P-1)}{k-P+1} = \frac{P}{k} \frac{1}{\frac{1}{P-1} - \frac{1}{k}}.$$

For the four following examples, an appeal to length-biased sampling distributions from $\mathbf{S}_\infty(\gamma)$ is required.

(iv) With $\alpha \in (0, 1)$, take $\phi(x) = 1 - (1-x)^\alpha \in \mathcal{S}$, resulting in ξ being a Poisson sum of extended negative binomial increments δ (also called a Poisson-Pascal random variable). Here $\phi_1 = \alpha$, $\phi_m = \alpha(1-\alpha)_{m-1}$, $m \geq 1$ and the weight of large clusters is smaller than in Example (i) where $\phi_m = (m-1)!$. We therefore expect small clusters sizes to be enhanced. In this case, Y_γ is a damped α -stable subordinator with Lévy-measure: $d\pi(t) = \alpha/\Gamma(1-\alpha) \cdot t^{-(\alpha+1)}e^{-t}dt$. The Laplace exponent of Y_γ is $\psi(x) = (1+x)^\alpha - 1$, in accordance with $\psi(x) = -\phi(-x)$.

Because $\bar{\pi}(t) \sim 1/\Gamma(1-\alpha) \cdot t^{-\alpha}$ as $t \rightarrow 0$, $N_+(t) := \#\{k : \Delta_{(k)}(\gamma) > t\}$ grows like $\gamma/\Gamma(1-\alpha) \cdot t^{-\alpha}$ as $t \rightarrow 0$. Besides,

$$S_{(k),\gamma} \sim \left(\frac{\gamma}{\Gamma(1-\alpha)} \right)^{1/\alpha} Y_\gamma^{-1} k^{-1/\alpha} \text{ as } k \rightarrow \infty$$

and the ordered frequencies only decay algebraically fast with k . Species with small frequency are long-tailed (there are many small size groups or rare species in the Engen model, compared to the Ewens model).

In this model, $\phi_\bullet = \alpha(1-\alpha)_{\bullet-1}$. Because $\phi_1 = \alpha$ and $\phi_{m+1} = \phi_m(m-\alpha)$, $m \geq 1$, it follows from (3, 4) that $\sigma_{k+1}(\theta) = (\theta\alpha + k)\sigma_k(\theta) - \theta\alpha\sigma'_k(\theta)$. Thus, the Bell coefficients $B_{k,p}(\phi_\bullet)$, appearing in (16), again obey a simple 3-term recurrence

$$B_{k+1,p}(\phi_\bullet) = \alpha B_{k,p-1}(\phi_\bullet) + (k-p\alpha) B_{k,p}(\phi_\bullet).$$

They constitute generalized Stirling numbers studied by [9]. It can be checked that $\sigma_k(\theta) \notin ZR_-$.

This model is amenable to similar conclusions as the ones from the previous example with recursion now given by

$$\mathbf{P}^*(P_{k+1} = p) = \frac{\sigma_k(\gamma)}{\sigma_{k+1}(\gamma)} (\alpha\gamma\mathbf{P}^*(P_k = p-1) + (k-p\alpha)\mathbf{P}^*(P_k = p)).$$

Equation (32) with $\phi_\bullet = \alpha(1-\alpha)_{\bullet-1}$ is the Engen's extended negative binomial sampling formula [24]. The particular case $\alpha = 1/2$ is studied in [25]. The micro-canonical distribution (33) coincides when $\phi_\bullet = \alpha(1-\alpha)_{\bullet-1}$ with the one occurring in the Pitman sampling formula ([24], Remark 3).

(v) Let $\phi(x)$ solve the functional equation $\phi(x) = x \exp \phi(x)$. Then $\phi(x) = \sum_{m \geq 1} \frac{\phi_m}{m!} x^m$ with $\phi_m = m^{m-1}$ is the Cayley generating function appearing in the enumeration of rooted labeled trees with m nodes. The convergence radius of this series is $x_0 = e^{-1}$ with $\phi(x_0) = 1$ and $\phi'(x_0) = \infty$. Clearly ϕ_m is log-convex, it is a Stieltjes moment sequence and $\phi \in \mathcal{S}$. The associated Laplace exponent

$\psi(x) = -\phi(-x)$ is the Lambert function. Because $\psi(x) \sim \log x$ as $x \rightarrow \infty$, $\bar{\pi}(t) \sim -\log t$ as $t \rightarrow 0$ and $N_+(t) := \#\{k : \Delta_{(k)}(\gamma) > t\}$ grows like $-\gamma \log t$ as $t \rightarrow 0$. Besides, like in Example (i)

$$-\log S_{(k),\gamma} \sim k/\gamma \text{ as } k \rightarrow \infty.$$

The partition function $Z_\theta(x) = \exp \theta \phi(x)$ occurs in the enumeration of forests of Cayley trees. The Bell coefficients are $B_{k,p}(\phi_\bullet) = \binom{k-1}{p-1} k^{k-p}$ (number of forests with k nodes and p trees) in accordance with the global weights $\sigma_k(\theta) = \theta(k+\theta)^{k-1}$. So $\sigma_k(\theta) \in ZR_-$. Assuming θ known, the MLE estimator of n in the finitely many species model is given implicitly by $P = \hat{n} \left(1 - \frac{\sigma_k((\hat{n}-1)\theta)}{\sigma_k(\hat{n}\theta)}\right)$, so here

$$P = \hat{n} - (\hat{n} - 1) \left(1 - \frac{\theta}{k + \hat{n}\theta}\right)^{k-1}.$$

The MLE estimator of γ in the infinitely many species model is given by $P = \hat{\gamma} \frac{\sigma'_k(\hat{\gamma})}{\sigma_k(\hat{\gamma})}$, so here explicit

$$\hat{\gamma} = \frac{k(P-1)}{k-P}.$$

The alternative (biased) estimator is $\tilde{\gamma} = \frac{B_{k,P-1}(\phi_\bullet)}{B_{k,P}(\phi_\bullet)}$. Thus

$$\tilde{\gamma} = \frac{k(P-1)}{k-P+1} = \frac{1}{\frac{1}{P-1} - \frac{1}{k}};$$

it is also explicit and very close to $\hat{\gamma}$.

(vi) As a next example, let $\phi(x)$ solve the functional equation $\phi(x) = xg(\phi(x))$ where $g(x) = (1+bx)^a$ with either $b > 0$ and $a \geq 1$ or a and b both negative. $\phi(x)$ is the generating function appearing in the enumeration of rooted trees when the generating function g of the offspring is either (generalized) binomial or negative binomial. Then $\phi_m = (m-1)! \binom{am}{m-1} b^{m-1}$ are non-negative numbers. We conjecture that $\phi \in \mathcal{S}$. It holds [28] that $x_0 = (ab)^{-1} (1-1/a)^{a-1}$ with $\phi(x_0) = 1/(b(a-1))$ and $\phi'(x_0) = \infty$. For this tree model first discussed in [4], the Lagrange inversion formula gives [1]

$$B_{k,p}(\phi_\bullet) = \binom{k-1}{p-1} \{ak\}_{k-p} b^{k-p},$$

where $\{a\}_l := a(a-1)\dots(a-l+1)$. Recalling $\tilde{\gamma} = \frac{B_{k,P-1}(\phi_\bullet)}{B_{k,P}(\phi_\bullet)}$, we get

$$\tilde{\gamma} = \frac{b(P-1)}{k-P+1} ((a-1)k+P) = \frac{b}{\frac{1}{P-1} - \frac{1}{k}} \left(a-1 + \frac{P}{k}\right),$$

which is explicit. Again, would $1/k$ be close to $1/(P-1)$, then $\tilde{\gamma}$ would be estimated to be large. Would $a \rightarrow \pm\infty$, $b \rightarrow \pm 0$ while $ab \rightarrow 1$, we recover the results just obtained for Cayley trees (consistently with $g(x) = (1+bx)^a \rightarrow e^x$). If $a = b = 1$, we recover Example (iii) with $\alpha = 1$. When k is large, the minimum of $B_{k,p}^2(\phi_\bullet) / (B_{k,p-1}(\phi_\bullet) B_{k,p+1}(\phi_\bullet))$ is attained when $p = [\lambda k]$ for some $\lambda \in (0, 1)$, with value

$$\frac{\lambda}{1-\lambda} \frac{(1-\lambda)k+1}{\lambda k-1} \frac{(a-1+\lambda)k+1}{(a-1+\lambda)k} \xrightarrow{k \rightarrow \infty} 1$$

and the sequence $B_{k,p}(\phi_\bullet)$ is p -log-concave.

(vii) Let $\alpha > 0$ and let $\phi(x) = \sum_{m \geq 1} m^{-\alpha} x^m$ be the polylog function. The convergence radius of this series is $x_0 = 1$ with $\phi(x_0) < \infty$ iff $\alpha > 1$ and $\phi'(x_0) < \infty$ iff $\alpha > 2$. $\phi(x)$ is defined for $x < x_0$ and $\phi(x) \rightarrow -\infty$ as $x \rightarrow -\infty$. We have $\phi_m = m!m^{-\alpha}$ and $(\phi_m)_{m \geq 1}$ constitutes a log-convex sequence because for all $m \geq 2$,

$$\begin{aligned} \phi_{m+1}\phi_{m-1} &= (m+1)!(m-1)!(m^2-1)^{-\alpha} \\ &> (m+1)!(m-1)!m^{-2\alpha} > m!^2m^{-2\alpha} = \phi_m^2. \end{aligned}$$

The sequence ϕ_m satisfies Carleman's condition $\sum_{m \geq 1} \phi_m^{-1/(2m)} = \infty$. Thus $\phi \in \mathcal{S}$ and $\psi(x) = -\phi(-x)$, $x > -1$, is the Laplace exponent of some polylog subordinator with Lévy measure $d\pi$. Because $\phi(x) \sim -[\log(-x)]^\alpha / \Gamma(1+\alpha)$ as $x \rightarrow -\infty$, [11], $-\phi(-x) =: \psi(x) \rightarrow \infty$ as $x \rightarrow \infty$ and π has infinite total mass. In this example, when $\alpha > 1$, the weight of large clusters ϕ_m is smaller than in Example (i) where $\phi_m = (m-1)!$. When $\alpha > 1$, we therefore expect small clusters sizes to be enhanced as in Example (iv), but to a lesser extent. Because indeed $\bar{\pi}(t) \sim [-\log t]^\alpha / \Gamma(1+\alpha)$ as $t \rightarrow 0$, $N_+(t) := \#\{k : \Delta_{(k)}(\gamma) > t\}$ grows like $\gamma[-\log t]^\alpha / \Gamma(1+\alpha)$ as $t \rightarrow 0$. Besides,

$$-\log S_{(k),\gamma} \sim (\Gamma(1+\alpha)/\gamma)^{1/\alpha} k^{1/\alpha} \text{ as } k \rightarrow \infty$$

and the ordered frequencies decay exponentially fast, but now with $k^{1/\alpha}$ (in a 'stretched exponential' Weibull way).

(viii) As another example with $\phi \in \mathcal{S}$ but with π integrable, consider the Mittag-Leffler function $\phi(x) = \sum_{m \geq 1} \frac{1}{\Gamma(1+m\alpha)} x^m$, where $\alpha \in (0, 1)$. We have $\psi(x) := -\phi(-x) =: 1 - \varphi(x)$ where

$$\varphi(x) := \sum_{m \geq 0} \frac{1}{\Gamma(1+m\alpha)} (-x)^m.$$

$\varphi(x)$ is the Mittag-Leffler LST of the random variable $S_\alpha^{-\alpha}$ where S_α is an α -stable random variable with LST $\mathbf{E}(e^{-xS_\alpha}) = e^{-x^\alpha}$, [38]. Here $\phi_\bullet = \frac{\Gamma(1+\bullet)}{\Gamma(1+\alpha\bullet)}$ and because of the latter link with the Mittag-Leffler LST, the ϕ_\bullet sequence is log-convex and $\phi \in \mathcal{S}$. For this model, the discrete abundance ξ is thus a Poisson sum of discrete Mittag-Leffler increments δ with

$$\mathbf{P}(\delta = m) = \frac{1}{\Gamma(1+m\alpha)} \frac{x^m}{\phi(x)}, \quad m \geq 1.$$

In the size-bias sampling from a random partition point of view, the Lévy-measure corresponding to Y_γ is $d\pi(t) = f_\alpha(t) dt$ where $f_\alpha(t)$ is the density of $S_\alpha^{-\alpha}$. The Laplace exponent of Y_γ is $\psi(x) = -\phi(-x)$. Because π is integrable with mass 1, Y_γ is a subordinator in the compound Poisson class (a Poisson(γ) sum of iid positive jumps with Mittag-Leffler density $f_\alpha(t)$). In the Mittag-Leffler case, the size-bias sampling is again from a finite random partition of unity, as in Example (iii). Note that as $\alpha \rightarrow 0$, $\phi(x) \sim (1-x)^{-1} - 1$ (which is a particular case of (iii)) whereas when $\alpha \rightarrow 1$, $\phi(x) \sim e^x - 1$ which is the Bell model, also in the \mathcal{S} class.

(ix) Let $\phi(x)$ solve the functional equation $\phi(x) = xg(\phi(x))$ where $g(x) = 1 + x^2/2$. Then $\phi(x) = (1 - \sqrt{1 - 2x^2})/x$ is the generating function appearing in the enumeration of rooted binary labeled trees. Only the odd ϕ_m 's are non-zero. The convergence radius of this series is $x_0 = 1/\sqrt{2}$ with $\phi(x_0) = \sqrt{2}$ and $\phi'(x_0) = \infty$. Clearly $\phi \notin \mathcal{S}$ because ϕ is only defined on $|x| \leq x_0$, so not absolutely monotone on $(-\infty, x_0)$.

6. A NEW ENGEN-LIKE EXAMPLE

We end up giving a new example of ξ sharing some common issues with the Engen's model.

Preliminaries. Previously, let us start with a general fact. Let $\phi^*(x)$ be some 'local' generating function with non-negative coefficients ϕ_m^* . Define $Z_1^*(x) = \exp \phi^*(x)$, together with $\sigma_k^*(\theta)$, the Bell polynomials associated to $\phi^*(x)$: $Z_1^*(x)^\theta = 1 + \sum_{k \geq 1} \frac{\sigma_k^*(\theta)}{k!} x^k$. Define now the new generating functions

$$\phi(x) = xZ_1^*(x) \text{ and } Z_\theta(x) = \exp(\theta\phi(x)).$$

The Taylor coefficients of ϕ are: $\phi_m = m\sigma_{m-1}^*(1)$. The Bell polynomials now associated to $\phi(x)$ are: $Z_\theta(x) = 1 + \sum_{k \geq 1} \frac{\sigma_k(\theta)}{k!} x^k$, with

$$\sigma_k(\theta) = \sum_{p=1}^k B_{k,p}(\bullet\sigma_{\bullet-1}^*(1)) \theta^p.$$

Because $\sigma_k^*(\theta)$ are binomial convolution polynomials, the following identity holds, [1]

$$(55) \quad B_{k,p}(\bullet\sigma_{\bullet-1}^*(1)) = \binom{k}{p} \sigma_{k-p}^*(p).$$

Three simple examples are:

- $\phi^*(x) = \alpha x$, $\alpha > 0$. Then $\sigma_k^*(\theta) = \alpha^k \theta^k$ leading to: $B_{k,p}(\bullet\alpha^{\bullet-1}) = \binom{k}{p} (\alpha p)^{k-p}$.

- $\phi^*(x) = e^{\alpha x} - 1$, $\alpha > 0$. Then $\sigma_k^*(\theta) = \alpha^k \sum_{p=1}^k S_{k,p} \theta^p$ (where $S_{k,p}$ are the second kind Stirling numbers), leading to: $B_{k,p}(\alpha^{\bullet-1} B_{\bullet-1}) = \binom{k}{p} \alpha^{k-p} \sum_{q=1}^{k-p} S_{k-p,q} p^q$ where $B_k = \sum_{p=1}^k S_{k,p}$ are the Bell numbers.

- $\phi^*(x)$ solves $\phi^*(x) = x \exp(\alpha \phi^*(x))$, $\alpha > 0$. Then $\sigma_k^*(\theta) = \sum_{p=1}^k B_{k,p}(\phi_\bullet^*) \theta^p$ with $B_{k,p}(\phi_\bullet^*) = \binom{k-1}{p-1} (\alpha k)^{k-p}$, leading to

$$\sigma_k^*(\theta) = \theta (\theta + \alpha k)^{k-1}.$$

We conclude that, with $\phi_\bullet = \bullet(1 + \alpha(\bullet - 1))^{\bullet-2}$

$$B_{k,p}(\phi_\bullet) = \binom{k}{p} p(p + \alpha(k-p))^{k-p-1}.$$

If $\alpha = 1$, $\phi_\bullet = \bullet(1 + \alpha(\bullet - 1))^{\bullet-2} = \bullet^{\bullet-1}$ and we recover $B_{k,p}(\bullet^{\bullet-1}) = \binom{k}{p} p k^{k-p-1} = \binom{k-1}{p-1} k^{k-p}$.

The example.

Let $\phi^*(x) = -\alpha \log(1-x)$, $\alpha > 0$. Then $\sigma_k^*(\theta) = (\alpha\theta)_k$. Looking at $\phi(x) = x \exp \phi^*(x)$ and

$$Z_\theta(x) = \exp(\theta\phi(x)) = e^{\theta x(1-x)^{-\alpha}},$$

with $\phi_\bullet = \bullet(\alpha)_{\bullet-1}$, we get $\sigma_k(\theta) = \sum_{p=1}^k B_{k,p}(\phi_\bullet) \theta^p$ where

$$(56) \quad B_{k,p}(\bullet(\alpha)_{\bullet-1}) = \binom{k}{p} (\alpha p)_{k-p}.$$

Proposition 15. *The new model $\phi(x) = x(1-x)^{-\alpha} \in \mathcal{S}$ iff $\alpha \in [0, 1]$.*

Proof: First, the convergence radius of ϕ is $x_0 = 1$.

We have $\phi'(x) = (1-x)^{-(\alpha+1)}(1-x(1-\alpha))$ and $\phi' > 0$ for all $x < x_0$ only if $\alpha \in [0, 1]$. Let then $\alpha \in [0, 1]$. Then $\phi^{(k)}(x) = (1-x)^{-(\alpha+k)}(a_k - xb_k)$ and suppose both a_k and b_k are positive with $a_k/b_k > 1$ in such a way that $\phi^{(k)} > 0$ for all $x < x_0$. Then

$$\phi^{(k+1)}(x) = (1-x)^{-(\alpha+k+1)}((\alpha+k)a_k - xb_k(\alpha+k-1))$$

with $a_{k+1} = (\alpha+k)a_k$ and $b_{k+1} = b_k(\alpha+k-1)$. Both a_{k+1} and b_{k+1} are positive with $a_{k+1}/b_{k+1} > a_k/b_k > 1$. So $\phi^{(k+1)} > 0$ for all $x < x_0$. \diamond

Corollary 16. *When $\alpha \in (0, 1)$, in the infinitely many species context, sampling from a discrete abundance model ξ built on $\phi(x) = x(1-x)^{-\alpha}$ interprets as size-bias sampling from a random partition of unity $\mathbf{S}_\infty(\gamma)$ with ordered frequencies decaying algebraically fast with k . The Laplace exponent associated to Y_γ is $\psi(x) = -\phi(-x) = x(1+x)^{-\alpha}$, $x > -1$. The estimator $\tilde{\gamma}$ of the biodiversity parameter γ is explicitly given by*

$$(57) \quad \tilde{\gamma} = \frac{P}{k-P+1} \frac{(\alpha(P-1))_{k-P+1}}{(\alpha P)_{k-P}}.$$

Proof: Clearly $\psi(x) \sim x^{1-\alpha} \rightarrow \infty$ as $x \rightarrow \infty$ and the corresponding Lévy measure π has infinite mass.

We have $\bar{\pi}(t) \sim t^{-(1-\alpha)} \rightarrow \infty$ as $t \rightarrow 0$ so that $N_+(t) := \#\{k : \Delta_{(k)}(\gamma) > t\}$ grows like $\gamma t^{-(1-\alpha)}$ as $t \rightarrow 0$ and

$$S_{(k),\gamma} \sim Y_\gamma^{-1}(k/\gamma)^{-1/(1-\alpha)} \text{ as } k \rightarrow \infty.$$

Like in the Engen model, the ordered frequencies decay algebraically fast with k .

The expression of $\tilde{\gamma}$ in (57) follows from (56). \diamond

When both k and P are large, together with $k-(1-\alpha)P$, using a simple asymptotic form for (56)

$$\tilde{\gamma} \sim \frac{P(k-(1-\alpha)P)}{k-P+1} \left(1 + \frac{\alpha+k-P}{\alpha(P-1)}\right)^{-\alpha}.$$

Acknowledgments: T.H. acknowledges partial support from the ANR Modélisation Aléatoire en Écologie, Génétique et Évolution (ANR-Manège-09-BLAN-0215 project) and from the labex MME-DII (Modèles Mathématiques et Économiques de la Dynamique, de l'Incertitude et des Interactions). Part of this work was done while S.M. was visiting Professor at the University of Cergy-Pontoise. Both authors thank

support from Basal CONICYT project PFB-03.

REFERENCES

- [1] Abbas, M.; Bouroubi, S. On new identities for Bell's polynomials. *Discrete Mathematics*, Volume 293, Issues 13, Pages 5-10 (2005).
- [2] Arratia, R.; Goldstein, L. Size bias, sampling, the waiting time paradox, and infinite divisibility: when is the increment independent? [arXiv:1105.2834](https://arxiv.org/abs/1105.2834) (2010).
- [3] Bahls, P.; Devitt-Ryder, R.; Nguyen, T. On the location of roots of logarithmically concave polynomials. Preprint available at <http://facstaff.unca.edu/~pbahls/papers/BahlsDevittRyderNguyenV2.pdf> (2010).
- [4] Berestycki, N.; Pitman, J. Gibbs distributions for random partitions generated by a fragmentation process. *Journal of Statistical Physics*, Volume 127, Number 2, 381-418 (2007).
- [5] Bernstein, S. Sur les fonctions absolument monotones. *Acta Math.* 52, no. 1, 1-66 (1929).
- [6] Bertoin J., Lévy processes. Cambridge University Press, Cambridge, (1996).
- [7] Blackwell, D.; MacQueen, J.B. Ferguson distributions via Pólya urn schemes. *The Annals of Statistics*, 1, 353–355 (1973).
- [8] Bunge, J.; Fitzpatrick, M. Estimating the Number of Species: A Review. *Journal of the American Statistical Association*, Vol. 88, No. March 1998, pp. 364-37 (1998).
- [9] Charalambides, Ch. A.; Singh, J. A review of the Stirling numbers, their generalizations and statistical applications. *Comm. Statist. Theory Methods*, 17, no. 8 (1988).
- [10] Comtet, L. *Analyse combinatoire*. Tomes 1 et 2. Presses Universitaires de France, Paris, (1970).
- [11] Costin, O.; Garoufalidis, S. Resurgence of the fractional polylogarithms. *Math. Res. Lett.* 16, no. 5, 817-826 (2009).
- [12] Darroch, J. N. On the distribution of the number of successes in independent trials. *Ann. Math. Statist.* 35, 1317-1321 (1964).
- [13] Davenport, H.; Pólya, G. On the product of two power series. *Canadian J. Math.*, no 1, 1-5 (1949).
- [14] Donnelly, P. Partition structures, Pólya urns, the Ewens sampling formula and the age of alleles. *Theoretical Population Biology*, 30, 271-288 (1986).
- [15] Engen, S. On species frequency models. *Biometrika*, no 61, 263-270 (1974).
- [16] Engen, S. *Stochastic abundance models*. Monographs on Applied Probability and Statistics, Chapman and Hall, London, (1978).
- [17] Ewens, W.J. Some remarks on the law of succession. *Athens Conference on Applied Probability and Time Series Analysis* (1995), Vol. I, 229–244, *Lecture Notes in Statistics*, 114, Springer, New York (1996).
- [18] Ewens, W.J. The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, 3, 87-112 (1972).
- [19] Ewens, W.J. Population genetics theory - the past and the future. In: *Mathematical and Statistical Developments of Evolutionary Theory*, S. Lessard Edt., Kluwer, Dordrecht, (1990).
- [20] Fisher, R. A.; Corbet, A. S.; Williams, C. B. The relation between the number of species and the number of individuals in a random sample of an animal population. *Journal of Animal Ecology*, 12, 42-58 (1943).
- [21] Garibaldi, U.; Scalas, E. *Finitary probabilistic methods in econophysics*. Cambridge University Press, Cambridge (2010).
- [22] Hardy, G. H.; Littlewood, J. E.; Pólya, G. *Inequalities*. 2d ed. Cambridge, at the University Press (1952).
- [23] Holst, L. The Poisson-Dirichlet distribution and its relatives revisited. available at: <http://www.math.kth.se/matstat/fofu/reports/PoiDir.pdf> (2001).
- [24] Hoshino, N. Engen's extended negative binomial model revisited. *Ann. Inst. Statist. Math.*, 57, No. 2, 369–387 (2005).
- [25] Hoshino, N. Random clustering based on the conditional inverse Gaussian-Poisson distribution. *J. Japan Statist. Soc.*, 33, No. 1, 105–117 (2003).
- [26] Hubbell, S. P. The neutral theory of biodiversity and biogeography and Stephen Jay Gould. *Paleobiology* 31, 122-123 (2005).

- [27] Huillet, T. Unordered and ordered sample from Dirichlet distribution. *Ann. Inst. Statist. Math.*, Vol 57, Issue 3, 597-616 (2005).
- [28] Huillet, T.; Möhle, M. Asymptotics of symmetric compound Poisson population models. Preprint available at hal-00730734 (2012).
- [29] Keener, R; Rothman, E.; Starr, N. Distributions on partitions. *Ann. Statist.* 15, no. 4, 1466-1481 (1987).
- [30] Kingman, J.F.C. Random discrete distributions. *Journal of the Royal Statistical Society. Series B*, 37, 1–22 (1975).
- [31] Kingman, J. F. C. Mathematics of genetic diversity. CBMS-NSF Regional Conference Series in Applied Mathematics, 34. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, Pa. (1980).
- [32] Kingman, J.F.C. Poisson processes. Clarendon Press, Oxford (1993).
- [33] Kolchin, V. F. Random mappings. Translated from the Russian. With a foreword by S. R. S. Varadhan. Translation Series in Mathematics and Engineering. Optimization Software, Inc., Publications Division, New York (1986).
- [34] Neveu, J. Processus ponctuels. École d' Été de Probabilités de Saint-Flour, VI 1976, pp. 249-445. *Lecture Notes in Math.*, Vol. 598, Springer-Verlag, Berlin (1977).
- [35] Pitman, J. Random discrete distributions invariant under size-biased permutation. *Advances in Applied Probability*, 28, 525-539 (1996).
- [36] Pitman, J. Exchangeable and partially exchangeable random partitions. *Probability Theory and Related Fields*, 102, 145-158 (1995).
- [37] Pitman, J.; Yor, M. The two parameter Poisson-Dirichlet distribution derived from a stable subordinator. *Annals of Probability*, 25, 855-900 (1997).
- [38] Pollard, H. The completely monotonic character of the Mittag-Leffler function $Ea(-x)$. *Bull. Amer. Math. Soc.* 54, 1115-1116 (1948).
- [39] Schoenberg, I. J. On the zeros of the generating functions of multiply positive sequences and functions. *Ann. of Math.* (2), 62, 447-471 (1955).
- [40] Steutel, F. W.; van Harn, K. Infinite divisibility of probability distributions on the real line. *Monographs and Textbooks in Pure and Applied Mathematics*, 259. Marcel Dekker, Inc., New York (2004).
- [41] Tavaré, S.; Ewens, W.J. Multivariate Ewens distribution. Chapter 41 in *Discrete Multivariate Distributions*, N.L. Johnson, S. Kotz and N. Balakrishnan Edts, Wiley, New York, 232-246 (1997).
- [42] Watterson, G. A. The sampling theory of selectively neutral alleles. *Advances in Appl. Probability* 6, 463-488 (1974).
- [43] Yang, S.L. Some identities involving the binomial sequences. *Discrete Mathematics*, Volume 308, Pages 51-58 (2008).

¹LABORATOIRE DE PHYSIQUE THÉORIQUE ET MODÉLISATION, UNIVERSITÉ DE CERGY-PONTOISE, CNRS UMR-8089, SITE DE SAINT MARTIN, 2 AVENUE ADOLPHE-CHAUVIN, 95302 CERGY-PONTOISE, FRANCE, ²DEPTO. INGENIERIA MATEMATICA AND CENTRO MODELAMIENTO MATEMATICO, UNIVERSIDAD DE CHILE, UMI 2071, UCHILE-CNRS, CASILLA 170-3 CORREO 3, SANTIAGO, CHILE, E-MAIL: HUILLET@U-CERGY.FR, SMARTINE@DIM.UCHILE.CL