

## CONDITIONALLY IDENTICALLY DISTRIBUTED SPECIES SAMPLING SEQUENCES

FEDERICO BASSETTI,\* *University of Pavia*

IRENE CRIMALDI,\*\* *University of Bologna*

FABRIZIO LEISEN,\*\*\* *University of Navarra*

### Abstract

In this paper the theory of species sampling sequences is linked to the theory of conditionally identically distributed sequences in order to enlarge the set of species sampling sequences which are mathematically tractable. The conditional identity in distribution (see Berti, Pratelli and Rigo (2004)) is a new type of dependence for random variables, which generalizes the well-known notion of exchangeability. In this paper a class of random sequences, called *generalized species sampling sequences*, is defined and a condition to have conditional identity in distribution is given. Moreover, two types of generalized species sampling sequence that are conditionally identically distributed are introduced and studied: the *generalized Poisson–Dirichlet sequence* and the *generalized Ottawa sequence*. Some examples are discussed.

*Keywords:* Conditional identity in distribution; Poisson–Dirichlet sequence; random partition; random probability measure; randomly reinforced urn; species sampling sequence; stable convergence

2010 Mathematics Subject Classification: Primary 60F05; 60G57; 60B10

### 1. Introduction

A sequence  $(X_n)_{n \geq 1}$  of random variables defined on a probability space  $(\Omega, \mathcal{A}, P)$  taking values in a Polish space, is a *species sampling sequence* if (a version of) the conditional distribution of  $X_{n+1}$  given  $X(n) := (X_1, \dots, X_n)$  is the transition kernel

$$K_{n+1}(\omega, \cdot) := \sum_{k=1}^n \tilde{p}_{n,k}(\omega) \delta_{X_k(\omega)}(\cdot) + \tilde{r}_n(\omega) \mu(\cdot), \quad (1.1)$$

where  $\tilde{p}_{n,k}(\cdot)$  and  $\tilde{r}_n(\cdot)$  are real-valued measurable functions of  $X(n)$  and  $\mu$  is a probability measure. See Pitman (1996).

As explained in Hansen and Pitman (2000), a species sampling sequence  $(X_n)_{n \geq 1}$  can be interpreted as the sequential random sampling of individuals' species from a possibly infinite population of individuals belonging to several species. If, for the sake of simplicity, we assume that  $\mu$  is diffuse (i.e. without atoms), then the interpretation is as follows. We assign a random tag

Received 2 October 2008; revision received 25 January 2010.

\* Postal address: Department of Mathematics, University of Pavia, via Ferrata 1, 27100 Pavia, Italy.

Email address: federico.bassetti@unipv.it

\*\* Postal address: Department of Mathematics, University of Bologna, Piazza di Porta San Donato 5, 40126 Bologna, Italy. Email address: crimaldi@dm.unibo.it

\*\*\* Postal address: Faculty of Economics, University of Navarra, Campus Universitario, Edificio de Biblioteca (Entrada Este), 31008, Pamplona, Spain. Email address: fabrizio.leisen@unimore.it

$X_1$ , distributed according to  $\mu$ , to the species of the first individual. Given the tags  $X_1, \dots, X_n$  of the first  $n$  individuals observed, the species of the  $(n + 1)$ th individual is a new species with probability  $\tilde{r}_n$  and it is equal to the observed species  $X_k$  with probability  $\sum_{j=1}^n \tilde{p}_{n,j} \mathbf{1}_{\{X_j=X_k\}}$ .

The concept of a species sampling sequence is naturally related to that of random partition induced by a sequence of observations (see Pitman (2006, p. 40)). Given a random vector  $X(n) = (X_1, \dots, X_n)$ , we denote by  $L_n$  the (random) number of distinct values of  $X(n)$  and by  $X^*(n) = (X_1^*, \dots, X_{L_n}^*)$  the random vector of the distinct values of  $X(n)$  in the order in which they appear. The *random partition induced by  $X(n)$*  is the random partition of the set  $\{1, \dots, n\}$  given by  $\pi^{(n)} = [\pi_1^{(n)}, \dots, \pi_{L_n}^{(n)}]$ , where

$$i \in \pi_k^{(n)} \iff X_i = X_k^*.$$

Two distinct indices  $i$  and  $j$  clearly belong to the same block  $\pi_k^{(n)}$  for a suitable  $k$  if and only if  $X_i = X_j$ . It follows that the *prediction rule* (1.1) can be rewritten as

$$K_{n+1}(\omega, \cdot) = \sum_{k=1}^{L_n(\omega)} \tilde{p}_{n,k}^*(\omega) \delta_{X_k^*(\omega)}(\cdot) + \tilde{r}_n(\omega) \mu(\cdot),$$

where

$$\tilde{p}_{n,k}^* := \sum_{j \in \pi_k^{(n)}} \tilde{p}_{n,j}.$$

In Hansen and Pitman (2000) it was proved that, if  $\mu$  is diffuse and  $(X_n)_{n \geq 1}$  is an exchangeable sequence, the coefficients  $\tilde{p}_{n,k}^*$  are almost surely equal to some function of  $\pi^{(n)}$  and they must satisfy a suitable recurrence relation. Although there are only a few explicit prediction rules which give rise to exchangeable sequences, this kind of prediction rule is appealing for many reasons. Indeed, exchangeability is a very natural assumption in many statistical problems, in particular from the Bayesian viewpoint, as well as for many stochastic models. Moreover, remarkable results are known for exchangeable sequences. For instance, such sequences satisfy a strong law of large numbers and they can be completely characterized by the well-known de Finetti representation theorem. See, e.g. Aldous (1985, Sections 3 and 7). Furthermore, for an exchangeable sequence, the empirical mean  $\sum_{k=1}^n f(X_k)/n$  and the predictive mean, i.e.  $E[f(X_{n+1}) \mid X_1, \dots, X_n]$ , converge to the same limit as the number of observations goes to  $\infty$ . This fact can be invoked to justify the use of the empirical mean in the place of the predictive mean, which is usually harder to compute. Nevertheless, in some situations the assumption of exchangeability can be too restrictive. For instance, instead of a classical Pólya urn scheme, it may be useful to deal with the so-called randomly reinforced urn scheme. See, for example, Aletti *et al.* (2009), Bai and Hu (2005), Berti *et al.* (2004), Berti *et al.* (2009), Crimaldi (2009), Crimaldi and Leisen (2008), May and Flournoy (2009), Janson (2006), May *et al.* (2005), Pemantle (2007), and the references therein. Such processes fail to be exchangeable. Our purpose is to introduce and study a class of *generalized species sampling sequences*, which are generally not exchangeable but which still have interesting mathematical properties.

We thus need to recall the notion of *conditional identity in distribution*, introduced and studied in Berti *et al.* (2004). Such a form of dependence generalizes the notion of exchangeability, preserving some of its nice predictive properties. We say that a sequence  $(X_n)_{n \geq 1}$ , defined on  $(\Omega, \mathcal{A}, P)$  and taking values in a measurable space  $(E, \mathcal{E})$ , is *conditionally identically distributed* with respect to a filtration  $\mathcal{G} = (\mathcal{G}_n)_{n \geq 0}$  (in the sequel,  $\mathcal{G}$ -CID for short), whenever

$(X_n)_{n \geq 1}$  is  $\mathcal{G}$ -adapted and, for each  $n \geq 0, j \geq 1$ , and every bounded measurable real-valued function  $f$  on  $E$ ,

$$E[f(X_{n+j}) \mid \mathcal{G}_n] = E[f(X_{n+1}) \mid \mathcal{G}_n].$$

This means that, for each  $n \geq 0$ , all the random variables  $X_{n+j}$ , with  $j \geq 1$ , are identically distributed conditionally on  $\mathcal{G}_n$ . It is clear that every exchangeable sequence is a CID sequence with respect to its natural filtration, but a CID sequence is not necessarily exchangeable. Moreover, it is possible to show that a  $\mathcal{G}$ -adapted sequence  $(X_n)_{n \geq 1}$  is  $\mathcal{G}$ -CID if and only if, for each bounded measurable real-valued function  $f$  on  $E$ ,

$$V_n^f := E[f(X_{n+1}) \mid \mathcal{G}_n]$$

is a  $\mathcal{G}$ -martingale; see Berti *et al.* (2004). Hence, the sequence  $(V_n^f)_{n \geq 0}$  converges almost surely to a random variable  $V_f$ . One of the most important features of CID sequences is the fact that this random variable  $V_f$  is also the almost-sure limit of the empirical means. More precisely, CID sequences satisfy the following strong law of large numbers: for each bounded measurable real-valued function  $f$  on  $E$ , the sequence  $(M_n^f)_{n \geq 1}$ , defined by

$$M_n^f := \frac{1}{n} \sum_{k=1}^n f(X_k),$$

converges almost surely to  $V_f$ . It also follows that the predictive mean  $E[f(X_{n+1}) \mid X_1, \dots, X_n]$  converges almost surely to  $V_f$ . In other words, CID sequences share with exchangeable sequences the remarkable fact that the predictive mean and the empirical mean merge when the number of observations diverges. Unfortunately, while, for an exchangeable sequence, we have  $V_f = E[f(X_1) \mid \mathcal{T}] = \int f(x)m(\omega, dx)$ , where  $\mathcal{T}$  is the tail  $\sigma$ -field and  $m$  is the directing random measure of the sequence, it is difficult to characterize explicitly the limit random variable  $V_f$  for a CID sequence. Indeed, no representation theorems are available for CID sequences.

This paper is organized as follows. In Section 2 we state our definition of a generalized species sampling sequence and we give a condition under which a generalized species sampling sequence is CID with respect to a suitable filtration  $\mathcal{G}$ . After recalling the notion of stable convergence in Section 3, we introduce and analyze two types of generalized species sampling sequences which are CID: the *generalized Poisson–Dirichlet sequence* (see Section 4) and the *generalized Ottawa sequence* (see Section 5). We study the asymptotic behavior of the length  $L_n$  of the random partition induced at time  $n$ , i.e. the random number of distinct values until time  $n$ , we give some central limit theorems in the sense of stable convergence, and we discuss some examples. The paper closes with a section devoted to proofs.

## 2. Generalized species sampling sequences

The Blackwell–MacQueen urn scheme provides the most famous example of exchangeable prediction rule, that is,

$$P[X_{n+1} \in \cdot \mid X_1, \dots, X_n] = \sum_{i=1}^n \frac{1}{\theta + n} \delta_{X_i}(\cdot) + \frac{\theta}{\theta + n} \mu(\cdot),$$

where  $\theta$  is a strictly positive parameter and  $\mu$  is a probability measure; see, e.g. Blackwell and MacQueen (1973) and Pitman (1996). This prediction rule determines an exchangeable

sequence  $(X_n)_{n \geq 1}$  whose directing random measure is a Dirichlet process with parameter  $\theta \mu(\cdot)$ ; see Ferguson (1973). According to this prediction rule, if  $\mu$  is diffuse, a new species is observed with probability  $\theta / (\theta + n)$  and an old species  $X_j^*$  is observed with probability proportional to the cardinality of  $\pi_j^{(n)}$ , a sort of *preferential attachment principle*. In terms of random partitions, this rule corresponds to the so-called *Chinese restaurant process*; see Pitman (2006, Chapter 3) and the references therein.

A randomly reinforced prediction rule of the same kind could work as follows:

$$P[X_{n+1} \in \cdot \mid X_1, \dots, X_n, Y_1, \dots, Y_n] = \sum_{i=1}^n \frac{Y_i}{\theta + \sum_{j=1}^n Y_j} \delta_{X_i}(\cdot) + \frac{\theta}{\theta + \sum_{j=1}^n Y_j} \mu(\cdot), \tag{2.1}$$

where  $\mu$  is a probability measure and  $(Y_n)_{n \geq 1}$  is a sequence of independent positive random variables. If  $\mu$  is diffuse then we have the following interpretation: each individual has a random positive weight  $Y_i$  and, given the first  $n$  tags  $X(n) = (X_1, \dots, X_n)$  together with the weights  $Y(n) = (Y_1, \dots, Y_n)$ , it is supposed that the species of the next individual is a new species with probability  $\theta / (\theta + \sum_{j=1}^n Y_j)$  and one of the species observed so far, say  $X_i^*$ , with probability  $\sum_{i \in \pi_j^{(n)}} Y_i / (\theta + \sum_{j=1}^n Y_j)$ . Again a preferential attachment principle. Note that, in this case, instead of describing the law of  $(X_n)_{n \geq 1}$  with the sequence of the conditional distributions of  $X_{n+1}$  given  $X(n)$ , we have a latent process  $(Y_n)_{n \geq 1}$  and we characterize  $(X_n)_{n \geq 1}$  with the sequence of the conditional distributions of  $X_{n+1}$  given  $(X(n), Y(n))$ .

Now that we have given an idea, let us formalize what we mean by a *generalized species sampling sequence*. Let  $(\Omega, \mathcal{A}, P)$  be a probability space, and let  $E$  and  $S$  be two Polish spaces, endowed with their Borel  $\sigma$ -fields  $\mathcal{E}$  and  $\mathcal{S}$ , respectively. In the sequel,  $\mathcal{F}^Z = (\mathcal{F}_n^Z)_{n \geq 0}$  will stand for the natural filtration associated with any sequence of random variables  $(Z_n)_{n \geq 1}$  on  $(\Omega, \mathcal{A}, P)$  and we set  $\mathcal{F}_\infty^Z = \bigvee_{n \geq 0} \mathcal{F}_n^Z$ . Finally,  $\mathcal{P}_n$  will denote the set of all partitions of  $\{1, \dots, n\}$ .

We shall say that a sequence  $(X_n)_{n \geq 1}$  of random variables on  $(\Omega, \mathcal{A}, P)$ , with values in  $E$ , is a generalized species sampling sequence if:

- (h<sub>1</sub>)  $X_1$  has distribution  $\mu$ ;
- (h<sub>2</sub>) there exists a sequence  $(Y_n)_{n \geq 1}$  of random variables with values in  $(S, \mathcal{S})$  such that, for each  $n \geq 1$ , a version of the conditional distribution of  $X_{n+1}$  given

$$\mathcal{F}_n := \mathcal{F}_n^X \vee \mathcal{F}_n^Y$$

is

$$K_{n+1}(\omega, \cdot) = \sum_{i=1}^n p_{n,i}(\pi^{(n)}(\omega), Y(n)(\omega)) \delta_{X_i(\omega)}(\cdot) + r_n(\pi^{(n)}(\omega), Y(n)(\omega)) \mu(\cdot)$$

with  $p_{n,i}(\cdot, \cdot)$  and  $r_n(\cdot, \cdot)$  suitable measurable functions defined on  $\mathcal{P}_n \times S^n$  with values in  $[0, 1]$ ;

- (h<sub>3</sub>)  $X_{n+1}$  and  $(Y_{n+j})_{j \geq 1}$  are conditionally independent given  $\mathcal{F}_n$ .

**Example 2.1.** Let  $\mu$  be a probability measure on  $E$ , let  $(v_n)_{n \geq 1}$  be a sequence of probability measures on  $S$ , and let  $(r_n)_{n \geq 1}$  and  $(p_{n,i})_{n \geq 1, 1 \leq i \leq n}$  be measurable functions such that

$$r_n : \mathcal{P}_n \times S^n \rightarrow [0, 1], \quad p_{n,i} : \mathcal{P}_n \times S^n \rightarrow [0, 1],$$

and

$$\sum_{i=1}^n p_{n,i}(q_n, y_1, \dots, y_n) + r_n(q_n, y_1, \dots, y_n) = 1 \tag{2.2}$$

for each  $n \geq 1$  and each  $(q_n, y_1, \dots, y_n)$  in  $\mathcal{P}_n \times S^n$ . By the Ionescu Tulcea theorem, there are two sequences of random variables,  $(X_n)_{n \geq 1}$  and  $(Y_n)_{n \geq 1}$ , defined on a suitable probability space  $(\Omega, \mathcal{A}, P)$ , taking values in  $E$  and  $S$ , respectively, such that conditions  $(h_1)$ ,  $(h_2)$ , and the following condition are satisfied:

- for each  $n$ , the random variable  $Y_{n+1}$  has distribution  $\nu_{n+1}$  and it is independent of the  $\sigma$ -field

$$\mathcal{F}_n \vee \sigma(X_{n+1}) = \mathcal{F}_{n+1}^X \vee \mathcal{F}_n^Y.$$

This last condition implies that, for each  $n$ ,  $(Y_{n+j})_{j \geq 1}$  is independent of  $\mathcal{F}_{n+1}^X \vee \mathcal{F}_n^Y$ . It follows, in particular, that  $(Y_n)_{n \geq 1}$  is a sequence of independent random variables. Next we prove that  $(h_3)$  also holds. For each bounded  $\mathcal{F}_n$ -measurable real-valued random variable  $V$ , each bounded Borel real-valued function  $f$  on  $E$ , each  $j \geq 1$ , and each bounded Borel real-valued function  $h$  on  $S^j$ , we have

$$\begin{aligned} E[Vf(X_{n+1})h(Y_{n+1}, \dots, Y_{n+j})] &= E[Vf(X_{n+1}) E[h(Y_{n+1}, \dots, Y_{n+j}) \mid \mathcal{F}_n \vee \sigma(X_{n+1})]] \\ &= E[Vf(X_{n+1}) E[h(Y_{n+1}, \dots, Y_{n+j})]] \\ &= E[V E[f(X_{n+1}) \mid \mathcal{F}_n] E[h(Y_{n+1}, \dots, Y_{n+j}) \mid \mathcal{F}_n]]. \end{aligned}$$

Hence,

$$E[f(X_{n+1})h(Y_{n+1}, \dots, Y_{n+j}) \mid \mathcal{F}_n] = E[f(X_{n+1}) \mid \mathcal{F}_n] E[h(Y_{n+1}, \dots, Y_{n+j}) \mid \mathcal{F}_n].$$

This fact is sufficient to verify that assumption  $(h_3)$  also holds.

In order to state our first result concerning generalized species sampling sequences, we need some further notation. Set

$$p_{n,j}^*(\pi^{(n)}) = p_{n,j}^*(\pi^{(n)}, Y(n)) := \sum_{i \in \pi_j^{(n)}} p_{n,i}(\pi^{(n)}, Y(n)) \quad \text{for } j = 1, \dots, L_n$$

and

$$r_n := r_n(\pi^{(n)}, Y(n)).$$

Given a partition  $\pi^{(n)}$ , denote by  $[\pi^{(n)}]_{j+}$  the partition of  $\{1, \dots, n + 1\}$  obtained by adding the element  $(n + 1)$  to the  $j$ th block of  $\pi^{(n)}$ . Finally, denote by  $[\pi^{(n)}; (n + 1)]$  the partition obtained by adding a block containing  $(n + 1)$  to  $\pi^{(n)}$ . For instance, if  $\pi^{(3)} = [(1, 3); (2)]$  then  $[\pi^{(3)}]_{2+} = [(1, 3); (2, 4)]$  and  $[\pi^{(3)}; (4)] = [(1, 3); (2); (4)]$ .

**Theorem 2.1.** *A generalized species sampling sequence  $(X_n)_{n \geq 1}$  with  $\mu$  diffuse is a CID sequence with respect to the filtration  $\mathcal{G} = (\mathcal{G}_n)_{n \geq 0}$  with  $\mathcal{G}_n := \mathcal{F}_n^X \vee \mathcal{F}_\infty^Y$  if and only if, for each  $n$ , the following condition holds  $P$ -almost surely:*

$$p_{n,j}^*(\pi^{(n)}) = r_n p_{n+1,j}^*([\pi^{(n)}; (n + 1)]) + \sum_{l=1}^{L_n} p_{n+1,j}^*([\pi^{(n)}]_{l+}) p_{n,l}^*(\pi^{(n)}) \tag{2.3}$$

for  $1 \leq j \leq L_n$ .

In the following sections, we shall introduce and study two types of generalized species sampling sequences that are CID.

We conclude this section with some remarks on the length  $L_n$  of the random partition induced by a generalized species sampling sequence at time  $n$ , i.e. the random number of distinct values of a generalized species sampling sequence until time  $n$ .

Let  $A_0 := E$  and  $A_n(\omega) := E \setminus \{X_1(\omega), \dots, X_n(\omega)\} = \{y \in E : y \notin \{X_1(\omega), \dots, X_n(\omega)\}\}$  for  $n \geq 1$ , and set  $s_0 := 1$  and  $s_n := r_n(\pi^{(n)}, Y(n))\mu(A_n) = r_n\mu(A_n)$  for each  $n \geq 1$ . (If the probability measure  $\mu$  is diffuse then  $s_n = r_n$ .) Reconsidering the species interpretation, given  $X(n) = (X_1, \dots, X_n)$  and  $Y(n) = (Y_1, \dots, Y_n)$ , the species of the  $(n + 1)$ th individual is a new species with probability  $s_n$ , that is,

$$P[L_{n+1} = L_n + 1 \mid \mathcal{F}_n] = s_n = r_n\mu(A_n).$$

Moreover, setting  $B_n = \{L_n = L_{n-1} + 1\} \in \mathcal{F}_n$  for each  $n \geq 1$  (with  $L_0 = 0$ ), we have

$$L_n = \sum_{k=1}^n I_{B_k} \quad \text{and} \quad \sum_{k \geq 1} P[B_k \mid \mathcal{F}_{k-1}] = \sum_{k \geq 1} s_{k-1}.$$

Then, by Lévy’s extension of Borel–Cantelli lemmas (see, for instance, Williams (1991, Section 12.15)), we can obtain the following simple, but useful, result.

**Proposition 2.1.** *Let  $(X_n)_{n \geq 1}$  be a generalized species sampling sequence. Then*

- (i)  $\sum_{k \geq 0} s_k < +\infty$  almost surely (a.s.) implies that  $L_n \xrightarrow{\text{a.s.}} L$ , where  $L$  is a random variable with  $P(L < +\infty) = 1$ ;
- (ii)  $\sum_{k \geq 0} s_k = +\infty$  a.s. implies that  $L_n / \sum_{k=1}^n s_{k-1} \xrightarrow{\text{a.s.}} 1$ .

In particular, in case (ii), if there exists a sequence  $(h_n)_{n \geq 1}$  of positive numbers and a random variable  $L$  such that

$$h_n \uparrow +\infty \quad \text{and} \quad \frac{1}{h_n} \sum_{k=1}^n s_{k-1} \xrightarrow{\text{a.s.}} L,$$

then  $L_n/h_n \xrightarrow{\text{a.s.}} L$ .

### 3. Stable convergence

Since in the sequel we shall deal with stable convergence, we briefly recall here this form of convergence.

Stable convergence has been introduced in Rényi (1963) and subsequently studied by various authors; see, for example, Aldous and Eagleson (1978), Jacod and Mémín (1981), Hall and Heyde (1980, p. 56). A detailed treatment, including some strengthened forms of stable convergence, can be found in Crimaldi *et al.* (2007).

Given a probability space  $(\Omega, \mathcal{A}, P)$  and a Polish space  $E$  (endowed with its Borel  $\sigma$ -field  $\mathcal{E}$ ), recall that a kernel  $K$  on  $E$  is a family  $K = (K(\omega, \cdot))_{\omega \in \Omega}$  of probability measures on  $E$  such that, for each bounded Borel real-valued function  $g$  on  $E$ , the function

$$K(g)(\omega) = \int g(x)K(\omega, dx)$$

is measurable with respect to  $\mathcal{A}$ . Given a sub- $\sigma$ -field  $\mathcal{H}$  of  $\mathcal{A}$ , we say that the kernel  $K$  is  $\mathcal{H}$ -measurable if, for each bounded Borel real-valued function  $g$  on  $E$ , the random variable  $K(g)$  is measurable with respect to  $\mathcal{H}$ . In the following, the symbol  $\mathcal{N}$  will denote the sub- $\sigma$ -field generated by the  $P$ -negligible events of  $\mathcal{A}$ . Given a sub- $\sigma$ -field  $\mathcal{H}$  of  $\mathcal{A}$  and an  $(\mathcal{H} \vee \mathcal{N})$ -measurable kernel  $K$  on  $E$ , a sequence  $(Z_n)_{n \geq 1}$  of random variables on  $(\Omega, \mathcal{A}, P)$  with values in  $E$  converges  $\mathcal{H}$ -stably to  $K$  if, for each bounded continuous real-valued function  $g$  on  $E$  and each bounded  $\mathcal{H}$ -measurable real-valued random variable  $W$ ,

$$E[g(Z_n)W] \rightarrow E[K(g)W].$$

If  $(Z_n)_{n \geq 1}$  converges  $\mathcal{H}$ -stably to  $K$  then, for each  $A \in \mathcal{H}$  with  $P(A) \neq 0$ , the sequence  $(Z_n)_{n \geq 1}$  converges in distribution under the probability measure  $P_A = P(\cdot | A)$  to the probability measure  $P_A K$  on  $E$  given by

$$P_A K(B) = P(A)^{-1} E[\mathbf{1}_A K(\cdot, B)] = \int K(\omega, B) P_A(d\omega) \quad \text{for each } B \in \mathcal{E}. \tag{3.1}$$

In particular, if  $(Z_n)_{n \geq 1}$  converges  $\mathcal{H}$ -stably to  $K$  then  $(Z_n)_{n \geq 1}$  converges in distribution to the probability measure  $P K$  on  $E$  given by

$$P K(B) = E[K(\cdot, B)] = \int K(\omega, B) P(d\omega) \quad \text{for each } B \in \mathcal{E}. \tag{3.2}$$

Moreover, if all the random variables  $Z_n$  are  $\mathcal{H}$ -measurable then the  $\mathcal{H}$ -stable convergence obviously implies the  $\mathcal{A}$ -stable convergence.

Throughout the paper, if  $U$  is a positive random variable, we shall call the *Gaussian kernel* associated with  $U$  the family

$$\mathcal{N}(0, U) = (\mathcal{N}(0, U(\omega)))_{\omega \in \Omega}$$

of Gaussian distributions with zero mean and variance equal to  $U(\omega)$  (with  $\mathcal{N}(0, 0) := \delta_0$ ). Note that, in this case, the probability measures defined in (3.1) and (3.2) are mixtures of Gaussian distributions.

### 4. Generalized Poisson–Dirichlet sequences

Let  $\alpha \geq 0$  and  $\theta > -\alpha$ . Moreover, let  $\mu$  be a probability measure on  $E$  and  $(\nu_n)_{n \geq 1}$  be a sequence of probability measures on  $[\alpha, +\infty)$ . Consider the following sequence of functions

$$p_{n,i}(q_n, y(n)) := \frac{y_i - \alpha / C_i(q_n)}{\theta + \sum_{j=1}^n y_j},$$

$$r_n(q_n, y(n)) := \frac{\theta + \alpha L(q_n)}{\theta + \sum_{j=1}^n y_j},$$

where  $y(n) = (y_1, \dots, y_n) \in [\alpha, +\infty)^n$ ,  $q_n \in \mathcal{P}_n$ ,  $C_i(q_n)$  is the cardinality of the block in  $q_n$  which contains  $i$ , and  $L(q_n)$  is the number of blocks of  $q_n$ . It is easy to see that such functions satisfy (2.2). Hence, by Example 2.1, there exists a generalized species sampling sequence  $(X_n)_{n \geq 1}$  for which

$$P[X_{n+1} \in \cdot | X(n), Y(n)] = \sum_{l=1}^{L_n} \frac{(\sum_{i \in \pi_l^{(n)}} Y_i) - \alpha}{\theta + S_n} \delta_{X_l^*(\cdot)} + \frac{\theta + \alpha L_n}{\theta + S_n} \mu(\cdot),$$

where  $(Y_n)_{n \geq 1}$  is a sequence of independent random variables such that each  $Y_n$  has law  $\nu_n$  and

$S_n := \sum_{j=1}^n Y_j$  (with  $S_0 = 0$ ). If  $\mu$  is diffuse, we can easily check that (2.3) of Theorem 2.1 holds and so  $(X_n)_{n \geq 1}$  is a CID sequence with respect to  $\mathcal{G} = (\mathcal{F}_n^X \vee \mathcal{F}_\infty^Y)_{n \geq 1}$ .

It is worthwhile noting that if  $\mu$  is diffuse,  $Y_n = 1$  for every  $n \geq 1$ ,  $\alpha \in [0, 1)$ , and  $\theta > -\alpha$ , then we get an exchangeable sequence directed by the well-known two parameters Poisson–Dirichlet process, i.e. an exchangeable sequence described by the prediction rule

$$P[X_{n+1} \in \cdot \mid X_1, \dots, X_n] = \sum_{l=1}^{L_n} \frac{\text{card}(\pi_l^{(n)}) - \alpha}{\theta + n} \delta_{X_l^*}(\cdot) + \frac{\theta + \alpha L_n}{\theta + n} \mu(\cdot). \tag{4.1}$$

See, e.g. Pitman and Yor (1997) and Pitman (2006, Section 3.2).

The case  $\alpha = 0$  essentially reduces to the randomly reinforced urn model that has been deeply studied by many authors (see, for instance, Aletti *et al.* (2009), Bai and Hu (2005), Berti *et al.* (2009), Crimaldi (2009), Flournoy and May (2009), Janson (2006), May *et al.* (2005), Pemantle (2007), and the references therein; see also Section 5 of this paper). The case when  $\mu$  is discrete and  $\alpha > 0$  has been treated in Berti *et al.* (2009). Here, we present some results for the case when  $\mu$  is diffuse and  $\alpha > 0$ .

**Proposition 4.1.** *If  $\sup_n E[Y_n^2] < +\infty$  and  $\lim_n E[Y_n] = m$ , then*

$$r_n = \frac{\theta + \alpha L_n}{\theta + S_n} \xrightarrow{\text{a.s.}} R \quad \text{and} \quad \frac{L_n}{n} \xrightarrow{\text{a.s.}} R,$$

where  $R$  is a random variable such that  $P(0 \leq R \leq 1) = 1$ . In particular, if  $m > \alpha$ , we have  $P(R = 0) = 1$ .

Later on we shall see some examples in which  $P(R > 0) > 0$ .

Let us take  $A \in \mathcal{E}$  and set  $V_n^A := P[X_{n+1} \in A \mid \mathcal{F}_n]$ . Since  $(X_n)_{n \geq 1}$  is CID, there exists a random variable  $V_A$  such that

$$V_n^A \xrightarrow{\text{a.s.}} V_A \quad \text{and} \quad M_n^A := \frac{1}{n} \sum_{k=1}^n I_A(X_k) \xrightarrow{\text{a.s.}} V_A.$$

We shall prove the following central limit theorem.

**Theorem 4.1.** *Let us assume the following conditions:*

- (i)  $\sup_n E[Y_n^u] < +\infty$  for some  $u > 2$ ;
- (ii)  $m = \lim_n E[Y_n]$ ,  $q = \lim_n E[Y_n^2]$ .

Then

$$(\sqrt{n}(M_n^A - V_n^A), \sqrt{n}(V_n^A - V_A)) \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, U_A) \times \mathcal{N}(0, \Sigma_A),$$

where

$$U_A := \left( \frac{q}{m^2} - 1 \right) V_A(1 - V_A) + \frac{\alpha^2}{m^2} R\mu(A)[1 - \mu(A)],$$

$$\Sigma_A := \frac{q}{m^2} V_A(1 - V_A) + \frac{\alpha}{m} \left( \frac{\alpha}{m} - 2 \right) R\mu(A)[1 - \mu(A)].$$

In particular,  $\sqrt{n}(M_n^A - V_A) \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, U_A + \Sigma_A)$ . Moreover,

$$E[g(\sqrt{n}(V_n^A - V_A)) \mid \mathcal{F}_n] \xrightarrow{\text{a.s.}} \mathcal{N}(0, \Sigma_A)(g)$$

for each  $g \in \mathcal{C}_b(\mathbb{R})$ .



**4.1. Case  $m > \alpha$**

By Proposition 4.1 we have  $P(R = 0) = 1$ , and so

$$U_A = \left(\frac{q}{m^2} - 1\right)V_A(1 - V_A) \quad \text{and} \quad \Sigma_A = \frac{q}{m^2}V_A(1 - V_A).$$

Taking into account the analogy with randomly reinforced Pólya urns, it is natural to think that the random variable  $V_A$  is generally not degenerate (see Aletti *et al.* (2009)). This fact implies that  $\Sigma_A$  is not degenerate, while  $U_A$  is 0 if and only if  $Y_n$  converges in  $L^2$  to the constant  $m$ . This happens, for example, in the classical case (see (4.1)) studied in Pitman and Yor (1997) and Pitman (2006, Section 3.2).

**4.2. Case  $m = \alpha$**

If  $m = \alpha$  and  $q = \alpha^2$  (i.e.  $Y_n \xrightarrow{L^2} \alpha$ ), then

$$U_A = R\mu(A)[1 - \mu(A)], \quad \Sigma_A = V_A(1 - V_A) - R\mu(A)[1 - \mu(A)],$$

and  $U_A + \Sigma_A = V_A(1 - V_A)$ .

The following examples show that, if  $m = \alpha$ , we can have  $P(R > 0) > 0$ .

**Example 4.1.** Let us take  $\alpha > 0$  and  $-\alpha < \theta \leq 0$ . Setting

$$W_n := \frac{\alpha L_n}{\alpha + \theta + S_{n-1}},$$

we have (see Lemma 6.2, below)

$$\begin{aligned} \Delta_n &:= E[W_{n+1} - W_n \mid \mathcal{F}_n] \\ &= \frac{(\alpha - Y_n)W_n}{\theta + S_n} + \frac{\alpha\theta}{(\theta + S_n)(\alpha + \theta + S_n)} \\ &\geq \frac{(\alpha - Y_n)\alpha n}{(\theta + \alpha n)^2} + \frac{\alpha\theta}{(\theta + \alpha n)(\alpha + \theta + \alpha n)}. \end{aligned}$$

Therefore, we have

$$\begin{aligned} E[W_{n+1} \mid \mathcal{F}_n] - W_1 &= \sum_{k=1}^n E[W_{k+1} - W_k \mid \mathcal{F}_k] \\ &= \sum_{k=1}^n \Delta_k \\ &\geq \alpha \sum_{k=1}^n \frac{(\alpha - Y_k)k}{(\theta + \alpha k)^2} + \alpha\theta \sum_{k=1}^n \frac{1}{(\theta + \alpha k)(\alpha + \theta + \alpha k)}, \end{aligned}$$

and so

$$\begin{aligned} E[W_{n+1}] &\geq \frac{\alpha}{\alpha + \theta} + \alpha \sum_{k=1}^n \frac{E[\alpha - Y_k]k}{(\theta + \alpha k)^2} + \alpha\theta \sum_{k=1}^n \frac{1}{(\theta + \alpha k)(\alpha + \theta + \alpha k)} \\ &= \frac{\alpha}{\alpha + \theta} + \alpha \sum_{k=1}^n \frac{E[\alpha - Y_k]k}{(\theta + \alpha k)^2} + \theta \sum_{k=1}^n \left[ \frac{1}{\theta + \alpha k} - \frac{1}{\theta + \alpha(k + 1)} \right] \\ &= 1 - \frac{\theta}{\theta + \alpha(n + 1)} + \alpha \sum_{k=1}^n \frac{E[\alpha - Y_k]k}{(\theta + \alpha k)^2}. \end{aligned}$$

Letting  $n \rightarrow +\infty$ , we obtain

$$E[R] \geq 1 + \alpha \sum_{k=1}^{\infty} \frac{E[\alpha - Y_k]k}{(\theta + \alpha k)^2}.$$

Therefore, if

$$\alpha \sum_{k=1}^{\infty} \frac{E[Y_k - \alpha]k}{(\theta + \alpha k)^2} < 1$$

then  $E[R] > 0$  and so  $P(R > 0) > 0$ .

**Example 4.2.** Let us take  $\alpha > 0, \theta = 0, P(Y_1 > \alpha) > 0$ , and  $Y_n = \alpha$  for each  $n \geq 2$ . Then, using the same notation as in the previous example, we have  $\Delta_n = 0$  for each  $n \geq 2$  and so  $E[R] = E[W_2]$ . On the other hand, we have

$$0 < E[W_2] \leq E\left[\frac{2\alpha}{\alpha + Y_1}\right] < 1.$$

Then we get  $\min[P(R > 0), P(R < 1)] > 0$ . Moreover, since it must be that  $P(\Sigma_A \geq 0) = 1$ , we find that if  $0 < \mu(A) < 1$  then  $P(V_A = 0, R > 0) = 0$  and  $P(V_A = 1, R > 0) = 0$ .

If we take  $\alpha > 0, \theta = 0$ , and  $Y_n = \alpha$  for each  $n \geq 1$ , we have  $E[R] = E[W_1] = 1$ , and so  $P(R = 1) = 1$ . Therefore, the random variable  $U_A$  is degenerate and if  $0 < \mu(A) < 1$  then  $\max[P(V_A = 0), P(V_A = 1)] = 0$ .

### 5. Generalized Ottawa sequences

We shall say that a generalized species sampling sequence  $(X_n)_{n \geq 1}$  is a *generalized Ottawa sequence* or, more briefly, a GOS, if the following conditions are satisfied for every  $n \geq 1$ .

- The functions  $r_n$  and  $p_{n,i}$  (for  $i = 1, \dots, n$ ) do not depend on the partition; hence,

$$K_{n+1}(\omega, \cdot) = \sum_{i=1}^n p_{n,i}(Y(n)(\omega))\delta_{X_i(\omega)}(\cdot) + r_n(Y(n)(\omega))\mu(\cdot).$$

- The functions  $r_n$  are strictly positive and

$$r_n(Y_1, \dots, Y_n) \geq r_{n+1}(Y_1, \dots, Y_n, Y_{n+1})$$

almost surely.

- The functions  $p_{n,i}$  satisfy, for each  $y(n) = (y_1, \dots, y_n) \in S^n$ , the equalities

$$p_{n,i}(y(n)) = \frac{r_n(y(n))}{r_{n-1}(y(n-1))} p_{n-1,i}(y(n-1)) \quad \text{for } i = 1, \dots, n-1,$$

$$p_{n,n}(y(n)) = 1 - \frac{r_n(y(n))}{r_{n-1}(y(n-1))},$$

with  $r_0 = 1$ .

For simplicity, from now on, we shall denote by  $r_n$  and  $p_{n,i}$  the  $\mathcal{F}_n^Y$ -measurable random variables  $r_n(Y(n))$  and  $p_{n,i}(Y(n))$ , that is,  $r_n := r_n(Y(n))$  and  $p_{n,i} := p_{n,i}(Y(n))$ .

First of all let us stress that any GOS is a CID sequence with respect to the filtration  $\mathcal{G} = (\mathcal{F}_n^X \vee \mathcal{F}_\infty^Y)_{n \geq 0}$ . Indeed, since  $\mathcal{G}_n = \mathcal{F}_n \vee \sigma(Y_{n+j} : j \geq 1)$ , condition (h<sub>3</sub>) implies that

$$E[f(X_{n+1}) \mid \mathcal{G}_n] = E[f(X_{n+1}) \mid \mathcal{F}_n]$$

for each bounded Borel real-valued function  $f$  on  $E$  and, hence, by (h<sub>2</sub>), we obtain

$$V_n^f := E[f(X_{n+1}) \mid \mathcal{G}_n] = \sum_{i=1}^n p_{n,i} f(X_i) + r_n E[f(X_1)].$$

Since the random variables  $p_{n+1,i}$  are  $\mathcal{G}_n$ -measurable, it follows that

$$\begin{aligned} E[V_{n+1}^f \mid \mathcal{G}_n] &= \sum_{i=1}^n p_{n+1,i} f(X_i) + p_{n+1,n+1} E[f(X_{n+1}) \mid \mathcal{G}_n] + r_{n+1} E[f(X_1)] \\ &= \frac{r_{n+1}}{r_n} \sum_{i=1}^n p_{n,i} f(X_i) + V_n^f - \frac{r_{n+1}}{r_n} V_n^f + r_{n+1} E[f(X_1)] \\ &= \frac{r_{n+1}}{r_n} V_n^f - r_{n+1} E[f(X_1)] + V_n^f - \frac{r_{n+1}}{r_n} V_n^f + r_{n+1} E[f(X_1)] \\ &= V_n^f. \end{aligned}$$

Some examples follow.

**Example 5.1.** Consider a GOS for which  $Y_n = a_n$ , where  $(a_n)_{n \geq 0}$  is a decreasing deterministic sequence with  $a_0 = 1$ ,  $a_n > 0$ , and  $r_n(y_1, \dots, y_n) = y_n$ .

If  $\mu$  is diffuse, by Proposition 2.1, we can say that  $L_n$  converges a.s. to an integrable random variable if and only if  $\sum_k a_k < +\infty$ .

**Example 5.2.** Consider a GOS for which  $(Y_n)_{n \geq 1}$  is a sequence of random variables taking values in  $(0, 1)$  and

$$r_n(y_1, \dots, y_n) = \prod_{i=1}^n y_i.$$

Note that in this case

$$P[X_{n+1} \in \cdot \mid X(n), Y(n)] = \sum_{j=1}^n \left[ (1 - Y_j) \prod_{i=j+1}^n Y_i \right] \delta_{X_j}(\cdot) + \left[ \prod_{i=1}^n Y_i \right] \mu(\cdot).$$

Assume that  $\mu$  is diffuse and that  $(Y_j)_{j \geq 1}$  is a sequence of independent random variables distributed according to a beta distribution of parameter  $(j, 1 - \alpha)$  with  $\alpha$  in  $[0, 1)$ . That is, each  $Y_j$  has density (with respect to the Lebesgue measure) on  $[0, 1]$  given by

$$x \mapsto \frac{\Gamma(j + 1 - \alpha)}{\Gamma(j)\Gamma(1 - \alpha)} \frac{x^{j-1}}{(1 - x)^\alpha},$$

where  $\Gamma(z) = \int_0^{+\infty} x^{z-1} e^{-x} dx$ . Set  $m_{1,n} := E[L_n]$  and  $m_{2,n} := E[L_n^2]$ . Note that

$$m_{1,n+1} = m_{1,n} + E[r_n] = \sum_{j=0}^n E[r_j] \tag{5.1}$$

and

$$m_{2,n+1} = 3m_{1,n+1} - 2 + 2 \sum_{j=2}^n \sum_{i=1}^{j-1} E[r_i r_j]. \tag{5.2}$$

If  $\alpha = 0$ , (5.1) gives

$$m_{1,n+1} = 1 + \sum_{j=1}^n \prod_{i=1}^j \frac{i}{1+i} = \sum_{j=0}^n \frac{1}{1+j},$$

and, after some computations, from (5.2), we also obtain

$$\begin{aligned} m_{2,n+1} &= 1 + 3 \sum_{j=1}^n \frac{1}{1+j} + 2 \sum_{j=2}^n \sum_{i=1}^{j-1} \prod_{h=1}^i \frac{h}{h+2} \prod_{k=i+1}^j \frac{k}{k+1} \\ &= 1 + 3 \sum_{j=1}^n \frac{1}{1+j} + 4 \sum_{j=3}^{n+1} \frac{1}{j} \sum_{i=3}^j \frac{1}{i}. \end{aligned}$$

Now recall that

$$\lim_{n \rightarrow +\infty} \frac{1}{\log n} \sum_{j=1}^n \frac{1}{j} = 1, \tag{5.3}$$

and, moreover, observe that

$$\lim_{n \rightarrow +\infty} \frac{1}{\log^2(n)} \sum_{j=1}^n \frac{1}{j} \sum_{i=1}^j \frac{1}{i} = \frac{1}{2}.$$

This shows that the mean of  $L_n$  diverges as  $\log n$  and the second moment diverges as  $\log^2(n)$ . More precisely,

$$\lim_{n \rightarrow +\infty} \frac{m_{1,n+1}}{\log n} = \lim_{n \rightarrow +\infty} \frac{m_{2,n+1}}{2 \log^2(n)} = 1.$$

If  $\alpha \neq 0$ , in the same way, we obtain

$$m_{1,n+1} = 1 + \Gamma(2 - \alpha) \sum_{j=1}^n \frac{\Gamma(j + 1)}{\Gamma(j + 2 - \alpha)}.$$

Now recall that

$$\frac{\Gamma(j + 1)}{\Gamma(j + 2 - \alpha)} = \frac{1}{j^{1-\alpha}} \left( 1 + O\left(\frac{1}{j}\right) \right)$$

for  $j \rightarrow +\infty$  and that

$$\lim_{n \rightarrow +\infty} \frac{1}{n^\alpha} \sum_{j=1}^n \frac{1}{j^{1-\alpha}} = \frac{1}{\alpha} \quad \text{for } \alpha \in (0, 1). \tag{5.4}$$

Hence, when  $\alpha \neq 0$ , we have

$$\lim_{n \rightarrow +\infty} \frac{m_{1,n+1}}{n^\alpha} = \frac{\Gamma(2 - \alpha)}{\alpha}.$$

**Example 5.3.** Consider a GOS for which  $(Y_n)_{n \geq 1}$  is a sequence of random variables taking values in  $\mathbb{R}_+$  and

$$r_n(y_1, \dots, y_n) = \frac{\theta}{\theta + \sum_{j=1}^n y_j}$$

with  $\theta > 0$ . Note that the randomly reinforced Blackwell–McQueen urn scheme (described by (2.1)) gives rise to a GOS. This example will be reconsidered later.

For the length  $L_n$  of the random partition induced by a GOS, we shall prove the following central limit theorem.

**Theorem 5.1.** *Let  $(X_n)_{n \geq 1}$  be a GOS with  $\mu$  diffuse, and suppose that there exists a sequence  $(h_n)_{n \geq 1}$  of positive numbers and a positive random variable  $\Lambda$  such that the following properties hold:*

$$h_n \uparrow +\infty \quad \text{and} \quad \Lambda_n := \frac{\sum_{j=1}^n r_{j-1}(1 - r_{j-1})}{h_n} \xrightarrow{\text{a.s.}} \Lambda.$$

Then, setting  $R_n := \sum_{j=1}^n r_{j-1}$ , we have

$$T_n := \frac{L_n - R_n}{\sqrt{h_n}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, \Lambda).$$

**Corollary 5.1.** *Under the same assumptions of Theorem 5.1, if  $\mathbb{P}(\Lambda > 0) = 1$  then we have*

$$\frac{T_n}{\sqrt{\Lambda_n}} = \frac{L_n - R_n}{\sqrt{\sum_{j=1}^n r_{j-1}(1 - r_{j-1})}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, 1).$$

**Example 5.4.** Let us consider a GOS with  $\mu$  diffuse and

$$r_n = \frac{\theta}{\theta + \sum_{i=1}^n Y_i},$$

where  $\theta > 0$  and the random variables  $Y_n$  are independent positive random variables such that  $\sum_n \mathbb{E}[Y_n^2]/n^2 < +\infty$  and  $\lim_n \mathbb{E}[Y_n] = m > 0$ . Then  $s_n = r_n$ , and (see, for instance, Lemma 3 of Berti *et al.* (2009)) we have

$$\left(\frac{\theta}{j} + \frac{1}{j} \sum_{i=1}^j Y_i\right)^{-1} \xrightarrow{\text{a.s.}} \frac{1}{m}.$$

Recall that

$$\frac{1}{h_n} \sum_{j=1}^n a_j b_j \rightarrow b, \tag{5.5}$$

provided that  $a_j \geq 0$ ,  $\sum_{j=1}^n a_j/h_n \rightarrow 1$ , and  $b_n \rightarrow b$  as  $n \rightarrow +\infty$ . Setting  $h_n = \log n$  and  $L = \theta/m$ , by (5.3) and (5.5), we obtain

$$\frac{1}{\log n} R_n = \frac{1}{\log n} + \frac{\theta}{\log n} \sum_{j=1}^{n-1} \frac{1}{\theta + \sum_{i=1}^j Y_i} \sim \frac{\theta}{\log n} \sum_{j=1}^{n-1} \frac{1}{j} \left(\frac{\theta}{j} + \frac{1}{j} \sum_{i=1}^j Y_i\right)^{-1} \xrightarrow{\text{a.s.}} \frac{\theta}{m},$$

and so, by Proposition 2.1, we can conclude that  $L_n/\log n \xrightarrow{\text{a.s.}} \theta/m$ .

Moreover, by (5.5) we have

$$\begin{aligned} \frac{\sum_{j=0}^{n-1} r_j(1 - r_j)}{\log n} &= \frac{\theta}{\log n} \sum_{j=1}^{n-1} \frac{\sum_{i=1}^j Y_i}{(\theta + \sum_{i=1}^j Y_i)^2} \\ &= \frac{\theta}{\log n} \sum_{j=1}^{n-1} \left( \frac{\sum_{i=1}^j Y_i/j}{\theta/j + \sum_{i=1}^j Y_i/j} \right)^2 \frac{j}{\sum_{i=1}^j Y_i} \frac{1}{j} \\ &\rightarrow \frac{\theta}{m}. \end{aligned}$$

Therefore, by Theorem 5.1 we obtain

$$T_n = \frac{L_n - R_n}{\sqrt{\log n}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}\left(0, \frac{\theta}{m}\right),$$

and so

$$\frac{L_n - R_n}{\sqrt{(\theta/m) \log n}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, 1).$$

If we take  $Y_i = 1$  for all  $i$ , we get the well-known results for the asymptotic distribution of the length of the random partition obtained for the Blackwell–McQueen urn scheme. Indeed, since  $\sum_{j=1}^n j^{-1} - \log n = \gamma + O(1/n)$ , we obtain

$$\frac{L_n - \theta \log n}{\sqrt{\theta \log n}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, 1).$$

See, for instance, Pitman (2006, pp. 68–69).

The partition structure introduced in the next example is interesting if compared to the well-known results concerning the classical two parameters Poisson–Dirichlet process. Let us recall that, for the classical Poisson–Dirichlet process, the length of the partition rescaled by  $n^{-\alpha}$  converges a.s. to a strictly positive (nondegenerate) random variable. See Pitman (2006, Theorem 3.8, p. 68). In point of fact, we give an example of a partition where the predictive structure is similar to the Chinese restaurant process while the limit behaviour is close to the classical two parameters Poisson–Dirichlet process. Indeed, the length of the partition studied below scales as  $n^{-\alpha}$ , although the limit, in this case, is a constant and not a random variable.

**Example 5.5.** Let us consider Example 5.1 with  $\mu$  diffuse and

$$a_n = \frac{\theta}{\theta + n^{1-\alpha}},$$

where  $\theta > 0$  and  $0 < \alpha < 1$ . We have  $s_n = r_n = a_n$  and, setting  $h_n = n^\alpha$  and  $L = \theta/\alpha$ , from (5.4) we obtain

$$\frac{1}{n^\alpha} R_n = \frac{1}{n^\alpha} \sum_{j=0}^{n-1} \frac{\theta}{\theta + j^{1-\alpha}} \rightarrow \frac{\theta}{\alpha}.$$

Thus, by Proposition 2.1, we obtain  $L_n/n^\alpha \xrightarrow{\text{a.s.}} \theta/\alpha$ . Furthermore, by (5.5), it is easy to see that

$$\frac{\sum_{j=0}^{n-1} r_j(1 - r_j)}{n^\alpha} = \frac{\theta}{n^\alpha} \sum_{j=1}^{n-1} \frac{j^{1-\alpha}}{(\theta + j^{1-\alpha})^2} = \frac{\theta}{n^\alpha} \sum_{j=1}^{n-1} \left( \frac{j^{1-\alpha}}{\theta + j^{1-\alpha}} \right)^2 \frac{1}{j^{1-\alpha}} \rightarrow \frac{\theta}{\alpha}.$$

Therefore, by Theorem 5.1 we obtain

$$T_n = \frac{L_n - R_n}{n^{\alpha/2}} \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}\left(0, \frac{\theta}{\alpha}\right).$$

We recall that, since a GOS  $(X_n)_{n \geq 1}$  is CID, then, for each bounded Borel real-valued function  $f$  on  $E$ , we have

$$V_n^f = E[f(X_{n+1}) \mid \mathcal{F}_n] \xrightarrow{\text{a.s.}} V_f \quad \text{and} \quad M_n^f = \frac{1}{n} \sum_{k=1}^n f(X_k) \xrightarrow{\text{a.s.}} V_f.$$

Inspired by Theorem 3.3 of Berti *et al.* (2004) and the results in Crimaldi (2009), we conclude this section with the statements of some central limit theorems for a GOS.

**Theorem 5.2.** *Let  $(X_n)_{n \geq 1}$  be a GOS. For each bounded Borel real-valued function  $f$  and each  $n \geq 1$ , let us set*

$$C_n^f := \sqrt{n}(M_n^f - V_n^f)$$

and, for  $1 \leq j \leq n$ ,

$$Z_{n,j}^f := \frac{1}{\sqrt{n}}[f(X_j) - jV_j^f + (j-1)V_{j-1}^f] = \frac{1}{\sqrt{n}}(1 + jp_{j,j})[f(X_j) - V_{j-1}^f].$$

Suppose that

- (a)  $U_n^f := \sum_{j=1}^n (Z_{n,j}^f)^2 \xrightarrow{P} U_f$ ;
- (b)  $(Z_n^f)^* := \sup_{1 \leq j \leq n} |Z_{n,j}^f| \xrightarrow{L^1} 0$ .

Then the sequence  $(C_n^f)_{n \geq 1}$  converges  $\mathcal{A}$ -stably to the Gaussian kernel  $\mathcal{N}(0, U_f)$ .

In particular, conditions (a) and (b) are satisfied if the following conditions hold:

- (a1)  $U_n^f \xrightarrow{\text{a.s.}} U_f$ ;
- (b1)  $\sup_{n \geq 1} E[(C_n^f)^2] < +\infty$ .

**Theorem 5.3.** *Let  $(X_n)_{n \geq 1}$  be a GOS, and let  $f$  be a bounded Borel real-valued function. Using the previous notation, for  $n \geq 0$ , set*

$$Q_n := p_{n+1,n+1} = 1 - \frac{r_{n+1}}{r_n} \quad \text{and} \quad D_n^f := \sqrt{n}(V_n^f - V_f).$$

Suppose that the following conditions are satisfied:

- (i)  $n \sum_{k \geq n} Q_k^2 \xrightarrow{\text{a.s.}} H$ , where  $H$  is a positive real random variable;
- (ii)  $\sum_{k \geq 0} k^2 E[Q_k^4] < \infty$ .

Then

$$E[g(D_n^f) \mid \mathcal{F}_n] \xrightarrow{\text{a.s.}} \mathcal{N}(0, H(V_{f^2} - V_f^2))(g)$$

for each  $g \in \mathcal{C}_b(\mathbb{R})$ . In particular, we have  $D_n^f \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, H(V_{f^2} - V_f^2))$ .

**Corollary 5.2.** *Using the notation of Theorem 5.3, let us set, for  $k \geq 0$ ,*

$$\rho_k := \frac{1}{r_{k+1}} - \frac{1}{r_k},$$

and assume that the following conditions hold:

- (a)  $r_k \leq c_k$  a.s. with  $\sum_{k \geq 0} k^2 c_{k+1}^4 < \infty$  and  $kr_k \xrightarrow{\text{a.s.}} \alpha$ , where  $c_k$  and  $\alpha$  are strictly positive constants;
- (b) the random variables  $\rho_k$  are independent and identically distributed with  $E[\rho_k^4] < \infty$ .

Finally, let us set  $\beta := E[\rho_k^2]$  and  $h := \alpha^2 \beta$ .

Then, the conclusion of Theorem 5.3 holds with  $H$  equal to the constant  $h$ .

Furthermore, if the assumptions of both Theorems 5.2 and 5.3 hold, then, by Lemma 1 of Berti et al. (2009), we obtain

$$(C_n^f, D_n^f) \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, U_f) \times \mathcal{N}(0, H(V_{f^2} - V_f^2)).$$

In particular,  $\sqrt{n}(M_n^f - V_f) = C_n^f + D_n^f \xrightarrow{\mathcal{A}\text{-stably}} \mathcal{N}(0, U_f + H(V_{f^2} - V_f^2))$ .

Since the proofs of these results are essentially the same as those in Berti et al. (2004), Crimaldi (2009), and Berti et al. (2009), we shall omit them. The interested reader can find all the details and some simple examples in the preprint version of this paper, Bassetti et al. (2008).

### 6. Proofs

This section contains all the proofs of the paper. Recall that

$$\mathcal{F}_n = \mathcal{F}_n^X \vee \mathcal{F}_n^Y \quad \text{and} \quad \mathcal{G}_n = \mathcal{F}_n^X \vee \mathcal{F}_\infty^Y = \mathcal{F}_n \vee \sigma(Y_{n+j} : j \geq 1).$$

So condition  $(h_3)$  of the definition of generalized species sampling sequence implies that

$$V_n^g := E[g(X_{n+1}) \mid \mathcal{G}_n] = E[g(X_{n+1}) \mid \mathcal{F}_n]$$

for each bounded Borel real-valued function  $g$  on  $E$ .

#### 6.1. Proof of Theorem 2.1

We start with a useful lemma.

**Lemma 6.1.** *If  $(X_n)_{n \geq 1}$  is a generalized species sampling sequence then we have*

$$\begin{aligned} P[(n+1) \in \pi_l^{(n+1)} \mid \mathcal{G}_n] &= P[X_{n+1} = X_l^* \mid \mathcal{F}_n] \\ &= \sum_{j \in \pi_l^{(n)}} p_{n,j}(\pi^{(n)}, Y(n)) + r_n(\pi^{(n)}, Y(n)) \mu(\{X_l^*\}) \end{aligned}$$

for each  $l = 1, \dots, L_n$ . Moreover, for each bounded Borel real-valued function  $f$  on  $E$ ,

$$\begin{aligned} E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{G}_n] &= E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{F}_n] \\ &= r_n(\pi^{(n)}, Y(n)) \int_{A_n} f(y) \mu(dy) \end{aligned}$$

holds with  $A_0 := E$  and  $A_n$  the random ‘set’ defined by

$$A_n(\omega) := E \setminus \{X_1(\omega), \dots, X_n(\omega)\} = \{y \in E : y \notin \{X_1(\omega), \dots, X_n(\omega)\}\} \quad \text{for } n \geq 1.$$



In particular, we have

$$P[L_{n+1} = L_n + 1 \mid \mathcal{G}_n] = P[L_{n+1} = L_n + 1 \mid \mathcal{F}_n] = r_n(\pi^{(n)}, Y(n))\mu(A_n) := s_n(\pi^{(n)}, Y(n)).$$

If  $\mu$  is diffuse, we have

$$P[(n + 1) \in \pi_l^{(n+1)} \mid \mathcal{G}_n] = P[X_{n+1} = X_l^* \mid \mathcal{F}_n] = \sum_{j \in \pi_l^{(n)}} p_{n,j}(\pi^{(n)}, Y(n))$$

for each  $l = 1, \dots, L_n$  and

$$\begin{aligned} E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{G}_n] &= E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{F}_n] \\ &= r_n(\pi^{(n)}, Y(n)) E[f(X_1)], \end{aligned}$$

$$P[L_{n+1} = L_n + 1 \mid \mathcal{G}_n] = P[L_{n+1} = L_n + 1 \mid \mathcal{F}_n] = r_n(\pi^{(n)}, Y(n)).$$

*Proof.* Since  $\mathcal{G}_n = \mathcal{F}_n \vee \sigma(Y_{n+j} : j \geq 1)$ , condition  $(h_3)$  implies that

$$P[(n + 1) \in \pi_l^{(n+1)} \mid \mathcal{G}_n] = P[X_{n+1} = X_l^* \mid \mathcal{G}_n] = P[X_{n+1} = X_l^* \mid \mathcal{F}_n].$$

Hence, by condition  $(h_2)$  we have

$$\begin{aligned} P[X_{n+1} = X_l^* \mid \mathcal{F}_n] &= \sum_{i=1}^n p_{n,i}(\pi^{(n)}, Y(n))\delta_{X_i}(X_l^*) + r_n(\pi^{(n)}, Y(n))\mu(\{X_l^*\}) \\ &= \sum_{j \in \pi_l^{(n)}} p_{n,j}(\pi^{(n)}, Y(n)) + r_n(\pi^{(n)}, Y(n))\mu(\{X_l^*\}) \end{aligned}$$

for each  $l = 1, \dots, L_n$ . If  $\mu$  is diffuse, we obtain

$$P[X_{n+1} = X_l^* \mid \mathcal{F}_n] = \sum_{j \in \pi_l^{(n)}} p_{n,j}(\pi^{(n)}, Y(n))$$

for each  $l = 1, \dots, L_n$ .

Now, we observe that

$$\mathbf{1}_{\{L_{n+1}=L_n+1\}} = \mathbf{1}_{B_{n+1}}(X_1, \dots, X_n, X_{n+1}),$$

where  $B_{n+1} := \{(x_1, \dots, x_{n+1}) : x_{n+1} \notin \{x_1, \dots, x_n\}\}$ . Thus, by conditions  $(h_3)$  and  $(h_2)$ , we have

$$\begin{aligned} E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{G}_n] &= E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{F}_n] \\ &= \int \mathbf{1}_{B_{n+1}}(X_1, \dots, X_n, y) f(y) K_{n+1}(\cdot, dy) \\ &= \sum_{i=1}^n p_{n,i}(\pi^{(n)}, Y(n)) \int_{A_n} f(y) \delta_{X_i}(dy) + r_n(\pi^{(n)}, Y(n)) \int_{A_n} f(y) \mu(dy) \\ &= r_n(\pi^{(n)}, Y(n)) \int_{A_n} f(y) \mu(dy). \end{aligned}$$

If we take  $f = 1$ , we obtain

$$P[L_{n+1} = L_n + 1 \mid \mathcal{G}_n] = P[L_{n+1} = L_n + 1 \mid \mathcal{F}_n] = r_n(\pi^{(n)}, Y(n))\mu(A_n).$$

Finally, if  $\mu$  is diffuse then  $\mu(A_n(\omega)) = 1$  for each  $\omega$ , and so we have

$$\int_{A_n} f(y)\mu(dy) = E[f(X_1)].$$

*Proof of Theorem 2.1.* Let us fix a bounded Borel real-valued function  $f$  on  $E$ . Using the given prediction rule, we have

$$\begin{aligned} V_n^f &= \sum_{i=1}^n p_{n,i}(\pi^{(n)}, Y(n))f(X_i) + r_n(\pi^{(n)}, Y(n))E[f(X_1)] \\ &= \sum_{j=1}^{L_n} p_{n,j}^*(\pi^{(n)})f(X_j^*) + r_n E[f(X_1)]. \end{aligned}$$

The sequence  $(X_n)$  is  $\mathcal{G}$ -CID if and only if, for each bounded Borel real-valued function  $f$  on  $E$ , the sequence  $(V_n^f)_{n \geq 0}$  is a  $\mathcal{G}$ -martingale. We observe that we have (for the sake of simplicity, we omit the dependence on  $(Y_n)_{n \geq 1}$ )

$$\begin{aligned} E[V_{n+1}^f | \mathcal{G}_n] &= \sum_{i=1}^n f(X_i)E_i + E[p_{n+1,n+1}(\pi^{(n+1)})f(X_{n+1}) | \mathcal{G}_n] + E[r_{n+1} | \mathcal{G}_n]\bar{f} \\ &= \sum_{j=1}^{L_n} f(X_j^*) \sum_{i \in \pi_j^{(n)}} E_i + E[p_{n+1,n+1}(\pi^{(n+1)})f(X_{n+1}) | \mathcal{G}_n] \\ &\quad + E[r_{n+1} | \mathcal{G}_n]\bar{f}, \end{aligned}$$

where  $E_i = E[p_{n+1,i}(\pi^{(n+1)}) | \mathcal{G}_n]$  and  $\bar{f} = E[f(X_1)]$ .

Now we compute the various conditional expectations which appear in the second member of the above equality. Since  $\mu$  is diffuse, using Lemma 6.1, we have

$$\begin{aligned} E_i &= E[p_{n+1,i}(\pi^{(n+1)}) | \mathcal{G}_n] \\ &= \sum_{l=1}^{L_n} E[\mathbf{1}_{\{(n+1) \in \pi_l^{(n+1)}\}} p_{n+1,i}(\pi^{(n+1)}) | \mathcal{G}_n] + E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} p_{n+1,i}(\pi^{(n+1)}) | \mathcal{G}_n] \\ &= \sum_{l=1}^{L_n} p_{n+1,i}([\pi^{(n)}]_{l+}) E[\mathbf{1}_{\{(n+1) \in \pi_l^{(n+1)}\}} | \mathcal{G}_n] \\ &\quad + E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} | \mathcal{G}_n] p_{n+1,i}([\pi^{(n)}; n+1]) \\ &= \sum_{l=1}^{L_n} p_{n+1,i}([\pi^{(n)}]_{l+}) \sum_{j \in \pi_l^{(n)}} p_{n,j}(\pi^{(n)}) + r_n p_{n+1,i}([\pi^{(n)}; n+1]) \\ &= \sum_{l=1}^{L_n} p_{n+1,i}([\pi^{(n)}]_{l+}) p_{n,l}^*(\pi^{(n)}) + r_n p_{n+1,i}([\pi^{(n)}; n+1]) \end{aligned}$$

and so

$$\begin{aligned} \sum_{i \in \pi_j^{(n)}} E_i &= \sum_{l=1, l \neq j}^{L_n} p_{n+1,j}^*([\pi^{(n)}]_{l+}) p_{n,l}^*(\pi^{(n)}) + \sum_{i \in \pi_j^{(n)}} p_{n+1,i}([\pi^{(n)}]_{j+}) p_{n,j}^*(\pi^{(n)}) \\ &\quad + r_n p_{n+1,j}^*([\pi^{(n)}; n+1]) \\ &= \sum_{l=1}^{L_n} p_{n+1,j}^*([\pi^{(n)}]_{l+}) p_{n,l}^*(\pi^{(n)}) - p_{n+1,n+1}([\pi^{(n)}]_{j+}) p_{n+1,j}^*(\pi^{(n)}) \\ &\quad + r_n p_{n+1,j}^*([\pi^{(n)}; n+1]). \end{aligned}$$

Moreover, using Lemma 6.1 again, we have

$$\begin{aligned} &E[p_{n+1,n+1}(\pi^{(n+1)}) f(X_{n+1}) \mid \mathcal{G}_n] \\ &= \sum_{l=1}^{L_n} E[\mathbf{1}_{\{(n+1) \in \pi_l^{(n+1)}\}} p_{n+1,n+1}(\pi^{(n+1)}) f(X_{n+1}) \mid \mathcal{G}_n] \\ &\quad + E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} p_{n+1,n+1}(\pi^{(n+1)}) f(X_{n+1}) \mid \mathcal{G}_n] \\ &= \sum_{l=1}^{L_n} E[\mathbf{1}_{\{(n+1) \in \pi_l^{(n+1)}\}} \mid \mathcal{G}_n] p_{n+1,n+1}([\pi^{(n)}]_{l+}) f(X_l^*) \\ &\quad + E[\mathbf{1}_{\{L_{n+1}=L_n+1\}} f(X_{n+1}) \mid \mathcal{G}_n] p_{n+1,n+1}([\pi^{(n)}; n+1]) \\ &= \sum_{l=1}^{L_n} \left( \sum_{k \in \pi_l^{(n)}} p_{n,k}(\pi^{(n)}) \right) p_{n+1,n+1}([\pi^{(n)}]_{l+}) f(X_l^*) + r_n p_{n+1,n+1}([\pi^{(n)}; n+1]) \bar{f} \\ &= \sum_{l=1}^{L_n} p_{n,l}^*(\pi^{(n)}) p_{n+1,n+1}([\pi^{(n)}]_{l+}) f(X_l^*) + r_n p_{n+1,n+1}([\pi^{(n)}; n+1]) \bar{f}. \end{aligned}$$

Finally, we have

$$\begin{aligned} E[r_{n+1} \mid \mathcal{G}_n] &= 1 - \sum_{i=1}^{n+1} E[p_{n+1,i}(\pi^{(n+1)}) \mid \mathcal{G}_n] \\ &= 1 - \sum_{i=1}^n E_i - E_{n+1} \\ &= 1 - \sum_{i=1}^n E_i - \sum_{l=1}^{L_n} p_{n,l}^*(\pi^{(n)}) p_{n+1,n+1}([\pi^{(n)}]_{l+}) \\ &\quad - r_n p_{n+1,n+1}([\pi^{(n)}; n+1]). \end{aligned}$$

Thus, we obtain

$$E[V_{n+1}^f \mid \mathcal{G}_n] = \sum_{j=1}^{L_n} c_{n,j} f(X_j^*) + \left( 1 - \sum_{j=1}^{L_n} c_{n,j} \right) \bar{f},$$

where

$$\begin{aligned}
 c_{n,j} &= \sum_{i \in \mathcal{X}_j^{(n)}} E_i + p_{n+1,n+1}([\pi^{(n)}]_{j^+}) p_{n,j}^*(\pi^{(n)}) \\
 &= r_n p_{n+1,j}^*([\pi^{(n)}; n+1]) + \sum_{l=1}^{L_n} p_{n+1,j}^*([\pi^{(n)}]_{l^+}) p_{n,l}^*(\pi^{(n)}).
 \end{aligned}$$

We can conclude that  $(X_n)_{n \geq 1}$  is  $\mathcal{G}$ -CID if and only if we have, for each bounded Borel real-valued function  $f$  on  $E$  and each  $n$ ,

$$\sum_{j=1}^{L_n} p_{n,j}^* f(X_j^*) + r_n \bar{f} = \sum_{j=1}^{L_n} c_{n,j} f(X_j^*) + \left(1 - \sum_{j=1}^{L_n} c_{n,j}\right) \bar{f} \quad \text{P-a.s.}$$

Since  $E$  is a Polish space, we may affirm that  $(X_n)_{n \geq 1}$  is  $\mathcal{G}$ -CID if and only if, for each  $n$ , we have, P-a.s.,

$$\sum_{j=1}^{L_n} p_{n,j}^* \delta_{X_k^*}(\cdot) + r_n \mu(\cdot) = \sum_{j=1}^{L_n} c_{n,j} \delta_{X_k^*}(\cdot) + \left(1 - \sum_{j=1}^{L_n} c_{n,j}\right) \mu(\cdot).$$

But this last equality holds if and only if, for each  $n$ , we have, P-a.s.,

$$p_{n,j}^* = c_{n,j} \quad \text{for } 1 \leq j \leq L_n;$$

that is,

$$p_{n,j}^*(\pi^{(n)}) = r_n p_{n+1,j}^*([\pi^{(n)}; \{n+1\}]) + \sum_{l=1}^{L_n} p_{n+1,j}^*([\pi^{(n)}]_{l^+}) p_{n,l}^*(\pi^{(n)}).$$

This is exactly the condition in the statement of Theorem 2.1.

**6.2. Proofs of Section 4**

We need the following preliminary lemma.

**Lemma 6.2.** *Let us set  $S_n := \sum_{j=1}^n Y_j$ . Then*

$$W_n := \frac{\alpha L_n}{\alpha + \theta + S_{n-1}} \xrightarrow{\text{a.s./}L^1} R,$$

where  $R$  is a random variable such that  $P(0 \leq R \leq 1) = 1$ .

*Proof.* We have

$$E[L_{n+1} | \mathcal{F}_n] = L_n + r_n = \frac{(\alpha + \theta + S_n)L_n + \theta}{\theta + S_n}.$$

Hence, we obtain

$$\begin{aligned}
 \Delta_n &:= E[W_{n+1} | \mathcal{F}_n] - W_n \\
 &= \frac{\alpha L_n}{\theta + S_n} - \frac{\alpha L_n}{\alpha + \theta + S_{n-1}} + \frac{\alpha \theta}{(\theta + S_n)(\alpha + \theta + S_n)} \\
 &= \frac{\alpha - Y_n}{\theta + S_n} \frac{\alpha L_n}{\alpha + \theta + S_{n-1}} + \frac{\alpha \theta}{(\theta + S_n)(\alpha + \theta + S_n)} \\
 &= \frac{(\alpha - Y_n)W_n}{\theta + S_n} + \frac{\alpha \theta}{(\theta + S_n)(\alpha + \theta + S_n)}.
 \end{aligned}$$

If  $-\alpha < \theta \leq 0$  then  $\Delta_n \leq 0$  for each  $n$  and so  $(W_n)_n$  is a positive  $\mathcal{F}$ -supermartingale. Therefore, it converges a.s. to a random variable  $R$ .

If  $\theta > 0$ , let us set  $Z_n := W_n + \theta/(\theta + S_{n-1})$ . For each  $n$ , we have

$$\begin{aligned} E[Z_{n+1} \mid \mathcal{F}_n] - Z_n &= \Delta_n - \frac{\theta Y_n}{(\theta + S_{n-1})(\theta + S_n)} \\ &= \frac{(\alpha - Y_n)W_n}{\theta + S_n} + \frac{\alpha\theta}{\theta + S_n} \left[ \frac{1}{\alpha + \theta + S_n} - \frac{Y_n}{\alpha(\theta + S_{n-1})} \right] \\ &= \frac{(\alpha - Y_n)W_n}{\theta + S_n} + \frac{\theta(\alpha - Y_n)}{(\theta + S_{n-1})(\theta + S_n)} - \frac{\alpha\theta(Y_n + \alpha)}{(\theta + S_{n-1})(\theta + S_n)(\alpha + \theta + S_n)} \\ &\leq 0. \end{aligned}$$

Therefore, the sequence  $(Z_n)_n$  is a positive  $\mathcal{F}$ -supermartingale and so it converges a.s. to a random variable  $R$ . Since  $S_n$  goes to  $+\infty$ , we obtain

$$W_n = Z_n - \frac{\theta}{\theta + S_{n-1}} \xrightarrow{\text{a.s.}} R.$$

Finally, we observe that  $0 \leq W_n \leq \alpha n/(\theta + \alpha n) \rightarrow 1$ .

*Proof of Proposition 4.1.* It is easy to verify that  $S_n/n \xrightarrow{\text{a.s.}} m$  (see, for instance, Lemma 3 of Berti *et al.* (2009)), and so by Lemma 6.2 we obtain

$$r_n = \frac{\theta + \alpha L_n}{\theta + S_n} = \frac{\theta}{\theta + S_n} + W_n \frac{\alpha + \theta + S_{n-1}}{\theta + S_n} \xrightarrow{\text{a.s.}} R.$$

Moreover, we have

$$\sum_{k \geq 1} r_{k-1} \geq 1 + (\theta + \alpha) \sum_{k \geq 1} \frac{1}{\theta + S_k} \stackrel{\text{a.s.}}{\sim} \frac{\alpha + \theta}{m} \sum_{k \geq 1} \frac{1}{k} = \infty.$$

Then, by Proposition 2.1 we find that

$$\frac{L_n}{\sum_{k=1}^n r_{k-1}} \xrightarrow{\text{a.s.}} 1.$$

Since Cesaro’s lemma implies that

$$\frac{1}{n} \sum_{k=1}^n r_{k-1} \xrightarrow{\text{a.s.}} R,$$

we obtain  $L_n/n \xrightarrow{\text{a.s.}} R$ . On the other hand, we have

$$\frac{L_n}{n} \stackrel{\text{a.s.}}{\sim} \frac{m}{\alpha} r_n \xrightarrow{\text{a.s.}} \frac{m}{\alpha} R.$$

Therefore, we have  $(m/\alpha)R \stackrel{\text{a.s.}}{\sim} R$  and so, if  $m \neq \alpha$ , it must be that  $P(R = 0) = 1$ .

*Proof of Theorem 4.1.* As already observed, assumption (ii) implies that  $S_n/n \xrightarrow{\text{a.s.}} m$  and  $r_n \xrightarrow{\text{a.s.}} R$ . After some calculations, we find that

$$V_{n+1}^A - V_n^A = \frac{Y_{n+1}[\mathbf{1}_A(X_{n+1}) - V_n^A]}{\theta + S_{n+1}} + \alpha \frac{[\mu(A) - \mathbf{1}_A(X_{n+1})]}{\theta + S_{n+1}} \mathbf{1}_{\{L_{n+1}=L_n+1\}}.$$

We want to apply Lemma 1, Theorem 2, and Remark 4 of Berti *et al.* (2009). Therefore, we have to prove the following conditions:

1.  $(1/\sqrt{n}) E[\max_{1 \leq k \leq n} k |V_{k-1}^A - V_k^A|] \rightarrow 0$ ;
2.  $E[\sup_{k \geq 1} \sqrt{k} |V_{k-1}^A - V_k^A|] < +\infty$ ;
3.  $n \sum_{k \geq n} (V_{k-1}^A - V_k^A)^2 \xrightarrow{\text{a.s.}} \Sigma_A$ ;
4.  $(1/n) \sum_{k=1}^n \{\mathbf{1}_A(X_k) - V_{k-1}^A + k(V_{k-1}^A - V_k^A)\}^2 \xrightarrow{P} U_A$ .

Conditions 1 and 2 hold. We observe that

$$|V_{n+1}^A - V_n^A| \leq \frac{Y_{n+1} + \alpha}{\theta + \alpha(n+1)}.$$

This inequality and assumption (i) imply that

$$\frac{1}{n^{u/2}} \left( E \left[ \max_{1 \leq k \leq n} k |V_{k-1}^A - V_k^A| \right] \right)^u \leq \frac{1}{n^{u/2}} \sum_{k=1}^n k^u \frac{E[(Y_k + \alpha)^u]}{(\theta + \alpha k)^u} \rightarrow 0$$

and

$$E \left[ \left( \sup_{k \geq 1} \sqrt{k} |V_{k-1}^A - V_k^A| \right)^u \right] \leq \sum_{k \geq 1} k^{u/2} \frac{E[(Y_k + \alpha)^u]}{(\theta + \alpha k)^u} < +\infty.$$

Condition 3 holds. We observe that

$$\begin{aligned} & n \sum_{k \geq n} (V_{k-1}^A - V_k^A)^2 \\ &= n \sum_{k \geq n} \frac{Y_k^2 [\mathbf{1}_A(X_k) - V_{k-1}^A]^2}{(\theta + S_k)^2} + n\alpha^2 \sum_{k \geq n} \frac{[\mu(A) - \mathbf{1}_A(X_k)]^2}{(\theta + S_k)^2} \mathbf{1}_{\{L_k=L_{k-1}+1\}} \\ &+ 2\alpha n \sum_{k \geq n} \frac{Y_k [\mathbf{1}_A(X_k) - V_{k-1}^A] [\mu(A) - \mathbf{1}_A(X_k)]}{(\theta + S_k)^2} \mathbf{1}_{\{L_k=L_{k-1}+1\}} \\ &\stackrel{\text{a.s.}}{\sim} \frac{n}{m^2} \sum_{k \geq n} \frac{Y_k^2 [\mathbf{1}_A(X_k) - V_{k-1}^A]^2}{k^2} \\ &+ \frac{\alpha^2}{m^2} n \sum_{k \geq n} \frac{[\mu(A) - \mathbf{1}_A(X_k)]^2}{k^2} \mathbf{1}_{\{L_k=L_{k-1}+1\}} \\ &+ 2 \frac{\alpha}{m^2} n \sum_{k \geq n} \frac{Y_k [\mathbf{1}_A(X_k) - V_{k-1}^A] [\mu(A) - \mathbf{1}_A(X_k)]}{k^2} \mathbf{1}_{\{L_k=L_{k-1}+1\}}. \end{aligned}$$

We want to use Lemma 3 of Berti *et al.* (2009). Therefore, we observe that

$$E[Y_k^2 [\mathbf{1}_A(X_k) - V_{k-1}^A]^2 | \mathcal{F}_{k-1}] = E[Y_k^2] E[[\mathbf{1}_A(X_k) - V_{k-1}^A]^2 | \mathcal{F}_{k-1}] \xrightarrow{\text{a.s.}} q V_A (1 - V_A),$$

and so, by a suitable truncation technique (see the proof of Corollary 7 of Berti *et al.* (2009) for details), the first term converges a.s. to

$$\frac{q}{m^2} V_A(1 - V_A).$$

Moreover, we observe that we have

$$E[[\mu(A) - \mathbf{1}_A(X_k)]^2 \mathbf{1}_{\{L_k=L_{k-1}+1\}} \mid \mathcal{F}_{k-1}] = r_{k-1}\mu(A)[1 - \mu(A)] \xrightarrow{\text{a.s.}} R\mu(A)[1 - \mu(A)],$$

and so the second term converges a.s. to

$$\frac{\alpha^2}{m^2} R\mu(A)[1 - \mu(A)]. \tag{6.1}$$

Finally, we have

$$\begin{aligned} & E[Y_k[\mathbf{1}_A(X_k) - V_{k-1}^A][\mu(A) - \mathbf{1}_A(X_k)] \mathbf{1}_{\{L_k=L_{k-1}+1\}} \mid \mathcal{F}_{k-1}] \\ &= -E[Y_k]r_{k-1}\mu(A)[1 - \mu(A)] \\ &\xrightarrow{\text{a.s.}} -mR\mu(A)[1 - \mu(A)], \end{aligned}$$

and so the third term converges a.s. to

$$-2\frac{\alpha}{m} R\mu(A)[1 - \mu(A)].$$

*Condition 4 holds.* We observe that

$$\begin{aligned} & \frac{1}{n} \sum_{k=1}^n \{ \mathbf{1}_A(X_k) - V_{k-1}^A + k(V_{k-1}^A - V_k^A) \}^2 \\ &= \frac{1}{n} \sum_{k=1}^n \left\{ [\mathbf{1}_A(X_k) - V_{k-1}^A] \left[ 1 - \frac{kY_k}{\theta + S_k} \right] + \frac{k\alpha[\mu(A) - \mathbf{1}_A(X_k)]}{\theta + S_k} \mathbf{1}_{\{L_k=L_{k-1}+1\}} \right\}^2 \\ &\stackrel{\text{a.s.}}{\sim} \frac{1}{n} \sum_{k=1}^n \left\{ [\mathbf{1}_A(X_k) - V_{k-1}^A] \left[ 1 - \frac{Y_k}{m} \right] + \frac{\alpha[\mu(A) - \mathbf{1}_A(X_k)]}{m} \mathbf{1}_{\{L_k=L_{k-1}+1\}} \right\}^2 \\ &= \frac{1}{n} \sum_{k=1}^n [\mathbf{1}_A(X_k) - V_{k-1}^A]^2 \left[ 1 - \frac{Y_k}{m} \right]^2 + \frac{\alpha^2}{m^2 n} \sum_{k=1}^n [\mu(A) - \mathbf{1}_A(X_k)]^2 \mathbf{1}_{\{L_k=L_{k-1}+1\}} \\ &\quad + \frac{2\alpha}{m} \frac{1}{n} \sum_{k=1}^n [\mathbf{1}_A(X_k) - V_{k-1}^A] \left[ 1 - \frac{Y_k}{m} \right] [\mu(A) - \mathbf{1}_A(X_k)] \mathbf{1}_{\{L_k=L_{k-1}+1\}}. \end{aligned}$$

We want to use Lemma 3 of Berti *et al.* (2009) once again. Therefore, we observe that the second term converges a.s. to (6.1). Moreover, we have

$$E \left[ [\mathbf{1}_A(X_k) - V_{k-1}^A]^2 \left[ 1 - \frac{Y_k}{m} \right]^2 \mid \mathcal{F}_{k-1} \right] = E \left[ \left[ 1 - \frac{Y_k}{m} \right]^2 \right] E \left[ [\mathbf{1}_A(X_k) - V_{k-1}^A]^2 \mid \mathcal{F}_{k-1} \right],$$

and so, by the same suitable truncation technique mentioned above, the first term converges a.s. to

$$\left( \frac{q}{m^2} - 1 \right) V_A(1 - V_A).$$

Finally, we have

$$\begin{aligned} & \mathbb{E} \left[ \left[ 1 - \frac{Y_k}{m} \right] [\mathbf{1}_A(X_k) - V_{k-1}^A][\mu(A) - \mathbf{1}_A(X_k)] \mathbf{1}_{\{L_k=L_{k-1}+1\}} \mid \mathcal{F}_{k-1} \right] \\ &= \mathbb{E} \left[ 1 - \frac{Y_k}{m} \right] \mathbb{E} [[\mathbf{1}_A(X_k) - V_{k-1}^A][\mu(A) - \mathbf{1}_A(X_k)] \mathbf{1}_{\{L_k=L_{k-1}+1\}} \mid \mathcal{F}_{k-1}], \end{aligned}$$

and so the third term converges a.s. to 0.

**6.3. Proof of Theorem 5.1**

It will be useful to introduce the sequence of increments

$$U_1 := L_1 = 1 \quad \text{and} \quad U_n := L_n - L_{n-1} \quad \text{for } n \geq 2.$$

We need a preliminary lemma.

**Lemma 6.3.** *If  $(X_n)_{n \geq 1}$  is a GOS with  $\mu$  diffuse then, for each fixed  $k$ , a version of the conditional distribution of  $(U_j)_{j \geq k+1}$  given  $\mathcal{G}_k$  is the kernel  $Q_k$  defined as*

$$Q_k(\omega, \cdot) := \bigotimes_{j=k+1}^{\infty} \mathcal{B}(1, r_{j-1}(\omega)),$$

where  $\mathcal{B}(1, r_{j-1}(\omega))$  denotes the Bernoulli distribution with parameter  $r_{j-1}(\omega)$ .

*Proof.* It is enough to verify that, for each  $n \geq 1$ , each  $\varepsilon_{k+1}, \dots, \varepsilon_{k+n} \in \{0, 1\}$ , and each bounded  $\mathcal{G}_k$ -measurable real-valued random variable  $Z$ , we have

$$\mathbb{E}[Z \mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}, \dots, U_{k+n}=\varepsilon_{k+n}\}}] = \mathbb{E} \left[ Z \prod_{j=k+1}^{k+n} r_{j-1}^{\varepsilon_j} (1 - r_{j-1})^{1-\varepsilon_j} \right]. \tag{6.2}$$

We continue with the proof by induction on  $n$ . For  $n = 1$ , by Lemma 6.1 we have

$$\mathbb{E}[Z \mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}\}}] = \mathbb{E}[Z \mathbb{E}[\mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}\}} \mid \mathcal{G}_k]] = \mathbb{E}[Z r_k^{\varepsilon_{k+1}} (1 - r_k)^{1-\varepsilon_{k+1}}].$$

Assume that (6.2) is true for  $n - 1$ , and let us prove it for  $n$ . Let us fix a bounded  $\mathcal{G}_k$ -measurable real-valued random variable  $Z$ . By Lemma 6.1 we have

$$\begin{aligned} & \mathbb{E}[Z \mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}, \dots, U_{k+n}=\varepsilon_{k+n}\}}] \\ &= \mathbb{E}[Z \mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}, \dots, U_{k+n-1}=\varepsilon_{k+n-1}\}} \mathbb{E}[U_{k+n} = \varepsilon_{k+n} \mid \mathcal{G}_{k+n-1}]] \\ &= \mathbb{E}[Z r_{k+n-1}^{\varepsilon_{k+n}} (1 - r_{k+n-1})^{1-\varepsilon_{k+n}} \mathbf{1}_{\{U_{k+1}=\varepsilon_{k+1}, \dots, U_{k+n-1}=\varepsilon_{k+n-1}\}}]. \end{aligned}$$

We have done because the random variable  $Z r_{k+n-1}^{\varepsilon_{k+n}} (1 - r_{k+n-1})^{1-\varepsilon_{k+n}}$  is also  $\mathcal{G}_k$ -measurable and (6.2) is true for  $n - 1$ .

We also need the following known result.

**Theorem 6.1.** *Let  $(Z_{n,i})_{n \geq 1, 1 \leq i \leq k_n}$  be a triangular array of square integrable centred random variables on a probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ . Suppose that, for each fixed  $n$ ,  $(Z_{n,i})_{1 \leq i \leq k_n}$  is*



independent ('row-independence property'). Moreover, set

$$\lambda_{n,i} := E[Z_{n,i}^2] = \text{var}[Z_{n,i}], \quad \lambda_n := \sum_{i=1}^{k_n} \lambda_{n,i},$$

$$V_n := \sum_{i=1}^{k_n} Z_{n,i}^2, \quad Z_n^* := \sup_{1 \leq i \leq k_n} |Z_{n,i}|,$$

and assume that  $(V_n)_{n \geq 1}$  is uniformly integrable,  $Z_n^* \xrightarrow{P} 0$ , and  $\lambda_n \rightarrow \lambda$ .

Then  $\sum_{i=1}^{k_n} Z_{n,i} \rightarrow \mathcal{N}(0, \lambda)$  in law.

*Proof.* In Hall and Heyde (1980, pp. 53–54) it was proved that, under the uniform integrability of  $(V_n)_n$ , the convergence in probability to 0 of  $(Z_n^*)_n$  is equivalent to the Lindeberg condition. Hence, it is possible to apply Corollary 3.1 of Hall and Heyde (1980, pp. 58–59) with  $\mathcal{F}_{n,i} = \sigma(Z_{n,1}, \dots, Z_{n,i})$ .

*Proof of Theorem 5.1.* Without loss of generality, we can assume that  $h_n > 0$  for each  $n$ . In order to prove the desired  $\mathcal{A}$ -stable convergence, it is enough to prove the  $(\mathcal{F}_\infty^X \vee \mathcal{F}_\infty^Y)$ -stable convergence of  $(T_n)_n$  to  $\mathcal{N}(0, \Lambda)$ . But, in order to prove this last convergence, since we have  $\mathcal{F}_\infty^X \vee \mathcal{F}_\infty^Y = \bigvee_k \mathcal{G}_k$ , it suffices to prove that, for each  $k$  and  $A$  in  $\mathcal{G}_k$  with  $P(A) \neq 0$ , the sequence  $(T_n)_n$  converges in distribution under  $P_A$  to the probability measure  $P_A \mathcal{N}(0, \Lambda)$ . In other words, it is sufficient to fix  $k$  and to verify that  $(T_{k+n})_n$  (and so  $(T_n)_n$ ) converges  $\mathcal{G}_k$ -stably to  $\mathcal{N}(0, \Lambda)$ . (Note that the kernel  $\mathcal{N}(0, \Lambda)$  is  $(\mathcal{G}_k \vee \mathcal{N})$ -measurable for each fixed  $k$ .) To this end, we observe that we have

$$T_{k+n} = \frac{\sum_{j=1}^{k+n} (U_j - r_{j-1})}{\sqrt{h_{k+n}}} = \frac{\sum_{j=1}^k (U_j - r_{j-1})}{\sqrt{h_{k+n}}} + \frac{\sum_{j=k+1}^{k+n} (U_j - r_{j-1})}{\sqrt{h_{k+n}}}.$$

Obviously, for  $n \rightarrow +\infty$ , we have

$$\frac{\sum_{j=1}^k (U_j - r_{j-1})}{\sqrt{h_{k+n}}} \xrightarrow{\text{a.s.}} 0.$$

Therefore, we have to prove that

$$\frac{\sum_{j=k+1}^{k+n} (U_j - r_{j-1})}{\sqrt{h_{k+n}}} \xrightarrow{\mathcal{G}_k\text{-stably}} \mathcal{N}(0, \Lambda). \tag{6.3}$$

From Lemma 6.3 we know that a version of the conditional distribution of  $(U_j)_{j \geq k+1}$  given  $\mathcal{G}_k$  is the kernel  $Q_k$  defined by

$$Q_k(\omega, \cdot) := \bigotimes_{j=k+1}^{\infty} \mathcal{B}(1, r_{j-1}(\omega)).$$

On the canonical space  $\mathbb{R}^{\mathbb{N}^*}$  let us consider the canonical projections  $(\xi_j)_{j \geq k+1}$ . Then, for each  $n \geq 1$ , a version of the conditional distribution of

$$\frac{\sum_{j=k+1}^{k+n} (U_j - r_{j-1})}{\sqrt{h_{k+n}}}$$

given  $\mathcal{G}_k$  is the kernel  $N_{k+n}$  so characterized: for each  $\omega$ , the probability measure  $N_{k+n}(\omega, \cdot)$  is the distribution, under the probability measure  $Q_k(\omega, \cdot)$ , of the random variable (which is defined on the canonical space)

$$\frac{\sum_{j=k+1}^{k+n} (\xi_j - r_{j-1}(\omega))}{\sqrt{h_{k+n}}}.$$

On the other hand, for almost every  $\omega$ , under  $Q_k(\omega, \cdot)$ , the random variables

$$Z_{n,i} := \frac{\xi_{k+i} - r_{k+i-1}(\omega)}{\sqrt{h_{k+n}}} \quad \text{for } n \geq 1, 1 \leq i \leq n$$

form a triangular array which satisfies the assumptions of Theorem 6.1. Indeed, we have the row-independence property and

$$E^{Q_k(\omega, \cdot)}[Z_{n,i}] = 0, \quad E^{Q_k(\omega, \cdot)}[Z_{n,i}^2] = \frac{r_{k+i-1}(\omega)(1 - r_{k+i-1}(\omega))}{h_{k+n}}.$$

Therefore, by assumption, for  $n \rightarrow +\infty$ , we have, for almost every  $\omega$ ,

$$\begin{aligned} \sum_{i=1}^n E^{Q_k(\omega, \cdot)}[Z_{n,i}^2] &= \frac{\sum_{i=1}^n r_{k+i-1}(\omega)(1 - r_{k+i-1}(\omega))}{h_{k+n}} \\ &= \Lambda_{k+n}(\omega) - \frac{h_{k-1}\Lambda_{k-1}(\omega)}{h_{k+n}} \\ &\rightarrow \Lambda(\omega). \end{aligned}$$

Moreover, under  $Q_k(\omega, \cdot)$ , we have  $Z_n^* := \sup_{1 \leq i \leq n} Z_{n,i} \leq 2/\sqrt{h_{k+n}} \rightarrow 0$ . Finally, we observe that, setting  $V_n := \sum_{i=1}^n Z_{n,i}^2$ , we have

$$E^{Q_k(\omega, \cdot)}[V_n^2] = \text{var}^{Q_k(\omega, \cdot)}[V_n] + \left( \Lambda_{k+n}(\omega) - \frac{h_{k-1}\Lambda_{k-1}(\omega)}{h_{k+n}} \right)^2$$

with

$$\begin{aligned} \text{var}^{Q_k(\omega, \cdot)}[V_n] &= \sum_{i=1}^n \text{var}^{Q_k(\omega, \cdot)}[Z_{n,i}^2] \\ &\leq \sum_{i=1}^n E^{Q_k(\omega, \cdot)}[Z_{n,i}^4] \\ &\leq 4 \left( \Lambda_{k+n}(\omega) - \frac{h_{k-1}\Lambda_{k-1}(\omega)}{h_{k+n}} \right) \frac{1}{h_{k+n}}. \end{aligned}$$

Since, for almost every  $\omega$ , the sequence  $(\Lambda_n(\omega))_n$  is bounded and  $h_n \uparrow +\infty$ , it follows that, for almost every  $\omega$ , the sequence  $(V_n)_n$  is bounded in  $L^2$  under  $Q_k(\omega, \cdot)$  and so uniformly integrable. Theorem 6.1 assures that, for almost every  $\omega$ , the sequence of probability measures

$$(N_{k+n}(\omega, \cdot))_{n \geq 1}$$

weakly converges to the Gaussian distribution  $\mathcal{N}(0, \Lambda(\omega))$ . This fact implies that, for each bounded continuous real-valued function  $g$ , we have

$$E \left[ g \left( \frac{\sum_{j=k+1}^{k+n} (U_j - r_{j-1})}{\sqrt{h_{k+n}}} \right) \middle| \mathcal{G}_k \right] \xrightarrow{\text{a.s.}} \mathcal{N}(0, \Lambda)(g).$$

The  $\mathcal{G}_k$ -stable convergence (6.3) then obviously follows.

## Acknowledgements

This research work was supported by funds of GNAMPA 2009. Irene Crimaldi would like to thank Luca Pratelli for useful discussions on generalized Poisson–Dirichlet sequences. Fabrizio Leisen would like to thank Patrizia Berti for the support during and after the PhD thesis that inspired this work.

## References

- ALDOUS, D. J. (1985). Exchangeability and related topics. In *École d'Été de Probabilités de Saint-Flour, XIII-1983* (Lecture Notes Math. **1117**), Springer, Berlin, pp. 1–198.
- ALDOUS, D. J. AND EAGLESON, G. K. (1978). On mixing and stability of limit theorems. *Ann. Prob.* **6**, 325–331.
- ALETTI, G., MAY, C. AND SECCHI, P. (2009). A central limit theorem, and related results, for a two-color randomly reinforced urn. *Adv. Appl. Prob.* **41**, 829–844.
- BAI, Z.-D. AND HU, F. (2005). Asymptotics in randomized URN models. *Ann. Appl. Prob.* **15**, 914–940.
- BASSETTI, F., CRIMALDI, I. AND LEISEN, F. (2008). Conditionally identically distributed species sampling sequences. Preprint. Available at <http://arxiv.org/abs/0806.2724>.
- BERTI, P., PRATELLI, L. AND RIGO, P. (2004). Limit theorems for a class of identically distributed random variables. *Ann. Prob.* **32**, 2029–2052.
- BERTI, P., CRIMALDI, I., PRATELLI, L. AND RIGO, P. (2009). A central limit theorem and its applications to multicolor randomly reinforced urns. Preprint. Available at <http://arxiv.org/abs/0904.0932>.
- BLACKWELL, D. AND MACQUEEN, J. B. (1973). Ferguson distributions via Pólya urn schemes. *Ann. Statist.* **1**, 353–355.
- CRIMALDI, I. (2009). An almost sure conditional convergence result and an application to a generalized Pólya urn. *Internat. Math. Forum* **4**, 1139–1156.
- CRIMALDI, I. AND LEISEN, F. (2008). Asymptotic results for a generalized Pólya urn with ‘multi-updating’ and applications to clinical trials. *Commun. Statist. Theory Meth.* **37**, 2777–2794.
- CRIMALDI, I., LETTA, G. AND PRATELLI, L. (2007). A strong form of stable convergence. In *Séminaire de Probabilités XL* (Lecture Notes Math. **1899**), Springer, Berlin, pp. 203–225.
- FERGUSON, T. S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209–230.
- HALL, P. AND HEYDE, C. C. (1980). *Martingale Limit Theory and Its Application*. Academic press, New York.
- HANSEN, B. AND PITMAN, J. (2000). Prediction rules for exchangeable sequences related to species sampling. *Statist. Prob. Lett.* **46**, 251–256.
- JACOD, J. AND MÉMIN, J. (1981). Sur un type de convergence intermédiaire entre la convergence en loi et la convergence en probabilité. In *Séminaire de Probabilités XV* (Lecture Notes Math. **850**), Springer, Berlin, pp. 529–546.
- JANSON, S. (2006). Limit theorems for triangular urn schemes. *Prob. Theory Relat. Fields* **134**, 417–452.
- MAY, C. AND FLOURNOY, N. (2009). Asymptotics in response-adaptive designs generated by a two-color, randomly reinforced urn. *Ann. Statist.* **37**, 1058–1078.
- MAY, C., PAGANONI, A. AND SECCHI, P. (2005). On a two-color generalized Pólya urn. *Metron* **63**, 115–134.
- PEMANTLE, R. (2007). A survey of random processes with reinforcement. *Prob. Surveys* **4**, 1–79.
- PITMAN, J. (1996). Some developments of the Blackwell–MacQueen urn scheme. In *Statistics, Probability and Game Theory* (IMS Lecture Notes Monogr. Ser. **30**), eds T. S. Ferguson *et al.*, Institute of Mathematical Statistics, Hayward, CA, pp. 245–267.
- PITMAN, J. (2006). *Combinatorial Stochastic Processes* (Lecture Notes Math. **1875**), Springer, Berlin.
- PITMAN, J. AND YOR, M. (1997). The two-parameter Poisson–Dirichlet distribution derived from a stable subordinator. *Ann. Prob.* **25**, 855–900.
- RÉNYI, A. (1963). On stable sequences of events. *Sankhyā A* **25**, 293–302.
- WILLIAMS, D. (1991). *Probability with Martingales*. Cambridge University Press.