

## The Bernoulli sieve: an overview

Alexander Gnedin, Alexander Iksanov, Alexander Marynych

► **To cite this version:**

Alexander Gnedin, Alexander Iksanov, Alexander Marynych. The Bernoulli sieve: an overview. 21st International Meeting on Probabilistic, Combinatorial, and Asymptotic Methods in the Analysis of Algorithms (AofA'10), 2010, Vienna, Austria. pp.329-342. hal-01185568

**HAL Id: hal-01185568**

**<https://hal.inria.fr/hal-01185568>**

Submitted on 20 Aug 2015

**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.

# The Bernoulli sieve: an overview

Alexander Gnedin<sup>1</sup>   Alexander Iksanov<sup>2</sup>   Alexander Marynych<sup>2</sup>

<sup>1</sup>*Department of Mathematics, Utrecht University, Postbus 80010, 3508 TA Utrecht, The Netherlands*

<sup>2</sup>*Faculty of Cybernetics, National T. Shevchenko University of Kiev, 01033 Kiev, Ukraine*

The Bernoulli sieve is a version of the classical balls-in-boxes occupancy scheme, in which random frequencies of infinitely many boxes are produced by a multiplicative random walk, also known as the residual allocation model or stick-breaking. We give an overview of the limit theorems concerning the number of boxes occupied by some balls out of the first  $n$  balls thrown, and present some new results concerning the number of empty boxes within the occupancy range.

**Keywords:** the occupancy problem, random partitions

## 1 Introduction

In a classical occupancy scheme  $n$  balls are thrown independently in an infinite array of boxes with probability  $p_k$  of hitting box  $k = 1, 2, \dots$ , where  $(p_k)_{k \in \mathbb{N}}$  is a fixed sequence of positive frequencies summing up to one. The quantities of traditional interest are

- $K_n$  the number of boxes occupied by at least one of  $n$  balls,
- $K_{n,r}$  the number of boxes occupied by exactly  $r$  out of  $n$  balls,
- $M_n$  the range of occupancy, equal to the maximal index of occupied box,
- $L_n := M_n - K_n$  the number of empty boxes within the occupancy range,
- $Z_n$  the number of balls in the  $M_n$ th box.

In applications ‘boxes’ are clusters, species, types of data etc. The quantities in the list characterise the sample variability, which for large  $n$  is dominantly determined by the boxes occupied by a few balls, thus determined by the way the frequencies  $p_k$  approach zero as  $k \rightarrow \infty$ . The first two variables are functionals of the induced partition of  $n$ , defined as the unordered collection of positive occupancy counts.

The *Bernoulli sieve* is a version of the occupancy scheme with random frequencies

$$p_k := W_1 W_2 \cdots W_{k-1} (1 - W_k), \quad k \in \mathbb{N}, \quad (1)$$

where  $(W_k)_{k \in \mathbb{N}}$  are independent copies of a random variable  $W$  taking values in  $(0, 1)$ . The name derives from the following recursive construction based on i.i.d.  $q_k =_d 1 - W$ : at round 1 a coin with probability

$q_1$  for heads is flipped for each of  $n$  balls and every time it turns heads the ball is put in box 1, then at round 2 a coin with probability  $q_2$  for heads is flipped for each of the remaining balls and every time it turns heads the ball is sent to box 2, and so on until all balls are allocated in boxes.

It is useful to identify frequencies (1) with the lengths of component intervals induced by splitting  $[0, 1]$  at points visited by a multiplicative random walk  $(Q_k)_{k \in \mathbb{N}_0}$ , where

$$Q_0 := 1, \quad Q_j := \prod_{i=1}^j W_i, \quad j \in \mathbb{N}.$$

In the spirit of Kingman's 'paintbox representation' of exchangeable partitions [18], we may identify the boxes with open intervals  $(Q_k, Q_{k-1})$ , and mark the balls by independent points  $U_1, \dots, U_n$  sampled from the uniform  $[0, 1]$  distribution, independently of  $(Q_k)$ . The event  $U_i \in (Q_{k-1}, Q_k)$  then means that ball  $i$  falls in box  $k$ . Keep in mind that in the natural order the intervals are indexed from the right to the left, thus the occupancy range is determined by the interval containing the leftmost mark  $\min(U_1, \dots, U_n)$ .

The Bernoulli sieve has nonrandom frequencies only when the law of  $W$  is a Dirac mass  $\delta_p$  located at some  $p \in (0, 1)$ , the frequencies  $p_k$  comprise then a geometric distribution. Results for this case can be readily recast from the numerous studies on sampling from the geometric distribution [5, 6, 19, 25] and related models like the leader election algorithms [3, 11, 22, 28], absorption sampling [7, 24] etc. It is known that asymptotic expansions of the moments of  $K_n, M_n$  and many other quantities have a component that oscillates periodically on the  $\log n$ -scale with a small amplitude [11, 27]. The same applies to distributions of the  $L_n$ 's [19, 26]. There are some peculiarities in the symmetric case  $p = 1/2$  [11, 28].

The best analytically tractable case involves random factors having beta( $\theta, 1$ ) density  $\mathbb{P}\{W \in dx\} = \theta x^{\theta-1} dx$  on  $(0, 1)$  with parameter  $\theta > 0$ . In this case the Bernoulli sieve may be viewed as a way to generate a random partition of  $n$  which follows the multivariate distribution known as the Ewens sampling formula [1]. This model has been widely studied in connection with problems of combinatorics, statistics and biology. In particular, the case  $\theta = 1$  of uniform factors is related to records and cycle patterns of random permutations under the uniform distribution on the symmetric group. It is well known [1] that  $(K_n - \theta \log n) / (\theta \log n)^{1/2}$  is asymptotically normal, and that the  $K_{n,r}$ 's converge jointly to independent Poisson( $\theta/r$ ) random variables. These classical results are complemented by the observation that  $M_n$  exhibits the same asymptotics of moments and distribution as  $K_n$ , and the number of empty boxes has the following surprising limit law:

**Theorem 1.1** [16] *If  $W$  has beta( $\theta, 1$ ) distribution then  $L_n \rightarrow_d L_\infty$ , where  $L_\infty$  has probability generating function*

$$\mathbb{E}_s^{L_\infty} = \frac{\Gamma(1+\theta)\Gamma(1+\theta-\theta s)}{\Gamma(1+2\theta-\theta s)}, \quad s \in [0, 1],$$

*which corresponds to a mixed Poisson distribution with the parameter distributed like  $\theta |\log(1-W)|$ .*

Throughout we shall use the following notation for the moments

$$\mu := \mathbb{E}|\log W|, \quad \sigma^2 := \text{Var}(\log W), \quad \nu := \mathbb{E}|\log(1-W)|,$$

which may be finite or infinite. The *standing assumption* for what follows is that the distribution of  $|\log W|$  is non-lattice. In particular, the case of sampling from the geometric distribution will be excluded.

## 2 Markov chains and distributional recursions

A random combinatorial structure which captures the occupancy of boxes by  $n$  indistinguishable balls is the *weak composition*  $C_n^*$  comprised of nonnegative integer parts summing up to  $n$ . The term *weak composition* means that zero parts are allowed, for instance, the sequence  $(2, 3, 0, 1, 0, 0, 1, 0, 0, 0, \dots)$  (padded by infinitely many 0's) is a possible value of  $C_7^*$ . A related structure which contains less information is a composition  $C_n$  obtained by discarding zero parts of  $C_n^*$ . Discarding further the order of parts in  $C_n$  yields a random partition of  $n$ . The parts of  $C_n^*$  can be represented (see [18, p. 452]) as the magnitudes of jumps of a time-homogeneous nonincreasing Markov chain  $Q_n^* = (Q_n^*(k))_{k \in \mathbb{N}_0}$  on integers, which starts at  $n$  and moves from  $n$  to  $m$  with transition probabilities

$$q^*(n, m) = \binom{n}{m} \mathbb{E}(1 - W)^{n-m} W^m, \quad m = 0, \dots, n.$$

In the same direction, parts of the composition  $C_n$  are the magnitudes of jumps of a Markov chain  $Q_n = (Q_n(k))_{k \in \mathbb{N}_0}$  with transition probabilities

$$q(n, m) = \binom{n}{m} \frac{\mathbb{E}(1 - W)^{n-m} W^m}{1 - \mathbb{E}W^n}, \quad m = 0, \dots, n - 1.$$

This Markovian realisation implies the following distributional recursions (see [16, Section 3]):

$$\begin{aligned} M_0 &= 0, \quad M_n =_d M_{Q_n^*(1)} + 1, \quad n \in \mathbb{N}, \\ K_0 &= 0, \quad K_n =_d K_{Q_n(1)} + 1, \quad n \in \mathbb{N}, \\ L_0 &= 0, \quad L_n =_d L_{Q_n^*(1)} + 1_{\{Q_n^*(1)=n\}}, \quad n \in \mathbb{N}, \end{aligned} \tag{2}$$

where in the right-hand side  $Q_n^*(1)$  is assumed independent of  $\{M_n : n \in \mathbb{N}\}$  and  $\{L_n : n \in \mathbb{N}\}$ , and  $Q_n(1)$  independent of  $\{K_n : n \in \mathbb{N}\}$ . Analysis of the recursions by known direct methods is difficult, as these impose restrictive conditions on the moments of  $Q_n(1)$  or  $Q_n^*(1)$ . Nevertheless, coupling with the multiplicative random walk allows to gain a lot of information about the compositions. For instance, let  $g(n, m)$  be the potential function, equal to the probability that  $Q_n$  ever visits state  $m$ ,

$$g(n, m) = \sum_{j=0}^{\infty} \mathbb{P}\{Q_n(j) = m\}.$$

The coupling implies that ([12, Proposition 5])

$$\lim_{n \rightarrow \infty} g(n, m) = \frac{1 - \mathbb{E}W^m}{\mu m}, \tag{4}$$

which is 0 if  $\mu = \infty$ .

The coupling readily implies stochastic subadditivity  $M_{n+m} <_d M_n + M'_m$  where the terms in the right-hand side are independent. Indeed, note first that  $M_n$  is nondecreasing. Now, when  $n$  balls have been allocated within the range  $M_n$ , adding  $m$  new balls leads to (stochastically) maximal increase of the occupancy range when all  $m$  fall outside the old range  $M_n$ , in which event the new range of occupancy is distributed like  $M_n + M'_m$ . With analogous notation,  $L_{n+m} <_d L_n + L'_m$  for exactly the same reason (although  $L_n$  is not monotone).

### 3 Asymptotics of $M_n$

Passing from the multiplicative to conventional (additive) random walk we introduce

$$S_0 := 0, \quad S_k := |\log W_1| + \dots + |\log W_k|, \quad k \in \mathbb{N}. \tag{5}$$

In this scenario the Bernoulli sieve can be defined as allocation of balls with exponentially distributed marks  $E_j = -\log U_j$ ,  $1 \leq j \leq n$ , in boxes  $(S_k, S_{k+1})$ ,  $k \in \mathbb{N}_0$ . Define

$$N_t := \inf\{k \geq 1 : S_k > t\}, \quad t \geq 0, \tag{6}$$

which is the first time  $(S_k)$  enters  $(t, \infty)$ . From the extreme-value theory we know that the maximum statistic  $T_n := \max(E_1, \dots, E_n)$  satisfies  $T_n - \log n \rightarrow_d T$ , where  $T$  has the standard Gumbel distribution  $\mathbb{P}\{T \leq x\} = \exp(-e^{-x})$ ,  $x \in \mathbb{R}$ . A key observation is that

$$M_n = N_{T_n},$$

thus the asymptotic behaviour of  $M_n$  is very much the same as that of  $N_{\log n}$ , and the latter can be concluded by means of the renewal theory. A complete description of possible limit laws and scaling/centering constants for the number of renewals  $N_t$  [16, Proposition A.1] leads to the following classification of possible limit laws for  $M_n$ .

**Theorem 3.1** [16] *The following assertions are equivalent:*

- (i) *There exist sequences  $\{a_n, b_n : n \in \mathbb{N}\}$  with  $a_n > 0$  and  $b_n \in \mathbb{R}$  such that, as  $n \rightarrow \infty$ , the variable  $(M_n - b_n)/a_n$  converges weakly to some non-degenerate and proper distribution.*
- (ii) *The distribution of  $|\log W|$  either belongs to the domain of attraction of a stable law, or the function  $\mathbb{P}\{|\log W| > x\}$  slowly varies at  $\infty$ .*

Accordingly, there are five possible modes of convergence:

- (a) *If  $\sigma^2 < \infty$  then, with constants  $b_n = \mu^{-1} \log n$  and  $a_n = (\mu^{-3} \sigma^2 \log n)^{1/2}$ , the limiting distribution of  $(M_n - b_n)/a_n$  is standard normal.*
- (b) *If  $\sigma^2 = \infty$ , and*

$$\int_0^x y^2 \mathbb{P}\{|\log W| \in dy\} \sim L(x) \quad x \rightarrow \infty,$$

*for some function  $L$  slowly varying at  $\infty$ , then, with  $b_n = \mu^{-1} \log n$  and  $a_n = \mu^{-3/2} c_{[\log n]}$ , where  $c(x)$  is any positive function satisfying  $\lim_{x \rightarrow \infty} xL(c(x))/c^2(x) = 1$ , the limiting distribution of  $(M_n - b_n)/a_n$  is standard normal.*

- (c) *If*

$$\mathbb{P}\{|\log W| > x\} \sim x^{-\alpha} L(x), \quad x \rightarrow \infty, \tag{7}$$

*for some  $L$  slowly varying at  $\infty$  and  $\alpha \in (1, 2)$  then, with  $b_n = \mu^{-1} \log n$  and  $a_n = \mu^{-\frac{\alpha+1}{\alpha}} c_{\log n}$ , where  $c(x)$  is any positive function satisfying  $\lim_{x \rightarrow \infty} xL(c(x))/c^\alpha(x) = 1$ , the limiting distribution of  $(M_n - b_n)/a_n$  is  $\alpha$ -stable with characteristic function*

$$t \mapsto \exp\{-|t|^\alpha \Gamma(1 - \alpha)(\cos(\pi\alpha/2) + i \sin(\pi\alpha/2) \operatorname{sgn}(t))\}, \quad t \in \mathbb{R}.$$

- (d) Assume that the relation (7) holds with  $\alpha = 1$ . Let  $r : \mathbb{R} \rightarrow \mathbb{R}$  be any nondecreasing function such that  $\lim_{x \rightarrow \infty} x \mathbb{P}\{|\log W| > r(x)\} = 1$  and set

$$m(x) := \int_0^x \mathbb{P}\{|\log W| > y\} dy, \quad x > 0.$$

Then, with  $b_n = \log n / (m(\log n / r(m(\log n))))$  and

$$a_n := \frac{r(\log n / m(\log n))}{m(\log n)},$$

the limiting distribution of  $(M_n - b_n) / a_n$  is 1-stable with characteristic function

$$t \mapsto \exp\{-|t|(\pi/2 - i \log |t| \operatorname{sgn}(t))\}, \quad t \in \mathbb{R}.$$

- (e) If the relation (7) holds for  $\alpha \in [0, 1)$  then, with  $b_n \equiv 0$  and  $a_n := \log^\alpha n / L(\log n)$ , the limiting distribution of  $M_n / a_n$  is the Mittag-Leffler law  $\theta_\alpha$  with moments

$$\int_0^\infty x^k \theta_\alpha(dx) = \frac{k!}{\Gamma^k(1 - \alpha) \Gamma(1 + \alpha k)}, \quad k \in \mathbb{N}.$$

## 4 Asymptotics of $K_n$

Loosely speaking,  $\nu$  controls the mean number of empty boxes, so that  $\nu < \infty$  implies  $\lim_{n \rightarrow \infty} \mathbb{E}L_n = \nu / \mu < \infty$  (Theorem 7.1 to follow). Thus when  $\nu < \infty$  the identity  $K_n = M_n - L_n$  suggests that  $K_n$  does not differ much from  $M_n$ . A first result of this kind was obtained in [12]: assuming  $\nu < \infty$  and  $\sigma^2 < \infty$  it was shown that

$$(K_n - \mu^{-1} \log n) / \sqrt{\sigma^2 \mu^{-3} \log n} \rightarrow_d \text{normal}(0, 1), \quad n \rightarrow \infty.$$

The proof was based on a careful analysis of the recursion (2) to conclude on the asymptotics of  $\operatorname{Var} K_n$  and to eventually prove the normal limit.

The similarity between  $M_n$  and  $K_n$  was justified in full generality in [16], where it was shown that under the assumption  $\nu < \infty$  Theorem 3.1 remains valid if  $M_n$  is replaced by  $K_n$ .

Another approach which allows one to treat the cases of finite and infinite  $\nu$  in a unified way was proposed in [15]. It was suggested to approximate  $K_n$  by  $N^*(\log n)$ , where

$$\begin{aligned} N^*(x) &:= \#\{k \in \mathbb{N} : p_k \geq e^{-x}\} \\ &= \#\{k \in \mathbb{N} : W_1 \cdots W_{k-1} (1 - W_k) \geq e^{-x}\}, \quad x > 0. \end{aligned}$$

The connection exemplifies the general idea that the variability of  $K_n$  stems from randomness in frequencies  $(p_k)$  superposed with randomness in sampling, and the first often plays a dominating role through the conditional law of large numbers  $K_n \sim \mathbb{E}(K_n | (p_k))$  a.s. (see [23]). Thus we believe that the approach based on  $N^*(x)$  offers a natural and the most adequate way to study the asymptotics of  $K_n$ . The following result was proved in [15].

**Theorem 4.1** *If there exist functions  $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  and  $g : \mathbb{R}_+ \rightarrow \mathbb{R}$  such that  $(N_t - g(t))/f(t)$  converges weakly (as  $t \rightarrow \infty$ ) to some non-degenerate and proper distribution, then also  $(K_n - b_n)/a_n$  converge weakly (as  $n \rightarrow \infty$ ) to the same distribution, where the constants are given by*

$$b_n = \int_0^{\log n} g(\log n - y) \mathbb{P}\{|\log(1 - W)| \in dy\}, \quad a_n = f(\log n).$$

As in [16], the convergence criterion for  $N_t$  leads to a complete characterisation of possible normalisations and limiting laws for  $K_n$ , see Corollary 1.1 in [15]. But Theorem 4.1 says more: if  $\nu = \infty$  the behaviour of  $L_n$  may affect the asymptotics of  $K_n = M_n - L_n$ . The following example illustrates the phenomenon.

**Example 4.1** Assume that, for some  $\gamma \in (0, 1/2)$ ,

$$\mathbb{P}\{W > x\} = \frac{1}{1 + |\log(1 - x)|^\gamma}, \quad x \in [0, 1).$$

Then

$$\mathbb{E} \log^2 W < \infty \quad \text{and} \quad \mathbb{P}\{|\log(1 - W)| > x\} \sim x^{-\gamma} \quad \text{as } x \rightarrow \infty,$$

and in this case,

$$a_n = \text{const } \log^{1/2} n \quad \text{and} \quad d_n = \mu^{-1}(\log n - (1 - \gamma)^{-1} \log^{1-\gamma} n + o(\log^{1-\gamma} n)).$$

Thus we see that the second term  $d_n - \mu^{-1} \log n$  of centering cannot be ignored. Moreover, one can check that

$$\mathbb{E} L_n \sim \frac{1}{\mu} \sum_{k=1}^n \frac{\mathbb{E} W^k}{k} \sim b_n - \mu^{-1} \log n \sim \frac{1}{\mu(1 - \gamma)} \log^{1-\gamma} n,$$

which reveals the indispensable contribution of  $L_n$ .

## 5 Weak convergence of $K_{n,r}$

Assume  $\mu < \infty$ . For  $B := \{\prod_{i=1}^k W_i : k \in \mathbb{N}_0\}$  the set of sites visited by the multiplicative random walk, consider a point process with unit atoms located at points of  $-\log B$  (which are the sites visited by  $S_k$ ,  $k \in \mathbb{N}_0$ ). By the renewal theorem the point process  $-\log B - \log n$  vaguely converges to a shift-invariant renewal process  $\mathcal{P}$  on the whole line. Therefore, the point process  $nB$  converges vaguely to a point process  $\mathcal{B} := \exp(-\mathcal{P})$  on  $\mathbb{R}_+$ . Think of intervals between consecutive points of  $\mathcal{B}$  as a series of boxes. Note that the process is self-similar, meaning that  $c\mathcal{B} =_d \mathcal{B}$  for every  $c > 0$ , and has the intensity measure  $(\mu x)^{-1} dx$ , so the atoms accumulate at 0 and  $\infty$ . In the role of balls assume the points of a unit Poisson process  $\mathcal{U}$  independent of  $\mathcal{B}$ . A well-known fact of extreme value theory is that  $\mathcal{U}$  is the vague limit of the point process with unit atoms located at  $nU_j$ ,  $1 \leq j \leq n$ . The location of the leftmost atom of  $\mathcal{U}$ , say  $Y$ , has exponential distribution. For  $r \geq 0$  define  $\hat{K}_r$  to be the number of component intervals of  $(Y, \infty) \setminus \mathcal{B}$  that contain exactly  $r$  atoms of  $\mathcal{U}$ . The existence of weak limits for the occupancy counts is read off from the convergence of point processes:

**Theorem 5.1** [17] *As  $n \rightarrow \infty$  we have the joint convergence in distribution*

$$(L_n, K_{n,1}, K_{n,2}, \dots) \rightarrow_d (\hat{K}_0, \hat{K}_1, \hat{K}_2, \dots)$$

along with

$$\mathbb{E}K_{n,r} \rightarrow \mathbb{E}\hat{K}_r = \frac{1}{r\mu}, \quad r > 0.$$

When  $W =_d \text{beta}(\theta, 1)$  the process  $\mathcal{B}$  is Poisson with intensity  $\theta x^{-1}dx$ . By self-similarity, the partition induced by allocation of  $n$  leftmost atoms of  $\mathcal{U}$  is the Ewens partition. The theorem allows to re-prove the results on asymptotics of the Ewens partition mentioned in Introduction, along with Theorem 1.1. Except the  $\text{beta}(\theta, 1)$  case no explicit formulas for the distribution of occupancy counts are known; in general the  $\hat{K}_r$ 's are neither independent, nor Poisson. See more on self-similar partitions in [13, Section 5].

## 6 Asymptotics of $Z_n$

The variable  $Z_n$  is analogous to the number of winners in the leader election algorithm [4, 5, 6, 25].

**Theorem 6.1** [16] *The number of balls in the last occupied box satisfies:*

- (1) *If  $\mu < \infty$  then  $Z_n \rightarrow_d Z$ ,  $n \rightarrow \infty$ , where the variable  $Z$  has distribution*

$$\mathbb{P}\{Z = k\} = \frac{\mathbb{E}(1 - W)^k}{\mu k}, \quad k \in \mathbb{N}.$$

- (2) *If (7) holds with  $\alpha \in [0, 1)$  then*

$$\frac{\log Z_n}{\log n} \rightarrow_d Z^{(\alpha)}, \quad n \rightarrow \infty,$$

where the law of  $Z^{(0)}$  is  $\delta_1$ , while for  $\alpha \in (0, 1)$  we have  $Z^{(\alpha)} =_d \text{beta}(1 - \alpha, \alpha)$ .

- (3) *If (7) holds with  $\alpha = 1$  and  $\mu = \infty$ , then*

$$\frac{m(\log Z_n)}{m(\log n)} \rightarrow_d Z^{(1)}, \quad n \rightarrow \infty,$$

where  $m(x) = \int_0^x \mathbb{P}\{|\log W| > y\}dy$ , and  $Z^{(1)} =_d \text{uniform}[0, 1]$ .

The case  $\mu < \infty$  is quite elementary, as is seen from

$$\mathbb{P}\{Z_n = m\} = g(n, m)\mathbb{P}\{Q_m(1) = 0\} = g(n, m) \frac{\mathbb{E}(1 - W)^m}{1 - \mathbb{E}W^m} \tag{8}$$

and (4). In the case  $\mu = \infty$  the result follows from the known limit distribution of the undershoot  $U(z) = z - S_{N(z)-1}$  (see [8, 10]) and the representation

$$\mathbb{P}\{Z_n > k\} = \mathbb{P}\{U(E_{n,n}) > E_{n,n} - E_{n-k,n}\}, \quad k \in \mathbb{N},$$

where  $E_{1,n} \leq \dots \leq E_{n,n} = T_n$  are the order statistics of the exponential variables  $E_j$ ,  $1 \leq j \leq n$ .



## 7 Asymptotics of $L_n$

Although there is an explicit formula

$$\mathbb{E}L_n = \sum_{k=1}^n (-1)^{k+1} \binom{n}{k} \frac{1 - \mathbb{E}(1 - W)^k}{1 - \mathbb{E}W^k}, \quad (9)$$

it does not seem possible to employ it in order to conclude on the asymptotic behaviour of  $\mathbb{E}L_n$  without restrictive additional assumptions.

Using a different approach we arrived at

**Theorem 7.1** *The expectation  $\mathbb{E}L_n$  exhibits the following asymptotic behaviour:*

(i) *If  $\mu = \infty$  and  $\nu = \infty$  then*

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}W^n}{\mathbb{E}(1 - W)^n} \leq \liminf_{n \rightarrow \infty} \mathbb{E}L_n \leq \limsup_{n \rightarrow \infty} \mathbb{E}L_n \leq \limsup_{n \rightarrow \infty} \frac{\mathbb{E}W^n}{\mathbb{E}(1 - W)^n}.$$

*In particular,*

$$\lim_{n \rightarrow \infty} \frac{\mathbb{E}W^n}{\mathbb{E}(1 - W)^n} = \gamma_0 \in [0, \infty]$$

*implies  $\lim_{n \rightarrow \infty} \mathbb{E}L_n = \gamma_0$ .*

(ii) *If  $\nu < \infty$  and  $\mu \leq \infty$  then*

$$\lim_{n \rightarrow \infty} \mathbb{E}L_n = \nu/\mu.$$

(iii) *If  $\mu < \infty$  and  $\nu = \infty$  then, as  $n \rightarrow \infty$ ,*

$$\mathbb{E}L_n \sim \frac{1}{\mu} \int_1^n \frac{\mathbb{E}e^{-y(1-W)}}{y} dy.$$

**Proof:** Part (i). Set  $s_m = \frac{\mathbb{E}W^m}{\mathbb{E}(1-W)^m}$ . We will use the representation

$$\mathbb{E}L_n = \mathbb{E}s_{Z_n}, \quad (10)$$

which follows from (8). The array  $c_{n,m} := \mathbb{P}\{Z_n = m\}$  verifies the conditions of Lemma 8.1 in Appendix, in particular by the assumption  $\mu = \infty$ . Hence the lemma can be applied to  $t_n = \mathbb{E}L_n$ , whence the assertion. When  $\gamma_0$  is well defined the proof is simpler, as in this case the statement follows from (10), divergence of  $Z_n$ , and by using dominated convergence in the case  $\gamma_0 < \infty$ , respectively using Fatou's lemma in the case  $\gamma_0 = \infty$ .

See [16] and [17] for (ii).

For part (iii) we use the poissonised version of the Bernoulli sieve, in which balls are thrown one-by-one at the epochs of a unit Poisson process  $(\Pi_t)_{t \geq 0}$ , independent of  $W_k$ 's. One can check that

$$\mathbb{E}(L_{\Pi_t} | (W_k)_{k \in \mathbb{N}}) = \sum_{k=1}^{\infty} \left( e^{-tW_1 \dots W_{k-1}(1-W_k)} - e^{-tW_1 \dots W_{k-1}} \right).$$

Recalling definitions (5),(6) and setting  $\varphi(t) := \mathbb{E}e^{-t(1-W)}$ ,  $U(x) := \mathbb{E}N_x = \sum_{k=1}^{\infty} \mathbb{P}\{S_{k-1} \leq x\}$ , we have

$$\begin{aligned} \mathbb{E}L_{\Pi_t} &= \mathbb{E} \sum_{k=1}^{\infty} \left( \varphi(te^{-S_{k-1}}) - \exp(-te^{-S_{k-1}}) \right) \\ &= \int_0^{\infty} \left( \varphi(te^{-x}) - \exp(-te^{-x}) \right) U(dx) \end{aligned} \tag{11}$$

$$= \sum_{k=1}^{\infty} (-1)^{k+1} \frac{t^k}{k!} \frac{(1 - \mathbb{E}(1 - W)^k)}{(1 - \mathbb{E}W^k)}, \tag{12}$$

where the familiar formula for Laplace transform of the potential measure,

$$\int_0^{\infty} e^{-sx} U(dx) = \frac{1}{1 - \mathbb{E}W^s}, \quad s > 0,$$

has been utilised. Note that (12) is an obvious counterpart of (9).

Set  $K(t) = \varphi(e^t) - \exp(-e^t)$ ,  $t \in \mathbb{R}$ . Since  $\nu = \infty$  and

$$\int_0^{\infty} \frac{e^{-z(1-W)} - e^{-z}}{z} dz = |\log(1 - W)|,$$

we conclude that

$$\lim_{t \rightarrow \infty} \int_{-\infty}^t K(z) dz = \infty. \tag{13}$$

Applying a minor extension of [29, Theorem 5] to the equality

$$\mathbb{E}L_{\Pi_{e^t}} = \int_0^{\infty} K(t - x) U(dx), \tag{14}$$

which is equivalent to (11), yields

$$\mathbb{E}L_{\Pi_{e^t}} \sim \frac{1}{\mu} \int_0^t \varphi(e^x) dx \sim \frac{1}{\mu} \int_1^{e^t} \frac{\varphi(x)}{x} dx.$$

The asymptotics of  $\mathbb{E}L_n$  is now obtained by the depoissonisation Lemma 8.2 in Appendix. The lemma is applicable because  $\mathbb{E}L_{\Pi(t)}$  is slowly varying. Indeed, slow variation of  $\int_1^t \varphi(u) du/u$  is checked straightforwardly from  $\varphi(t) \downarrow 0$  and the divergence of the integral for  $t = \infty$ .  $\square$

Similarly to the above, the proof of the next theorem is based on the poissonisation technique.

**Theorem 7.2** [16] *If  $\mu < \infty$  and  $\nu < \infty$  then  $L_n \rightarrow_d L_{\infty}$  as  $n \rightarrow \infty$  for some random variable  $L_{\infty}$  whose distribution satisfies*

$$\mathbb{P}\{L_{\infty} \geq i\} = \frac{1}{\mu} \sum_{j=1}^{\infty} \frac{\mathbb{E}W^j}{j} \mathbb{P}\{L_j = i\}, \quad i \in \mathbb{N}.$$

Moreover, the convergence of all moments holds, i.e.  $\mathbb{E}L_n^k \rightarrow \mathbb{E}L_{\infty}^k < \infty$  for  $k \in \mathbb{N}$ .

It is also known that if  $\mu < \infty$  and  $\nu = \infty$  then  $L_n \rightarrow_d \infty$  (see [17]), and that  $L_n \rightarrow_P 0$  if  $\nu < \infty$  and  $\mu = \infty$ . In the cases not covered by these results the question about the weak convergence of  $L_n$  is open.

Note that Theorem 7.2 only gives implicit specification of the limit law through distributions of  $L_n$ 's, which are not easy to determine, with one remarkable exception. Obviously from the recursive construction of the Bernoulli sieve, the distribution of  $L_1$  is geometric with parameter  $\mathbb{E}W$ . Curiously, the same is true for all  $n$  provided the law of  $W$  is symmetric about the midpoint  $1/2$ .

**Proposition 7.1** *If  $W =_d 1 - W$  then  $L_n$  is geometrically distributed with parameter  $1/2$  for all  $n \in \mathbb{N}$ .*

**Proof:** The argument is based on the recursion (3) for marginal distributions of the  $L_n$ 's. The symmetry  $W =_d 1 - W$  yields  $\mathbb{E}W^k = \mathbb{E}(1 - W)^k$  for all  $k \in \mathbb{N}$  and

$$\mathbb{P}\{Q_n^*(1) = n\} = \mathbb{P}\{Q_n^*(1) = 0\} \quad (15)$$

for all  $n \in \mathbb{N}$ . We will show by induction on  $n$  that  $\mathbb{P}\{L_n = k\} = 2^{-k-1}$  for all  $k \in \mathbb{N}_0$ . Using (3) and (15) we obtain

$$\begin{aligned} \mathbb{P}\{L_n = 0\} &= \mathbb{P}\{Q_n^*(1) = 0\} + \sum_{k=1}^{n-1} \mathbb{P}\{L_k = 0\} \mathbb{P}\{Q_n^*(1) = k\} \\ &= \mathbb{P}\{Q_n^*(1) = 0\} + \frac{1}{2} \left(1 - 2\mathbb{P}\{Q_n^*(1) = 0\}\right) = \frac{1}{2}, \end{aligned}$$

by the induction hypothesis. Assuming now that  $\mathbb{P}\{L_n = i\} = 2^{-i-1}$  for all  $i < k$  we have

$$\begin{aligned} \mathbb{P}\{L_n = k\} &= \sum_{j=1}^{n-1} \mathbb{P}\{Q_n^*(1) = j\} \mathbb{P}\{L_j = k\} + \mathbb{P}\{Q_n^*(1) = n\} \mathbb{P}\{L_n = k - 1\} \\ &= 2^{-k-1} \left(1 - 2\mathbb{P}\{Q_n^*(1) = 0\}\right) + \mathbb{P}\{Q_n^*(1) = 0\} 2^{-k} = 2^{-k-1}, \end{aligned}$$

and the proof is complete.  $\square$

Alternatively, one can use a representation of  $L_n$  through the sojourns of the Markov chain  $Q_n^*$  in positive states. Indeed, recall that  $L_1$  has geometric distribution with parameter  $\mathbb{E}W$ . Then using (15) and induction it can be checked that the distribution of  $L_n$  does not depend on  $n \geq 1$ .

## 8 Appendix.

For ease of reference we include a result due to Toeplitz and Schur (see [20], Theorem 2 on p. 43 and Theorem 9 on p. 52). We rewrite it in a form suitable for our purposes.

**Lemma 8.1** *Let  $\{s_n, n \in \mathbb{N}\}$  be any sequence of real numbers and let  $\{c_{nm}, n, m \in \mathbb{N}\}$  be a nonnegative array. Define another sequence  $\{t_n, n \in \mathbb{N}\}$  by  $t_n = \sum_{m=1}^n c_{nm} s_m$ . If*

(i)  $\lim_{n \rightarrow \infty} c_{nm} = 0$  for all  $m$ ,

(ii)  $\lim_{n \rightarrow \infty} \sum_{m=1}^n c_{nm} = 1$ ,

then

$$\liminf_{n \rightarrow \infty} s_n \leq \liminf_{n \rightarrow \infty} t_n \leq \limsup_{n \rightarrow \infty} t_n \leq \limsup_{n \rightarrow \infty} s_n \leq +\infty.$$

Now we address the issue of depoissonisation. Recall that the function  $\mathbb{E}L_{\Pi(t)}$  is slowly varying.

**Lemma 8.2** *If  $\lim_{t \rightarrow \infty} \mathbb{E}L_{\Pi_t} = +\infty$  then  $\mathbb{E}L_n \sim \mathbb{E}L_{\Pi_n}$ , as  $n \rightarrow \infty$ .*

**Proof:** For any fixed  $\varepsilon \in (0, 1)$ ,

$$\mathbb{E}L_{\Pi_t} = \mathbb{E}L_{\Pi_t} 1_{\{|\Pi_t - t| > \varepsilon t\}} + \mathbb{E}L_{\Pi_t} 1_{\{|\Pi_t - t| \leq \varepsilon t\}} =: A(t) + B(t).$$

Sublinearity of  $\mathbb{E}L_{\Pi_t}$  and the elementary large deviation bound for the Poisson distribution [2],

$$\mathbb{P}\{|\Pi_t - t| > \varepsilon t\} < c_1 e^{-c_2 t}, \quad t > 0 \quad (16)$$

with some  $c_1, c_2 > 0$ , yield  $A(t) \rightarrow 0$ .

It remains to evaluate  $B(t)$ . Since both  $M_n$  and  $K_n$  are non-decreasing, we have

$$B(t) = \mathbb{E}(M_{\Pi_t} - K_{\Pi_t}) 1_{\{|\Pi_t - t| \leq \varepsilon t\}} \leq \mathbb{E}L_{[(1-\varepsilon)t]} + \mathbb{E}(M_{[(1+\varepsilon)t]} - M_{[(1-\varepsilon)t]}).$$

Similarly,  $B(t) \geq \mathbb{E}L_{[(1+\varepsilon)t]} - \mathbb{E}(M_{[(1+\varepsilon)t]} - M_{[(1-\varepsilon)t]})$ . First step is to prove that

$$\lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \mathbb{E}(M_{[(1+\varepsilon)n]} - M_n) = 0. \quad (17)$$

Recalling the notation  $T_n = \max(E_1, \dots, E_n)$  and using the subadditivity of  $U(\cdot)$  we obtain

$$\begin{aligned} D(n) &:= \mathbb{E}\left(M_{[(1+\varepsilon)n]} - M_n\right) = \mathbb{E}\left(U(T_{[(1+\varepsilon)n]}) - U(T_n)\right) \\ &\leq \mathbb{E}U(T_{[(1+\varepsilon)n]} - T_n) 1_{\{T_{[(1+\varepsilon)n]} > T_n\}}. \end{aligned}$$

An appeal to estimate  $U(x) < ax + b$  (with some  $a, b > 0$ ) allows us to conclude that

$$D(n) \leq \mathbb{E}\left(a(T_{[(1+\varepsilon)n]} - T_n) + b\right) 1_{\{T_{[(1+\varepsilon)n]} > T_n\}} \leq a(H_{[(1+\varepsilon)n]} - H_n) + b\mathbb{P}\{T_{[(1+\varepsilon)n]} > T_n\},$$

where the equality  $\mathbb{E}T_n = H_n := \sum_{k=1}^n \frac{1}{k}$  has been utilised. Since

$$\lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} (H_{[(1+\varepsilon)n]} - H_n) = 0$$

and by exchangeability

$$\mathbb{P}\{T_{[(1+\varepsilon)n]} > T_n\} = 1 - \frac{n}{[(1+\varepsilon)n]} \rightarrow \varepsilon,$$

as  $n \rightarrow \infty$ , we arrive at (17).

We are ready to finish the proof. Divide the inequality

$$\mathbb{E}L_{\Pi_n/(1-\varepsilon)} \leq A(n/(1-\varepsilon)) + \mathbb{E}L_n + \mathbb{E}(M_{[\frac{1+\varepsilon}{1-\varepsilon}n]} - M_n), \quad (18)$$

by  $\mathbb{E}L_{\Pi_n}$ . Letting  $n \rightarrow \infty$  then  $\varepsilon \rightarrow 0$  and using the slow variation of  $\mathbb{E}L_{\Pi_n}$  we obtain

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}L_n}{\mathbb{E}L_{\Pi_n}} \geq 1.$$

The upper bound follows in the same way from the inequality

$$\mathbb{E}L_{\Pi_n/(1+\varepsilon)} \geq \mathbb{E}L_n - \mathbb{E}(M_n - M_{\lfloor \frac{1-\varepsilon}{1+\varepsilon}n \rfloor}). \quad (19)$$

□

## References

- [1] R. Arratia, A. Barbour, and S. Tavaré. *Logarithmic combinatorial structures*. European Mathematical Society, 2003.
- [2] R. R. Bahadur. Some limit theorems in statistics. *CBMS Regional conference series in applied mathematics. Philadelphia: SIAM*, 4, 1971.
- [3] Y. Baryshnikov, B. Eisenberg, and G. Stengle. A necessary and sufficient condition for the existence of the limiting probability of a tie for first place. *Statistics and Probability Letters*, 23(3):203–209, 1995.
- [4] J. Brands, F. Steutel, and R. Wilms. On the number of maxima in a discrete sample. *Statistics and Probability Letters*, 20:209–217, 1994.
- [5] F. T. Bruss and R. Grübel. On the multiplicity of the maximum in a discrete random sample. *Ann. Appl. Probab.*, 13(4):1252–1263, 2003.
- [6] F. T. Bruss and C. A. O’Cinneide. On the maximum and its uniqueness for geometric random samples. *J. Appl. Probab.*, 27:598–610, 1990.
- [7] C. F. Dunkl. The absorption distribution and the  $q$ -binomial theorem. *Communications in Statistics - Theory and Methods*, 10(19):1915–1920, 1981.
- [8] E. B. Dynkin. Some limit theorems for sums of independent random variables with infinite mathematical expectations. *Selected Transl. in Math. Statist. and Probability*, 1:171–189, 1961.
- [9] B. Eisenberg, G. Stengle, and G. Strang. The asymptotic probability of a tie for first place. *Ann. Appl. Probab.*, 3:731–745, 1993.
- [10] K. B. Erickson. Strong renewal theorems with infinite mean. *Trans. Amer. Math. Soc.*, 151:263–291, 1970.
- [11] J. Fill, H. Mahmoud, and W. Szpankowski. On the distribution for the duration of a randomized leader election algorithm. *Ann. Appl. Probab.*, 6:1260–1283, 1996.
- [12] A. Gnedin. The Bernoulli sieve. *Bernoulli*, 10:79–96, 2004.

- [13] A. Gnedin. Regeneration in random combinatorial structures. *Probability Surveys*, 7:105–156, 2010.
- [14] A. Gnedin, A. Hansen, and J. Pitman. Notes on the occupancy problem with infinitely many boxes: general asymptotics and power laws. *Probability Surveys*, 4:146–171, 2007.
- [15] A. Gnedin, A. Iksanov, and A. Marynych. Limit theorems for the number of occupied boxes in the Bernoulli sieve. *Submitted to Theory of Stochastic Processes*, 2010.
- [16] A. Gnedin, A. Iksanov, P. Negadajlov, and U. Rösler. The Bernoulli sieve revisited. *Ann. Appl. Probab.*, 19:1634–1655, 2009.
- [17] A. Gnedin, A. Iksanov, and U. Rösler. Small parts in the Bernoulli sieve. *Discrete Mathematics and Theoretical Computer Science*, Proceedings Series Volume AI:239–246, 2008.
- [18] A. Gnedin and J. Pitman. Regenerative composition structures. *Ann. Probab.*, 33:445–479, 2005.
- [19] W. M. Y. Goh and P. Hitczenko. Gaps in samples of geometric random variables. *Discrete Mathematics*, 22:2871–2890, 2007.
- [20] G. H. Hardy. *Divergent Series*. AMS Bookstore, 2000.
- [21] P. Hitczenko and A. Knopfmacher. Gap-free compositions and gap-free samples of geometric random variables. *Discrete Mathematics*, 294(3):225–239, 2005.
- [22] S. Janson and W. Szpankowski. Analysis of an asymmetric leader election algorithm. *Electron. J. Combin.*, 4(1), Art. #R17, 1997.
- [23] S. Karlin. Central limit theorems for certain infinite urn schemes. *J. Math. Mech.*, 17:373–401, 1967.
- [24] A. W. Kemp. Absorption sampling and the absorption distribution. *J. Appl. Probab.*, 35:489–494, 1998.
- [25] P. Kirschenhofer and H. Prodinger. The number of winners in a discrete geometrically distributed sample. *Ann. Appl. Probab.*, 6:687–694, 1996.
- [26] G. Louchard and H. Prodinger. On gaps and unoccupied urns in sequence of geometrically distributed random variables. *Discrete Mathematics*, 308(9):1538–1562, 2008.
- [27] G. Louchard and H. Prodinger. The asymmetric leader election algorithm: another approach. *Annals of Combinatorics*, 12(4):449–478, 2009.
- [28] H. Prodinger. How to select a loser. *Discrete Math.*, 120:149–159, 1993.
- [29] M. S. Sgibnev. Renewal theorem in the case of an infinite variance. *Siberian Math. J.*, 22:787–796, 1981.

