



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 1999

On Poisson-Dirichlet limits for random decomposable combinatorial structures

Arratia, R ; Barbour, A D ; Tavaré, S

Abstract: We prove a joint local limit law for the distribution of the r largest components of decomposable logarithmic combinatorial structures, including assemblies, multisets and selections. Our method is entirely probabilistic, and requires only weak conditions that may readily be verified in practice.

DOI: <https://doi.org/10.1017/S0963548399003788>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-22118>

Journal Article

Originally published at:

Arratia, R; Barbour, A D; Tavaré, S (1999). On Poisson-Dirichlet limits for random decomposable combinatorial structures. *Combinatorics, Probability Computing*, 8(3):193-208.

DOI: <https://doi.org/10.1017/S0963548399003788>

On Poisson–Dirichlet Limits for Random Decomposable Combinatorial Structures

RICHARD ARRATIA^{1†}, A.D. BARBOUR^{2‡}
and SIMON TAVARÉ^{1†}

¹ Department of Mathematics, University of Southern California,
Los Angeles, CA 90089-1113, USA
(e-mail: rarratia@math.usc.edu stavare@gnome.usc.edu)

² Abteilung für Angewandte Mathematik, Universität Zürich,
Winterthurerstrasse 190, CH-8057, Zürich, Switzerland
(e-mail: adb@amath.unizh.ch)

Received 20 June 1997; revised 16 March 1998

We prove a joint local limit law for the distribution of the r largest components of decomposable logarithmic combinatorial structures, including assemblies, multisets and selections. Our method is entirely probabilistic, and requires only weak conditions that may readily be verified in practice.

1. Introduction

The proportion of integers in the largest, second largest, ... cycles of a random permutation of n objects have, asymptotically as $n \rightarrow \infty$, the Poisson–Dirichlet distribution $PD(\theta)$, with parameter $\theta = 1$. The result for the largest cycle appears in Goncharov [8], that for the k th largest in Shepp and Lloyd [17], and the joint distributional result in Kingman [12] and Vershik and Schmidt [19]. Related limit laws are now known for a variety of decomposable combinatorial structures. For example, the relative sizes of the largest, second largest, ... components of a random mapping have asymptotically the $PD(\theta)$ law with $\theta = 1/2$. Convergence for the marginal distributions appears in Kolchin [14, 15], and for the joint distribution in Aldous [1]. An analogous result holds for the ordered degree sequence of the factors of a polynomial over $GF(q)$; in this case $\theta = 1$, just as for permutations [3].

† Supported in part by NSF grant DMS 90-05833 and DMS 96-26412.

‡ Supported in part by Schweizerischer NF Projekt Nr 20-43453.95.

More recently, Hansen [9] provided a unified approach to the asymptotics of the order statistics of the component sizes of multisets and assemblies, which are families of decomposable combinatorial structures. She establishes weak convergence to the PD(θ) distribution, and shows how to identify the appropriate value of θ . Her arguments rely in part on complex analysis, and involve conditions on generating functions which are not always easy to verify.

In this paper we establish a local limit theorem for the joint distribution of the large components, refining previous results, under very weak conditions. Our method uses probabilistic, as opposed to complex analytic, arguments, and the conditions become correspondingly transparent. We consider randomly chosen decomposable combinatorial structures of total size n whose component counts $(C_1(n), \dots, C_n(n))$, where $C_i(n)$ denotes the number of components of size i , have joint distribution determined by the conditioning relation

$$(CR) \quad \mathcal{L}(C_1(n), \dots, C_n(n)) = \mathcal{L}(Z_1, \dots, Z_n \mid T_n = n), \quad (1.1)$$

where $(Z_i, i \geq 1)$ are independent random variables over \mathbb{Z}_+ , and $T_n = \sum_{i=1}^n iZ_i$. Such structures include assemblies (for which the Z_i are Poisson-distributed), multisets (for which the Z_i are negative binomially distributed), and selections (for which the Z_i are binomially distributed): see [5]. We require only that the Z_i satisfy the logarithmic condition LC:

$$(LC) \quad \lim_{i \rightarrow \infty} i\mathbb{P}(Z_i = 1) = \theta = \lim_{i \rightarrow \infty} i\mathbb{E}Z_i \quad (1.2)$$

for some $\theta > 0$, together with the additional tail condition

$$\mathbb{P}(Z_i \geq 2) = O(i^{-2}), \quad (1.3)$$

conditions weaker than those in [9]. In fact, (1.3) is implied by (1.2) for the examples we consider in Section 5. More general structures can also be analysed, under the mild extra condition given in (5.14), since the approach also ties in with arguments using Stein's method, discussed in detail in [4].

The simplest example of such a construct is that in which the Z_i are Poisson-distributed with means *exactly* θ/i . In this case, the $C_i(n)$ have as joint distribution the Ewens Sampling Formula $\text{ESF}_n(\theta)$ given in (2.1). In Section 2, we collect known facts about PD(θ), $\text{ESF}_n(\theta)$, and the limit distribution P_θ of $n^{-1}T_n$ under $\text{ESF}_n(\theta)$. The local limit approximation to the joint distribution of the sizes of the large cycles by PD(θ), when the $C_i(n)$ are distributed according to $\text{ESF}_n(\theta)$, is then established in Section 3. The argument used is readily generalized in Section 4 to arbitrary combinatorial structures that satisfy (LC), and for which T_n/n admits the local limit approximation (LLA) given in (4.6). In Section 5 it is shown that (LLA) holds for assemblies, multisets and selections satisfying (LC).

2. The Ewens Sampling Formula

Let \mathcal{S}_n denote the set of permutations of $\{1, 2, \dots, n\}$. We write $\pi \in \mathcal{S}_n$ as an ordered product of cycles. The integer 1 starts the first cycle, followed by the image of 1, the image of that point and so on. The smallest integer not in the first cycle begins the second cycle, followed by its images. In this way, π is decomposed into an ordered product of cycles.

We consider random permutations with distribution determined by

$$\mathbb{P}(\pi) = \frac{\theta^{|\pi|}}{\theta_{(n)}}, \pi \in \mathcal{S}_n,$$

where $|\pi|$ denotes the number of cycles in π , $\theta \in (0, \infty)$, and

$$\theta_{(n)} = \theta(\theta + 1) \cdots (\theta + n - 1).$$

Let $C^{(n)} \equiv (C_1(n), C_2(n), \dots, C_n(n))$ be the counts of cycles of sizes $1, 2, \dots, n$ in such a θ -biased random permutation of size n . The distribution of $C^{(n)}$ is given by the Ewens Sampling Formula [7] $\text{ESF}_n(\theta)$: for any $a \in \mathbb{Z}_+^n$,

$$\mathbb{P}(C^{(n)} = a) = \mathbb{1}\left(\sum_{i=1}^n ia_i = n\right) \frac{n!}{\theta_{(n)}} \prod_{j=1}^n \left(\frac{\theta}{j}\right)^{a_j} \frac{1}{a_j!}, \tag{2.1}$$

$\mathbb{1}(A)$ denoting the indicator of A . Let Z_1, Z_2, \dots be independent Poisson random variables with means $\mathbb{E}Z_i = \theta/i, i = 1, 2, \dots, n$, and let $Z[1, n] = (Z_1, \dots, Z_n)$. It is well known [21, 2] that

$$\mathbb{P}(C^{(n)} = a) = \mathbb{P}(Z[1, n] = a | T_n = n), \tag{2.2}$$

where

$$T_n = Z_1 + 2Z_2 + \cdots + nZ_n,$$

so we have a combinatorial structure satisfying (1.1), with Z_i having a Poisson distribution with mean θ/i , for which the conditions (LC) and (1.3) are clearly satisfied. Understanding the asymptotic behaviour of (2.2) requires knowledge of the asymptotic behaviour of T_n , which we now review.

2.1. The limit T of T_n/n

The density of T_n satisfies the recursion

$$k\mathbb{P}(T_n = k) = \theta \sum_{j=1}^n \mathbb{P}(T_n = k - j), \quad k = 1, 2, \dots \tag{2.3}$$

It follows from this that

$$k\mathbb{P}(T_n = k) = (k - 1 + \theta)\mathbb{P}(T_n = k - 1), \quad k = 1, 2, \dots, n,$$

so that

$$\mathbb{P}(T_n = k) = \frac{\theta_{(k)}}{k!} \mathbb{P}(T_n = 0) = \exp(-\theta h(n)) \frac{\theta_{(k)}}{k!}, \quad k \leq n,$$

where $h(n) = \sum_{j=1}^n 1/j$. Hence

$$\lim_{n \rightarrow \infty} n\mathbb{P}(T_n = k) = \frac{e^{-\gamma\theta} x^{\theta-1}}{\Gamma(\theta)} \quad \text{if } k \leq n, \quad k/n \rightarrow x \in (0, 1]. \tag{2.4}$$

Using (2.3) and (2.4), we conclude that

$$\lim_{n \rightarrow \infty} \mathbb{P}(T_n/n \leq x) = \frac{x^\theta e^{-\gamma\theta}}{\Gamma(\theta + 1)} \tag{2.5}$$

if $x \in (0, 1]$. In fact, $\lim_{n \rightarrow \infty} \mathbb{P}(T_n/n \leq x)$ exists for all $x > 0$.

Theorem 2.1. As $n \rightarrow \infty$, the random variable T_n/n converges in distribution to a random variable T whose distribution P_θ has Laplace transform given by

$$\mathbb{E}e^{-sT} = \exp\left(-\int_0^1 (1 - e^{-sx}) \frac{\theta}{x} dx\right). \quad (2.6)$$

Proof. Let μ_n be the measure that puts mass n^{-1} at points in^{-1} , $i = 1, 2, \dots, n$ and note that μ_n converges weakly to Lebesgue measure. The Laplace transform of the random variable T_n/n is

$$\begin{aligned} \mathbb{E}e^{-sT_n/n} &= \exp\left(-\sum_{i=1}^n \frac{\theta}{i} (1 - e^{-si/n})\right) \\ &= \exp\left(-\int_0^1 (1 - e^{-sx}) \frac{\theta}{x} \mu_n(dx)\right) \\ &\rightarrow \exp\left(-\int_0^1 (1 - e^{-sx}) \frac{\theta}{x} dx\right), \end{aligned}$$

the last step following by dominated convergence. \square

It follows from (2.5) that the density g_θ of T satisfies

$$g_\theta(x) = \frac{e^{-\gamma\theta} x^{\theta-1}}{\Gamma(\theta)}, \quad 0 \leq x \leq 1,$$

so that

$$g_\theta(1) = \frac{e^{-\gamma\theta}}{\Gamma(\theta)}. \quad (2.7)$$

An expression for $g_\theta(x)$ for $x > 1$ is given in [20]; it satisfies the integral equation

$$xg_\theta(x) = \theta \int_{x-1}^x g_\theta(u) du, \quad x > 0, \quad (2.8)$$

with $g_\theta(x) = 0$ if $x < 0$.

2.2. The Poisson–Dirichlet distribution

The Poisson–Dirichlet distribution, denoted by $\text{PD}(\theta)$, was defined by Kingman [11] to be the distribution of the normalized points $\sigma_1 > \sigma_2 > \dots$ of a Poisson process with intensity $\theta e^{-x}/x$, $x > 0$:

$$\text{PD}(\theta) = \mathcal{L}((\sigma_1/\sigma, \sigma_2/\sigma, \dots)), \quad (2.9)$$

where $\sigma = \sigma_1 + \sigma_2 + \dots$. Other representations of the Poisson–Dirichlet distribution may be found in [13, 4], for example.

The density $f_\theta^{(r)}$ of the first r coordinates of $\text{PD}(\theta)$ was found by Watterson [22] in the form

$$f_\theta^{(r)}(x_1, \dots, x_r) = \frac{e^{\gamma\theta} \theta^r \Gamma(\theta) x_r^{\theta-1}}{x_1 x_2 \cdots x_r} g_\theta\left(\frac{1 - x_1 - \cdots - x_r}{x_r}\right), \quad (2.10)$$

for $r \geq 1$, and x_1, \dots, x_r satisfying $0 < x_r < \cdots < x_1 < 1$ and $0 < x_1 + \cdots + x_r < 1$.

The Poisson–Dirichlet distribution with parameter θ arises as the limit law of the renormalized sizes $L_1(n), L_2(n), \dots$ of the largest, second largest, ... cycles of a θ -biased permutation.

Theorem 2.2 ([12]). *As $n \rightarrow \infty$,*

$$n^{-1}(L_1(n), L_2(n), \dots) \Rightarrow (L_1, L_2, \dots),$$

where (L_1, L_2, \dots) has the $PD(\theta)$ distribution.

3. A local limit law for large cycles under $ESF_n(\theta)$

3.1. Point probabilities for T_n

Theorem 2.1 extends the convergence in (2.5) from $x \in (0, 1]$ to all $x > 0$. This suggests that the same extension may also be feasible for (2.4), as is shown in the following lemma.

Lemma 3.1. *Suppose that $m/n \rightarrow y \in (0, \infty)$ as $n \rightarrow \infty$. Then*

$$\lim_{n \rightarrow \infty} n\mathbb{P}(T_n = m) = g_\theta(y). \tag{3.1}$$

Proof. Equation (2.3) gives

$$m\mathbb{P}(T_n = m) = \theta\mathbb{P}(m - n \leq T_n < m).$$

Multiplying by n/m and using (2.8) and the fact that $T_n/n \Rightarrow T$, which has a continuous distribution function, shows that

$$\begin{aligned} \lim_{n \rightarrow \infty} n\mathbb{P}(T_n = m) &= \frac{\theta}{y}\mathbb{P}(y - 1 \leq T \leq y) \\ &= g_\theta(y), \end{aligned}$$

completing the proof. □

The next result uses elementary arguments to derive bounds for the point probabilities $\mathbb{P}(T_{bn} = m)$, where

$$T_{bn} = \sum_{j=b+1}^n jZ_j, \quad 0 \leq b < n.$$

Lemma 3.2. *Write $\bar{\theta} = \min(1, \theta)$. Then*

$$\max_{k \geq 0} \mathbb{P}(T_{bn} = k) \leq e^{-\bar{\theta}(h(n) - h(b))}. \tag{3.2}$$

Proof. First consider the case $0 < \theta \leq 1$. We use the fact ([5]) that the point probabilities for T_{bn} satisfy

$$m\mathbb{P}(T_{bn} = m) = \theta\mathbb{P}(m - n \leq T_{bn} < m - b). \tag{3.3}$$

Hence, for $m \geq 1$,

$$\mathbb{P}(T_{bn} = m) \leq \frac{1}{m} \sum_{j=0}^{m-1} \mathbb{P}(T_{bn} = j).$$

Thus $\mathbb{P}(T_{bn} = m)$ is at most the average of the previous m values, and so, by induction, $\max_{k \geq 0} \mathbb{P}(T_{bn} = k) \leq \mathbb{P}(T_{bn} = 0) = e^{-\theta(h(n)-h(b))}$.

For the case $\theta > 1$, let $\tilde{Z}_j, j \geq 1$ be independent Poisson random variables with $\mathbb{E}\tilde{Z}_j = 1/j$, and define $\tilde{T}_{bn} = \sum_{j=b+1}^n j\tilde{Z}_j$. Define $T'_{bn} = \sum_{j=b+1}^n jZ'_j$, where the Z'_j are independent Poisson random variables with mean $(\theta - 1)/j$, independent of the \tilde{Z}_j . Then we can write

$$T_{bn} = \tilde{T}_{bn} + T'_{bn},$$

with independent summands. It follows that

$$\begin{aligned} \mathbb{P}(T_{bn} = m) &= \sum_{j=0}^m \mathbb{P}(\tilde{T}_{bn} = j) \mathbb{P}(T'_{bn} = m - j) \\ &\leq \max_{0 \leq j \leq m} \mathbb{P}(\tilde{T}_{bn} = j) \\ &\leq e^{-(h(n)-h(b))}, \end{aligned}$$

the last step following from the case proved earlier. \square

3.2. The local limit theorem

In this section we derive a joint local limit law for the distribution of the r largest cycle lengths $L_1(n), \dots, L_r(n)$ under $\text{ESF}_n(\theta)$.

Theorem 3.3. *For $r \geq 1$, suppose that $0 < x_r < x_{r-1} < \dots < x_1 < 1$ satisfy $0 < x_1 + \dots + x_r < 1$. Then*

$$\lim_{n \rightarrow \infty} n^r \mathbb{P}(L_i(n) = \lfloor nx_i \rfloor, 1 \leq i \leq r) = f_\theta^{(r)}(x_1, \dots, x_r), \quad (3.4)$$

where the density $f_\theta^{(r)}$ is given in (2.10).

Proof. First assume that integers m_1, m_2, \dots, m_r satisfy the conditions

$$1 \leq m_r < m_{r-1} < \dots < m_1 < n, \quad m \equiv m_1 + \dots + m_r \leq n,$$

and let $A_n(C^{(n)}) = A_n(C^{(n)}; m_1, m_2, \dots, m_{r-1}, m_r)$ denote the event

$$\left\{ \begin{aligned} C_n(n) = 0, \dots, C_{m_1+1}(n) = 0, C_{m_1}(n) = 1, C_{m_1-1}(n) = 0, \dots, \\ C_{m_2+1}(n) = 0, C_{m_2}(n) = 1, C_{m_2-1}(n) = 0, \dots, C_{m_{r-1}+1}(n) = 0, \\ C_{m_{r-1}}(n) = 1, C_{m_{r-1}-1}(n) = 0, \dots, C_{m_r+1}(n) = 0 \end{aligned} \right\}.$$

Then

$$\mathbb{P}(L_1(n) = m_1, \dots, L_r(n) = m_r) = \mathbb{P}(A_n(C^{(n)}), C_m(n) \geq 1).$$

This last probability can be written

$$\mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = 1) + \sum_{l \geq 2} \mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = l).$$

The first term is, using (2.2),

$$\mathbb{P}(A_n(Z), Z_{m_r} = 1 | T_n = n) = \mathbb{P}(A_n(Z)) \mathbb{P}(Z_{m_r} = 1) \frac{\mathbb{P}(T_{m_r-1} = n - m)}{\mathbb{P}(T_n = n)}, \tag{3.5}$$

which reduces to

$$\frac{\mathbb{P}(T_{m_r-1} = n - m)}{\mathbb{P}(T_n = n)} \frac{\theta^r e^{-\theta(h(n) - h(m_r-1))}}{m_1 \cdots m_r} \tag{3.6}$$

Applying the result of Lemma 3.1 and simplifying shows that

$$\lim_{n \rightarrow \infty} n^r \mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = 1) = f_\theta^{(r)}(x_1, \dots, x_r).$$

It remains to show that $\sum_{l \geq 2} \mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = l) = o(n^{-r})$. But this probability is just

$$\begin{aligned} \mathbb{P}(A_n(Z)) \sum_{l \geq 2} \mathbb{P}(Z_{m_r} = l) \frac{\mathbb{P}(T_{m_r-1} = n - m - (l - 1)m_r)}{\mathbb{P}(T_n = n)} \\ \leq \mathbb{P}(A_n(Z)) \frac{e^{-\bar{\theta}h(m_r-1)}}{\mathbb{P}(T_n = n)} \mathbb{P}(Z_{m_r} \geq 2), \end{aligned}$$

using Lemma 3.2. Since $\mathbb{P}(T_n = n) \sim n^{-1}g_\theta(1)$, $\mathbb{P}(A_n(Z)) \leq \theta^{r-1}/(m_1 \cdots m_{r-1})$, and $\mathbb{P}(Z_{m_r} \geq 2) \leq \theta^2/(2m_r^2)$, we see that this term is of order $O(n^{-r-1} \cdot n \cdot n^{-\bar{\theta}}) = O(n^{-r}n^{-\bar{\theta}}) = o(n^{-r})$, as required. \square

Remarks.

- (1) Theorem 2.2 follows from Theorem 3.3 using Scheffé’s Theorem [16].
- (2) It is crucial in the hypothesis $x_1 + \cdots + x_r < 1$ to have strict inequality. To see this, take $r = 1$ and $x_1 = 1$, and note that

$$n\mathbb{P}(L_1(n) = n) \sim \Gamma(\theta + 1)n^{1-\theta}.$$

4. Combinatorial structures

In this section, we show that a joint local limit law like that in Theorem 3.3 is true for a large class of decomposable combinatorial structures. $C_i(n)$ now denotes the number of components of size i , $i = 1, 2, \dots, n$, and we consider structures that satisfy the relation (2.2) for independent random variables Z_i taking values in \mathbb{Z}_+ . However, the Z_i no longer satisfy $Z_i \sim \text{Po}(\theta/i)$; instead, we merely require the ‘logarithmic condition’ (LC), repeated here for convenience:

$$\lim_{i \rightarrow \infty} i\mathbb{P}(Z_i = 1) = \theta = \lim_{i \rightarrow \infty} i\mathbb{E}Z_i \tag{LC}$$

for some $\theta \in (0, \infty)$, and the tail condition (1.3). Note that (LC) implies that

$$\tilde{\theta} \equiv \sup_{i \geq 1} i\mathbb{E}Z_i < \infty. \tag{4.1}$$

We begin with some preliminaries, the first of which requires no proof.

Lemma 4.1. *Suppose the Z_i satisfy (LC). Then as $i \rightarrow \infty$,*

$$\mathbb{P}(Z_i \geq 2) = o(i^{-1}), \quad (4.2)$$

and

$$\mathbb{P}(Z_i = 0) = 1 - \theta i^{-1} + o(i^{-1}). \quad (4.3)$$

We see from this that, for large i , the distribution of Z_i is indeed close to Poisson with mean θ/i .

Corollary 4.2. *Let Z_i^* be independent Poisson random variables with $\mathbb{E}Z_i^* = \theta/i$, $i \geq 1$. There is a sequence $\epsilon(i) \downarrow 0$ as $i \rightarrow \infty$ such that*

$$d_{TV}(Z_i, Z_i^*) \leq \epsilon(i)i^{-1}.$$

Proof.

$$\begin{aligned} 2d_{TV}(Z_i, Z_i^*) &= \sum_{j \geq 0} |\mathbb{P}(Z_i = j) - \mathbb{P}(Z_i^* = j)| \\ &\leq |\mathbb{P}(Z_i = 0) - \mathbb{P}(Z_i^* = 0)| + |\mathbb{P}(Z_i = 1) - \mathbb{P}(Z_i^* = 1)| \\ &\quad + \mathbb{P}(Z_i \geq 2) + \mathbb{P}(Z_i^* \geq 2). \end{aligned}$$

The result now follows from Lemma 4.1. □

In order to prove a local limit theorem for the r largest component sizes $L_1(n), \dots, L_r(n)$, analogous to Theorem 3.3, we use the same recipe. The first ingredient is the counterpart of Theorem 2.1; an alternative proof may be found in [5].

Theorem 4.3. *For $i = 1, 2, \dots$, let Z_i be independent random variables taking values in \mathbb{Z}_+ and satisfying (LC). Then, as $n \rightarrow \infty$,*

$$n^{-1}T_n \Rightarrow T. \quad (4.4)$$

Proof. Let Z_i^* be independent Poisson random variables with $\mathbb{E}Z_i^* = \theta/i$, and write $T_{bn}^* = \sum_{j=b+1}^n jZ_j^*$, $Z^*(b, n) = (Z_{b+1}^*, \dots, Z_n^*)$. Corollary 4.2 shows that $d_{TV}(Z_i, Z_i^*) \leq \epsilon(i)i^{-1}$. Choose any sequence $b_n = o(n)$ such that $\epsilon(b_n) \log(n/b_n) \rightarrow 0$ as $n \rightarrow \infty$. Then we immediately find that

$$\begin{aligned} d_{TV}(T_{b_n n}, T_{b_n n}^*) &\leq d_{TV}(Z(b_n, n), Z^*(b_n, n)) \\ &\leq \sum_{j=b_n+1}^n \epsilon(j)j^{-1} \leq \epsilon(b_n) \log(n/b_n). \end{aligned} \quad (4.5)$$

Since $\mathbb{E}n^{-1}T_{b_n}^* = n^{-1}\theta b_n \rightarrow 0$, it follows that $n^{-1}T_{b_n}^* \Rightarrow 0$. Hence, since $n^{-1}T_n^* = n^{-1}T_{b_n}^* + n^{-1}T_{b_n n}^*$, $n^{-1}T_{b_n n}^* \Rightarrow T$ by Theorem 2.1, and it then follows from (4.5) that $n^{-1}T_{b_n n} \Rightarrow T$.

Finally,

$$n^{-1}\mathbb{E}T_{b_n} = n^{-1} \sum_{j=1}^{b_n} j\mathbb{E}Z_j \leq \tilde{\theta}n^{-1}b_n \rightarrow 0,$$

so that $n^{-1}T_n \Rightarrow T$ also. □

The second ingredient is a bound on point probabilities, echoing Lemma 3.2.

Theorem 4.4. *As $n \rightarrow \infty$, $\max_{k \geq 0} \mathbb{P}(T_n = k) \rightarrow 0$.*

Proof. As in the proof of Theorem 4.3, let Z_i^* be independent Poisson random variables with $\mathbb{E}Z_i^* = \theta/i$, and choose any sequence $b_n = o(n)$ such that (4.5) holds. Since $T_n = T_{b_n} + T_{b_n n}$ and the two summands are independent,

$$\max_{k \geq 0} \mathbb{P}(T_n = k) \leq \max_{k \geq 0} \mathbb{P}(T_{b_n n} = k).$$

Now, from (4.5),

$$\mathbb{P}(T_{b_n n} = k) \leq \mathbb{P}(T_{b_n n}^* = k) + \epsilon(b_n) \log(n/b_n),$$

and by Lemma 3.2, defining $\bar{\theta} = \min(1, \theta)$,

$$\max_{k \geq 0} \mathbb{P}(T_{b_n n}^* = k) \leq e^{-\bar{\theta}(h(n)-h(b_n))} = O((b_n/n)^{\bar{\theta}}).$$

Hence $\max_{k \geq 0} \mathbb{P}(T_{b_n n} = k) \rightarrow 0$ as $n \rightarrow \infty$. □

The final ingredient is that T_n should satisfy a local limit approximation:

$$(LLA) \quad n\mathbb{P}(T_n = m) \sim g_\theta(y), \quad \text{as } n \rightarrow \infty, m/n \rightarrow y \in (0, \infty). \quad (4.6)$$

The proofs of (LLA) are somewhat different for the various classes of combinatorial structure we consider, and a detailed treatment is given in the next section. Whenever (LC) and (LLA) hold, we have the following joint local limit law.

Theorem 4.5. *Suppose that a combinatorial structure $C^{(n)}$ has distribution given by (CR), and satisfies (LC), (LLA), and the tail condition (1.3). For $r \geq 1$, suppose that $0 < x_r < x_{r-1} < \dots < x_1 < 1$ satisfy $0 < x_1 + \dots + x_r < 1$. Then*

$$\lim_{n \rightarrow \infty} n^r \mathbb{P}(L_i(n) = \lfloor nx_i \rfloor, 1 \leq i \leq r) = f_\theta^{(r)}(x_1, \dots, x_r),$$

where the density $f_\theta^{(r)}$ is given in (2.10); hence also, as $n \rightarrow \infty$,

$$n^{-1}(L_1(n), L_2(n), \dots) \Rightarrow \text{PD}(\theta).$$

Proof. The proof mimics that of Theorem 3.3 down to (3.5). Expression (3.6) is now replaced by

$$\frac{\mathbb{P}(T_{m_{r-1}} = n - m)}{\mathbb{P}(T_n = n)} \prod_{i=m_r}^n \mathbb{P}(Z_i = 0) \prod_{s=1}^r \frac{\mathbb{P}(Z_{m_s} = 1)}{\mathbb{P}(Z_{m_s} = 0)}.$$

Under (LLA), the first term is asymptotic to $x_r^{-1}g_\theta((1-x_1-\cdots-x_r)/x_r)/g_\theta(1)$. Using (4.3), the first product is asymptotic to $x_r^{-\theta}$, while from (LC) the second product is asymptotic to $n^{-r}\theta^r x_1^{-1}\cdots x_r^{-1}$. Combining these terms and using (2.10) shows that

$$\lim_{n \rightarrow \infty} n^r \mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = 1) = f_\theta^{(r)}(x_1, \dots, x_r).$$

To show that $\sum_{l \geq 2} \mathbb{P}(A_n(C^{(n)}), C_{m_r}(n) = l) = o(n^{-r})$, note that the left side is just

$$\begin{aligned} \mathbb{P}(A_n(Z)) \sum_{l \geq 2} \mathbb{P}(Z_{m_r} = l) \frac{\mathbb{P}(T_{m_r-1} = n - m - (l-1)m_r)}{\mathbb{P}(T_n = n)} \\ \leq \frac{\mathbb{P}(A_n(Z))\mathbb{P}(Z_{m_r} \geq 2)}{\mathbb{P}(T_n = n)} \max_{k \geq 0} \mathbb{P}(T_{m_r-1} = k). \end{aligned}$$

Since

$$\mathbb{P}(A_n(Z)) \leq \mathbb{P}(Z_{m_1} = 1) \cdots \mathbb{P}(Z_{m_{r-1}} = 1) = O(n^{-(r-1)}),$$

and $\mathbb{P}(Z_{m_r} \geq 2) = O(n^{-2})$, we see from (LLA) that the first factor is of order n^{-r} , whereas the second term tends to 0 by Theorem 4.4. \square

5. Verifying the local limit approximation

This section is devoted to establishing (LLA) for a wide variety of combinatorial models. Once done, Theorem 4.5 can then be applied.

5.1. Assemblies

Random assemblies are decomposable combinatorial structures for which the counts $C_j(n)$ of components of size j satisfy the conditioning relation (2.2) for Poisson-distributed Z_j with means

$$\mathbb{E}Z_j \equiv \lambda_j = \frac{m_j x^j}{j!}, \text{ for some } x > 0.$$

In these models, the integers m_j are prescribed in advance, and, for a satisfying $a_1 + 2a_2 + \cdots + na_n = n$, the probabilities

$$\begin{aligned} \mathbb{P}(C^{(n)} = a) &= \mathbb{P}(Z[1, n] = a | T_n = n) \\ &= \mathbb{P}(Z[1, n] = a) / \mathbb{P}(T_n = n) \\ &= \frac{\prod_{i=1}^n (m_i x^i / i!)^{a_i} / a_i!}{\sum_{\{d_1+2d_2+\cdots+nd_n=n\}} \prod_{i=1}^n (m_i x^i / i!)^{d_i} / d_i!} \end{aligned}$$

are the same for any arbitrary value of x . Hence, to be in the logarithmic class, it is enough that $m_j \sim \theta(j-1)!y^j$ for some $y > 0$, since we can take $x = y^{-1}$. Condition (LC) then reduces to the requirement that $j\lambda_j \rightarrow \theta$, in which case

$$\mathbb{P}(Z_j \geq 2) = [1 - e^{-\lambda_j}(1 + \lambda_j)] \leq \frac{1}{2}\lambda_j^2 = O(j^{-2}),$$

so that the tail condition (1.3) is satisfied. Among the examples are permutations for which $m_j = (j-1)!, x = 1, \theta = 1$, and random mappings for which $m_j = (j-1)! \sum_{l=0}^{j-1} l! / l!, x = e^{-1}, \theta = 1/2$. Many other examples are described in [5]. We note in passing that the

following proofs make no use of the m_j being integers; an application in the more general setting appears in Section 5.4.

Most of the results depend on the analogue of equation (2.3) for the density of T_n . Using [5], it takes the form

$$k\mathbb{P}(T_n = k) = \sum_{j=1}^n \mathbb{P}(T_n = k - j)j\lambda_j, \quad k = 0, 1, \dots \tag{5.1}$$

with $j\lambda_j$ in the place of θ . Intuitively, this should make little difference for large n , because $j\lambda_j \rightarrow \theta$.

To verify (LLA), note that, according to equation (5.1),

$$\begin{aligned} k\mathbb{P}(T_n = k) &= \sum_{j=1}^n \mathbb{P}(T_n = k - j)j\lambda_j \\ &= \theta\mathbb{P}(k - n \leq T_n < k) + r_n(k), \end{aligned}$$

where

$$r_n(k) = \sum_{j=1}^n \mathbb{P}(T_n = k - j)(j\lambda_j - \theta).$$

The remainder of the proof follows just as for Lemma 3.1, but now using Theorem 4.3 instead of Theorem 2.1, if we can show that $|r_n(k)| \rightarrow 0$ as $n \rightarrow \infty$ when $k/n \rightarrow y > 0$. To do this, let $\epsilon > 0$ be arbitrary, and choose $j_0 = j_0(\epsilon)$ such that $|j\lambda_j - \theta| < \epsilon$ for all $j > j_0$. Then, for $n > j_0$,

$$\begin{aligned} |r_n(k)| &\leq \sum_{j=1}^{j_0} \mathbb{P}(T_n = k - j)|j\lambda_j - \theta| + \epsilon \sum_{j > j_0} \mathbb{P}(T_n = k - j) \\ &\leq \max_{j \geq 1} |j\lambda_j - \theta| \mathbb{P}(k - j_0 \leq T_n < k - 1) + \epsilon. \end{aligned}$$

Hence

$$\limsup_{n \rightarrow \infty} |r_n(k)| \leq \max_{j \geq 1} |j\lambda_j - \theta| \limsup_{n \rightarrow \infty} \sup_{x > 0} \mathbb{P}(x - j_0/n \leq n^{-1}T_n < x) + \epsilon = \epsilon,$$

because T_n/n converges in distribution to T , which has continuous distribution function.

5.2. Multisets

For combinatorial multisets, the Z_i have negative binomial distributions $\text{NB}(m_i, x^i)$, with

$$\mathbb{P}(Z_i = k) = \binom{m_i + k - 1}{k} (1 - x^i)^{m_i} x^{ik}, \quad k = 0, 1, \dots,$$

for any $x \in (0, 1)$; once again, the integers m_i are prescribed in the structure, and the joint distribution of the component sizes is the same for any choice of x . We have

$$\mathbb{E}Z_i = \frac{m_i x^i}{1 - x^i}, \quad \text{Var}Z_i = \frac{m_i x^i}{(1 - x^i)^2};$$

in the logarithmic class are those structures for which

$$m_i \sim \frac{\theta y^i}{i}, \quad \text{for some } y > 1, \theta \in (0, \infty),$$

when we take $x = y^{-1}$, and record that then

$$\lim_{i \rightarrow \infty} i \mathbb{E}Z_i = \lim_{i \rightarrow \infty} i m_i x^i = \theta. \quad (5.2)$$

Furthermore, the tail condition (1.3) is also satisfied, since

$$\mathbb{P}(Z_i \geq 2) = 1 - (1 - x^i)^{m_i} [1 + m_i x^i] \leq \frac{m_i(m_i + 1)x^{2i}}{2(1 - x^i)^2} = O(i^{-2}).$$

We note once more that, in the proofs below, the m_i need not be integers: see the examples in Section 5.4.

The recursion analogous to (5.1) for the distribution of T_n is (see [5])

$$k \mathbb{P}(T_n = k) = \sum_{j=1}^k g_n(j) \mathbb{P}(T_n = k - j), \quad (5.3)$$

where

$$g_n(j) = x^j \sum_{l=1; l|j}^n l m_l. \quad (5.4)$$

This is already substantially more complicated than (5.1). However, we note that, for $j \leq n$,

$$g_n(j) = g(j) \equiv x^j \sum_{l=1; l|j}^j l m_l,$$

and that

$$\lim_{i \rightarrow \infty} g(i) = \theta. \quad (5.5)$$

On the other hand, for $j > n$ we have

$$\begin{aligned} g_n(j) &= x^j \sum_{l=1; l|j}^n l m_l \leq x^j \sum_{l=1}^n l m_l \\ &= x^{j-n} \sum_{l=1}^n (l m_l x^l) x^{n-l} \leq \tilde{\theta} x^{j-n} \sum_{l=0}^{n-1} x^l \leq \frac{\tilde{\theta} x^{j-n}}{1-x}, \end{aligned} \quad (5.6)$$

where

$$\tilde{\theta} = \sup_{j \geq 1} j m_j x^j < \infty$$

under assumption (5.2).

Applying Theorem 4.4 when $k > n$ and using (5.6) shows that

$$\begin{aligned} \sum_{i=n+1}^k g_n(i) \mathbb{P}(T_n = k - i) &\leq \sum_{i=n+1}^k \frac{\tilde{\theta} x^{i-n}}{1-x} \mathbb{P}(T_n = k - i) \\ &= \frac{\tilde{\theta}}{1-x} \max_{l \geq 0} \mathbb{P}(T_n = l) \sum_{l=0}^{k-n-1} x^{k-n-l} \\ &\leq \max_{l \geq 0} \mathbb{P}(T_n = l) \frac{\tilde{\theta} x}{(1-x)^2} \\ &= o(1), \end{aligned} \quad (5.7)$$

uniformly in $k > n$.

This can be exploited to verify (LLA) as follows. The bound (5.7) shows that, for $k = 0, 1, \dots$,

$$k\mathbb{P}(T_n = k) = \sum_{i=1}^n g(i)\mathbb{P}(T_n = k - i) + o(1), \tag{5.8}$$

uniformly in $k \geq 0$. The method of proof in the previous section, together with (5.5), then shows that

$$k\mathbb{P}(T_n = k) = \theta\mathbb{P}(k - n \leq T_n < k) + r_n(k),$$

where $r_n(k) \rightarrow 0$ as $n \rightarrow \infty$, uniformly in k . The result follows from Theorem 4.3.

5.3. Selections

The next case we consider is the case of combinatorial selections, for which the Z_j are binomially distributed with

$$\mathbb{P}(Z_i = k) = \binom{m_i}{k} \left(\frac{x^i}{1+x^i}\right)^k \left(\frac{1}{1+x^i}\right)^{m_i-k}, \quad k = 0, 1, \dots, m_i,$$

for any $0 < x < 1$. Once more, the assumption that

$$m_i \sim \frac{\theta y^i}{i}$$

is necessary. In this case, we take $x = y^{-1} \in (0, 1)$, and (LC) and the tail condition hold.

To verify (LLA), the method of the previous section can be used, but this time based on the recurrence (see [5]) in (5.3), where

$$g_n(j) = x^j \sum_{l=1; l|j}^n (-1)^{j/l-1} l m_l. \tag{5.9}$$

The steps that lead to (5.6) and (5.7) follow immediately, with appropriate modification for the alternating nature of the $g_n(j)$.

5.4. Biased combinatorial structures

The preceding results are applicable to combinatorial structures that are not chosen uniformly from the set of possible structures of weight n , but rather with probability proportional to $\kappa^{\#\text{components}}$, for some $\kappa > 0$. According to the results in Section 8 of [5], these models also satisfy the identity (2.2); for assemblies the Z_j are Poisson with mean $\kappa m_j x^j / j!$, for multisets they are negative binomial with parameters m_j and κx^j , and for selections, binomial with parameters m_j and $\kappa x^j / (1 + \kappa x^j)$. It follows that if the uniform structure satisfies the conditions in (LC), then so do the biased structures, with θ replaced by $\kappa\theta$. Theorem 4.5 then follows from the earlier results from this section. For further examples of biasing, see [5].

5.5. General combinatorial structures

Now suppose the Z_i are arbitrary \mathbb{Z}_+ -valued random variables, with means $\mathbb{E}Z_i$ satisfying (LC). In combinatorial settings we are aware of (for example, [10]), Z_j can be decomposed

into the sum of m_j i.i.d. random variables Y_{j1}, \dots, Y_{jm_j} , each with p.g.f. $\phi_j(s)$ and means

$$\mathbb{E}Y_{j1} = y_j,$$

with the y_j eventually decreasing, and such that

$$j\mathbb{E}Z_j = jm_j y_j \rightarrow \theta \in (0, \infty).$$

Theorems 4.3 and 4.4 continue to hold, and

$$\begin{aligned} \mathbb{P}(Z_j \geq 2) &\leq \mathbb{E} \left\{ \sum_{i=1}^{m_j} \mathbb{1}(Y_{ji} \geq 2) + \sum_{1 \leq i < l \leq m_j} \mathbb{1}(Y_{ji} \geq 1) \mathbb{1}(Y_{jl} \geq 1) \right\} \\ &\leq m_j \mathbb{P}(Y_{j1} \geq 2) + \frac{1}{2} m_j^2 y_j^2 \end{aligned}$$

is of order $O(j^{-2})$ under (LC) if also

$$m_j \mathbb{P}(Y_{j1} \geq 2) = O(j^{-2});$$

this is typically the case. However, to get further we need a recursion for the point probabilities $\mathbb{P}(T_n = k)$. Since

$$\mathbb{E}s^{T_n} = \prod_{j=1}^n (\phi_j(s^j))^{m_j},$$

logarithmic differentiation leads to

$$k\mathbb{P}(T_n = k) = \sum_{l=1}^k g_n(l)\mathbb{P}(T_n = k - l),$$

where

$$g_n(l) = \sum_{j=1}^n jm_j [s^{l-j}] \frac{\phi'_j(s^j)}{\phi_j(s^j)}, \tag{5.10}$$

$[x^l]f(x)$ denoting the coefficient of x^l in $f(x)$. The following example shows that this recursion need not be easy to use.

Example 5.1. Suppose that a combinatorial structure is conditioned to have at most one component of each size [18]. If the original structure $\tilde{C}^{(n)}$ satisfies a conditioning relation like (2.2), then

$$\begin{aligned} \mathbb{P}(\tilde{C}^{(n)} = a | \tilde{C}^{(n)} \leq \mathbf{1}) &= \mathbb{P}(\tilde{Z}[1, n] = a | \tilde{Z}[1, n] \leq \mathbf{1}, \tilde{T}_n = n) \\ &= \mathbb{P}(Z[1, n] = a | T_n = n), \end{aligned}$$

where $\mathbf{1} = (1, \dots, 1)$ and $Z = (Z_1, \dots)$ is a vector of independent Bernoulli random variables satisfying

$$\mathbb{P}(Z_j = a) = \mathbb{P}(\tilde{Z}_j = a | \tilde{Z}_j \leq 1), \quad a = 0, 1. \tag{5.11}$$

If the original \tilde{Z}_j satisfy the (LC), then so too do the Z_j ; that is,

$$\pi_j \equiv \mathbb{P}(Z_j = 1) \sim \frac{\theta}{j}, \quad \theta \in (0, \infty), \tag{5.12}$$

and then $j\mathbb{E}Z_j = j\mathbb{P}(Z_j = 1) \rightarrow \theta$ automatically, and the tail condition is trivially satisfied. The point probabilities $\mathbb{P}(T_n = k)$ satisfy an equation of the form (5.3). Adapting (5.10) to the present setting with $m_j = 1$, $y_j = \pi_j$ and $\phi_j(s) = 1 - \pi_j + \pi_j s$ leads, after some simplification, to the fact that

$$g_n(l) = - \sum_{j=1, j|l}^n (-1)^{l/j} j h_j^{l/j}, \quad h_j \equiv \frac{\pi_j}{1 - \pi_j}. \quad (5.13)$$

It seems difficult to make progress with this approach in general, although in special cases verification of (LLA) should be possible. In [4], we have developed an alternative approach based on Stein’s method for compound Poisson approximation (*cf.* [6]). This leads to recursions for point probabilities that are easier to handle. If, for example, the Z_i satisfy (LC), together with the mild additional condition that

$$\sum_{i \geq 1} i \mathbb{E}(Z_i \mathbb{1}[Z_i \geq r]) < \infty, \quad \text{for some } r, \quad (5.14)$$

then (LLA) follows. Note that (5.14) clearly holds with $r = 2$ for Example 5.1. The alternative approach also provides bounds on the accuracy of the approximations in this paper.

Acknowledgement

We thank Dr Jennie Hansen for helpful comments on an earlier draft.

References

- [1] Aldous, D. J. (1985) *Exchangeability and Related Topics*, Vol. 1117 of *Lecture Notes in Mathematics*, Springer, New York, pp. 1–198.
- [2] Arratia, R., Barbour, A. D. and Tavaré, S. (1992) Poisson process approximations for the Ewens Sampling Formula. *Ann. Appl. Probab.* **2** 519–535.
- [3] Arratia, R., Barbour, A. D. and Tavaré, S. (1993) On random polynomials over a finite field. *Math. Proc. Camb. Phil. Soc.* **114** 347–368.
- [4] Arratia, R., Barbour, A. D. and Tavaré, S. (1998) *Logarithmic Combinatorial Structures*. In preparation.
- [5] Arratia, R. and Tavaré, S. (1994) Independent process approximations for random combinatorial structures. *Adv. Math.* **104** 90–154.
- [6] Barbour, A. D., Chen, L. H. Y. and Loh, W.-L. (1992) Compound Poisson approximation for nonnegative random variables via Stein’s method. *Ann. Probab.* **20** 1843–1866.
- [7] Ewens, W. J. (1972) The sampling theory of selectively neutral alleles. *Theoret. Pop. Biol.* **3** 87–112.
- [8] Goncharov, V. L. (1944) Some facts from combinatorics. *Izvestia Akad. Nauk. SSSR, Ser. Mat.* **8** 3–48. See also: On the field of combinatory analysis. *Trans. Amer. Math. Soc.* **19** 1–46.
- [9] Hansen, J. C. (1994) Order statistics for decomposable combinatorial structures. *Random Struct. Alg.* **5** 517–533.
- [10] Hansen, J. C. and Schmutz, E. (1993) How random is the characteristic polynomial of a random matrix? *Math. Proc. Camb. Phil. Soc.* **114** 507–515.
- [11] Kingman, J. F. C. (1975) Random discrete distributions. *J. Royal Statist. Soc.* **37** 1–22.

- [12] Kingman, J. F. C. (1977) The population structure associated with the Ewens sampling formula. *Theoret. Pop. Biol.* **11** 274–283.
- [13] Kingman, J. F. C. (1993) *Poisson Processes*, Oxford University Press, Oxford.
- [14] Kolchin, V. F. (1976) A problem of the allocation of particles in cells and random mappings. *Theor. Probab. Appl.* **21** 48–63.
- [15] Kolchin, V. F. (1986) *Random Mappings*, Optimization Software, Inc., New York.
- [16] Scheffé, H. (1947) A useful convergence theorem for probability distributions. *Ann. Math. Stat.* **18** 434–438.
- [17] Shepp, L. A. and Lloyd, S. P. (1966) Ordered cycle lengths in a random permutation. *Trans. Amer. Math. Soc.* **121** 340–357.
- [18] Stark, D. (1994) Total variation distance for independent process approximations of random combinatorial objects. PhD thesis, University of Southern California.
- [19] Vershik, A. M. and Schmidt, A. A. (1977) Limit measures arising in the theory of groups I. *Theor. Probab. Appl.* **22** 79–85.
- [20] Vervaat, W. (1972) *Success Epochs in Bernoulli Trials with Applications in Number Theory*. Vol. 42 of *Mathematical Center Tracts*, Mathematisch Centrum, Amsterdam.
- [21] Watterson, G. A. (1974) The sampling theory of selectively neutral alleles. *Adv. Appl. Probab.* **6** 463–488.
- [22] Watterson, G. A. (1976) The stationary distribution of the infinitely-many-alleles diffusion model. *J. Appl. Probab.* **13** 639–651.