# Asymptotic Normality Through Factorial Cumulants and Partition Identities

**Konstancja Bobecka**[1,†], **Paweł Hitczenko**[2,‡], **Fernando López-Blázquez**[3,†], **Grzegorz Rempała**[4,§], and **Jacek Wesołowski**[1,¶]

[1]Wydział Matematyki i Nauk Informacyjnych, Politechnika Warszawska, Warszawa, Poland, (bobecka@mini.pw.edu.pl, wesolo@mini.pw.edu.pl)

[2]Department of Mathematics, Drexel University, Philadelphia, USA, (phitczenko@math.drexel.edu)

[3]Facultad de Matemáticas Universidad de Sevilla, Sevilla, Spain, (lopez@us.es)

[4]Department of Biostatistics, Georgia Health University, Augusta, USA, (grempala@georgiahealth.edu)

## Abstract

In the paper we develop an approach to asymptotic normality through factorial cumulants. Factorial cumulants arise in the same manner from factorial moments as do (ordinary) cumulants from (ordinary) moments. Another tool we exploit is a new identity for 'moments' of partitions of numbers. The general limiting result is then used to (re-)derive asymptotic normality for several models including classical discrete distributions, occupancy problems in some generalized allocation schemes and two models related to negative multinomial distribution.

## 1. Introduction

Convergence to the normal law is one of the most important phenomena of probability. As a consequence, a number of general methods, often based on transforms of the underlying sequence, have been developed as techniques for establishing such convergence. One of these methods, called the method of moments, rests on the fact that the standard normal random variable is uniquely determined by its moments, and that for such a random variable $X$, if $(X_n)$ is a sequence of random variables having all moments and $\mathbb{E} X_n^k \to \mathbb{E} X^k$ for all $k = 1, 2, \ldots$, then $X_n \overset{d}{\to} X$: see, *e.g.*, [22, Theorem 2.22] or [2, Theorem 30.2]. Here and throughout the paper we use '$\overset{d}{\to}$' to denote the convergence in distribution.

Since moments are not always convenient to work with, one can often use some other characteristics of random variables to establish the convergence to the normal law. For example, in one classical situation we consider a sequence of cumulants (we recall the definition in the next section) rather than moments. On the one hand, since the $k$th cumulant is a continuous function of the first $k$ moments, to prove that $X_n \overset{d}{\to} X$ instead of

convergence of moments one can use convergence of cumulants. On the other hand, all cumulants of the standard normal distribution are zero except for the cumulant of order 2, which equals 1. This often makes it much easier to establish the convergence of cumulants of $(X_n)$ to the cumulants of the standard normal random variable. We refer the reader to [8, Section 6.1], for example, for a more detailed discussion.

In this paper we develop an approach to asymptotic normality that is based on factorial cumulants. They will be discussed in the next section. Here we just indicate that factorial cumulants arise in the same manner from factorial moments as do (ordinary) cumulants from (ordinary) moments. The motivation for our work is the fact that quite often one encounters situations in which properties of random variables are naturally expressed through factorial (rather than ordinary) moments. As we will see below, this is the case, for instance, when random variables under consideration are sums of indicator variables.

In developing our approach we first provide a simple and yet quite general sufficient condition for the Central Limit Theorem (CLT) in terms of factorial cumulants (see Proposition 2.1 below). Further, as we will see in Theorem 4.1, we show that the validity of this condition can be verified by controlling the asymptotic behaviour of factorial moments. This limiting result will be then used in Section 5 to (re-)derive asymptotic normality for several models including classical discrete distributions, occupancy problems in some generalized allocation schemes (GAS), and two models related to a negative multinomial distribution; they are examples of what may be called generalized inverse allocation schemes (GIAS). Generalized allocation schemes were introduced in [13], and we refer the reader to chapters in books [14, 15] by the same author for more details, properties, and further references. The term 'generalized inverse allocation schemes' does not seem to be commonly used in the literature; in our terminology the word 'inverse' refers to inverse sampling, a sampling method proposed in [6]. A number of distributions with 'inverse' in their names reflecting the inverse sampling are discussed in the first ([10]) and the fourth ([9]) volumes of the Wiley Series in Probability and Statistics.

We believe that our approach may turn out to be useful in other situations when the factorial moments are natural and convenient quantities to work with. We wish to mention, however, that although several of our summation schemes are, in fact, a GAS or a GIAS, we do not have a general statement that would give reasonably general sufficient conditions under which a GAS or a GIAS exhibits asymptotic normality. It may be an interesting question, worth further investigation.

Aside from utilizing factorial cumulants, another technical tool we exploit is an identity for 'moments' of partitions of natural numbers (see Proposition 3.1). As far as we know this identity is new, and may be of independent interest to the combinatorics community. As of now, however, we do not have any combinatorial interpretation, either for its validity or for its proof.

## 2. Factorial cumulants

Let *X* be a random variable with the Laplace transform

$$\phi_X(t) = \mathbb{E}e^{tX} = \sum_{k=0}^{\infty} \mu_k \frac{t^k}{k!}$$

and the cumulant transform

$$\psi_X(t) = \log(\phi_X(t)) = \sum_{k=0}^{\infty} \kappa_k \frac{t^k}{k!}.$$

Then $\mu_k = \mathbb{E}X^k$ and $\kappa_k$ are, respectively, the $k$th moment and the $k$th cumulant of $X$, $k = 0, 1, \ldots$.

One can view the sequence $\underline{\kappa} = (\kappa_k)$ as obtained by a transformation $f = (f_k)$ of the sequence $\underline{\mu} = (\mu_k)$, that is, $\underline{\kappa} = f(\underline{\mu})$, where the $f_k$ are defined recursively by $f_1(\underline{x}) = x_1$ and

$$f_k(\mathrm{x}) = \mathrm{x_k} - \sum_{\mathrm{j}=1}^{k-1} \binom{k-1}{j-1} \mathrm{f_j(x) x_{k-j}}, \mathrm{k} > 1. \quad (2.1)$$

The Laplace transform can also be expanded in the form

$$\phi_X(t) = \sum_{k=0}^{\infty} \nu_k \frac{(e^t - 1)^k}{k!}, \quad (2.2)$$

where $\nu_k$ is the $k$th factorial moment of $X$, that is,

$$\nu_0 = 1 \text{ and } \nu_k = \mathbb{E}(X)_k = \mathbb{E}X(X-1) \cdots (X-k+1), k = 1, \ldots.$$

Here and below we use the Pochhammer symbol $(x)_k$ for the falling factorial $x(x-1) \cdots (x-(k-1))$.

In analogy to (2.1), one can view the sequence $\underline{\mu}$ as obtained by a transformation $g = (g_k)$ of the sequence $\underline{\nu} = (\nu_k)$, that is, $\underline{\mu} = g(\nu)$, where the $g_k$ are defined by

$$g_k(\mathrm{x}) = \sum_{\mathrm{j}=1}^{k} \mathrm{S_2(k, j) x_j}, \mathrm{k} \geq 1, \quad (2.3)$$

where $S_2(k, j)$ are the Stirling numbers of the second kind defined by the identity

$$y^k = \sum_{j=1}^{k} S_2(k, j)(y)_j,$$

holding for any $y \in \mathbb{R}$ (see, *e.g.*, [5, (6.10)]).

Following this probabilistic terminology for any sequence of real numbers $\underline{a} = (a_k)$, one can define its cumulant and factorial sequences, $\underline{b} = (b_k)$ and $\underline{c} = (c_k)$, respectively, by

$$\sum_{k=0}^{\infty} a_k \frac{t^k}{k!} = \exp\left\{\sum_{k=0}^{\infty} b_k \frac{t^k}{k!}\right\} = \sum_{k=0}^{\infty} c_k \frac{(e^t - 1)^k}{k!}. \quad (2.4)$$

The first relation is known in combinatorics by the name 'exponential formula', and its combinatorial interpretation when both $(a_n)$ and $(b_n)$ are non-negative integers may be found, for example, in [20, Section 5.1].

Note that if $\underline{a}$, $\underline{b}$ and $\underline{c}$ are related by (2.4), then $\underline{b} = f(\underline{a})$ and $\underline{a} = g(\underline{c})$, where $f$ and $g$ are given by (2.1) and (2.3), respectively.

Let the sequence $\underline{d} = (d_k)$ be defined by

$$\exp\left\{\sum_{k=0}^{\infty} d_k \frac{(e^t - 1)^k}{k!}\right\} = \sum_{k=0}^{\infty} c_k \frac{(e^t - 1)^k}{k!}.$$

Then, regarding $e^t - 1$ as a new variable, we see that $\underline{d} = f(\underline{c})$. Since $\underline{c}$ is a factorial sequence for $\underline{a}$ and $\underline{d}$ is a cumulant sequence for $\underline{c}$, we call $\underline{d} = f(\underline{c})$ the *factorial cumulant sequence* for $\underline{a}$.

Observe that, by (2.4),

$$\sum_{k=0}^{\infty} d_k \frac{(e^t - 1)^k}{k!} = \sum_{k=0}^{\infty} b_k \frac{t^k}{k!},$$

and thus $\underline{b} = g(\underline{d}) = g(f(\underline{c}))$. That is,

$$b_k = \sum_{j=1}^{k} S_2(k, j) f_j(c), \mathrm{k} = 1, 2, \dots. \quad (2.5)$$

This observation is useful for proving convergence in law to the standard normal variable.

**Proposition 2.1.** *Let $(S_n)$ be a sequence of random variables having all moments. Assume that*

$$\mathrm{Var} S_n \xrightarrow{n \to \infty} \infty \quad (2.6)$$

*and*

$$\frac{\mathbb{E} S_n}{\mathrm{Var}^{\frac{3}{2}} S_n} \xrightarrow{n \to \infty} 0. \quad (2.7)$$

*For any $n = 1, 2, \dots,$ let $\underline{c}_n = (c_{k,n})_{k=1, \dots}$ be the sequence of factorial moments of $S_n$, that is, $c_{k,n} = \mathbb{E}(S_n)_k$, $k = 1, 2, \dots,$ and let $f_{J,n} = f_J(\underline{c}_n)$ (where $f_J$ is defined by (2.1)) be the Jth factorial cumulant of $S_n$, $J = 1, 2, \dots.$ Assume that*

$$\frac{f_{J,n}}{\mathrm{Var}^{\frac{J}{2}} S_n} \xrightarrow{n \to \infty} 0, \textit{for } J \geq 3. \quad (2.8)$$

*Then*

$$U_n = \frac{S_n - \mathbb{E} S_n}{\sqrt{\mathrm{Var} S_n}} \xrightarrow{d} \mathcal{N}(0, 1). \quad (2.9)$$

**Proof.** We will use the cumulant convergence theorem (see, *e.g.*, [8, Theorem 6.14]). Let $\kappa_{J,n}$ denote the Jth cumulant of $U_n$ and recall that all cumulants of the standard normal distribution are zero except for the cumulant of order 2, which is 1. It is obvious that $\kappa_{1,n} = 0$ and $\kappa_{2,n} = 1$. Therefore, to prove (2.9) it suffices to show that $\kappa_{J,n} \to 0$ for $J \quad 3$. By (2.5),

$$\kappa_{J,n} = \frac{\sum_{j=1}^{J} S_2(J,j) f_{j,n}}{\mathrm{Var}^{\frac{J}{2}} S_n}.$$

Fix arbitrary $J \geq 3$. To prove that $\kappa_{J,n} \to 0$, it suffices to show that

$$\frac{f_{j,n}}{\mathrm{Var}^{\frac{J}{2}} S_n} \to 0, \text{for} \quad \text{all} j = 1, 2, \ldots, J. \quad (2.10)$$

Note first that by (2.1)

$$f_{1,n} = \mathbb{E} S_n \text{and} f_{2,n} = \mathrm{Var}(S_n) - \mathbb{E} S_n.$$

Therefore the assumptions (2.6) and (2.7) imply (2.10) for $j = 1, 2$.

If $j \in \{3, \ldots, J-1\}$, we write

$$\frac{f_{j,n}}{\mathrm{Var}^{\frac{J}{2}} S_n} = \frac{f_{j,n}}{\mathrm{Var}^{\frac{j}{2}} S_n} \frac{1}{\mathrm{Var}^{\frac{J-j}{2}} S_n}.$$

By (2.8) the first factor tends to zero and by (2.6) the second factor tends to zero as well.

Finally, for $j = J$ the conditions (2.10) and (2.8) are identical.

The above result is particularly useful when the factorial moments of $S_n$ are available in a nice form. We will now describe a general situation when this happens.

For any numbers $\delta_j$, $j = 1, \ldots, N$, assuming values 0 or 1 we have

$$x^{\sum_{j=1}^{N} \delta_j} = \sum_{m=0}^{\sum_{j=1}^{N} \delta_j} \binom{\sum_{j=1}^{N} \delta_j}{m} (x-1)^m = 1 + \sum_{m=1}^{N} (x-1)^m \sum_{1 \leq j_1 < \cdots < j_m \leq N} \delta_{j_1} \cdots \delta_{j_m}.$$

Therefore, if $(\varepsilon_1, \ldots, \varepsilon_N)$ is a random vector valued in $\{0, 1\}^N$ and $S = \sum_{i=1}^{N} \varepsilon_i$, then

$$\mathbb{E} e^{tS} = 1 + \sum_{m=1}^{\infty} (e^t - 1)^m \sum_{1 \leq j_1 < \cdots < j_m \leq N} \mathbb{P}(\varepsilon_{j_1} = \varepsilon_{j_2} = \cdots = \varepsilon_{j_m} = 1).$$

Comparing this formula with (2.2), we conclude that factorial moments of $S$ have the form

$$\mathbb{E}(S)_k = k! \sum_{1 \leq j_1 < \cdots < j_k \leq N} \mathbb{P}(\varepsilon_{j_1} = \varepsilon_{j_2} = \cdots = \varepsilon_{j_k} = 1) =: c_k, k = 1, 2, \ldots. \quad (2.11)$$

If, in addition, the random variables $(\varepsilon_1, \ldots \varepsilon_N)$ are exchangeable, then the above formula simplifies to

$$\mathbb{E}(S)_k = (N)_k \mathbb{P}(\varepsilon_1 = \cdots = \varepsilon_k = 1) =: c_k, k = 1, 2, \ldots. \quad (2.12)$$

As we will see in Section 5, our sufficient condition for asymptotic normality will work well for several set-ups falling within such a scheme. This will be preceded by a derivation of new identities for integer partitions, which will give us a major enhancement of the tools we will use to prove limit theorems.

## 3. Partition identities

Recall that if $\underline{b} = (b_n)$ is a cumulant sequence for a sequence of numbers $\underline{a} = (a_n)$, that is, $\underline{b} = f(\underline{a})$ with $f$ given by (2.1), then for $J \geq 1$

$$b_J = \sum_{\pi \subset J} D_\pi \prod_{i=1}^{J} a_i^{m_i}, \text{ where } D_\pi = \frac{(-1)^{\sum_{i=1}^{J} m_i - 1} J!}{\prod_{i=1}^{J} (i!)^{m_i} \sum_{i=1}^{J} m_i} \binom{\sum_{i=1}^{J} m_i}{m_1, \ldots, m_J}, \quad (3.1)$$

and where the sum is over all partitions $\pi$ of a positive integer $J$, *i.e.*, over all vectors $\pi = (m_1, \ldots, m_J)$ with non-negative integer components which satisfy $\sum_{i=1}^{J} i m_i = J$ (for a proof, see, *e.g.*, [12, Section 3.14]).

Note that for $J \geq 2$

$$\sum_{\pi \subset J} D_\pi = 0. \quad (3.2)$$

This follows from the fact that all the cumulants, except for the first one, of the constant random variable $X = 1$ a.s. are zero.

For $\pi = (m_1, \ldots m_J)$ we denote $H_\pi(s) = \sum_{i=1}^{J} i^s m_i, s = 0, 1, 2, \cdots$. The main result of this section is the identity which considerably extends (3.2).

**Proposition 3.1.** *Assume $J \geq 2$. Let $I \geq 1$ and $s_i \geq 1, i = 1, \ldots, I$, be such that*

$$\sum_{i=1}^{I} s_i \leq J + I - 2.$$

*Then*

$$\sum_{\pi \subset J} D_\pi \prod_{i=1}^{I} H_\pi(s_i) = 0. \quad (3.3)$$

**Proof.** We use induction with respect to $J$. Note that if $J = 2$ then $I$ may be arbitrary, but $s_i = 1$ for all $i = 1, \ldots, I$. Thus the identity (3.3) follows from (3.2) since $H_\pi(1) = J$ for any $J$ and any $\pi \subset J$.

Now assume that the result holds true for $J = 2, \ldots, K - 1$, and consider the case of $J = K$. That is, we want to study

$$\sum_{\pi \subset K} D_\pi \prod_{i=1}^{I} H_\pi(s_i)$$

under the condition $\sum_{i=1}^{I} S_i \leq K+I-2$.

Let us introduce functions $g_i$, $i = 1, \ldots, K$, by letting

$$g_i(m_1, \ldots, m_K) = (\tilde{m}_1, \ldots, \tilde{m}_{K-1}) = \begin{cases} (m_1, \ldots, m_{i-1}+1, m_i-1, \ldots, m_{K-1}) & \text{if } i \neq 1, K, \\ (m_1-1, m_2, \ldots, m_{K-1}) & \text{if } i=1, \\ (m_1, \ldots, m_{K-2}, m_{K-1}+1) & \text{if } i=K. \end{cases}$$

Note that

$$g_i : \{\pi \subset K : m_i \geq 1\} \rightarrow \{\tilde{\pi} \subset (K-1) : \tilde{m}_{i-1} \geq 1\}, \quad i=1, \ldots, K,$$

are bijections. Here for consistency we assume $\tilde{m_0} = 1$.

Observe that for any $s$, any $\pi \subset K$ such that $m_i \geq 1$, and for $\tilde{\pi} = g_i(\pi) \subset (K-1)$ we have

$$H_{\tilde{\pi}}(s) = H_\pi(s) - i^s + (i-1)^s = H_\pi(s) - 1 - A_s(i-1),$$

where

$$A_s(i-1) = \sum_{k=1}^{s-1} \binom{s}{k} (i-1)^k$$

is a polynomial of degree $s - 1$ in the variable $i - 1$ with constant term equal to zero. In particular, for $i = 1$ we have $H_\pi(s) = H_\pi(s) - 1$. Therefore, expanding $H_\pi(s_1)$ we obtain

$$\sum_{\pi \subset K} D_\pi \prod_{i=1}^{I} H_\pi(s_i) = \sum_{i=1}^{K} i^{s_1} \sum_{\substack{\pi \subset K \\ m_i \geq 1}} m_i D_\pi \prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j) + 1 + A_{s_j}(i-1)].$$

Note that if $m_i \geq 1$ then

$$\frac{1}{K!} i m_i D_\pi = \begin{cases} \dfrac{(-1)^{M_\pi-1} M_\pi!}{M_\pi(m_1-1)! \prod_{k=2}^{K-1} m_k!(k!)^{m_k}} = -\dfrac{1}{(K-1)!} \sum_{j=1}^{K-1} D_{\tilde{\pi}} \tilde{m}_j & \text{if } i=1, \\ \dfrac{(-1)^{M_\pi-1} M_\pi!}{M_\pi(m_i-1)!(i-1)!(i!)^{m_i-1} \prod_{\substack{k=2 \\ k \neq i}}^{K} m_k!(k!)^{m_k}} = \dfrac{1}{(K-1)!} D_{\tilde{\pi}} \tilde{m}_{i-1} & \text{if } i=2, \ldots, K, \end{cases}$$

where $\tilde{\pi} = (\tilde{m_1}, \ldots, \tilde{m_{K-1}}) = g_i(\pi)$, respectively. Therefore,

$$\frac{1}{K} \sum_{\pi \subset K} D_\pi \prod_{i=1}^{I} H_\pi(s_i) = -\sum_{i=1}^{K-1} \sum_{\tilde{\pi} \subset (K-1)} D_{\tilde{\pi}} \tilde{m}_i \prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j) + 1] + \sum_{i=2}^{K} \sum_{\substack{\tilde{\pi} \subset (K-1) \\ \tilde{m}_{i-1} \geq 1}} D_{\tilde{\pi}} \tilde{m}_{i-1} i^{s_1-1} \prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j) + 1 + A_{s_j}(i-1)]. \quad (3.4)$$

The second term in the above expression can be written as

$$\sum_{i=1}^{K-1} \sum_{\tilde{\pi} \subset (K-1)} D_{\tilde{\pi}} \tilde{m}_i (i+1)^{s_1-1} \prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j)+1+A_{s_j}(i)]. \quad (3.5)$$

Note that

$$(i+1)^{s_1-1} \prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j)+1+A_{s_j}(i)] = \sum_{r=0}^{s_1-1} \binom{s_1-1}{r} \sum_{M=0}^{I-1} \sum_{2 \le u_1 < \cdots < u_M \le I} i^r \prod_{h=1}^{M} (H_{\tilde{\pi}}(s_{u_h})+1) \prod_{\substack{2 \le j \le I \\ j \notin \{u_1,\ldots,u_m\}}} A_{s_j}(i).$$

The term with $r = 0$ and $M = I - 1$ in the above expression is $\prod_{j=2}^{I} [H_{\tilde{\pi}}(s_j)+1]$, so this term of the sum (3.5) cancels with the first term of (3.4).

Hence, we only need to show that for any $r \in \{1, \ldots, s_1 - 1\}$, any $M \in \{0, \ldots, I - 1\}$, and any $2 \le u_1 < \cdots < u_M \le I$,

$$\sum_{\tilde{\pi} \subset (K-1)} D_{\tilde{\pi}} \prod_{h=1}^{M} (H_{\tilde{\pi}}(s_{u_h})+1) \sum_{i=1}^{K-1} \tilde{m}_i i^r \prod_{\substack{2 \le j \le I \\ j \notin \{u_1,\ldots,u_M\}}} A_{s_j}(i) = 0. \quad (3.6)$$

Observe that the expression

$$\sum_{i=1}^{K-1} \tilde{m}_i i^r \prod_{\substack{2 \le j \le I \\ j \notin \{u_1,\ldots,u_M\}}} A_{s_j}(i)$$

is a linear combination of $H_{\tilde{\pi}}$ functions with the largest value of an argument equal to

$$\sum_{\substack{2 \le j \le I \\ j \notin \{u_1,\ldots,u_M\}}} (s_j - 1) + r.$$

Therefore the left-hand side of (3.6) is a respective linear combination of terms of the form

$$\sum_{\tilde{\pi} \subset (K-1)} D_{\tilde{\pi}} \prod_{w=1}^{W} H_{\tilde{\pi}}(t_w), \quad (3.7)$$

where

$$\sum_{w=1}^{W} t_w \le \sum_{j=1}^{M} s_{u_j} - (M-W+1) + \sum_{\substack{2 \le j \le I \\ j \notin \{u_1,\ldots,u_M\}}} (s_j-1) + s_1 - 1 \le \sum_{i=1}^{I} s_i - (M-W+1) - (I-1-M) - 1.$$

But we assumed that $\sum_{i=1}^{I} s_i \le K+I - 2$. Therefore

$$\sum_{w=1}^{W} t_w \le K+I-2-(M-W+1)-(I-1-M)-1=(K-1)+W-2.$$

Now, by the inductive assumption it follows that any term of the form (3.7) is zero, thus proving (3.6).

Note that in a similar way one can prove that (3.3) is no longer true when

$$\sum_{i=1}^{I} s_i = J+I-1.$$

**Remark 3.2.** Richard Stanley [21] provided the following combinatorial description of the left-hand side of (3.3). Put

$$F_n(x)=\sum_k S_2(n,k)x^k,$$

and let $(s_i)$ be a sequence of positive integers. Then the left-hand side of (3.3) is a coefficient of $x^J/J!$ in

$$\sum_{\mathscr{P}}(-1)^{|\mathscr{P}|-1}(|\mathscr{P}|-1)!\prod_{B\in\mathscr{P}}F_{\sigma B}(x),$$

where the sum ranges over all partitions $\mathscr{P}$ of a set $\{1, \ldots, I\}$ into $|\mathscr{P}|$ non-empty pairwise disjoint subsets, and where for any such subset $B \in \mathscr{P}$

$$\sigma B=\sum_{i\in B}s_i.$$

In the simplest case when $I = 1$ for any positive integer $s_1$, the left-hand side of equation (3.3) is equal to $J!S_2(s_1, J)$. Since this is the number of surjective maps from an $s_1$-element set to a $J$-element set, it must be 0 for $s_1 < J$, which is exactly what Proposition 3.1 asserts. It is not clear to us how easy it would be to show that (3.3) holds for the larger values of $I$.

## 4. Central Limit Theorem: general set-up

To illustrate and motivate how our approach is intended to work, consider a sequence $(S_n)$ of Poisson random variables, where $Sn \sim \text{Poisson}(\lambda_n)$. As is well known, if $\lambda_n \to \infty$ then $U_n=(S_n-\mathbb{E}S_n)/\sqrt{\text{var}S_n}$ converges in distribution to $\mathscr{N}0, 1)$. To see how it follows from our approach, note that $\mathbb{E}S_n = \text{Var } S_n = \lambda_n$, and therefore the assumptions (2.6) and (2.7) of Proposition 2.1 are trivially satisfied. Moreover, the factorial moments of $S_n$ are $c_{i,n}=\lambda_n^i$. Consequently, $\prod_{i=1}^{J} c_{i,n}^{m_i}=\lambda_n^J$ for any partition $\pi = (m_1, \ldots, m_J)$ of $J \in \mathbb{N}$, and thus

$$f_{J,n}=f_J(c_n)=\sum_{\pi\subset J}D_\pi\prod_{i=1}^{J}c_{i,n}^{m_i}=\lambda_n^J\sum_{\pi\subset J}D_\pi.$$

It now follows from the simplest case (3.2) of our partition identity that $f_{J,n} = 0$ as long as $J$ 2. Hence the assumption (2.8) of Proposition 2.1 is also trivially satisfied and we conclude the asymptotic normality of $(U_n)$.

The key feature in the above argument was, of course, the very simple form of the factorial moments $c_{i,n}$ of $S_n$, which resulted in factorization of the products of $c_{i,n}^{m_i}$ in the expression for $f_{J,n}$. It is not unreasonable, however, to expect that if the expression for moments does not depart too much from the form it took for the Poisson variable, then with the full strength of Proposition 3.1 one might be able to prove the CLT. This is the essence of condition (4.1) in the next result. This condition, when combined with the extension of (3.2) given in Proposition 3.1, allows us to greatly refine Proposition 2.1 towards a possible use in schemes of summation of indicators, as was suggested in the final part of Section 2.

**Theorem 4.1.** *Let $(S_n)$ be a sequence of random variables with factorial moments $c_{i,n}$, $i$, $n = 1, 2, \ldots$. Assume that (2.6) and (2.7) are satisfied and suppose that $c_{i,n}$ can be decomposed into the form*

$$c_{i,n} = L_n^i \exp\left(\sum_{j \geq 1} \frac{Q_{j+1}^{(n)}(i)}{jn^j}\right), i, n = 1, 2, \ldots, \quad (4.1)$$

*where $(L_n)$ is a sequence of real numbers and $Q_j^{(n)}$ is a polynomial of degree at most j such that*

$$|Q_j^{(n)}(i)| \leq (Ci)^j, \text{ for all } i \in \mathbb{N}, \quad (4.2)$$

*with $C > 0$ a constant not depending on n or j. Assume further that for all $J$ 3*

$$\frac{L_n^J}{n^{J-1}\text{Var}^{\frac{J}{2}}S_n} \xrightarrow{n \to \infty} 0. \quad (4.3)$$

*Then*

$$U_n = \frac{S_n - \mathbb{E}S_n}{\sqrt{\text{Var}S_n}} \xrightarrow{d} \mathcal{N}(0, 1), \text{ as } n \to \infty. \quad (4.4)$$

**Proof.** Due to Proposition 2.1 we only need to show that (2.8) holds. The representation (4.1) implies

$$f_{J,n} = \sum_{\pi \subset J} D_\pi \prod_{i=1}^{J} c_{i,n}^{m_i} = L_n^J \sum_{\pi \subset J} D_\pi e^{z_\pi(J,n)},$$

where

$$z_\pi(J,n) = \sum_{j \geq 1} \frac{A_\pi^{(n)}(j)}{jn^j} \text{ with } A_\pi^{(n)}(j) = \sum_{i=1}^{J} m_i Q_{j+1}^{(n)}(i).$$

Fix arbitrary $J$ 3. To prove (2.8), in view of (4.3) it suffices to show that $\sum_{\pi \subset J} D_\pi e^{z_\pi(J,n)}$ is of order $n^{-(J-1)}$. To do that, we expand $e^{z_\pi(J,n)}$ into power series to obtain

$$\sum_{\pi \subset J} D_\pi e^{z_\pi(J,n)}$$

$$=\sum_{\pi \subset J} D_\pi e^{\sum_{j \geq 1} \frac{1}{jn^j} A_\pi^{(n)}(j)}$$

$$=\sum_{\pi \subset J} D_\pi \sum_{s=0}^{\infty} \frac{1}{s!} \left( \sum_{j \geq 1} \frac{1}{jn^j} A_\pi^{(n)}(j) \right)^s$$

$$=\sum_{s \geq 1} \frac{1}{s!} \sum_{l \geq s} \frac{1}{n^l} \sum_{\substack{j_1,\ldots,j_s \geq 1 \\ \sum_{k=1}^s j_k = l}} \frac{1}{\prod_{k=1}^s j_k} \sum_{\pi \subset J} D_\pi \prod_{k=1}^s A_\pi^{(n)}(j_k).$$

We claim that whenever $\sum_{k=1}^s j_k \leq j - 2$ then

$$\sum_{\pi \subset J} D_\pi \prod_{k=1}^s A_\pi^{(n)}(j_k)=0. \quad (4.5)$$

To see this, note that by changing the order of summation in the expression for $A_\pi^{(n)}(j)$ we can write it as

$$A_\pi^{(n)}(j)=\sum_{k=0}^{j+1} \alpha_{k,j+1}^{(n)} H_\pi(k),$$

where $(\alpha_{k,j}^{(n)})$ are the coefficients of the polynomial $Q_j^{(n)}$, that is,

$$Q_j^{(n)}(x)=\sum_{k=0}^{j} \alpha_{k,j}^{(n)} x^k.$$

Consequently, (4.5) follows from identity (3.3).

To handle the terms for which $\sum_{k=1}^s j_k > J - 2$, note that

$$|A_\pi^{(n)}(j)| \leq (CJ)^{j+1} \sum_{i=1}^{J} m_i < K(CJ)^j,$$

where $K > 0$ is a constant depending only on $J$ (and not on the partition $\pi$). Hence,

$$\left| \sum_{\pi \subset J} D_\pi \prod_{k=1}^s A_\pi^{(n)}(j_k) \right| \leq \sum_{\pi \subset J} |D_\pi| \prod_{k=1}^s K(CJ)^{j_k} \leq \tilde{C} K^s (CJ)^{\sum_{k=1}^s j_k}, \quad (4.6)$$

where $\tilde{C} = \Sigma_{\pi \subset J} |D_\pi|$ is a constant depending only on $J$. Therefore, restricting the sum according to (4.5) and using (4.6), we get

$$\left|\sum_{\pi \subset J} D_\pi e^{z_\pi(J,n)}\right| \leq \sum_{s \geq 1} \frac{1}{s!} \sum_{1 \geq \max\{s, J-1\}} \frac{1}{n^l} \sum_{\substack{j_1,\ldots,j_s \geq 1 \\ \sum_{k=1}^s j_k = l}} \frac{1}{\prod_{k=1}^s j_k} \left|\sum_{\pi \subset J} D_\pi \prod_{k=1}^s A_\pi^{(n)}(j_k)\right| \leq \tilde{C} \sum_{s \geq 1} \frac{K^s}{s!} \sum_{l \geq J-1} \frac{1}{n^l} l^s (CJ)^l.$$

Here we used the inequality

$$\sum_{\substack{j_1,\ldots,j_s \geq 1 \\ \sum_{k=1}^s j_k = l}} \frac{1}{\prod_{k=1}^s j_k} < l^s,$$

which may be seen by trivially bounding the sum by the number of its terms. Now we change the order of summations, arriving at

$$\left|\sum_{\pi \subset J} D_\pi e^{z_\pi(J,n)}\right| \leq \tilde{C} \sum_{l \geq J-1} \left(\frac{CJ}{n}\right)^l \sum_{s \geq 1} \frac{(lK)^s}{s!} \leq \tilde{C} \sum_{l \geq J-1} \left(\frac{CJe^K}{n}\right)^l = \tilde{C} \left(\frac{CJe^K}{n}\right)^{J-1} \sum_{l \geq 0} \left(\frac{CJe^K}{n}\right)^l.$$

The result follows, since for $n$ sufficiently large (such that $CJe^K < n$), the series in the last expression converges.

**Remark 4.2.** A typical way Theorem 4.1 will be applied is as follows. Assume that $\mathbb{E}S_n$ and $\text{Var } S_n$ are of the same order $n$. Then obviously, (2.6) and (2.7) are satisfied. Assume also that (4.1) and (4.2) hold and that $L_n$ is also of order $n$. Then clearly (4.3) is satisfied and thus (4.4) holds true.

## 5. Applications

In this section we show how the tools developed in the previous section, and in particular the decomposition (4.1) together with the condition (4.3), can be conveniently used for proving CLTs in several situations, mostly in summation schemes of {0, 1}-valued random variables, as was indicated in Section 2. First, four more or less standard limit results for the binomial, negative binomial, hypergeometric and negative hypergeometric schemes will be re-proved. Then more involved schemes of allocation problems for distinguishable balls, indistinguishable balls, coloured balls, and random forests will be considered. The CLTs for the number of boxes with exactly $r$ balls in the first two problems and for the number of trees with exactly $r$ non-root vertices in the third problem will be derived. While the CLT in the case of distinguishable balls has been known in the literature for years (see, *e.g.*, [16]), the main focus in the other two cases appears to be on the local limit theorems (see, *e.g.*, [14, 15, 19]). We have not found any references for the asymptotic normality results for the problems we consider in GIAS models.

The models in Sections 5.2.1–5.2.4 are examples of *generalized allocation schemes* (GAS), that is,

$$\mathbb{P}(\xi_1^{(n)} = k_1, \ldots, \xi_N^{(n)} = k_N) = \mathbb{P}(\eta_1 = k_1, \ldots, \eta_N = k_N \mid \eta_1 + \cdots + \eta_N = n), \quad (5.1)$$

where $\eta_1, \ldots, \eta_N$ are independent random variables.

On the other hand the models in Sections 5.3.1 and 5.3.2 are examples of what may be called *generalized inverse allocation schemes* (GIAS), that is,

$$\mathbb{P}(\xi_1^{(n)}=k_1,\ldots,\xi_N^{(n)}=k_N)=C\mathbb{P}(\eta_1=k_1,\ldots,\eta_N=k_N|\eta_0+\eta_1+\cdots+\eta_N=n+k_1+\cdots+k_N), \quad (5.2)$$

where $\eta_0$, $\eta_1$, ..., $\eta_N$ are independent random variables, $C$ is a proportionality constant, and the equality is understood to hold whenever the right-hand side is summable. This last requirement is not vacuous: if, *e.g.*, $N = 1$ and $\eta_0 = n$ a.s., then trivially the probability on the right-hand side of (5.2) is 1 regardless of the value of $k_1$, and hence these probabilities are not summable as long as $\eta_1$ takes on infinitely many values.

In practical situations of GAS models the $(\eta_j)$ are identically distributed, and in the case of GIAS we assume that the $\eta_j$ have the same distribution for $j = 1$, ..., $N$, which may differ from the distribution of $\eta_0$.

In the derivations below we will often use the expansion

$$\left(1-\frac{a}{b}\right)^c=e^{c\log(1-\frac{a}{b})}=e^{-c\sum_{j=1}^{\infty}\frac{a^j}{jb^j}}, \quad (5.3)$$

which is valid for any $0 < |a| < |b|$ and any real $c$. We also recall (see, *e.g.*, [5, Chapter 6.5]) that

$$\mathcal{Q}_{j+1}(M) := \sum_{k=1}^{M-1}k^j=\frac{1}{j+1}\sum_{l=1}^{j+1}\binom{j+1}{l}B_{j+1-l}M^l, \quad (5.4)$$

where $(B_k)$ are the Bernoulli numbers. Clearly, $\mathcal{Q}_j$ is a polynomial of degree $j$ satisfying (4.2) with $C = 1$. For notational convenience we let

$$T(m,t)=\prod_{k=1}^{m-1}\left(1-\frac{k}{t}\right)$$

for $t > 0$ and integer $m > 0$. It follows from (5.3) and (5.4) that for $t > m$

$$T(m,t)=e^{-\sum_{j\geq 1}\frac{1}{jt^j}\mathcal{Q}_{j+1}(m)}. \quad (5.5)$$

## 5.1. Classical discrete distributions

In this subsection we re-derive asymptotic normality of

$$\frac{S_n - \mathbb{E}S_n}{\sqrt{\mathrm{Var}S_n}}$$

for laws of $S_n$ belonging to four classical families of discrete distributions: binomial, negative binomial, hypergeometric and negative hypergeometric.

**5.1.1. Binomial scheme**—Let $(\varepsilon_i)$ be a sequence of i.i.d. Bernoulli random variables, $P(\varepsilon_1 = 1) = p = 1 - P(\varepsilon_1 = 0)$. Then $S_n=\sum_{i=1}^{n}\varepsilon_i$ has the binomial $b(n, p)$ distribution. To see how Theorem 4.1 allows us to re-derive the de Moivre–Laplace theorem,

$$\frac{S_n - np}{\sqrt{np(1-p)}} \xrightarrow{d} \mathcal{N}(0,1), \quad (5.6)$$

in a simple way, we first set $L_n = np$. Then $\mathbb{E}S_n = L_n$ and Var $S_n = L_n(1 - p)$. Furthermore, $\mathbb{P}(\varepsilon_1 = \cdots = \varepsilon_i = 1) = p^i$, and thus by (2.12) it follows that the $i$th factorial moment of $S_n$ is

$$c_{i,n} = \mathbb{E}(S_n)_i = (n)_i p^i = L_n^i T(i,n). \quad (5.7)$$

Thus (5.5) implies representation (4.1) with $Q_{j+1} = -\mathcal{Q}_{j+1}$ and (5.6) follows from Remark 4.2.

**5.1.2. Negative binomial scheme**—Let $S_n$ denote the number of failures until the $n$th success in Bernoulli trials, with $p$ being the probability of success in a single trial, that is, $S_n$ is negative binomial $nb(n, p)$ with

$$\mathbb{P}(S_n = k) = \binom{n+k-1}{k}(1-p)^k p^n, k = 0, 1, \ldots.$$

We will show how Theorem 4.1 allows us to re-derive the CLT for $(S_n)$ in a simple way, which states that for $n \to \infty$

$$\frac{pS_n - n(1-p)}{\sqrt{n(1-p)}} \xrightarrow{d} \mathcal{N}(0,1). \quad (5.8)$$

Set $L_n = n(1 - p)/p$ so that $\mathbb{E}S_n = L_n$ and Var $S_n = L_n/p$. Furthermore, the $i$th factorial moment of $S_n$ is easily derived as

$$c_{i,n} = \mathbb{E}(S_n)_i = L_n^i T(i, -n).$$

Hence (4.1) holds with $Q_{j+1} = (-1)^{j+1}\mathcal{Q}_{j+1}$, and thus (5.8) follows from Remark 4.2.

**5.1.3. Hypergeometric scheme**—From an urn containing $N$ white and $M$ black balls we draw subsequently without replacement $n$ balls ($n \leq \min\{M,N\}$). For $i = 1, \ldots, n$, let $\varepsilon_i = 1$ if a white ball is drawn at the $i$th drawing and let $\varepsilon_i = 0$ otherwise. Then $S_n = \sum_{i=1}^n \varepsilon_i$ has a hypergeometric distribution $Hg(N,M; n)$, that is,

$$\mathbb{P}(S_n = k) = \frac{\binom{N}{k}\binom{M}{n-k}}{\binom{N+M}{n}}, k = 0, 1, \ldots, n.$$

Using Theorem 4.1 again, we will show that under the assumptions $N = N(n) \to \infty$, $M = M(n) \to \infty$, and $N/(N + M) \to p \in (0, 1)$ with $n \to \infty$,

$$\frac{(N+M)S_n - nN}{\sqrt{nNM(N+M-n)/(N+M-1)}} \xrightarrow{d} \mathcal{N}(0,1). \quad (5.9)$$

Setting $L_n = nN/(N + M)$ we have

$$\mathbb{E}S_n = L_n \text{ and } \operatorname{Var}S_n = L_n \frac{M(N+M-n)}{(N+M)(N+M-1)}. \quad (5.10)$$

Moreover, on noting that $(\varepsilon_1, \ldots, \varepsilon_n)$ is exchangeable by (2.12), we get

$$c_{i,n} = \mathbb{E}(S_n)_i = (n)_i \mathbb{P}(\varepsilon_1 = \cdots = \varepsilon_i = 1) = \frac{(n)_i (N)_i}{(N+M)_i} = L_n^i \frac{T(i,n)T(i,N)}{T(i,N+M)}.$$

As in earlier schemes we obtain representation (4.1) with

$$Q_{j+1} = \left[ -1 - \left( \frac{n}{N} \right)^j + \left( \frac{n}{N+M} \right)^j \right] \mathscr{Q}_{j+1}.$$

Moreover, the condition (4.3) is fulfilled since $\mathbb{E}S_n$, $\operatorname{Var}S_n$ and $L_n$ are all of order $n$. See again Remark 4.2 to conclude that (5.9) holds true.

**5.1.4. Negative hypergeometric scheme**—Let $S_n$ be a random variable with negative hyper-geometric distribution of the first kind, that is,

$$\mathbb{P}(S_n = k) = \left( \begin{array}{c} n \\ k \end{array} \right) \frac{B(\alpha_n + k, \beta_n + n - k)}{B(\alpha_n, \beta_n)}, k = 0, 1, \ldots, n,$$

with $\alpha_n = n\alpha$ and $\beta n = n\beta$. The CLT for Sn states that for $n \to \infty$

$$\frac{(\alpha+\beta)^{\frac{3}{2}} S_n - n\alpha \sqrt{\alpha+\beta}}{\sqrt{n\alpha\beta(1+\alpha+\beta)}} \xrightarrow{d} \mathscr{N}(0,1). \quad (5.11)$$

To quickly derive it from Theorem 4.1, let $L_n = n\alpha/(\alpha + \beta)$ and note that

$$\mathbb{E}S_n = L_n \text{ and } \operatorname{Var}S_n = L_n \frac{n\beta(1+\alpha+\beta)}{(\alpha+\beta)^2(n\alpha+n\beta+1)}.$$

Further, the $i$th factorial moment of $S_n$ is easily derived as

$$c_{i,n} = \mathbb{E}(S_n)_i = L_n^i \frac{T(i,n)T(i,-\alpha n)}{T(i,-(\alpha+\beta)n)}.$$

Thus, again due to (5.5) we conclude that representation (4.1) holds with

$$Q_{j+1} = \left( -1 - \frac{(-1)^j}{\alpha^j} + \frac{(-1)^j}{(\alpha+\beta)^j} \right) \mathscr{Q}_{j+1}(i).$$

The final result follows by Remark 4.2.

### 5.2. Asymptotics of occupancy in generalized allocation schemes (GAS)

In this subsection we will derive asymptotics for

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i^{(n)} = r)$$

in several generalized allocation schemes as defined at the beginning of Section 5. As we will see, when $n \to \infty$ and $N/n \to \lambda \in (0, \infty]$ the order of $\mathbb{E}S_n^{(r)}$ is $n^r/N^{r-1}$ for any $r = 0, 1,$ ..., and the order of $\mathrm{Var}S_n^{(r)}$ is the same for $r \geq 2$. When $\lambda = \infty$ and $r = 0$ or 1, the order of $\mathrm{Var}S_n^{(r)}$ is $n^2/N$. Consequently, we will derive asymptotic normality of

$$\frac{S_n - \mathbb{E}S_n^{(r)}}{\sqrt{n^r/N^{r-1}}}$$

when either

    **a.**   $r \geq 0$ and $\lambda < \infty$ or

    **b.**   $r \geq 2, \lambda = \infty$

and $n^r/N^{r-1} \to \infty$, and asymptotic normality of

$$\sqrt{N}\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{n}$$

when $\lambda = \infty$, $n^2/N \to \infty$ and $r = 0, 1$.

Although in all the cases the results look literally the same (with different asymptotic expectations and variances and with different proofs), for the sake of precision we decided to repeat formulations of theorems in each of the subsequent cases.

**5.2.1. Indistinguishable balls**—Consider a scheme of a random distribution of $n$ indistinguishable balls into $N$ distinguishable boxes, such that all distributions are equiprobable. That is, if $\xi_i = \xi_i^{(n)}$ denotes the number of balls which fall into the $i$th box, $i = 1, \ldots, N$, then

$$\mathbb{P}(\xi_1 = i_1, \ldots, \xi_N = i_N) = \binom{n+N-1}{n}^{-1}$$

for any $i_k \geq 0$, $k = 1, \ldots, N$, such that $i_1 + \cdots + i_N = n$. Note that this is a GAS and that (5.1) holds with $\eta_i \sim \mathrm{Geom}(p)$, $0 < p < 1$.

Let

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i = r)$$

denote the number of boxes with exactly $r$ balls. Note that the distribution of $(\xi_1, \ldots, \xi_N)$ is exchangeable. Moreover,

$$\mathbb{P}(\xi_1 = \cdots = \xi_i = r) = \frac{\binom{n - ri + N - i - 1}{n - ri}}{\binom{n + N - 1}{n}}.$$

Therefore, by (2.12) we get

$$c_{i,n} = \mathbb{E}(S_n^{(r)})_i = \frac{(N)_i (N-1)_i (n)_{ir}}{(n+N-1)_{i(r+1)}}. \quad (5.12)$$

Consequently,

$$\mathbb{E}S_n^{(r)} = c_{1,n} = \frac{N(N-1)(n)_r}{(n+N-1)_{r+1}},$$

and since $\mathrm{Var}S_n^{(r)} = c_{2,n} - c_{1,n}^2 + c_{1,n}$ we have

$$\mathrm{Var}S_n^{(r)} = \frac{N(N-1)^2(N-2)(n)_{2r}}{(n+N-1)_{2r+2}} - \left(\frac{N(N-1)(n)_r}{(n+N-1)_{r+1}}\right)^2 + \frac{N(N-1)(n)_r}{(n+N-1)_{r+1}}.$$

In the asymptotics below we consider the situation when $n \to \infty$ and $N/n \to \lambda \in (0, \infty]$. Then for any integer $r \quad 0$

$$\frac{N^{r-1}}{n^r}\mathbb{E}S_n^{(r)} \to \left(\frac{\lambda}{1+\lambda}\right)^{r+1} (=1 \text{ for } \lambda = \infty). \quad (5.13)$$

It is also elementary but more laborious to prove that, for any $r \quad 2$ and $\lambda \in (0, \infty]$ or $r = 0, 1$ and $\lambda \in (0, \infty)$,

$$\frac{N^{r-1}}{n^r}\mathrm{Var}S_n^{(r)} \to \left(\frac{\lambda}{1+\lambda}\right)^{r+1}\left(1 - \frac{\lambda(1+\lambda+(\lambda r - 1)^2)}{(1+\lambda)^{r+2}}\right) =: \sigma_r^2 (=1 \text{ for } \lambda = \infty). \quad (5.14)$$

Further, for $\lambda = \infty$,

$$\frac{N}{n^2}\mathrm{Var}S_n^{(0)} \to 1 =: \tilde{\sigma}_0^2 \text{ and } \frac{N}{n^2}\mathrm{Var}S_n^{(1)} \to 4 =: \tilde{\sigma}_1^2.$$

Similarly, in this case

$$\frac{N}{n^2}\mathrm{Cov}(S_n^{(0)}, S_n^{(1)}) \to -2 \text{ and } \frac{N}{n^2}\mathrm{Cov}(S_n^{(0)}, S_n^{(2)}) \to 1,$$

and thus for the correlation coefficient we have

$$\rho(S_n^{(0)}, S_n^{(1)}) \to -1 \text{and} \rho(S_n^{(0)}, S_n^{(2)}) \to 1. \quad (5.15)$$

Now we are ready to deal with CLTs.

**Theorem 5.1.** *Let $N/n \to \lambda \in (0, \infty]$. Let either*

    **a.**    *$r$   $0$ and $\lambda < \infty$, or*

    **b.**    *$r$   $2, \lambda = \infty$ and $n^r/N^{r-1} \to \infty$.*

*Then*

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n^r/N^{r-1}}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2).$$

**Proof.** Note that (5.12) can be rewritten as

$$c_{i,n} = L_n^i \frac{T(i,N)T(i,N-1)T(ir,n)}{T(i(r+1), n+N-1)} \text{with} L_n = \frac{N(N-1)n^r}{(n+N-1)^{r+1}}.$$

Therefore, as in the previous cases, using (5.5) we conclude that representation (4.1) holds with

$$Q_{j+1}(i) = -\left[\left(\frac{n}{N}\right)^j + \left(\frac{n}{N-1}\right)^j\right]\mathcal{Q}_{j+1}(i) - \mathcal{Q}_{j+1}(ri) + \left(\frac{n}{n+N-1}\right)^j \mathcal{Q}_{j+1}((r+1)i).$$

To conclude the proof we note that $\mathbb{E}S_n^{(r)}$, $\mathrm{Var}S_n^{(r)}$ and $L_n$ are of the same order, $n^r/N^{r-1}$, and we use Remark 5.2 stated below.

**Remark 5.2.** If $\mathbb{E}S_n^{(r)}$ and $\mathrm{Var}S_n^{(r)}$ are of the same order and diverge to $\infty$, then (2.6) and (2.7) hold. Moreover, if $L_n$ and $\mathrm{Var}S_n^{(r)}$ are both of order $n^r/N^{r-1}$ then the left-hand side of (4.3) is of order

$$\frac{1}{n^{J-1}}\left(\frac{n^r}{N^{r-1}}\right)^{\frac{J}{2}} = \left(\frac{n}{N}\right)^{\frac{J}{2}(r-1)}\frac{1}{n^{\frac{J}{2}-1}}.$$

That is, when either $\lambda \in (0, \infty)$ and $r = 0, 1, \ldots$ or $\lambda = \infty$ and $r = 2, 3, \ldots$, the condition (4.3) is satisfied.

In the remaining cases we use asymptotic correlations.

**Theorem 5.3.** *Let $N/n \to \infty$ and $n^2/N \to \infty$. Then for $r = 0, 1$*

$$\sqrt{N}\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{n} \xrightarrow{d} \mathcal{N}(0, \tilde{\sigma}_r^2).$$

**Proof.** Due to the second equation in (5.15), it follows that

$$\sqrt{N}\frac{S_n^{(0)} - \mathbb{E}S_n^{(0)}}{\tilde{\sigma}_0 n} - \sqrt{N}\frac{S_n^{(2)} - \mathbb{E}S_n^{(2)}}{\sigma_2 n} \xrightarrow{L^2} 0.$$

Therefore the result for $r = 0$ holds. Similarly, for $r = 1$ it suffices to observe that the first equation in (5.15) implies

$$\sqrt{N}\frac{S_n^{(0)} - \mathbb{E}S_n^{(0)}}{\tilde{\sigma}_0 n} + \sqrt{N}\frac{S_n^{(1)} - \mathbb{E}S_n^{(1)}}{\tilde{\sigma}_1 n} \xrightarrow{L^2} 0.$$

### 5.2.2. Distinguishable balls

Consider a scheme of a random distribution of $n$ distinguishable balls into $N$ distinguishable boxes, such that any such distribution is equally likely. Then, if $\xi_i = \xi_i^{(n)}$ denotes the number of balls which fall into the $i$th box, $i = 1, \ldots, N$,

$$\mathbb{P}(\xi_1 = i_1, \ldots, \xi_N = i_N) = \frac{n!}{i_1! \cdots i_N!} N^{-n}$$

for any $i_l \geq 0$, $l = 1, \ldots, N$, such that $i_1 + \cdots + i_N = n$. This is a GAS with $\eta_i \sim \text{Poisson}(\lambda)$, $\lambda > 0$, in (5.1).

For a fixed non-negative integer $r$ let

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i = r)$$

be the number of boxes with exactly $r$ balls. Obviously, the distribution of $(\xi_1, \ldots, \xi_N)$ is exchangeable, and

$$\mathbb{P}(\xi_1 = \cdots = \xi_i = r) = \frac{n!}{(r!)^i (n-ir)!} N^{-ir} \left(1 - \frac{i}{N}\right)^{n-ir}.$$

Therefore, by (2.12) we get

$$c_{i,n} = \mathbb{E}(S_n^{(r)})_i = \frac{(N)_i (n)_{ir}}{(r!)^i N^{ri}} \left(1 - \frac{i}{N}\right)^{n-ir}. \quad (5.16)$$

Consequently, for any $r = 0, 1, \ldots,$

$$\mathbb{E}S_n^{(r)} = c_{1,n} = \frac{(n)_r \left(1 - \frac{1}{N}\right)^{n-r}}{r! N^{r-1}}$$

and

$$\text{Var}S_n^{(r)} = c_{2,n} - c_{1,n}^2 + c_{1,n} = \frac{(N-1)(n)_{2r}(1-\frac{2}{N})^{n-2r}}{(r!)^2 N^{2r-1}} - \frac{(n)_r^2(1-\frac{1}{N})^{2(n-r)}}{(r!)^2 N^{2(r-1)}} + \frac{(n)_r(1-\frac{1}{N})^{n-r}}{r! N^{r-1}}.$$

In the asymptotics below we consider the situation when $n \to \infty$ and $N/n \to \lambda \in (0, \infty]$. Then, for any integer $r \geq 0$,

$$\lim_{n\to\infty} \frac{N^{r-1}}{n^r} \mathbb{E}S_n^{(r)} = \frac{1}{r!} e^{-\frac{1}{\lambda}} \left( = \frac{1}{r!} \text{for} \lambda = \infty \right). \quad (5.17)$$

It is also elementary but more laborious to check that, for any fixed $r \geq 2$ and $\lambda \in (0, \infty]$ or $r = 0, 1$ and $\lambda \in (0,\infty)$,

$$\lim_{n\to\infty} \frac{N^{r-1}}{n^r} \text{Var}S_n^{(r)} = \frac{e^{-\frac{1}{\lambda}}}{r!} \left( 1 - \frac{e^{-\frac{1}{\lambda}}(\lambda + (\lambda r - 1)^2)}{r! \lambda^{r+1}} \right) := \sigma_r^2 \left( = \frac{1}{r!} \text{for} \lambda = \infty \right). \quad (5.18)$$

Further, for $\lambda = \infty$,

$$\frac{N}{n^2} \text{Var}S_n^{(0)} \to \frac{1}{2} =: \tilde{\sigma}_0^2 \text{and} \frac{N}{n^2} \text{Var}S_n^{(1)} \to 2 =: \tilde{\sigma}_1^2.$$

Similarly one can prove that

$$\frac{N}{n^2} \text{Cov}(S_n^{(0)}, S_n^{(1)}) \to -1.$$

Therefore, for the correlation coefficients we have

$$\rho(S_n^{(0)}, S_n^{(1)}) \to -1. \quad (5.19)$$

Since

$$\frac{N}{n^2} \text{Var}S_n^{(2)} \to \frac{1}{2} = \sigma_2^{(2)} \text{and} \frac{N}{n^2} \text{Cov}(S_n^{(0)}, S_n^{(2)}) \to \frac{1}{2}$$

we also have

$$\rho(S_n^{(0)}, S_n^{(2)}) \to 1. \quad (5.20)$$

We consider the cases when $r \geq 2$ and $\lambda = \infty$ or $r \geq 0$ and $\lambda \in (0, \infty)$.

**Theorem 5.4.** *Let $N/n \to \lambda \in (0, \infty]$. Let either*

    **a.**   *$r \geq 0$ and $\lambda < \infty$, or*

    **b.**   *$r \geq 2$, $\lambda = \infty$ and $n^r/N^{r-1} \to \infty$.*

*Then*

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n^r/N^{r-1}}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2).$$

**Proof.** Write $c_{i,n}$ as

$$c_{i,n} = L_n^i e^{i\frac{n}{N}} \left(1 - \frac{i}{N}\right)^{n-ir} T(i, N) T(ir, n), \text{ where } L_n = \frac{n^r e^{-\frac{n}{N}}}{r! N^{r-1}}.$$

Then, the representation (4.1) holds with

$$Q_{j+1}(i) = \left(r - \frac{j}{j+1}\frac{n}{N}\right)\left(\frac{n}{N}\right)^j i^{j+1} - \left(\frac{n}{N}\right)^j \mathscr{Q}_{j+1}(i) - \mathscr{Q}_{j+1}(ri).$$

Since $\mathbb{E}S_n^{(r)}$, $\mathrm{Var}\, S_n^{(r)}$ and $L_n$ are of order $n^r/N^{r-1}$, the final result follows by Remark 5.2.

As in the case of indistinguishable balls, using (5.20) and (5.19) we get the following.

**Theorem 5.5.** *Let $N/n \to \infty$ and $n^2/N \to \infty$. Then for $r = 0, 1$*

$$\sqrt{N}\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{n} \xrightarrow{d} \mathcal{N}(0, \tilde{\sigma}_r^2).$$

**5.2.3. Coloured balls**—An urn contains $NM$ balls, $M$ balls of each of $N$ colours. From the urn a simple random sample of $n$ elements is drawn. We want to study the asymptotics of the number of colours with exactly $r$ balls in the sample. More precisely, let $\xi_i = \xi_i^{(n)}$ denote the number of balls of colour $i$, $i = 1, \ldots, N$. Then

$$\mathbb{P}(\xi_1 = k_1, \ldots, \xi_N = k_N) = \frac{\prod_{i=1}^{N} \binom{M}{k_i}}{\binom{NM}{n}}$$

for all integers $k_i \geq 0$, $i = 1, \ldots, N$, such that $\sum_{i=1}^{N} k_i = n$. Obviously, the random vector $(\xi_1, \ldots, \xi_N)$ is exchangeable and the GAS equation (5.1) holds with $\eta_i \sim b(M, p)$, $0 < p < 1$.

For an integer $r \geq 0$ we want to study the asymptotics of

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i = r).$$

For the $i$th factorial moment we get

$$c_{i,n}=(N)_i\mathbb{P}(\xi_1=\cdots\xi_i=r)=(N)_i\frac{\dbinom{M}{r}^i\dbinom{(N-i)M}{n-ri}}{\dbinom{NM}{n}}. \quad (5.21)$$

Consequently, for any $r = 0, 1, \ldots,$

$$\mathbb{E}S_n^{(r)}=c_{1,n}=N\frac{\dbinom{M}{r}\dbinom{(N-1)M}{n-r}}{\dbinom{NM}{n}}$$

and

$$\operatorname{Var}S_n^{(r)}=c_{2,n}-c_{1,n}^2+c_{1,n}=N(N-1)\frac{\dbinom{M}{r}^2\dbinom{(N-2)M}{n-2r}}{\dbinom{NM}{n}}$$

$$-N^2\frac{\dbinom{M}{r}^2\dbinom{(N-1)M}{n-r}^2}{\dbinom{NM}{n}^2}$$

$$+N\frac{\dbinom{M}{r}\dbinom{(N-1)M}{n-r}}{\dbinom{NM}{n}}.$$

In the asymptotics below we consider the situation when $n \to \infty$, $N/n \to \lambda \in (0, \infty]$, and $M = M(n)$ $n$.

Although the computational details are different, asymptotic formulas for $\mathbb{E}S_n^{(r)}, \operatorname{Var}S_n^{(r)}, \operatorname{Cov}(S_n^{(0)}, S_n^{(1)})$ and $\operatorname{Cov}(S_n^{(0)}, S_n^{(2)})$ are literally the same as for their counterparts in the case of occupancy for distinguishable balls studied in Section 5.2.2.

First we will consider the limit result in the case $r$ $2, \lambda = \infty$, and $r$ $0, \lambda \in (0, \infty)$.

**Theorem 5.6.** *Let $N/n \to \lambda \in (0, \infty]$ and $M = M(n)$ $n$. Let either*

   **a.** $r$ $0$ *and* $\lambda < \infty$, *or*

   **b.** $r$ $2, \lambda = \infty$ *and* $n^r/N^{r-1} \to \infty$.

*Then*

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n^r/N^{r-1}}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2).$$

**Proof.** Rewrite the formula (5.21) as

$$c_{i,n} = L_n^i \frac{T((M-r)i, NM-n)T(i,N)T(ri,n)}{T(Mi, NM)} \text{ with } L_n = Nn^r \binom{M}{r} \frac{(1 - \frac{n}{NM})^M}{(NM-n)^r}.$$

Thus the representation (4.1) holds with

$$Q_{j+1}(i) = \left(\frac{n}{NM}\right)^j \mathcal{Q}_{j+1}(Mi) - \left(\frac{n}{NM-n}\right)^j \mathcal{Q}_{j+1}((M-r)i) - \left(\frac{n}{N}\right)^j \mathcal{Q}_{j+1}(i) - \mathcal{Q}_{j+1}(ri).$$

We need to see that the polynomials $Q_j$ satisfy bound (4.2). This is clearly true for each of the last two terms in the above expression for $Q_{j+1}$. For the first two terms we have

$$\left| \left(\frac{n}{NM}\right)^j \mathcal{Q}_{j+1}(Mi) \right.$$
$$- \left(\frac{n}{NM-n}\right)^j \mathcal{Q}_{j+1}((M-r)i) \right|$$
$$= \left| \left(\frac{n}{NM}\right)^j \sum_{k=1}^{Mi-1} k^j \right.$$
$$- \left(\frac{n}{NM-n}\right)^j \sum_{k=1}^{(M-r)i-1} k^j \left| \le \right| \left(\frac{n}{NM}\right)^j$$
$$- \left(\frac{n}{NM-n}\right)^j \mathcal{Q}_{j+1}(Mi)$$
$$+ \left(\frac{n}{NM-n}\right)^j \sum_{k=(M-r)i}^{Mi-1} k^j.$$

Since

$$\left| \left(\frac{n}{NM}\right)^j - \left(\frac{n}{NM-n}\right)^j \right| \le \left(\frac{n}{NM-n}\right)^j \frac{jn}{NM} \le \left(\frac{2n}{NM-n}\right)^j \frac{n}{NM},$$

and

$$\sum_{k=(M-r)i}^{Mi-1} k^j \le rM^j i^{j+1},$$

and $\mathcal{Q}_{j+1}(M_i) \quad M^{j+1} i^{j+1}$, we conclude that the $Q_j$ do satisfy (4.2).

Clearly, $\mathbb{E}S_n^{(r)}$, $\operatorname{Var}S_n^{(r)}$ and $L_n$ are of order $n^r/N^{r-1}$, and again we conclude the proof by referring to Remark 5.2.

Asymptotic normality for $S_n^{(1)}$ and $S_n^{(0)}$ for $\lambda = \infty$ also holds with an identical statement to that of Theorem 5.5 for distinguishable balls.

**5.2.4. Rooted trees in random forests**—Let $\mathscr{T}N, n)$ denote a forest with $N$ roots (that is, $N$-rooted trees) and $n$ non-root vertices. Consider a uniform distribution on the set of such $\mathscr{T}N, n)$ forests. Let $\xi_i = \xi_i^{(n)}$, $i = 1, \ldots, N$, denote the number of non-root vertices in the $i$th tree. Then (see, *e.g.*, [3], [18] or [19]), for any $k_i \geq 0$ such that $\sum_{i=1}^{N} k_i = n$,

$$\mathbb{P}(\xi_1 = k_1, \ldots, \xi_N = k_N) = \frac{n!}{\prod_{i=1}^{N} k_i} \frac{\prod_{i=1}^{N}(k_i+1)^{k_i-1}}{N(N+n)^{n-1}}.$$

Note that this distribution is exchangeable and that it is a GAS with $\eta_i$ in (5.1) given by

$$\mathbb{P}(\eta_i = k) = \frac{\lambda^k(k+1)^{k-1}}{k!} e^{-(k+1)\lambda}, k = 0, 1, \ldots, \lambda > 0.$$

We mention in passing that the distribution of $\eta_i$ may be identified as an Abel distribution discussed in [17] with (in their notation) $p = 1$ and $\theta = \ln\lambda - \lambda$. We refer to [17, Example D] for more information on Abel distributions, including further references.

For a fixed number $r \geq 0$ we are interested in the number $S_n^{(r)}$ of trees with exactly $r$ non-root vertices:

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i = r).$$

Since the $i$th factorial moment of $S_n^{(r)}$ is of the form

$$c_{i,n} = (N)_i \mathbb{P}(\xi_1 = \cdots = \xi_i = r),$$

we have to find the marginal distributions of the random vector $(\xi_1, \ldots, \xi_N)$. From the identity

$$s\sum_{k=0}^{m} \binom{m}{k}(k+1)^{k-1}(m-k+s)^{m-k-1} = (s+1)(m+1+s)^{m-1},$$

which is valid for any natural $m$ and $s$, we easily obtain that, for $k_j \geq 0$ such that $\sum_{j=1}^{i+1} k_j = n$,

$$\mathbb{P}(\xi_1 = k_1, \ldots, \xi_i = k_i) = \frac{n!}{\prod_{j=1}^{i+1} k_j!} \frac{(N-i)(k_{i+1}+N-i)^{k_{i+1}-1}\prod_{j=1}^{i}(k_j+1)^{k_j-1}}{N(N+n)^{n-1}}.$$

Therefore

$$c_{i,n} = \frac{(r+1)^{i(r-1)}}{(r!)^i} \frac{N-i}{N} \frac{(N)_i (n)_{ri}}{(n+N-(r+1)i)^{ri}} \left(1 - \frac{(r+1)i}{n+N}\right)^{n-1}. \quad (5.22)$$

Hence

$$\mathbb{E}S_n^{(r)} = c_{1,n} = \frac{(r+1)^{r-1}}{r!} \frac{(N-1)(n)_r}{(n+N-r-1)^r} \left(1 - \frac{r+1}{n+N}\right)^{n-1}.$$

Thus, if $N/n \to \lambda \in (0, \infty]$ we have

$$\frac{N^{r-1}}{n^r} \mathbb{E}S_n^{(r)} \to \frac{(r+1)^{r-1}}{r!} \left(\frac{\lambda}{\lambda+1}\right)^r e^{-\frac{r+1}{\lambda+1}} \left(= \frac{(r+1)^{r-1}}{r!} \text{for} \lambda = \infty\right).$$

Since $\text{Var}S_n^{(r)} = c_{2,n} - c_{1,n}^2 + c_{1,n}$, elementary but cumbersome computations lead to

$$\frac{N^{r-1}}{n^r} \text{Var}S_n^{(r)} \to \sigma_r^2 = \frac{e^{-\frac{r+1}{\lambda+1}}}{(r+1)!} \left(\frac{\lambda(r+1)}{\lambda+1}\right)^r \left[1 - \frac{e^{-\frac{r+1}{\lambda+1}}}{(r+1)!} \left(\frac{r+1}{\lambda+1}\right)^r\right] - \lambda \left(\frac{e^{-\frac{r+1}{\lambda+1}}(\lambda r - 1)}{(r+1)!(\lambda+1)}\right)^2 \left(\frac{\lambda(r+1)^2}{(\lambda+1)^2}\right)^r$$

for $r \geq 2$ and $\lambda \in (0, \infty]$ and for $r = 0, 1$ and $\lambda \in (0, \infty)$. For $r = 0, 1$ and $\lambda = \infty$, and $n^2/N \to \infty$, we have

$$\frac{N}{n^2} \text{Var}S_n^{(0)} \to \frac{3}{2} = \tilde{\sigma}_0^2 \text{and} \frac{N}{n^2} \text{Var}S_n^{(1)} \to 6 = \tilde{\sigma}_1^2.$$

Similarly one can prove that

$$\frac{N}{n^2} \text{Cov}(S_n^{(0)}, S_n^{(1)}) \to -3.$$

Therefore, for the correlation coefficients we have

$$\rho(S_n^{(0)}, S_n^{(1)}) \to -1. \quad (5.23)$$

Since

$$\frac{N}{n^2} \text{Var}S_n^{(2)} \to \frac{3}{2} = \sigma_2^{(2)} \text{and} \frac{N}{n^2} \text{Cov}(S_n^{(0)}, S_n^{(2)}) \to 1,$$

we also have

$$\rho(S_n^{(0)}, S_n^{(2)}) \to 1. \quad (5.24)$$

**Theorem 5.7.** *Let $N/n \to \lambda \in (0, \infty]$.* *Let either*

**a.**   $r$   0 *and* $\lambda < \infty$, *or*

**b.**   $r$   2, $\lambda = \infty$ *and* $n^r/N^{r-1} \to \infty$.

*Then*

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n^r/N^{r-1}}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2).$$

**Proof.** Since the asymptotics of $\mathrm{Var}S_n^{(r)}$ and $\mathbb{E}S_n^{(r)}$ is of the same order as in Theorem 5.1, the conditions (2.6) and (2.7) are satisfied. Using (5.22) we write

$$c_{i,n} = L_n^i e^{i\frac{(n-1)(r+1)}{n+N}} \left(1 - \frac{(r+1)i}{n+N}\right)^{n-1-ri} T(i+1, N)T(ri, n),$$

where

$$L_n = \frac{N(r+1)^{r-1}}{r!} \left(\frac{n}{n+N}\right)^r e^{-\frac{(n-1)(r+1)}{n+N}}.$$

Thus the representation (4.1) holds true with

$$Q_{j+1}(i) = \left(r - \frac{j(r+1)(n-1)}{(j+1)(n+N)}\right) \left(\frac{(r+1)(n-1)}{n+N}\right)^j i^{j+1} - \left(\frac{n}{N}\right)^j \mathscr{Q}_{j+1}(i+1) - \mathscr{Q}_{j+1}(ri).$$

The final result follows again by Remark 5.2, on noting that $\mathbb{E}S_n^{(r)}$, $\mathrm{Var}S_n^{(r)}$ and $L_n$ are of order $n^r/N^{r-1}$.

Again, as in previous cases we use (5.23) and (5.24) to obtain the following result.

**Theorem 5.8.** *Let* $N/n \to \infty$ *and* $n^2/N \to \infty$. *Then for* $r = 0, 1$

$$\sqrt{N}\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{n} \xrightarrow{d} \mathcal{N}(0, \tilde{\sigma}_r^2).$$

## 5.3. Asymptotics in generalized inverse allocation schemes (GIAS)

Our final two settings are examples of the GIAS as defined in (5.2). As in the case of the GAS for

$$S_n^{(r)} = \sum_{i=1}^{N} I(\xi_i^{(n)} = r),$$

we will obtain asymptotic normality of

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n}}$$

when $N/n \to \lambda \in (0, \infty)$.

### 5.3.1. Exchangeable negative multinomial model

Let $(\xi_i) = \xi_i^{(n)}$ be a random vector with a negative multinomial distribution, that is,

$$\mathbb{P}(\xi_1 = k_1, \ldots, \xi_N = k_N) = \frac{(n + \sum_{j=1}^N k_j)!}{n! \prod_{j=1}^N k_j!} p^{\sum_{j=1}^N k_j} (1 - Np)^{n+1}. \quad (5.25)$$

Note that this is an exchangeable case of a model usually referred to as the Bates–Neyman model, introduced in [1]. We refer to [9, Chapter 36, Sections 1–4] for a detailed account of this distribution, its properties, applications, and further references. Here, we just note that this is a GIAS for which (5.2) holds with $\eta_0 \sim \text{Poisson}(\lambda_0)$, $\eta_i \sim \text{Poisson}(\lambda_1)$, $i = 1, \ldots, N$, and $C = \lambda_0/(\lambda_0 + N\lambda_1)$. Thus (5.2) implies (5.25) with $p = \lambda_1/(\lambda_0 + N\lambda_1)$.

For a fixed integer $r$ we are interested in the asymptotics of

$$S_n^{(r)} = \sum_{j=1}^N I(\xi_j = r).$$

Denoting $\beta_n = (Np)^{-1} - 1$, we obtain

$$c_{i,n} = (N)_i \frac{(n + ri)!}{n!(r!)^i} \frac{(N\beta_n)^{n+1}}{(N\beta_n + i)^{n+1+ri}}. \quad (5.26)$$

To study the asymptotic properties of $S_n^{(r)}$ we will assume that $N/n \to \lambda \in (0, \infty)$. Moreover we let $p = p_n$ depend on $n$ in such a way that $Np_n \to \alpha \in (0, 1)$, *i.e.*, $\beta_n \to \alpha^{-1} - 1$. Consequently, setting $\Delta := \alpha/(\lambda(1 - \alpha))$, for any $r = 0, 1, \ldots$,

$$\lim_{n \to \infty} \frac{1}{n} \mathbb{E}S_n^{(r)} = \frac{\lambda \Delta^r}{r!} e^{-\Delta}$$

and

$$\lim_{n \to \infty} \frac{1}{n} \text{Var} S_n^{(r)} = \frac{\lambda \Delta^r e^{-\Delta}}{r!} \left[ 1 - \frac{\Delta^r e^{-\Delta}}{r!} (1 - \lambda(r - \Delta)^2) \right] =: \sigma_r^2. \quad (5.27)$$

**Theorem 5.9.** *Let $N/n \to \lambda \in (0, \infty)$ and $Np_n \to \alpha \in (0, 1)$. Then, for any $r = 0, 1, \ldots$,*

$$\frac{S_n^{(r)} - \mathbb{E}S_n^{(r)}}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2),$$

*with $\sigma_r^2$ defined in (5.27).*

**Proof.** Write (5.26) as

$$c_{i,n} = L_n^i e^{i \frac{n}{N\beta_n}} \frac{T(ri+1, -n)T(i+1, N)}{(1+\frac{i}{N\beta_n})^{n+1+ri}} \text{ with } L_n = \frac{n^r e^{-\frac{n}{N\beta_n}}}{r! N^{r-1} \beta_n^r}.$$

Thus the representation (4.1) holds true with

$$Q_{j+1}(i) = -\frac{i}{N\beta_n} + \left(r - \frac{j(n+1)}{(j+1)N\beta_n}\right)\left(-\frac{n}{N\beta_n}\right)^j i^{j+1} + (-1)^{j+1} \mathcal{Q}_{j+1}(ri+1) - \left(\frac{n}{N}\right)^j \mathcal{Q}_{j+1}(i).$$

Moreover, $\mathbb{E}S_n^{(r)}, \operatorname{Var}S_n^{(r)}$ and $L_n$ are all of order $n$ and thus the final conclusion follows from Remark 4.2.

**5.3.2. Dirichlet negative multinomial model**—Finally we consider an exchangeable version of what is known as the 'Dirichlet model of buying behaviour', introduced in a seminal paper by Goodhart, Ehrenberg and Chatfield [4] and subsequently studied in numerous papers up to the present. This distribution is also mentioned in [9, Chapter 36, Section 6]. Writing, as usual, $\xi_i^{(n)} = \xi_i$, the distribution under consideration has the form

$$\mathbb{P}(\xi_1 = k_1, \ldots, \xi_N = k_N) = \frac{(n+\sum_{i=1}^N k_i)!}{n! \prod_{i=1}^N k_i!} \frac{\Gamma(Na+b)}{\Gamma^N(a)\Gamma(b)} \frac{\Gamma(n+1+b)\prod_{i=1}^N \Gamma(k_i+a)}{\Gamma(Na+b+n+1+\sum_{i=1}^N k_i)}$$

for any $k_i = 0, 1, \ldots, i = 1, \ldots, N$. Here $n > 0$ is an integer and $a, b > 0$ are parameters. This is again a GIAS, for which (5.2) holds with $\eta_i \sim nb(a, p)$, $i = 1, \ldots, N$, $\eta_0 \sim nb(b + 1, p)$, $0 < p < 1$, and $C = b/(Na + b)$. When $a$ and $b$ are integers we recall a nice interpretation of $(\xi_1, \ldots, \xi_N)$ via the Pólya urn scheme. An urn contains $b$ black balls and $a$ balls in each of $N$ non-black colours. In subsequent steps a ball is drawn at random and returned to the urn together with one ball of the same colour. The experiment is continued until the $n$th black ball is drawn. Then $\xi_i$ is the number of balls of the $i$th colour at the end of experiment, $i = 1, \ldots, N$. This distribution can also be viewed as multivariate negative hypergeometric law of the second kind.

From the fact that $c_{i,n} = (N)_i \mathbb{P}(\xi_1 = \cdots = \xi_i = r)$, we get

$$c_{i,n} = (N)_i \frac{(n+ri)!}{n!(r!)^i} \frac{\Gamma(ia+b)}{\Gamma^i(a)\Gamma(b)} \frac{\Gamma(n+1+b)\Gamma^i(r+a)}{\Gamma((r+a)i+n+1+b)}. \quad (5.28)$$

To study the asymptotic behaviour of $S_n^{(r)}$ we will assume that $N/n \to \lambda \in (0, \infty)$ and that $b = b_n$ depends on $n$ in such a way that $b_n/n \to \beta > 0$.

Below we use the following product representation of the Gamma function:

$$\Gamma(x) = \frac{1}{xe^{\gamma x}} \prod_{k \geq 1} \frac{e^{\frac{x}{k}}}{1+\frac{x}{k}}, \quad (5.29)$$

where $\gamma$ is the Euler constant and $x > 0$. We also recall (see, *e.g.*, [23, Section 12.16]) that for a digamma function $\Psi(x) = d \ln(\Gamma(x))/dx$ we have

$$\Psi(x+1) = -\gamma + \sum_{k \geq 1} \left( \frac{1}{k} - \frac{1}{k+x} \right), x \neq -1, -2, \ldots.$$

Then, for any $0 < x < y$ we can write

$$\frac{\Gamma(x+y)}{\Gamma(y)} = e^{x \Psi(y+1)} \frac{y}{x+y} e^{\sum_{k \geq 1} \left( \frac{x}{k+y} - \log\left(1 + \frac{x}{k+y}\right) \right)}, \quad (5.30)$$

and the series

$$\sum_{k \geq 1} \left( \frac{x}{k+y} - \log\left(1 + \frac{x}{k+y}\right) \right)$$

converges. Note that

$$\Psi(y) - \ln y \to 0, \text{ as } y \to \infty, \quad (5.31)$$

so that, if $\alpha_n/n \to \alpha$, then for any $x > 0$

$$n^{-x} \frac{\Gamma(\alpha_n + x)}{\Gamma(\alpha_n)} \to \alpha^x.$$

Consequently,

$$\lim_{n \to \infty} \frac{1}{n} \mathbb{E} S_n^{(r)} = \frac{\lambda \Gamma(a+r)}{r! \Gamma(a)} \frac{\beta^a}{(1+\beta)^{a+r}}.$$

Similarly,

$$\lim_{n \to \infty} \frac{1}{n} \mathrm{Var} S_n^{(r)} = \frac{\lambda \Gamma(r+a) \beta^a}{r! \Gamma(a) (1+\beta)^{a+r}} \left( 1 - \frac{\Gamma(r+a) \beta^a [1 + \lambda(\frac{(a+r)^2}{\beta+1} - \frac{\alpha^2}{\beta} - r^2)]}{r! \Gamma(a)(1+\beta)^{a+r}} \right) =: \sigma_r^2. \quad (5.32)$$

**Theorem 5.10.** *Let $N/n \to \lambda \in (0, \infty)$ and $b_n/n \to \beta \in (0, \infty)$. Then, for any $r = 0, 1, \ldots,$*

$$\frac{S_n^{(r)} - \mathbb{E} S_n^{(r)}}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, \sigma_r^2),$$

with $\sigma_r^2$ defined in (5.32).

**Proof.** Note that (5.28) can be written as

$$c_{i,n}=\left(\frac{N\Gamma(r+a)n^r}{r!\Gamma(a)}\right)^i T(ir+1,-n)T(i,N)\frac{\Gamma(b_n+ia)}{\Gamma(b_n)}\frac{\Gamma(n+1+b_n)}{\Gamma(n+1+b_n+i(r+a))}.$$

Moreover, setting

$$h_j(x)=\sum_{k\geq 1}\frac{1}{(k+x)^{j+1}}, x>0, j\geq 2,$$

we see that (5.30) can be developed into

$$\frac{\Gamma(x+y)}{\Gamma(y)}=e^{x\Psi(y+1)}e^{\sum_{j\geq 1}\left(\frac{(-x)^j}{jy^j}+\frac{(-x)^{j+1}}{j+1}h_j(y)\right)}.$$

Therefore, taking $(x, y) = (ia, b_n)$ and $(x, y) = (i(r + a), n + 1 + b_n)$, we decompose $c_{i,n}$ according to (4.1), where

$$L_n=\frac{N\Gamma(r+a)n^r}{r!\Gamma(a)}e^{a\Psi(1+b_n)-(r+a)\Psi(2+n+b_n)}$$

and

$$Q_{j+1}(i)=\left[\left(\frac{n(r+a)}{b_n+n+1}\right)^j\right.$$
$$-\left(\frac{na}{b_n}\right)^j(-1)^{j+1}i^j+\frac{j(-n)^j}{j+1}[(r+a)^{j+1}h_j(b_n$$
$$+n+1)-a^{j+1}h_j(b_n)]i^{j+1}$$
$$+(-1)^{j+1}\mathcal{Q}_{j+1}(ir$$
$$\left.+1)-\left(\frac{n}{N}\right)^j\mathcal{Q}_{j+1}(i).$$

On noting that $a_n/n \to a$ implies that $n^j h_j(a_n) < c^j(a)$ uniformly in $n$, we conclude that polynomials $Q_j$ satisfy condition (4.2). Moreover, (5.31) yields that $L_n$ is of order $n$. Since $\mathbb{E}S_n^{(r)}$ and $\mathrm{Var}S_n^{(r)}$ are of order $n$ too, the result follows by Remark 4.2.

## Acknowledgments

## References

1. Bates GE, Neyman J. Contributions to the theory of accident proneness I: An optimistic model of the correlation between light and severe accidents. Univ. California Publ. Statist. 1952; 1:215–253.

2. Billingsley, P. Probability and Measure. third edition. Wiley; 1995.

3. Chuprunov A, Fazekas I. An inequality for moments and its applications to the generalized allocation scheme. Publ. Math. Debrecen. 2010; 76:271–286.

4. Goodhart GJ, Ehrenberg ASC, Chatfield C. The Dirichlet: A comprehensive model of buying behaviour. J. Royal Statist. Soc. Section A. 1984; 147:621–655.

5. Graham, RL.; Knuth, DE.; Patashnik, O. Concrete Mathematics. Addison-Wesley; 1994.

6. Haldane JBS. On a method of estimating frequencies. Biometrika. 1945; 33:222–225. [PubMed: 21006837]

7. Hwang H-K, Janson S. Local limit theorems for finite and infinite urn models. Ann. Probab. 2008; 36:992–1022.

8. Janson, S.; Łuczak, T.; Ruci ski, A. Random Graphs. Wiley-Interscience Series in Discrete Mathematics and Optimization, Wiley-Interscience; 2000.

9. Johnson, NL.; Kotz, S.; Balakrishnan, N. Discrete Multivariate Distributions. Wiley Series in Probability and Statistics, Wiley; 1997.

10. Johnson, NL.; Kotz, S.; Kemp, AW. Univariate Discrete Distributions. second edition. Wiley Series in Probability and Statistics, Wiley; 1992.

11. Keener, RW.; Wu, WB. Random Walk, Sequential Analysis and Related Topics. World Scientific; 2006. On Dirichlet multinomial distributions; p. 118-130.

12. Kendall, MG.; Stuart, A. The Advanced Theory of Statistics, Vol. 1: Distribution Theory. Addison-Wesley; 1969.

13. Kolchin VF. A certain class of limit theorems for conditional distributions. Litovsk. Mat. Sb. 1968; 8:53–63.

14. Kolchin, VF. Random Graphs. Vol. 53 of Encyclopedia of Mathematics and its Applications. Cambridge University Press; 1999.

15. Kolchin, VF. Random Mappings. Optimization Software; 1986.

16. Kolchin, VF.; Sevastyanov, BA.; Chistyakov, VP. Random Allocations. Winston; 1978.

17. Letac G, Mora M. Natural real exponential families with cubic variance functions. Ann. Statist. 1990; 18:1–37.

18. Pavlov YL. Limit theorems for the number of trees of a given size in a random forest. Mat. Sbornik. 1977; 103:335–345.

19. Pavlov, YL. Random Forests. VSP; 2000.

20. Stanley, RP. Enumerative Combinatorics. Vol. Vol. 2. Cambridge University Press; 1999.

21. Stanley, RP. Personal communication. 2009.

22. van der Vaart, AW. Asymptotic Statistics. Cambridge University Press; 1998.

23. Whittaker, ET.; Watson, GN. A Course of Modern Analysis. fourth edition. Cambridge University Press; 1996.