# A DISCRETE MULTIVARIATE DISTRIBUTION RESULTING FROM THE LAW OF SMALL NUMBERS

NOBUAKI HOSHINO,* *Kanazawa University*

## Abstract

In the present article we derive a new discrete multivariate distribution using a limiting argument that is essentially the same as the law of small numbers. The distribution derived belongs to an exponential family, and randomly partitions positive integers. The facts shown about the distribution are useful in many fields of application involved with count data. The derivation parallels that of the Ewens distribution from the gamma distribution, and the new distribution is produced from the inverse Gaussian distribution. The method employed is regarded as the discretization of an infinitely divisible distribution over nonnegative real numbers.

*Keywords:* Conditional inverse Gaussian–Poisson distribution; extended negative binomial model; frequencies of frequencies; infinite divisibility; random clustering; species abundance

2000 Mathematics Subject Classification: Primary 62E10; 60E07
Secondary 62P99

## 1. Introduction

In the following, $\mathbb{R}_+$ denotes the set of nonnegative real numbers, and $\mathbb{N}_0$ and $\mathbb{N}$ are respectively the sets of nonnegative integers and positive integers.

Let us denote the set of all unordered partitions of a positive integer $n$ by

$$\mathscr{S}_n := \left\{ \boldsymbol{s}_n := (s_1, s_2, \ldots, s_n) : s_i \in \mathbb{N}_0, \ i = 1, 2, \ldots, n, \ \sum_{i=1}^{n} i s_i = n \right\}.$$

For a random vector $\boldsymbol{S}_n := (S_1, S_2, \ldots, S_n)$, the Ewens distribution is defined for $\theta > 0$ as

$$P(\boldsymbol{S}_n = \boldsymbol{s}_n) = \frac{\theta^u}{\theta^{[n]}} \frac{n!}{\prod_{i=1}^{n} i^{s_i} s_i!}, \qquad \boldsymbol{s}_n \in \mathscr{S}_n,$$

where $u = \sum_{i=1}^{n} s_i$ and $\theta^{[n]} = \theta(\theta + 1) \ldots (\theta + n - 1)$. This distribution constitutes random partitioning of a positive integer and has applications in many fields: genetics (Ewens (1972)), ecology (Aoki (2000)), linguistics (Sibuya (1991)), and statistical disclosure control (Hoshino (2001)), to name a few. See Chapter 41 of Johnson *et al.* (1997) for a review of this distribution.

A derivation of the Ewens distribution can be made in the following way. Suppose that random variables $F_1, F_2, \ldots, F_J$ are independently and identically distributed according to the negative binomial distribution. Let

$$S_i := \sum_{j=1}^{J} 1(F_j = i), \qquad i \in \mathbb{N}_0,$$

where $1(\cdot)$ denotes the indicator function, and write $\boldsymbol{S} := (S_1, S_2, \dots)$. Then $\boldsymbol{S}$ is multinomially distributed and $P(\boldsymbol{S} = \boldsymbol{s})$ is defined over

$$\mathscr{S}_\infty(J) := \left\{ \boldsymbol{s} := (s_1, s_2, \dots) : s_i \in \mathbb{N}_0, \ i = 1, 2, \dots, \ \sum_{i=1}^\infty s_i \le J \right\}.$$

The total sum $F_1 + F_2 + \cdots + F_J$ is denoted by

$$N := \sum_{i=1}^\infty i \, S_i. \tag{1}$$

Then the conditional distribution of $\boldsymbol{S}$ given $N$ becomes the negative multivariate hypergeometric distribution or Dirichlet–multinomial mixture. Given that $N = n$, we have $S_i = 0$ for $i > n$. Hence, we consider $P(\boldsymbol{S}_n = \boldsymbol{s}_n \mid N = n)$, which is defined over

$$\mathscr{S}_n(J) := \left\{ \boldsymbol{s}_n : s_i \in \mathbb{N}_0, \ i = 1, 2, \dots, n, \ \sum_{i=1}^n i s_i = n, \ \sum_{i=1}^n s_i \le J \right\}.$$

Taking $J \to \infty$, while the distribution of $N$ remains unchanged, the limiting distribution of the Dirichlet–multinomial distribution is the Ewens distribution.

Hoshino and Takemura (1998) pointed out that in this derivation the order of the conditioning on $N$ and the taking of the limit as $J \to \infty$ is exchangeable. In the same way, we start from independent, identically distributed $F_j$. However, first taking $J \to \infty$, while the distribution of $N$ remains unchanged, the limiting distribution of $\boldsymbol{S}$ is Anscombe's (1950) logarithmic series model, in which each $S_i$ is independently Poisson distributed. Then, conditioning $\boldsymbol{S}$ on $N$, we obtain the Ewens distribution. See Figure 1.

The above construction of the Ewens distribution can be regarded as a discretization of the gamma distribution, because the negative binomial distribution is the Poisson distribution mixed with the gamma distribution. Actually, any infinitely divisible distribution over $\mathbb{R}_+$ can be discretized in the same way.

Let us denote the Poisson distribution with mean $\lambda$ by $Po(\lambda)$. If $\lambda$ is distributed according to a distribution over $\mathbb{R}_+$, we say that $\lambda$ is mixed with this distribution, and $E(Po(\lambda))$ is a mixed Poisson distribution. Mixing $\lambda$ of $Po(\lambda)$ with an infinitely divisible distribution over $\mathbb{R}_+$ results in an infinitely divisible distribution over $\mathbb{N}_0$ (see Steutel and van Harn (2004, p. 368)); any infinitely divisible distribution over $\mathbb{N}_0$ produces a random partitioning distribution by the method based on Hoshino (2005). Summarizing the argument, we have the following discretization.

**Construction 1.** By mixing $\lambda$ of $Po(\lambda)$ with an infinitely divisible distribution over $\mathbb{R}_+$ we obtain an infinitely divisible distribution, say $F$, over $\mathbb{N}_0$. Let $F_1, F_2, \dots, F_J$ be independent and identically distributed with distribution $F$. Then the conditional distribution of $\boldsymbol{S}$ given $N = n$ converges in distribution to a random partitioning distribution of $n$ as $J \to \infty$, while the distribution of $N$ remains unchanged.

It might be noteworthy that the order of the conditioning and the taking of the limit is exchangeable in Construction 1. For an infinitely divisible distribution $Q$ over $\mathbb{R}_+$, let us denote a random partitioning distribution generated in this way by $\mathscr{D}(Q)$. Denoting the gamma distribution by Ga, the Ewens distribution is expressed as $\mathscr{D}(Ga)$.
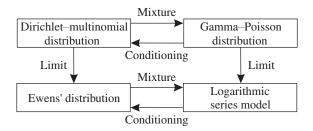
FIGURE 1: Relationships among distributions derived from the gamma–Poisson mixture.

The present article introduces $\mathcal{D}(\mathrm{IG})$, where IG denotes the inverse Gaussian distribution, which is infinitely divisible over $\mathbb{R}_+$. The random partitioning distribution $\mathcal{D}(\mathrm{IG})$ seems new and tractable. In particular, it belongs to an exponential family. Hence, it affords us a practical tool to describe count data in the aforementioned fields.

The organization of the present article is as follows. In Section 2 we provide notation and definitions of relevant distributions. In Section 3 we derive $\mathcal{D}(\mathrm{IG})$ and clarify its properties. The effect of the limiting argument will turn out to be the same as that of the law of small numbers. In Section 4 we discuss the parameter estimation of $\mathcal{D}(\mathrm{IG})$. In Section 5 we apply $\mathcal{D}(\mathrm{IG})$ to a typical data set to exemplify the usefulness of the distribution.

## 2. Preliminary results

In the statistical literature, $(S_0, S_1, \dots)$ are called size indices (Sibuya (1993)) or frequencies of frequencies (Good (1953)). In the following we denote the number of positive $F_j$ by

$$U := \sum_{i=1}^{\infty} S_i. \qquad (2)$$

This variate is of practical importance. For example, in ecology each $F_j$ represents the number of individuals of a species and $U$ corresponds to the total number of observable species. See Bunge and Fitzpatrick (1993) for a survey on the estimation of $U$. When $N = n$ is given, we write, in particular, $U_n := \sum_{i=1}^{n} S_i$, since $S_i$ has to be 0 for $i \geq n + 1$.

Recall that $u := \sum_{i=1}^{\infty} s_i$. Using this, define $s_0 := J - u$.

From now on, we follow Construction 1 for the IG distribution. Mixing $\lambda$ of $\mathrm{Po}(\lambda)$ with the IG distribution, we have the inverse Gaussian–Poisson (IGP) distribution, which is well reviewed in Chapter 7.1 of Seshadri (1999).

Assume that random variables $F_j$, $j = 1, \dots, J$, are independent and identically distributed with an IGP distribution. Then, for $\theta$, $0 < \theta \leq 1$, and $\alpha > 0$,

$$\mathrm{P}(\boldsymbol{S} = \boldsymbol{s}) = J! \prod_{i=0}^{\infty} \left\{ \sqrt{\frac{2\alpha}{\pi}} \exp(\alpha\sqrt{1-\theta}) \frac{(\alpha\theta/2)^i}{i!} K_{i-1/2}(\alpha) \right\}^{s_i} \frac{1}{s_i!}, \qquad \boldsymbol{s} \in \mathcal{S}_{\infty}(J), \quad (3)$$

where $K_{i-1/2}(\cdot)$ is the modified Bessel function of the third kind, of order $i - \frac{1}{2}$.

Equations (4) to (6) are taken directly from Watson (1944, pp. 79–80) for convenience. The argument of the modified Bessel function is always real and positive in the present article. We

have

$$K_{-1/2}(\alpha) = K_{1/2}(\alpha) = \sqrt{\frac{\pi}{2}} \alpha^{-1/2} \exp(-\alpha) \tag{4}$$

and

$$K_{y-1/2}(\alpha) = \sqrt{\frac{\pi}{2\alpha}} \exp(-\alpha) \sum_{i=0}^{y-1} \frac{(y-1+i)!}{(y-1-i)!\,i!} (2\alpha)^{-i}, \qquad y \in \mathbb{N}, \tag{5}$$

which can be shown by the recurrence formula

$$K_{\gamma+1}(\alpha) = \frac{2\gamma}{\alpha} K_{\gamma}(\alpha) + K_{\gamma-1}(\alpha). \tag{6}$$

Equation (5) can be expressed in another way, using the C-number $C(n, l, z)$ that is defined for any real (or, more generally, complex) numbers $z$ and $t$ by

$$(zt)^{(n)} = \sum_{l=0}^{n} C(n, l, z) t^{(l)}, \qquad n \in \mathbb{N}_0,$$

where $t^{(l)} = t(t-1)\cdots(t-l+1)$. See Charalambides and Singh (1988) for a comprehensive review of the C-number. As evaluated by, e.g. Yamato *et al.* (2001, p. 25),

$$C(y, l, \tfrac{1}{2}) = (-1)^{l-y} \frac{(2y-l-1)!}{(l-1)!\,(y-l)!} \left(\frac{1}{2}\right)^{2y-l}, \qquad l = 1, 2, \ldots, y,$$

and, thus,

$$K_{y-1/2}(\alpha) = \sqrt{\frac{\pi}{2\alpha}} \exp(-\alpha) \sum_{l=1}^{y} C(y, l, \tfrac{1}{2}) 2^y \left(\frac{-1}{\alpha}\right)^{y-l}, \qquad y \in \mathbb{N}. \tag{7}$$

When the order of the function is large, Ismail (1977) showed an asymptotic formula using which the computations become easy:

$$K_{\gamma}(\alpha) \sim 2^{\gamma} \gamma^{\gamma} \exp(-\gamma) \alpha^{-\gamma} \sqrt{\frac{\pi}{2\gamma}} \qquad \text{as } \gamma \to \infty. \tag{8}$$

As noted in Hoshino (2003), under (3) $N$ is distributed as

$$P(N = n) = \sqrt{\frac{2J\alpha}{\pi}} \exp(J\alpha\sqrt{1-\theta}) \frac{(J\alpha\theta/2)^n}{n!} K_{n-1/2}(J\alpha), \qquad n \in \mathbb{N}_0. \tag{9}$$

Hence, the conditional distribution of $S$ given $N = n$ is, for $\alpha > 0$,

$$P(S_n = s_n \mid N = n)$$
$$= \left(\frac{2\alpha}{\pi}\right)^{(J-1)/2} \frac{J!\,n!}{J^{n+1/2} K_{n-1/2}(J\alpha)} \prod_{i=0}^{n} \left\{ \frac{K_{i-1/2}(\alpha)}{i!} \right\}^{s_i} \frac{1}{s_i!}, \qquad s_n \in \mathcal{S}_n(J), \tag{10}$$

which is the conditional inverse Gaussian–Poisson (CIGP) distribution (Hoshino (2003)).

From (9) we note that the distribution of $N$ remains unchanged as

$$J \to \infty \quad \text{and} \quad \alpha \to 0 \qquad \text{such that } J\alpha = \mu > 0. \tag{11}$$

This property is the key for the exchangeability between the conditioning and the taking of the limit in the derivation of $\mathcal{D}(\text{IG})$ in the next section.

## 3. Main results

This section derives $\mathcal{D}(\mathrm{IG})$ and investigates its properties. The proofs of the theorems in this section are all provided in Appendix A.

**Theorem 1.** *By applying (11), the CIGP distribution (10) converges in distribution to*

$$\mathrm{P}(S_n = s_n \mid N = n) = \sqrt{\frac{\pi}{2\mu}} \frac{n! \exp(-\mu)}{\mu^{n-u} K_{n-1/2}(\mu)} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!}, \qquad s_n \in \mathcal{S}_n, \quad (12)$$

*where* $(-1)!! = 1$ *and* $(2i-3)!! = (2i-3)(2i-5)\cdots 1$.

The distribution (12) is referred to as $\mathcal{D}(\mathrm{IG})$ or the limiting CIGP (LCIGP) distribution. To understand the meaning of the limiting argument (11), the following result is useful.

**Theorem 2.** *By applying (11), the IGP model (3) converges in distribution to*

$$\mathrm{P}(S = s) = \prod_{i=1}^{\infty} \frac{\exp(-\tau(i;\mu,\theta))\tau(i;\mu,\theta)^{s_i}}{s_i!}, \qquad s \in \mathbb{N}_0^{\infty}, \quad (13)$$

*where* $0 < \theta \le 1$ *and*

$$\tau(i;\mu,\theta) = \mu\left(\frac{\theta}{2}\right)^i \frac{(2i-3)!!}{i!} = \frac{\mu\theta^i}{2\sqrt{\pi}} \frac{\Gamma(i-1/2)}{\Gamma(i+1)}, \qquad i \in \mathbb{N}.$$

*That is, each* $S_i$ *is independently distributed according to* $\mathrm{Po}(\tau(i;\mu,\theta))$, $i = 1, 2, \ldots$.

The conditional distribution of (13) given $N = n$ is (12) and, conversely, the mixture of (12) with the distribution of $N$, (9), results in (13).

Theorem 2 indicates that the effect of (11) is essentially that of the law of small numbers. In (3), $S$ is multinomially distributed and each marginal $S_i$, except for $S_0$, becomes independent and Poisson distributed as $J \to \infty$. In this limit, almost every $F_j$ is 0 and the rare event of having a positive $F_j$ is described by the limiting distribution.

In (13), $\tau(i;\mu,\theta)$ is proportional to the following special case of the extended truncated negative binomial distribution (Engen (1974)):

$$\mathrm{P}(X = i) = \frac{1}{1 - \sqrt{1-\theta}} \frac{\theta^i (2i-3)!!}{2^i i!}, \qquad i \in \mathbb{N},\ 0 < \theta \le 1. \quad (14)$$

This is so because the IGP distribution is the compound Poisson distribution defined by (14). A distribution over $\mathbb{N}$ is used to define the compound Poisson distribution, and the class of compound Poisson distributions coincides with the class of infinitely divisible distributions over $\mathbb{N}_0$; see Steutel and van Harn (2004, Theorem 3.2). Therefore, an infinitely divisible distribution over $\mathbb{N}_0$ can be expressed as a compound Poisson distribution whose defining distribution over $\mathbb{N}$ determines the limiting distribution of $S$ in Construction 1; see Hoshino (2005, Theorem 2.1).

Construction 1 for the IG distribution has been completed by Theorem 1 and Theorem 2. The relationships derived are illustrated in Figure 2, which parallels Figure 1. Because Hoshino's (2005) discussion of Engen's extended negative binomial model deals with (13) as a special case, our attention henceforth centers upon (12), that is, $\mathcal{D}(\mathrm{IG})$.

To begin with, we point out that $\mathcal{D}(\mathrm{IG})$ does not belong to an existing class of random partitioning distributions. Pitman (2003) discussed the discretization of an infinitely divisible
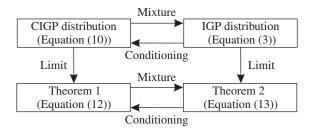
FIGURE 2: Relationships among distributions derived from the IGP mixture.

distribution over $\mathbb{R}_+$. This method produces a random partitioning distribution called the exchangeable partition probability function (EPPF). The Ewens distribution is the EPPF generated from the gamma distribution. However, the EPPF generated from the inverse Gaussian distribution differs from $\mathcal{D}(\text{IG})$. Hence, Construction 1 is a different discretization than Pitman's.

To confirm the difference, we note the partition structure (Kingman (1978)), where a given partition of $n$ elements results from the deletion of one element uniformly at random from a partition of $n + 1$ elements. More precisely, a distribution that has the partition structure satisfies, for $n = 1, 2, \ldots,$

$$
\begin{aligned}
&\mathrm{P}(S_1 = s_1, \ S_2 = s_2, \ \ldots \ | \ N = n) \\
&\quad = \mathrm{P}(S_1 = s_1 + 1, \ S_2 = s_2, \ \ldots \ | \ N = n + 1)\frac{s_1 + 1}{n + 1} \\
&\qquad + \sum_{r=2}^{n+1} \mathrm{P}(S_1 = s_1, \ \ldots, \ S_{r-1} = s_{r-1} - 1, \ S_r = s_r + 1, \ \ldots \ | \ N = n + 1)\frac{r(s_r + 1)}{n + 1}.
\end{aligned}
$$

Every EPPF has this partition structure, and the Ewens distribution is an example of this type. However, $\mathcal{D}(\text{IG})$ contradicts this rule: for example, it is easy to verify that

$$
\begin{aligned}
&\mathrm{P}(S_1 = 2, \ S_2 = 0 \ | \ N = 2) \\
&\quad \neq \mathrm{P}(S_1 = 3, \ S_2 = 0, \ S_3 = 0 \ | \ N = 3) + \tfrac{2}{3}\mathrm{P}(S_1 = 1, \ S_2 = 1, \ S_3 = 0 \ | \ N = 3)
\end{aligned}
$$

under (12). Hence, Remark 1 is justified, and $\mathcal{D}(\text{IG})$ is not the EPPF generated from the IG distribution. The conditional distribution $\mathrm{P}(S_n = s_n \ | \ N = n, \ U_n = u)$ of (12) does coincide with that of an EPPF, however.

**Remark 1.** The LCIGP distribution does not possess Kingman's partition structure.

Although we have shown the difference in the partition structure, $\mathcal{D}(\text{IG})$ is an analogue of the Ewens distribution. Hence, it is natural to expect that the LCIGP distribution has properties similar to those of the Ewens distribution. The following investigation will first show results closely similar to Proposition 2.2 of Sibuya (1993), who considered the same limit of the Ewens distribution as in Theorem 3.

**Theorem 3.** *Let m be a finite, fixed positive integer. Suppose that $S_n$ is distributed as in (12). Then, as $n \to \infty$, the limiting distribution of the first m components $(S_1, S_2, \ldots, S_m)$ is the joint distribution of independent* $\mathrm{Po}(\tau(i; \mu, 1))$, $i = 1, 2, \ldots, m$.

In the limit as $n \to \infty$, the difference between (12) and (13) is diminished in the lower tail (i.e. $S_1, S_2, \ldots, S_m$). Moreover, $\theta$ of (13) has to be unity in the limit, because a heavy upper tail is required to describe an infinite number of individuals.

Next, using (5) we rewrite (12) as

$$P(S_n = s_n \mid N = n)$$
$$= \exp\left( u \log \mu - \log\left( \mu^n \sum_{i=0}^{n-1} \frac{(n-1+i)!}{(n-1-i)!\, i!} (2\mu)^{-i} \right) \right) n! \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!},$$

which illustrates the following fact.

**Remark 2.** The LCIGP distribution belongs to an exponential family, and $U_n$ is its sufficient statistic.

The Ewens distribution also belongs to an exponential family and has $U_n$ as its sufficient statistic; see Sibuya (1991). It should be remembered that $U_n$ is an important variate in applications. Therefore, a natural goal is to find the distribution of $U_n$.

**Theorem 4.** *Suppose that the size indices are distributed as in (12). Then*

$$P(U_n = v) = \sqrt{\frac{\pi}{2\mu}} \frac{\exp(-\mu)}{K_{n-1/2}(\mu)} \left(\frac{1}{2\mu}\right)^{n-v} \frac{(2n-v-1)!}{(v-1)!\,(n-v)!} \tag{15}$$

$$= \frac{1}{\sum_{l=1}^{n} (-1/\mu)^{n-l} C(n,l,\frac{1}{2})} \left(\frac{-1}{\mu}\right)^{n-v} C(n,v,\tfrac{1}{2}), \qquad v = 1,2,\ldots,n. \tag{16}$$

Let us regard $U_n$ as the number of urns and $n$ as the number of balls. Considering the process of increasing the number of balls, Sibuya (1993) noted that the Ewens distribution has the following urn model implication:

$$P_{\text{Ewens}}(U_{n+1} = v) = P_{\text{Ewens}}(U_{n+1} = v \mid U_n = v) P_{\text{Ewens}}(U_n = v)$$
$$+ P_{\text{Ewens}}(U_{n+1} = v \mid U_n = v-1) P_{\text{Ewens}}(U_n = v-1).$$

Here the number of urns is increased by adding a new urn that contains one ball. If a ball is put into an existing urn then the number of urns is unchanged. As regards the LCIGP distribution, (6) leads to the recurrence formula:

$$P(U_{n+1} = v) = \frac{K_{n-1/2}(\mu)}{K_{n+1-1/2}(\mu)} \frac{2n-1}{\mu} P(U_n = v) + \frac{K_{n-1-1/2}(\mu)}{K_{n+1-1/2}(\mu)} P(U_{n-1} = v-2),$$

which can be interpreted as

$$P(U_{n+1} = v) = P(U_{n+1} = v \mid U_n = v) P(U_n = v)$$
$$+ P(U_{n+1} = v \mid U_{n-1} = v-2) P(U_{n-1} = v-2). \tag{17}$$

Equation (17) implies that the number of urns is increased by adding new two urns and two balls. If one new ball added, it is put into an existing urn. Hence, it seems that the LCIGP distribution is not constructed by a sequence of one-by-one addition of balls. This view is consistent with Remark 1.

The moments of $U_n$ are functions of the moments of the size indices, since $U_n = \sum_{i=1}^{n} S_i$. The joint factorial moments of size indices, which are vital in practice, are given in the next theorem.

TABLE 1: Expectations of size indices, for $N = 1000$ (see Theorem 5).

| $\mu$ | $E(S_1)$ | $E(S_2)$ | $E(S_3)$ | $E(S_4)$ | $E(S_5)$ |
|---|---|---|---|---|---|
| 0.1 | 0.05 | 0.01 | 0.01 | 0.00 | 0.00 |
| 1.0 | 0.50 | 0.13 | 0.06 | 0.04 | 0.03 |
| 10.0 | 5.01 | 1.25 | 0.63 | 0.39 | 0.28 |
| 100.0 | 49.95 | 12.47 | 6.23 | 3.89 | 2.72 |
| 300.0 | 146.98 | 85.93 | 17.64 | 10.80 | 7.41 |
| 500.0 | 236.37 | 55.87 | 26.41 | 15.60 | 10.32 |
| 700.0 | 315.59 | 71.13 | 32.05 | 18.05 | 11.39 |
| 900.0 | 384.13 | 81.95 | 34.96 | 18.64 | 11.12 |
| 10 000.0 | 905.08 | 40.92 | 3.70 | 0.42 | 0.05 |

**Theorem 5.** *Suppose that the size indices are distributed as in (12). Then, for all $r_1, \ldots, r_n \in \mathbb{N}_0$ such that $R := \sum_{i=1}^{n} i r_i \leq n$, the factorial moments are*

$$E\left( \prod_{i=1}^{n} S_i^{(r_i)} \right) = \frac{K_{n-R-1/2}(\mu)\mu^{r-R} n!}{K_{n-1/2}(\mu)(n-R)!} \prod_{i=1}^{n} \left( \frac{(2i-3)!!}{i!} \right)^{r_i}, \tag{18}$$

*where $r = \sum_{i=1}^{n} r_i$.*

In particular,

$$E(S_i) = \frac{K_{n-i-1/2}(\mu)\mu^{1-i}(2i-3)!!\, n!}{K_{n-1/2}(\mu) i!\, (n-i)!}, \qquad i = 1, 2, \ldots, n. \tag{19}$$

Table 1 summarizes values of $E(S_i)$ for $i = 1, 2, \ldots, 5$ and various parameter values, given that $N = 1000$. Because $E(U_n) = \sum_{i=1}^{n} E(S_i)$, the following proposition follows from (19). Higher moments of $U_n$ can be obtained in an analogous way.

**Proposition 1.** *Suppose that a random variable $U_n$ is distributed as in (15). Then its expectation is*

$$E(U_n) = \frac{n!}{K_{n-1/2}(\mu)} \sum_{i=1}^{n} \frac{K_{n-i-1/2}(\mu)\mu^{1-i}(2i-3)!!}{i!\, (n-i)!}.$$

The moments shown become simple as $n \to \infty$, because the limiting distribution of a size index is Poisson, by Theorem 3. Also, the limiting distribution of $U_n$ can be shown to be shifted Poisson, as follows.

**Theorem 6.** *Suppose that $U_n$ is subject to (15). Then, as $n \to \infty$, $U_n$ converges in distribution to $1 + X$, where $X$ is Poisson distributed with mean $\mu$.*

## 4. Parameter estimation

This section deals with the parameter estimation of the LCIGP distribution. After constructing the maximum likelihood estimation, we consider an approximate moment estimator.

The loglikelihood of (12) is denoted by

$$L(\mu) = -\mu - (n - u + \tfrac{1}{2}) \log \mu - \log K_{n-1/2}(\mu) + \text{const.}$$

In the following we will use the notation $R_\gamma(\alpha) = K_{\gamma+1}(\alpha)/K_\gamma(\alpha)$. As given by Seshadri (1999, p. 125) for instance, we have

$$\frac{\partial \log K_\gamma(\alpha)}{\partial \alpha} = -R_\gamma(\alpha) + \frac{\gamma}{\alpha}.$$

The first derivative of $L$,

$$\frac{dL}{d\mu} = -1 - (2n-u)\frac{1}{\mu} + R_{n-1/2}(\mu),$$

is hence easy to calculate. The maximum likelihood estimator is the unique solution to $dL/d\mu = 0$; Remark 1 justifies the claim of uniqueness, based on the property of the exponential family discussed, for example, by Lehmann (1991, p. 417).

A numerical method is required to solve the likelihood equation. The second derivative,

$$\frac{d^2 L}{d\mu^2} = R_{n-1/2}^2(\mu) - \frac{2n}{\mu}R_{n-1/2}(\mu) + (2n-u)\frac{1}{\mu^2} - 1,$$

enables us to employ the fast convergent iteration of the Newton–Raphson method.

A moment estimator is also useful, to obtain the starting value of the iteration procedure, but an exact one is inconvenient to compute because of the modified Bessel function. Therefore, we use an approximate estimator. Under distribution (13), $E(U) = \mu(1 - \sqrt{1-\theta})$ and $\theta = 4E(S_2)/E(S_1)$. By substituting the latter equation into the former, we have

$$\mu = \frac{E(U)}{1 - \sqrt{1 - 4E(S_2)/E(S_1)}},$$

which leads to the following approximate moment estimator:

$$\tilde{\mu} = \frac{U}{1 - \sqrt{1 - 4S_2/S_1}}. \tag{20}$$

However, $1 - 4s_2/s_1$ can be negative. In such a situation, it may be possible to use $u$ as an estimate of $\mu$.

Asymptotically, as $n \to \infty$, the parameter estimation becomes easy according to Theorem 6: the estimate of $\mu$ should be $u - 1$ when $n$ is sufficiently large.

## 5. An application result

In this section we demonstrate the applicability of the LCIGP distribution by fitting it to the frequency data of South and Edelsten (1939). The same data were analyzed by Engen (1978, p. 109).

The data describe frequencies of genera with different numbers of species in British macro-lepidoptera, excluding butterflies. The total number of genera, $u$, is 357, and the total number of species, $n$, is 756. The parameter $\mu$ of the LCIGP distribution is estimated, using the maximum likelihood estimator, to be $\hat{\mu} = 514.982$ and, using the approximate estimator (20), to be $\tilde{\mu} = 578.306$. Under the maximum likelihood estimate, the expectations of size indices are tabulated in the fourth column of Table 2. The first column corresponds to '$i$' of $E(S_i)$, and the second column contains the observed size indices. For comparison, the Ewens distribution is also fitted using maximum likelihood estimation and summarized in the third column; see, e.g. Hoshino and Takemura (1998) for details of the fitting of the Ewens distribution.

In this example, the fit of the Ewens distribution is not satisfactory, but the LCIGP distribution is a reasonable fit, and thus broadens the scope of count data modeling.

TABLE 2: Frequencies of genera with different number of species.

| Species per genus | Observed size indices | Ewens' distribution | LCIGP distribution |
|---|---|---|---|
| 1 | 239 | 195.69 | 233.39 |
| 2 | 51 | 72.59 | 52.87 |
| 3 | 25 | 35.89 | 23.95 |
| 4 | 12 | 19.96 | 13.56 |
| 5 | 8 | 11.83 | 8.59 |
| 6 | 3 | 7.30 | 5.84 |
| 7 | 4 | 4.64 | 4.15 |
| 8 | 3 | 3.00 | 3.05 |
| 9 | 2 | 1.98 | 2.30 |
| 10+ | 10 | 4.13 | 9.30 |

## Appendix A.

Theorem 1 and 2 are actually special cases of Hoshino (2005, Propositions 2.2 and 2.1). However, analytical proofs are provided here to support the general argument in Hoshino (2005).

*Proof of Theorem 1.* In the following, the probability function (10) is shown to converge to (12). Let us rewrite the right-hand side of (10) as $C_1 \times C_2 \times C_3$, where

$$C_1 = \frac{J!}{(J-u)! \, J^u}, \qquad C_2 = \frac{n!}{K_{n-1/2}(\mu)} \prod_{i=1}^{n} \left(\frac{1}{i!}\right)^{s_i} \frac{1}{s_i!},$$

$$C_3 = \left(\frac{1}{J}\right)^{n-u+1/2} \left(\sqrt{\frac{\pi}{2\alpha}}\right)^{1-u} \exp(-\mu + \alpha u) \prod_{i=1}^{n} K_{i-1/2}(\alpha)^{s_i}.$$

Because $C_1 \to 1$ as $J \to \infty$, it suffices to show that $C_2 \times C_3$ converges to the desired limit.

As stated in Jørgensen (1982, p. 171), for $\gamma > 0$ we have

$$K_\gamma(\alpha) \sim \Gamma(\gamma) 2^{\gamma-1} \alpha^{-\gamma} \tag{21}$$

as $\alpha \to 0$. Using this result, as $\alpha \to 0$ we have

$$\prod_{i=1}^{n} K_{i-1/2}(\alpha)^{s_i} \sim 2^{n-3u/2} \alpha^{u/2-n} \prod_{i=1}^{n} \Gamma\left(i - \frac{1}{2}\right)^{s_i}$$

$$= 2^{n-3u/2} \alpha^{u/2-n} \prod_{i=1}^{n} (2^{1-i} \sqrt{\pi} (2i-3)!!)^{s_i}$$

$$= 2^{-u/2} (\sqrt{\pi})^u \alpha^{u/2-n} \prod_{i=1}^{n} \{(2i-3)!!\}^{s_i}$$

by (1) and (2). Therefore, by (11),

$$C_3 \sim \sqrt{\frac{\pi}{2}}\left(\frac{1}{J}\right)^{n-u+1/2} \alpha^{u-n-1/2} \exp(-\mu + \alpha u) \prod_{i=1}^{n}\{(2i-3)!!\}^{s_i}$$

$$\rightarrow \sqrt{\frac{\pi}{2}}\left(\frac{1}{\mu}\right)^{n-u+1/2} \exp(-\mu) \prod_{i=1}^{n}\{(2i-3)!!\}^{s_i}.$$

Multiplying this limit by $C_2$ completes the proof.

*Proof of Theorem 2.* We will obtain the required convergence in distribution by showing that the probability generating function (PGF) of size indices converges to the PGF of the joint distribution of independent Poisson variables.

First we derive the PGF, $G(z_1, z_2, \ldots, z_l)$, of the joint distribution of $(S_1, S_2, \ldots, S_l)$ under (3). Let us denote the PGF for $J = 1$ by

$$G_1(z_1, z_2, \ldots, z_l) = \mathrm{E}\left(\prod_{i=1}^{l} z_i^{S_i}\right)$$

$$= \sum_{i=1}^{l}(z_i - 1)\,\mathrm{P}(F_1 = i) + 1$$

$$= \sum_{i=1}^{l}(z_i - 1)\sqrt{\frac{2\alpha}{\pi}}\exp(\alpha\sqrt{1-\alpha})\frac{(\alpha\theta/2)^i}{i!}K_{i-1/2}(\alpha) + 1.$$

By the independence of the $F_j$, the PGF for a general $J$ is expressed as $G(z_1, z_2, \ldots, z_l) = G_1(z_1, z_2, \ldots, z_l)^J$. Now we consider the limit as $J \to \infty$, for $J\alpha = \mu$. Using (21) we obtain

$$\left[\sum_{i=1}^{l}(z_i - 1)\sqrt{\frac{2\alpha}{\pi}}\exp(\alpha\sqrt{1-\alpha})\frac{(\alpha\theta/2)^i}{i!}K_{i-1/2}(\alpha) + 1\right]^J$$

$$\rightarrow \left[1 + \frac{1}{J}J\sum_{i=1}^{l}(z_i - 1)\sqrt{\frac{2\alpha}{\pi}}\frac{(\alpha\theta/2)^i}{i!}\Gamma(i-\tfrac{1}{2})2^{i-3/2}\alpha^{-i+1/2}\right]^J$$

$$\rightarrow \exp\left(\sum_{i=1}^{l}(z_i - 1)\mu\left(\frac{\theta}{2}\right)^i\frac{(2i-3)!!}{i!}\right)$$

$$= \exp\left(\sum_{i=1}^{l}(z_i - 1)\tau(i; \mu, \theta)\right). \tag{22}$$

Equation (22) coincides with the PGF of the joint distribution of independent Poisson variables $S_i$, $i = 1, 2, \ldots, l$, with mean $\mathrm{E}(S_i) = \tau(i; \mu, \theta)$. This argument holds for any $l \in \mathbb{N}$, and the sequence of the joint distribution of $(S_1, S_2, \ldots, S_l)$ determines the limiting distribution of $\boldsymbol{S}$ as $l \to \infty$.

The distribution of $N$ under (13) is that of the IGP model. Hence, the conditional distribution of (13) given $N$ is the result of dividing (13) by (9), which equals (12) using (1) and (2). Conversely, (12) multiplied by (9) is (13).

*Proof of Theorem 3.* Following Sibuya (1993), we adopt the method of moments to show the required convergence in distribution; this proof depends on the joint factorial moments in (18), which equation will be proved below. Assuming that (18) is correct, the components have the joint factorial moments

$$M(r_1, r_2, \ldots, r_m) = \frac{K_{n-R-1/2}(\mu)\mu^{r-R}n!}{K_{n-1/2}(\mu)(N-R)!} \prod_{i=1}^{m} \left( \frac{(2i-3)!!}{i!} \right)^{r_i}.$$

From (8), as $n \to \infty$ we have

$$M(r_1, r_2, \ldots, r_m) \to 2^{-R} \left( 1 - \frac{R}{n-1/2} \right)^{n-R-1/2} \frac{1}{(n-1/2)^R} \exp(R)$$

$$\times \sqrt{\frac{2n-1}{2n-R-1}} \mu^r \frac{n!}{(n-R)!} \prod_{i=1}^{m} \left( \frac{(2i-3)!!}{i!} \right)^{r_i}$$

$$\to \prod_{i=1}^{m} \{\tau(i; \mu, 1)\}^{r_i}$$

$$< \infty.$$

The $r$th factorial moment of $\mathrm{Po}(\tau(i; \mu, 1))$ is $\tau(i; \mu, 1)^r$. Hence, it is observable that $M(r_1, r_2, \ldots, r_m)$ converges to the joint factorial moments of $\mathrm{Po}(\tau(i; \mu, 1))$, $i = 1, 2, \ldots, m$, for any $(r_1, r_2, \ldots, r_m)$. Consequently, we have proved the theorem.

*Proof of Theorem 4.* From definition (12),

$$K_{n-1/2}(\mu)\,\mathrm{P}(U_n = v) = \sum_{s_1+\cdots+s_n=v} \sqrt{\frac{\pi}{2\mu}} \frac{n! \exp(-\mu)}{\mu^{n-v}} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!}.$$

Thus, using (5), we obtain

$$\left( \sum_{i=0}^{n-1} \frac{(n-1+i)!}{(n-1-i)!\,i!} (2\mu)^{-i} \right) \mathrm{P}(U_n = v) = \sum_{s_1+\cdots+s_n=v} \frac{n!}{\mu^{n-v}} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!}.$$

Because $\sum_{v=1}^{n} \mathrm{P}(U_n = v) = 1$, we have

$$\sum_{i=0}^{n-1} \frac{(n-1+i)!}{(n-1-i)!\,i!} (2\mu)^{-i} = \sum_{v=1}^{n} \sum_{s_1+\cdots+s_n=v} \frac{n!}{\mu^{n-v}} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!}.$$

By comparing the coefficients of $\mu$ on the left- and right-hand sides of this equation, we have

$$\sum_{s_1+\cdots+s_n=v} n! \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!} = \frac{(n-1+n-v)!}{(n-1-n+v)!\,(n-v)!} 2^{v-n},$$

which leads to (15).

The above proof depends on the fact that the sum of probabilities is unity, but Professor H. Yamato has suggested a direct derivation of $\mathrm{P}(U_n)$. The following proof implies that (12) has the same probability of partitioning $n$ given $U_n$ as does a special case of Pitman's (1995)

sampling formula; see Hoshino (2005) for further discussion. Equation (7) results in another expression for (12):

$$P(S_n = s_n \mid N = n) = \frac{n! \, \mu^{u-n}}{\sum_{l=1}^{n} C(n, l, \frac{1}{2}) 2^n (-1/\mu)^{n-l}} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!}.$$

Because

$$(2i-3)!! = 2^i (-1)^{i-1} \binom{\frac{1}{2}}{i} i!,$$

the distribution can be further rewritten as

$$P(S_n = s_n \mid N = n) = \frac{n!}{\sum_{l=1}^{n} C(n, l, \frac{1}{2})(-1/\mu)^{n-l}} \left( \frac{-1}{\mu} \right)^{n-u} \prod_{i=1}^{n} \left\{ \binom{\frac{1}{2}}{i} \right\}^{s_i} \frac{1}{s_i!}.$$

Charalambides and Singh (1988, Equation 3.24) noted that

$$\sum_{s_1 + s_2 + \cdots + s_n = u} n! \prod_{i=1}^{n} \binom{\frac{1}{2}}{i}^{s_i} \frac{1}{s_i!} = C(n, u, \tfrac{1}{2}).$$

Substituting this into the previous equation leads to (16).

*Proof of Theorem 5.* We can prove (18) using the fact that $\sum_{s_n \in \mathcal{S}_n} P(S_n = s_n) = 1$, as follows:

$$
\begin{aligned}
E\left( \prod_{i=1}^{N} S_i^{(r_i)} \right) &= \sum_{s_n \in \mathcal{S}_n} \sqrt{\frac{\pi}{2\mu}} \frac{n! \exp(-\mu)}{\mu^{n-u} K_{n-1/2}(\mu)} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i - r_i} \frac{1}{(s_i - r_i)!} \left\{ \frac{(2i-3)!!}{i!} \right\}^{r_i} \\
&= \frac{n! \, \mu^{r-R} K_{n-R-1/2}(\mu)}{(n-R)! \, K_{n-1/2}(\mu)} \sum_{s_n \in \mathcal{S}_n} \sqrt{\frac{\pi}{2\mu}} \frac{(n-R)! \exp(-\mu)}{\mu^{n-R-(u-r)} K_{n-R-1/2}(\mu)} \\
&\qquad \times \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i - r_i} \frac{1}{(s_i - r_i)!} \left\{ \frac{(2i-3)!!}{i!} \right\}^{r_i} \\
&= \frac{n! \, \mu^{r-R} K_{n-R-1/2}(\mu)}{(n-R)! \, K_{n-1/2}(\mu)} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{r_i} \\
&\qquad \times \sum_{s_n - R \in \mathcal{S}_{n-R}} \sqrt{\frac{\pi}{2\mu}} \frac{(n-R)! \exp(-\mu)}{\mu^{n-R-u} K_{n-R-1/2}(\mu)} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{s_i} \frac{1}{s_i!} \\
&= \frac{n! \, \mu^{r-R} K_{n-R-1/2}(\mu)}{(n-R)! \, K_{n-1/2}(\mu)} \prod_{i=1}^{n} \left\{ \frac{(2i-3)!!}{i!} \right\}^{r_i}.
\end{aligned}
$$

**Remark 3.** Given the result in (19), the relation $\sum_{i=1}^{n} i \, E(S_i) = n$ is tantamount to the following recurrence formula:

$$K_{n-1/2}(\mu) = \sum_{i=1}^{n} K_{n-i-1/2}(\mu) \mu^{1-i} \frac{(2i-3)!! \, (n-1)!}{(i-1)! \, (n-i)!}.$$

*Proof of Theorem 6.* We show that the PGF of $U_n$ converges to that of $1 + \mathrm{Po}(\mu)$. Since the distribution of $U_n$ is given by (15), its PGF is evaluated as

$$
\begin{aligned}
\mathrm{E}(z^{U_n}) &= z^n \sqrt{\frac{\pi}{2\mu}} \frac{\exp(-\mu)}{K_{n-1/2}(\mu)} \sum_{v=1}^{n} \left(\frac{1}{2\mu z}\right)^{n-v} \frac{(2n-v-1)!}{(v-1)!\,(n-v)!} \\
&= z^n \sqrt{\frac{\pi}{2\mu}} \frac{\exp(-\mu)}{K_{n-1/2}(\mu)} \sqrt{\frac{2\mu z}{\pi}} \frac{K_{n-1/2}(\mu z)}{\exp(-\mu z)} \\
&= z^{n+1/2} \exp(\mu(z-1)) \frac{K_{n-1/2}(\mu z)}{K_{n-1/2}(\mu)},
\end{aligned}
$$

using the fact that $\sum_{v=1}^{n} \mathrm{P}(U_n = v) = 1$. Then, from (8), $K_{n-1/2}(\mu z)/K_{n-1/2}(\mu) \sim z^{-n+1/2}$ as $n \to \infty$. Hence the limit of the PGF is $z \exp(\mu(z-1))$.

## Acknowledgements

## References

ANSCOMBE, F. J. (1950). Sampling theory of the negative binomial and logarithmic series distributions. *Biometrika* **37,** 358–382.

AOKI, M. (2000). Cluster size distributions of economic agents of many types in a market. *J. Math. Anal. Appl.* **249,** 32–52.

BUNGE, J. AND FITZPATRICK, M. (1993). Estimating the number of species: a review. *J. Amer. Statist. Assoc.* **88,** 364–373.

CHARALAMBIDES, C. A. AND SINGH, J. (1988). A review of the Stirling numbers, their generalizations and statistical applications. *Commun. Statist. Theory Meth.* **17,** 2533–2595.

ENGEN, S. (1974). On species frequency models. *Biometrika* **61,** 263–270.

ENGEN, S. (1978). *Stochastic Abundance Models.* Chapman and Hall, London.

EWENS, W. J. (1972). The sampling theory of selectively neutral alleles. *Theoret. Pop. Biol.* **3,** 87–112.

GOOD, I. J. (1953). The population frequencies of species and the estimation of population parameters. *Biometrika* **40,** 237–264.

HOSHINO, N. (2001). Applying Pitman's sampling formula to microdata disclosure risk assessment. *J. Official Statist.* **17,** 499–520.

HOSHINO, N. (2003). Random clustering based on the conditional inverse Gaussian–Poisson distribution. *J. Japan Statist. Soc.* **33,** 105–117.

HOSHINO, N. (2005). Engen's extended negative binomial model revisited. *Ann. Inst. Statist. Math.* **57,** 369–387.

HOSHINO, N. AND TAKEMURA, A. (1998). Relationship between logarithmic series model and other superpopulation models useful for microdata disclosure risk assessment. *J. Japan Statist. Soc.* **28,** 125–134.

ISMAIL, M. E. H. (1977). Integral representations and complete monotonicity of various quotients of Bessel functions. *Canad. J. Math.* **29,** 1198–1207.

JOHNSON, N. L., KOTZ, S. AND BALAKRISHNAN, N. (1997). *Discrete Multivariate Distributions.* John Wiley, New York.

JØRGENSEN, B. (1982). *Statistical Properties of the Generalized Inverse Gaussian Distribution* (Lecture Notes Statist. **9**). Springer, New York.

KINGMAN, J. F. C. (1978). Random partitions in population genetics. *Proc. R. Soc. London A* **361,** 1–20.

LEHMANN, E. L. (1991). *Theory of Point Estimation.* Wadsworth and Brooks, Pacific Grove, CA.

PITMAN, J. (1995). Exchangeable and partially exchangeable random partitions. *Prob. Theory Relat. Fields* **102,** 145–158.

PITMAN, J. (2003). Poisson–Kingman partitions. In *Science and Statistics: A Festschrift for Terry Speed* (IMS Lecture Notes Monogr. Ser. **40**), Institute of Mathematical Statistics, Hayward, CA, pp. 1–34.

SESHADRI, V. (1999). *The Inverse Gaussian Distribution* (Lecture Notes Statist. **137**). Springer, New York.

SIBUYA, M. (1991). A cluster-number distribution and its application to the analysis of homonyms. *Japanese J. Appl. Statist.* **20,** 139–153 (in Japanese).

SIBUYA, M. (1993). A random clustering process. *Ann. Inst. Statist. Math.* **45,** 459–465.

SOUTH, R. AND EDELSTEN, H. M. (1939). *The Moths of the British Isles*. Frederick Warne, London.

STEUTEL, F. W. AND VAN HARN, K. (2004). *Infinite Divisibility of Probability Distributions on the Real Line*. Marcel Dekker, New York.

WATSON, G. N. (1944). *A Treatise on the Theory of Bessel Functions*, 2nd edn. Cambridge University Press.

YAMATO, H., SIBUYA, M. AND NOMACHI, T. (2001). Ordered sample from two-parameter GEM distribution. *Statist. Prob. Lett.* **55,** 19–27.