

Applying Pitman's sampling formula to microdata disclosure risk assessment

Nobuaki Hoshino*

Faculty of Economics, Kanazawa University

Abstract

Bethlehem et al. (1990) proposed a superpopulation model called Poisson-gamma model to assess microdata disclosure risk. Takemura (1999) introduced the Ewens sampling formula (Ewens (1972)) studied in statistical ecology to the disclosure problem as a limiting form of a conditional Poisson-gamma model. Pitman (1995) considered an extension of the Ewens sampling formula in a different context, and in this paper we assess usefulness of the Pitman sampling formula in the disclosure field. After giving some theoretical properties of the Pitman model, we compare various superpopulation models based on the Akaike Information Criterion (AIC) by applying them to real data sets from the Japanese labor force survey. Our comparison strongly supports the Pitman model. This result suggests that our superpopulation model based approach is very promising for the microdata disclosure problem as well as for statistical ecology.

keywords: Privacy, Uniqueness, Species abundance, Superpopulation, Random clustering

1 Introduction

In releasing a microdata set, the statistical agency must eliminate records identifiable to a particular individual. A record is composed of fields that correspond to categorized attributes of an individual. Attackers might identify an individual using information on records. In practical sense, we may consider individuals that are unique in the population with respect to the categorization in sample data to be identifiable. The number of population uniques is thus an important control object in the context of microdata disclosure, and it is important to estimate the number of population uniques from sample data at hand. We may regard a data set that contains many population uniques as risky.

**Address for correspondence* : Nobuaki Hoshino, Faculty of Economics, Kanazawa University, Kakuma-machi, Kanazawa-shi, Ishikawa, Japan. E-mail: hoshino@kenroku.kanazawa-u.ac.jp

To estimate population uniques, Bethlehem et al. (1990) introduced the Poisson-gamma model, which is the first application of the superpopulation model in the field of the microdata disclosure problem. Under the superpopulation model based approach, we assume that the population is generated by an appropriate (prior) distribution. By means of the assumption on the prior distribution, the risk inference is reduced to the problem of parameter estimation. We should be pragmatic since it is impossible to know the true mechanism of generating population. Here we adopt empirical Bayes methods; what is required is a prior distribution flexible enough to describe various populations.

We briefly survey various superpopulation models used in literature. Several authors apply the Poisson-gamma model to actual data sets, but insufficient fits are reported. See Skinner (1992) or Skinner and Holmes (1993). Skinner and Holmes (1993) applied the Poisson-lognormal model and the logarithmic series distribution to US and Italian data sets. These models have mainly been studied in ecology, where frequencies of species are estimated from sample frequency structure. The stochastic abundance model (Engen (1978)) is used for modeling the populations consisting of large number of species in statistical ecology. Hoshino and Takemura (1998) clarified relations between various superpopulation models and revealed that the superpopulation model based approach in the disclosure problem has a connection with the stochastic abundance models. The Poisson-lognormal model is studied, for example, in Bulmer (1974) or Aitchison and Ho (1989). Fisher's classical logarithmic series model (Fisher et al. (1943)) leads to many versions of superpopulation models; see Section 3.2 of Engen (1978) and Johnson et al. (1993). Hoshino and Takemura (1998), based on an interpretation of Anscombe (1950), noted that a limiting Poisson-gamma model becomes a logarithmic series model different from Skinner and Holmes' (1993). Takemura (1999) considered a sampling distribution from the Poisson-gamma model and derived the Dirichlet-multinomial model. Takemura (1999) also identified that the Ewens sampling formula originally developed in genetics is a limiting form of the Dirichlet-multinomial model. See Ewens (1990), Sibuya (1993) and Johnson et al. (1997) for the Ewens distribution. In Hoshino and Takemura (1998), we showed that the Ewens model is derived from the logarithmic series model by the same conditioning argument as the Dirichlet-multinomial model is derived from the Poisson-gamma model. Watterson (1973) referred to the Ewens distribution as a version of the logarithmic series distribution.

Pitman (1995) considered the random partition of the positive integers, and obtained a new generalization of the Ewens distribution. See Pitman (1996), Pitman and Yor (1997), Yamato

et al. (1999) for the context. The obtained distribution is the Pitman sampling formula, which contains the Ewens model and the Dirichlet-multinomial model as special cases. Thus the fit of the Pitman model is at least as good as that of these models, though the degree of freedom decreases. If the Pitman model greatly improves prediction on the disclosure risk, then the superpopulation model based approach becomes much more relevant not only for the disclosure problem but also for the stochastic abundance model fitting. It is important to apply above superpopulation models to actual data sets and compare each model on the same appropriate criterion. The conventional χ^2 type criterion is not desirable for the comparison of models that have different number of parameters. Therefore we choose the Akaike Information Criterion (AIC), which adjusts the differences in the number of parameters on the comparison. It seems that tools of testing hypotheses are still commonly used for the purpose of model selection in the area of stochastic abundance models. We later discuss the model selection criterion.

The organization of this paper is as follows. In Section 2 we derive some relevant moments of the Pitman model. Estimation problems concerning the Pitman model are discussed in Section 3. We compare the Pitman model with other superpopulation models in Section 4, applying the models to Japanese labor force survey data sets. Section 5 offers motivation for the Pitman model and the conclusion. In the rest of this section we fix notation and define existing superpopulation models compared in Section 4.

1.1 Notation and summary of existing superpopulation models

Consider a discrete population of size N . Let K denote the total number of the cells and let F_j , $j = 1, \dots, K$, denote the size of the j -th cell. Under the superpopulation model approach we consider F_j , $j = 1, \dots, K$, as random variables; the population size $N = \sum_{j=1}^K F_j$ may or may not be a random variable. Let S_i denote the number of the cells of size i . In terms of the indicator function

$$I(F_j = i) = \begin{cases} 1, & F_j = i, \\ 0, & F_j \neq i, \end{cases}$$

the number of the cells of size i is expressed as

$$S_i = \sum_{j=1}^K I(F_j = i), \quad i = 0, 1, \dots,$$

which are called size indices (Sibuya (1993)) or frequencies of frequencies (Good (1965)). These ideas correspond to equivalence class (Greenberg and Zayatz (1992)) in the context of the microdata disclosure problem. In disclosure risk assessment, the number of population uniques S_1 is of particular importance.

Obviously

$$\sum_{i=0}^{\infty} S_i = K, \quad \sum_{i=1}^{\infty} i \cdot S_i = N.$$

Here K is the total number of the cells including the number of the empty cells S_0 . In the following we denote the number of the non-empty cells by

$$U = K - S_0 = \sum_{i=1}^{\infty} S_i.$$

One important difference between the disclosure problem and statistical ecology is the handling of U and K . In statistical ecology usually we only consider the marginal distribution of (S_1, \dots) , and do not include K in the models. The reason is that species of frequency zero in population have little meaning and there exists no means to specify S_0 in statistical ecology. However, as far as the microdata problem is concerned, we can set K as the product of the number of categories in variables assessed. Generally K becomes huge. The limiting process of $K \rightarrow \infty$ is thus reasonable.

In the following we summarize existing superpopulation models. We classify these models by paying attention to the following two points: (a) whether the population size N is a random variable or not, and (b) whether S_0 is defined or not. Models in which S_0 is not defined are described without explicit dependence on K .

Poisson-gamma model : The population size N is a random variable having the negative binomial distribution, and S_0 is defined. Under the Poisson-gamma model, F_j is the Poisson random variable with mean $N_0\mu$ and μ has the gamma distribution with parameters γ and β , which are assumed to satisfy the restriction $\gamma\beta = 1/K$. The unconditional distribution of F_j becomes the negative binomial distribution. Furthermore F_j , $j = 1, \dots, K$, are assumed to be independently and identically distributed. In summary the Poisson-gamma model is defined by the joint probability function of F_j 's as

$$P(F_1, \dots, F_K) = \prod_{j=1}^K \frac{\Gamma(F_j + \gamma)}{\Gamma(\gamma) F_j!} p^\gamma q^{F_j}, \quad q = \frac{N_0\beta}{N_0\beta + 1}, \quad p = 1 - q, \quad \gamma\beta = \frac{1}{K}. \quad (1)$$

The expected population size is $E(N) = KE(F_j) = KN_0\gamma\beta = N_0$.

Poisson-lognormal model : The population size N is a random variable, and S_0 is defined. In the Poisson-lognormal model, F_j is the Poisson random variable with mean λ , and $\log \lambda$ is normally distributed with mean M and variance V . As in the Poisson-gamma model, F_j , $j = 1, \dots, K$, are assumed to be independently and identically distributed. The Poisson-lognormal model is defined by

$$P(F_1, \dots, F_K) = \prod_{j=1}^K \frac{1}{F_j! \sqrt{2\pi V}} \int_0^\infty \lambda^{F_j-1} \exp(-\lambda - (\log \lambda - M)^2/2V) d\lambda. \quad (2)$$

The expected population size becomes $K \exp(M + V/2)$. In this paper we restrict the model such that $K \exp(M + V/2) = N_0$. Namely $M = \log N_0 - \log K - V/2$, and V is the unique parameter.

Dirichlet-multinomial model : The population size N is fixed, and S_0 is defined. The Dirichlet-multinomial model is the conditional Poisson-gamma model given N and defined by

$$P(S_0, \dots, S_N) = \frac{N! K! \Gamma(K\gamma)}{\Gamma(K\gamma + N)} \prod_{i=0}^N \left(\frac{\Gamma(\gamma + i)}{\Gamma(\gamma) i!} \right)^{S_i} \frac{1}{S_i!}. \quad (3)$$

Logarithmic series model : The population size N is a random variable, and S_0 is not defined. Fisher's logarithmic series model is defined in terms of the joint distribution of size indices S_i , $i \geq 1$. Let

$$\lambda_i = N_0 \frac{p \cdot q^{i-1}}{i}, \quad i = 1, 2, \dots,$$

where $N_0 > 0$, $0 < p < 1$ and $q = 1 - p$. Here S_i is independent Poisson random variable with mean λ_i . The joint probability function of the size indices (S_1, S_2, \dots) becomes

$$P(S_1, S_2, \dots) = \prod_{i=1}^{\infty} \frac{\lambda_i^{S_i} \exp(-\lambda_i)}{S_i!}. \quad (4)$$

Here only finite number of S_i 's are nonzero. This model is the limiting form of the Poisson-gamma model as $K \rightarrow \infty$ with $K\gamma$ fixed.

Ewens model : The population size N is fixed, and S_0 is not defined. Applying the limiting process $K \rightarrow \infty$ to the Dirichlet-multinomial model with $K\gamma = \theta$ fixed, we obtain the Ewens model with parameter θ :

$$P(S_1, \dots, S_N) = \frac{\theta^U}{\theta^{[N]}} \frac{N!}{\prod_{i=1}^N i^{S_i} S_i!}, \quad (5)$$

where $\theta^{[N]} = \theta(\theta + 1)(\theta + 2) \cdots (\theta + N - 1)$, $U = \sum_{i=1}^N S_i$.

2 Some theoretical results on the Pitman sampling formula

In this section we introduce the Pitman model and derive some moments of the model needed for the estimation.

Sibuya (1993) describes the urn scheme construction of the Ewens sampling formula. To begin with, we explain the urn model implication of the Pitman model, based on Proposition 9 in Pitman (1995). Let us consider the following process: Suppose that n balls are distributed over u urns such that no empty urn exists; we put a new ball to a new urn with the probability of

$$\frac{\theta + u\alpha}{\theta + n}, \quad (6)$$

or put the ball to one of the existing u urns with the probability of

$$\frac{j - \alpha}{\theta + n}, \quad (7)$$

where j is the number of balls in the urn. Starting the process with $n = u = 0$, we obtain the Pitman sampling formula.

For each pair of real parameters α and θ , such that either $0 \leq \alpha < 1$ and $\theta > -\alpha$, or $\alpha < 0$ and $\theta = -m\alpha$ for some natural number m , the Pitman model is defined by

$$P(S_1, \dots, S_N) = N! \frac{\theta^{[U:\alpha]}}{\theta^{[N]}} \prod_{j=1}^N \left(\frac{(1 - \alpha)^{[j-1]}}{j!} \right)^{S_j} \frac{1}{S_j!}, \quad (8)$$

where $\theta^{[U:\alpha]} = \theta(\theta + \alpha) \cdots (\theta + (U - 1)\alpha)$, $\theta^{[N]} = \theta(\theta + 1) \cdots (\theta + N - 1)$. If α equals zero, (8) amounts to the Ewens model (5). Assuming that $\alpha < 0$, let $\theta = -K\alpha > 0$, $\gamma = -\alpha > 0$. Then (8) amounts to the Dirichlet-multinomial model (3).

Write

$$S_{N,U} = \{S = (S_1, \dots, S_N) \mid \sum_{i \geq 1} i S_i = N, \sum_{i \geq 1} S_i = U\}.$$

In the following we consider $U_N = \sum_{i=1}^N S_i$ as a random variable given N . Namely

$$P(U_N = U) = \begin{cases} \sum_{S \in \mathcal{S}_{N,U}} N! \frac{\theta^{[U:\alpha]}}{\theta^{[N]}} \prod_{j=1}^N \left(\frac{(1-\alpha)^{[j-1]}}{j!} \right)^{S_j} \frac{1}{S_j!} & \text{if } U = 1, \dots, N, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Yamato et al. (1999) give the explicit form of (9), though it is complicated.

Theorem 1 *Suppose that size indices are distributed according to (8). Then*

$$E(S_i) = \frac{(1-\alpha)^{[i-1]}}{i!} (\theta + \alpha \cdot E(U_{N-i})) \prod_{j=1}^i \frac{N-j+1}{\theta+N-j}, \quad (10)$$

where

$$E(U_{N-i}) = \frac{\theta}{\theta+N-i-1} + \sum_{l=0}^{N-i-2} \frac{\theta}{\theta+l} \prod_{j=l+1}^{N-i-1} \left(1 + \frac{\alpha}{\theta+j} \right). \quad (11)$$

Proof First we show (10).

$$\begin{aligned} E(S_i) &= \sum_{U=1}^N \sum_{S \in \mathcal{S}_{N,U}} N! \frac{\theta^{[U:\alpha]}}{\theta^{[N]}} \prod_{j=1}^N \left(\frac{(1-\alpha)^{[j-1]}}{j!} \right)^{S_j} \frac{1}{S_j!} S_i \\ &= \sum_{U=1}^N \sum_{S \in \mathcal{S}_{N,U}} N! \frac{\theta^{[U:\alpha]}}{\theta^{[N]}} \prod_{j=1(j \neq i)}^N \left(\frac{(1-\alpha)^{[j-1]}}{j!} \right)^{S_j} \frac{1}{S_j!} \left(\frac{(1-\alpha)^{[i-1]}}{i!} \right)^{S_i} \frac{1}{(S_i-1)!} \\ &= \sum_{U=1}^N \sum_{S \in \mathcal{S}_{N-i,U-1}} N! \frac{\theta^{[U:\alpha]}}{\theta^{[N]}} \prod_{j=1}^N \left(\frac{(1-\alpha)^{[j-1]}}{j!} \right)^{S_j} \frac{1}{S_j!} \frac{(1-\alpha)^{[i-1]}}{i!} \\ &= \sum_{U=1}^N \frac{(1-\alpha)^{[i-1]}}{i!} \prod_{l=1}^i \frac{N-l+1}{\theta+N-l} (\theta + (U-1)\alpha) P(U_{N-i} = U-1). \end{aligned} \quad (12)$$

Since $\sum_{U=1}^N (U-1) P(U_{N-i} = U-1) = E(U_{N-i})$, (10) is proved.

To prove (11) we utilize a recurrence relation:

$$E(U_{N+1}) = E(U_N) + \frac{\theta + E(U_N)\alpha}{\theta + N}, \quad (13)$$

where $E(U_0) = 0$. Assuming that (13) is true, we can easily prove (11) by induction. The relation of (13) holds from the fact that

$$P(U_{N+1} = U) = P(U_{N+1} = U | U_N = U) P(U_N = U) + P(U_{N+1} = U | U_N = U-1) P(U_N = U-1)$$

and

$$P(U_{N+1} = U + 1 | U_N = U) = \frac{\theta + U\alpha}{\theta + N}, \quad P(U_{N+1} = U | U_N = U) = \frac{N - U\alpha}{\theta + N} \quad (14)$$

derived from (6) and (7).

Q.E.D.

We can similarly calculate the factorial moment:

$$E\left(\prod_{j=1}^N S_j^{(r_j)}\right) = \sum_{U=1}^N \frac{N^{(R)}\theta^{[U:\alpha]}}{(\theta + N - 1)^{(R)}\theta^{[U-r:\alpha]}} \prod_{j=1}^N \left(\frac{(1-\alpha)^{[j-1]}}{j!}\right)^{r_j} P(U_{N-R} = U - r), \quad (15)$$

where $r = \sum_{j=1}^N r_j$, $R = \sum_{j=1}^N j r_j$ and $n^{(R)} = n(n-1)\cdots(n-R+1)$. The higher moments of U_N are evaluated through recurrence relations like (13). For example

$$E(U_{N+1}^2) = \frac{N + \theta + 2\alpha}{\theta + N} E(U_N^2) + \frac{2\theta + \alpha}{\theta + N} E(U_N) + \frac{\theta}{\theta + N}.$$

In particular we obtain

$$E(S_1) = \frac{N\theta + N\alpha E(U_{N-1})}{\theta + N - 1}, \quad (16)$$

$$E(S_2) = E\left[\frac{N^{(2)}(\theta + \alpha U_{N-2})}{(\theta + N - 1)^{(2)}}\right] \frac{1 - \alpha}{2}, \quad (17)$$

$$E(S_1^{(2)}) = E\left[\frac{N^{(2)}(\theta + \alpha U_{N-2})(\theta + \alpha(U_{N-2} + 1))}{(\theta + N - 1)^{(2)}}\right] \quad (18)$$

from (10) and (15). These expectations lead to moment estimators discussed in Section 3.

3 Estimation with the Pitman model

In this section we consider the estimation of the disclosure risk under the Pitman model. We denote sample size by n and sample size indices by $s = (s_1, \dots, s_n)$. The total number of nonempty cells or clusters is $u = \sum_{i=1}^n s_i$. Suppose that n individuals are drawn simply at random without replacement.

3.1 The estimation of the parameters

The Pitman model enjoys the property of exchangeability, with respect to individuals in population, assumed in Lemma 1 of Takemura (1999). Accordingly the marginal distribution of

sample individuals coincides with the prior distribution of values of n individuals directly drawn from the superpopulation. That is to say,

$$P(s_1, \dots, s_n) = n! \frac{\theta^{[u:\alpha]}}{\theta^{[n]}} \prod_{j=1}^n \left(\frac{(1-\alpha)^{[j-1]}}{j!} \right)^{s_j} \frac{1}{s_j!} \quad (19)$$

is obtained from replacing N and U in (8) by n and u . We can show the result in another way. Suppose that N objects are partitioned into classes according to a probability distribution p_N . A partition structure (Kingman (1978)) is a sequence p_1, p_2, \dots of distributions wherein, assuming that an object is deleted uniformly at random from the N objects, the partition of the $N - 1$ remaining objects is distributed according to p_{N-1} . The Pitman sampling formula is known to have a partition structure, with the result that (19) holds.

We then construct the Maximum Likelihood Estimator (MLE) of θ and α . Let the logarithm of the right hand side of (19) be L . The MLE is the solution of

$$\frac{\partial L}{\partial \theta} = \sum_{i=1}^{u-1} \frac{1}{\theta + i\alpha} - \sum_{i=1}^{n-1} \frac{1}{\theta + i} = 0$$

and

$$\frac{\partial L}{\partial \alpha} = \sum_{i=1}^{u-1} \frac{i}{\theta + i\alpha} - \sum_{i=2}^n s_i \sum_{j=1}^{i-1} \frac{1}{j - \alpha} = 0.$$

These simultaneous equations can be solved by the Newton-Raphson method using second derivatives:

$$\begin{aligned} \frac{\partial^2 L}{(\partial \theta)^2} &= - \sum_{i=1}^{u-1} \frac{1}{(\theta + i\alpha)^2} + \sum_{i=1}^{n-1} \frac{1}{(\theta + i)^2}, \\ \frac{\partial^2 L}{(\partial \alpha)^2} &= - \sum_{i=1}^{u-1} \frac{i^2}{(\theta + i\alpha)^2} - \sum_{i=2}^n s_i \sum_{j=1}^{i-1} \frac{1}{(j - \alpha)^2}, \\ \frac{\partial^2 L}{\partial \theta \partial \alpha} &= - \sum_{i=1}^{u-1} \frac{i}{(i\alpha + \theta)^2}. \end{aligned}$$

To solve the ML estimation, we investigate approximate moment estimators for the starting values of the Newton-Raphson procedure. Our moment estimators are

$$\hat{\theta} = \frac{nuc - s_1(n-1)(2u+c)}{2s_1u + s_1c - nc}, \quad (20)$$

$$\hat{\alpha} = \frac{\hat{\theta}(s_1 - n) + (n - 1)s_1}{nu}, \quad (21)$$

where $c = s_1(s_1 - 1)/s_2$. The derivation is given in Appendix. In six of seven cases in Section 4, except for Case 1, these estimators gave convergences in the Newton-Raphson procedure. In Case 1 the author reached the solution by random starting value generation.

3.2 Risk inference

In the following we discuss some statistics concerning the disclosure risk under the Pitman model. We state three propositions useful for the disclosure problem. All the proofs are given in Appendix.

As regards the risk inference, we shall evaluate the expectation of the number of population uniques $E(S_1)$ with the ML estimates of the parameters. However necessary moments given in Section 3 are not in convenient forms to compute. From Theorem 1 we investigate simple forms of the moments.

Proposition 1 *If $\alpha \neq 0$, the expectation of U_N under (8) is reduced to*

$$E(U_N) = \frac{\theta}{\alpha} \left(\frac{(\theta + \alpha)^{[N]}}{\theta^{[N]}} - 1 \right). \quad (22)$$

The result of Yamato and Sibuya (1999) coincides with (22). We can rewrite (22) using the gamma function. Based on the asymptotic property of the gamma function, we find a useful approximation of $E(U_N)$, which is a special case of Lemma 2 in Yamato and Sibuya (1999). If N is sufficiently large

$$E(U_N) = \frac{\theta}{\alpha} \left(\frac{\Gamma(\theta + \alpha + N)\Gamma(\theta)}{\Gamma(\theta + N)\Gamma(\theta + \alpha)} - 1 \right) \approx \frac{\Gamma(\theta + 1)}{\alpha\Gamma(\theta + \alpha)} N^\alpha \quad (23)$$

for $\alpha \neq 0$. Our expression of $E(S_1)$ depending on $E(U_{N-1})$ then becomes simpler. We obtain

$$E(S_1) = \frac{N\Gamma(\theta + \alpha + N - 1)\Gamma(\theta + 1)}{\Gamma(\theta + N)\Gamma(\theta + \alpha)} \approx \frac{\Gamma(\theta + 1)}{\Gamma(\theta + \alpha)} N^\alpha. \quad (24)$$

As a result, the evaluation of $E(S_1)$ is not very hard, once the ML estimates are obtained.

The following propositions may have interesting implications on the disclosure problem.

Proposition 2 *Suppose that size indices are distributed according to (8). For $\alpha \geq 0$*

$$\lim_{N \rightarrow \infty} \frac{E(S_1)}{E(U_N)} = \alpha.$$

Proposition 2 suggests that the ratio of population uniques to the number of non-zero frequency groups is α , which is smaller than unity; the implication is consistent with the author's experience that the Ewens model (i.e. $\alpha = 0$) tends to underestimate the number of population uniques. Since the Ewens model is a limiting form of the conditional Poisson-gamma model, these models give similar population unique estimates as can be seen in Section 4. It is suggested that the poor performance of the Poisson-gamma model and related models, including the logarithmic series distribution, occurs when population uniques constitute no negligible proportion of the population. In other words these models might be suitable only for safe data sets.

Based on Proposition 2, we propose a simple estimator of α :

$$\tilde{\alpha} = \frac{s_1}{u}. \quad (25)$$

We could replace the previous moment estimator (21) by (25); see Table 7 of Section 4.

Proposition 3 *Let $n/N = f$ be fixed. If we assume (8) and (19) then*

$$\lim_{N \rightarrow \infty} \frac{E(S_1)}{E(s_1)} f = f^{1-\alpha}. \quad (26)$$

The left hand side of (26) is interpreted as the ratio of population uniques in the sample to sample uniques s_1 ; this ratio is often an index of the disclosure risk. Combining the simple estimate of (25), we can roughly evaluate the risk of a data set by $f^{1-\tilde{\alpha}}$, where the sampling ratio f is known. This simple procedure is useful because the data editing for anonymization requires repeated trial and error.

4 An application to Japanese labor force survey data

In this section we examine performances of the Pitman model and other superpopulation models. Takemura (1998) gives some size indices data of the Japanese labor force survey. We apply the Pitman model (8), the Ewens model (5), the Poisson-gamma model (the Dirichlet-multinomial

model (3)) and the Poisson-lognormal model (29) to the data, and compare each model by its Akaike Information Criterion (AIC) value.

A χ^2 type statistic like

$$\chi^2 = \sum_{i=1}^n \frac{(s_i - E(s_i))^2}{E(s_i)} \quad (27)$$

is conventionally used to evaluate the goodness of stochastic abundance model fitting. If (s_1, \dots, s_n) is multinomially distributed given u , then (27) is the classical χ^2 test. The symmetrical model description in terms of independent F_j 's can be converted in terms of S_i 's, where

$$P(S_0, \dots) = \binom{K}{S_0 S_1 \dots} \prod_{i=0}^{\infty} P(F = i)^{S_i} \quad (28)$$

is in the form of multinomial distribution. Since the marginal distribution of the multinomial distribution becomes multinomial, the χ^2 type statistic might be suitable for the Poisson-gamma model and the Poisson-lognormal model. However the assumption seems to be inappropriate for the other superpopulation models. As described in Section 7.2 of Engen (1978), we can only use (27) to “form a picture of the similarity” between s_i 's and $E(s_i)$'s. In the disclosure context, Zayatz (1991) used the Kolmogorov-Smirnov goodness of fit test for the Poisson-gamma model and found a significant lack of fit at the .01 level. Skinner and Holmes (1993) calculated (27) and likelihood ratio statistics for the logarithmic series distribution and the Poisson-lognormal model. These ideas are based on the theory of testing hypotheses and we can, only at best, tell whether the model assumption is acceptable or not. In other words these statistics are not comparable between different models.

Let the number of parameters in a model be λ . Let the log likelihood of the model maximized with respect to the parameters be denoted by L . The AIC selects the model that has the lowest $-2L + 2\lambda$. See Atkinson (1980) or Konishi and Kitagawa (1996), for example.

The likelihood depends on the sampling mechanism. We assume simple random sampling without replacement for the Pitman model, the Ewens model and the Poisson-gamma model. The sampling distributions of the Pitman model and the Ewens model are again the Pitman model and the Ewens model as in (19). Under simple random sampling without replacement, the sampling distribution of the Fisher's logarithmic series model is the Ewens model (Hoshino and Takemura (1998)). Thus we only evaluate the logarithmic series model through the Ewens model. See Sibuya (1991) or Hoshino and Takemura (1998) for the parameter estimation of the

Ewens model. The sampling distribution $P(s_0, \dots | n)$ of the Poisson-gamma model becomes the Dirichlet-multinomial model. Thus we evaluate the MLE of the Poisson-gamma model by the Dirichlet-multinomial model.

Under simple random sampling without replacement, the sampling distribution $P(s_0, \dots | n)$ of the Poisson-lognormal model is hard to manipulate. Therefore we assume the Bernoulli sampling (Särndal et al. (1992)) in which each individual is drawn if a coin with some success probability results in head. This scheme is a convenient approximation to simple random sampling without replacement, but it is more natural in ecological sampling than simple random sampling without replacement. When the success probability is n/N_0 , we obtain the sampling distribution $P(s_0, \dots)$ of the form (28) replacing N_0 by n . Another approximation we use is the normal approximation of the sample size distribution. The variance of the sample size becomes $T = K(\exp(M + V/2) + \exp(2M + 2V) - \exp(2M + V))$, and the expected sample size is set to n . Therefore we set the probability of the sample size to be n as $1/\sqrt{2\pi T}$. Consequently the conditional Poisson-lognormal model $P(s_0, \dots)/P(\sum is_i = n)$ is approximated by

$$P(s_0, \dots | n) = \binom{K}{s_0 \cdots s_n} \prod_{i=0}^n \left\{ \frac{1}{i! \sqrt{2\pi V}} \int_0^\infty \lambda^{i-1} \exp(-\lambda - (\log \lambda - M)^2/2V) d\lambda \right\}^{s_i} \sqrt{2\pi T}, \quad (29)$$

where $M = \log n - \log K - V/2$. We need numerical integration to evaluate the model. A transformation suited to the Hermitian integration is discussed in Aitchison and Ho (1989). The author programmed the numerical integration with GNU C compiler, checking results against Grundy (1951).

Now we sketch the data to be assessed. The purpose of the labor force survey is to elucidate the current state of employment and unemployment in Japan. The data at hand was collected in January 1995. The population of the survey is composed of all persons 15 years old and over usually residing in 47 prefectures of Japan. However the data at hand consists of persons only in nine prefectures near Tokyo. The sample size n is 27230. Corresponding population size N is about 35.85 million. For simplicity we assume that individuals are drawn simply at random without replacement, although the actual sampling scheme is more complicated. Seven different combinations of “global recoding” and “global suppression” are applied to the data. See Willenborg and de Waal (1996) for these techniques of anonymization. The size indices are enumerated with respect to the categorization of nine (Case 1–2), eight (Case 3–6) and seven (Case 7) variables. These variables are geographical codes, classified number of persons

in household, relationship to head of household, sex, age, and marital status. Table 1 provides more information on the categorization. The results of our model fitting are tabulated in Table 2 to 5.

The Pitman model highly dominates in all the cases, and the Poisson-lognormal model shows the least performance. Engen (1978) and Skinner and Holmes (1993) reported relatively good fits of the Poisson-lognormal model based on the χ^2 type statistic in (27). It might be the case that for the Poisson-lognormal model the maximization of the marginal likelihood on (s_1, \dots, s_n) gives a different estimate. Figure 1 illustrates the fits of the Ewens model, the Pitman model and the Poisson-lognormal model in Case 7. The vertical axis corresponds to $E(s_i)$'s under the ML estimates of the parameters, and the horizontal axis corresponds to $i = 1, \dots, 15$. The actual sample size indices are plotted in the same scale. Under the Poisson-lognormal model, $E(s_1)$ shows huge overshoots. It seems that the inclusion of zero frequency groups causes the lowest fit of the Poisson-lognormal model. Note that the Pitman model ignores the restriction that K is finite. Thus for fairer comparison, we explore $K = K^*$ in which the Poisson-lognormal model attains the smallest AIC value; it is a kind of marginal fitting. The results are provided in Table 6. The Pitman model still dominates except for Case 3. We accordingly observe the strong support of the Pitman model.

However we further analyze the results for more detailed evaluation. Case 3 shows some peculiarity when we apply the simple estimator (25) of α . Table 7 lists differences between the MLE and (25); there is a large difference in Case 3 of Table 7. We may be able to regard (25) as a model check. Since there exists no all-purpose estimation procedure, we should probably examine the possibility of an alternative approach.

Let us then return to Table 6. We first realize that K is much greater than K^* , and $E(S_1)$ with K^* is very small. These facts suggest that the fit of the Poisson-lognormal model has no robustness in withstanding changes of K . The presence of structural zeros, for example caused by the cross-classification of age and marital status, may lead us to accept a claim that true K is smaller than the product of the number of categories in variables assessed. Although the author considers that structural zeros are also realizations from a superpopulation, let us consider the possibility of such decrease in K . We may regard K^* as an estimate of true K , but it seems that K^* is too small; a great underestimate of K implies an underestimate of population uniques. Concerning the risk inference, an underestimate should be more heavily penalized than an overestimate. It is thus not persuasive to believe that K^* equals true K . Note that

an arbitrariness can not be totally eliminated in determining true K . It is therefore preferable that the risk inference does not depend on K . However, this independence does not seem to hold in the Poisson-lognormal model. Moreover Engen (1978) provides an example in which the estimated parameter of the Poisson-lognormal model by marginal (excluding zero groups) fitting varies with respect to the size of the sample from the same population. It suggests that the use of K^* by marginal fitting leads to the erroneous estimate of population uniques.

Figure 1 clearly represents a typical tendency of model fitting in the disclosure field. We often observe a great difference between s_1 and s_2 . The author considers that “shape” parameter is required to describe this kind of non-smoothness. In view of the urn model implication, the Pitman parameter α specially adjusts the rate of unique cells. This fact would be the reason why the Pitman model dominates.

5 Discussion

5.1 The Pitman model and the lognormal distribution

Construction 16 of Pitman (1995) provides another derivation of the Pitman sampling formula. In this section we observe that it gives a justification similar to that of the lognormal model for the Pitman model. This interpretation may motivate the Pitman model.

The lognormal distribution has long been used to describe various populations of species, savings in households, mineral gains and numerous seemingly unrelated objects. Halmos (1944) gave the following justification over the wide applicability of the lognormal distribution. Let $\mathbf{W} = (W_1, W_2, \dots)$ be a sequence of random variables, where $0 \leq W_i \leq 1, i = 1, 2, \dots$. Define $\bar{W}_i = 1 - W_i$. Let

$$P_i = \bar{W}_1 \cdots \bar{W}_{i-1} W_i, \quad i = 1, 2, \dots \quad (30)$$

Then $\mathbf{P} = (P_1, P_2, \dots)$ constitutes a random classification where the i -th group has proportion P_i . The equation (30) implies that $P_i = (1 - P_1 - P_2 - \cdots - P_{i-1})W_i$. Hence the process allocates the residual. The logarithm of (30) equals

$$\log P_i = \log \bar{W}_1 + \cdots + \log \bar{W}_{i-1} + \log W_i.$$

The right hand side is a sum of random variables. Thus under appropriate regularity conditions the central limit theorem holds, and $\log P_i$ is normally distributed. Namely P_i is subject to the lognormal distribution in many cases.

Assume that $X_j, j = 1, \dots, N$, are independently identically distributed given \mathbf{P} with $P(X_j = i|\mathbf{P}) = P_i, i = 1, 2, \dots$. Here X_j is the j -th sample from the infinite population of individuals. We can interpret P_i as the long run relative frequency of the i -th group. The marginal distribution of the frequency $F_i = \sum_{j=1}^N I(X_j = i)$ given \mathbf{P} is the binomial distribution:

$$P(F_i = y|\mathbf{P}) = \binom{N}{y} P_i^y (1 - P_i)^{N-y}, \quad y = 0, 1, \dots, N.$$

It is well known that the binomial distribution above is approximated by the Poisson distribution with mean NP_i . If $\log NP_i = \log N + \log P_i$ is subject to the normal distribution, then the marginal frequency of the i -th group is approximately the Poisson-lognormal.

We now turn to the Pitman model. Let us suppose that W_i of (30) independently possesses the beta distribution with parameters $(1 - \alpha, \theta + i\alpha)$, where $0 \leq \alpha < 1, \theta > -\alpha$. Now we can explicitly derive the distribution of the samples X_1, \dots, X_N . According to Pitman (1995), the size indices of the samples are then subject to (8): the Pitman model, in fact.

We have observed that the same residual allocation structure induces the Pitman model and the lognormal distribution. The Pitman model may consequently have a motivation for use in the disclosure field, because the Poisson-lognormal model has been used to measure the risk. Yamato et al. (1999) clearly explain the derivation of the Pitman model from the beta distribution. The corresponding derivation of the Ewens model is given in Johnson et al. (1997). The process of (30) with independent W_i is known as the residual allocation model. See Pitman (1996) for a survey.

5.2 Concluding remarks

We saw that the Pitman model fits well to the Japanese labor force survey data set in comparison with other existing superpopulation models. It should also be emphasized that the computation on the Pitman model is not so heavy compared to the Poisson-lognormal model. Thus it seems that the Pitman model is a promising tool for the disclosure risk assessment. This section appends a few arguments in this regard.

We generally face difficulties in estimation problems concerning the tail of a distribution, such as population uniques. For instance, even the approximation of simple random sampling without replacement by the Bernoulli sampling might considerably affect the distribution of the tail; preferable models are those that employ no approximation in the sampling scheme. The

Pitman model, due to its partition structure, is consistent with the sampling scheme in the disclosure field. With stratified sampling structure, we may apply the Pitman model in each stratum.

However we should note that the Pitman model ignores the restriction of K . If there is a great difference between the sample size and the population size, then the disregard of K may cause the overestimate of the risk; it is, in the extreme, possible that U becomes larger than the disregarded K . In applying the Pitman model, we should check whether U is too large compared to K .

Appendix

We first derive the moment estimators (20) and (21). For simplicity we use $E(s_1)$ and $E(s_1(s_1 - 1))/E(s_2)$ to estimate θ and α . We denote the total number of clusters given $n - 1$ and $n - 2$ by u_{n-1} and u_{n-2} . By (16)

$$\alpha = \frac{(\theta + n - 1)E(s_1) - n\theta}{nE(u_{n-1})}. \quad (31)$$

Referring to (17) and (18), we derive

$$\begin{aligned} C = \frac{E(s_1(s_1 - 1))}{E(s_2)} &= \frac{2E[(\theta + \alpha u_{n-2})(\theta + \alpha(u_{n-2} + 1))]}{(1 - \alpha)E(\theta + \alpha u_{n-2})} \\ &\approx \frac{2E(\theta + \alpha(u_{n-2} + 1))}{(1 - \alpha)}. \end{aligned}$$

Then

$$\alpha = \frac{C - 2\theta}{2E(u_{n-2}) + C + 2}. \quad (32)$$

From (31) and (32),

$$\theta = \frac{nE(u_{n-1})C - (n - 1)E(s_1)(2E(u_{n-2}) + C + 2)}{(E(s_1) - n)(2E(u_{n-2}) + C + 2) + 2nE(u_{n-1})}. \quad (33)$$

Now we give the moment estimator of θ . Ignoring the relation $E(u_{n-1}) = (\theta + n)/(\theta + n + \alpha)E(u_n) - \theta/(\theta + n + \alpha)$, we replace $E(u_{n-1})$ of (33) by u and $E(u_{n-2})$ by $u - 1$, whereby the estimator becomes simpler. Let $C = c = s_1(s_1 - 1)/s_2$. Substituting $E(s_1)$ of (33) by s_1 , we obtain (20). (21) is a direct consequence of (31).

In the following we show the propositions stated in Section 3.

Proof of Proposition 1 We derive (22) from (11). The proposition is shown by the relation that

$$\begin{aligned}
1 + \frac{\alpha E(U_N)}{\theta} &= 1 + \frac{\alpha}{\theta + N - 1} + \sum_{l=0}^{N-2} \frac{\alpha}{\theta + l} \prod_{j=l+1}^{N-1} \left(\frac{\theta + j + \alpha}{\theta + j} \right) \\
&= \frac{\theta + N - 1 + \alpha}{\theta + N - 1} \left\{ 1 + \sum_{l=0}^{N-3} \frac{\alpha}{\theta + l} \prod_{j=l+1}^{N-2} \left(\frac{\theta + j + \alpha}{\theta + j} \right) + \frac{\alpha}{\theta + N - 2} \right\} \\
&= \frac{(\theta + N - 1 + \alpha)(\theta + N - 2 + \alpha)}{(\theta + N - 1)(\theta + N - 2)} \left\{ 1 + \sum_{l=0}^{N-4} \frac{\alpha}{\theta + l} \prod_{j=l+1}^{N-3} \left(\frac{\theta + j + \alpha}{\theta + j} \right) + \frac{\alpha}{\theta + N - 3} \right\} \\
&\quad \vdots \\
&= \frac{(\theta + \alpha)^{[N]}}{\theta^{[N]}}.
\end{aligned}$$

Q.E.D.

To prove Proposition 2, we need the lemma below.

Lemma 1 For $\alpha \geq 0$,

$$\lim_{N \rightarrow \infty} E(U_N) = \infty.$$

Proof With respect to nonnegative α , $E(U_N)$ is monotonically increasing. Thus it suffices to show that $E(U_N)$ diverges at $\alpha = 0$. When α equals zero

$$E(U_N) = \sum_{l=0}^{N-1} \frac{\theta}{\theta + l}$$

from (11), and it is well known that the right hand side diverges as $N \rightarrow \infty$.

Q.E.D.

Proof of Proposition 2 From (10) we obtain

$$\frac{E(S_1)}{E(U_N)} = \frac{N\alpha}{\theta + N - 1} \frac{E(U_{N-1})}{E(U_N)} + \frac{N\theta}{\theta + N - 1} \frac{1}{E(U_N)}.$$

Since

$$E(U_N) = E(U_{N-1}) \left(1 + \frac{\alpha}{\theta + N} \right) + \frac{\theta}{\theta + N}$$

from (13), $E(U_N)/E(U_{N-1}) \rightarrow 1$ as $N \rightarrow \infty$. Also $1/E(U_N) \rightarrow 0$ by Lemma 1. We consequently obtain the formula.

Q.E.D.

Proof of Proposition 3 By the relation of (24),

$$\frac{E(S_1)}{E(s_1)} = \frac{N\Gamma(\theta + \alpha + N - 1)\Gamma(\theta + n)}{n\Gamma(\theta + N)\Gamma(\theta + \alpha + n - 1)}.$$

The formula then holds because of the asymptotic relation of the gamma function. Q.E.D.

Acknowledgment

The author would like to thank Prof. Akimichi Takemura for many useful comments and Prof. Masaaki Sibuya for stimulating discussion.

References

- [1] Aitchison, J. and Brown, J.A.C. (1957). *The Lognormal Distribution*, Cambridge University Press, Cambridge.
- [2] Aitchison, J. and Ho, C.H. (1989). The multivariate Poisson-log normal distribution. *Biometrika*, **76**, 4, 643–53.
- [3] Anscombe, F.J. (1950). Sampling theory of the negative binomial and logarithmic series distributions. *Biometrika*, **37**, 358–382.
- [4] Atkinson, A.C. (1980). A note on the generalized information criterion for choice of a model. *Biometrika*, **67**, 2, 413–418.
- [5] Bethlehem, J.G., Keller, W.J. and Pannekoek, J. (1990). Disclosure control of microdata. *Journal of the American Statistical Association*, **85**, 38–45.
- [6] Bulmer, M.G. (1974). On fitting the Poisson lognormal distribution to species-abundance data. *Biometrics*, **30**, 101–110.
- [7] Engen, S. (1978). *Stochastic Abundance Models*. Chapman and Hall, London.
- [8] Ewens, W.J. (1972). The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, **3**, 87–112, and erratum, p. 376.
- [9] Ewens, W.J. (1990). Population genetics theory – the past and the future. in *Mathematical and Statistical Development of Evolutionary Theory*, S. Lessard ed., 177–227, Kluwer, Dordrecht.

- [10] Fisher, R.A., Corbet, A.S. and Williams, C.B. (1943). The relation between the number of species and the number of individuals in a random sample of an animal population. *Journal of Animal Ecology*, **12**, 42–58.
- [11] Good, I.J. (1965). *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. MIT Press, Cambridge, Massachusetts.
- [12] Greenberg, B.V. and Zayatz, L.V. (1992). Strategies for measuring risk in public use microdata file. *Statistica Neerlandica*, **46**, 33–48.
- [13] Grundy, P.M. (1951). The expected frequencies in a sample of an animal population in which the abundances of species are log-normally distributed. *Biometrika*, **38**, 427–434.
- [14] Halmos, P.R. (1944). Random alms. *Annals of Mathematical Statistics*, **15**, 182–189.
- [15] Hoshino, N. and Takemura, A. (1998). Relationship between logarithmic series model and other superpopulation models useful for microdata disclosure risk assessment. *Journal of the Japan Statistical Society*, **28**, 2, 125–134.
- [16] Johnson, N.L., Kotz, S. and Kemp, A.W. (1993). *Univariate Discrete Distributions*, 2nd ed., Chap. 7, Wiley, New York.
- [17] Johnson, N.L., Kotz, S. and Balakrishnan, N. (1997). *Discrete Multivariate Distributions*, Chap. 41, Wiley, New York.
- [18] Kingman, J.F. (1978). Random partitions in population genetics. *Proceedings of the Royal Society of London, A*, **361**, 1–20.
- [19] Konishi, S. and Kitagawa, G. (1996). Generalised information criteria in model selection. *Biometrika*, **83**, 4, 875–890.
- [20] Pitman, J. (1995). Exchangeable and partially exchangeable random partitions. *Probability Theory and Related Fields*, **102**, 145–158.
- [21] Pitman, J. (1996). Random discrete distributions invariant under size-biased permutation. *Advances in Applied Probability*, **28**, 525–539.
- [22] Pitman, J. and Yor, M. (1997). The two-parameter Poisson-Dirichlet distribution derived from a stable subordinator. *Annals of Probability*, **25**, 855–900.

- [23] Särndal, C.E., Swensson, B. and Wretman, J. (1992). *Model Assisted Survey Sampling*. Springer, New York.
- [24] Sibuya, M. (1991). A cluster-number distribution and its application to the analysis of homonyms. *Japanese Journal of Applied Statistics*, **20**, 139–153 (in Japanese).
- [25] Sibuya, M. (1993). A random clustering process. *Annals of Institute of Statistical Mathematics*, **45**, 459–465.
- [26] Skinner, C.J. (1992). On identification disclosure and prediction disclosure for microdata. *Statistica Neerlandica*, **46**, 21–32.
- [27] Skinner, C.J. and Holmes, D.J. (1993). Modelling population uniqueness. in *Proceedings of the International Seminar on Statistical Confidentiality*, Dublin.
- [28] Takemura, A. (1998). Size indices data from the Japanese labor force survey. *Report of the research supported by the Grant No. 09206102 for Scientific Research of the Ministry of Education*, 95–104 (In Japanese).
- [29] Takemura, A. (1999). Some superpopulation models for estimating the number of population uniques. in *Statistical data protection - Proceedings of the conference, Lisbon, 25 to 27 March 1998 - 1999 edition*, 45–58, Office for Official Publications of the European Communities, Luxembourg.
- [30] Watterson, G.A. (1973). Models for the logarithmic species abundance distributions. *Theoretical Population Biology*, **6**, 217–250.
- [31] Willenborg, L. and de Waal, T. (1996). *Statistical Disclosure Control in Practice*, Lecture Notes in Statistics 111, Springer, New York.
- [32] Yamato, H. and Sibuya, M. (1999). Moments and asymptotic distribution of number of components of partition associated with Pitman sampling formula. (submitted).
- [33] Yamato, H., Sibuya, M. and Nomachi, T. (1999). Ordered sample from two-parameter GEM distribution. (submitted).
- [34] Zayatz, L.V. (1991). Estimation of the percent of unique population elements in a microdata file using the sample. *Statistical Research Division Report RR-91/08*, U.S. Bureau of the Census.

Variable	Case 1	Case 2	Case 3
(A) Prefecture code	(9)	(9)	(9)
(B) — zone code	824	824	824
(C) Persons 15 years old & over	8	<u>3</u>	<u>3</u>
(D) — under 15 years (male)	6	<u>3</u>	<u>2</u>
(E) — under 15 years (female)	5	<u>3</u>	(D)
(F) Relationship to the head	12	<u>5</u>	<u>5</u>
(G) Sex	2	2	2
(H) Age	100	<u>20</u>	<u>10</u>
(I) Marital status	4	4	4
K (Cell total)	1898496000	17798400	1977600

Variable	Case 4	Case 5	Case 6	Case 7
(A) Prefecture code	9	9	9	9
(B) — zone code	<u>x</u>	<u>x</u>	<u>x</u>	<u>x</u>
(C) Persons 15 years old & over	8	8	8	8
(D) — under 15 years (male)	6	6	<u>4</u>	<u>6</u>
(E) — under 15 years (female)	5	5	<u>4</u>	(D)
(F) Relationship to the head	12	12	12	12
(G) Sex	2	2	2	2
(H) Age	100	<u>20</u>	<u>20</u>	<u>20</u>
(I) Marital status	4	4	4	4
K (Cell total)	20736000	4147200	2211840	829440

Table 1: The number of categories in each variable (Case 1–7 of Section 4).

NOTE: “x” implies that the variable is suppressed. The underline implies that the number of categories is smaller than that of Case 1. ”(D)” implies that the information on the variable is represented by the variable (D).

	Case 1	Case 2
Total cell number (K)	1898496000	17798400
Total non-empty cell number (u)	25923	21851
Sample uniques (s_1)	25046	18275
Maximum cell size	28	28
Ewens parameter θ by MLE	280628.969879	52004.115657
Log likelihood (AIC)	-518.9 (1039.7)	-245.3 (492.6)
Estimated population uniques $E(S_1)$	278449.3	51928.8
Dirichlet-multi parameter γ by MLE	0.000148	0.002928
Log likelihood (AIC)	-518.9 (1039.8)	-246.3 (494.6)
Estimated population uniques $E(S_1)$	278244.1	51044.0
Pitman parameters α, θ by MLE	0.917448, 16389.753923	0.520587, 21297.598824
Log likelihood (AIC)	-111.8 (227.6)	-100.9 (205.8)
Estimated population uniques $E(S_1)$	19000174.4	1017904.0
Poisson-lognormal parameter V by MLE	10.530755	8.524957
Log likelihood (AIC)	-547296.6 (1094595.2)	-5750.6 (11503.1)
Estimated population uniques $E(S_1)$	11950655.9	1773952.1

Table 2: Case 1–2 of Section 4.

	Case 3	Case 4
Total cell number (K)	1977600	20736000
Total non-empty cell number (u)	18221	12390
Sample uniques (s_1)	12919	8049
Maximum cell size	41	54
Ewens parameter θ by MLE	24249.278863	8804.206385
Log likelihood (AIC)	-142.0 (286.1)	-686.0 (1374.0)
Estimated population uniques $E(S_1)$	24232.9	8802.0
Dirichlet-multi parameter γ by MLE	0.012424	0.000425
Log likelihood (AIC)	-143.8 (289.6)	-686.7 (1375.4)
Estimated population uniques $E(S_1)$	22427.7	8769.8
Pitman parameters α, θ by MLE	0.140768, 19948.932049	0.501239, 2585.173765
Log likelihood (AIC)	-131.4 (266.8)	-101.9 (207.7)
Estimated population uniques $E(S_1)$	57260.1	308054.4
Poisson-lognormal parameter V by MLE	5.166813	14.268244
Log likelihood (AIC)	-4235.1 (8472.2)	-14134.2 (28270.4)
Estimated population uniques $E(S_1)$	357874.6	495826.6

Table 3: Case 3–4 of Section 4.

	Case 5	Case 6
Total cell number (K)	4147200	2211840
Total non-empty cell number (u)	6657	6653
Sample uniques (s_1)	3813	3805
Maximum cell size	154	154
Ewens parameter θ by MLE	2813.718472	2810.978767
Log likelihood (AIC)	-997.1 (1996.2)	-993.1 (1988.2)
Estimated population uniques $E(S_1)$	2813.5	2810.8
Dirichlet-multi parameter γ by MLE	0.000679	0.001271
Log likelihood (AIC)	-998.6 (1999.3)	-996.0 (1993.9)
Estimated population uniques $E(S_1)$	2795.9	2778.0
Pitman parameters α, θ by MLE	0.505272, 523.377001	0.504301, 525.742679
Log likelihood (AIC)	-219.7 (443.3)	-219.7 (443.5)
Estimated population uniques $E(S_1)$	145294.2	144053.2
Poisson-lognormal parameter V by MLE	14.370145	13.053184
Log likelihood (AIC)	-13291.8 (26585.5)	-10398.8 (20799.7)
Estimated population uniques $E(S_1)$	320562.6	199167.2

Table 4: Case 5–6 of Section 4.

	Case 7
Total cell number (K)	829440
Total non-empty cell number (u)	5682
Sample uniques (s_1)	2974
Maximum cell size	154
Ewens parameter θ by MLE	2188.670938
Log likelihood (AIC)	-759.9 (1521.7)
Estimated population uniques $E(S_1)$	2188.5
Dirichlet-multi parameter γ by MLE	0.002646
Log likelihood (AIC)	-764.8 (1531.6)
Estimated population uniques $E(S_1)$	2138.9
Pitman parameters α, θ by MLE	0.443278, 524.588977
Log likelihood (AIC)	-228.2 (460.3)
Estimated population uniques $E(S_1)$	72949.3
Poisson-lognormal parameter V by MLE	9.209263
Log likelihood (AIC)	-10761.7 (21525.4)
Estimated population uniques $E(S_1)$	108974.3

Table 5: Case 7 of Section 4.

	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case7
K^*	657385	198613	71065	36594	10258	10254	7883
Log likelihood	-635.2	-157.5	-122.5	-864.9	-2153.4	-2148.4	-2103.6
AIC	1272.4	316.9	246.9	1731.9	4308.7	4298.7	4209.3
$E(S_1)$ with K^*	3894.5	431.9	1.5	5.3	0.0	0.0	0.0
$E(S_1)$ with K	33902689.2	4076872.5	104255.6	4273132.0	480176.7	147338.8	8979.4

Table 6: Poisson-lognormal model fits with optimized K (Section 4).

	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case7
s_1/u	0.97	0.84	0.71	0.65	0.57	0.57	0.52
α by MLE	0.92	0.52	0.14	0.50	0.50	0.50	0.44

Table 7: The simple estimates of the Pitman parameter α (Section 4).

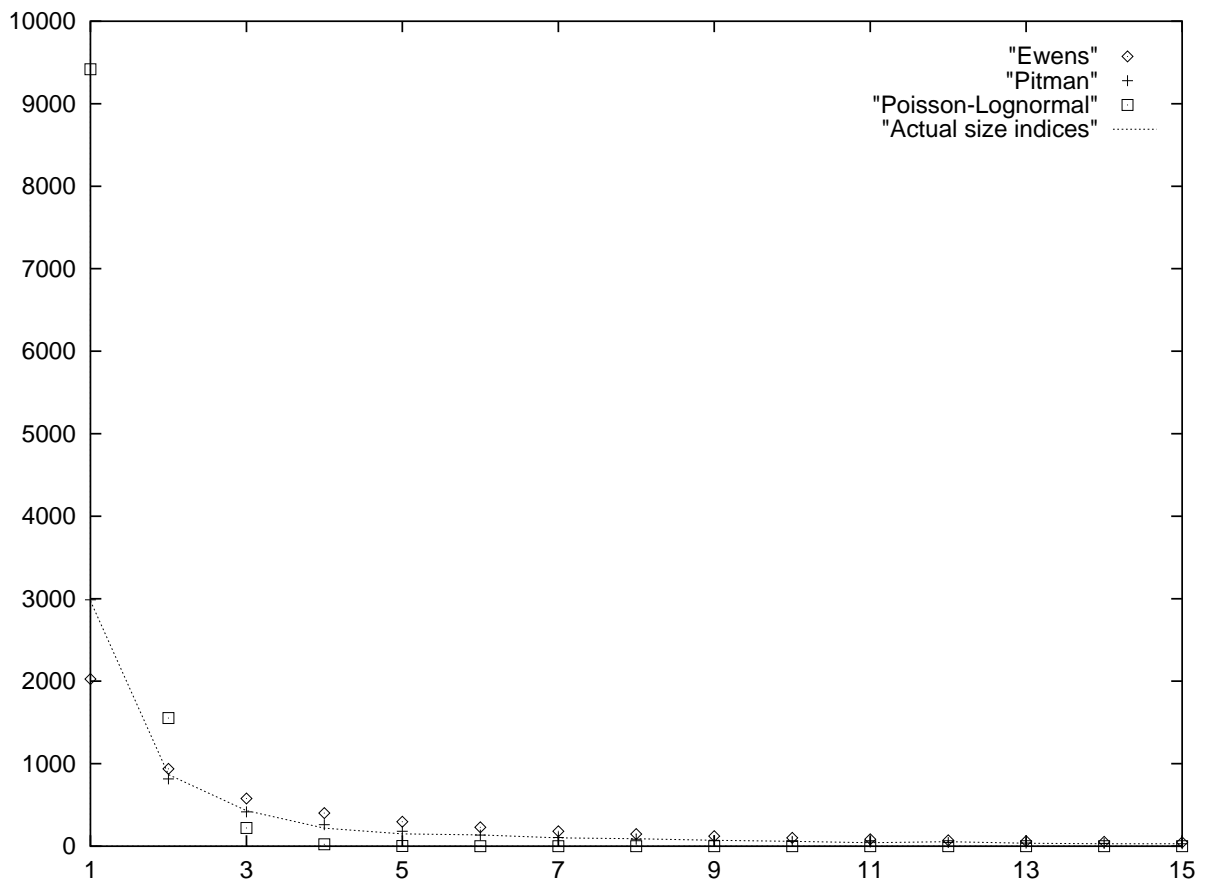


Figure 1: Fits of $E(s_i)$'s under ML estimates (Case 7, Section 4).