Contents lists available at ScienceDirect



# **Computational Statistics and Data Analysis**



journal homepage: www.elsevier.com/locate/csda

# An extension of Chao's estimator of population size based on the first three capture frequency counts

# K. Lanumteang\*, D. Böhning

Department of Mathematics and Statistics, The University of Reading, Philip Lyle Building, Reading, RG6 6BX, UK

# ARTICLE INFO

Article history: Received 20 November 2009 Received in revised form 20 January 2011 Accepted 22 January 2011 Available online 22 February 2011

Keywords: Estimation of population size Capture-recapture methods Horvitz-Thompson estimator Chao's estimator Homogeneous and heterogeneous Poisson models Negative binomial

# 1. Introduction

# ABSTRACT

A new estimator for estimating the size of an elusive target population is presented using frequency counts from capture-recapture sampling. The proposed estimator is developed by extending the idea of Chao's estimator using monotonicity of ratios of neighbouring frequency counts under a specific Poisson mixture sampling framework, the Poisson-Gamma mixture or negative binomial. The new estimator is achieved using a simple linear model on the basis of the log-ratio of neighbouring frequency counts as dependent variable which is valid under the Poisson-Gamma mixture. A simulation study is provided to study the performance of the proposed estimator under a variety of heterogeneous Poisson capture probabilities. Confidence interval estimation is done by means of an approximating normal approach and a modified bootstrap method, and was found to perform well. A variety of real data sets were also examined in order to illustrate the use of the proposed method.

© 2011 Elsevier B.V. All rights reserved.

Estimation of the size of an elusive target population is of considerable interest in several fields. For example, ecologists commonly consider how to estimate the number of species in a wildlife population. In social sciences, there is major concern about certain social problems and determining its amount in a target population such as illicit drug users, violators of a law or the number of illegal immigrants. In medicine, there is wide interest in estimating the hidden disease occurrence, the unobserved part of the disease iceberg (Woodward, 1999). In public health and epidemiology, there is the frequent problem of determining the completeness of a disease registry (e.g. Corrao et al., 2000; Gallay et al., 2000; Hook and Regal, 1995; Nardone et al., 2003).

Capture-recapture models have been ordinarily used to estimate animal abundance or population size in the ecological sciences (see, for a review, Chao and Bunge, 2002; Darroch, 1958; Eberhardt, 1969; Edwards and Eberhardt, 1967; McDonald and Palanacki, 1989; North, 1981; Pollock, 2000). The origin of capture-recapture modelling goes back to Petersen and Lincoln (Seber, 2002), who used the independent information of two identifying sources or lists to construct an estimator of population size.

Capture-recapture models currently tend to be generally applied in a variety of applications including estimation of the size of a human target population, usually defined by a specific disease experiencing potential severe undercount (e.g. Böhning et al., 2004; Corrao et al., 2000; Gallay et al., 2000; Hay et al., 2009; Hook and Regal, 1995; Nardone et al., 2003; Smit et al., 2002; van Hest et al., 2008), as well as estimation of an elusive target population in the social sciences such as illegal gun owners or car drivers without licence (e.g. Carothers, 1973; Chang et al., 1999; Hay, 1997; Hope et al., 2005; van der Heijden et al., 2003a,b).

Corresponding author. E-mail addresses: k.lanumteang@pgr.reading.ac.uk, klanumteang@hotmail.com (K. Lanumteang), d.a.w.bohning@reading.ac.uk (D. Böhning).

<sup>0167-9473/\$ -</sup> see front matter © 2011 Elsevier B.V. All rights reserved. doi:10.1016/j.csda.2011.01.017

Several estimators have been proposed to estimate the size of a target population when several identifications of the same unit are available. These include maximum likelihood methods, Zelterman's estimator (Zelterman, 1988), and Chao's lower bound estimators (Chao, 1987). However, several aspects of these estimators are critical. The maximum likelihood estimator is usually efficient only under Poisson homogeneity, whereas Chao's lower bound estimator – although developed under Poisson heterogeneity – uses only part of the available information and, hence, suffers under a lack of efficiency. To be more precise, let  $f_1$  denote the frequencies of individuals which have been identified exactly once in the capture–recapture study,  $f_2$  the number of individuals with exactly two identifications, and so forth, with *m* being the largest number of reidentifications. Then,  $n = f_1 + f_2 + \cdots + f_m$  is the size of the observed sample. Chao's estimator is given as  $\hat{N} = n + f_1^2/(2f_2)$  and it is clear from its form that it uses only part of the available information, namely the proportion  $(f_1 + f_2)/n$ .

In this paper we propose a modification of this estimator, namely  $\hat{N} = n + (3f_1f_3)/(2f_2^2) \times f_1^2/(2f_2)$  which extends the estimator of Chao by incorporating the adjustment factor  $\hat{\gamma} = (3f_1f_3)/(2f_2^2)$ . The central point of the paper is to show that this adjustment improves bias and efficiency of Chao's estimator under a wide class of models allowing heterogeneity.

#### 2. The proposed estimator

The purpose of a capture–recapture model is to provide an estimator of the population size N or, equivalently, of the frequency of unobserved individuals  $f_0$ . From the individual capture–recapture history we can determine the count X of repeated identifications per individual. Let  $f_1, f_2, f_3, \ldots, f_m$  denote the frequencies of distinct individuals identified exactly 1, 2, 3, ..., m times during the period of study, and  $f_0$  is the frequency of individuals that were never identified in the study period and hence remain unobserved. Consequently, the total number of population size N can be written as  $N = f_0 + f_1 + f_2 + \cdots + f_m = f_0 + n$ , where  $n = \sum_{j=1}^m f_j$  is the total number of distinct individuals observed. Furthermore, let  $p_0$  be the probability that an individual remains unobserved, so that  $E(f_0) = Np_0$ . Therefore, we can also write the expected population size as  $N = Np_0 + N(1 - p_0)$ . Estimating  $N(1 - p_0)$  with n leads to  $\hat{N} = \frac{n}{1-p_0}$ , the Horvitz–Thompson estimator (Horvitz and Thompson, 1952). The key issue is to estimate  $p_0$ .

Let  $p_j$  denote the probability for identifying an individual exactly j times, j = 0, 1, 2, ..., m. Under the Poisson distribution these probabilities are given as  $p_0 = e^{-\lambda}$ ,  $p_1 = e^{-\lambda}\lambda$ ,  $p_2 = \frac{e^{-\lambda}\lambda^2}{2!}$ , ...,  $p_m = \frac{e^{-\lambda}\lambda^m}{m!}$  and  $\frac{p_1}{p_0} = 2\frac{p_2}{p_1}$ . Replacing the unknown Poisson probabilities by observed frequencies provides  $f_1/f_0 = 2f_2/f_1$  as an estimating equation for  $f_0$  and Chao's estimator  $\hat{f}_0 = f_1^2/(2f_2)$  follows. However, Poisson homogeneity is rarely met in practice and it is more appropriate to incorporate heterogeneity of the identifying probability it is more reasonable to assume that the actual target population may consist of a variety of subgroups. This leads to a Poisson mixture model of the form

$$p_j = \int_0^\infty \frac{\mathrm{e}^{-\lambda} \lambda^j}{j!} f(\lambda) \mathrm{d}\lambda,\tag{1}$$

where  $f(\lambda)$  represents the heterogeneity distribution of the model parameter in the population. A prominent example for  $f(\lambda)$  is the Gamma distribution  $f(\lambda) = \lambda^{k-1} \exp(-\lambda/\theta')/(\theta'^k \Gamma(k))$  with parameters  $\theta', k > 0$ , so that  $p_j$  is Poisson–Gamma mixture, or if the marginal is worked out, the *negative binomial* distribution. Let  $r_j = \frac{jp_j}{p_{j-1}}$ , where  $p_j = \frac{\Gamma(k+j)}{\Gamma(j+1)\Gamma(k)}\theta^k(1-\theta)^j$  with  $\theta' = (1-\theta)/\theta$ , then we achieve  $r_j = (k+j-1)(1-\theta)$ . This clearly implies that there is a linear relationship  $r_j = (k-1)(1-\theta) + (1-\theta)j$  between  $r_j$  and j. Plotting  $r_j$  against j leads to the *ratio plot*, and specific patterns indicate a certain distribution, such as linearity indicates a negative binomial, a horizontal line means the presence of a Poisson distribution and a line passing through the origin indicates a geometric distribution.

To derive our estimator we consider a Taylor expansion of  $\log r_i$  around (k - 1) so that

$$\log r_j = \log(k+j-1) + \log(1-\theta) \approx \underbrace{\log(1-\theta) + \log(k-1)}_{\alpha} + \underbrace{\frac{\beta}{1-1}}_{k-1} j.$$
(2)

The motivation for the approximation (2) is as follows. Using a logarithmic transformation will guarantee that our population size estimate is feasible (which is not necessarily so when working on the  $r_j$  scale). Now, for j = 2 or j = 3 in (2) we get  $\log(r_2) = \log(\frac{2f_2}{f_1}) = \alpha + 2\beta$  and  $\log(r_3) = \log(\frac{3f_3}{f_2}) = \alpha + 3\beta$ . Solving these equations in  $\alpha$  and  $\beta$  can easily be achieved as  $\hat{\alpha} = 3\log(\frac{2f_2}{f_1}) - 2\log(\frac{3f_3}{f_2})$  and  $\hat{\beta} = \log(\frac{3f_3}{f_2}) - \log(\frac{2f_2}{f_1})$ . Then, plugging  $\hat{\alpha}$  and  $\hat{\beta}$  into (2) and using j = 1, (2) provides  $\log(r_1) = \log(\frac{f_1}{f_2}) = \alpha + \beta$ , or

$$\log\left(\frac{f_1}{f_0}\right) = 3\log\left(\frac{2f_2}{f_1}\right) - 2\log\left(\frac{3f_3}{f_2}\right) + \log\left(\frac{3f_3}{f_2}\right) - \log\left(\frac{2f_2}{f_1}\right) = 2\log\left(\frac{2f_2}{f_1}\right) - \log\left(\frac{3f_3}{f_2}\right)$$

Finally, we achieve that  $\log(f_0) = \log(f_1) - \log(\frac{4f_2^2}{f_1^2}) + \log(\frac{3f_3}{f_2}) = \log(\frac{3f_1^3f_3}{4f_2^3})$ . Hence, our estimator for  $f_0$  and N, respectively, is

$$\hat{f}_{0New} = \frac{3f_1^3 f_3}{4f_2^3} \quad \text{and} \quad \hat{N}_{New} = n + \frac{3f_1^3 f_3}{4f_2^3}.$$
 (3)

### 3. Properties of the proposed estimator

In this section we summarize some properties of the new estimator. Firstly, it should be noted that (3) is closely associated with Chao's estimator  $\hat{N}_{Chao} = n + \frac{f_1^2}{2f_2}$  (Chao, 1987) in that we can think of (3) as an adjusted Chao estimator of the form  $\hat{N}_{New} = n + \frac{f_1^2}{2f_2}\hat{\gamma}$ ; where  $\hat{\gamma} = \frac{3f_1f_3}{2f_2^2}$ . Hence, we investigate the effect of this adjustment factor.

**Theorem 1.** Under arbitrary mixing in (1) we have that

$$\lim_{N \to \infty} \frac{E(N_{New})}{N} \ge \lim_{N \to \infty} \frac{E(N_{Chao})}{N}$$

and

$$\lim_{N\to\infty} E(\hat{\gamma}) = \lim_{N\to\infty} E\left(\frac{3f_1f_3}{2f_2^2}\right) = \frac{3p_1p_3}{2p_2^2} \ge 1.$$

**Proof.** As a consequence of the Cauchy–Schwarz inequality we have for arbitrary mixing that the ratios of neighbouring count probabilities experience a monotonicity property as follows (see Chao, 1987)

$$\frac{p_1}{p_0} \le \frac{2p_2}{p_1} \le \frac{3p_3}{p_2} \le \frac{4p_4}{p_3} \le \cdots,$$

so that in particular  $\frac{2p_2}{p_1} \leq \frac{3p_3}{p_2}$ . Now,  $E(\hat{f}_{0New})/N = E(\frac{3f_1^3f_3}{4f_2^3})/N \rightarrow \frac{3}{4}(\frac{p_1^3p_3}{p_2^3}) = \frac{3}{2}(\frac{p_1p_3}{p_2})(\frac{p_1^2}{2p_2})$  and  $E(\hat{f}_{0Chao})/N \rightarrow \frac{p_1^2}{2p_2}$  for  $N \rightarrow \infty$ . It remains to show that  $\frac{3}{2}\frac{p_1p_3}{p_2^2} \geq 1$ . The latter follows from  $\frac{2p_2}{p_1} \leq \frac{3p_3}{p_2}$  which also implies the second part of the theorem, and this ends the proof.  $\Box$ 

Chao's estimator is a lower bound estimator in the sense that  $E(\hat{N}_{Chao})/N \leq 1$  for  $N \to \infty$  using that  $\frac{p_1}{p_0} \leq \frac{2p_2}{p_1}$ . Hence typically Chao's estimator will underestimate the population size. The property in Theorem 1 is remarkable since it guarantees that the asymptotically expected value of (3) is larger than that of Chao's estimator – under fairly general conditions. Next we show that (3) is asymptotically unbiased under Poisson homogeneity – as is Chao's estimator.

**Theorem 2.** Under Poisson homogeneity  $p_i = e^{-\lambda} \lambda^j / j!$  we have that

$$\lim_{N\to\infty}\frac{E(\hat{N}_{New})}{N}\to 1.$$

**Proof.**  $E(f_j/N)$  converges with  $N \to \infty$  to  $p_j$ . Hence,  $E(\frac{\hat{f}_{0New}}{N}) = E(\frac{3(f_1/N)^3(f_3/N)}{4(f_2/N)^3})$  converges to  $\frac{3p_1^3p_3}{4p_2^3} = e^{-\lambda}$ . Finally,  $E(\hat{N}_{New}/N) = E(n + \hat{f}_{0New})/N$  converges to  $(1 - e^{-\lambda}) + e^{-\lambda} = 1$  and ends the proof.  $\Box$ 

The next result compares the asymptotic biases for the new and Chao's estimator.

**Theorem 3.** Under Poisson heterogeneity according to a Gamma distribution, e.g. (1) is the negative binomial  $p_j = \frac{\Gamma(k+j)}{\Gamma(j+1)\Gamma(k)}$  $\theta^k (1-\theta)^j$  for j = 0, 1, 2, ... we have that

$$\lim_{n \to \infty} \frac{E(\hat{N}_{New})}{N} = 1 - \frac{\theta^k}{(k+1)^2}$$

and

N

$$\lim_{N \to \infty} \frac{E(\hat{N}_{Chao})}{N} = 1 - \frac{\theta^k}{k+1},$$
  
with  $1 - \frac{1}{k+1} \le 1 - \frac{1}{(k+1)^2} \le 1.$ 

2304

**Proof.** We have for  $N \to \infty$  that

$$\begin{split} E(\hat{f}_{0New})/N &= E\left(\frac{3f_1^3f_3}{4f_2^3}\right)/N \\ & \to \frac{3}{4}\left(\frac{\frac{k!^3}{(k-1)!^3}\theta^{3k}(1-\theta)^3\frac{(k+2)!^3}{3!(k-1)!}\theta^k(1-\theta)^3}{\frac{(k+1)!^3}{2!^3(k-1)!^3}\theta^{3k}(1-\theta)^6}\right) \\ & = \frac{k(k+2)}{(k+1)^2}\theta^k, \end{split}$$

so that  $E(\hat{N}_{New})/N \to (1 - \theta^k) + \frac{k(k+2)}{(k+1)^2} \theta^k = (1 - \theta^k + \frac{k(k+2)}{(k+1)^2} \theta^k) = 1 - \theta^k/(k+1)^2$ . On the other hand,

$$E(\hat{f}_{0Chao})/N = E\left(\frac{f_{1}^{2}}{2f_{2}}\right)/N$$
  

$$\rightarrow \frac{1}{2}\left(\frac{\frac{k!^{2}}{(k-1)!^{2}}\theta^{2k}(1-\theta)^{2}}{\frac{(k+1)!}{2!(k-1)!}\theta^{k}(1-\theta)^{2}}\right)$$
  

$$= \frac{k}{k+1}\theta^{k},$$

and  $\hat{N}_{Chao}/N \rightarrow (1 - \theta^k) + \frac{k}{k+1}\theta^k = (1 - \theta^k + \frac{k}{k+1}\theta^k) = 1 - \frac{\theta^k}{(k+1)}$ . The result in Theorem 3 indicates the large potential of reducing bias with the new estimator. To explore this a bit further we consider exponential mixing in (1).  $\Box$ 

**Corollary 1.** Let the mixing density  $f(\lambda)$  in (1) be the exponential, k = 1, so that the marginal (1) is the geometric. Then:

$$\lim_{N\to\infty}\frac{E(\hat{N}_{New})}{N}=1-\frac{\theta}{4}\quad and\quad \lim_{N\to\infty}\frac{E(\hat{N}_{Chao})}{N}=1-\frac{\theta}{2}.$$

The condition in Corollary 1 might appear difficult to be checked. However, exponential mixing means that the shape parameter k equals one which implies that the line in the ratio plot passes through the origin. This can be simply diagnosed and formally tested. An asymptotic unbiased Chao-type estimator for this case (k = 1) is provided as  $n + f_1^2/f_2$  and an asymptotic unbiased estimator incorporating the first three capture frequency counts is also available as  $n + f_1^3 f_3 / f_2^3$ .

Note that (3) is only well-defined as long as  $f_2$  is positive. Therefore, we suggest to use a modification of (3) which allows  $f_2 = 0$ , as follows

$$\hat{N}_{NewMo} = n + \frac{3}{4} \frac{f_1(f_1 - 1)(f_1 - 2)f_3}{(f_2 + 1)(f_2 + 2)(f_2 + 3)}.$$
(4)

The modification (4) has been suggested by an anonymous reviewer for which we are grateful. The reviewer argued that  $E(\frac{f_3f_1(f_1-1)(f_1-2)}{(f_2+1)(f_2+2)(f_2+3)}) \approx \frac{E(f_3)(E(f_1))^3}{(E(f_2))^3}$ . Indeed a small simulation confirm this results as follows. Since the original bias-corrected version of Chao's estimator  $\hat{N}_{Chao} = n + \frac{f_1^2}{2f_2}$  is  $\hat{N}_{BChao} = n + \frac{f_1(f_1-1)}{2(f_2+1)}$  we looked at a bias-corrected version (modification) of our new estimator  $\hat{N}_{New} = n + \frac{f_1^2}{2f_2} \times \frac{3f_1f_3}{2f_2^2}$  in the form  $\hat{N}_{New_{Mo}} = n + \frac{3}{4}\frac{f_1(f_1-1)}{(f_2+1)} \times \frac{(f_1-C_1)f_3}{(f_2+C_2)(f_2+C_3)}$  with constants  $(C_1, C_2, C_3)$  to be chosen. Notice that the denominator form arises to avoid a division by 0. As it turned out in our simulation the values providing the smallest bias were frequently given by the combination  $C_1 = 2$ ,  $C_2 = 2$ ,  $C_3 = 3$  (or  $C_1 = 2$ ,  $C_2 = 3$ ,  $C_3 = 2$ giving the identical form).

In addition, we consider the following truncated version of  $\hat{N}_{New}$  to improve its variance. It can be seen from Theorem 3 (and by replacing  $f_j$  by their theoretical value  $p_j$ ) that the expected value of  $\hat{\gamma} = \frac{3f_1f_3}{2f_2^2}$  approaches

$$\frac{3\Gamma(k+1)\Gamma(k+3)\Gamma(3)^{2}\Gamma(k)^{2}}{2\Gamma(2)\Gamma(k)\Gamma(4)\Gamma(k)\Gamma(k+2)^{2}} = \frac{k+2}{k+1}$$

for *N* becoming large assuming the negative binomial distribution for the count probabilities  $p_i, j = 0, ..., m$ . Note that

$$1 \le \frac{k+2}{k+1} \le 2$$

for  $0 \le k \le \infty$ . Hence, truncation at the upper and lower bound of the asymptotically expected value of the multiplier  $\hat{\gamma}$  appears reasonably and leads to an adjusted form  $\hat{N}_{New}$  as follows:

$$\hat{N}_{NewAdj} = \begin{cases} n + \frac{f_1^2}{2f_2}, & \text{if } \frac{3f_1f_3}{2f_2^2} \le 1\\ n + \frac{f_1^2}{2f_2} \left(\frac{3f_1f_3}{2f_2^2}\right), & \text{if } 1 < \frac{3f_1f_3}{2f_2^2} < 2\\ n + \frac{f_1^2}{f_2}, & \text{if } \frac{3f_1f_3}{2f_2^2} \ge 2. \end{cases}$$

$$(5)$$

When  $f_2 = 0$  we replace  $f_2$  by  $f_2 + 1$ , so that (5) is always well-defined. The adjusted form (5) can be expected to show an improved performance in terms of reducing the variance while retaining the reduction in bias, which will be seen in the simulation study section.

# 4. Variance estimator and confidence interval

# 4.1. Variance estimator

In order to investigate the variance of the proposed estimator we simply derive it by conditioning. It can be noted that the variation of  $\hat{N}_{New} = n + \frac{3f_1^3 f_3}{4f_2^3}$  is arising from two sources, the random variation of sampling *n* individuals from *N* and the random variation with respect to estimation of  $\hat{\lambda}_0$  where  $\hat{\lambda}_0 = \frac{3f_1^3 f_3}{4f_2^3}$ . Böhning (2008) provided a simple formula for variance computation of population size which can be also applied to derived the variance approximation of the new proposed estimator as follows:

$$Var_{\hat{\lambda_0}|n}(n+\hat{\lambda_0}) = \underbrace{E_n\{Var_{\hat{\lambda_0}|n}(n+\hat{\lambda_0})\}}_{[1]} + \underbrace{Var_n\{E_{\hat{\lambda_0}|n}(n+\hat{\lambda_0})\}}_{[2]}, \tag{6}$$

where  $E_n$  and  $Var_n$  refer to the marginal distribution of n and  $\hat{\lambda_0} = \frac{3}{4} \frac{f_1^3 f_3}{f_2^3}$ . Assuming that  $E_{\hat{\lambda_0}|n}(n + \hat{\lambda_0})$  in the second term [2] of (6) can be estimated by  $n + \hat{\lambda_0}$  we have that

$$Var_{n}\{E_{\hat{\lambda}_{0}|n}(n+\hat{\lambda}_{0})\} = \widehat{Var_{n}}\{n+\hat{\lambda}_{0}\} = Var_{n}\{n\} = Np_{0}(1-p_{0}).$$
(7)

Since  $\hat{p_0} = \frac{\hat{f_0}}{n+\hat{f_0}}$  and  $N(1-p_0) = n$ , (7) can be estimated by

$$\widehat{Var_n}\{E_{\hat{\lambda_0}|n}(n+\hat{\lambda_0})\} = \frac{\frac{3n}{4}f_1^3f_3}{nf_2^3 + \frac{3}{4}f_1^3f_3}.$$
(8)

Now, consider the first term in (6),  $E_n\{Var_{\hat{\lambda}_0|n}(n+\hat{\lambda}_0)\}$ , and assume again that  $E_n\{Var_{\hat{\lambda}_0|n}(n+\hat{\lambda}_0)\}$  can be estimated by  $Var_{\hat{\lambda}_0|n}(n+\hat{\lambda}_0) = Var_{\hat{\lambda}_0|n}(\frac{3}{4}\frac{f_1^3f_3}{f_2^3})$ . Using the multivariate delta-method (see Bishop et al., 1975) we are able to achieve that

$$\operatorname{Var}_{\hat{\lambda}_0|n} = \nabla g \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix}^T \operatorname{Cov} \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix} \nabla g \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix}, \tag{9}$$

where  $g(f_1, f_2, f_3) = \frac{3}{4} \frac{f_1^2 f_3}{f_2^3}$  and  $\nabla_i g(f_1, f_2, f_3) = \frac{\partial}{\partial f_i} g(f_1, f_2, f_3)$ . It is easy to see that

$$\nabla g \begin{pmatrix} f_1 \\ f_2 \\ f_3 \end{pmatrix} = \left( \frac{9}{4} \frac{f_1^2 f_3}{f_2^3} - \frac{9}{4} \frac{f_1^3 f_3}{f_2^4} - \frac{3}{4} \frac{f_1^3}{f_2^3} \right)^T.$$
(10)

Recall that the covariance matrix of the multinomial vector  $(f_1, f_2, f_3)^T$  is estimated by

$$\hat{Cov}\begin{pmatrix} f_1\\ f_2\\ f_3 \end{pmatrix} = \begin{pmatrix} f_1\left(1 - \frac{f_1}{n}\right) & -\frac{f_1f_2}{n} & -\frac{f_1f_3}{n} \\ -\frac{f_1f_2}{n} & f_2\left(1 - \frac{f_2}{n}\right) & -\frac{f_2f_3}{n} \\ -\frac{f_1f_3}{n} & -\frac{f_2f_3}{n} & f_3\left(1 - \frac{f_3}{n}\right) \end{pmatrix}.$$
(11)

Hence (9) becomes ultimately

$$Var_{\hat{\lambda}_0|n}\left(\frac{3}{4}\frac{f_1^3f_3}{f_2^3}\right) = \left(\frac{9}{4}\right)^2 \frac{f_1^5f_3^2}{f_2^6} \left\{\frac{f_1}{f_2} + 1\right\} + \left(\frac{3}{4}\right)^2 \frac{f_1^6f_3}{f_2^6} \left\{1 - \frac{f_3}{n}\right\}.$$
(12)

Substituting (8) and (12) into (6), we finally have that

$$Var_{\hat{\lambda}_0|n}\left(n+\frac{3}{4}\frac{f_1^3f_3}{f_2^3}\right) = \left(\frac{9}{4}\right)^2 \frac{f_1^5f_3^2}{f_2^6} \left\{\frac{f_1}{f_2}+1\right\} + \left(\frac{3}{4}\right)^2 \frac{f_1^6f_3}{f_2^6} \left\{1-\frac{f_3}{n}\right\} + \frac{\frac{3n}{4}f_1^3f_3}{nf_2^3+\frac{3}{4}f_1^3f_3}.$$
(13)

It is seen from (13) that the first term  $(\frac{9}{4})^2 \frac{f_1^5 f_2^3}{f_2^6} \{\frac{f_1}{f_2} + 1\} + (\frac{3}{4})^2 \frac{f_1^{16} f_3}{f_2^6} \{1 - \frac{f_3}{n}\}$  is estimating the random variation stemming from sampling *n* units from the target population and the second term  $\frac{\frac{3n}{2}f_1^3f_3}{nf_1^3 + \frac{3}{2}f_1^3f_3}$  is the approximating the random variation due to estimating the number of unobserved cases  $f_0$ .

# 4.2. Confidence interval

Once we have provided an estimator of the variance of the estimator of interest, a confidence interval of the population size N can be constructed using the normal approximation as  $\hat{N} \pm 1.96Se(\hat{N})$ , where  $Se(\hat{N})$  is the estimated standard error of  $\hat{N}$ . Alternatively, we can also investigate the confidence interval by using the bootstrap method. The main benefit of using bootstrap method is that it does not require a formula for a variance estimator and might be preferable for small sizes. We focus here on the percentile bootstrap method. The procedure for constructing 95% confidence intervals using the percentile bootstrap method is as follows:

- (1) A sample of size  $\hat{N}$  is drawn with replacement from the data set which contains both observed individuals (*n* counts of 1, 2, 3, ..., m with associated frequencies  $f_1, f_2, \ldots, f_m$ ) and estimated unobserved frequency  $\hat{f}_0$  of individuals with zero counts.  $\hat{N}$  is determined according to the estimator under investigation. We do not only bootstrap from the observed sample of size *n* because the variance of estimating *N* arises from two sources, the random variation due to drawing *n* individuals from the target population of size *N* and the random variation from estimating the parameter of interest from the observed *n* units, as just mentioned in Section 4.1 (see, for a review, van der Heijden et al., 2003a; Böhning, 2008).
- (2) Then, the resampled zero counts of individuals never identified are omitted. Then, using only the new sample of size  $n^*$ a new estimate  $\hat{N}^*$  is calculated.
- (3) Steps (1) and (2) are repeated *B* times. This provides  $\hat{N}_1^*, \hat{N}_2^*, \hat{N}_3^*, \dots, \hat{N}_B^*$ . (4) The lower and upper bound of the 95% confidence interval are calculated from  $P_{2.5}$  and  $P_{97.5}$ , the 2.5th and 97.5th percentile of the data set obtained in (3), respectively.
- (5) The standard error of estimate is now found from the sample  $\hat{N}_1^*, \hat{N}_2^*, \hat{N}_3^*, \dots, \hat{N}_n^*$ .

# 5. Real data examples and empirical applications

There exist a variety of published studies applying the ideas of capture-recapture to the estimation of the total number (adjusted for hidden cases) of patients with infectious diseases such as tuberculosis, HIV/AIDS, legionnaires disease, or malaria (e.g. Gallay et al., 2000; Nardone et al., 2003; van Hest et al., 2008). However, most studies use frequency data from multiple sources with problems of matching and potentially different target areas. Here, we illustrate the use of our proposed estimator in particular data sets with repeated identifications from only one source which is the underlying assumption to apply the model in (1). To achieve a better judgement of the proposed estimator we include the following estimators in the comparison:

• Chao: 
$$\hat{N}_{Chao} = n + f_1^2/(2f_2),$$
  
 $Var(\hat{N}_{Chao}) = \frac{1}{4}\frac{f_1^4}{f_2^3} + \frac{f_1^3}{f_2^2} + \frac{1}{2}\frac{f_1^2}{f_2} - \frac{1}{4}\frac{f_1^4}{nf_2^2} - \frac{1}{2}\frac{f_1^4}{f_2(2nf_2+f_2)}$ 

• MLE:  $\hat{N}_{MLE} = \frac{n}{1 - \exp(-\hat{\lambda})}$  where  $\hat{\lambda}$  is the maximum likelihood estimator under Poisson homogeneity,

$$\widehat{ar(\hat{N}_{MLE})} = \frac{\hat{N}_{MLE}}{(\exp(\frac{\sum if_j}{\hat{N}_{MLE}}) - \frac{\sum if_j}{\hat{N}_{MLE}} - 1)}$$

V

• Zel: 
$$N_{Zel} = \frac{n}{1 - \exp(-\frac{2f_2}{f_1})},$$
  
 $\widehat{Var(\hat{N}_{Zel})} = n(\frac{\exp(-\frac{2f_2}{f_1})}{(1 - \exp(-\frac{2f_2}{f_1}))^2})\{1 + n(\frac{\exp(-\frac{2f_2}{f_1})}{(1 - \exp(-\frac{2f_2}{f_1}))^2})(\frac{2f_2}{f_1})^2(\frac{1}{f_1} + \frac{1}{f_2})\}$ 

the latter being suggested by Zelterman (1988). We apply these estimators to studies from illicit drug use and biodiversity. We have also computed confidence intervals according to both, the approximate normal and Bootstrap method, outlined in the previous section.

#### Table 1

	The frequency of individual	contacts at Scottish needle exchange, in 1997; $n = 647$
--	-----------------------------	--

j	1	2	3	4	5	6	7	8	9	10	11	12	13	14
$f_j$	175	85	50	47	37	38	32	16	17	17	15	11	9	12
j	15	16	17	18	19	20	21	22	23	24	25	26	27	28+
$f_j$	13	7	6	2	3	5	8	2	6	1	2	3	3	25

#### Table 2

Estimated total number of Scottish drug injectors in 1997.

Method	Ñ	$\widehat{Se(\hat{N})}$	95% Cl (Approximate normal)	$Se(\hat{N})_{BT}$	95% CI (Bootstrap percentile)
MLE	648	0.67	646-650	1.00	646-649
Chao	828	34.85	759–897	36.91	763–907
New	975	137.99	704-1246	150.94	788-1379
NewAdj	975	-	_	103.76	779-1169
NewMo	948	-	_	145.78	774-1326
Zel	1042	85.25	874-1210	87.44	909-1246

#### Table 3

The frequency count of times that heroin users contacted health treatment centres in Bangkok, Thailand in 2002; n = 9, 302.

$j$ $f_j$	1	2	3	4	5	6	7	8	9	10	11
	2176	1600	1278	976	748	570	455	368	281	254	188
j	12	13	14	15	16	17	18	19	20	21	
fj	138	99	67	44	34	17	3	3	2	1	

#### Table 4

Estimated total number of heroin users in Bangkok, Thailand 2000.

Method	Ñ	$\widehat{Se(\hat{N})}$	95% CI (Approximate normal)	$Se(\hat{N})_{BT}$	95% CI (Bootstrap percentile)
MLE	9,454	12.84	9,428-9,480	13.40	9,518-9,573
Chao	10,782	80.21	10,624-10,940	85.71	10,625-10,945
New	11,714	250.16	11,223-12,205	265.07	11,256-12,279
NewAdj	11,714	-	_	249.39	11,257-12,241
NewMo	11,701	-	_	255.71	11,250-12,216
Zel	12,078	184.54	11,716–12,440	188.45	11,728–12,476

# 5.1. Drug use

#### 5.1.1. Scottish drug users

Hay and Smit (2003) provide data on drug user contact to a Scottish needle exchange programme in 1997. As the authors say

Data were collated on individuals who have visited a Scottish needle exchange in the year 1997. We prefer not to explicitly state the needle exchange from which we obtained these data; however the data were collated during a programme of drug misuse prevalence research in Scotland and was the only one operating in that area at that time. The needle exchange assigns a unique identifier number to each individual accessing the service, thus enabling it to produce statistics on the number of people who had contacted the service over a given period.

The system provided a record of the number of individuals accessing the service over the time period from January to December 1997. The number of visited drug users over this 12 months was 647 and the frequency count of contacting this treatment centre is shown in Table 1, with a maximum number of contacts of 105. Here, only the frequency count up to 28 is shown. We are able to compute all estimators of the total estimate of drug users for this data set. The result is shown in Table 2.

#### 5.1.2. Bangkok heroin users

We are interested in estimating the total number of heroin users in Bangkok (Thailand) in 2002. The data was collected by the Office of the Narcotics Control Board (ONCB), Ministry of the Prime Minister, in cooperation with the Drug Abuse Prevention and Treatment Division, Health Department and Medical Service Department, Bangkok Metropolitan Administration. The database recorded all replicated treatment contacts of drug users from the 61 health treatment centres in Bangkok metropolis. The treatment episodes for heroin users are shown in Table 3 (Source: Viwatwongkasem et al., 2008). We use this frequency table as basis for all estimators considered as provided in Table 4.

watayati	i butterny ua	ata (Fisher d	et al., 1945)										
j	1	2	3	4	5	6	7	8	9	10	11	12	13
$f_j$	118	74	44	24	29	22	20	19	20	15	12	14	6
j	14	15	16	17	18	19	20	21	22	23	24	24+	п
$f_j$	12	6	9	9	6	10	10	11	5	3	3	119	620

# Table 5 Malayan butterfly data (Ficher et al. 1943)

#### Table 6

Estimated total number of Malayan butterfly species.

Method	Ñ	$\widehat{Se(\hat{N})}$	95% CI (Approximate normal)	$Se(\hat{N})_{BT}$	95% CI (Bootstrap percentile)
MLE	624	1.80	620-628	2.23	616-624
Chao	715	22.07	671–759	24.78	672-766
New	754	63.49	629-879	85.33	659-1017
NewAdj	754	-	_	66.17	670-919
NewMo	741	-	_	62.80	664-898
Zel	868	67.04	736–1000	64.61	711–973

As can be seen in Table 4, the estimated number of heroin users from our method is between the estimators obtained using Chao's and Zelterman's methods. However, similar to the previous application, the proposed estimator shows larger variation.

# 5.2. Butterfly data

The Malayan butterfly data go back to Fisher et al. (1943) (see also Chao and Bunge, 2002) and have been frequently serving as test data for estimators under development. The frequency count of identifying distinct species is shown in Table 5. There are in 620 observed distinct species. Table 6 shows the result of estimated numbers of Malayan butterfly species.

In this example, the overall impression is that all estimators show similar results in terms of estimated numbers of species for both point and interval estimations. As expected, the MLE provides not only the smallest estimate (underestimation bias), but also gives the least variation. In contrast, our proposed estimator and Zelterman's method yields a larger estimate and variation. In addition, the new estimator provides a similar estimate for this data as the Poisson–Gamma-based estimator suggested by Chao and Bunge (2002) who also used the Malayan butterfly data to illustrate their estimator.

#### 6. Simulation study

# 6.1. Simulation scenarios

A simulation experiment is undertaken to study the performance of the proposed estimator and some competitors such as maximum likelihood, Zelterman's and Chao's estimator. The scope of study covers a variety of situations of heterogeneity in the capture probabilities. Counts have been sampled from the following distributions: negative binomial  $(k, \theta)$ ;  $k = 2, 3, \theta = 0.4, 0.6$  and  $k = 4, 5, \theta = 0.6, 0.8$ , Geometric  $(\theta)$ ;  $\theta = 0.3, 0.4, 0.5$  and two-component Poisson Mixture;  $0.5Poi(\lambda)+0.5Poi(\mu), \lambda = 0.5, 1, \mu = 2, 3, 4$ . We used a population size of 100, 1000, 10,000 and 100,000, respectively, and each scenario is repeated 10,000 times. To evaluate performance of estimation, we look at relative bias ( $RBias = \frac{E(\hat{N})-N}{N}$ ) and relative mean square error ( $RMSE = \frac{E(\hat{N}-N)^2}{N^2}$ ). Furthermore, both simulated approximation variance of the new proposed estimator and bootstrap percentile method (using a resample size of 1000) is investigated.

#### 6.2. Simulation results

We start with an illustration and show the results of estimating population size *from one sample* from a population with a known capture probability and a known population size. The artificial data set of frequency counts of identifications of distinct individuals was generated from  $f_j = E(f_j) = Np_j$ . We assume that N = 1000 and  $p_j$  corresponds to a negative binomial  $NB(4, \theta = 0.6, 0.7, 0.8)$ , a Geometric  $Geo(\theta = 0.3, 0.4, 0.5)$  and a two-component Poisson mixture  $0.5Poi(0.5) + 0.5Poi(\mu; \mu = 1, 2, 3, 4)$ . To illustrate the behaviour of the estimators a sample was generated from each of the above three scenarios, for example, for N = 1000 and NB(4, 0.7), we got  $f_0 = 240, f_1 = 288, f_2 = 216, f_3 = 130, f_4 = 68, f_5 = 33, f_6 = 15, f_7 = 6, f_8 = 3$  and  $f_9 = 1$ . Then,  $f_0$  was omitted and only the remaining zero-truncated frequencies  $f_1, f_2, \ldots, f_m$  with  $n = \sum_{j=1}^m f_j$  were used to estimate  $f_0$  and N. The results for all estimators are shown in Table 7 (see Supplementary data). It is clear that if heterogeneity becomes more pronounced our proposed estimator noticeably provides the most accurate results. However, these are the results from only one simulated sample. We now undertake a more profound simulation investigation.

### 6.2.1. Heterogeneity in identification

As a summary result it can be said that under a negative binomial the MLE and Chao's estimator show a clear underestimation of population size, whereas Zelterman, the new estimator and the adjusted form tend to overestimate for a small population size; see Table 8 in Supplementary data. It can also be seen from Table 8 that the proposed estimator and its adjusted form perform similarly in cases of large population size. In addition, the adjusted form also significantly reduces the variance if compared with its original form, in particular for small size. Furthermore, the proposed estimators show a good performance for estimating population size as does Chao's and Zelterman's estimator, in particular they provide smallest *RBias* and *RMSE* for the larger sample size. In summary, it is reasonable to state that the proposed estimators (in particular the adjusted versions) are suitable under the negative binomial distributional model.

For the case that the identification probabilities arise from the Geometric distribution, the new estimator generally shows a good performance in terms of accuracy as it gives on average the smallest *RBias* in almost all cases; see Table 9 in Supplementary data. According to *RMSE*, although the new estimator seems to be of lack of precision for the small population size, it shows excellent performance against the other methods for larger size; see Table 9.

Similar to the results under a negative binomial distribution, Zelterman and the new estimator seem to provide overestimation of population size, in contrast to the MLE and Chao's estimator which always show underestimation for all two-component Poisson Mixture scenarios; see Table 10 in Supplementary data. If large population sizes are considered under a discrete Poisson mixture, our proposed estimators not only shows high performance in terms of accuracy, but it also performs similar to the other methods in terms of precision. However, the new proposed estimator is less satisfactory for smaller size as well as it shows high variance; see Table 10.

# 6.2.2. Variance approximation

This section presents the results on variance approximation of the new estimator. We compared the variance approximation of the new estimator in Eq. (13) with estimating variance using the bootstrap and simulation methods. To define the investigated statistics in Table 11 (see Supplementary data),  $Se(\hat{N})$  denotes standard error of the new estimator computing based 10,000 repeated simulation samples, whereas  $Mean(Se(\hat{N}))$  and  $Mean(Se(\hat{N})_{Bt})$  represent an average estimated standard error from Eq. (13) and the bootstrap percentile method, respectively. It is seen from Table 11 that,  $Se(\hat{N})$ ,  $Mean(Se(\hat{N}))$  and  $Mean(Se(\hat{N})_{Bt})$  are quite similar in their values.  $Mean(Se(\hat{N}))$  is slightly smaller than  $Se(\hat{N})$  and  $Mean(Se(\hat{N})_{Bt})$ . As a result, it is reasonable to state that the variance approximation of the new estimator in Eq. (13) can be utilized to represent the true variance.

According to the bootstrap percentile confidence interval for *N*, the coverage probabilities of our proposed estimator and its modified forms are close (sometime very close) to the desired confidence level; see Table 12 in Supplementary data. We have noted previously that the newly proposed estimators  $\hat{N}_{New}$ ,  $\hat{N}_{NewMo}$ , and  $\hat{N}_{NewAdj}$  experience increased variance. This – evidently unfavourable property – turns out to be beneficial when it comes to coverage probability, as we see here. In contrast, the other estimators, in particular MLE and Chao's estimator, show lack of sufficient coverage probability due to the fact that the MLE and Chao's estimator are underestimating and yield also a small variance. In short, with respect to coverage probability in confidence interval estimation, the new estimator(s) tend to perform noticeably well.

# 7. Conclusions and discussion

A diversity of estimators in the capture-recapture field exists such as the estimators of Chao (1987) and Zelterman (1988), being widely applied in many areas of interest, especially in public health and social sciences. Here, we have introduced a new method of estimating the population size under a specific form of heterogeneity for the identification probability of distinct individuals. We have also been able to see how accurate and precise the method is performing when it is compared to other frequently used estimators. Overall, the proposed estimator is more accurate as well as providing small bias in the homogeneous Poisson case which asymptotically disappears. It is also found that the new estimator compares well with Chao's estimator since its expected value is equal to or greater than the one of Chao's estimator. Hence, it improves Chao's estimator which is known to provide a lower bound. In a simulation study, the new estimator tended to overestimate, whereas all the other methods under consideration provided the known underestimation phenomenon in almost all scenarios of heterogeneous identification probabilities. However, although the proposed estimator showed superior performance in terms of accuracy, it evidently gave also the largest variation. The reason for this increased variation is likely the use of high-order frequencies which are generally not as stable as the frequencies of the first two orders (e.g. number of singletons and doubletons). Hence, the new method has lack of precision; nonetheless, the variation of the new estimators considerably decreased for large population size (1000 and more) which is typically the case for situations of interest in surveillance and public health. In addition, the adjusted forms of the new estimator can be utilized for sample sizes below 1000 which significantly reduces the variance. We also currently view the modified estimators (4) and (5) as the better choices and the ones to be used in practice. We also provided a formula of variance approximation of the new estimator. This variance formula is not only useful to determine the efficiency of estimating, but it can be also used to construct confidence intervals. In short, the new estimator can be an alternative form of population size estimation especially for large populations and heterogeneous capturing probabilities.

# Acknowledgements

Both authors would like to thank the reviewers for the very helpful comments which lead to considerable improvement of the paper as well as to interesting and promising modifications of the proposed estimator. Both authors are grateful to the Ministry of Science & Technology, the Royal Thai Government for providing Ph.D.-funding for the first author. The authors thank Dr. Chukiat Viwatwongkasem and Busaba Suppawattanabodee for their cooperation in providing public access to the surveillance data on drug users in Bangkok Metropolis.

# Appendix. Supplementary data

Supplementary material related to this article can be found online at doi:10.1016/j.csda.2011.01.017.

# References

Bishop, Y.Y.M., Firnberg, S.E., Holland, P.W., 1975. Discrete Multivariate Analysis: Theory and Practice. MIT Press, Cambridge.

Böhning, D., 2008. A simple variance formula for population size estimators by conditioning. Statistical Methodology 5, 410-423.

Böhning, D., Suppawattanabodee, B., Kusolvisitkul, W., Viwatwongkasem, C., 2004. Estimating the number of drug users in Bangkok 2001: a capture-recapture approach using repeated entries in one list. European Journal of Epidemiology 19, 1075–1083.

Carothers, A.D., 1973. Capture-recapture methods applied to a population with known parameters. The Journal of Animal Ecology 42, 125-146.

Chang, Y.F., LaPorte, R.E., Aaron, D.J., Songer, T.J., 1999. The importance of source selection and pilot study in the capture-recapture application. Journal of Clinical Epidemiology 52, 927–928.

Chao, A., 1987. Estimating the population size for capture-recapture data with unequal catchability. Biometrics 43, 783-791.

Chao, A., Bunge, J., 2002. Estimating the number of species in a stochastics abudance model. Biometrics 58, 531–539.

Corrao, G., Bagnardi, V., Vittadini, G., Favilli, S., 2000. Capture-recapture methods to size alcohol related problems in a population. Journal of Epidemiology and Community Health 54, 603–610.

Darroch, J.N., 1958. The multiple-recapture census: I. Estimation of a closed population. Biometrika 45, 343–359.

Eberhardt, L., 1969. Population estimation from recapture frequencies. The Journal of Wildlife Management 33, 28–39.

Edwards, W.R., Eberhardt, L., 1967. Estimating cottontail abundance from livetrapping data. The Journal of Wildlife Management 31, 87–96.

Fisher, R.A., Steven, A., Williams, C.B., 1943. The relation between the number of species and the number of individuals in a random smaple of an animal population. Journal of Animal Ecology 12, 42–58.

Gallay, A., Vaillant, V., Bouvet, P., Grimont, P., Desenclos, J.C., 2000. How many foodborne outbreaks of salmonella infection occurred in France in 1955?: application of the capture-recapture method to three surveillance systems. American Journal of Epidemiology 152, 171–177.

Hay, G., 1997. The selection from multiple data source in epidemiological capture-recapture studies. The Statistician 46, 515-520.

Hay, G., Gannon, M., MacDougall, J., Eastwood, C., Williams, K., Millar, T., 2009. Capture-recapture and anchored prevalence estimation of injecting drug users in England: national and regional estimates. Statitical Methods in Medical Research 18, 323–339.

Hay, G., Smit, F., 2003. Estimating the number of drug injectors from needle exchange data. Addiction Research and Theory 11, 235-243.

Hook, E.B., Regal, R., 1995. Capture-recapture methods in epidemiology: methods and limitations. Epidemiologic Reviews 17, 243–264.

Hope, V.D., Hickman, M., Tilling, K., 2005. Capturing crack cocaine use: estimating the prevalence of crack cocaine use in London using capture-recapture with covariates. Addiction 100, 1701–1708.

Horvitz, D.G., Thompson, D.J., 1952. A generalization of sampling without replacement from a finite universe. Journal of the American Statistical Association 47, 663–685.

McDonald, J.F., Palanacki, D., 1989. Interval estimation of the size of a small population from a mark-recapture experiment. Biometrics 45, 1223-1231.

Nardone, A., Decludt, B., Jarraud, S., Etienne, J., Hubert, B., Infuso, A., Gallay, A., Desenclos, J.C., 2003. Repeat capture-recapture studies as part of the evaluation of the surveillance of legionnaires's disease in France. Epidemiology and Infection 131, 647–654.

North, P.M., 1981. Models for a simple capture-recapture procedure to estimate population size. Biometrics 37, 661–672.

Pollock, K.H., 2000. Capture-recapture models. Journal of the American Statistical Association 95, 293-296.

Seber, G.A.F., 2002. The Estimation of Animal Abundance and Related Parameter, 2nd ed. Blackburn Press, Caldwell, NI.

Smit, F., Reinking, D., Reijerse, M., 2002. Estimating the number of people eligible for health service use. Evaluation and Program Planning 25, 101–105.

van der Heijden, P.G.M., Bustami, R., Cruyff, M.J.L.F., Engbersen, G., van Houwelingen, H.C., 2003a. Point and interval estimation of the population size using the truncated poisson regression model. Statistical Modelling 3, 305–322.

van der Heijden, P.G.M., Cruyff, M.J.L.F., van Houwelingen, H.C., 2003b. Estimating the size of a criminal population from police records using the truncated poisson regression model. Statistica Neerlandica 57, 289–304.

van Hest, N.A.H., Grant, A.D., Smit, F., Story, A., Richardus, J.H., 2008. Estimating infectious diseases incidence: validity of capture-recapture analysis and truncated models for incomplete count data. Epidemiology and Infection 136, 14–22.

Viwatwongkasem, C., Kuhnert, R., Satitvipawee, P., 2008. A comparison of population size estimators under the truncated count model with and without allowance for contaminations. Biometrical Journal 50, 1006–1021.

Woodward, M., 1999. Epidemiology: Study Design and Data Analysis. Chapman and Hall, CRC, London, Boca Raton, FL.

Zelterman, D., 1988. Robust estimation in truncated discrete distributions with application to capture–recapture experiments. Journal of Statistical Planning and Inference 18, 225–237.