



ESTIMATING UNKNOWN POPULATION SIZES: TECHNIQUES AND APPROACHES

Mohammad Saleh Al-Yousef and Mohammad Faisal Al-Saleh

Department of Statistics, Yarmouk University, Jordan

Currently at: Department of Mathematics and Statistics, Jordan University of Science and Technology, Jordan

Abstract: Statistics, as the science of drawing meaningful conclusions about populations from representative samples, employs various sampling techniques, including simple, stratified, systematic, and cluster random sampling. The choice of sampling method hinges on research objectives, available population information, and budget constraints. From these samples, valuable insights into population parameters are derived, with a focus on mean (μ), proportion (p), variance (σ^2), and total (τ). Additionally, the population size (N) serves as a crucial yet less common parameter.

In the context of finite populations, estimating these parameters and their variances relies on a known value for N . However, when N is unknown, it becomes necessary to estimate it beforehand to accurately assess population totals and estimator variances.

This paper delves into the intricacies of statistics, emphasizing the significance of appropriate sampling techniques and parameter estimation, particularly in scenarios where the population size remains uncertain.

Keywords: Statistics, sampling techniques, parameter estimation, population parameters, population size.

estimated.

Assume that we have a population of size N (known). Let o_1, \dots, o_N be the population measurements.

The population mean, total and variance are:

N

$$\text{Population Mean: } \bar{o} = \frac{1}{N} \sum_{i=1}^N o_i$$

1. Introduction and Some Related Topics

Statistics is the science of making inference about a population using the information contained in a sample selected from it. The sample can be chosen by one of the techniques such as simple, stratified, systematic, cluster random sampling, etc. The choice of the technique depends on the objectives of the study, information available about the population of interest and the budget. The information obtained from the chosen sample is used to estimate the population parameters. For finite population, the main population parameters are the mean (\bar{o}), proportion (p), variance (σ^2), total (τ). Less common parameter is the population size (N). Estimation of the parameters and their variances depend on N , which is usually known. If N is unknown, then it has to be estimated first so that the population total and the variances of the estimators can be

$\square \quad i \square 1 \quad N$

A simple random sample, SRS, of size n is a sample obtained from the population in such a way so that all possible samples of size n have equal chance of being the chosen sample, i.e. $P(\text{a subset of size } n \text{ from a population of size } N \text{ is the chosen sample}) = 1/\binom{N}{n}$. Let x_1, \dots, x_n be a SRS of size n from this population.

The usual estimators of the population mean & to talare, respectively:

$$\frac{\sum x_i}{n} \quad \bar{x}, \quad \frac{\sum x_i^2}{n} \quad s^2.$$

It is well known that \bar{x} & s^2 are unbiased estimators of μ & σ^2 , respectively. Their variances are:

$$\frac{2\sigma^2}{n}$$

$$Var(\bar{x}) = \frac{\sigma^2}{n} \quad \text{and} \quad Var(s^2) = \frac{2\sigma^4}{n(n-1)}.$$

The main concern in this work is the estimation of the population total μ when N is unknown. If N is known, then some estimators of N can be used as a guard against unsuitable estimates of μ . Capture-Recapture technique is the main method used to estimate N . There are two main procedures of this technique; Capture- Recapture with Direct Sampling and Capture- Recapture with Indirect (or Inverse) Sampling:

$\square \quad \text{Direct Sampling}$

The direct Capture- Recapture sampling, known as *Petersen's method*, goes back to 1894. Assume that there is a closed population with equal chance of each member to be selected. A random sample of m elements is drawn, tagged (marked) and then released back into the population. Then, after some enough period of time necessary for the marked units to mix with the remaining elements of the population, a second random sample of size n is drawn. Let T be the number of recaptured elements in the second sample. Then, the Petersen estimator of the population size is:

$$\hat{N}_P = \frac{nm}{T}$$

$$\bar{T}$$

An approximate estimator of the variance of \hat{N}_P is (Sekar and Deming, 1949):

$$Var(\hat{N}_P) \approx mn(m-1)(n-1)(n-2)/[T(T-1)(T-2)]$$

$$\bar{T}$$

Actually, \hat{N}_P is the maximum likelihood estimator (MLE) and also the method of moments estimator (MME) of N . A modified estimator of N was proposed by Chapman (1951):

$$\hat{N}_C = \frac{(m+1)(n+1)}{T+1} - 1,$$

with variance estimated by $Var(\hat{N}_C) \approx (m-1)((Tn-1)-1)(2m(T-1)(T-2))(n-1)$.

$$\bar{k}$$

(Scheaffer et al., 1995). This estimator has the advantage of being valid even when $T=0$.

$\square \quad \text{Inverse Sampling (Indirect Sampling)}$

Inverse sampling is another method for estimating N . In this method, a random sample of size m is chosen, marked and released. Later, elements are randomly selected from the population until k (fixed in advance) elements are being recaptured, then

$$\hat{N} = Tm/k$$

$$\bar{k}$$

Where, T^{\square} is the total number of elements selected in the second random sample to obtain k previously captured elements. The variance of N^{\wedge} is estimated by

$$m^2 T^{\square} (T^{\square} k) \\ Var^{\wedge} (N^{\wedge}) \square k 2(k \square 1) , \text{ (Scheaffer et al., 1995).}$$

□

Capture-Recapture technique is an old method used to estimate the size of fish and wildlife population. Later on, the method was used for estimating other population sizes. Azevedo-Silva et al. (2009) analyzed the number of cases and incidence of childhood acute lymphoblastic leukemia by using two source capture-recapture procedures in three different cities in Brazil. Estimating of birth and death rates in India was considered by SeKar and Deming (1949). Estimating the population size of Injecting Drug Users (IDU) was discussed by Luan et al.

(2005). Estimating the number of people eligible for health service was studied by Smith et al. (2002). In their graduation project, Mohammad and Abdullah (2007) compared some Capture-Recapture techniques and cluster sampling for estimating the total number of times the word "Allah" "الله" appears in the Holy Quran. For more details about the estimation of population total and size, see also Gutierrez and Breidt (2009), Otieno et al. (2005), and Arnab (2004).

In many situations, the ratio estimator is used to estimate \square of the variable of interest for a population of size N (unknown). One way to overcome the difficulty of not knowing N is to use a suitable auxiliary variable. Let O_1, \dots, O_N be the population measurements of the main variable of interest (O) and the corresponding values of an auxiliary variable (V) be V_1, \dots, V_N ; the population measurements are $(O_1, V_1), \dots, (O_n, V_n)$. Assume that there is a fair degree of association between O & V . Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be the elements of a SRS, from this population. Now, using the relation

$$\square o \square \square o,$$

$$\square V \square V$$

X

we can estimate $\square o$ by $\square \hat{o} \square \square v$. $Y X$

Let $r \square -$, an estimate of the variance of $\square \hat{o}$ is given by:

$$Y \square =$$

n

N

$$Var^{\wedge} (\square \hat{o}) \square \square v^2 \square n 1 Sr^2 , \text{ where } Sr^2 \square \square i \square 1 (Y_i \square r X_i) \square 2 , \text{ (Scheaffer et al., 1995).}$$

$$\overline{N} n \square \overline{V}^2 \square n \square 1$$

Ahmad et al. (2000) introduced another method to estimate the population total and the population size. They used sequential sampling with replacement until a fixed (k) members are repeated.

In this paper, the estimation of population total (\square) utilizing estimators of the population size is considered. In Section 2, we consider the estimation of the population total when N is unknown using Direct Sampling. Two estimators are suggested for \square ; one is based on Chapman estimator of N and the other is based on a suggested modified estimator of N . In Section 3, \square is estimated using Indirect Sampling. The suggested estimators are compared. Concluding remarks and suggested future works are outlined in Section 4.

2. Estimation of \square When N is Unknown Using Capture-Recapture- Direct Sampling

Assume that we have a closed population with an equal chance of each member to be selected in a random sample. In Direct Sampling, a SRS of n_1 elements is drawn, marked(tagged) and the value of the random variable(r.v.) of interest(Y) is noted for each element, the n_1 items are released back into the population. After waiting enough period of time so that the marked elements mixed with the

remaining population elements(this necessary when the population is a mobile one), a second SRS of n_2 elements is drawn. Let T be the number of recaptured elements in the second sample. The values of the variable Y are noted for each of the $n_2 \square T$ elements. From the first and second sample we obtain a net random sample of size n , where

$$n \square n_1 \square n_2 \square T. \quad (2.1)$$

Note that n is a r.v. (not fixed). T has a hypergeometric distribution with probability function:

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$f(t) = P(T = t) = \frac{\binom{n_1}{t} \binom{n_2}{t} \binom{N-n}{N-t}}{\binom{N}{t}}, \quad t = 0, 1, 2, \dots, \min(n_1, n_2).$$

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

Actually, the smallest value of t is $\max(0, n_1 - n_2 + N)$, but in practice N is large and $\max(0, n_1 - n_2 + N) = 0$. Now, using the properties of hypergeometric distribution we have:

nn

$$E(n) = E(n_1 + n_2 - T) = n_1 + n_2 - \frac{1}{N}n^2 \quad (2.2)$$

$$Var(n) = Var(n_1 + n_2 - T) = n_2 \frac{N-n_1}{N} \frac{N-n_2}{N} \frac{N-n}{N}$$

$$(2.3)$$

Petersen estimator of the population size (N) is:

$$\hat{N}_P = \frac{n_1 n_2}{T}, \quad T = 0, 1, \dots, \min(n_1, n_2). \quad (2.4)$$

Note that T may equal zero with positive probability; in this case, \hat{N}_P is ∞ . To overcome this difficulty, an alternative estimator of N was proposed by Chapman (1951) as:

$$\hat{N}_C = \frac{(n_1 + 1)(n_2 + 1)}{T + 1} \quad (2.5)$$

Now,

$$\hat{N}_C = \frac{(n_1 + 1)(n_2 + 1)}{T + 1} = \frac{(n_1 + 1)(n_2 + 1)}{\min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - T + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1}$$

$$\begin{aligned} & \frac{T + 1}{\min(n_1, n_2) + 1} = \frac{T + 1}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{T + 1}{n_1 + n_2 - \min(n_1, n_2) + 1} \\ & \frac{n_1 + n_2 - T + 1}{\min(n_1, n_2) + 1} = \frac{n_1 + n_2 - T + 1}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{n_1 + n_2 - T + 1}{n_1 + n_2 - \min(n_1, n_2) + 1} \\ & \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - T + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} \\ & \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} \\ & \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} = \frac{(n_1 + 1)(n_2 + 1)}{n_1 + n_2 - \min(n_1, n_2) + 1} \end{aligned} \quad (2.6)$$

Clearly, \hat{N}_C is negatively biased. Also,

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$\frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N}$$

$$Var(\hat{N}_C) = E(\hat{N}_C^2) - E(\hat{N}_C)^2 = \frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N} \frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N} \frac{n_1}{N} \frac{n_2}{N} \frac{N-n}{N} \quad (2.7)$$

Using the two estimators (2.4) and (2.5), we suggest the following new estimator of N :

$$\hat{N}_S = \begin{cases} \hat{N}_P & \text{if } T = 0 \\ \hat{N}_C & \text{if } T \neq 0 \end{cases} \quad (2.8)$$

$$\hat{N}_S = \begin{cases} \hat{N}_P & \text{if } T = 0 \\ \hat{N}_C & \text{if } T \neq 0 \end{cases}$$

Now,

$$\hat{E}(N_S) = \frac{1}{n_1 n_2} \min(\lceil n \rceil, n_1, n_2) \cdot \frac{1}{n_1 n_2} \sum_{t=1}^{\min(n_1, n_2)} P(T=t) \cdot \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} (n_1 - i + 1)(n_2 - j + 1) P(T=t)$$

But,

2

$\neg \Box N^s \square E(N^s) \square \square \square \square N^p I(T \square o) \square E(N^s) I(T \square o) \square \square \square N^c I(T \square o) \square E(N^s) I(T \square o) \square \square \square N^p I(T \square o) \square E(N^s) I(T \square o) \square^2 \square \square N^c I(T \square o) \square E(N^s) I(T \square o) \square^2 \square \square N^p \square E(N^s) \square^2 I(T \square o) \square \square N^c \square E(N^s) \square^2 I(T \square o).$

Thus,

□ The above results are summarized in the following lemma:

Lemma (2.1)

The expected value and variance of the suggested estimator, N_s , of the population size N are given by
 $E(N_s) = n_1 N / (n_1 + 1)$ $V(N_s) = n_1 N^2 / ((n_1 + 1)^2 (n_1 + 2))$

$\min(n_1, n_2) \leq N^{\frac{1}{2}} S^{\frac{1}{2}} \sqrt{n_1 n_2} \leq \sqrt{1 + 2} \sqrt{n_1 n_2} \leq \sqrt{n_1 n_2} \leq t \sqrt{n_1 n_2} \leq N^{\frac{1}{2}} n_1 n_2 \leq t \sqrt{n_1 n_2} \leq (n_1 + 1)(n_2 + 1) \leq 1 + \sqrt{n_1 n_2} \leq nnN^{22} \leq n^2 N^{22}.$

t_1 n_2 n_2

Given n , let Y_1, Y_2, \dots, Y_n be the values of the variable Y for the sample elements. The suggested estimators of the population total (\square) are: –

$$\hat{C} \hat{N} CY, \quad (2.11)$$

and

$$\square^* S \square N^* SY, \quad (2.12)$$

We conjecture here that given n , Y_1, Y_2, \dots, Y_n is a SRS from the population. We have not been able to prove this conjecture yet.

Now,

$$E(\bar{N}^s) = E(E(\bar{N}^s | n)) = E(E(N^s Y | n)) = E(\bar{N}^s E(Y | n)) = E(N^s)$$

$$= \frac{n_1}{N} N = n_1$$

$$\therefore \bar{N}^s = n_1$$

\square^t

$$n_2 \min(n_1, n_2) = n_2 \min(n_1, n_2) \quad (2.13)$$

$$\bar{N}^s = n_1 \quad \therefore n_2 = n_2 \quad \therefore n_2 = n_2$$

\square

Similarly, for \bar{N}^c we have

$$(n_2 - 1) \bar{N}^c = n_2 - 1$$

$$E(\bar{N}^c) = E(N^c) = N - 1$$

$$= (n_2 - 1) - n_2 = -1$$

(2.14)

$$\square N$$

$$= (n_2 - 1) - n_2 = -1$$

$$\square N$$

$$= (n_2 - 1) - n_2 = -1$$

$$\square$$

1 2 $t \square 1 \square \square \square \square t2(n1 \square n2 \square t) \square \square \square \square nN2 \square \square \square \square$ $\square \square \square \square$ $n1 \square n2$ \square
 $\square \square \square nN2 \square \square \square \square$

Thus,

\square \square $\square n_1 \square \square N \square n_1 \square \square$ $\square N \square n_1 \square$

1 $\boxed{} \quad \boxed{} \quad \boxed{t} \boxed{} \boxed{} \boxed{} \boxed{n} \boxed{t} \boxed{} \boxed{} \boxed{} \boxed{} N \boxed{(n \boxed{1})} (n \boxed{1}) \boxed{1} \boxed{} \boxed{} \boxed{n} \quad \boxed{} \boxed{} \quad \boxed{} \quad Var(N^{\wedge})$
 $\boxed{} \boxed{} E(N^{\wedge}) \boxed{} \boxed{}$

$\square \quad \square$
 $Var(\hat{S}) \leq N \leq 2(1 + \min(tn_{11}, n_2) + t(n_1 + n_2)t) \leq N_2 \leq 1 + n_1 + n_2 \leq N_2 \leq S^2 Var(N).$

$$\overbrace{\square \square \quad \square \quad \square} n_2 \square \square \quad \square \square \quad \square \square n_2 \square \square$$

(2.15)

Similarly for $\square^{\wedge}c$, we have

The above results are given in the following lemma:

Lemma (2.2)

The expected value and the variance of the estimators of the population total are given by

$$E(\hat{N}_S) = E(N_S), E(\hat{N}_C) = E(N_C)$$

$$Var(\bar{X}) = \frac{N^2}{n} \left(\frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j=1, j \neq i}^n (x_{ij} - \bar{x}_i)^2 \right) = \frac{N^2}{n} \left(\frac{1}{n(n-1)} \sum_{i=1}^n n(n-1) \bar{s}_{i-}^2 \right) = \frac{N^2}{n} \bar{s}_{-}^2$$

$$S = \boxed{} \quad 1 \quad \boxed{0} \quad \boxed{0} t_2(n_1 \boxed{} n_2 \boxed{} t) \boxed{0} \quad nN \boxed{0} \quad \boxed{0} \quad n_1 \boxed{} n_2 \quad \boxed{0} \quad nN_2 \quad \boxed{0}$$

[View Details](#) [Edit](#) [Delete](#)

t_1

□ □ □ □ □ □ □ □ □ □

| □ 2 □ □

$$\square \text{Var}(N^\wedge s) \square \square E(N^\wedge s) \square^2 \square \square \square \square^2 \text{Var}(N^\wedge s)$$

$$N^2 \leq \min(n_1, n_2) N^2 \leq nt_1 N^2 \leq n_2 nt_1 N^2 \leq Var(N^C) E(N^C)^2 \leq 2Var(N^C).$$

$$Var(\hat{C}) = \frac{n}{N} \left(\frac{1}{n} \sum_{i=1}^n \hat{C}_i - \bar{\hat{C}} \right)^2$$

$$Var(\square_C) = \frac{N-n}{N}$$

— □□ *to* □□ □□□ □ □

Now, if N is known, then \hat{N} can be estimated based on SRS of size n : $n_1 \hat{n}_2 \hat{T}$ by $\hat{N} = NY$, (2.17)

with

$$E(\Box^{\wedge}) \Box E(E(\Box^{\wedge} \mid n)) \Box N \Box \Box \Box.$$

Thus, $\hat{\mu}$ is an unbiased estimator of μ .

$$Var(\square^{\wedge}) \square Var(NY) \square E(Var(NY | n)) \square Var(E(NY | \bar{n}))$$

E \square \square \square $N2$ \square $n2$ \square \square \square NN \square \square $1n$ \square \square \square \square \square \square \square Var \square N \square \square \square $NN2$ \square \square 12

□ □ □ □ *NE* □ □ □ *In* □ □ □ □ 1 □ □ □ □

$$\begin{aligned} & \square \quad \square n_1 \square \square N \square n_1 \square \quad \square \\ & N 2 \square 2 \square \min(n, n) \square \quad 1 \quad \square \square \\ & \square N \square 1 \square \square \square N \square t \square 10 \square 2 \square \square \square \end{aligned}$$

$$\boxed{} \boxed{} n_1 \boxed{} n_2 \boxed{} t \boxed{} \quad \boxed{} nN_2 \boxed{} \quad \boxed{} \quad \boxed{}$$

$$K^2.$$

where,

where,

A horizontal row containing two small, solid black squares. They are positioned side-by-side with a small gap between them.

1

The efficiency of \square^S with respect to \square^* is

$$Eff(\hat{\square_s}; \hat{\square}) \leq MSE(\hat{\square}),$$

$$\overline{MSE(\square^{\wedge}_S)}$$

where,

where,
 $MSE(\hat{\square}) \equiv Var(\hat{\square}), MSE(\hat{\square}_S) \equiv bias(\hat{N}_S)^2 + Var(\hat{\square}_S).$

Now, let $\text{Var}(\hat{S}) = L^2 \text{Var}(N^{\wedge} S) / 2$, where

$$\begin{aligned} \text{Now, let } \text{Var}(\square B) = L \square 2 \square \text{Var}(B \square S) \square 2, \text{ where} \\ \square & \quad 2 \quad \min(n, n') L \square \quad \square \square Nn \quad t \square 11 \ 2 \ \square \square t \ 2 (n1 \square n2 \square t) \quad \square \quad t) \square \square \square \square \quad N \\ & \quad 1 \ \square 1n)(1n \square 2 \ n \square 21) \square 1 \square 2 \ P(T \ \square \ o) \square \square \square \\ 1 & \quad 1 \ n22 \ \square \ \square \ 1 \quad \quad P(T \ \square \ \square \ (n \end{aligned}$$

1

then,

$$Eff(\square^{\wedge}S; \square^{\wedge}) \leq \overline{MSE(N^{\wedge})} \leq \overline{MSE(N^{\wedge}S)} \leq \overline{L}.$$

The efficiency can be rewritten in terms of the coefficient of variation (CV) given by $CV = \sqrt{\frac{S^2}{\bar{X}^2}}$ as

$$\begin{aligned} K & \quad 2 \\ \square CV(y) \square & \\ MSE(N^s) & \quad (2.19) \\ Eff(\square^s; \square^c) \square & \end{aligned}$$

$$\begin{aligned} L & \quad 2 \\ 1 \square & \quad \square CV(y) \square \\ MSE(N^s) & \end{aligned}$$

The efficiency of \square^s and \square^c w.r.t. \square^c are given in Tables (2.1) and (2.2) respectively. Also, the efficiency

of \square^s w.r.t. \square^c is given in Table (2.3). Based on these tables, we can see that \square^s is more efficient than \square^c for small expected sample size. \square^s & \square^c are more efficient than \square^c when $E(n)$ is small and for large CV .

3. Estimation of \square When N is Unknown Using Capture-Recapture- Indirect Sampling

Indirect sampling is another Capture -Recapture method for estimating \square . In this method, sampling continues until a fixed number (T) of recaptured elements are obtained. So, a random sample of size n_1 is chosen, marked and released. Later, we select elements randomly from the population until T elements are being recaptured. Let n_2 be the total number of elements selected in the second random sample to obtain T previously captured elements, then from the first and the second sample we obtain a random sample of size n elements, where

$$n \square n_1 \square n_2 \square T. \quad (3.1)$$

Here, n_2 is a random variable (not fixed); it has the negative hypergeometric distribution with probability function given by:

$$\begin{aligned} \square n_2 \square 1 \square \square N \square n_2 \square \\ \square \square T \square 1 \square \square \square \square \square \square n_1 \square T \square \square \square, \quad n \end{aligned}$$

$$\frac{\square}{\square} f(n_2) \square \quad 2 \square T, T \square 1, \dots, N \square n_1 \square T. \quad (3.2)$$

$$\begin{aligned} \square N \square \\ \square \square \square n_1 \square \square \square \end{aligned}$$

Also, $E(n_2) \square T(N \square 1) \quad (\text{Balakrishnan 2003}) \quad n_1 \square 1$

and $Var(n_2) \square T(N \square (1n)(N1 \square) 2n(1n)(1n \square 12 \square) 1 \square T) \quad (\text{Khan, 1994}).$

$$\begin{aligned} \square \\ \text{Now,} \\ E(n) \square E(n_1 \square n_2 \square T) \square n_1 \square T \square T(nN_1 \square \square 1 1) \square n_1 \square T \square \square \square Nn_1 \square \square n_{11} \square \square \square \square \quad (3.3) \end{aligned}$$

$$\begin{aligned} \square \\ \text{Also,} \\ Var(n) \square Var(n_2) \square T(N \square (1n)(1 \square N1 \square) 2n(1n)(1n \square 12 \square) 1 \square T) \quad (3.4) \end{aligned}$$

Table (3.1) contains the expected value of the sample size for different values of n_1, T, N . It can be seen that $E(n)$ is increasing in T for fixed n_1 and decreasing in n_1 for fixed T .

The estimator of the population size N is $N_I \hat{\square} n^I T^{\frac{1}{2}}$, $1 \square T \square n_1, T$ is integer. (3.5)

Now,

$$E(N_I) \square E \square n^I n^{\frac{1}{2}} \square \square n^I (N \square 1) \quad (3.6)$$

$$\square T \square \quad n_1 \square 1$$

Clearly, $E(N_I)$ does not depend on T and increasing in n_1 , N_I is negatively biased. Now,

$$\begin{aligned} & \text{Var}(N^{\wedge} I) \square \text{Var} \square \square n_1 n_2 \square \square \square T n \underline{12} 2 \text{Var}(n2) \\ & \square T \square \\ & n_1^2 (N \square 1)(N \square n_1)(n_1 \square 1 \square T) \end{aligned}$$

$$\frac{\square \quad T \quad (n_1 \square 1 \quad 2(n_1 \square 2))}{\square} \quad . \quad (3.7)$$

)

Hence,

$$1 \text{ } MSE(N^*) \square (n_1 \square 2 \square \square n_{T^{12}}(N \square 1)(N_{(n} \square_1 n_{\square)})(_2)n^1 \square 1 \square T) \square \square n_1 \square N \square^2 \square \square \square . \quad (3.8)$$

I

1) □

It can be seen from (3.6) that this estimator of N can be corrected to be unbiased estimator of N as follows:

$$E(N^+ I) \quad \square n_1(N^- \square 1) \quad , n_1 \square 1$$

which gives

$$E(N^\wedge \ln^1 \square^1) \square N \square 1.$$

n 1

So, the estimator

$$N^{\wedge} \square \square \square n_1 n \square_1 1 \square \square \square \square \square \square_1 \square n_2 (n T_1 \square 1) \square_1$$

is an unbiased estimator of N . Let

$$N_I^* \square n^2(n^1 \square 1) \quad \square 1, \quad (3.9) T$$

then

$$E(N^{\wedge I^*}) \square N.$$

Therefore, the mean square error of $N^* I^*$ is

(

$$MSE(N^* I^*) \square Var(N^* I^*) \square n_{1T} \square_{21)2} Var(n_2) \square (N \square_1)(NT(\square n_1 n_1 \square)(2n)_1 \square_1 \square T). \quad (3.10)$$

2

Clearly, $\text{Var}(N^I) \leq n_1 n_1 \text{Var}(N^{I^*})$, thus, for any values of N, n_1 and T , $\text{Var}(N^I) \leq \text{Var}(N^{I^*})$, however,

the $MSE(\hat{N}^I)$ is not necessarily less than $MSE(\hat{N}^I^*)$.

The efficiency of \hat{N}_I w.r.t. \hat{N}_I^* is:

The efficiency of IV T.w.r.t. IV T is:

Table (3.2) contains some numerical values of the efficiency of \hat{N}_I w.r.t. \hat{N}_{I^*} . Clearly, the two estimators almost have the same performance.

Given T , Y_1, Y_2, \dots, Y_n are the values of variable for the sample element. Two estimators of the population total \square are: $\square \hat{Y}^I \square N^I Y$ and $\square \hat{Y}^{II} \square N^{II} Y$

Now,

$$\begin{array}{c} E\Box\Box^I \Box\Box E\Box\Box^m n_2\Box\Box\Box | E\Box E\Box N^I Y n_2\Box\Box\Box E\Box N^I E\Box Y n_2\Box\Box \\ n(N) \\ \Box\Box E\Box N^I \Box\Box 1 \quad \Box 1) \Box. \end{array}$$

Also, $E(\bar{x}^i \bar{r}^*) = N\bar{x}\bar{r}$. The variance of \bar{x}^i and \bar{r}^* can be derived as follow

$$n_{12} \square_2 E \square \square n_{22}(N \square (n_1 \square n_2 \square T)) \square \square$$

$T^2(N \square 1) \square \square n_1 \square n_2 \square T \square \square \square \square 2Var \square N^\wedge I \square$
 $\square n_2 \square 1 \square \square N \square n_2 \square \square$

$$\square B \square^2 \square \square^2 Var(\hat{N}_I), \quad (3.11)$$

where,

n_2 1 N n_2

$$\overline{B \square n_{12} N \square \square 2n \square 1 \square T \square \square n_{22}(Nn \square 1 \square (nn_{12} \square \square nT_2 \square T)) \square \square T \square 1 \square \square \square N \square \underline{n} \square 1 \square T \square \square \square \square}. \quad (3.12)$$

$$T(N \square 1) n T \square \square \square \square \square \square n1 \square \square \square \quad \square \square$$

Now,

$$Var(\square^{\wedge I^*} n_2) \quad | \quad \square E(Var(\square^{\wedge I^*} n_2)) \square Var(E(\square^{\wedge I^*} n_2)) \quad \square E(Var(\bar{N}^{\wedge I^*} Y n_2)) \quad \square Var(E(N^{\wedge I^*} Y n_2))$$

$$\square \boxed{E(N^* 2Var(Y n2))} \square \boxed{2Var(N^* I^*)} \quad \square \boxed{E \overline{\square \square \square \square}} \square \boxed{n2(nT1 \square 1)} \overline{\square 1 \square \square \square 2 n1} \quad \square \boxed{n22 T N} \\ \square \boxed{nN1} \square \boxed{n12} \square \boxed{T} \square \boxed{\square \square \square \square}$$

I

$$\square_1(N \square_1)(N \square n_1)(n_1 \square_1 \square T) \square_2$$

$T(n_1 \square 2)$

$$\square^2 \quad \square \quad {}_2\square \quad N \quad \square\square$$

$$\square T\, {}_2(N\square 1)\, E\square\square\square\square n_2(n_1\square 1)\square T\square\square\square n_1\square n_2\square T\square 1\square\square\square\square\square\square$$

$$\begin{aligned}
 & (N\Box 1)(N\Box n)(n \\
 & \Box 1 \quad 1 \quad \Box 1\Box T)\Box 2 \\
 & T(n_1\Box 2) \\
 & \Box Z\Box^2 \Box J\Box^2,
 \end{aligned} \tag{3.18}$$

where

Note that, the previous derivations depend on the fact that n_2 has a negative hypergeometric distribution.

Now, if N is known, then \underline{N} can be estimated based on a simple random sample of size n $\square n_1 \square n_2 \square t$ by $\hat{N} \square NY$.

$$\begin{aligned} & \text{by } \square \square NY : \\ & \frac{2}{2} \quad N \square n \square \square Var \square N \square \square \\ & Var(\square \wedge) \square Var(NY) \square E(Var(NY \mid n)) \square Var(E(NY \mid n)) \square E \square \square \square N^2 \square n \square \square \square N \square 1 \square \square \square \square \square \end{aligned}$$

$\square N_2 \Box_2 \Box N \Box \Box E \Box \Box_1 \Box N \Box_1 \Box n \Box$

$\square N_2 \square_2 \bar{E} \square \square N \square \square \square N_2 \square_2 \square \quad N_3 \square_2 E \square \square_1 \square \square \square \quad N_2 \square_2$

$$\overline{N \square 1} \quad \overline{\square n \square} \quad \overline{\overline{N \square 1}} \quad \overline{N \square 1} \quad \overline{\square n \square} \quad \overline{N \square 1}$$

$$\square \quad \square n_2 \square 1 \square \square N \square n_2 \square \square$$

1 \square
 $\square NN_3 \square \square 12 Nn \square \square 2n \square 1 T \square T \square \square \square \square n_1 \square n_2 \square t \square \square T \square 1 \square \square \square \square \underline{N} \square \square \underline{n} \square 1 \square T$
 $\square \square \square \square \square \square \square \square NN_2 \square \square 1 2 \square H \square 2,$

where,

where,

$\boxed{}_n \boxed{} \quad \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} \boxed{} n_1 \boxed{} \boxed{}$

The efficiency of $\hat{\square}_I$ with respect to $\hat{\square}$ (obtained for a sample size equal the expected sample size) is:

$$Eff(\hat{\square}_I; \hat{\square}) = \frac{MSE(\hat{\square}_I)}{MSE(\hat{\square})},$$

$$\begin{aligned} & \overline{MSE(\hat{\mu}_I)} \\ & MSE(\hat{\mu}_I) \leq bias(\hat{\mu}_I)^2 + Var(\hat{\mu}_I) \leq bias(N^*_{\cdot I})^2 + B^2 Var(N^*_{\cdot I}) \\ & \leq B^2 MSE(N^*_{\cdot I})^2, \\ & Eff(\hat{\mu}_I; \hat{\mu}) = B^{-2} MSE(\hat{\mu}_{\cdot I})^2. \end{aligned}$$

The efficiency can be rewritten in terms of the coefficient of variation (CV), given by $CV \square$, (assume

\square
 $\square \square o$, as:

$$\begin{aligned} H^2 &= \frac{\square CV(y)^2}{MSE(N^I)} \\ Eff(\hat{N}_I; \hat{N}) &= \frac{1}{B} \frac{\square CV(y)^2}{MSE(N^I)} \end{aligned}$$

Some values of the efficiency of \hat{N}_I w.r.t. \hat{N} are given in Table (3.3). Similarly,
 $Eff(\hat{N}_I^*; \hat{N}_I) = \frac{B}{\square CV(y)^2} \frac{1}{MSE(N^I)}$

$$\underline{Z \square CV(y) \square \square J}$$

The efficiency of \hat{N}_I^* w.r.t. \hat{N}_I is given in Table (3.4). Based on the tables, we have the following conclusions:

1. From Table (3.3), the efficiency of \hat{N}_I w.r.t. \hat{N} increases when the CV increases.
2. For small sample size, the efficiency of \hat{N}_I w.r.t. \hat{N} 's is greater than or close to one. Also, the efficiency of \hat{N}_I w.r.t. \hat{N}_c is greater than one for small sample size and small value of the CV.
3. For large sample size, the efficiency of \hat{N}_I w.r.t. \hat{N} 's is less than one.
4. For large value of the coefficient of variation, the efficiency of \hat{N}_I w.r.t. \hat{N}_c is less than one, and it decreases when the sample size increases.
5. \hat{N}_I^* is more efficient than \hat{N}_I for large value of the expected sample size and CV.

4. Conclusions and Suggestions for Further Research

In this paper, four different estimators for the population total are discussed and from the results obtained we can conclude the following:

1. For large expected sample size, when we use Direct Sampling we found that the estimator of the population total \hat{N}_c based on Chapman estimator N_c is better than the estimator \hat{N}_s based on the suggested modified estimator N_s . On the other hand, if the expected sample size is small then \hat{N}_s is more efficient than \hat{N}_c .
2. For small expected sample size and CV, we found that it is better to use Indirect Sampling to estimate \hat{N} than using Direct Sampling.
3. The bias in N_I can be corrected to obtain an unbiased estimator of N , N_I^* , and also an unbiased estimator of \hat{N} .

Suggestions for Future Work

- If N is known, then an estimator of N (pretending it is unknown) can be used as a guard against unsuitable or insufficient sample. So, we may suggest estimator of \hat{N} conditioning on N to be between $N \square \square$ and $N \square \square$ for some \square .
- Estimation of \hat{N} based on other sampling techniques when N is unknown can also be considered next.

References

Ahmad, M., Alalouf, S. and Chaubey, Y. (2000). Estimation of the population total when the population size is unknown. Statistics and Probability Letters **49**, 211-216.

Arnab, R. (2004). Optimum estimation of a finite population total in PPS sampling with replacement for multicharacter surveys. Journal of the Indian Society of Agricultural Statistics **58**, 231-243.

Azevedo-Silva, F., Reis, R., Santos, M., Luiz, R.and Pombo-de-oliveira, M. (2009).

Evaluation of childhood acute leukemia incidence and underreporting in Brazil by capture-recapture methodology. *Cancer Epidemiology* **33**, 403-405.

Balakrishnan, N., Charalambides, C.and Papadatos, N. (2003). Bounds on expectation of order statistics from a finite population. *Journal of Statistical Planning and Inference* **113**, 569 – 588.

Chapman, D. (1951). Some properties of the Hyper-geometric distribution with application to zoological censuses. University of California Publication in Statistics.

Gutierrez, H.and Breidt, J. (2009). Estimation of the population total using the generalized difference estimator and Wilcoxon ranks. *Revista Colombiana de Estadística* **32**, 123-143.

Khan, R. A.(1994). A Note on the generating function of a negative hypergeometric distribution. *Sankhya*, 56, 309-313.

Luan, R., Liang, B., Yuan, P.,Fan, L., Huang, Y., Zeng, G., Wang, L., and Wang, S. (2005). A study on the Capture-Recapture method for estimating the population size of injecting drug users in Southwest China. *Journal of Health Science* **51**, 405-409.

Mohammad, M.and Abdullah, M. (2007). Undergraduate graduation project. Qatar University.

Otieno, R., Mwita, P., Obudho, G. and Wafula, C. (2005). Model-based estimation of finite population total in stratified sampling. *East African Journal of Stat.* **1**, 23-40.

Scheaffer, Mendenhall and Ott (1995). Elementary survey sampling. Fifth edition, Duxbury Press.

Sekar, C. and Deming, W. (1949). On a method of estimation birth and death rates and the extent of registration. *JASA* **44**, 101-115.

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

Table(2.1): $Eff(N^S; N^C)$

N	$E(n)$	$MSE(N^S)$	$MSE(N^C)$	$Eff(N^S; N^C)$
1000	51.9	728205.9	732439.8804	1.005814
1000	57.6	332593.5929	402913.9604	1.21143
1000	61.4	177120.7456	289004.9081	1.631683
1000	67.1	84362.2724	213053.84	2.525464
1000	72.8	117772.4201	201699.6164	1.712622
1000	73.8	132095.56	203386.25	1.53969
1000	80.4	275199.96	229010.9225	0.832162
1000	92.8	598569.24	286622.6784	0.478846
1000	116.5	918583.24	311101.1376	0.338675
5000	227.6	16746267.56	7491881	0.447376
5000	237.3	19768303.04	7903480.49	0.399806
5000	247.0	22055514.24	8134961	0.36884
5000	252.82	23062512.25	8193942.25	0.355293

5000	253.79	23204389.44	8198340.81	0.35331
5000	266.4	24406592.01	8135765.61	0.333343
5000	290.65	23995928.16	7566230.24	0.315313
5000	314.9	21416904.64	6725680.81	0.314036
5000	344.0	17364010.41	5689186.49	0.327642

Table (2.2): $Eff(\square^C; \square^C)$

N	E(n)	CV \square 0.5	CV \square 1	CV \square 5	CV \square 10
1000	51.9	0.006241	0.024954	0.615966	2.370284
1000	57.6	0.010147	0.040398	0.878713	2.499718
1000	61.4	0.01319	0.052189	0.968865	2.147734
1000	67.1	0.016208	0.063373	0.920717	1.595043
1000	72.8	0.015642	0.060687	0.773051	1.220907
1000	73.75	0.015294	0.059299	0.748042	1.17425
1000	80.4	0.012367	0.04792	0.598649	0.93414
1000	92.75	0.008463	0.032934	0.440701	0.718821
1000	116.5	0.006056	0.023686	0.346428	0.603322
5000	227.6	0.003486	0.01378	0.250864	0.542584
5000	237.3	0.003163	0.012512	0.230536	0.506165
5000	247.0	0.002947	0.011661	0.216907	0.482064
5000	252.82	0.002855	0.011299	0.211125	0.471948
5000	253.79	0.002842	0.011248	0.21031	0.470536
5000	266.4	0.002721	0.010773	0.202774	0.457653
5000	290.65	0.002668	0.010566	0.199714	0.453306
5000	314.9	0.002756	0.010912	0.205493	0.464102
5000	344.0	0.002964	0.011728	0.218439	0.486251

Table (2.3): $Eff(\square^S; \square^S)$

N	E(n)	CV \square 0.5	CV \square 1	CV \square 5	CV \square 10
1000	51.9	0.006277	0.025099	0.619341	2.380999
1000	57.6	0.012284	0.0488	1.002176	2.573034
1000	61.4	0.021431	0.083741	1.20335	2.066933
1000	67.1	0.040074	0.147658	1.047878	1.294508
1000	72.8	0.026292	0.096826	0.684161	0.844183
1000	73.75	0.023162	0.08576	0.634498	0.793077
1000	80.4	0.010235	0.039044	0.393256	0.548862
1000	92.75	0.004054	0.015796	0.215683	0.356763
1000	116.5	0.002056	0.008094	0.134857	0.264126
5000	227.6	0.00156	0.006172	0.114769	0.254933
5000	237.3	0.001265	0.005014	0.096394	0.223911
5000	247.0	0.001088	0.004314	0.084914	0.204024
5000	252.82	0.001015	0.004028	0.080156	0.195766
5000	253.79	0.001005	0.003988	0.079481	0.194601
5000	266.4	0.000908	0.003606	0.073122	0.183947
5000	290.65	0.000842	0.003348	0.069215	0.179724
5000	314.9	0.000867	0.003445	0.071808	0.188995
5000	344.0	0.000973	0.003866	0.080548	0.211892

Table (2.4): $Eff(\hat{s}; \hat{c})$

N	$E(n)$	$CV \hat{0.5}$	$CV \hat{1}$	$CV \hat{5}$	$CV \hat{10}$
1000	51.9	1.005811	1.005801	1.005478	1.004521
1000	57.6	1.210561	1.207977	1.140504	1.02933
1000	61.4	1.624743	1.604566	1.24202	0.962378
1000	67.1	2.472407	2.329969	1.138111	0.811582
1000	72.8	1.680817	1.595495	0.885014	0.69144
1000	73.75	1.514448	1.446229	0.848211	0.675391
1000	80.4	0.827601	0.814762	0.656905	0.587559
1000	92.75	0.479045	0.47962	0.489408	0.496317
1000	116.5	0.339446	0.341712	0.389279	0.437785
5000	227.6	0.447514	0.44792	0.457494	0.46985
5000	237.3	0.400046	0.400759	0.418129	0.442368
5000	247.0	0.36913	0.36999	0.391477	0.42323
5000	252.82	0.355601	0.356517	0.379659	0.414804
5000	253.79	0.353621	0.354545	0.377924	0.413574
5000	266.4	0.333682	0.334688	0.36061	0.401936
5000	290.65	0.315694	0.316825	0.34657	0.396475
5000	314.9	0.314463	0.315735	0.349444	0.407227
5000	344.0	0.328139	0.329615	0.368744	0.435766

Table (3.1): Expectation of n

N	T	n_1	$E(n)$
1000	1	50	68.627
1000	3	50	105.88
1000	4	50	124.51
1000	5	100	144.55
1000	10	100	189.11
1000	15	100	233.66
1000	25	250	324.7
1000	63	250	438.25
1000	75	250	474.1
5000	1	100	148.51
5000	2	100	197.03
5000	8	400	491.77

5000	32	400	767.08
5000	63	400	812.97
5000	50	500	949.1
5000	65	500	1083.8

Table(3.2): $Eff(N^I; N^{*I})$

N	$E(n)$	$Eff(N^I, N^{*I})$
1000	68.627	0.96155
1000	105.88	0.962362
1000	124.51	0.962765
1000	144.55	0.980766
1000	189.11	0.981287
1000	233.66	0.981863
1000	324.7	0.992361
1000	438.25	0.993052
1000	474.1	0.993324
5000	148.51	0.980392
5000	197.03	0.980514
5000	491.77	0.995068
5000	767.08	0.995215
5000	812.97	0.995449
5000	949.1	0.99621
5000	1083.8	0.996266

Table(3.3): $Eff(\square^I; \square^*)$

N	$E(n)$	$CV \square 0.5$	$CV \square 1$	$CV \square 5$	$CV \square 10$
1000	68.627	0.004084	0.016098	0.275107	0.553307
1000	105.88	0.008043	0.030201	0.254925	0.332162
1000	124.51	0.009028	0.032967	0.217609	0.263777
1000	144.55	0.0087	0.030742	0.162453	0.187566
1000	189.11	0.012318	0.036371	0.096949	0.102273
1000	233.66	0.013651	0.034478	0.067367	0.069437

1000	324.7	0.013146	0.026527	0.039341	0.039944
1000	438.25	0.011182	0.014516	0.016046	0.016099
1000	474.1	0.01013	0.012492	0.013499	0.013533
5000	148.51	0.001876	0.007452	0.1538	0.398104
5000	197.03	0.002887	0.011358	0.186294	0.359178
5000	491.77	0.003979	0.014462	0.092053	0.110597
5000	767.08	0.007933	0.018096	0.030669	0.031349
5000	812.97	0.008032	0.0173	0.027427	0.027938
5000	949.1	0.007608	0.014631	0.019868	0.020093
5000	1083.8	0.008022	0.012369	0.015381	0.015499