

Total Estimators in Inverse Sampling

L. Greco

Dipartimento di Metodi Quantitativi, Università degli Studi di Siena
greco2@unisi.it

S. Naddeo

Dipartimento di Metodi Quantitativi, Università degli Studi di Siena
naddeo@unisi.it

Abstract. Inverse sampling is an efficient design when the population of interest is subdivided into two groups, one of which contains only few units. With this design it is possible to obtain efficient parameter estimators for the whole population as well as for its subgroups. We consider unequal selection probabilities in sampling with replacement, as well as equal selection probabilities in sampling with and without replacement.

Keywords. sequential sampling, rare populations, unequal selection probabilities, stopping rules.

1. Introduction

Inverse sampling is an efficient design for estimating the parameters of a population in which only few units display the characteristic of interest or when the variable of interest tends to be at or near zero for many of the population units and only one subgroup exhibits values different from zero. This design is also efficient when the aim of the survey is to estimate the parameters of a subgroup formed by few units, in addition to the parameters of the whole population. In many real situations the population may be considered as subdivided into two subgroups, and the subgroup to which a unit belongs is not known until the unit is sampled. In these cases the classical fixed sample-size designs may not give efficient estimates for the population parameters. Instead, in inverse sampling designs the selection of units continues until a prefixed number of units with the characteristic of interest is observed in the sample.

Recently, Christman and Lan (2001) considered inverse sampling in simple random sampling with and without replacement, when all the population units have equal selection probabilities. They provided some unbiased estimators of the population total and their variances, but not the variance estimators.

Salehi and Seber (2001, 2004) considered inverse sampling with equal selection probabilities by using Murthy's method, showing that this method is an application of the Rao-Blackwell theorem. With this approach they obtained the unbiased estimator of the population total using simple random sampling without replacement, with the same results given by Christman and Lan (2001). Salehi and Seber (2001, 2004) also obtained a variance estimator which is unbiased on the basis of Murthy's method.

In this paper we consider inverse sampling when the units have unequal selection probabilities, as happens in environmental surveys when the selection probabilities are evaluated at the moment in which the units are selected. We derive unbiased estimators of the totals of the two subgroups, their variance and the corresponding variance estimators in inverse sampling with replacement (ISWR). Our estimator of the whole population total is equivalent to that given by Christman and Lan (2001) when the population units have equal selection probabilities.

Finally, we obtain similar results in inverse sampling without replacement (ISWOR) when the population units have equal selection probabilities, and we show that the variance estimator reported by Salehi and Seber (2001, 2004) is the unbiased estimator of the variance given by Christman and Lan (2001).

2. Unequal selection probabilities

Let us consider a population of N units divided into two groups of N_1 and N_2 units respectively, with $N_1+N_2=N$, and let y_{1i} ($i=1,2,\dots,N_1$) and y_{2i} ($i=1,2,\dots,N_2$) be the values of the variable of interest Y in the first and second group. For example, the population can be subdivided into two subgroups according to whether the y -values satisfy some condition C . Moreover, let T_1 and T_2 be the totals of the two groups and $T=T_1+T_2$ the population total, p_{1i} the selection probability of the i -th unit of the first group and p_{2i} the selection probability of the i -th unit of the second group.

Finally, let $P = \sum_{i=1}^{N_1} p_{1i}$ be the selection probability of the first group and $1 - P = \sum_{i=1}^{N_2} p_{2i}$ the selection probability of the second group.

In inverse sampling designs, the selection of units continues until k units of the first group are observed in the sample, so that the sample size ν is a random variable which can assume the values $k, k+1, \dots$. Let us consider the event $\nu=n$ and all the possible samples of size n . Among them the only samples taken into account with inverse sampling are formed by k units coming from the first group and $n-k$ units from the second group.

In ISWR the total probability of all samples is

$$P^k (1 - P)^{n-k},$$

so that the probability of selecting a sample of n elements in a given order is

$$\frac{\prod_{i=1}^{N_1} p_{1i}^{m_{1i}} \prod_{i=1}^{N_2} p_{2i}^{m_{2i}}}{P^k (1 - P)^{n-k}} = \prod_{i=1}^{N_1} \left(\frac{p_{1i}}{P} \right)^{m_{1i}} \prod_{i=1}^{N_2} \left(\frac{p_{2i}}{1 - P} \right)^{m_{2i}},$$

where m_{ji} counts how many times the i -th unit of j -th group is in the sample ($j=1,2; i=1,\dots, N_j$), with

$$\sum_{i=1}^{N_1} m_{1i} = k \quad \text{and} \quad \sum_{i=1}^{N_2} m_{2i} = n - k.$$

For a given n the selection probability of the i -th unit in the first group is $\frac{p_{1i}}{P}$, while $\frac{p_{2i}}{1 - P}$ is the selection probability of the i -th unit in the second group. Thus, conditionally to $v=n$ the selection procedure with inverse sampling is exactly like the selection procedure with stratified sampling, where the two strata correspond to the two subgroups and k and $n-k$ units are respectively selected from the first and from the second stratum. Obviously, the sample results in the two groups are independent.

Let Y_{ji} be the random variable “value of Y ” on the i -th sample unit from the j -th group, so that the estimator of the general total may be written as

$$\hat{T} = \hat{T}_1 + \hat{T}_2 = \frac{P}{k} \sum_{i=1}^k \frac{Y_{1i}}{p_{1i}} + \frac{1 - P}{v - k} \sum_{i=1}^{v-k} \frac{Y_{2i}}{p_{2i}}. \quad (1)$$

Usually, the value of P is not known, but an unbiased estimator is

$$\hat{P} = \frac{k - 1}{v - 1}, \quad (2)$$

as can be derived from the negative binomial distribution of v . Hence, if we denote by

$$W_{ji} = \frac{Y_{ji}}{p_{ji}}$$

the sample i.i.d. random variables, an estimator of the total is

$$\tilde{T} = \tilde{T}_1 + \tilde{T}_2 = \frac{\hat{P}}{k} \sum_{i=1}^k W_{1i} + \frac{1-\hat{P}}{v-k} \sum_{i=1}^{v-k} W_{2i} = \hat{P}\bar{W}_1 + (1-\hat{P})\bar{W}_2, \quad (3)$$

where the means of \bar{W}_1 and \bar{W}_2 are respectively

$$E(\bar{W}_1) = W_1 = \frac{T_1}{P}$$

and

$$E(\bar{W}_2) = W_2 = \frac{T_2}{1-P}.$$

Since

$$E_v[E_Y(\tilde{T})] = PW_1 + (1-P)W_2 = T_1 + T_2,$$

estimator \tilde{T} is unbiased. As regards its variance, we may observe that

$$V_v[E_Y(\tilde{T})] = (W_1 - W_2)^2 = V_v(\hat{P}),$$

and

$$\begin{aligned} E_v[V_Y(\tilde{T})] &= \frac{\sigma_{1w}^2}{k} E_v(\hat{P}^2) + \sigma_{2w}^2 E_v\left[\frac{(1-\hat{P}^2)}{v-k}\right] = \\ &= \frac{\sigma_{1w}^2}{k} E_v(\hat{P}^2) + \frac{\sigma_{2w}^2}{k-1} E_v[\hat{P}(1-\hat{P})] \end{aligned}$$

where

$$\sigma_{jw}^2 = \sum_{i=1}^{N_j} (W_j - W_j)^2 p_{ji} \quad j=1,2,$$

so that

$$V(\tilde{T}) = (W_1 - W_2)^2 V_v(\hat{P}) + \frac{\sigma_{1w}^2}{k} E_v(\hat{P}^2) + \frac{\sigma_{2w}^2}{k-1} E_v[\hat{P}(1-\hat{P})].$$

In order to obtain an unbiased variance estimator, let us note that an unbiased estimator of W_1^2 is

$$\hat{W}_{1q} = \bar{W}_1^2 - \frac{\hat{S}_{1w}^2}{k},$$

while the corresponding unbiased estimator of W_2^2 is

$$\hat{W}_{2q} = \bar{W}_2^2 - \frac{\hat{S}_{2w}^2}{v-k},$$

where \hat{S}_{jw}^2 is the unbiased sample variance of the W_{ji} s in the j -th group.

Moreover, an unbiased estimator of P^2 is

$$\hat{P}_q = \frac{(k-1)(k-2)}{(v-1)(v-2)} \quad (4)$$

so that an unbiased estimator of the variance of \hat{P} is

$$\hat{V}(\hat{P}) = \hat{P}^2 - \frac{(k-1)(k-2)}{(v-1)(v-2)} = \frac{\hat{P}(1-\hat{P})}{v-2}.$$

and finally

$$\hat{V}(\tilde{T}) = (\bar{W}_1 - \bar{W}_2)^2 \frac{\hat{P}(1-\hat{P})}{v-2} + \frac{\hat{S}_{1w}^2}{k} \hat{P}_q + \frac{\hat{S}_{2w}^2}{k-1} \left(\hat{P} - \frac{k-1}{k-2} \hat{P}_q \right). \quad (5)$$

As regards the two subgroups, observe that the variances of the two unbiased estimators of the totals \tilde{T}_1 and \tilde{T}_2 are respectively

$$V(\tilde{T}_1) = W_1^2 V_v(\hat{p}) + \frac{\sigma_{1w}^2}{k} E_v(\hat{p}^2)$$

and

$$V(\tilde{T}_2) = W_2^2 V_v(\hat{p}) + \frac{\sigma_{2w}^2}{k-1} E_v[\hat{p}(1-\hat{p})]$$

and their unbiased estimators may be obtained from expression (5).

3. Stopping rules

With inverse sampling it is possible to run out of resources (e.g., money or time) prior to selecting k units from the first subgroup. In this case it is still possible to get unbiased estimators of the subgroup totals and their variances.

Thus, we consider a sampling design in which the selection procedure stops if: (i) k units of the first subgroup are observed in the sample; (ii) the maximum sample size M is obtained without getting k units of the first subgroup. In the first case, the total estimator is given in expression (3), while in the second case the estimator is the classical estimator of the totals of the two subgroups used in the fixed sample-size designs when M units are randomly selected from the whole population. Obviously, in the last case the number of selected units from the first subgroup is a random variable X with a binomial distribution of parameters P and M .

By using this stopping rule, the unbiased estimator of the population total is

$$\tilde{T}_s = \begin{cases} \tilde{T}_a = \hat{P}_a \overline{W}_{1a} + (1 - \hat{P}_a) \overline{W}_{2a} & \text{case (i)} \\ \tilde{T}_b = \hat{P}_b \overline{W}_{1b} + (1 - \hat{P}_b) \overline{W}_{2b} & \text{case (ii)} \end{cases} \quad (6)$$

where \tilde{T}_a is used when the k -th unit of the first subgroup is selected at the n -th replication (with $k \leq n \leq M$), \tilde{T}_b is the estimator which is computed when $n=M$ and the number of sampled units of the

first subgroup is x (with $0 \leq x < k$), $\hat{P}_a = \frac{k-1}{v-1}$, $\hat{P}_b = \frac{X}{M}$, and finally \bar{W}_{1a} , \bar{W}_{2a} , \bar{W}_{1b} , \bar{W}_{2b} are the sample means in the two cases.

The variance of estimator (6) is

$$V(\tilde{T}_s) = \sum_{n=k}^M E_Y \left[(\tilde{T}_a - T)^2 \right] f(n; k, P) + \sum_{x=0}^{k-1} E_Y \left[(\tilde{T}_b - T)^2 \right] g(x; M, P),$$

where $f(n; k, P)$ denotes the negative binomial distribution of parameters k and P , while $g(x; M, P)$ is the binomial distribution of parameters M and P . An unbiased variance estimator is

$$\hat{V}(\tilde{T}_s) = \begin{cases} \hat{V}_a(\tilde{T}_s) = (\bar{W}_{1a} - \bar{W}_{2a})^2 (\hat{P}_a^2 - \hat{P}_{aq}) + \frac{\hat{S}_{1w}^2}{k} \hat{P}_{aq} + \frac{\hat{S}_{2w}^2}{v-k} (1 - 2\hat{P}_a + \hat{P}_{aq}) & \text{case (i)} \\ \hat{V}_b(\tilde{T}_s) = (\bar{W}_{1b} - \bar{W}_{2b})^2 (\hat{P}_b^2 - \hat{P}_{bq}) + \frac{\hat{S}_{1w}^2}{X} \hat{P}_{bq} + \frac{\hat{S}_{2w}^2}{M-X} (1 - 2\hat{P}_b + \hat{P}_{bq}) & \text{case (ii)} \end{cases} \quad (7)$$

where \hat{P}_{aq} is given by (4), $\hat{P}_{bq} = \frac{X}{M} \frac{X-1}{M-1}$ and $\hat{V}_a(\tilde{T}_s)$ is equivalent to (5).

Estimators (6) and (7) are unbiased, as is shown in the Appendix.

The unbiased estimators of T_1 and T_2 and their unbiased variance estimators may be easily obtained from the previous paragraph and from (7).

A more complex design considers an initial sample of a prefixed size, say m . In this case the selection procedure stops if: (i) at least k units of the first subgroup are selected in the initial sample; (ii) exactly k units of the first subgroup are observed in a sample of size n (with $m < n \leq M$); (iii) the maximum sample size M is obtained without obtaining k units of the first subgroup.

Obviously, previous results apply in cases (ii) and (iii), while in case (i) the expressions of the totals estimators, of their variances and of the variances estimators may be easily obtained from the previous expressions when the binomial distribution applies, substituting M with m . Also with this stopping rule, the total estimators of the population and of the two subgroups are unbiased, as well as the estimators of the variances of T , T_1 and T_2 , as is shown in Appendix.

Finally, it is worth noting that if the procedure selection considers only cases (i) and (ii), i.e. a maximum sample size M is not fixed, the corresponding estimators are still unbiased.

4. Equal selection probabilities

If the selection probabilities are equal for every population unit so that the selection probability of the i -th unit in the j -th group ($j=1,2; i=1,2,\dots,N_j$) is $1/N$, estimator (3) under ISWR assumes the form

$$\tilde{T} = N[\hat{P}\bar{Y}_1 + (1 - \hat{P})\bar{Y}_2], \quad (8)$$

where \bar{Y}_1 and \bar{Y}_2 are the sample means in the two subgroups. Its variance is

$$V(\tilde{T}) = N^2 \left\{ (\mu_1 - \mu_2)^2 V_v(\hat{P}) + \frac{\sigma_1^2}{k} E_v(\hat{P}^2) + \frac{\sigma_2^2}{k-1} E_v[\hat{P}(1 - \hat{P})] \right\} \quad (9)$$

where μ_j and σ_j^2 are the mean and the variance of the variable of interest in the j -th group ($j=1,2$) respectively. It is worth noting that the estimator of the population total and its variance are equivalent to the expressions given by Christman and Lan (2001).

An unbiased estimator of (9) is

$$\hat{V}(\tilde{T}) = N^2 \left\{ (\bar{Y}_1 - \bar{Y}_2)^2 \frac{\hat{P}(1 - \hat{P})}{v-2} + \frac{\hat{S}_1^2}{k} \hat{P}_q + \frac{\hat{S}_2^2}{k-1} \left[\hat{P} - \frac{k-1}{k-2} \hat{P}_q \right] \right\}$$

where \hat{S}_j^2 is the unbiased sample variance in the j -th group.

If we consider ISWOR and equal selection probabilities, from the hypergeometric distribution it turns out that the first order inclusion probabilities are the same as in the previous case. In fact, conditionally to $v=n$, the number of samples formed by k units from the first subgroup and $n-k$ units from the second subgroup is

$$\binom{N_1}{k} \binom{N_2}{n-k},$$

while the number of samples containing, for example, the i -th unit of the first subgroup is

$\binom{N_1-1}{k-1} \binom{N_2}{n-k}$. Thus, its inclusion probability turns out to be $\frac{k}{N_1}$. Similarly, the inclusion

probability of any unit coming from the second subgroup is $\frac{n-k}{N_2}$. The second order inclusion probability of two units in the first subgroup is $\frac{k}{N_1} \frac{k-1}{N_1-1}$, the corresponding probability for two units in the second group is $\frac{n-k}{N_2} \frac{n-k-1}{N_2-1}$, while it is $\frac{k}{N_1} \frac{n-k}{N_2}$ if the two units come from different subgroups.

Practically, these results show that, conditionally to $v=n$, the selection procedure with inverse sampling is exactly like the selection procedure with stratified sampling, where k and $n-k$ units are respectively selected without replacement in the first stratum and, independently, in the second stratum.

Since in this case v has a negative hypergeometric distribution, an unbiased estimator of P is still given by (2) and the estimator of the population total is given by expression (8). Its variance is

$$V(\tilde{T}) = N^2 \left\{ (\mu_1 - \mu_2)^2 V_v(\hat{P}) + \frac{S_1^2}{k} \left(1 - \frac{k}{N_1}\right) E_v(\hat{P}^2) + \frac{S_2^2}{k-1} E_v \left[\hat{P}(1-\hat{P}) \left(1 - \frac{v-k}{N_2}\right) \right] \right\},$$

where

$$S_j^2 = \frac{1}{N_j - 1} \sum_{i=1}^{N_j} (y_{ji} - \mu_j)^2, \quad j = 1, 2,$$

and it is similar to the variance of \tilde{T} in ISWR with finite population correction factors. Also in this case the estimator of the total and its variance correspond to the expressions given by Christman and Lan (2001).

To obtain an unbiased variance estimator, let us note that an unbiased estimator of μ_1^2 is

$$\hat{\mu}_{1q} = \bar{Y}_1^2 - \left(1 - \frac{k}{N_1}\right) \frac{\hat{S}_1^2}{k},$$

while the corresponding unbiased estimator of μ_2^2 is

$$\hat{\mu}_{2q} = \bar{Y}_2^2 - \left(1 - \frac{v-k}{N_2}\right) \frac{\hat{S}_2^2}{v-k}.$$

An unbiased estimator of P^2 is

$$\hat{P}_q = \frac{k-1}{v-1} \frac{k-2}{v-2} + \frac{1}{N} \frac{k-1}{v-1} \left(1 - \frac{k-2}{v-2}\right),$$

so that we have

$$\hat{V}(\hat{P}) = \left(\frac{k-1}{v-1}\right)^2 - \hat{P}_q = \frac{\hat{P}(1-\hat{P})}{v-2} \left(1 - \frac{v-1}{N}\right).$$

Thus, it is easy to obtain the unbiased estimator (10) of the variance of \tilde{T}

$$\hat{V}(\tilde{T}) = N^2 \left[(\bar{Y}_1 - \bar{Y}_2)^2 \hat{V}(\hat{P}) + \frac{\hat{S}_1^2}{k} \left(\hat{P}_q - \frac{k}{N} \hat{P}\right) + \frac{\hat{S}_2^2}{v-k} \left(\hat{P}_q - \frac{v-k}{N} (1-\hat{P})\right) \right],$$

which is equivalent to the expression given by Salehi and Seber (2004).

The estimators of the two subgroup totals and of their variances in ISWOR may be easily obtained from the results given in paragraph 2. Finally, Salehi and Seber (2004) give the unbiased estimators of the population total and its variance for stopping rules similar to those considered in this paper.

References

- Christman, M.C. and Lan, F. (2001). Inverse Adaptive Cluster Sampling. *Biometrics* **57**, 1096-1105.
- Salehi, M. M. and Seber, G.A.F. (2001). A new proof of Murthy's estimator which applies to sequential sampling. *Australian and New Zealand Journal of Statistics* **43**, 281-286.
- Salehi, M. M. and Seber, G.A.F. (2004). A General Inverse Sampling Scheme and its Application to Adaptive Cluster Samplig. *Australian and New Zealand Journal of Statistics* **46**, 483-494.

Appendix

It is worth pointing out the connection between the negative binomial and the binomial distributions, from which the probability that more than M replications are necessary to get k units from a given subgroup is equal to the probability that less than k units from the same subgroup are selected in M replications. Thus,

$$\sum_{n=M+1}^{\infty} f(n; k, P) = \sum_{x=0}^{k-1} g(x; M, P),$$

where $f(n; k, P)$ denotes the negative binomial distribution of parameters k and P , while $g(x; M, P)$ is the binomial distribution of parameters M and P , from which it turns out

$$\sum_{n=k}^M f(n; k, P) + \sum_{x=0}^{k-1} g(x; M, P) = 1. \quad (\text{A1})$$

From the previous equality it is apparent that estimator (6) is unbiased. In fact

$$\begin{aligned} E(\tilde{T}_s) &= \sum_{n=k}^M E_y(\tilde{T}_a) f(n; k, P) + \sum_{x=0}^{k-1} E_y(\tilde{T}_b) g(x; M, P) = \\ &= \sum_{n=k}^M [\hat{p}_a (W_1 - W_2) + W_2] f(n; k, P) + \sum_{x=0}^{k-1} [\hat{p}_b (W_1 - W_2) + W_2] g(x; M, P) = \\ &= P(W_1 - W_2) + W_2 = T_1 + T_2, \end{aligned}$$

where $\hat{p}_a = \frac{k-1}{n-1}$, and $\hat{p}_b = \frac{x}{M}$, as results from (A1) and the following equality

$$\sum_{n=k}^M \hat{p}_a f(n; k, P) + \sum_{x=0}^{k-1} \hat{p}_b g(x; M, P) = P. \quad (\text{A2})$$

As regards the variance estimator (7), observe that

$$\begin{aligned}
V(\tilde{T}_s) &= \sum_{n=k}^M E_y \left[(\tilde{T}_a - T)^2 \right] f(n; k, P) + \sum_{x=0}^{k-1} E_y \left[(\tilde{T}_b - T)^2 \right] g(x; M, P) = \\
&= \sum_{n=k}^M \left[(W_1 - W_2)^2 (\hat{P}_a - P)^2 + \frac{\sigma_{1w}^2}{k} \hat{P}_a^2 + \frac{\sigma_{2w}^2}{n-k} (1 - \hat{P}_a)^2 \right] f(n; k, P) + \\
&+ \sum_{x=0}^{k-1} \left[(W_1 - W_2)^2 (\hat{P}_b - P)^2 + \frac{\sigma_{1w}^2}{M} \hat{P}_b^2 + \frac{\sigma_{2t}^2}{M} (1 - \hat{P}_b)^2 \right] g(x; M, P)
\end{aligned}$$

and

$$\begin{aligned}
E[\hat{V}(\tilde{T}_s)] &= \sum_{n=k}^M E_y \left[\hat{V}_a(\tilde{T}_s) \right] f(n; k, P) + \sum_{x=0}^{k-1} E_y \left[\hat{V}_b(\tilde{T}_s) \right] g(x; M, P) = \\
&= \sum_{n=k}^M (W_1 - W_2)^2 (\hat{P}_a^2 - \hat{P}_{aq}) + \frac{\sigma_{1w}^2}{k} \hat{P}_a^2 + \frac{\sigma_{2w}^2}{n-k} (1 - \hat{P}_a)^2 f(n; k, P) + \\
&+ \sum_{x=0}^{k-1} \left[(W_1 - W_2)^2 (\hat{P}_b^2 - \hat{P}_{bq}) + \frac{\sigma_{1w}^2}{M} \hat{P}_b^2 + \frac{\sigma_{2t}^2}{M} (1 - \hat{P}_b)^2 \right] g(x; M, P).
\end{aligned}$$

From (A1), (A2) and

$$\sum_{n=k}^M \hat{P}_{aq} f(n; k, P) + \sum_{x=0}^{k-1} \hat{P}_{bq} g(x; M, P) = P^2, \quad (A3)$$

where $\hat{P}_{aq} = \frac{k-1}{n-1} \frac{k-2}{n-2}$ and $\hat{P}_{bq} = \frac{x}{M} \frac{x-1}{M-1}$, it turns out that $E[\hat{V}(\tilde{T}_s)] = V(\tilde{T}_s)$.

When at least m and at most M replications are performed, the equality analogous to (A1) is

$$\sum_{x=k}^m g(x; m, P) + \sum_{n=m+1}^M f(n; k, P) + \sum_{x=0}^{k-1} g(x; M, P) = 1. \quad (A4)$$

From (A4), equalities analogous to (A2) and (A3) may be easily obtained and these equalities are sufficient to show that the totals estimators are unbiased, as well as the estimators of their variances.