

On studying the relationship between the ratio and the regression estimators for a population proportion

Muñoz, Juan F.

University of Granada, Department of Quantitative Methods in Economic and Business

Facultad de Ciencias Económicas y Empresariales, Campus Cartuja s/n

Granada (18071), Spain

E-mail: jfmunoz@ugr.es

Arcos, Antonio

University of Granada, Department of Statistics and O.R.

Facultad de Farmacia, Campus Cartuja s/n

Granada (18071), Spain

E-mail: arcos@ugr.es

Álvarez-Verdejo, Encarnación

University of Granada, Department of Quantitative Methods in Economic and Business

Facultad de Ciencias Económicas y Empresariales, Campus Cartuja s/n

Granada (18071), Spain

E-mail: encarniav@ugr.es

Abstract

Estimation of a proportion is commonly used in many areas and disciplines. Traditional estimators and confidence intervals for a population proportion does not assume auxiliary information at the estimation stage. Assuming auxiliary information, the ratio and regression methods are known techniques for the purpose of estimation of a population mean. They can be adapted to the problem of the estimation of the population proportion, although some differences can be observed. For example, a ratio estimator for a population proportion can be defined via a ratio estimator for the complementary of the population proportion. We define an optimum ratio estimator for a population proportion and show that this estimator coincides with the regression estimator. Results derived for the population proportion can be easily extended to the problem of estimating the distribution function.

Keywords: Auxiliary information, sample survey, variance, inclusion probability.

Introduction

Traditional estimators and confidence intervals for a population proportion (Blyth and Still, 1983; Cohen and Yang, 1994; Fleiss et al., 2003; Newcombe, 1998) do not assume auxiliary information at the estimation stage. Assuming auxiliary information, the ratio and regression methods are known techniques for the purpose of estimation of a population mean. They can be adapted to the problem of the estimation of the population proportion (see, for example, Rueda et al., 2011a), although some differences can be observed. For example, a ratio estimator for a population proportion can be defined via a ratio estimator for the complementary of the population proportion. Furthermore, one should be aware of the risks when confidence intervals are constructed for a population proportion, since limits outside $[0, 1]$ could be achieved.

We consider the scenario of a finite population $U = \{1, \dots, N\}$ containing N units. Let A_1, \dots, A_N denote the values of a attribute of interest A , where $A_i = 1$ if i th unit possesses the attribute A and $A_i = 0$ otherwise. Let B denote an auxiliary attribute associated with A and values given by B_1, \dots, B_N . We also assume that a sample s , of size n , is selected from U according to an

arbitrary sampling design with first and second order inclusion probabilities given by π_i and π_{ij} .

The aim is to estimate the population proportion of individuals that possess the attribute A , i.e. $P_A = N^{-1} \sum_{i=1}^N A_i$. The naive estimator of P_A , which makes no use of the auxiliary information, is given by $\hat{p}_A = N^{-1} \sum_{i \in s} d_i A_i$, where $d_i = \pi_i^{-1}$. Note that \hat{p}_A can take values larger than 1. An alternative estimator that avoids values larger than 1 is the Hájek type estimator, which is given by $\hat{p}_{A.H} = \hat{N}^{-1} \sum_{i \in s} d_i A_i$, where $\hat{N} = \sum_{i \in s} d_i$.

Note that results related to the population proportion can be easily extended to the problem of estimating the distribution function $F(t)$ of a variable of interest y , since $F(t)$ is defined as the proportion of individuals in the population which take values of y less or equal than a given argument t .

Assuming a general sampling design, Rueda et al. (2011a) defined the ratio type estimators $\hat{p}_r = \hat{R}P_B$ and $\hat{p}_{r.q} = 1 - \hat{q}_r = \hat{R}_c Q_B$, where $\hat{R} = \hat{p}_A / \hat{p}_B$ is the design-based estimator of the population ratio $R = P_A / P_B$, $\hat{p}_B = N^{-1} \sum_{i \in s} d_i B_i$, $P_B = N^{-1} \sum_{i=1}^N B_i$ is the population proportion of individuals that possess the attribute B , $\hat{R}_c = \hat{q}_A / \hat{q}_B$, and the complementary proportions are defined as $\hat{q}_A = 1 - \hat{p}_A$, $\hat{q}_B = 1 - \hat{p}_B$ and $Q_B = 1 - P_B$. It is assumed that P_B is known from a census or estimated without error. Rueda et al. (2011b) also defined a regression type estimator for the population proportion under a general sampling design.

Assuming a general sampling design, an optimum ratio estimator for the population proportion P_A is defined. This estimator is based on a linear combination between the ratio estimators \hat{p}_r and $\hat{p}_{r.q}$. We also show that the proposed optimum estimator coincides with the regression type estimator.

The optimum ratio estimator

Assuming a general sampling design and the ratio estimators \hat{p}_r and $\hat{p}_{r.q}$ previously defined, the following class of estimator can be defined

$$(1) \quad \hat{p}_{r.w} = w\hat{p}_r + (1 - w)\hat{p}_{r.q}.$$

It can be easily seen that the optimum value for w in the sense of minimum variance into the class of estimators $\hat{p}_{r.w}$ is

$$(2) \quad w_{opt} = \frac{AV(\hat{p}_{r.q}) - cov(\hat{p}_r, \hat{p}_{r.q})}{AV(\hat{p}_r) + AV(\hat{p}_{r.q}) - 2cov(\hat{p}_r, \hat{p}_{r.q})},$$

where the asymptotic variance of \hat{p}_r is given by

$$AV(\hat{p}_r) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i - RB_i}{\pi_i} \frac{A_j - RB_j}{\pi_j}.$$

Analogously, the asymptotic variance of $\hat{p}_{r.q}$ is given by

$$AV(\hat{p}_{r.q}) = AV(\hat{q}_r) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i^c - R_c B_i^c}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j}.$$

A^c and B^c denote the complementary attribute of A and B . The covariance between the ratio estimators \hat{p}_r and $\hat{p}_{r.q}$ is

$$cov(\hat{p}_r, \hat{p}_{r.q}) = -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i - RB_i}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j}.$$

Proof

Using Taylor series (see Särndal et al., 1992, pg 178) under a general sampling design, the ratios \widehat{R} and \widehat{R}_c can be approximated as

$$\widehat{R} \cong R + \frac{1}{NP_B} \sum_{i \in s} \frac{A_i - RB_i}{\pi_i}$$

and

$$\widehat{R}_c \cong R_c + \frac{1}{NQ_B} \sum_{i \in s} \frac{A_i^c - R_c B_i^c}{\pi_i}.$$

From previous expressions we obtain

$$\widehat{p}_r \cong P_A + \frac{1}{N} \sum_{i \in s} \frac{A_i - RB_i}{\pi_i}$$

and

$$\widehat{q}_r \cong Q_A + \frac{1}{N} \sum_{i \in s} \frac{A_i^c - R_c B_i^c}{\pi_i},$$

which can be used to obtain the covariance between \widehat{p}_r and $\widehat{p}_{r,q}$:

$$\begin{aligned} cov(\widehat{p}_r, \widehat{p}_{r,q}) &= cov(\widehat{p}_r, 1 - \widehat{q}_r) = -cov(\widehat{p}_r, \widehat{q}_r) = \\ &= -cov\left(P_A + \frac{1}{N} \sum_{i \in s} \frac{A_i - RB_i}{\pi_i}, Q_A + \frac{1}{N} \sum_{i \in s} \frac{A_i^c - R_c B_i^c}{\pi_i}\right) = \\ &= -\frac{1}{N^2} cov\left(\sum_{i \in s} \frac{A_i - RB_i}{\pi_i}, \sum_{i \in s} \frac{A_i^c - R_c B_i^c}{\pi_i}\right) = \\ &= -\frac{1}{N^2} cov\left(\sum_{i=1}^N \frac{A_i - RB_i}{\pi_i} I_i, \sum_{i=1}^N \frac{A_i^c - R_c B_i^c}{\pi_i} I_i\right) = \\ &= -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N cov\left(\frac{A_i - RB_i}{\pi_i} I_i, \frac{A_j^c - R_c B_j^c}{\pi_j} I_j\right) = \\ &= -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \frac{A_i - RB_i}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j} cov(I_i, I_j) = \\ &= -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i - RB_i}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j}, \end{aligned}$$

where I_i is the sample indicator, i.e., $I_i = 1$ if $i \in s$ and $I_i = 0$ otherwise.

We observe that w_{opt} depends on unknown quantities. An estimator of w_{opt} is given by

$$\widehat{w}_{opt} = \frac{\widehat{V}(\widehat{p}_{r,q}) - \widehat{cov}(\widehat{p}_r, \widehat{p}_{r,q})}{\widehat{V}(\widehat{p}_r) + \widehat{V}(\widehat{p}_{r,q}) - 2\widehat{cov}(\widehat{p}_r, \widehat{p}_{r,q})},$$

where

$$\widehat{V}(\widehat{p}_r) = \frac{1}{N^2} \sum_{i \in s} \sum_{j \in s} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \frac{A_i - RB_i}{\pi_i} \frac{A_j - RB_j}{\pi_j},$$

$$\widehat{V}(\widehat{p}_{r,q}) = AV(\widehat{q}_r) = \frac{1}{N^2} \sum_{i \in s} \sum_{j \in s} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \frac{A_i^c - R_c B_i^c}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j}$$

and

$$\widehat{cov}(\widehat{p}_r, \widehat{p}_{r,q}) = -\frac{1}{N^2} \sum_{i \in s} \sum_{j \in s} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_{ij}} \frac{A_i - RB_i}{\pi_i} \frac{A_j^c - R_c B_j^c}{\pi_j}.$$

The estimator of the optimum weight w_{opt} can be used to define an optimum estimator for P_A , which is given by

$$\hat{p}_{r.opt} = \hat{w}_{opt}\hat{p}_r + (1 - \hat{w}_{opt})\hat{p}_{r.q}.$$

An alternative expression of the optimum weight w_{opt} is

$$(3) \quad w_{opt} = \frac{R_c - \beta}{R_c - R},$$

where $R_c = P_A/P_B$,

$$\beta = \frac{cov(\hat{p}_A, \hat{p}_B)}{V(\hat{p}_B)},$$

$$cov(\hat{p}_A, \hat{p}_B) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i}{\pi_i} \frac{B_j}{\pi_j}$$

and

$$V(\hat{p}_B) = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\pi_{ij} - \pi_i \pi_j) \frac{A_i}{\pi_i} \frac{A_j}{\pi_j}.$$

Proof

Variances and the covariance in (2) can be expressed as

$$V_1 = AV(\hat{p}_r) = V(\hat{p}_A) + R^2V(\hat{p}_B) - 2Rcov(\hat{p}_A, \hat{p}_B),$$

$$V_2 = AV(\hat{p}_{r.q}) = V(\hat{q}_A) + R_c^2V(\hat{q}_B) - 2R_c cov(\hat{q}_A, \hat{q}_B) =$$

$$= V(\hat{p}_A) + R_c^2V(\hat{p}_B) - 2R_c cov(\hat{p}_A, \hat{p}_B)$$

and

$$C = cov(\hat{p}_r, \hat{p}_{r.q}) = V(\hat{p}_A) + RR_cV(\hat{p}_B) - (R + R_c)cov(\hat{p}_A, \hat{p}_B).$$

The numerator and the denominator given by (2) can be expressed as

$$V_2 - C = V(\hat{p}_B)(R_c^2 - RR_c) - cov(\hat{p}_A, \hat{p}_B)[2R_c - (R - R_c)] =$$

$$= V(\hat{p}_B)R_c(R_c - R) - cov(\hat{p}_A, \hat{p}_B)(R_c - R) =$$

$$= (R_c - R)[V(\hat{p}_B)R_c - cov(\hat{p}_A, \hat{p}_B)]$$

and

$$V_1 + V_2 - 2C = V(\hat{p}_B)(R^2 + R_c^2 - 2RR_c) - cov(\hat{p}_A, \hat{p}_B)[2R + 2R_c - 2(R + R_c)] =$$

$$= V(\hat{p}_B)(R_c - R)^2$$

By substituting these expressions into (2) we obtain

$$w_{opt} = \frac{V_2 - C}{V_1 + V_2 - 2C} = \frac{(R_c - R)[V(\hat{p}_B)R_c - cov(\hat{p}_A, \hat{p}_B)]}{V(\hat{p}_B)(R_c - R)^2} =$$

$$= \frac{V(\hat{p}_B)R_c - cov(\hat{p}_A, \hat{p}_B)}{(R_c - R)V(\hat{p}_B)} = \frac{R_c - \beta}{R_c - R}.$$

An estimator of expression (3) is

$$\hat{w}_{opt} = \frac{\hat{R}_c - \hat{\beta}}{\hat{R}_c - \hat{R}},$$

where

$$\hat{\beta} = \frac{\widehat{cov}(\hat{p}_A, \hat{p}_B)}{\widehat{V}(\hat{p}_B)}.$$

An alternative and simpler expression for the optimum ratio estimator $\hat{p}_{r.opt}$ is

$$(4) \quad \hat{p}_{r.opt} = \hat{p}_A + \hat{\beta}(P_B - \hat{p}_B).$$

Proof

$$\begin{aligned} \hat{p}_{r.opt} &= \hat{w}_{opt}\hat{p}_r + (1 - \hat{w}_{opt})\hat{p}_{r.q} = \frac{\hat{R}_c - \hat{\beta}}{\hat{R}_c - \hat{R}}\hat{R}P_B + \left(1 - \frac{\hat{R}_c - \hat{\beta}}{\hat{R}_c - \hat{R}}\right)(1 - \hat{R}_cQ_B) = \\ &= \frac{\hat{R}_c\hat{R}P_B - \hat{\beta}\hat{R}P_B}{\hat{R}_c - \hat{R}} + \frac{\hat{R}_c - \hat{R} - \hat{R}_c + \hat{\beta}}{\hat{R}_c - \hat{R}}(1 - \hat{R}_c(1 - P_B)) = \\ &= \frac{\hat{R}_c\hat{R}P_B - \hat{\beta}\hat{R}P_B - \hat{R} + \hat{R}\hat{R}_c - \hat{R}_c\hat{R}P_B + \hat{\beta} - \hat{\beta}\hat{R}_c + \hat{\beta}\hat{R}_cP_B}{\hat{R}_c - \hat{R}} = \\ &= \frac{\hat{\beta}P_B(\hat{R}_c - \hat{R}) + (\hat{R}_c - 1)(\hat{R} - \hat{\beta})}{\hat{R}_c - \hat{R}} = \hat{\beta}P_B + \frac{(\frac{\hat{q}_A}{\hat{q}_B} - 1)(\frac{\hat{p}_A}{\hat{p}_B} - \hat{\beta})}{\frac{\hat{q}_A}{\hat{q}_B} - \frac{\hat{p}_A}{\hat{p}_B}} = \\ &= \hat{\beta}P_B + \frac{1 - \hat{p}_A - 1 + \hat{p}_B \hat{p}_A - \hat{\beta}\hat{p}_B}{(1 - \hat{p}_A)\hat{p}_B - \hat{p}_A(1 - \hat{p}_B)} \frac{\hat{p}_B}{(1 - \hat{p}_B)\hat{p}_B} = \\ &= \hat{\beta}P_B + \frac{(\hat{p}_B - \hat{p}_A)(\hat{p}_A - \hat{\beta}\hat{p}_B)}{\hat{p}_B - \hat{p}_A\hat{p}_B - \hat{p}_A + \hat{p}_A\hat{p}_B} = \\ &= \hat{p}_A + \hat{\beta}(P_B - \hat{p}_B). \end{aligned}$$

We observe that the optimum regression estimator given by (4) coincides with the regression type estimator given in Rueda et al. (2011b). This relationship between the ratio and regression estimators show that the regression type estimator is more efficient than the ratio estimators \hat{p}_r and $\hat{p}_{r.q}$.

A variance estimator of the optimum ratio estimator is given by

$$\hat{V}(\hat{p}_{r.opt}) = \frac{1}{N^2} \sum_{i \in s} \sum_{j \in s} \frac{\pi_{ij} - \pi_i\pi_j}{\pi_{ij}} \frac{A_i - \hat{A}_i}{\pi_i} \frac{A_j - \hat{A}_j}{\pi_j},$$

where $\hat{A}_i = \hat{p}_A + \hat{\beta}(B_i - \hat{p}_B)$.

Acknowledgment

This work is supported by the project PYR-2010-20 (PYR-GENIL) of the "Campus de Excelencia Internacional CEI BioTIC GENIL (CEB09-0010)", which belongs to the CEI programme of the "Ministerio de Ciencia e Innovación" in Spain.

REFERENCES

Blyth, C. R., Still, H. A. (1983). Binomial confidence intervals. J. Am. Stat. Assoc. 78: 108-116.
 Cohen, G. R., Yang, S. Y. (1994). Mid-p confidence intervals for the Poisson expectation. Stat. Med. 13: 2189-2203.
 Fleiss, J. L., Levin, B., Paik, M. C. (2003). Statistical methods for rates and proportions 3rd edn. Wiley, New Jersey.

Newcombe, R. G. (1998). Two-sided confidence intervals for the single proportion: comparison of seven methods. *Stat. Med.* 17: 857-872.

Rueda, M.M., Muñoz, J.F., Arcos, A., Álvarez, E. (2011a). Estimators and confidence intervals for the proportion using binary auxiliary information with applications to the estimation of prevalences. *J. Bioph. Stat.* 21, 526-554.

Rueda, M.M., Muñoz, J.F., Arcos, A., Álvarez, E. (2011b). Indirect estimation of proportions in natural resource surveys. *Mathematics and Computers in Simulation*. In press.

Särndal, C.E., Swensson, B. Wretman, J. (1992) *Model Assisted Survey sampling*. Springer Verlag.