

## ON ASYMPTOTIC PROPERTIES OF A GENERALISED PREDICTOR OF FINITE POPULATION VARIANCE

By P. MUKHOPADHYAY

*Indian Statistical Institute*

*SUMMARY.* A predictor of a finite population variance under probability sampling suggested by a multiple regression model is shown to be asymptotically design unbiased and consistent.

### 1. INTRODUCTION

We consider estimating a finite population variance through probability sampling. Let  $U$  denote a finite population of  $N$  identifiable units labelled  $1, 2, \dots, N$  and  $y$  the character of interest taking value  $y_i$  on unit  $i, i = 1, \dots, N$ . Its variance is

$$V(\mathbf{y}) = a_1 \sum_1^N y_i^2 - a_2 \sum_{i \neq i'=1}^N y_i y_{i'} \quad \dots \quad (1.1)$$

where  $a_1 = \frac{1}{N} \left( 1 - \frac{1}{N} \right)$ ,  $a_2 = \frac{1}{N^2}$  and  $\mathbf{y} = (y_1, \dots, y_N)$ . Let a sample  $s$  be selected from  $U$  following a design  $p$ , having inclusion—probabilities  $\pi_i = \sum_{s \ni i} p(s)$ ,  $\pi_{ii'} = \sum_{s \ni i, i'} p(s)$ , etc. Let  $I_i, I_{ii'}$  denote indicator random variables with  $I_i = 1(0)$  according as unit  $i \in (\phi)s$  and  $I_{ii'} = 1(0)$  according as the pair  $(i, i') \in (\phi)s$ . Suppose auxiliary variable  $x_j$  with  $x_{ij}$  its value on unit  $i$  is available. Also assume that  $y_i$  is the realised value of a random variable  $Y_i, i = 1, \dots, N$ . We propose a predictor of  $V(\mathbf{Y})$  where  $\mathbf{Y} = (Y_1, \dots, Y_N)$  as

$$v_G(\mathbf{Y}) = a_1 \sum_{i=1}^N \frac{I_i Y_i^2}{\pi_i} - a_2 \sum_{i \neq i'=1}^N \frac{I_{ii'} Y_i Y_{i'}}{\pi_{ii'}} + \sum_{j=1}^k \beta_j \left\{ a_1 \sum_{i=1}^N \left( \frac{I_i}{\pi_i} - 1 \right) x_{ij}^2 - a_2 \sum_{i \neq i'=1}^N \left( \frac{I_{ii'}}{\pi_{ii'}} - 1 \right) x_{ij} x_{i'j} \right\}, \dots \quad (1.2)$$

here  $\beta_j$  is a function of  $\mathbf{I}, \mathbf{Y}$  and  $\mathbf{X}, \mathbf{I} = (I_1, \dots, I_N)'$ ,  $\mathbf{X} = ((x_{ij}))$  an  $N \times k$  matrix such that  $\beta_j$  when suitably assigned is computable given the data stated above. The multiple-regression model-based form (1.2) is suggested following Särndal (1980).

*AMS (1980) subject classification :* 62D05.

*Key words and phrases :* Multiple regression model, Generalised predictor, Asymptotic design unbiasedness, Asymptotic design consistency.

Following Isaki and Fuller (1982) and Robinson and Särndal (1983) we show that  $v_G(\mathbf{Y})$  is asymptotically design unbiased and consistent for  $V(\mathbf{Y})$  under conditions which do not require any modelling.

2. ASYMPTOTIC DESIGN UNBIASEDNESS AND CONSISTENCY OF THE GENERALISED PREDICTOR

Following Robinson and Särndal (1983) we define a sequence of populations  $U_t$  of increasing sizes  $N_1 < N_2 < N_3 < \dots$  such that  $U_1 \subset U_2 \subset U_3 \dots$ . Let  $\{s_t\}$  denote a sequence of samples  $s_t$  of effective size  $n_t$  drawn from  $U_t$  using sampling design  $p_t$ ,  $t = 1, 2, 3, \dots$  with  $n_1 < n_2 < n_3 < \dots$ . Let  $\pi_{it}$ ,  $\pi_{i'it}$  etc. denote inclusion-probabilities for  $p_t$ . Let also  $I_{it}$  and  $I_{i'it}$  denote corresponding indicator variables. Then we have a sequence of population values  $\{\mathbf{y}^t, \mathbf{X}^t\}$  where  $\mathbf{y}^t = (y_1, \dots, y_{N_t})$ ,  $\mathbf{X}^t = ((x_{ij}))$  is an  $N_t \times k$  matrix, a sequence of population parameters  $\{V_t(\mathbf{y}^t)\}$  and a sequence of predictors  $\{v_{G_t}(\mathbf{Y}^t)\}$  where

$$\begin{aligned}
 v_{G_t}(\mathbf{Y}^t) = & a_{1t} \sum_{i=1}^{N_t} \frac{I_{it} Y_i^2}{\pi_{it}} - a_{2t} \sum_{i \neq i'=1}^{N_t} \frac{I_{i'it} Y_i Y_{i'}}{\pi_{i'it}} \\
 & + \sum_{j=1}^k \beta_{jt} \left\{ a_1 \sum_{j=1}^{N_t} \left( \frac{I_{ij}}{\pi_{it}} - 1 \right) x_{ij}^2 \right. \\
 & \left. - a_{2t} \sum_{i \neq i'=1}^{N_t} \left( \frac{I_{i'it}}{\pi_{i'it}} - 1 \right) x_{ij} x_{i'j} \right\}, \quad \dots \quad (2.1)
 \end{aligned}$$

$a_{1t} = \frac{1}{N_t} \left( 1 - \frac{1}{N_t} \right)$ ,  $a_{2t} = \frac{1}{N_t^2}$ ,  $\beta_{jt}$  is a function of  $\mathbf{I}_t$ ,  $\mathbf{Y}^t$  and  $\mathbf{X}^t$  with  $\mathbf{I}_t = (I_{1t}, \dots, I_{N_t t})'$  and  $\mathbf{Y}^t = (Y_1, \dots, Y_{N_t})$ .

For the asymptotic analysis let  $N_t \rightarrow \infty$  as  $t \rightarrow \infty$ . Let  $\xi$  be the probability distribution of the infinite dimensional random vector  $(Y_1, Y_2, \dots)$ .

*Definition 1.*  $\{v_{G_t}\}$  is asymptotically design unbiased (ADU) if

$$\lim_{t \rightarrow \infty} E\{(v_{G_t} | \mathbf{Y}^t) - V(\mathbf{Y}^t)\} = 0$$

with  $\xi$ -probability one.

*Definition 2.*  $\{v_{G_t}\}$  is asymptotically design consistent (ADC) for  $V_t$  if given any  $\epsilon > 0$ ,

$$\lim_{t \rightarrow \infty} P\{|v_{G_t} - V_t| > \epsilon | \mathbf{Y}^t\} = 0$$

with  $\xi$ -probability one.

Here  $E$  denotes design expectation. By Markov's inequality if  $v_{Gt}$  is  $ADU$  it must be  $ADC$ .

Theorem : Under assumptions (a. 1)–(a. 9) below,  $v$  is  $ADU$  and  $ADC$ . The assumptions are :

$$(a. 1) \overline{\lim}_{t \rightarrow \infty} \frac{1}{N_t} \sum_{i=1}^{N_t} Y_i^4 < \infty \text{ with } \xi\text{-probability one.}$$

$$(a. 2) \overline{\lim}_{t \rightarrow \infty} \phi_1(t) = \infty \text{ where } \phi_1(t) = N_t \min_{1 \leq i \leq N_t} \pi_{it}.$$

$$(a. 3) \lim_{t \rightarrow \infty} \psi_1(t) = 0 \text{ where } \psi_1(t) = \max_{1 \leq i \neq i' \leq N_t} \left| \frac{\pi_{ii't}}{\pi_{ii'}\pi_{i't}} - 1 \right|$$

$$(a. 4) \lim_{t \rightarrow \infty} \phi_2(t) = \infty \text{ where } \phi_2(t) = N_t^2 \min_{1 \leq i \neq i' \leq N_t} \pi_{ii't}$$

$$(a. 5) \lim_{t \rightarrow \infty} \psi_2(t) = 0 \text{ where } \psi_2(t) = \frac{1}{N_t} \max_{1 \leq i \neq i' \neq i'' \leq N_t} \left| \frac{\pi_{ii'i''t}}{\pi_{ii'}\pi_{i'i''t}} - 1 \right|$$

$$(a. 6) \lim_{t \rightarrow \infty} \psi_3(t) = 0 \text{ where } \psi_3(t) = \max_{1 \leq i \neq i' \neq i'' \neq i''' \leq N_t} \left| \frac{\pi_{ii'i''i'''t}}{\pi_{ii'}\pi_{i'i''i'''t}} - 1 \right|$$

$$(a. 7) \lim_{t \rightarrow \infty} \psi_4(t) = 0 \text{ where } \psi_4(t) = \max_{1 \leq i \neq i' \neq i'' \leq N_t} \left| \frac{\pi_{ii'i''t}}{\pi_{ii'}\pi_{i'i''t}} - 1 \right|$$

$$(a. 8) \overline{\lim}_{t \rightarrow \infty} \frac{1}{N_t} \sum_{i=1}^{N_t} x_{ij}^4 < \infty \text{ for } j = 1, 2, \dots, k.$$

$$(a. 9) \overline{\lim}_{t \rightarrow \infty} E \left( \sum_{j=1}^k \beta_{jt}^2 \right) < \alpha \text{ with } \xi\text{-probability one.}$$

Proof : We have

$$v_{Gt} - V_t = C_t(\mathbf{y}) + \sum_{j=1}^k \beta_{jt} C_t(x_j) \quad \dots \quad (2.2)$$

where

$$C_t(\mathbf{y}) = a_{1t} \sum_{i=1}^{N_t} Y_i^2 \left( \frac{I_{it}}{\pi_{it}} - 1 \right) - a_{2t} \sum_{i \neq i'}^{N_t} Y_i Y_{i'} \left( \frac{I_{ii't}}{\pi_{ii't}} - 1 \right)$$

and  $C_t(x_j)$  is defined similarly. Hence

$$E\{|v_{Gt} - V_t| | \mathbf{Y}^t\} \leq \sqrt{E(C_t^2(\mathbf{y}) | \mathbf{y}^t)} + \sqrt{E \left( \sum_{j=1}^k \beta_{jt}^2 | \mathbf{Y}^t \right) E \left( \sum_{j=1}^k C_t(x_j)^2 \right)} \quad \dots \quad (2.3)$$

Now

$$\begin{aligned}
 E(C_{it}^2(\mathbf{y}) | \mathbf{Y}^t) &= a_{1t}^2 \left[ \sum_{i=1}^{N_t} Y_i^4 \left( \frac{1}{\pi_{it}} - 1 \right) + \sum_{i \neq i'=1}^{N_t} Y_i^2 Y_{i'}^2 \right. \\
 &\quad \left. \left( \frac{\pi_{ii't}}{\pi_{it}\pi_{i't}} - 1 \right) \right] + a_{2t}^2 \left[ 2 \sum_{i \neq i'=1}^{N_t} Y_i^2 Y_{i'}^2 \left( \frac{1}{\pi_{ii't}} - 1 \right) \right. \\
 &\quad + 4 \sum_{i \neq i' \neq i''=1}^{N_t} Y_i^2 Y_{i'} Y_{i''} \left( \frac{\pi_{ii'i''t}}{\pi_{ii't} \pi_{i'i''t}} - 1 \right) \\
 &\quad + \sum_{i \neq i' \neq i'' \neq i'''=1}^{N_t} \sum \sum Y_i Y_{i'} Y_{i''} Y_{i'''} \left( \frac{\pi_{ii'i''i'''t}}{\pi_{ii't} \pi_{i'i'''t}} - 1 \right) \left. \right] \\
 &\quad - 2 a_{1t} a_{2t} \left[ 2 \sum_{i \neq i'=1}^{N_t} Y_i^3 Y_{i'} \left( \frac{1}{\pi_{it}} - 1 \right) \right. \\
 &\quad \left. + \sum_{i \neq i' \neq i''=1}^{N_t} \sum \sum Y_i^2 Y_{i'} Y_{i''} \left( \frac{\pi_{ii'i''t}}{\pi_{it}\pi_{i'i''t}} - 1 \right) \right] \dots \quad (2.4)
 \end{aligned}$$

The first term in (2.4) is dominated by

$$\frac{1}{N_t} \sum_{i=1}^{N_t} \frac{Y_i^4}{\phi_1(t)}$$

and  $\rightarrow 0$  as  $t \rightarrow \infty$  with  $\xi$ -probability one under assumptions (a. 1) and (a.2). The subsequent terms also tend to 0 with  $\xi$ -probability one as  $t \rightarrow \infty$  under (a. 1)–(a. 7). Hence  $E(C_{it}^2(\mathbf{y}) | \mathbf{Y}^t) \rightarrow 0$  with  $\xi$ -probability one as  $t \rightarrow \infty$ . Similarly under (a. 2)–(a. 8),  $E \sum_{j=1}^k C_{jt}^2(x_j) \rightarrow 0$  as  $t \rightarrow \infty$ . These coupled with the assumption (a. 9) prove that  $v_{Gt}$  is ADU and ADC.

*Note :* The assumptions (a. 2), (a. 4), (a. 5) imply  $n_t \rightarrow \infty$  as  $t \rightarrow \infty$ . All the assumptions (a. 2)–(a. 7) are satisfied for simple random sampling.

*Acknowledgement.* Thanks are due to a referee for some suggestions for improvement in presentation.

REFERENCES

ISAKI C. T. and FULLER, W. A. (1982): Survey design under the regression super-population model. *Jour. Am. Stat. Assoc.* **77**, 89-96.

MUKHOPADHYAY, P. (1986): Asymptotic properties of a generalised predictor of finite population variance under probability sampling. *Ind. Stat. Instt. Tech. Rep.* No. ASC/86/19.

ROBINSON, P. M. and SÄRNDALE, C. E. (1983): Asymptotic properties of the generalised regression estimator in probability sampling. *Sankhyā B*, **45**, 240-248.

SÄRNDALE, C. E. (1980): On  $\pi$ -inverse weighting versus best linear unbiased weighting in probability sampling. *Biometrika*, **67**, 639-650.

*Paper received :* October, 1987.

*Revised :* November, 1988.