

Efficiency of Standard Regression Model-based Ratio-synthetic Estimators in Sample Surveys Combining Time Series and Cross-sectional Data

P. D. Ghangurde

Consultant, 1370 Plante Drive, Ottawa, ON K1V9G3

Abstract

Efficiency of ratio-synthetic estimator compared to BLUP estimator, both based on cross-sectional data, depends on ratios of harmonic means of auxiliary variable values of sample units to sample sizes in respective domains, domain of interest U_i and complementary domain U_c . Assuming equal unit error variances in model-based domain estimation BLUP estimator of domain mean or total is less efficient when compared to ratio-synthetic estimator (Ghangurde, P. D.; JSM (2014)). In this paper efficiencies of ratio-synthetic estimators based on time series data as compared to the same based on cross-sectional data are analyzed for various values of correlation coefficient in AR(1) model for errors in the regression model assuming 5 and 10 time-points or occasions with diagonal covariance matrix for sampling errors. These efficiencies are not affected by small area domain effects in the mixed model assumed to derive BLUP estimator. Empirical results on gains in efficiencies of ratio-synthetic are compared to those of BLUP estimator (pages 159-62; Rao, J.N.K. (2003)). Empirical results can be similarly obtained for ratio-synthetic estimators in household surveys with rotation sample designs.

Key Words: Efficiency, ratio-synthetic, model-based, time-series, cross-sectional, data

1. Introduction

In BLUP extension of ratio-synthetic estimator $\hat{\bar{X}}_i$, \bar{X}_i is known mean of auxiliary variable (e.g. business or household income) values of units in the domain of interest U_i and $\hat{\beta}$ is estimator of β , regression coefficient in the mixed model. In household surveys the auxiliary variable is population of areal units such as primary sampling units or clusters of households within strata. Evaluation of variance of ratio-synthetic as compared to m.s.e. of BLUP estimator was done for several sample sizes and harmonic means in the domain of interest U_i and complementary domain U_c , unit error variances and variances of domain effects in the mixed model. The evaluation has provided new theoretical and empirical results on efficiencies of BLUP and

Prepared for presentation at the meetings of the Survey Research Methods section of the ASA August 8 - 13, 2015, Seattle, USA.

ratio-synthetic estimators assuming sample survey framework. The important conclusion from the study is that in domain estimation in sample surveys assuming two domains U_i and U_c , equal unit error variances in the model ratio-synthetic estimator is more efficient than BLUP estimator based on the

mixed model. In ratio-synthetic estimator \bar{X}_i , β , β and its variance are based on standard regression model. Alternative methods for obtaining total population in the domain of interest U_i in inter-census period have also been briefly reviewed in a paper on the study (Ghangurde, P. D.; JSM (2014)).

Whatever methodology is used in practice to obtain auxiliary variable and population totals the theoretical and empirical results on efficiency comparison of ratio-synthetic and BLUP estimators referred to above hold good, since the totals are assumed to be known.

Efficiencies of ratio-synthetic estimators are not affected by small area domain effects in the mixed model assumed in the case of BLUP estimators (pages 135–38; Rao, J.N.K. (2003)). In this paper efficiencies of ratio-synthetic estimators, based on time-series data vs cross-sectional data, have been analyzed by taking into consideration efficiencies of ratio-synthetic as compared to BLUP estimators both based on cross-sectional data (Table 3; Ghangurde, P. D; JSM(2014)).

2. Ratio-synthetic Estimation Using Time-series vs Cross-sectional Data

Let population U consist of N units; a sample of n units is drawn from U by simple random sampling without replacement. Let $n_i (> 0)$ and $n_c=(n - n_i)$ be units in the sample from domains U_i and U_c respectively. Let N_i and $N_c = N - N_i$ be total number of units in domains U_i and U_c respectively.

We assume that n_i sample units from U_i and $(n - n_i)$ sample units from U_c satisfy the model :

$$Y_{ijt} = \chi_{ijt} \beta + \epsilon_{ijt} ; t = 0, 1, \dots, T; j = 1, \dots, n_i; j \in U_i;$$

$$Y_{cjt} = \chi_{cjt} \beta + \epsilon_{cjt} ; t = 0, 1, \dots, T; j = (n_i + 1), \dots, n; j \in U_c, (2.1)$$

where β is regression coefficient, ϵ_{ijt} and ϵ_{cjt} are i. i. d. errors; χ_{ijt} and χ_{cjt} are x-values of j th sample unit for $t = 0, 1, \dots, T$; Y_{ijt} and Y_{cjt} are y-values of j th sample unit for $t = 0, 1, \dots, T$ in U_i and U_c respectively. We assume $Y_{ijt} \geq 0$ and $Y_{cjt} \geq 0$ and $\chi_{ijt} > 0$ and $\chi_{cjt} > 0$, which includes the case $\chi_{ijt} = \chi_{cjt} = 1$ for sample surveys. We assume that

$E(\epsilon_{ijt}) = E(\epsilon_{cjt}) = 0$; $V(\epsilon_{ijt}) = V(\epsilon_{cjt}) = \sigma^2$. Thus equal unit error variances are assumed in the model (2.1).

When auxiliary variable is quantitative (counts or continuous) unit variances are unequal and proportional to their values, which was the assumption made in derivation of model-based domain and ratio-synthetic estimators (Ghangurde P. D.; JSM (2012)).

The model (2.1) when $t = 0$ is the model for cross-sectional data. Given \bar{X}_i as known mean of auxiliary variable values for U_i the ratio-synthetic estimator is defined as :

$$\hat{\bar{X}}_i = \beta \bar{X}_i, \tag{2.2}$$

where $\hat{\beta}$ is estimator of regression coefficient β under model (2.1) based on n_i and $(n-n_i)$ sample units from U_i and U_c respectively for cross-sectional data. It is given by

$$\hat{\beta} = \frac{\sum X_i V_i X_i + \sum X_c V_c X_c}{\sum X_i V_i Y_i + \sum X_c V_c Y_c}. \tag{2.3}$$

The variance of $\hat{\beta}$ based on cross-sectional data is given by

$$V(\hat{\beta}) = \frac{\sum X_i V_i X_i + \sum X_c V_c X_c}{\sum X_i V_i X_i + \sum X_c V_c X_c}, \tag{2.4}$$

where $V_i = \frac{\sigma^2}{n_i}$; $V_c = \frac{\sigma^2}{n-n_i}$;

$$X_i = [\text{row}(\chi_{ij0})]_{1 \leq j \leq n_i}; \quad X_c = [\text{row}(\chi_{cj0})]_{(n_i+1) \leq j \leq n}. \tag{2.5}$$

Thus

$$V(\hat{\beta}) = \frac{\sigma^2}{\sum_{j=1}^{n_i} \chi_{ij0}^2 + \sum_{j=(n_i+1)}^n \chi_{cj0}^2}. \tag{2.6}$$

is the variance of $\hat{\beta}$ for cross-sectional data assuming auxiliary variable values greater than zero. The estimator of regression coefficient is given by

$$\hat{\beta} = \frac{\sum_{j=1}^{n_i} Y_{ij0} \chi_{ij0} + \sum_{j=(n_i+1)}^n Y_{cj0} \chi_{cj0}}{\sum_{j=1}^{n_i} \chi_{ij0}^2 + \sum_{j=(n_i+1)}^n \chi_{cj0}^2}. \tag{2.7}$$

However, for sample surveys $\chi_{ij0} = \chi_{cj0} = 1$ and $V(\hat{\beta}) = \sigma^2 / n$. This is the case of interest for evaluation of efficiency of ratio-synthetic estimator which uses time-series data in AR(1) model for sampling errors in model (2.1) as compared to that based on cross-sectional data.

For time-series data of T time points we will assume a first order autoregressive process for sampling errors $\epsilon_{ijt}, \epsilon_{cjt}; t=1, \dots, T$. The estimator $\hat{\beta}$ of β for occasion $t=0$ is based on the same sample units. Assuming values of correlation coefficient $\rho, |\rho| < 1$, in the AR(1) model for sampling errors in the standard regression model reduces variance of $\hat{\beta}$ and increases efficiency. The AR(1) model is an extension for sample units in U_i and U_c in sample survey framework of model (5.4.8) (see page 83; Rao, J.N.K.(2003)) and is as follows:

$$\begin{aligned} \epsilon_{ijt} &= \rho \epsilon_{ij,t-1} + e_{ijt}; t = 1, \dots, T; j=1, \dots, n_i; j \in U_i; \\ \epsilon_{cjt} &= \rho \epsilon_{cj,t-1} + e_{cjt}; t = 1, \dots, T; j=(n_i+1), \dots, n; j \in U_c, \end{aligned} \quad (2.8)$$

where $E(e_{ijt}) = 0; E(e_{cjt}) = 0; V(e_{ijt}) = V(e_{cjt}) = \sigma^2$. The covariance matrices of $\mathbf{\epsilon}_i = (\epsilon_{ij1}, \dots, \epsilon_{ijT})$ and $\mathbf{\epsilon}_c = (\epsilon_{cj1}, \dots, \epsilon_{cjT})$ are the same for each sample unit in U_i and U_c due to the same ρ assumed in U_i and U_c . Thus for n_i and $(n-n_i)$ sample units in U_i and U_c respectively we have block diagonal covariance matrices $\sigma^2 \Lambda_i = \sigma^2 \Lambda_c$, where $\Lambda_i = \Lambda_c = \Lambda$.

The (t, s) th element of Λ is given by

$$\rho^{|t-s|} / (1 - \rho^2); t = 1, \dots, T; s = 1, \dots, T. \quad (2.9)$$

Thus in the case of time-series of T time-points $t = 1, \dots, T$ variance of $\hat{\beta}$, as extension of (2.4) under model (2.1) with $t=0$ to AR(1) model (2.8), is:

$$V(\hat{\beta}) = \left[X_i \text{diag} \left[\chi_{ij} V_{ij} \chi_{ij}' \right] X_i' + X_c \text{diag} \left[\chi_{cj} V_{cj} \chi_{cj}' \right] X_c' \right], \quad (2.10)$$

$1 \leq j \leq n_i \qquad (n_i+1) \leq j \leq n$

where $X_i = [\text{row}(\chi_{ij0})]; X_c = [\text{row}(\chi_{cj0})]; \chi_{ij} = [\text{row}(\chi_{ijt})]; \chi_{cj} = [\text{row}(\chi_{cjt})];$
 $1 \leq j \leq n_i \qquad (n_i+1) \leq j \leq n \qquad 1 \leq t \leq T \qquad 1 \leq t \leq T$

$$V_{ij} = \sigma^2 [I + \Lambda]; \quad V_{cj} = \sigma^2 [I + \Lambda], \quad (2.11)$$

where V_{ij} and V_{cj} have been defined for AR(1) model (2.8) for sampling errors and I is $T \times T$ identity matrix. Let variances

$\chi_{ij}^{-1} V_{ij} \chi_{ij}^{-2} = \sigma^{-2} S(i,j,T, \rho)$ and $\chi_{cj}^{-1} V_{cj} \chi_{cj}^{-2} = \sigma^{-2} S(c,j,T, \rho)$, which are scalars based on data of T time-points, χ_{ij} and χ_{cj} . The notation $\sigma^{-2} S(i,j,T, \rho)$ and $\sigma^{-2} S(c,j,T, \rho)$ for these variances indicates that these are obtained for given T and ρ for j th sample unit in U_i and U_c respectively. Thus

$$V(\hat{\beta}) = [W_i X_i + W_c X_c]^{-1}$$

where $W_i = \sigma^{-2} [\text{row}(\chi_{ij}^{-1} S(i,j,T, \rho))]_{1 \leq j \leq n_i}$; $W_c = \sigma^{-2} [\text{row}(\chi_{cj}^{-1} S(c,j,T, \rho))]_{(n_i+1) \leq j \leq n}$.

Hence

$$V(\hat{\beta}) = \sigma^{-2} \left[\sum_{j=1}^{n_i} \chi_{ij}^{-1} S(i,j,T, \rho) + \sum_{j=(n_i+1)}^n \chi_{cj}^{-1} S(c,j,T, \rho) \right]. \quad (2.12)$$

3. Efficiency of Ratio-synthetic Estimator under AR(1) Model for Time-series Data as Compared to Cross-sectional Data

The efficiency of ratio-synthetic estimator \bar{X}_i based on time-series data as compared to that based on cross-sectional data for $t=0$, assuming known

\bar{X}_i , is defined as:

$$\text{Efficiency} = \frac{V(\hat{\beta})_0}{V(\hat{\beta})_T} \quad (3.1)$$

$$= \frac{[\sum_{j=1}^{n_i} \chi_{ij}^{-1} S(i,j,T, \rho) + \sum_{j=(n_i+1)}^n \chi_{cj}^{-1} S(c,j,T, \rho)]}{[\sum_{j=1}^{n_i} \chi_{ij}^{-1} + \sum_{j=(n_i+1)}^n \chi_{cj}^{-1}]} \quad (3.2)$$

If error variances are unequal with $V(\epsilon_{ijt}) = \chi_{ijt}^2 \sigma^2$ and $V(\epsilon_{cjt}) = \chi_{cjt}^2 \sigma^2$; $t=1, \dots, T$ for j th sample units in U_i and U_c respectively, $S(i, j, T, \rho)$ and $S(c, j, T, \rho)$ will have to be redefined for appropriate time-series model. The case of unequal unit error variances has not been considered in this paper.

In the case of sample surveys $\chi_{ij0} = 1$; $j = 1, \dots, n_i$; $j \in U_i$ and $\chi_{cj0} = 1$; $j = (n_i+1), \dots, n$; $j \in U_c$ under unified model (see Ghangurde, P.D.; JSM (2014)).

Also, $\chi_{ijt} = \chi_{cjt} = 1$; $t=1, \dots, T$; $j \in U_i$ and $j \in U_c$; $\chi_{ij} = \chi_{cj} = \mathbf{1}$ are row vectors and $\chi_{ij} = \chi_{cj} = \mathbf{1}$ are column vectors of T elements equal to 1.

Thus $S(i, j, T, \rho) = \sigma^2 [1 \ V_{ij} \ 1]^{-1}$ and $S(c, j, T, \rho) = \sigma^2 [1 \ V_{cj} \ 1]^{-1}$. Substituting these values in (3.2) efficiency of ratio-synthetic estimator for time series as compared to cross-sectional data is:

$$\text{Efficiency} = \frac{\sigma^2 [n_i [1 \ V_{ij} \ 1]^{-1} + (n - n_i) [1 \ V_{cj} \ 1]^{-1}}{[n_i + (n - n_i)]}, \quad (3.3)$$

where $V_{ij} = \sigma^2 [I + \Lambda]$ and $V_{cj} = \sigma^2 [I + \Lambda]$ for any j th sample unit in U_i and U_c . Thus (3.3) reduces as

$$\text{Efficiency} = \text{Total of values of terms in } [I + \Lambda]^{-1}, \quad (3.4)$$

where I is $T \times T$ identity matrix and Λ is as defined in (2.9).

For $T=5$ and $T=10$ substituting $\rho = 0.3, 0.5, 0.7$ and 0.9 efficiencies were obtained. These are presented in Table 1 below. The efficiencies are greater than those for BLUP estimator (pages 159 – 62; Rao, J. N. K.(2003)).

Table 1: Efficiencies of Ratio-synthetic Estimator under AR(1) Model in Sample Surveys for Time-series as Compared to Cross-sectional Data

ρ	T = 5	T = 10
0.3	1.81	3.45
0.5	1.28	2.26
0.7	0.72	1.14
0.9	0.22	0.26

The domain effects assumed in the model to derive BLUP estimator reduce efficiencies of BLUP estimator as compared to ratio-synthetic estimator in

sample surveys assuming $\sigma^2 = 5 \sigma_i^2 = 5 \sigma_c^2$. For sample size $n = 100$, sample sizes in U_i , $n_i = 2$ to 80 , efficiency of ratio-synthetic as compared to BLUP estimator varies from 13.15 to 1.37 (see Table 3; Ghangurde; P. D.; JSM (2014)).

If we assume $\sigma = \sigma_i = \sigma_c$ ratio-synthetic estimator for cross-sectional data is expected to be even more efficient. There are no studies on efficiency of ratio-synthetic estimator as compared to BLUP estimator in Rao, J.N.K.(2003).

The efficiencies are greater for $T = 10$ than for $T = 5$. The efficiencies decrease for increasing values of ρ and are less than 1 for $T = 5$ and $\rho = 0.7$ and 0.9 and $T = 10$ and $\rho = 0.9$.

Due to sample survey framework assumed in the models (2.1) and (2.8) efficiencies can be obtained for assumed values of ρ without simulation.

4. Concluding Remarks

The covariance matrix for errors in the standard regression model is diagonal due to assumption about errors made in model-based domain, synthetic, ratio-synthetic and BLUP estimation as in previous papers. By assuming sample survey framework the model for domain estimation gives basic theoretical and empirical results on relationship between estimators used in sample surveys (see Ghangurde, P. D. ; JSM (2012) and JSM (2014)).

In sample surveys efficiency is obtained by substituting values $\chi_{ijt} = \chi_{cjt} = 1$; $t=1, \dots, T$; $j \in U_i$ and $j \in U_c$ in (3.2). Thus we have simpler formula for efficiency (3.4). For rotation sample designs used in household surveys empirical values in covariance matrix with correlated errors can be used in place of matrix $[I + \Lambda]$ in (3.4) to obtain efficiencies. The AR(1) model-based covariance matrix will not be needed.

References

- [1] Ghangurde, P. D. (2012), Small area estimation in household surveys when auxiliary variable totals are known, Presented at the Joint Statistical Meetings of the ASA, San Diego, July 28 – August 2, 2012 .
- [2] Ghangurde, P. D. (2014), Evaluation of efficiency of standard regression mixed model-based BLUP estimators in household surveys assuming unequal error variances, Presented at the Joint Statistical Meetings of the ASA, Boston, August, 3 - 8, 2014.
- [3] Rao, J. N. K. (2003), Small Area Estimation, Wiley-Interscience .