

Standard Regression Model-based Ratio-synthetic Estimators Assuming Unequal or Equal Unit Error Variances and Their Use in Survey Practice

P. D. Ghangurde

Consultant, 1370 Plante Drive, Ottawa, ON K1V9G3

Abstract

Assuming sample survey framework of two domains, domain of interest U_i and complementary domain U_c in sample design strata, unequal unit error variances proportional to auxiliary variable values, an assumption appropriate in household surveys, ratio-synthetic estimator was proved to be more efficient than mixed model-based BLUP estimator. Two components of approximate efficiency of ratio-synthetic estimator were derived, assuming known domain population totals or quantitative auxiliary variable totals (Ghangurde, P. D. (2014)). In this paper unit error variances are assumed to be equal, an assumption appropriate in most other sample surveys. Approximate efficiency of ratio-synthetic estimator is derived by unconditional analysis. The results in the case of sample surveys under unified model are similar to those in earlier paper. Approximate efficiencies under two models are compared assuming that auxiliary variable has lognormal distribution. Ratio-synthetic is more efficient than BLUP in sample surveys under standard regression model assuming unequal or equal unit error variances; it is simple to use as compared to BLUP. Some applications in survey practice and methods to estimate domain totals and means are reviewed.

Key Words: Efficiency, ratio-synthetic, model-based, BLUP, survey, practice

1. Introduction

Standard regression mixed model-based BLUP estimator is an extension of regression-synthetic estimator, obtained by addition of estimator of small area domain effect η_i in the mixed model for domain of interest U_i . Thus BLUP

estimator of domain mean of y-variable $\bar{\mu}_i$ based on vector of p (≥ 1) domain

means of auxiliary variable values \bar{X}_i and p regression coefficients β_k is:

$$\hat{\bar{\mu}}_i = \bar{X}_i \hat{\beta} + \hat{\eta}_i; \quad i = 1, 2, \dots, m, \quad (1.1)$$

where $\hat{\beta}$ is vector of estimators $\hat{\beta}_k$; $k = 1, 2, \dots, p$ in the mixed model and $\hat{\eta}_i$ is estimator of η_i . In the framework of Type-B model (see pages 135–137; Rao, J.N.K.(2003)) there are m (≥ 2) domains called small areas unlike in the

Prepared for presentation at the meetings of the Survey Research Methods Section of the ASA, July 30 – August 4, 2016, Chicago, USA.

framework of domain estimation in sampling theory, where domain of interest U_i and complementary domain U_c are two domains within each stratum or strata enclosing U_i (Hartley, H. O. (1959); pages 34-38; Cochran, W.G. (1977)). The mixed model-based BLUP estimator has been used in survey practice since 1990's; however, results on relationship between standard regression model-based domain estimator, synthetic estimator and ratio-synthetic estimator derived from it were first proved by the author (Ghangurde, P. D. (2012) and (2013)). Under standard regression model ratio-synthetic estimator was obtained from synthetic estimator by a simple substitution. It was also proved that under the model ratio-synthetic estimation makes the same well-known assumption as made in synthetic estimation.

Assuming a random sample of size n units from a stratum of sample design of a household survey, $n_i (> 0)$ sample units from U_i and $(n-n_i)$ sample units from

U_c , known domain mean \bar{X}_i of one auxiliary variable, efficiency of ratio-synthetic as compared to BLUP estimator was first derived by the author assuming model with unequal unit error variances. The approximate efficiency is determined by ratios of harmonic means of auxiliary variable values of sample units to sample sizes and by ratio of total sample size to that in domain U_i (Ghangurde, P. D. (2014)).

For high values of ratios \bar{X}_{ih}/n_i and $\bar{X}_{ch}/(n-n_i)$, where \bar{X}_{ih} and \bar{X}_{ch} are harmonic means of x -values of sample units in U_i and U_c respectively, m.s.e of BLUP is close to variance of ratio-synthetic estimator but is still greater than the latter. For low values of these ratios of harmonic means, as in sample surveys, ratio-synthetic estimator is far more efficient than BLUP estimator (Ghangurde, P. D. (2014)). This also holds in other sample surveys of households (e.g. Survey of Household Spending and LFS sample design-based surveys). These household surveys use the same estimation and variance estimation methodology as used in the LFS.

Standard regression model with equal unit error variances, appropriate in model-based domain estimation in most other surveys, is introduced in Section 2. In Section 3 approximate efficiency of ratio-synthetic estimator in sample surveys derived under unified model has the same expression as in Ghangurde, P. D. (2014), except that the second term is multiplied by square of inverse sampling ratio. When the auxiliary variable is quantitative the efficiency is obtained by assuming lognormal distribution and substituting its moments as approximations for means and means of squares of sample x -values from two domains. The differences in empirical results under these two models can be attributed to assumptions about errors made in the models. The use of lognormal as a distribution of auxiliary variable in model-based domain estimation gives these remarkable empirical results.

In view of efficiency of ratio-synthetic estimator as compared to BLUP estimator under standard regression model with unequal and also equal unit

error variances in sample surveys use of ratio-synthetic estimation in sample surveys is appropriate. In Section 4 its use in survey practice and methods involved are reviewed. In Section 5 some concluding remarks are given.

2. The Model and Model-based Estimators

We assume a population U of N units. A sample of size n units is drawn from U by simple random sampling without replacement. Let $n_i (>0)$ be number of sample units from N_i units in U_i and $(n-n_i)$ be number of sample units from $(N-N_i)$ units in U_c . The mixed model for domain estimation assuming one auxiliary variable is defined as :

$$\begin{aligned}
 Y_{ij} &= X_{ij} \beta + \eta_i + \epsilon_{ij}; j = 1, \dots, n_i; j \in U_i; \\
 Y_{cj} &= X_{cj} \beta + \eta_c + \epsilon_{cj}; j = (n_i+1), \dots, n; j \in U_c, \quad (2.1)
 \end{aligned}$$

where Y_{ij} and Y_{cj} are y -values of j th sample unit from U_i and U_c and X_{ij} and X_{cj} are x -values of j th sample unit from U_i and U_c respectively. We assume $Y_{ij} \geq 0$; $Y_{cj} \geq 0$; $X_{ij} > 0$; and $X_{cj} > 0$. It includes the case $X_{ij} = X_{cj} = 1$, which implies that j th units in U_i and U_c are in the sample. The assumptions $X_{ij} = 1$ and $X_{cj} = 1$ under unified model (Ghangurde, P. D. (2014)) make analysis of efficiency possible for surveys in which unit error variances are equal, following the same approach as in LFS sample design-based surveys in which unit error variances are unequal. The estimator of regression coefficient β and its variance are derived in this paper assuming that sampling errors ϵ_{ij} and ϵ_{cj} are independent of random domain effects η_i and η_c . We assume that

$$E(\eta_i) = E(\eta_c) = 0; V(\eta_i) = \sigma_i^2; V(\eta_c) = \sigma_c^2; E(\epsilon_{ij}) = E(\epsilon_{cj}) = 0; V(\epsilon_{ij}) =$$

$V(\epsilon_{cj}) = \sigma_e^2$. This was also the model for domain estimation in evaluation of efficiency of ratio-synthetic estimator in sample surveys combining time-series and cross-sectional data (see Ghangurde, P. D. (2015)). We now

derive estimator $\hat{\beta}$ and its variance $V(\hat{\beta})$ under the standard model (2.1)

assuming equal unit error variances and $\eta_i = \eta_c = 0$. Thus

$$\hat{\beta} = [X_i' V_i^{-1} X_i + X_c' V_c^{-1} X_c]^{-1} [X_i' V_i^{-1} Y_i + X_c' V_c^{-1} Y_c], \quad (2.2)$$

where

$$Y_i = [\text{row}(Y_{ij})]; \quad Y_c = [\text{row}(Y_{cj})]; \\
 1 \leq j \leq n_i \quad (n_i+1) \leq j \leq n$$

$$X_i = \begin{bmatrix} \text{row} (X_{ij}) \\ 1 \leq j \leq n_i \end{bmatrix}; \quad X_c = \begin{bmatrix} \text{row} (X_{cj}) \\ (n_i+1) \leq j \leq n \end{bmatrix};$$

$$V_i = \sigma_e \begin{bmatrix} I_{n_i} \\ -1 \quad -2 \end{bmatrix}; \quad V_c = \sigma_e \begin{bmatrix} I_{(n-n_i)} \\ -1 \quad -2 \end{bmatrix}.$$

Also,

$$V(\hat{\beta}) = \begin{bmatrix} X_i' V_i^{-1} X_i + X_c' V_c^{-1} X_c \\ 0 \end{bmatrix}^{-1} \quad (2.3)$$

Thus

$$\hat{\beta} = \left[\sum_{j=1}^{n_i} Y_{ij} X_{ij} + \sum_{j=(n_i+1)}^n Y_{cj} X_{cj} \right] / \left[\sum_{j=1}^{n_i} X_{ij}^2 + \sum_{j=(n_i+1)}^n X_{cj}^2 \right] \quad (2.4)$$

$$V(\hat{\beta}) = \sigma_e^2 \begin{bmatrix} \sum_{j=1}^{n_i} X_{ij}^2 + \sum_{j=(n_i+1)}^n X_{cj}^2 \\ 0 \end{bmatrix}^{-1} \quad (2.5)$$

assuming

$$X_{ij} > 0; j=1, \dots, n_i; j \in U_i; \quad X_{cj} > 0; j = (n_i+1), \dots, n; j \in U_c.$$

The ratio-synthetic estimator assuming one auxiliary variable and known

domain mean \bar{X}_i is:

$$\hat{\mu}_i = \bar{X}_i \hat{\beta} \quad (2.6)$$

and its variance under standard regression model (2.1) assuming $\eta_i = \eta_c = 0$ is:

$$\bar{X}_i^2 V(\hat{\beta}).$$

Under the mixed model (2.1) assuming N_i is large and n_i/N_i is small the BLUP estimator is given by

$$\hat{\mu}_i = \bar{X}_i \hat{\beta} + \hat{\eta}_i, \quad (2.7)$$

where $\hat{\beta}$ is estimator of β and $\hat{\eta}_i$ is estimator of η_i based on the mixed model (2.1).

We derive below m.s. error of $\hat{\mu}_i$ under the mixed model and then derive efficiency of ratio-synthetic estimator under the model (2.1) assuming $\eta_i = \eta_c = 0$.

The estimator of small area domain effect η_i for the case of equal unit variances is

$$\hat{\eta}_i = \xi_i (\bar{Y}_{ia} - \bar{X}_{ia} \hat{\beta}). \quad (2.8)$$

This is the case $k_{ij} = \delta_{ij} = 1$ (page 135-36; Rao, J.N.K. (2003)). The derivation of $\hat{\beta}$ and $V(\hat{\beta})$ is given later. Thus in (2.8)

$$\bar{Y}_{ia} = \left[\sum_{j=1}^{n_i} Y_{ij} \right] / n_i; \text{ since } a_{ij} = 1 \text{ in the model with equal unit error variances,}$$

$$a_{i.} = \sum_{j=1}^{n_i} a_{ij} = n_i. \text{ Similarly, } \bar{X}_{ia} = \left[\sum_{j=1}^{n_i} X_{ij} \right] / n_i \text{ and } \bar{X}_{ca} = \left[\sum_{j=(n_i+1)}^n X_{cj} \right] / (n - n_i).$$

The optimal weights derived for the two domains in the mixed model (2.1) are:

$$\xi_i = \frac{\sigma_i^2}{\sigma_i^2 + \sigma_e^2 / n_i} = 1 / \left[1 + (\sigma_e^2 / \sigma_i^2) / n_i \right]; \quad (2.9)$$

$$\xi_c = 1 / \left[1 + (\sigma_e^2 / \sigma_c^2) / (n - n_i) \right]. \quad (2.10)$$

Substituting for $\hat{\eta}_i$ from (2.8) in (2.7)

$$\hat{\mu}_i = (\bar{X}_{i.} - \xi_i \bar{X}_{ia}) \hat{\beta} + \xi_i \bar{Y}_{ia}. \quad (2.11)$$

Conditioning on sample x -values $E[\hat{Y}_{ij} - E(Y_{ij})] = E[\hat{\eta}_i + \epsilon_{ij}] = 0$.

Using the expression for $\hat{\beta}$ given later the two terms in (2.11) can be proved to have zero covariance.

$$\text{Also, } V[\hat{Y}_{ij}] = E[\hat{\eta}_i + \epsilon_{ij}]^2 = [\sigma_i^2 + \sigma_e^2 / n_i] = \sigma_e^2 \left[1 + (\sigma_i^2 / \sigma_e^2) / n_i \right].$$

Using the last result we have

$$\text{MSE}(\hat{\mu}_i) = (\bar{X}_{i.} - \xi_i \bar{X}_{ia})^2 V(\hat{\beta}) + \xi_i^2 \sigma_e^2 \left[1 + (\sigma_i^2 / \sigma_e^2) / n_i \right]. \quad (2.12)$$

The term (7.1.13), page 176; Rao, J.N.K and Molina, I.(2015) is incorrect.

$V(\hat{\beta})$, variance of $\hat{\beta}$ under the mixed model (2.1), is given by

$$V(\hat{\beta}) = \left[X_i' V_i^{-1} X_i + X_c' V_c^{-1} X_c \right]^{-1}, \quad (2.13)$$

where

$$V_i = \sigma_e^2 \left[I_{n_i} + \left(\frac{\sigma_i}{\sigma_e} \right)^2 \frac{1_{n_i} 1_{n_i}'}{n_i} \right];$$

$$V_c = \sigma_e^2 \left[I_{(n-n_i)} + \left(\frac{\sigma_c}{\sigma_e} \right)^2 \frac{1_{(n-n_i)} 1_{(n-n_i)'}}{(n-n_i)} \right],$$

where 1_{n_i} and $1_{(n-n_i)}$ are column vectors of ones. Using a result on matrix inversion :

$$[A + uv']^{-1} = A^{-1} - [A^{-1} u v' A^{-1}] / [1 + v' A^{-1} u],$$

$$V_i^{-1} = \sigma_e^{-2} \left[I_{n_i} - \xi_i \frac{1_{n_i} 1_{n_i}'}{n_i} \right];$$

$$V_c^{-1} = \sigma_e^{-2} \left[I_{(n-n_i)} - \xi_c \frac{1_{(n-n_i)} 1_{(n-n_i)'}}{(n-n_i)} \right].$$

We derive $\hat{\beta}$ under the mixed model (2.1) and prove $E(\hat{\beta}) = \beta$, as done in Ghangurde,P.D.(2014).

$$\hat{\beta} = [X_i V_i^{-1} X_i + X_c V_c^{-1} X_c]^{-1} [X_i V_i^{-1} Y_i + X_c V_c^{-1} Y_c]. \quad (2.14)$$

Substituting for V_i and V_c in $\hat{\beta}$ we have

$$\hat{\beta} = [X_i X_i' - \xi_i X_i' \frac{1_{n_i} 1_{n_i}'}{n_i} X_i + X_c X_c' - \xi_c X_c' \frac{1_{(n-n_i)} 1_{(n-n_i)'}}{(n-n_i)} X_c]^{-1} [X_i Y_i - \xi_i X_i' \frac{1_{n_i} 1_{n_i}'}{n_i} Y_i + X_c Y_c - \xi_c X_c' \frac{1_{(n-n_i)} 1_{(n-n_i)'}}{(n-n_i)} Y_c]$$

$$= \left[\sum_{j=1}^{n_i} X_{ij} X_{ij}' - \xi_i \sum_{j=1}^{n_i} X_{ij} X_{ia} X_{ia}' + \sum_{j=(n_i+1)}^n X_{cj} X_{cj}' - \xi_c \sum_{j=(n_i+1)}^n X_{cj} X_{ca} X_{ca}' \right] \text{ multiplied by}$$

$$\left[\sum_{j=1}^{n_i} X_{ij} Y_{ij} - \xi_i \sum_{j=1}^{n_i} X_{ij} X_{ia} Y_{ia} + \sum_{j=(n_i+1)}^n X_{cj} Y_{cj} - \xi_c \sum_{j=(n_i+1)}^n X_{cj} X_{ca} Y_{ca} \right]. \quad (2.15)$$

Conditioning on sample x-values it can be seen that $E(\hat{\beta}) = \beta$ and that covariance between two terms in (2.11) is zero. Also,

$$V(\hat{\beta}) = \sigma_e^2 \left[X_i X_i' - \xi_i X_i' \frac{1_{n_i} 1_{n_i}'}{n_i} X_i + X_c X_c' - \xi_c X_c' \frac{1_{(n-n_i)} 1_{(n-n_i)'}}{(n-n_i)} X_c \right]^{-1} \sigma_e^2 \left[\sum_{j=1}^{n_i} X_{ij} X_{ij}' - \xi_i \sum_{j=1}^{n_i} X_{ij} X_{ia} X_{ia}' + \sum_{j=(n_i+1)}^n X_{cj} X_{cj}' - \xi_c \sum_{j=(n_i+1)}^n X_{cj} X_{ca} X_{ca}' \right]. \quad (2.16)$$

3. Efficiency of Ratio-synthetic Estimators Assuming Equal (and Unequal) Unit Error Variances

The efficiency defined as in Ghangurde, P. D (2014) assuming known \bar{X}_i , is:

$$\begin{aligned} \text{MSE}(\hat{\mu}_i) / \bar{X}_i^2 V(\hat{\beta}) &= [1 - \xi_i (\bar{X}_{ia} / \bar{X}_i)]^2 \frac{V(\hat{\beta})}{m} \\ &+ \frac{\xi_i^2 \sigma_e^2 [1 + (\sigma_i / \sigma_e)^2]}{\bar{X}_i^2 V(\hat{\beta}) n_i} \end{aligned} \quad (3.1)$$

Given a partition of population U into two non-empty domains U_i and U_c

expectation of sample mean \bar{X}_{ia} is domain mean \bar{X}_i ; thus $\bar{X}_{ia} / \bar{X}_i = 1$ approximately. Substituting for

$V(\hat{\beta})$ and $V(\hat{\beta})$ approximate efficiency is:

$$\begin{aligned} [1 - \xi_i]^2 \frac{[\sum_{j=1}^{n_i} X_{ij}^2 + \sum_{j=(n_i+1)}^n X_{cj}^2]}{[\sum_{j=1}^{n_i} X_{ij}^2 - \xi_i n_i \bar{X}_{ia}^2 + \sum_{j=(n_i+1)}^n X_{cj}^2 - \xi_c (n-n_i) \bar{X}_{ca}^2]} \\ + \xi_i^2 \frac{[\sum_{j=1}^{n_i} X_{ij}^2 / \bar{X}_i^2 + \sum_{j=(n_i+1)}^n X_{cj}^2 / \bar{X}_i^2]}{n_i} [1 + (\sigma_i / \sigma_e)^2] \end{aligned} \quad (3.2)$$

We call $[1 - \xi_i]$ and ξ_i as weights and weighted components as terms. In case of surveys of incomes of households or businesses when interest is in fitting regression models to estimate expenditures on various categories from income (e.g. expenditures on food from household income or employee wages and salaries from business income) the auxiliary variable distribution is uni-modal and positively skew. In Economics and Business two or three parameter lognormal distributions have been studied as models for household and business income. The sample of n units in model (2.1) can be assumed to be drawn from

lognormal distribution with parameters μ and σ . This lognormal distribution has r th moment about the origin $E(x^r) = \exp(r\mu + [r\sigma]^2/2)$. For some μ and $\sigma = 1$, $E(x) = e^{\mu + 1/2}$ and $E(x^2) = e^{2\mu + 2}$ (see pages 3 - 9; Crow, E.L. and Shimizu, K. (1988)). Also, $E(n_i \bar{X}_{ia}) = n_i e^{\mu}$; $E([n-n_i] \bar{X}_{ca}) = [n-n_i] e^{\mu}$; $E[\sum_{j=1}^{n_i} X_{ij}] = e^{\mu} n_i$; $E[\sum_{j=(n_i+1)}^n X_{cj}] = e^{\mu} (n - n_i)$, where X_{ij} and X_{cj} are x -values of sample units from U_i and U_c respectively.

Substituting \bar{X}_{ia} for \bar{X}_i and $e^{2\mu+1}$ for \bar{X}_{ia} approximate efficiency by unconditional inference is:

$$[1 - \xi_i] \frac{e^{2\mu} [e^{n_i} + e^{(n-n_i)}]}{e^{2\mu} [e^{n_i} - \xi_i n_i e^{\mu}] + e^{2\mu} [e^{(n-n_i)} - \xi_c (n-n_i) e^{\mu}]}$$

$$+ \frac{\xi_i e^{2\mu} [e^{n_i} + e^{(n-n_i)}]}{n_i e^{2\mu+1}} [1 + (\sigma_i / \sigma_e)^2].$$

Also, $E(x^{-1}) = e^{-\mu - 1/2}$; $E(x^{-2}) = e^{-2\mu - 2}$; $E(\bar{X}_{ih}) = E(\bar{X}_{ch}) = e^{\mu - 1/2}$ and $E(\bar{X}_{ih}^2) = e^{2\mu - 2}$. These will be used in the model with unequal unit error variances. After simplification approximate efficiency under standard regression model with equal unit error variances assuming lognormal distribution with unspecified μ and $\sigma = 1$ is :

$$[1 - \xi_i] \frac{[n/n_i]}{[1 - (\xi_i / e) + [(n/n_i) - 1][1 - (\xi_c / e)]]} + \xi_i e^{2\mu} [n/n_i] [1 + (\sigma_i / \sigma_e)^2]. \tag{3.4}$$

The optimal weights ξ_i and ξ_c in the first term of (3.4) are reduced by factor $(1/e)$ and the second term is increased by factor e as compared to the same terms in approximate efficiency (3.5) below under standard regression model with unequal unit error variances. The efficiencies under

two models differ due to different error assumptions. The lognormal with $\sigma = 1$ has unimodal longtailed skew distribution appropriate for household or business income. The efficiency (3.5) depends on parameter μ since ξ_{iu} and ξ_{cu} depend on \bar{X}_{ih} and \bar{X}_{ch} and second term depends on $\bar{X}_{ih} / \bar{X}_{ih2}$. The approximate efficiency under standard regression model with unequal unit error variances is obtained by substituting for \bar{X}_{ih} , \bar{X}_{ch} and \bar{X}_{ih2} , the last one being harmonic mean of squares of x-values of sample units from U_i , their expectations assuming lognormal with $\mu = 11.0$ and $\sigma = 1$. Since for these values $e = 60,000$ it is realistic as distribution of household incomes and appropriate for computation of terms in (3.5). In Ghangurde, P. D.(2014) values were assumed for \bar{X}_{ih} and \bar{X}_{ch} to evaluate efficiencies. Now $\bar{X}_{ih} / \bar{X}_{ih2} = E(\bar{X}_{ih}) / E(\bar{X}_{ih2}) = e^{-9.5} = 0.000075$, giving approximate efficiency:

$$[1 - \xi_{iu}] \frac{[n/ni]}{[1 - \xi_{iu}] + [(n/ni) - 1][1 - \xi_{cu}]} + \xi_{iu} [n/ni] [1 + (\sigma_i / \sigma_e) (\bar{X}_{ih} / \bar{X}_{ih2})], \tag{3.5}$$

where optimal weights ξ_{iu} and ξ_{ic} are approximated assuming the lognormal;

$$\xi_{iu} = 1 / [1 + (\sigma_e / \sigma_i) (e / ni)] \text{ and } \xi_{cu} = 1 / [1 + (\sigma_e / \sigma_c) (e / (n - ni))] \text{ and } \sigma_e = 5\sigma_i = 5\sigma_c; (\sigma_i / \sigma_e) = 0.2. \text{ Thus second term of (3.5) reduces to } \xi_{iu} [n/ni] [1.000015].$$

Due to approximation of optimal weights ξ_{iu} , ξ_{cu} and $\bar{X}_{ih} / \bar{X}_{ih2}$ in (3.5) the differences in efficiencies based on two models can be attributed to differences in the assumptions about errors in these two models. The efficiencies based on (3.4) and (3.5) are presented in Tables 1 and 3.

We now consider efficiencies in sample surveys under unified model defined in Ghangurde, P. D.(2014). Assuming $X_{ij} = X_{cj} = 1$ for units in sample and zero outside sample term (3.2) is the same under two models. Also, in sample surveys optimal weights ξ_i and ξ_c are the same under two models.

Since U_i is a random domain from population U , $E[\bar{X}_i] = \bar{X} = [n / N]$ and second term (3.3) is $\xi_i [N / n] [n/ni] [1 + (\sigma_i / \sigma_e)]$ by unconditional inference.

We have assumed $\sigma_e = 5 \sigma_i$ giving approximate efficiency in sample surveys:

$$[1 - \xi_i] \frac{2}{[1 - \xi_i] + [1 - \xi_c] [(n/n_i) - 1]} \frac{n/n_i}{2} + \xi_i [N/n] \frac{2}{[n/n_i]} [1.2] \quad (3.6)$$

The second term in the model with unequal unit error variances, assuming

$\xi_{ij} = \xi_{cj} = 1$ in sample surveys is $\xi_i [N/n] \frac{2}{[n/n_i]} [1.2]$; see Ghangurde, P. D. (2014).

In (3.6) second term is greater due to factor $[N/n]$ making ratio-synthetic estimator more efficient than BLUP in sample surveys under model with equal unit error variances (see Tables 2 and 4).

In household surveys the use of ratio-synthetic estimator in survey practice assumes known domain totals of households and persons. In other sample surveys, where regression models are fitted for estimation, domain means of quantitative auxiliary variable (e.g. household income) are assumed to be known, which is the main issue in use of ratio-synthetic estimators in survey practice. In ratio-synthetic and BLUP estimators use of administrative data-based domain totals or means of auxiliary variable has been suggested in the literature on small area estimation. However, appropriate changes in processing of administrative data and other issues have not been addressed (see Ghangurde P.D. (2013); Rao J.N.K. and Molina, I. (2015)).

Since ratio-synthetic estimator is more efficient than BLUP under unified model for sample surveys assuming equal or unequal unit error variances, it should be used in survey practice. The use of estimates of household and population totals for domains based on counts of households in sample clusters and counts of persons in households, obtained in the first month in survey data collection, is outlined in brief in Section 4. This can be expanded based on survey practice.

4. Ratio-synthetic Estimation in Survey Practice

In household surveys with area sample designs clusters of households are sampling units within strata and households within clusters are the last stage sampling units. Weighted household and population totals within sample clusters in domains within strata can be used in ratio-synthetic estimation to obtain estimated domain totals X_i for households and persons. There are about 200 households per cluster in the LFS (see Methodology of the Canadian LFS; Statistics Canada (2008)). Cluster weights are multiplied by inverse sampling ratios within individual clusters to obtain household weights. In

strata with extreme growth or decline of population clusters may be re-designed and households listed. For details on procedure of household listing see Guide to the Labour Force Survey; Statistics Canada (2012).

In the LFS listing of households in clusters in rotation groups to be introduced in the survey and updating household weights is part of household survey operations. In domains of interest within strata listing of households in clusters not in the survey sample can be done by survey operations staff as an extension of sample cluster listing done a few months in advance of the first survey month when sample clusters in a rotation group are introduced. In the case of LFS a much better option would be matching clusters with areal units in the last Census of Population to obtain household counts for clusters in domains and updating these at the time of introduction of rotation groups in the sample by using Address Register (AR). The advantage of obtaining household totals for domains as part of ongoing survey operations or using AR to update cluster lists is that the statistical agency would be responsible for these activities and can ensure data quality.

Methods can be developed for multipliers to obtain population counts for domains in rural areas, suburbs, urban core and fringe areas from counts of LFS sample household members collected in the first month that rotation groups are in the survey. These methods can be refined over time and multipliers can be verified. These multipliers can be applied to weighted household counts in sample clusters to obtain estimated population in various areas identified within domain U_i enclosed by strata .

This approach would be more objective than using administrative data for related auxiliary variables from other government agencies without studies on changes in processing data to obtain estimates and whether increased correlations of these x-variables with y-variable are enough to offset substantial reduction in number of matched administrative records available (e.g increased correlation of unemployed with unemployment insurance recipients vs substantial reduction in persons eligible for unemployed insurance benefits). These issues with the use of administrative data have not been addressed in Rao, J.N.K.and Molina, I. (2015). There are other issues in the area of data quality control in agencies providing data (Section 4; Practical issues; Ghangurde, P. D. (2013)).

Domain totals X_i are usually not known in inter-census period. Ratio-synthetic estimation transforms the problem of updating synthetic weights in inter-census period to much simple one of estimation of domain totals for households and persons. Domains are parts of geographic strata. Ratio-synthetic estimation would be done within individual strata which enclose parts of domain of interest U_i . Estimation and variance estimation of other strata is not changed; also X_i are not needed for these strata, since there is no domain estimation in these strata. Household and population totals for these strata without domains would not need updating until re-design.

Thus estimation of domain mean $\bar{\mu}_i$, based on weighted data in a stratified sample design, would reduce to identifying strata enclosing parts of domain

and obtaining estimates $\hat{\beta}$ for the strata multiplied by estimated domain totals X_i based on the survey estimation methodology. The ratio-adjustment of sub-provincial estimates is a required step in household surveys based on LFS area sample design (see Ghangurde, P.D. and Gray, G.B. (1981)). Thus ratio-synthetic estimation would be more efficient and simple to implement as compared to BLUP as an integral part of estimation and variance estimation system for sub-provincial areas of household surveys.

Ratio-synthetic estimation can be extended to time-series data in household surveys with rotation sample designs; it is simpler than extension of BLUP. In a study on efficiency of ratio-synthetic estimators under standard regression model assuming equal unit error variances and 5 and 10 time-points or occasions AR(1) model was used (Ghangurde, P. D.(2015)). Time series models are useful in theoretical studies of ratio-synthetic and BLUP estimators. In practice, correlations in LFS survey estimates for characteristics over survey months can be based on matched households or persons in each rotation group; these correlations are higher for characteristics employed and in-labour force than for unemployed. Ratio-synthetic estimators for a domain total or mean and their standard errors can be obtained by using matched records over the last five months. The estimators can be combined for six rotation groups with proportion of matched households used as weights for estimates of each group providing efficiency gains in ratio-synthetic estimation in the LFS based on time-series of survey data of six months as compared to that for the current month.

The use of time-series models assuming equal or unequal unit error variances is not needed in the case of LFS, since correlations in survey data for a characteristic can be estimated from matched sample records of households and persons using survey variance estimation system. The same is true for other LFS area sample design-based surveys.

In a sample survey of construction industry in Nova Scotia, designed to obtain model-based estimates of total wages and salaries using gross business income as auxiliary variable ratio-synthetic estimator was more efficient than BLUP estimator (see example (7.3.2), pages 189 – 92; Rao, J.N.K. and Molina, I (2015)). The results in the study are based on simulation from a known population; thus domain total based on this population is known. However, in survey practice total income in a domain would have to be obtained from a recent census or survey.

In the case of some sample surveys an auxiliary variable highly correlated with estimation variable can be constructed as an index based on several correlated variables for which census data are available each year. National Agricultural Statistics Service (NASS) conducts annual survey of agricultural operations for estimation of average cash rental rates based on weighted survey data from a stratified sample survey. The construction of the index involved assumptions based on census data for these correlated auxiliary variables (see Berg, E., Cecere, W. and Ghosh, M. (2014)). The study used BLUP estimation for area level mixed model. Ratio-synthetic estimation can be used as efficient and simpler method of unit level domain estimation in a sample survey in which model with unequal unit or equal unit error variances is appropriate. The theory and empirical results to support the use of ratio-synthetic estimator instead of BLUP estimator in domain estimation at unit level are available in this paper and in Ghangurde, P. D. (2014).

5. Concluding Remarks

The important result proved in this paper shows that under standard regression model with equal unit error variances approximate efficiency of ratio-synthetic as compared to BLUP estimator in sample surveys is even greater than that under model with unequal unit error variances. The results on efficiencies in the case of quantitative auxiliary variable under models with equal and unequal unit error variances assuming lognormal distribution for quantitative variable are similar. Under model with equal unit error variances ratio-synthetic is more efficient than BLUP as compared to the model with unequal unit error variances in surveys in which regression models are fitted to estimate components of auxiliary variable total in a domain for several categories.

In both surveys the problem in survey practice is that of obtaining domain totals or means for the quantitative variable. Under unified model in the case of household surveys methods to estimate total households and population in domains in inter-census period were reviewed. However, in the case of other

surveys total for domains from recent census for quantitative auxiliary variable or estimate from recent survey can be used. Although ratio-synthetic estimator is more efficient and simple to use in survey practice than BLUP, for estimation of domain totals for auxiliary variables in sample surveys and domain means of quantitative auxiliary variables in surveys in which regression models are fitted for domain estimation, methods based on survey practice have to be developed.

Even so, ratio-synthetic estimation is simple and efficient method of providing survey sponsors and data analysts estimates and their standard errors for a range

of values of domain means \bar{X}_i and totals X_i . In most situations assuring one true value for mean or total seems not possible. We would be able to provide in household surveys estimate of persons in a domain U_i based on first month responses from households in clusters rotating in the sample. This is expected to be more accurate than estimates based on administrative data from local area jurisdictions. In case of counts of households in domains in strata, the counts based on sample clusters would be estimates. If the counts are listed by survey personnel or obtained from AR these would have bias due to under-coverage, which is called slippage in the LFS.

In practice, most surveys have stratified sample designs with different weights for sample units between strata. These weights are adjusted for different non-response rates between strata. Since lognormal distribution was used for approximations in both models results are comparable. Extension of results to weighted data for ratio-synthetic and BLUP estimators under standard regression model is only of theoretical interest. The use of ratio-synthetic estimation instead of BLUP in sample surveys is important because it is more efficient than BLUP and simple to use in survey practice. The domains of interest are likely to be enclosed by a few strata within urban or rural areas of a province. Thus estimation and variance estimation for ratio-synthetic estimates would make no assumptions than those made in methodology of estimation and variance estimation in these surveys and would be an integral part of the methodology for sub-provincial estimation in these surveys.

The assumption of lognormal distribution for auxiliary variable in obtaining efficiencies of ratio-synthetic estimators under models with unequal and equal unit error variances made it possible to attribute differences in efficiencies to different assumptions about errors. Efficiencies in the case of model with unequal unit error variances can be obtained by assuming a value of μ giving

median income e^{μ} close to that based on sample surveys as done in this study.

TABLE 1
Efficiencies based on sample size n ; auxiliary variable $U = 1$ with lognormal distribution; parameters unspecified μ and $\sigma = 1$
Standard regression model with equal unit error variances : $\sigma_e^2 = 5(\sigma_i^2 \text{ or } \sigma_c^2)$
Sample in U $n = 100$; unit error variance $\sigma_e^2 = 5(\sigma_i^2 \text{ or } \sigma_c^2)$

Sample in U_i : n_i	Optimal ξ_i	ξ_i / e	Optimal ξ_c	ξ_c / e	First term	Second term	Total Efficiency
2	0.28571	0.10513	0.95146	0.35002	0.7790	13.31	14.09
4	0.44444	0.16350	0.95050	0.34967	0.4692	16.11	16.58
6	0.54545	0.20066	0.94949	0.34930	0.3132	16.11	16.42
8	0.61538	0.22638	0.94845	0.34891	0.2238	15.44	15.66
10	0.66667	0.24525	0.94737	0.34859	0.1679	14.50	14.67
12	0.70588	0.25966	0.94623	0.34810	0.1306	13.54	13.67
50	0.90909	0.33441	0.90909	0.33443	0.0124	5.39	5.40
60	0.92308	0.33958	0.88889	0.32700	0.0089	4.64	4.65
80	0.94118	0.34624	0.80000	0.29430	0.0052	3.61	3.62

TABLE 2
Efficiencies based on sample size n in sample surveys; auxiliary variable = 1 for units in sample and = 0 for units not in sample; $U = 1$
Standard regression model with equal unit error variances : $\sigma_e^2 = 5(\sigma_i^2 \text{ or } \sigma_c^2)$
Sample in U $n = 100$; unit error variance $\sigma_e^2 = 5(\sigma_i^2 \text{ or } \sigma_c^2)$
Inverse Sampling Ratio $N/n = 10$

Sample in U_i : n_i	Optimal ξ_i	Optimal ξ_c	First term	Second term	Total Efficiency
2	0.28571	0.95146	8.25	490	498.25
4	0.44444	0.95050	4.43	580	584.43
6	0.54545	0.94949	2.76	595	597.76
8	0.61538	0.94845	1.89	568	569.89
10	0.66667	0.94737	1.38	533	534.38
12	0.70588	0.94623	1.05	498	499.05
50	0.90909	0.90909	0.09	198	198.09
60	0.92308	0.88889	0.06	170	170.06
80	0.94118	0.80000	0.04	133	133.04

TABLE 3

Efficiencies based on sample size n, with auxiliary variable with 2 lognormal distribution with parameters $\mu = 11.0$ and $\sigma = 1$; Standard regression model with unequal unit error variances :

$$\bar{X}_{ij}^2 \sigma_e^2 ; \bar{X}_{cj}^2 \sigma_e^2$$

Sample in U n = 100 ; unit error variance $\sigma_e^2 = 5 (\sigma_i^2 \text{ or } \sigma_c^2)$

Sam ple in U _i :n _i	Optimal ξ_{ju}	Optimal ξ_{cu}	First term	Second term	Total Efficiency
2	0.0000110	0.0005394	1.000506	0.000000006	1.000506
4	0.0000220	0.0005285	1.000464	0.000000012	1.000464
6	0.0000330	0.0005176	1.000426	0.000000018	1.000426
8	0.0000440	0.0005066	1.000384	0.000000024	1.000384
10	0.0000551	0.0004955	1.000342	0.000000030	1.000342
12	0.0000661	0.0004845	1.000300	0.000000036	1.000300
50	0.0002754	0.0002753	0.999750	0.000000150	0.999750
60	0.0003294	0.0002203	0.999700	0.000000180	0.999700
80	0.0004403	0.0001100	0.999520	0.000000240	0.999520

TABLE 4

x-variable = 1 for units in sample; x-variable = 0 for units not in sample; Unequal unit error variances model reducing to equal in sample surveys.

Sample in U n = 100 ; unit error variance $\sigma_e^2 = 5 (\sigma_i^2 \text{ or } \sigma_c^2)$

$$\bar{X}_{ih} = \bar{X}_{ch} = 1 \text{ due to } \bar{X}_{ij} = \bar{X}_{cj} = 1 \text{ for sample units } j$$

Sample in U _i : n _i	Optimal Weight ξ_i	Optimal Weight ξ_c	First term	Second term	Total Efficiency
2	0.28571	0.95146	8.25	4.90	13.15
4	0.44444	0.95050	4.43	5.80	10.23
6	0.54545	0.94949	2.76	5.95	8.71
8	0.61538	0.94845	1.89	5.68	7.57
10	0.66667	0.94737	1.38	5.33	6.71
12	0.70588	0.94623	1.05	4.98	6.03
50	0.90909	0.90909	0.09	1.98	2.07
60	0.92308	0.88889	0.06	1.70	1.76
80	0.94118	0.80000	0.04	1.33	1.37

References

- [1] Berg, E., Cecere, W. and Ghosh, M. (2014), Small area estimation for county level farm-land cash rental rates, *Journal of Survey Statistics and Methodology*, Volume 2 / Number 1.
- [2] Cochran, W.G. (1977), *Sampling Techniques*, Wiley.
- [3] Crow, E. L. and Shimizu, K. (1988), *Lognormal Distributions, Theory and Applications*, Marcel Dekker, Inc.
- [4] Ghangurde, P. D. and Gray, G.B. (1981), Estimation for small areas in household surveys, *Communications in Statistics, Series A* (10).
- [5] Ghangurde, P. D. (2012), Small area estimation in household surveys when auxiliary variable totals are known, Presented at the Joint Statistical Meetings of the American Statistical Association, San Diego, USA.
- [6] Ghangurde, P. D. (2013), Standard regression model-based small area domain estimation in household surveys, Presented at the Joint Statistical Meetings of the ASA, Montreal, Canada.
- [7] 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, USA.
- [8] Ghangurde, P. D. (2015), Efficiency of standard regression model-based ratio-synthetic estimators in sample surveys combining time series and cross-sectional data, Presented at the Joint Statistical Meetings of the ASA, Seattle, USA.
- [9] Hartley, H. O. (1959), Analytic studies of survey data, in a Volume in honor of Corrado Gini, Istituto di Statistica, Rome, Italy.
- [10] Rao, J. N. K. (2003), *Small Area Estimation*, Wiley Interscience.
- [11] Rao, J.N.K. and Molina, I. (2015), *Small Area Estimation*, Wiley.
- [12] Statistics Canada (2008), *Methodology of the Canadian Labour Force survey*, Catalogue No. 71-526-X.
- [13] Statistics Canada (2012), *Guide to the Labour Force Survey*, Catalogue No. 71-543-G.