Pak. J. Statist. 1995 Vol. 11(3), pp 173-189

ON GENERALIZED REGRESSION ESTIMATORS OF SMALL DOMAIN TOTALS - AN EVALUATION STUDY

A. Chaudhuri and A.K. Adhikary Indian Statistical Institute, Calcutta

(Received: July, 1994 Accepted: June, 1995)

Abstract

We consider drawing a sample from a survey population to estimate the totals of a variable of interest separately for its disjoint domains of varying sizes. Horvitz and Thompson's (HT in brief, 1952) method of estimation is of course applicable using 'inverse inclusion-probabilities' as weights for the observations on the sampled units. Assuming knowledge of population values of a related variable and postulating a linear regression through the origin two alternative estimators using further weights called 'g-weights' in two different forms may be employed for a possible improvement. One of them uses all sample values to estimate a common regression 'slope' and the other uses domain-specific values alone to estimate 'domain-wise' varying 'slopes'. These synthetic and non-synthetic versions respectively of generalized regression (greg, in brief) estimators have two alternative variance estimators each, respectively 'involving' and 'free of' the g-weights. All four of them are modifications of Yates and Grundy's (YG, in brief, 1953) variance estimator of an HT estimator. As it is difficult to theoretically compare the relative efficacies of these point estimators and corresponding interval estimators for the domain totals we undertake a numerical exercise using official records and simulations to empirically evaluate their performances. A fair conclusion tends to support the synthetic greg estimators coupled with variance estimators incorporating the g-weights.

Key Words

Empirical evaluation. Regression model; Small domain statistics.

1. INTRODUCTION

Developing resolve, timely and relevant statistics relating to 'small domains' is an important current area of research in survey sampling. Many new procedures for the purpose are rapidly emerging. We shall concentrate here only on three relatively simple 'design-based' methods of estimating domain specific totals of a variable of interest on drawing a sample from a population which is the union of several disjoint domains. If y be a variable of interest with a population total Y, the

Horvitz-Thompson (HT, 1952) estimator (HTE) for Y uses the 'reciprocals of the inclusion-probabilities of the units' as the weights for the sampled observations. If x be another variable well-correlated with y and its values be known for the entire population, then one may, instead, employ Särndal's (1980) generalized regression (greg, in brief) estimator as a possible improvement upon the HTE applying further 'multipliers', called 'q-weights' on the sample observations. A simple form of it postulates a linear regression of y on x through the origin. These x-values may or may not be well-associated with the inclusion-probabilities; in general, they will not, in multivariate surveys. If the same sample is intended to be used in such a situation in deriving estimators not only for Y but also for the totals of disjoint domains of the population the HTE can be applied with an obvious modification. But the 'greg' estimator may be employed in two alternative forms. If it is plausible to fit a single regression line for the entire population, then a single slope parameter is to be estimated using all the sample observations on y. As opposed to this resulting 'synthetic' greg estimator, the 'non-synthetic' alternative to it is based on postulation of separate regression lines for the respective domains and hence involves domain-specific y-values alone while using domain-wise slope-estimators. As the non-synthetic greg estimator uses auxiliary x-values while the HTE does not, the former may outperform the latter. But if a domain-size is small, both of them may turn out poor as the level of aggregation over y-values is small for both. But if the over-all sample-size is large enough then inspite of a small domain-size, the synthetic greg estimator may yet fare quite well especially if the postulation of a 'common' regression line for all the domains is not grossly untenable. A possible middle course of postulating 'distinct' common slopes for several disjoint subsets of domains, estimating them from respective sample values and then employing partially synthetic greg estimators borrowing strength across only 'like' domains with anticipated common slopes may also be tried. But this is not often put into practice because applying diagnostic tests for identification of domains with common slopes is not quite practicable in large-scale multi-subject surveys. It is simpler to postulate a common slope for all the domains at a time.

We shall throughout assume the units to be all distinct in our sample which is of a given size. For the HTE, the variance estimator is given by Yates and Grundy (YG, say for brevity, 1953). Särndal (1982) has given two modifications of the YG formula for variance estimators of the greg estimators of a population total. One alternative uses the 'g-weights' while the other does not. They extend easily to cover the synthetic and non-synthetic greg estimators described above. If as usual the 'pivotal' formed by the "estimator minus the parameter" divided by the estimated standard error be supposed to be distributed approximately like the standardized normal deviate τ , then one may construct confidence intervals with desired nominal confidence coefficients. It is of interest to ask how good are the confidence intervals that may be formed for the domain totals by the above noted alternative procedures. To reach conclusions on theoretical grounds seems difficult. So, we resort to a numerical exercise utilizing certain official records and carrying out simulations to empirically examine the relative performances of the confidence intervals

constructed employing the above procedures. The theory is presented briefly in Section 2. The live data we use are discussed in Section 3. Our numerical findings are presented in Section 4. In Section 5 we give our recommendations pointing out why one may be inclined to favour the use of the synthetic greg estimators of domain totals and the g-weighted variance estimators in comparable situations in practice in preference to other alternatives cited in this paper.

Incidentally, we may mention that in the literature the term 'synthetic' estimator often occurs but may not denote exclusively the one we have described. One may consult Särndal, Swensson and Wretman (1992). But in each case it involves y-values outside the domain for which the y-total is to be estimated.

2. THEORY OF ESTIMATION

Suppose a survey population U = (1, ..., i, ..., N) of N identifiable individuals labelled i = 1, ..., N, is divisible into D known disjoint segments $U_d(d = 1, ..., D)$ called domains. Let y be a real variable of interest with unknown values y_i and domain totals Y_d which we intend to estimate.

Let x and z be two positive-valued variables both well and positively correlated with y and respectively $x_i, z_i (i \in U)$ be their values with totals X and Z. Let the total of x for U_d be X_d and the values $p_i = z_i/Z$ be called normed size-measures of the units, $i \in U$.

To estimate Y_d 's, let a sample s of distinct units, n(< N) in number, be chosen from U with a probability p(s) admitting positive inclusion-probabilities π_i for i and π_{ij} for pairs (i, j). Sampling separately from respective domains is considered impractical. Let $\Delta_{ij} = (\pi_i \pi_j - \pi_{ij})/\pi_{ij}$; $I_{di} = 1$ if $i \in U_d$, but = 0, else; \sum' be sum over units k in s, $\sum' \sum'$ be sum over pairs of units k, k'(k < k') in s.

The HTE for Y_d is

 $t_{Hd} = \sum' \frac{y_i}{\pi_i} I_{di}$ and its YG form of variance estimator is $v_{YGd} = \sum' \sum' \Delta_{ij} \left(\frac{y_i I_{di}}{\pi_i} - \frac{y_j I_{dj}}{\pi_j} \right)^2.$

It is often feasible to postulate a model \underline{M}_d connecting y and x permitting one to write

$$y_i = \beta_d x_i + \epsilon_i$$
 for $i \in U_d, d = 1, \ldots, D_i$

Here β_d is an unknown constant and ϵ_i 's are uncorrelated random variables with means $E_m(\epsilon_i) = 0$ and variances

 $V_m(\epsilon_i) = \sigma_i^2, i \in U; \sigma_i (> 0)$

is an unknown constant for $i \in U$.

A special case of \underline{M}_d is \underline{M} for which

$$\beta_d = \beta$$
 for every $d = 1, \ldots, D$.

Let Q_i be an arbitrarily assignable positive constant and β_d be estimated by $\hat{\beta}_{Qd} = \sum' y_i x_i Q_i I_{di} / \sum' x_i^2 Q_i I_{di}$.

Writing $e_{di} = y_i - \hat{\beta}_{Qd} x_i$, Särndal's (1980) non-synthetic greg estimator of Y_d is

$$\begin{aligned} t_d &= X_d \hat{\beta}_{Qd} + \sum' e_{di} I_{di} / \pi_i = \sum' y_i \frac{I_{di}}{\pi_i} + \hat{\beta}_{Qd} \left(X_d - \sum' x_i \frac{I_{di}}{\pi_i} \right) \\ &= \sum' \frac{y_i}{\pi_i} I_{di} g_{sdi} \text{ where the "g-weight" is} \\ g_{sdi} &= 1 + \left(X_d - \sum' \frac{x_k}{\pi_k} I_{dk} \right) \frac{x_i Q_i \pi_i}{\sum' x_k^2 Q_k I_{dk}}. \end{aligned}$$

Estimating β by

 $\hat{\beta}_Q = \sum' y_i x_i Q_i / \sum' x_i^2 Q_i$ and writing $e_i = y_i - \hat{\beta}_Q x_i$, the synthetic greg estimator for Y_d is

$$t_{sd} = X_d \hat{\beta}_Q + \sum' e_i I_{di} / \pi_i$$

= $\sum' \frac{y_i}{\pi_i} I_{di} + \hat{\beta}_Q \left(X_d - \sum' \frac{x_i}{\pi_i} I_{di} \right)$
= $\sum' \frac{y_i}{\pi_i} g'_{sdi}$ with the "g-weights"
 $g'_{sdi} = I_{di} + \left(X_d - \sum' \frac{x_k}{\pi_k} I_{dk} \right) \frac{x_i Q_i \pi_i}{\sum' x_k^2 Q_k}.$

For t_d two variance estimators given by Särndal (1982) are

$$v_2 = \sum' \sum' \Delta_{ij} \left(\frac{e_{di}I_{di}g_{idi}}{\pi_i} - \frac{e_{dj}I_{dj}g_{idj}}{\pi_j} \right)$$

and v_1 which is v_2 with g_{sdi} replaced by unity. For t_{sd} , the two corresponding variance estimators are

$$v_{s2} = \sum' \sum' \Delta_{ij} \left(\frac{e_i g'_{sdi}}{\pi_i} - \frac{e_j g'_{sdj}}{\pi_j} \right)^2$$

and v_{s1} which is v_{s2} with g'_{sdi} replaced by I_{di} .

For Q_i four choices are usual; they are $1/x_i^g$, g = 1, 2, corresponding respectively to possible simple forms of σ_i^2 as $\sigma^2 x_i^g$ with $\sigma(>0)$ as an unknown constant; the other two are $1/p_{i_i}x_i$ and $(1 - 1)/\pi_i x_i$, respectively recommended by Hájek (1971) and Brewer (1979).

For any linear estimator e for a parameter θ having a positive-valued variance estimator v, for a large sample-size n, it is usual to regard the distribution of the pivotal quantity

$$(e-\theta)/\sqrt{v}$$

as close to that of the standard normal deviate τ . This helps construction of a $100(1-\alpha)$ per cent confidence interval for θ of the form $e \pm \tau_{\alpha/2}\sqrt{v}$, with α chosen in (0,1) and $\tau_{\alpha/2}$ the $100_{\alpha/2}$ per cent point on the right tail of the distribution of τ . With θ as Y_d we may construct such confidence intervals choosing (e, v) as $(t_{Hd}, v_{YGd}), (t_d, v_j)$ and $(t_{sd}, v_{sj}), j = 1, 2$. To evaluate relative efficacies of these various confidence intervals theoretically is difficult. So, we undertake a numerical investigation considering live data illustrated in Section 3 and by simulation. If a sample is drawn by the same method a number of times, say, R, then it is customary to evaluate the following three criteria described below and labelled I,II,III to discriminate among t_{Hd}, t_d, t_{sd} coupled respectively with $v_{YGd}, v_j, v_{sj}(j = 1, 2)$. By

 $\sum_{\mathbf{r}}$ we shall denote sums over the replicates of the simulated samples and let

$$PM(e_d) = \frac{1}{R} \sum_{r} (e_d - Y_d)^2$$

denote the Pseudo mean square error of e_d which stands for t_{Hd} , t_d and t_{sd} . By v_d we shall denote v_{YGd} , v_i , v_{sj} ; j = 1, 2. The criteria are

I. "The Actual Coverage Percentage" (ACP in brief),

II. "The average Coefficient of Variation" (ACV in brief), and

III. "The Relative Efficiency" (RE in brief)

of the estimator e_d , denoted $RE(e_d)$ for e_d as t_d and t_{sd} relative to t_{Hd} . The "Actual Coverage Percentage", ACP, is the percentage of replicated samples for which the confidence interval covers Y_d . The closer it is to $100(1 - \alpha)$ the better. The "Average Coefficient of Variation", ACV, is

$$\frac{1}{R}\sum_{r}\sqrt{v}/e.$$

This reflects the length of the confidence interval. The smaller its value the better the choice of (e, v). The "Relative Efficiency" $RE(e_d)$ is defined as $[PM(t_{Hd})/FM(e_d)]^1$ for e_d as t_d and t_{sd} .

The higher its value the better is e_d relative to t_{Hd} . The HTE t_{Hd} which does not use x_i 's and is not motivated either by \underline{M}_d or \underline{M} is taken as a basic estimator in terms of which we intend to judge the efficacies of t_d and t_{sd} respectively motivated by postulation of \underline{M}_d and \underline{M} .

3. DATA BASE

The Indian Statistical Institute (ISI), Calcutta, in April 1992, consisted of 39 administrative "units" which we shall refer to as "domains" and label them arbitrarily as $1, \ldots, d, \ldots, D = 39$. The respective roll strengths of the units or the domain sizes $N_1, \ldots, N_d, \ldots, N_D$ were respectively 73, 69, 21, 13, 4, 25, 6, 10, 25, 31, 7, 29, 5, 68, 69, 6, 35, 52, 50, 127, 3, 25, 10, 11, 13, 91, 9, 8, 22, 22, 3, 14, 4, 26, 46, 21, 69, 34 and 30. For every employee are ascertained from the Accounts Office for April, 1992 his/her dearness allowance (DA), gross pay and basic pay, respectively denoted by y, x and z. We shall illustrate application of the theory of Section 2 to estimate the 'total DA earned' by "all the respective "unit" - employees" for the 39 units. More current values could be utilized but for illustration we believe we need not mind using these slightly past data which were readily obtained during an investigation.

4. NUMERICAL FINDINGS

Out of the above 1186 workers we considered drawing samples of 200 workers. A worker's basic pay which varied from about 500 to 6,000 Indian rupees was available for use as the size-measure for sample selection. We employed two alternative schemes of sampling. For one due to Lahiri (1951), in choosing n units from a population of size N on the first draw one unit is chosen with probability $p_i = z_i/Z$ and followed up by a simple random sample (SRS) without replacement (WOR) of size (n-1) from the remaining (N-1) units. Then, the inclusion-probabilities turn out to be

$$\pi_i = \frac{n-1}{N-1} + \left(\frac{N-n}{N-1}\right) p_i, i \in U.$$

Though, p_i 's for i = 1, ..., 1186 for our example vary considerably among each other the term $\binom{n-1}{N-1}$ dominates the second term so appreciably that π_i for each i is close to the constant $\frac{n-1}{N-1}$. Yet, we find this scheme useful as is clarified below. We try a second competing scheme due to Hartley and Rao (1962). This scheme randomly arranges the units of $U = (1, \ldots, i, \ldots, N)$ and then chooses circular systematically a sample of n units with probabilities proportional to z_i so as to achieve the inclusion probabilities

$$\pi_i = np_i, i \in U.$$

Here π_i 's vary appreciably among one another. The formula for π_{ij} 's for both the schemes, with only approximations for the latter are available from the respective literature cited. For both the schemes we separately take R = 100 replicates of samples, calculate $(t_{Hd}, v_{YGd}), (t_d, v_i), (t_{sd}, v_{sj}), j = 1, 2$ taking Q_i separately as $1/x_i, 1/x_i^2, 1/\pi_i x_i$ and $(1-\pi_i)/\pi_i x_i$ and construct the confidence intervals based on them in manners described in Section 2 taking $\alpha = 0.05$. Since for the Lahiri (1951) scheme π_i 's are all close to $\frac{n-1}{N-1}$, we show our findings only for Q_i as $1/x_i$ and $1/x_1^2$ because the relevant values for the other two choices are almost the same as for $1/x_i$. We observe that the results for both the schemes are closely competitive and those for the Lahiri scheme are often more impressive. The main findings for the Hartley-Rao (HR in brief) scheme are presented in Table 1 and those for the Lahiri scheme in Table 2 in self-explanatory manners. The Lahiri scheme is easier to employ and its performance is not poorer. Though the inclusion-probabilities for this scheme are not proportional to the size measures, basing the greg estimator on this scheme is not inappropriate. Hence its utility. As the values of the performance criteria turn out rather poor for domains of sizes 15 or less, we do not show them in the Tables 1-4.

Domain	(t_{Hd}, v_{YGd})	(t_d, v_1)	(t_{sd}, v_{s1})
3	(95,30.9),	(89,8.0)/(89,8.0)/(89,7.1)/(89,7.1),	(96,9.4)/(88,8.7)/(90,8.6)/(89,8.5)
69	(99, 28, 2),	(79,4.3)/(66,4.2)/(65,4.0)/(65,4.0),	(86,5.6)/(75,4.0)/(77,4.8)/(76,4.6)
21	(82,50.7),	/(56,4.2)/(55,4.1)/	/(88,5.1)/
25	(88,54.1),	(59, 1.6)/(59, 1.6)/(59, 1.6)/(59, 1.6),	(93,8.9)/(93,6.1)/(93,6.7)/(93,6.1)
25	(93,53.9),	(75,11.7)/(74,11.1)/(75,11.0)/(75,10.7),	(94,29.7)/(93,46.2)/(93,40.7)/(93,44.8)
31	(97,46.6),	_	(96,7.1)/(94,9.3)/(95,8.8)/(95,9.3)
29	(93, 47.9),	(76,1.8)/(76,1.8)/(76,1.7)/(76,1.7),	(91,8.0)/(90,5.5)/(91,6.0)/(90,5.5)
68	(75,27.4),	(72,4.3)/(67,4.3)/(66,3.9)/(63,3.8),	(91,4.2)/(77,4.2)/(82,4.2)/(79,4.2)
69	(90,25.6),	(79,4.0)/(63,3.4)/(58,3.2)/(53,3.1),	(80,3.7)/(60,3.3)/(66,3.3)/(60,3.3)
35	(95, 46.1),	(557.5)/(75,9.3)/(76,7.4)/(76,7.5),	(87,8.5)/(95,10.7)/(94,10.2)/(94,10.7)
52	(95,33.7),	(91,4.7)/(92,4.4)/(90,4.3)/(89,4.3),	(95,5.4)/(96,4.6)/(96,4.7)/(97,4.6)
50	(85,30.1),	(81,4.0)/(75,3.0)/(74,3.5)/(71,3.4),	/(87,3.6)/
127	(90,23.9),	(94,2.5)/(94,2.5)/(94,2.4)/(94,2.4),	(91,5.7)/(92,4.4)/(92,4.7)/(92,4.4)
25	(88,42.8),	(71,9.9)/(65,11.0)/(61,8.8)/(58,7.9),	(82,13.1)/(66,14.2)/(69,13.7)/(66,13.9)
91	(98,26.4),	(92,2.6)/(91,2.5)/(90,2.3)/(90,2.3),	(97,3.8)/(97,2.8)/(97,3.0)/(97,2.8)
22	(88,57.9),	(72,4.3)/(71,4.0)/(71,4.0)/(71,4.0),	(89,4.7)/(84,4.9)/(90,4.7)/(86,4.9)
22	(82,56.3),	(61, 1.9)/(61, 1.8)/(61, 1.7)/(60, 1.8),	(85,7.0)/(85,4.4)/(85,4.9)/(85,4.4)
26	(92,41.6),	(83,5.7)/(83,4.9)/(82,4.7)/(81,4.6),	(95,6.0)/(93,5.1)/(95,5.2)/(92,5.1)
46	(95,45.5),	(82,9.0)/(81,10.0)/(80,8.6)/(80,8.7),	(91,34.1)/(91,46.2)/(91,41.3)/(91,43.3)
21	(83,47.5),	(73,6.0)/(73,5.5)/(73,5.4)/(72,5.3),	(91,34.1)/(91,46.2)/(91,41.3)/(91,43.3)
	(91,22.1),	(89,4.6)/(89,3.9)/(89,3.8)/(90,3.6),	(94,3.7)/(92,3.7)/(91,3.6)/(92,7.2)
69	(91,33.4),	(85,7.2)/(85,6.7)/(85,6.3)/(85,6.2),	(94,7.2)/(92,7.2)/(93,7.2)/(94,7.2)
69 34			

 Table 1

 Table 1

 Performances of procedures based on HR scheme in terms of (ACP, ACV). Slashes separate the values respectively for choice of Q_i as $1/x_i, 1/x_i^2, 1/\pi_i x_i$ and $(1 - \pi_i)/\pi_i x_i$. Commas after parentheses separate the values for rival procedures.

2

neralized Regression Estimators

ξų.

Table 1 (Continued)

Performances of procedures based on HR scheme in terms of (ACP, ACV). Slashes separate the values respectively for choice of Q_i as $1/x_i, 1/x_i^2, 1/\pi_i x_i$ and $(1 - \pi_i)/\pi_i x_i$. Commas after parentheses separate the values for rival procedures.

Domain size	(t_d, v_2)	(t_{sd}, v_{s2})
73	(83,8.0)/(82,8.1)/(82,5.9)/(81,5.8),	(92,10.4)/(89,10.3)/(89,10.0)/(88,10.0)
69	(65,4.1)/(76,3.7)/(74,3.6)/(74,3.6),	(87,5.5)/(88,4.5)/(90,4.7)/(86,4.5)
21	(52,5.3)/(57,4.7)/(57,4.7)/(59,4.6),	(87,5.9)/(88,4.8)/(87,5.1)/(87,4.9)
25	(67, 1.3)/(66, 1.3)/(66, 1.3)/(66, 1.3),	(95,9.3)/(95,6.5)/(95,7.1)/(95,6.6)
25	(70,8.6)/(72,7.3)/(72,7.0)/(73,7.0),	(91,26.6)/(91,37.7)/(91,34.0)/(91,36.2)
31	(80,4.2)/(77,4.1)/(78,3.9)/(78,3.9),	(98,6.9)/(93,9.0)/(96,8.4)/(93,8.9)
29	(80,1.5)/(78,1.4)/(78,1.4)/(78,1.4),	(92,7.9)/(92,5.5)/(92,6.0)/(92,5.5)
68	(81,5.1)/(74,5.1)/(73,4.6)/(71,4.4),	(91,4.3)/(74,4.2)/(76,4.1)/(72,4.1)
69	(82,4.6)/(75,4.0)/(73,3.9)/(71,4.4),	(91,4.3)/(74,4.2)/(72,3.3)/(62,3.3)
35	(74,6.6)/(72,7.7)/(73,6.4)/(73,6.4),	(92,8.1)/(94,9.1)/(100,9.2)/(91,9.3)
52	(88,4.7)/(892,4.4)/(88,4.3)/(89,4.3),	(91,5.3)/(93,4.6)/(93,4.7)/(92,4.6)
50	(84,4.9)/(79,4.6)/(78,4.5)/(77,4.3),	(95,4.3)/(83,3.8)/(86,3.8)/(84,3.8)
127	(91,2.0)/(90,2.0)/(91,1.9)/(92,1.9),	(96,5.7)/(97,4.4)/(97,4.6)/(97,4.4)
25	(73,14.1)/(68,17.1)/(66,10.4)/(689.5),	(81,21.2)/(65,23.6)/(66,22.5)/(63,22.9)
91	(93,2.4)/(91,2.3)/(87,2.0)/(87,2.1),	(95,3.8)/(96,2.8)/(80,8.6)/(84,8.3)
22	(68,4.2)/(69,3.8)/(68,3.7)/(68,3.7),	(91,5.1)/(76,5.0)/(89,5.0)/(88,4.6)
22	(69,1.9)/(69,1.8)/(69,1.8)/(69,1.8),	(90,6.9)/(88,4.5)/(81,5.0)/(88,4.6)
26	(74,5.9)/(78,5.3)/(79,5.1)/(82,5.0),	(93,6.0)/(92,5.0)/(94,5.2)/(92,5.1)
46	(79,7.9)/(79,8.6)/(79,7.8)/(79,7.8),	(93,30.9)/(93,38.8)/(93,36.2)/(93,37.7)
21	(80,5.9)/(81,5.4)/(81,5.4)/(81,5.3),	(93,6.2)/(93,6.7)/(95,6.5)/(94,6.7)
69	(95,4.9)/(95,4.2)/(95,4.1)/(97,3.9),	(96,3.8)/(94,3.7)/(95,3.7)/(93,3.7)
34	(88,7.3)/(88,6.8)/(88,6.4)/(84,6.3),	(93,6.8)/(93,6.7)/(93,6.7)/(91,6.7)
30	(59,8.4)/(60,8.5)/(60,7.9)/(60,8.0),	(76,11.5)/(71,13.2)/(71,12.8)/(71,13.1)

••

Table 2

Performances of procedures based on Lahiri's scheme in terms of (ACP, ACV). Slashes separate the values respectively for choice of Q_i as $1/x_i$, $1/x_i^2$. Commas after parentheses separate the values for rival procedures.

Domain size	(t_{Hd}, v_{YGd})	(t_d, v_1)	(t_{sd}, v_{s1})
73	(96,25.8),	(78, 4.2)/(80, 4.3),	(86,4.6)/(79,4.5)
69	(94,29.4),	(77, 5.4)/(80, 7.0),	(98,5.7)/(84,5.9)
21	(84,55.5),	(52, 7.0)/(55, 7.3),	(90,8.7)/(74,9.1)
25	(93,48.2),	(57, 1.2)/(56, 1.2),	(87, 5.6)/(87, 3.8)
25	(86, 51.0),	(70, 7.2)/(71, 7.5),	(87,13.2)/(88,15.4)
31	(90,42.6),	(73, 3.4)/(75, 3.5),	(93, 4.3)/(89, 5.6)
29	· (94,45.5),	(71, 1.3)/(74, 1.4),	(92,5.3)/(89,3.5)
68	(94,28.6),	(85, 6.1)/(87, 7.0),	(95,6.5)/(91,7.0)
69	(94,28.3),	(88, 5.2)/(90, 6.2),	(95, 6.0)/(92, 6.9)
35	(93,41.3),	(76, 5.0)/(75, 7.6),	(97, 5.9)/(95, 7.1)
52	(94,33.5),	(82, 4.7)/(84, 5.1),	(96,5.9)/(97,5.6)
50	(90,35.0),	(79, 5.4)/(82, 6.0),	(89,6.2)/(81,6.6)
127	(90,19.7),	(93, 1.5)/(93, 1.5),	(91,3.1)/(92,2.4)
25	(93,49.2),	(70,10.0)/(70,12.6),	(80,14.4)/(74,16.4)
91	(95,24.4),	(78, 2.5)/(78, 2.6),	(93, 3.0)/(88, 2.7)
22	(90,51.9),	(61, 4.2)/(63, 4.3),	(89,5.8)/(75,7.1)
22	(83,54.1),	(69, 1.8)/(69, 1.9),	(82,5.2)/(82,3.4)
26	(93,48.2),	(54, 5.8)/(56, 6.6),	(87, 7.8)/(62, 8.1)
46	(93,34.6),	(80, 5.9)/(81, 6.5),	(93,21.9)/(93,25.0)
21	(87,55.1),	(63, 6.4)/(64, 6.9),	(74,9.1)/(71,10.7)
69	(93,28.5),	(84, 6.0)/(89, 7.3),	(88,7.5)/(89,8.6)
34	(95,41.0),	(85,7.8)/(87,9.2),	(89,9.4)/(89,10.5)
30	(94,43.0),	(78, 7.6)/(80, 8.0),	(84,10.2)/(84,11.6)

i

Table 2 (Continued)

Performances of procedures based on Lahiri's scheme in terms of (ACP, ACV). Slashes separate the values respectively for choices of Q_i as $1/x_i$ and $1/x_i^2$. Commas after parentheses separate the values for rival procedures.

Domain size	(t_d, v_2)	$(t_s d, v_{s2})$
73	(76,3.6)/(77,3.9),	(85,4.6)/(79,4.5)
69	(78, 5.3)/(85, 7.3),	(98,5.7)/(86,5.9)
21	(52, 5.7)/(54, 6.4),	(93, 8.6)/(80, 9.0)
25	(62, 1.2)/(62, 1.2),	(89,5.6)/(87,3.8)
25	(71, 7.0)/(72, 7.6),	(88,13.2)/(88,15.4)
31	(78, 3.5)/(79, 3.6),	(95,4.3)/(89,3.6)
29	(77, 1.3)/(78, 1.4),	(92,5.3)/(89,3.6)
5	(12,5.6)/(12,10.5),	(81,20.8)/(38,27.8)
68	(89,6.5)/(91,7.3),	(95,6.5)/(91,7.0)
69	(92,5.3)/(95,6.2),	(95,6.0)/(93,6.9)
35	(81,5.5)/(83,8.3),	(98,6.0)/(94,7.1)
52	(87,4.7)/(89,5.1),	(96, 5.9)/(97, 5.6)
50	(79,5.4)/(81,6.0),	(89,6.2)/(81,6.6)
127	(95, 1.5)/(94, 1.5),	(91,3.1)/(92,2.4)
25	(73, 8.9)/(74, 11.8),	(80,14.3)/(76,16.3)
91	(82, 2.5)/(83, 2.6),	(93,3.0)/(89,2.7)
22	(65, 3.8)/(68, 4.5),	(92,5.7)/(76,6.7)
22	(71, 1.8)/(72, 1.9),	(82,5.2)/(82,3.4)
26	(55, 5.2)/(56, 6.3),	(87,7.8)/(62,8.1)
46	(81, 6.3)/(80, 7.0),	(93,21.8)/(94,24.9)
21	(64, 5.9)/(67, 6.5),	(74, 9.1)/(71, 10.7)
69	(87, 5.9)/(90, 7.2),	(88,7.4)/(89,8.6)
34	(82,7.5)/(82,9.0),	(40,9.3)/(89,10.5)
30	(77,7.6)/(78,7.9),	(85,10.2)/(84,11.6)

•...

183

Table 3 Relative efficiencies of t_d and t_{sd} for HR scheme. Slashes separate the values for respective choice of Q_i as $1/x_i$, $1/x_i^2$, $1/\pi_i x_i$ and $(1 - \pi_i)/\pi_i x_i$. Commas after parentheses separate the values for rival procedures.

Domain size	$RE(t_d)$	$\overline{RE(t_{sd})}$
73	(7.48/7.44/9.55/9.81),	(6.05/6.43/6.65/6.70)
69	(5.08/5.80/5.45/5.48),	(3.77/4.36/4.24/4.36)
21	(2.67/2.71/2.72/2.73),	(5.69/6.32/6.21/6.29)
25	(28.88/23.17/23.01/23.02),	(6.15/8.94/8.23/8.90)
25	(4.40/4.48/4.53/4.55),	(3.43/2.78/2.90/2.80)
31	(13.24/13.02/13.40/13.38),	(10.59/7.91/8.54/8.00
29	(25.42/25.97/25.77/25.81),	(6.06/8.83/8.10/8.76)
68	(4.47/4.03/4.24/4.09),	(5.44/4.94/5.11/4.98)
69	(3.47/3.28/3.22/3.11),	(3.57/3.30/3.35/3.28)
35	(6.69/6.87/7.03/7.06),	(5.87/5.16/5.30/5.16)
52	(10.44/10.32/10.15/9.98),	(7.34/9.20/8.79/9.03)
50	(5.27/5.41/5.41/5.36),	(6.42/6.55/6.61/6.55)
127	(21.10/21.23/22.33/22.43),	(5.33/7.27/6.84/7.30)
25	(2.54/2.33/2.65/2.66),	(2.31/2.14/2.19/2.16)
91	(9.84/10.56/11.66/11.73),	(7.03/9.32/9.08/9.48)
22	(7.63/7.86/7.85/7.87),	(8.51/8.48/8.62/8.43)
26	(5.63/6.23/6.35/6.45),	(6.54/7.93/7.73/7.96)
46	(4.68/4.43/4.56/4.52),	(1.79/1.55/1.60/1.56)
21	(2.57/2.59/2.59/2.60),	(7.25/6.38/6.58/6.34)
69	(3.69/4.17/4.27/4.49),	(4.56/4.64/6.69/6.69)
34	(3.58/3.65/3.77/3.81),	(4.02/3.86/3.94/3.90)
30	(3.29/3.22/3.31/3.31),	(3.96/3.41/3.52/3.42)

Relative efficiencies of t_d and t_{sd} for Lahiri's scheme. Slashes separate the values for respective
choice of Q_i as $1/x_i$ and $1/x_i^2$. Commas after parentheses separate the values for rival
procedures.

Table 4

Domain size	$RE(t_d)$	$RE(t_{sd})$	Domain Size	$RE(t_d)$	RE(t _{sd})
73	(5.65/5.35),	(4.97/4.80)	25	(2.96/2.78),	(3.01/2.75)
69	(4.05/3.20),	(4.59/4.31)	91	(7.05/7.02),	(7.11/7.55)
21	(3.61/3.56),	(.577/5.32)	22	(2.48/2.17),	(7.22/6.13)
25	(4.13/4.14),	(7.41/10.51)	22	(2.79/2.79),	(8.50/13.12)
25	(4.33/4.17),	(3.67/3.17)	26	(4.61/4.46),	(5.08/4.79)
31	(3.30/3.30),	(8.07/6.78)	46	(3.77/3.40),	(1.48/1.33)
29	(17.22/17.02),	(7,67/11.07)	21	(4.41/4.28),	(4.65/3.94)
68	(3.38/3.15),	(3.72/3.50)	69	(3.66/3.41),	(3.53/3.16)
69	(4.16/3.86),	(4.58/4.02)	34	(3.84/3.65),	(4.01/3.67)
35	(4.93/2.80),	(5.97/5.29)	30	(3.89/3.86),	(3.77/3.39)
52	(4.85/4.82),	(4.17/4.78)			
50	(5.21/4.95),	(5.46/4.89)			
127	(12.06/12.07),	, (6:17/7.96)			

Honouring the valuable suggestions from one of the referees we present a few summary measures of the performances of the above procedures in the tables below; they are:

- (i) Medians, first and third quartiles and the minimum and maximum values, respectively abbreviated as Med, Q_{1/4}, Q_{3/4}, Min and Max, of ACP, ACV and RE(.);
- (ii) Numbers of domains out of the 29 for which the values of ACP for the procedures are 90 or more;
- (iii) Number of domains for which ACP for t_{sd} is closer to 95 than that for t_d ;
- (iv) Numbers of instances in much the use of v_2, v_{s2} gives an ACP closer to 95 than that given by the use of v_1, v_{s1} ;
- v) Numbers of instances in which values of ACV are not more than the minimum. ACV plus 5;
- (vi) Numbers of instances in which t_{sd} gives a smaller ACV than t_d ;
- (vii) Numbers of instances in which the use of v_2, v_{s2} gives a smaller ACV than the use of v_1, v_{s1} ;
- (viii) Numbers of instances in which $RE(e_d)$ is greater than or equal to 5;
- (ix) Numbers of instances in which t_{sd} gives a larger RE than t_d ;

We present these values only for the HR scheme; as those for the Lahiri's scheme reveal roughly a similar pattern we do not show them here.

Max	Min	$Q_3/4$	Q1/4	Med	Criteria		Max	Min	$Q_{3/4}$	Q1/4	Med	Criteria	given su	Performar concernin
93.0,	22.1.	78.4,	33.7,	53.9,	(t_{Hd}, v_{YGd})	Lepeat	99,	33,	92,	72,	86.	(t_{Hd}, v_{YGd})	ccessively sepa	nces of procedu g ACP based o
(637.6/122.8/ 32.4/37.3)	(1.1/1.1/1.6/1.6)	(9.3/8.4/8.3/7.9)	(4.0/3.4/3.5/3.4)	(5.9/5.3/5.1/4.6)	(t_d, v_1)	Table 6 Lepeat of contents of Table 5 with ACP replaced by ACV.	(94/94/94/94),	(3/3/3/3),	(81/80/82/80),	(26/26/28/26),	(71/65/65/61).	(t_d, v_1)	given successively separated by slashes. Commas separate values for respective procedures	Performances of procedures in terms of summary measures "Med, $Q_1/4$, $Q_3/4$, Min and Max" concerning ACP based on HR scheme. Values for Q_1 as $\frac{1}{2}$, $\frac{1}{2}$, $\frac{1}{2}$, $\frac{1}{2}$, within parentheses
		-	, (5.9/4.9/4.9/4.9)	, (8.9/7.6/8.6/7.7)	(t_{sd}, v_{s1})	Table 6 sle 5 with ACP rej	(97/97/97/97),	(97/97/97/97),	(83/92/92/92),	(71/65/67/66),	(87/86/87/26).	(t_{3d}, v_{31})	Jommas separate v	Table 3 nmary measures " ues for Q_i as $\frac{1}{r_i}$,
_		15.3), (-	-		slaced by ACV.	(95/95/95/97),	(95/95/95/97),	(80/78/78/78),	(26/28/26/26),	(67/68/66/68).	(t_d, v_2)	alues for respect	Med, $Q_1/4, Q_3/4$
(629 4/93 n/	0.5/0.3/0.6/0.6)	(7.3/6.2/5.9/5.8),	(4.1/2.3/2.1/2.1),	(5.2/4.6/4.5/4.3),	(t_d, v_2)		(98/97/100/97	(98/97/100/97	(93/93/94/93)	(77/65/71/68)	(91/87/89/86	(tsd, Us2	tive procedures	, Min and Max thin parenthes
(201,1/84,1/	(37/28/21/20)	(11.2/10.6/11.4/11.4)	(5.9/4.8/5.0/4.9)	(9.1/6.7/7.2/7.2)	(t_{sd}, v_{s2})		/97)	/97)	93)	68)	86)			es X

185

¢

Table 7	
Repeat of contents of Table 5 with ACP replaced by RE (.) That v_d is in	relevant
here in assessing e_d may be noted. Also, only t_d and t_{sd} are relevant for	$RE(e_d)$
and t_{Hd} constitutes only the base.	

Criteria	$RE(t_d)$	$RE(t_{sd})$
Med	(3.15/3.20/3.22/3.11),	(6.04/6.43/6.58/6.44)
$Q_{1/4}$	(1.82/1.82/1.82/1.82),	(4.56/4.64/4.69/4.69)
$Q_{3/4}$	(5.63/6.87/6.35/6.45),	(7.05/8.48/8.23/8.43)
Min	(1.23/1.24/1.24/1.24),	(1.79/1.55/1.60/1.56)
Max	(25.42/25.97/25.77/25.81),	(10.59/13.67/12.22/13.69)

Table 8. For HR scheme

The numbers of instances with ACP as 90 or more for procedures. Values for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}, \frac{1-\pi_i}{\pi_i x_i}$ separated by slashes

$(t_{HD}, v_{YGD}),$	$(t_d, v_1),$	$(t_{sd}, v_{s1}),$	$(t_d, v_2),$	(t_{sd}, v_{s2})
16	(3/3/3/3),	(15/12/17/12),	(3/3/2/2),	(22/16/18/16)

For HR scheme, the number of instances in which ACP for t_{sd} is closer to 95 than that for t_D is 38.

Table 9. For HR scheme

The numbers of cases use of v_2, v_{s2} gives ACP closer to 95 than that of v_1, v_{s1} . Values are respectively given for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}, \frac{1-\pi_i}{\pi_i x_i}$ separated by slashes within parentheses for procedures separated by commas

$(t_d, v_2)vs(t_d, v_1),$	$(t_{sd}, v_{S2})vs(t_{sd}, v_{s1})$
(19/20/20/20),	(25/23/24/24)

Table 10. For HR scheme

Numbers of cases for which ACV does not exceed minimal ACV plus five. Values for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}, \frac{1-\pi_i}{\pi_i x_i}$ separated by slashes for procedures separated by commas.

(t_{Hd}, v_{YGd})	(t_d, v_1)	$(t_{sd}, v_{s1}),$	(t_d, v_2)	(t_{sd}, v_{s2})
0,	(17/18/19/21),	(6/10/10/10),	(18/21/23/24),	(5/12/10/10)

Table 11. (For HR scheme)

Number of cases with t_{sd} giving smaller ACV than t_d . Values for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}$ and $\frac{1-\pi_i}{\pi_i x_i}$ given respectively separated by slashes for procedures separated by commas

$(t_{sd},v_{s1})vs(t_d,v_1),$	$(t_{sd}, v_{s2})vs(t_d, v_2)$
(4/4/2/0),	(8/8/4/4)

Table 12. (For HR scheme)

Number of cases by use of v_2, v_{s2} yielding lesser ACV than that of v_1, v_{s1} . Values for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}$ and $\frac{1-\pi_i}{\pi_i x_i}$ respectively separated by slashes for procedures separated by commas

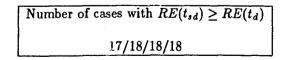
$(t_d, v_2)vs(t_d, v_1),$	$(t_{sd}, v_{s2})vs(t_{sd}, v_{s1})$	
(25/26/27/25),	(25/23/24/24)	

Table 13. (For HR scheme) Number of cases with $RE(e_d)$ greater than or equal to 5. Values separated by slashes for Q_i respectively as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}$ and $\frac{(1-\pi_i)}{\pi_i x_i}$ for procedures separated by commas

$RE(t_d)$	$RE(t_{sd})$
(12/12/12/12),	(28/27/29/27)

Table 14. (For HR scheme)

Number of instances for which $RE(t_{sd})$ exceeds $RE(t_d)$. Values separated by slashes for Q_i as $\frac{1}{x_i}, \frac{1}{x_i^2}, \frac{1}{\pi_i x_i}, \frac{(1-\pi_i)}{\pi_i x_i}$.



5. CONCLUDING REMARKS AND RECOMMENDATIONS

- (i) A procedure that fails to achieve a value of ACP at least 80 may not be acceptable. In the present example, very few cases fail by this criterion.
- (ii) For domains of size 15 or more, HTE for both the HR and Lahiri schemes, is adequate to achieve a desired ACP. But the ACV and hence the length of the confidence interval based on HTE is unacceptably poor. Also it is inefficient compared to t_d and t_{sd} .
- (iii) Although the non-synthetic estimator t_d achieves the best ACV, in most cases it ensures a poor level of ACP when coupled with either form of its variance estimator. It often turns out poorer than the HTE in terms of ACP although it is more efficient. Taking everything into consideration it need not be an improvement upon the HTE to a desired extent.
- (iv) The synthetic estimator t_{sd} is decidedly an improvement upon the HTE for both HR and Lahiri schemes in all the three respects, namely, ACP, ACV and RE. It is preferable to t_d except in terms of ACV. It combines better with the variance estimator that uses the g-weights.
- (v) For domains of sizes 15 or more, in the present example, the 'synthetic' greg estimator coupled with the g-weighted variance estimator turns out to be most appropriate for both HR and Lahiri schemes with Q_i chosen as $1/x_i$. The choice of Q_i as $1/x_i^2$ for both the schemes turns out poor in many situations and so should be avoided.
- (vi) There is not much to distinguish between these two schemes in using t_{sd} and among choices of Q_i as $\frac{1}{x_i}, \frac{1}{\pi_i x_i}, \frac{(1-\pi_i)}{\pi_i x_i}$ for both schemes.

ACKNOWLEDGEMENT

We are grateful to two referees and the editor whose suggestions helped us to substantially improve upon two earlier drafts. We are indebted to the authorities of Indian Statistical Institute, Calcutta who released official records for our use. Sri Milan Kumar Santra, Sri Arup Kumar Seal and Sri K.V.S. Ravi Kumar did the computational work to earn our thanks.

REFERENCES

- (1) Brewer, K.R.W. (1979). A class of robust sampling designs for large-scale surveys. Jour. Amer. Stat. Assoc. 74, 911-915.
- (2) Hajek, J. (1971). Comment on a paper by Basu, D. In Foundations of Statistical Inference. Ed. Godambe, V.P. and Sprott, D.A. Holt, Rinehart, Winston; Toronto, 203-242.

- (3) Hartley, H.O. and Rao, J.N.K. (1962). Sampling with unequal probabilities and without replacement. Ann. Math. Stat. 33, 350-374.
- Horvitz, D.G. and Thompson, D.J. (1952). A generalization of sampling without replacement from a finite universe. Jour. Amer. Stat. Assoc. 47, 663-685.
- (5) Lahiri, D.B. (1951). A method of sample selection providing unbiased ratio estimates. Bull. Int. Stat. Inst. 33, 133-140.
- (6) Särndal, C.E. (1980). On π-inverse weighting versus best linear weighting in probability sampling. *Biometrika*. 67, 639-650.
- (7) Särndal, C.E. (1982). Implications of survey design for generalized regression estimation of linear functions. Jour. Stat. Plan. and. Inf. 7, 155-170.
- (8) Särndal, C.E., Swensson, B.E. and Wretman, J.H. (1992). Model assisted survey sampling. Springer-Verlag, New York.
- (9) Yates, F. and Grundy, P.M. (1953). Selection without replacement from within strata with probability proportional to size. Jour. Roy. Stat. Soc. Ser. B, 15, 253-261.