# SAMPLING FOR THE COUNT IN A CENSUS

Cary T. Isaki, Julie H. Tsay, and Yves Thibaudeau, Bureau of the Census
Cary T. Isaki, Room 3132, Bldg. 4, Washington, D.C.   20233

KEY WORDS:   Small areas, bias, mean
square error

## I. Introduction

As part of a program of continuing research regarding the 2000 Decennial Census we conducted an empirical study concerning the possibility of sampling for the count. Sampling for the count represents a marked departure from traditional census methods in which every known unit is contacted for enumeration. In our study of sampling for the count we focussed on two sample designs - one that utilized the block as the sampling unit and another that used the housing unit. In the block design all housing units in the selected block are in sample. In the housing unit design, samples of housing units were selected from within each block.

The main purpose of this research was to construct several sample designs and provide empirical results concerning estimates of Voting Rights Act data at the block, address register area (ARA) and district office (DO) level. Voting Rights Act data are counts of persons 18 and older by Black, non Black Hispanic and the remainder. To give a rough perspective, blocks contain 50 housing units or so; ARAs contain roughly 30 blocks and DOs average 150,000 housing units. We used the 1990 Decennial Census data for our study. Biases and variances of estimators at the various geographic levels and for the Voting Rights Act data were computed for four DOs. Time limitations for research precluded processing of more DOs.

## II. Assumptions and Design

The sample design used in the study assumed that a frame of address consistent with that available in the 1990 Census would be used as a frame. For operational considerations we chose to sample within DOs. Furthermore, we considered only DOs that were covered entirely by mail in 1990. About five percent of the housing units in 1990 were covered by traditional census methods, termed list enumerate, while another comparable group was covered by an update/leave method. Both of these latter procedures utilized field enumerators.

We assumed that our frame of addresses were coded to the block level. Upon selection of the addresses, via a block sample or individual address

sample, we further assumed that all sampled units responded. We also assumed that the responses in the data file to be used in the study were that to be obtained under a sampling for the count scenario. On the data file, there existed a set of housing units coded to block within the DO that were not coded as being mailed a form yet were linked to data. We treated them as mail address cases for sampling purposes. In the following, we refer to the individual addresses as cases which in the survey process are converted to housing units (occupied or vacant) or deletes/kills (not housing units). This also removes some confusion among persons who view addresses as linked to structures and not units within, such as apartments in buildings. Finally, the group quarters universe (GQ) such as prisons are identified separately and were canvassed by a special operation in 1990. We excluded GQs from sampling in our work.

Several factors motivated consideration of a block sample design. If a Post Enumeration Survey (PES) were to be used for adjustment, then a subsample of blocks would support a PES as used in 1990. A block sample would also enhance the coverage of a case sample. Evidence of this surfaced when it was discovered (after completion of research) that some housing units on our data files that were not mailed forms were actually discovered via a field canvassing of units during the 1990 Census. These units are called census adds. Such a process could still be accomplished via a block sample design with field followup of nonrespondents, but not via a case sample design. Consequently, the case sample design is actually subject to coverage bias of census adds. The effect of the bias is trivial at the block level, but not so at the DO level.

The case sample design would provide unbiased estimators at the block level for cases on the list frame. It also provides samples in every block which some feel is a political advantage.

## III. Methodology

The survey designs under study user either the block or the case (mail address) as the sampling unit. It is a condition of the sampling for the count study that only mail addresses as available in the 1990 census could be assumed available for use in the construction of the survey designs that follow.

## A. Block Design

The block design consists of a stratification of blocks and a simple random sampling of blocks from within strata by DO. Certainly blocks consisting of large number of cases as well as those with ten or less cases are used. Within selected blocks all cases are mailed census forms except for a small number of cases that are canvassed by other means.

The estimator of block total for characteristic Y under the block design is $y_{ih}$ where

$$y_{ih} = \begin{cases} Y_{ih}, & \text{if the } i^{\text{th}} \text{ block in stratum h is in sample} \\ \hat{R}_h^y M_{ih}, & \text{otherwise} \end{cases}$$

where $Y_{ih}$ is the total for characteristic Y in the $i^{\text{th}}$ block in stratum h, $M_{ih}$ is the total number of cases in the $i^{\text{th}}$ block in stratum h (known for all blocks) and $\hat{R}_h^y = \sum_i^{n_h} Y_{ih} / \sum_i^{n_h} M_{ih}$ is the ratio of the sums of $n_h$ sample observations in stratum h.

The block estimator, $y_{ih}$, was motivated by the desire to use the sample observation for the block when it was selected. When it was not in sample, a ratio adjustment to its case count was used. The block estimator is biased. In fact it can be shown that any block estimator that uses its sample observation exclusively will be biased. Apart from ratio bias, the sum of $y_{ih}$ over blocks in the $h^{\text{th}}$ stratum is unbiased for the total in stratum h. To reduce the effect of ratio bias, stratum sample sizes were allocated to provide a minimum of approximately ten blocks per stratum.

It can be shown that the expected value of $y_{ih}$ is

$$E[Y_{ih}] = Y_{ih}(n_h/N_h) + R_{(i)h}^y M_{ih}(1 - n_h/N_h) \quad (1)$$

where $\hat{R}_{(i)h}^y = \overline{Y}_{(i)h}/\overline{M}_{(i)h}$ and $\overline{Y}_{(i)h}$ is the mean of $N_h - 1$ blocks in stratum h excluding block i.

From (1) we see that the bias of $y_{ih}$ is dictated by how close the term, $R_{(i)h}^y M_{ih}$, is to $y_{ih}$. The variance of $y_{ih}$, omitting the subscript h for brevity, is

$$\text{Var}[y_{ih}] = (n/N)(1 - n/N)[Y_i - R_{(i)}^y M_i]^2$$

$$+ n^{-1} M_i^2 \overline{M}^{-2}(1 - n/(N-1))(1 - n/N)$$

$$\sum_{\substack{j=1 \\ \neq i}}^{N-1} (Y_j - R_{(i)}^y M_j)^2 / (N-2) \quad (2)$$

Under the design it is clear that the covariance between $y_{ih}$ and $y_{kh'}$ is zero for $h \neq h'$.

Otherwise, an expression for the covariance between the two can be found in Isaki, et. al. (1993).

## B. Case Sample Design

The case design consists of simple random without replacement sampling of cases independently from within every block except for the blocks with less than or equal to ten cases where all cases are selected. Treatment of GQ persons, are as handled in the block design. A constant sampling rate is applied in all other blocks. The case design allows for unbiased estimators of total characteristic for every block. Estimators of total for ARAs and the DO are formed by summing over block estimators.

Under the case design we considered two estimators of block total. One estimator is the usual Horvitz-Thompson estimator which we call the Case 2 estimator and the other is a Royall type estimator which we term the Case 1 estimator. The Case 1 estimator, $y_{ih}$, for block i in stratum h is

$$y_{ih} = \sum_{j=1}^{n_{hi}} Y_{ijh} + \hat{R}_h(M_{ih} - n_{hi}) \quad \text{where} \quad (3)$$

$Y_{ijh}$ is the Y characteristic for case j in block i
$n_{hi}$ is the sample number of cases in block i

$$\hat{R}_h = \sum_i^{N_h} \sum_j^{n_{hi}} Y_{ijh} / \sum_i^{N_h} n_{hi} \quad ,$$

$M_{ih}$ is the number of cases in block i and
$N_h$ is the number of blocks in stratum h.

The subscript h denotes the $h^{\text{th}}$ stratum as defined in the block design. The blocks in the $h^{\text{th}}$ stratum are used in the second term in (3) to account for those cases not in the sample of $n_{hi}$ cases. Under the case sample design, it can be shown that the Case 1 estimator for block i in stratum h is biased with expectation

$$E[y_{ih}] = n_{hi} \overline{Y}_{ih} + R_h(M_{ih} - n_{hi}) \quad (4)$$

$$\text{where} \quad R_h = \sum_{i=1}^{N_h} n_{hi} \overline{Y}_{ih}/n_h \quad ,$$

$$n_h = \sum_i^{N_h} n_{hi}$$

Furthermore, it can be shown that the variance of $y_{ih}$

is

$$\text{Var}\,[y_{ih}] = M_{ih}\left[\frac{n_{hi}}{M_{ih}}\right]\left(1 - \frac{n_{hi}}{M_{ih}}\right)S_{ih}^2 +$$

$$(M_{ih} - n_{hi})^2\,V\,(\hat{R}_h)$$

$$+\,2\,n_{hi}^2(M_{ih} - n_{hi})\,n_h^{-1}\,V\left[\sum_j^{n_{hi}} Y_{ijh}/n_{hi}\right] \quad (5)$$

where $\quad V(\hat{R}_h) = \sum_i^{N_h} (n_{hi}/n_h)^2\,n_{hi}^{-1}\,(1 - n_{hi}/M_{ih})\,S_{ih}^2,$

$$n_h = \sum_i n_{hi} \quad \text{and}$$

$$S_{ih}^2 = \sum_j^{M_{ih}} (Y_{ijh} - \overline{Y}_{ih})^2/(M_{ih} - 1)$$

Similarly, the covariance between $y_{ih}$ and $y_{kh}$ is zero if i and k are in separate strata. When they are in the same strata, their covariance is

$$\text{Cov}\,(y_{ih}, y_{kh}) = n_{hi}\,(M_{kh} - n_{hk})\,\text{Cov}\left[\sum_j^{n_{hi}} Y_{ijh}/n_{hi}, \hat{R}_h\right]$$

$$+\,(M_{ih} - n_{hi})\,n_{hk}\,\text{Cov}\left[\sum_j^{n_{hk}} Y_{kjh}/n_{hk}, \hat{R}_h\right]$$

$$+\,(M_{ih} - n_{hi})\,(M_{kh} - n_{nk})\,\text{Var}\,(\hat{R}_h) \quad (6)$$

With the above expressions it is possible to derive the variances of block estimators, ARA estimators and of the DO estimator. We note that the Case 1 estimator is unbiased for stratum totals and hence, is unbiased for the DO. Because it is identically the sum of unbiased estimators for each stratum, its DO variance is the sum of the variances of the stratum totals. For each large certainty block under the Case 1 estimation scheme, a simple unbiased estimator was used. Inclusion of the Case 1 estimator in the study was motivated by the thought that it would provide a smaller variance than the Case 2 estimator. The effect of its bias was an open question.

## IV. Study Description

The study was based on a 1990 Census scenario in terms of the handling of the universe of mail addresses. The study utilized 1990 Census data as well with a modification for constructing strata and determining stratum sample sizes.

## A. Construction of Strata

Strata of blocks were constructed to support the block design and the same strata were also used in defining the Royall type estimator under the case sample design. The stratification process used in the study was based on a clustering algorithm that used four variables. The four variables used were the block's case count and the ARA's ratio of the number of Black 18[+], non Black Hispanic 18[+] and other 18[+] to the total case count which were applied to each member block's case count. It was felt that ARA race distributions were somewhat stable over time and hence, use of the 1990 distributions would provide a similar stratification result in the next census. The clustering algorithm identified certain blocks as outliers. Such blocks were treated as certainty and included in the group of blocks termed large certainty blocks.

## B. Allocation of Sample to Strata / Blocks

Given the strata and an overall sampling rate for mail universe cases we determined sample number of blocks per stratum in the block design and sample number of cases per block in the case design. For sample size allocation to strata of blocks we sought to minimize the variance of the stratified Horvitz-Thompson estimator of total 18[+]. The minimization was performed subject to a constraint on total number of cases, discounted for cases in certainty blocks and limited by the sampling rate specified. No limitation was placed on GQ cases. Hence, a specified sampling rate was maintained over the mail universe. Allocation of sample size to stratum was forced to be not less than ten blocks or so. This condition was imposed to control bias and variance of the block estimator.

For the case design, we assumed a proportional allocation of sample cases to each block. Hence, for all relevant blocks a fixed rate slightly reduced because cases in certainty blocks were canvassed completely, was applied in each block. This produced the obvious anomaly of some blocks being assigned smaller sample sizes than blocks of smaller size. We ignored this in the study.

## C. Summary Statistics

We created block and ARA files of expected values, variances and other statistics listed below by DO. The following summary statistics were computed.

Let

$Y_i$ be the census total characteristic for the $i^{th}$ area

$y_i$ be the estimator of total characteristic for the $i^{th}$ area

E $[y_i]$ be the expected value of $y_i$

Var $[y_i]$ be the variance of $y_i$

Bias $(y_i) = E\ [y_i] - Y_i$

MSE $(y_i) = Var[y_i] + [Bias\ (y_i)]^2$ be the mean square error of $y_i$

Relative Bias $(y_i) = Bias\ (y_i)\ /\ Y_i$

s.e. $(y_i) = [Var\ (y_i)]^{1/2}$

RMSE $(y_i) = [MSE\ (y_i)]^{1/2}$ be the root mean square error of $y_i$

The summary statistics by DO and by area (block, ARA) were -

i)   mean [s.e. $(y_i)$ / E $[y_i]$]

ii)  mean | Relative Bias $(y_i)$ |

iii) mean [RMSE $(y_i)$ / $Y_i$]

vi)  Sample size count of blocks and cases for - GQ, mail out, mail universe but not mailed out

## V.  Results

Table 1 in the Appendix summarizes the bias and mean squared error of the estimators at the block, ARA and DO level for total Black 18[+] persons in DO 2107 (South Boston). DO 2107 contained 2,636 blocks and 87 ARAs under study. There were 145,393 cases of which 9,975 were not mailed. The total number of persons in the DO was 339,389 of which about 39 percent were Black and seven percent were Hispanic.

Table 2 refers to Hispanics 18[+] in DO 2801 (North, D.C.). DO 2801 consisted of 2,210 blocks and 91 ARAs under study. There were 142,848 cases of which 6,595 were not mailed. The total number of persons in the DO was 309,391 of which about 87 percent were Black and three percent were Hispanic.

These DOs were selected because of the apparent anomalies in some statistics. The summary statistics contained within were typical of the results obtained for other DOs and race by age.

## VI.  Discussion of Some Tabular Results

The following presentation discusses some unusual observations in the tabular results and provides an explanation.

## A.  Increasing Variance of the Case 1 Estimator, as the Sampling Rate Increases

Consider the behavior of the average relative standard error of the Case 1 estimator, at the block level. This estimator yields an estimate of the number of individuals with a given characteristic (Black less than 18, black 18[+]...). Typically, the relative variance increases when the sampling rate of the process increases, for small values of the sampling rate.

This behavior is attributable to the nature of the variance of the Case 1 estimator, in terms of the sampling rate. The variance is the sum of three distinct parts corresponding to the variances and the covariance between the two terms involved in the Case 1 estimator. (See equation (5)). The first term is a second degree polynomial (in n/N) strictly increasing between 0 and .5. In the situations studied as part of this research it was found that the first component increases for <u>small values of the sampling rate</u> to such an extent that the variance increases. As the variance of the Case 1 estimator increases with increasing sampling rate, the bias decreases. The resulting relative M.S.E. is usually decreasing, as the sampling rate increases. The table below illustrates this situation with an example. For the count of Black 18[+] in D.O. 2107, it gives the average of the variances and of their three components over all the blocks with a true count greater than zero.

| Rate | 1st component | 2nd component | 3rd component | Variance | M.S.E. |
|---|---|---|---|---|---|
| 10% | 5.3722 | 2.3105 | 0.4373 | 8.1201 | 1646.27 |
| 20% | 9.7202 | 0.8695 | 0.4063 | 10.9961 | 588.294 |
| 30% | 12.4282 | 0.3140 | 0.2578 | 13.0001 | 375.936 |
| 40% | 13.7525 | 0.2885 | 0.3774 | 14.4184 | 199.916 |
| 50% | 13.9300 | 0.1181 | 0.2349 | 14.2832 | 138.459 |
| 60% | 12.6712 | 0.04638 | 0.1393 | 12.8569 | 78.5502 |
| 70% | 10.4599 | 0.01968 | 0.0919 | 10.5715 | 42.2967 |

The clear pattern in this table is responsible for inducing a similar pattern for the average of the relative errors of the Appendix.

## B.  Increasing Relative Bias from Block to D.O. Level (D.O. 2801)

Another interesting fact concerns the average relative bias at the ARA level for the Case 1 estimator, compared to the average relative bias at the block level for the same estimator. For instance, the relative bias of the Case 1 estimator for Hispanics 18[+], for DO 2801, is 1.1481 and at the ARA level, it is 4.3725. An ARA covers more cases than the typical block. Intuitively, the average relative bias should decrease. For the computation at the block level, the blocks with no cases of Hispanics 18[+] are ignored. Note however, some of these blocks without Hispanic 18[+] can have a large bias. On the other hand, for the computation at the ARA level, few ARA have no Hispanics 18[+] and all the blocks that are part of an ARA containing at

least one Hispanic $18^+$ are included. In particular some blocks with large bias, excluded from the block level calculation, are now included in the ARA analysis.

For example, a particular ARA (4052) in D.O. 2801, has only one Hispanic aged $18^+$. This individual is in one block, all the other blocks in the ARA were not part of the analysis at the block level, and their bias did not enter in the computation of the average relative bias for Hispanics $18^+$, at the block level. However, some of these "zero" blocks possessed a large absolute bias. One block (120), for instance, did not have any Hispanic $18^+$, but had an absolute bias of 47. In addition, several blocks with no Hispanic $18^+$ also exhibited moderate-to-large absolute biases. The total absolute bias for the ARA is about 78. Since there was only one Hispanic $18^+$ in this ARA, the relative bias of the ARA was 78. The relative bias of this ARA contributed heavily to the average relative bias at the ARA level.

Why block 120 in ARA (4052) possessed a large bias was found in the stratification. The rate of imputation embedded in the Case 1 estimator is a function of the stratification. At the 10% rate, the blocks were classified into four strata. Most of the blocks (77%) belong to strata 2 and 3. The associated imputation rates were .0916 and .0263. All the blocks in ARA 4052 were in stratum 3, except for block 120. Block 120 was in stratum 4 with a corresponding rate of imputation of .1349. In addition, there were 388 households in block 120. These circumstances together gave an estimated number of Hispanics $18^+$ of 47, while in fact, there were none.

## C. Conflicting Trends Between Absolute and Relative Measurements as the Sampling Rate Increases

The analysis of the block design also presented some unusual observations. In particular, in some cases, when the sampling rate increased, the total standard error decreased, but the total coefficient of variation increased. This phenomenon is recreated on a small scale in the two tables below, for a set of four blocks taken in D.O. 2107. These figures relate to the counts of Black $18^+$. The total for the standard error decreases when increasing the sampling rate from 50% to 60%, however, at the same time, the total for the coefficient of variation increased. The same opposite trends can be observed for the root mean square error vs. the relative root mean square error. In the following tables block 1, 2, 3, 4, refer to specific blocks in D.O. 2107 (South Boston). Block 1 is in fact block 303 in ARA 4037, block 2

is block 203 in ARA 4038, block 3 is block 305 in ARA 4038, and block 4 is block 103 in ARA 4040.

1. Sampling Rate 50%

| Block # | Y | E(Y) | S.E. | C.V. | R.M.S.E. | R.M.S.E. Y |
|---|---|---|---|---|---|---|
| 1 | 27 | 40.8 | 10.357 | .254 | 17.254 | .639 |
| 2 | 25 | 27.8 | 5.104 | .184 | 5.814 | .233 |
| 3 | 3 | 5.0 | 1.384 | .275 | 2.456 | .819 |
| 4 | 121 | 131.7 | 9.473 | .072 | 14.282 | .118 |
| Total | | | 26.318 | .785 | 39.806 | 1.809 |

2. Sampling Rate 60%

| Block # | Y | E(Y) | S.E. | C.V. | R.M.S.E. | R.M.S.E. Y |
|---|---|---|---|---|---|---|
| 1 | 27 | 19.0 | 5.599 | .294 | 9.728 | .360 |
| 2 | 25 | 27.3 | 3.518 | .129 | 4.223 | .169 |
| 3 | 3 | 11.0 | 5.409 | .490 | 9.687 | 3.229 |
| 4 | 121 | 131.7 | 8.934 | .068 | 14.002 | .116 |
| Total | | | 23.460 | .981 | 37.640 | 3.874 |

## VII. References

1. Dalzell, D. (1990), "1990 Census 21st Decennial Census 100% Detail File DCS*HEDF - < DO > draft", 13 pages.

2. Isaki, C.T., Tsay, J.H. and Thibaudeau, Y. (1993), "Empirical Results of Two Sample Design Options for Sampling for the Count", 23 pages plus two appendices and sets of tables.

3. Katzoff, E. and McLaughlin, G. (1992), "1990 Decennial Census - Census Data Organization Project Operation File [CDOPOP]", 22 pages.

496

4. Miskura, S. (1992), Memorandum of 7/6/92 to Distribution List, Subject: 2000 Census Research and Development Alternative Designs Program with Attachments.

5. Miskura, S., Woltman, H. and Thompson, J. (1984), "Research Plan Uses of Sampling for the Census Count," Proceedings of the Social Statistics Section, pg. 458-463, Meetings of the American Statistical Association.

6. Royal, R.M. (1970), "On Finite Population Sampling Theory Under Certain Linear Regression Models", Biometrika, 57, pg. 377-387.

Appendix. 2107-1

Table 1. Mean Summary Statistics by Design for Black 18+ by Sampling Rate and Geographic Level - DO 2107 Total is 86,355[1]

| | No. of Items[2] with Black 18+ > 0 | s.e. (y_i)/E[y_i] | | | Mean of Summary Statistics \|Rel. bias (y_i)\| | | | RMSE(y_i)/Y_i | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Block | Case 1 | Case 2 | Block | Case 1 | Case 2 | Block | Case 1 | Case 2 |
| Sampling Rate - 10% | | | | | | | | | | |
| Block | 1737 | .2387 | .1083 | 1.0531 | 1.2934 | 1.1844 | 0 | 1.4210 | 1.2177 | 1.0531 |
| ARA | 87 | .1840 | .0379 | .2445 | 1.2783 | 1.2916 | 0 | 1.3761 | 1.3003 | .2445 |
| DO | 1 | .0772 | .0123 | .0123 | 0 | 0 | 0 | .0772 | .0123 | .0123 |
| Sampling Rate - 30% | | | | | | | | | | |
| Block | | .2025 | .1448 | .6786 | .9147 | .4407 | 0 | 1.0877 | .9621 | .6786 |
| ARA | | .1677 | .0428 | .1575 | .4589 | .2392 | 0 | .6211 | .4712 | .1575 |
| DO | | .0540 | .0079 | .0079 | 0 | 0 | 0 | .0540 | .0079 | .0079 |
| Sampling Rate - 50% | | | | | | | | | | |
| Block | | .1714 | .1539 | .5126 | .8274 | .3605 | 0 | 1.0064 | .8870 | .5126 |
| ARA | | .1153 | .0418 | .1190 | .3971 | .2151 | 0 | .5471 | .4080 | .1190 |
| DO | | .0393 | .0060 | .0060 | 0 | 0 | 0 | .0393 | .0060 | .0060 |

1] Block ≡ block design
   Case 1 ≡ case design with Royall type estimator
   Case 2 ≡ case design with unbiased estimator
2] Excludes 40 blocks containing only GQ persons

| | | s.e. (y_i)/E[y_i] | | | \|Rel. bias (y_i)\| | | | RMSE(y_i)/Y_i | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Sampling Rate - 40% | | | | | | | | | | |
| Block | | .1723 | .1756 | .4088 | .5413 | .4716 | 0 | .6350 | .5661 | .4088 |
| ARA | | .1128 | .0467 | .0949 | .4582 | .3105 | 0 | .5147 | .3274 | .0949 |
| DO | | .0289 | .0048 | .0048 | 0 | 0 | 0 | .0289 | .0048 | .0048 |
| Sampling Rate - 50% | | | | | | | | | | |
| Block | | .1573 | .1718 | .3327 | .4384 | .4015 | 0 | .5204 | .4999 | .3327 |
| ARA | | .0846 | .0428 | .0772 | .4060 | .2496 | 0 | .4403 | .2650 | .0772 |
| DO | | .0186 | .0039 | .0039 | 0 | 0 | 0 | .0186 | .0039 | .0039 |
| Sampling Rate - 60% | | | | | | | | | | |
| Block | | .1529 | .1596 | .2711 | .4324 | .3207 | 0 | .5447 | .4142 | .2711 |
| ARA | | .0723 | .0406 | .0629 | .2771 | .1530 | 0 | .3064 | .1694 | .0629 |
| DO | | .0125 | .0032 | .0032 | 0 | 0 | 0 | .0125 | .0032 | .0032 |

2801-1

Table 2. Mean Summary Statistics by Design for Hisp 18+ by Sampling Rate and Geographic Level - DO 2801 Total is 13,333[1]

| | No. of Items[2] with Hisp 18+ > 0 | s.e. (y_i)/E[y_i] | | | Mean of Summary Statistics \|Rel. bias (y_i)\| | | | RMSE(y_i)/Y_i | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Block | Case 1 | Case 2 | Block | Case 1 | Case 2 | Block | Case 1 | Case 2 |
| Sampling Rate - 10% | | | | | | | | | | |
| Block | 866 | .5767 | .3939 | 2.2006 | 1.2132 | 1.1481 | 0 | 1.5383 | 1.2350 | 2.2006 |
| ARA | 86 | .5112 | .1305 | 1.0396 | 3.3232 | 4.3725 | 0 | 4.1558 | 4.4482 | 1.0396 |
| DO | 1 | .2226 | .0442 | .0442 | 0 | 0 | 0 | .2226 | .0442 | .0442 |
| Sampling Rate - 33% | | | | | | | | | | |
| Block | | .3999 | .4395 | 1.0144 | .6300 | .5944 | 0 | .7460 | .7684 | 1.0144 |
| ARA | | .2031 | .1259 | .4792 | 2.0297 | 1.9486 | 0 | 2.1660 | 2.0038 | .4792 |
| DO | | .0739 | .0204 | .0204 | 0 | 0 | 0 | .0739 | .0204 | .0204 |
| Sampling Rate - 50% | | | | | | | | | | |
| Block | | .3411 | .4317 | .7082 | .3982 | .3452 | 0 | .5076 | .5843 | .7082 |
| ARA | | .1823 | .1424 | .3346 | 1.1383 | .9842 | 0 | 1.2383 | 1.0582 | .3346 |
| DO | | .0312 | .0142 | .0142 | 0 | 0 | 0 | .0312 | .0142 | .0142 |

1] Block ≡ block design
   Case 1 ≡ case design with Royall type estimator
   Case 2 ≡ case design with unbiased estimator
2] Excludes 13 blocks containing only GQ persons