

ON THE STRATIFICATION OF SKEWED POPULATIONS

Pierre Lavallée and Michel A. Hidiroglou ,
Business Survey Methods Division, Statistics Canada, Tunney's Pasture, Ottawa, Ontario, Canada, K1A 0T6

ABSTRACT

The stratification of a highly skewed population requires that it be split into a take-all stratum and a number of take-some strata. This article presents an iterative algorithm which has, as objective function, the determination of stratification boundaries, such that, for a given allocation scheme and a level of precision, the resulting sample size is minimum. The sampling from the take-some strata is assumed to be simple random and without replacement.

KEYWORDS: Iterative algorithm; Optimum boundaries; Take-all; Take-some.

1. INTRODUCTION

Highly skewed populations such as those displayed by business surveys require that they be stratified into a take-all stratum and a number of take-some strata. Units in the take-all stratum are selected with certainty whereas units in the take-some strata are selected with a given probability mechanism. Approximate cut-off rules for stratifying a population into a take-all and a take-some stratum have been given by Glasser (1962) and Hidiroglou (1986). Glasser (1962) provided the cut-off value under the assumption that a fixed sample size was to be drawn from the take-all and take-some stratum, and that the take-some sampled units were to be selected without replacement using simple random sampling. Hidiroglou (1986) provided the cut-off value under the assumption that a required level of precision had to be satisfied. These two approaches are dual in the sense that Glasser's objective was to minimize sampling variance for fixed sample size, whereas Hidiroglou's objective was to minimize sample size for fixed sampling variance.

In this article, an algorithm for stratifying a highly skewed population into a take-all stratum and a number of take-some strata will be presented. The objective will be to minimize the overall sample size given reliability constraints and to satisfy the allocation scheme of the sample to the take-some strata. The strata boundaries will be derived in term of an auxiliary variable which is closely related to the information

being collected by the survey. The algorithm is a modification to Sethi's (1963) method for stratifying a population. The resulting boundaries, which are optimal, will provide the required minimum sample size. This method will be numerically compared, in terms of boundary values and sample size, to the Dalenius - Hodges (1959) cumulative square root f rule, as well as to a mixture of the Hidiroglou (1986) and the Dalenius - Hodges (1959) stratification methods. The algorithm, which is recursive in nature, is simple to program and converges rapidly to the optimum boundary points. It also offers substantial savings in terms of sample size for given reliability criteria.

2. THE PROBLEM

Consider a finite ordered population of N units:

$$y(1), y(2), \dots, y(N),$$

with $y(i) \leq y(i+1)$ for $i=1, 2, \dots, N-1$. This population is to be stratified into L strata, with the restriction that the first L-1 strata are to be take-some and that the Lth stratum is to be take-all. The number of units to be associated with each stratum is denoted as N_h , $h=1, 2, \dots, L$. The mean to be estimated is

$$\bar{Y} = \sum_{h=1}^L \frac{M_h}{\sum_{j=M_{h-1}+1}^{M_h}} y(j) / N \quad (2.1)$$

where $M_h = \sum_{i=1}^h N_i$ for $h=1, 2, \dots, L$ and M_0 is equal to zero.

The sampling scheme calls for n_h units to be drawn from each corresponding take-some stratum of size N_h ($h=1, 2, \dots, L-1$) without replacement, using simple random sampling, and N_L units to be selected with certainty from the Lth take-all stratum. Given this set up, the estimator of population mean \bar{Y} is

$$\hat{\bar{Y}} = \sum_{h=1}^{L-1} \frac{N_h}{n_h} \sum_{j=m_{h-1}+1}^{m_h} z_j + \sum_{j=M_{L-1}+1}^N y(j) / N \quad (2.2)$$

where $y_{M_{h-1}+1} \leq z_j \leq y_{M_h}$ for $j=m_{h-1}+1, \dots, m_h$
 ($h=1, 2, \dots, L-1$), $m_h = \sum_{i=1}^h n_i$ for $h=1, 2, \dots,$
 L and m_0 is equal to zero.

Assume that the desired level of precision for the estimated mean is specified by c (coefficient of variation) and that the proportion of sampled units to be allocated to each of the first $L-1$ strata is a_h ($h=1, 2, \dots, L-1$) where $\sum_{h=1}^{L-1} a_h = 1$. The term " a_h " is conveniently used to represent any type of allocation to the strata. For instance, in the case of N -proportional power allocation,

$$a_h = \frac{N_h^p}{\sum_{h=1}^{L-1} N_h^p} \quad (h=1, 2, \dots, L-1)$$

and in the case of Y -proportional power allocation,

$$a_h = \frac{Y_h^p}{\sum_{h=1}^{L-1} Y_h^p},$$

where $0 < p < \infty$. The power allocations have the particularity that under relatively simple assumptions and for a suitable choice of p , the coefficients of variation for the take-some strata tend to be equalized without a significant increase in the overall coefficient of variation. This equality of coefficients of variation is often asked by the users of the survey data.

Denoting the population variance of each stratum h as S_h^2 , the overall sample size which satisfies the above conditions is given by

$$n = N_L + \frac{\sum_{h=1}^{L-1} N_h^2 S_h^2 / a_h}{(N c \bar{Y})^2 + \sum_{h=1}^{L-1} N_h S_h^2}. \quad (2.3)$$

The problem is to find boundaries $b_{(1)}, b_{(2)}, \dots, b_{(L-1)}$ (where $y_{(1)} \leq \dots \leq b_{(1)} < \dots < b_{(L-1)} \leq y_{(N)}$) such that the overall sample size n is minimized, given the level of reliability c and the specific allocation scheme (represented by a_h).

3. THE ALGORITHM

The approach used in this paper, for obtaining stratification boundaries for a desired level of precision, has first been used by Dalenius (1950) in the

case of stratification boundaries for a given sample size. It is first assumed that the sampling is done from a population whose frequency distribution may with "sufficient accuracy" be represented by a continuous density $f(y)$. Then, for a given set of boundaries $b_{(1)}, \dots, b_{(L-1)}$ the following quantities are defined:

$$W_h = \int_{b_{(h-1)}}^{b_{(h)}} f(y) dy \quad (3.1)$$

$$\mu_h = \int_{b_{(h-1)}}^{b_{(h)}} y f(y) dy / W_h \quad (3.2)$$

$$\sigma_h^2 = \int_{b_{(h-1)}}^{b_{(h)}} y^2 f(y) dy / W_h - \mu_h^2 \quad (3.3)$$

for $h=1, \dots, L$.

Note that $b_{(0)}$ is defined as minus-infinity ($-\infty$) while $b_{(L)}$ is plus-infinity ($+\infty$). Based on these quantities, equation (2.3) can then be rewritten as

$$n = N W_L + \frac{N \left(\sum_{h=1}^{L-1} W_h^2 \sigma_h^2 / a_h \right)}{N c^2 \mu^2 + \sum_{h=1}^{L-1} W_h \sigma_h^2} \quad (3.4)$$

where

$$\mu = \int_{b_{(0)}}^{b_{(L)}} y f(y) dy.$$

It should be noted that even if the population is considered to be "large", the finite population correction (f.p.c.) factor is still present in equation (3.4). By definition, the take-all stratum needs to have a finite population in order to get a finite sample size. Also, ignoring the f.p.c. would not lead to a zero variance for the take-all stratum. Considering the f.p.c. in this kind of population representation has been previously used by Dalenius-Gurney (1951).

The a_h in equation (2.3) can also be represented using the quantities (3.1), (3.2) and (3.3). In the case of the N -proportional power allocation, we get:

$$a_h = \frac{W_h^p}{\sum_{h=1}^{L-1} W_h^p} \quad (3.5)$$

for $h=1, \dots, L-1$.

For the Y-proportional power allocation, the following is obtained:

$$a_h = \frac{(W_h \mu_h)^p}{\sum_{h=1}^{L-1} (W_h \mu_h)^p} \quad (3.6)$$

where $0 < p < \infty$.

In this paper, the Y-proportional power allocation will mainly be considered but the calculations can also be performed for the N-proportional power allocation and, in fact, for any kind of allocation represented by some a_h where $\sum_{h=1}^{L-1} a_h = 1$. Putting equation (3.6) into (3.4), we get

$$n = N W_L + \frac{N \left[\sum_{h=1}^{L-1} (W_h \sigma_h)^2 (W_h \mu_h)^{-p} \right] \left[\sum_{h=1}^{L-1} (W_h \mu_h)^p \right]}{N c^2 \mu^2 + \sum_{h=1}^{L-1} W_h \sigma_h^2} \quad (3.7)$$

In order to find the optimal boundaries $b(1), \dots, b(L-1)$ such that the sample size n will be minimum, the derivatives of equation (3.7) are taken with respect to $b(1), \dots, b(L-1)$, respectively, and equated to zero. The resulting equations are:

For $h=1, \dots, L-2$,

$$\begin{aligned} & [F T_h - F T_{h+1}] b_{(h)}^2 + \\ & [F K_h - 2\mu_h F T_h - F K_{h+1} + \\ & 2\mu_{h+1} F T_{h+1} + 2\mu_h AB - 2\mu_{h+1} AB] b_{(h)} + \\ & [F T_h \mu_h^2 + F T_h \sigma_h^2 - F T_{h+1} \mu_{h+1}^2 - \\ & F T_{h+1} \sigma_{h+1}^2 - AB\mu_h^2 + AB\mu_{h+1}^2] = 0 \end{aligned} \quad (3.8)$$

and for $h = L-1$,

$$\begin{aligned} & [F T_{L-1} - AB] b_{(L-1)}^2 + \\ & [F K_{L-1} - 2\mu_{L-1} F T_{L-1} + 2\mu_{L-1} AB] b_{(L-1)} + \\ & [F T_{L-1} \mu_{L-1}^2 + F T_{L-1} \sigma_{L-1}^2 - AB\mu_{L-1}^2 - F T_{L-1}^2] = 0 \end{aligned} \quad (3.9)$$

where

$$\begin{aligned} A &= \sum_{h=1}^{L-1} (W_h \mu_h)^p \\ B &= \sum_{h=1}^{L-1} (W_h \sigma_h)^2 (W_h \mu_h)^{-p} \\ F &= N c^2 \mu^2 + \sum_{h=1}^{L-1} W_h \sigma_h^2 \\ K_h &= B p (W_h \mu_h)^{p-1} - A p (W_h \sigma_h)^2 (W_h \mu_h)^{-p-1} \\ T_h &= A W_h (W_h \mu_h)^{-p}. \end{aligned}$$

Labeling the coefficient of $b_{(h)}^2$ as α_h , the coefficient of $b_{(h)}$ as β_h and the remaining terms as γ_h , equations (3.8) and (3.9) can be represented as quadratic equations of the form $\alpha_h b_{(h)}^2 + \beta_h b_{(h)} + \gamma_h = 0$. However, as pointed out by Sethi (1963), the terms α_h, β_h and γ_h are themselves functions of $b(1), \dots, b(L-1)$ through the integrals (3.1), (3.2) and (3.3). Using Sethi's (1963) approach, equations (3.8) and (3.9) can easily be solved using the following iterative method:

- STEP 1 : Start with some arbitrary boundaries $b'_1 < \dots < b'_{L-1}$.
- STEP 2 : Calculate the proportions W'_h , the means μ'_h and the variances σ'^2_h (from equations (3.1), (3.2) and (3.3), respectively) based on these boundaries, $h=1, \dots, L-1$.
- STEP 3 : Replace the initial set of boundaries by b''_1, \dots, b''_{L-1} where

$$b''_h = \frac{-\alpha'_h + \sqrt{\beta'^2_h - 4\alpha'_h \gamma'_h}}{2\alpha'_h}, \quad h=1, \dots, L-1. \quad (3.10)$$

- STEP 4 : Repeat steps 2 and 3 till two consecutive sets are either identical or differ by negligible quantities, i.e.

$$\max_{h=1}^{L-1} |b''_h - b'_h| < \epsilon \text{ for some } \epsilon > 0. \quad (3.11)$$

It should be noted that it can be proved that the sign before the square root ($\sqrt{\quad}$) is positive if we assume that b'_h lies between μ'_h and μ'_{h+1} .

The difficulty of using the above algorithm is that some knowledge of $f(y)$, the "approximate" density, is required. Since the population considered is finite, it is possible to overcome this difficulty by replacing the quantities (3.1), (3.2) and (3.3) by corresponding expressions based on the finite population property. Hence, proceeding as in Cochran (1977), the infinite population parameters given by expressions (3.1), (3.2) and (3.3) can be replaced by their finite population counterparts. That is:

$$W_h = \frac{N_h}{N} \quad (3.12)$$

$$\bar{y}_h = \frac{1}{N_h} \sum_{j=b_{(h-1)}+1}^{b_{(h)}} y_{(j)} \quad (3.13)$$

$$S_h^2 = \frac{1}{N_h-1} \sum_{j=b_{(h-1)}+1}^{b_{(h)}} y_{(j)}^2 - N_h \bar{y}_h^2 \quad (3.14)$$

for $h=1, \dots, L$.

Using these last quantities, the problem described in section 2 of finding boundaries $b_{(1)}, \dots, b_{(L-1)}$ such that the overall sample size n is minimized for a given level of reliability c and a specific allocation scheme can easily be solved by the following iterative method:

- STEP 0 : Sort the population y_1, \dots, y_N in ascending order and set $b_{(0)} = y_{(1)}$ and $b_{(L)} = y_{(L)}$.
- STEP 1 : Start with some arbitrary boundaries such that $b_{(0)} < b'_{(1)} < \dots < b'_{(L-1)} < b_{(L)}$.
- STEP 2 : Calculate the proportions W_h' , the mean \bar{y}_h' and the variance $S_h'^2$ (from equations (3.12), (3.13) and (3.14) respectively) based on these boundaries, $h=1, \dots, L-1$.
- STEP 3 : Replace the initial set of boundaries by $b''_{(1)}, \dots, b''_{(L-1)}$ where

$$b''_{(h)} = \frac{-\alpha_h' + \sqrt{\beta_h'^2 - 4\alpha_h'\gamma_h'}}{2\alpha_h'}, \quad h=1, \dots, L-1.$$

- STEP 4 : Repeat step 2 and 3 till two consecutive sets are either identical or differ by negligible quantities, i.e.

$$\max_{h=1}^{L-1} |b''_{(h)} - b'_{(h)}| < \epsilon \text{ for some } \epsilon > 0.$$

The use of this algorithm with real data will be compared to others in the next section.

4. SOME ILLUSTRATIONS

In order to display results given in Section 3, we will use data obtained from the Annual Retail Trade and Wholesale Trade Surveys conducted at Statistics Canada. These surveys measure the sales of companies whose principal business is retailing or wholesaling respectively. Three populations have been used to illustrate the algorithm. They are, respectively, other products in Wholesale in Quebec (Population 1), other foods in Wholesale in Manitoba (Population 2), and appliances, television, radio and stereo stores in Retail in Quebec (Population 3). Those populations have been chosen to reflect different combinations of population sizes: high, medium and low.

The numerical results provided by the algorithm will be compared to those obtained using two other methods. The first method is to simply stratify the population using the cumulative square root f rule given by Dalenius — Hodges (1959). The second method is to determine the cut-off boundary between take-all and take-some strata using the approximation given by Hidiroglou (1986) and then to apply the cumulative square root f rule to stratify the non take-all population into a number of take-some strata. The different methods will be labelled as i) Cum $f^{\frac{1}{2}}$ rule for the Dalenius — Hodges (1959) method, ii) mixture for the stratification using the Hidiroglou (1986) and Dalenius — Hodges (1959) method, and iii) optimum for the currently proposed algorithm. The sole use of the Dalenius — Hodges (1959) method is not realistic because it would, in practice, only be used after the take-all stratum had been identified using some given arbitrary rule. However, we display the sole use of this method to caution against its blind use in the context of highly skewed populations.

The Hidiroglou (1986) cut-off point is obtained via the following iterative process:

$$b_{TA}'' = \mu_{[N-t']} = \left\{ \frac{N-t'-1}{(N-t')^2} N^2 c^2 \bar{y}^2 + S_{[N-t']}^2 \right\}^{\frac{1}{2}} \quad (4.1)$$

where

$$\mu_{[N-t']} = \frac{1}{N-t'} \sum_{i=1}^{N-t'} Y(i) \quad (4.2)$$

and

$$S_{[N-t']}^2 = \frac{1}{N-t'-1} \sum_{i=1}^{N-t'} (Y(i) - \mu_{[N-t']})^2$$

The number of take-all units obtained for each step of this iterative process is t' . The starting point for this approximation is

$$b_{TA}' = \mu_{[N]} + \{N c^2 \bar{Y}^2 + S_{[N]}^2\}^{\frac{1}{2}} \quad (4.3)$$

The stopping point for (4.1) is reached when the following inequality is satisfied:

$$0 \leq 1 - n(t'')/n(t') < 0.10 \quad (4.4)$$

where

$$n(t') = t' + \frac{(N-t')^2 S_{[N-t']}^2}{(N c \bar{Y})^2 + (N-t') S_{[N-t']}^2}, \quad (4.5)$$

Tables 1 and 2 display the results for a large population (Population 1) and a small population (Population 2) for a number of different coefficients of variation and power allocations. Table 3 displays the results for the large population (Population 1) and a medium population (Population 3) by varying the number of strata. The allocation of the sample to the take-some strata is the power Y -proportional scheme, for the three tables. The contents of these tables is as follows:

1. the coefficient of variation, c
2. the power of the allocation, p
3. the stratum h population size, N_h
4. the stratum h sample size, n_h
5. the total sample size, "total"
6. the boundary between stratum h and $h+1$, $b_{(h)}$.

The following conclusions can be drawn from Tables 1 and 2. The use of the cumulative square root f rule to determine boundary points is very inefficient in the present context. Substantial gains, in terms of sample size reduction, are made by using the mixture rule. For the three strata used in those two tables, further reductions in sample size in the order of 20% can be achieved by using the optimum rule. For a given fixed coefficient of variation, the variation of the power "p"

has a minor impact on the resulting sample size. As expected, sample sizes increase when the required level of reliability, c , is decreased (for a fixed power allocation). The optimum method declares less take-all units (stratum 3) than the mixture method, or stated another way, the take-all-take-some boundary is higher for the optimum than for the corresponding boundary for the mixture. The cumulative square root rule loses its efficiency in the take-all-take-some boundary determination. It is readily observed that the boundary for this method is significantly higher than those obtained with the other methods.

In Table 3, we only compare the mixture and optimum methods for two populations, varying the number of strata, for a fixed coefficient of variation and Y -proportional power allocation. Similar conclusions as to those drawn from Tables 1 and 2 hold. The effect of increasing the number of strata is to reduce the number of samples units for both methods. However, the reduction becomes more pronounced for the optimum method as the number of strata increases.

5. CONCLUSION

The optimal stratification, of a skewed population into a take-all stratum and a number of take-some strata, has provided a substantial reduction in overall sample size for given reliability constraints. The method can be adopted to any type of allocation and to any number of strata. The take-all condition can also be excluded.

The method is dual, in the sense that, either the sampling variance can be minimized for a fixed sample size, or the sample size can be minimized for a fixed sampling variance.

The algorithm, which is recursive in nature, converges quickly. It is simple to implement on the computer using SAS, FORTRAN, or any other high level language.

REFERENCES

Cochran, W.G. (1977), "Sampling Techniques, 3rd edition", John Wiley & Sons, New York, 428 pages.
Dalenius, T. (1950), "The problem of optimum stratification", Skandinavisk Aktuarietidskrift, 33, 203-13.
Dalenius, T. and Gurney, M. (1951), "The problem of optimum stratification.II", Skandinavisk Aktuarietidskrift, 34, 133-48.
Dalenius, T., and Hodges, Jr. (1959), "Minimum

Variance Stratification", Skandinavisk Aktuarietidskrift, 54, 88-101.
Glasser, G. J. (1962), "On the Complete Coverage of Large Units in a Statistical Study", Review of the International Statistical Institute, 30, 28-32.
Hidioglou, M.A. (1986), "The Construction of a Self-Representing Stratum of Large Units in Survey Design", The American Statistician, 40, 27-31.
Sethi, V.K. (1963), "A Note on Optimum Stratification of Populations for Estimating the Population Means", Australian Journal of Statistics, 5, 20-33.

TABLE 1
Effect of Varying Coefficient of Variation and Power Allocation to Sample Sizes for Three Stratification Methods (Population 1 – Size = 1221)

| c | p | Strata | Stratification Method | | | | | | | | |
|------|------|--------|----------------------------|-------|------------|---------|-------|-----------|---------|-------|-----------|
| | | | Cum $f^{\frac{1}{2}}$ Rule | | | Mixture | | | Optimum | | |
| | | | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ |
| 0.05 | 0.25 | 1 | 1196 | 177* | | 1017 | 16 | | 891 | 11 | |
| | | 2 | 20 | 20 | 3,715,320 | 152 | 14 | 465,180 | 290 | 13 | 302,912 |
| | | 3 | 5 | 5 | 14,786,280 | 52 | 52 | 1,131,961 | 40 | 40 | 1,835,930 |
| | | Total | | 202 | | | 82 | | | 64 | |
| 0.05 | 0.50 | 1 | 1196 | 178* | | 1017 | 16 | | 863 | 10 | |
| | | 2 | 20 | 20 | 3,715,320 | 152 | 13 | 465,180 | 318 | 14 | 289,422 |
| | | 3 | 5 | 5 | 17,786,280 | 52 | 52 | 1,131,961 | 40 | 40 | 1,832,038 |
| | | Total | | 203 | | | 81 | | | 64 | |
| 0.01 | 1.00 | 1 | 1196 | 616* | | 751 | 37 | | 687 | 36 | |
| | | 2 | 20 | 20 | 3,715,320 | 215 | 34 | 196,840 | 374 | 78 | 162,068 |
| | | 3 | 5 | 5 | 14,786,280 | 255 | 255 | 383,033 | 160 | 160 | 564,076 |
| | | Total | | 641 | | | 326 | | | 274 | |
| 0.05 | 1.00 | 1 | 1196 | 180* | 3,715,320 | 1017 | 16 | | 858 | 8 | |
| | | 2 | 20 | 20 | 14,786,280 | 152 | 11 | 465,180 | 323 | 16 | 271,920 |
| | | 3 | 5 | 5 | | 52 | 52 | 1,131,961 | 40 | 40 | 1,867,254 |
| | | Total | | 205 | | | 79 | | | 64 | |
| 0.10 | 1.00 | 1 | 1196 | 56* | | 1073 | 7 | | 1007 | 7 | |
| | | 2 | 20 | 20 | 3,715,320 | 109 | 4 | 592,900 | 191 | 9 | 442,357 |
| | | 3 | 5 | 5 | 14,786,280 | 39 | 39 | 1,953,113 | 23 | 23 | 4,032,950 |
| | | Total | | 81 | | | 50 | | | 39 | |

* Requires over allocation to satisfy coefficient of variation.

TABLE 2

**Effect of Varying Coefficient of Variation and Power Allocation
to Sample Sizes for Three Stratification Methods
(Population 2 — Size = 44)**

| c | p | Strata | Stratification Method | | | | | | | | |
|------|------|--------|----------------------------|-----------------|-------------|---------|-----------------|------------|---------|-----------------|------------|
| | | | Cum $f^{\frac{1}{2}}$ Rule | | | Mixture | | | Optimum | | |
| | | | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ |
| 0.05 | 0.25 | 1 | 42 | 38 | | 32 | 1 | | 29 | 1 | |
| | | 2 | 1 | 1* | 137,939,900 | 6 | 1 | 4,708,409 | 11 | 1 | 3,029,455 |
| | | 3 | 1 | $\frac{1}{1}$ | 459,739,000 | 6 | $\frac{6}{6}$ | 10,622,301 | 4 | $\frac{4}{4}$ | 17,461,464 |
| | | Total | | $\frac{40}{40}$ | | | $\frac{8}{8}$ | | | $\frac{6}{6}$ | |
| 0.05 | 0.50 | 1 | 42 | 38 | | 32 | 1 | | 28 | 1 | |
| | | 2 | 1 | 1* | 137,939,900 | 6 | 1 | 4,708,409 | 12 | 1 | 2,582,819 |
| | | 3 | 1 | $\frac{1}{1}$ | 459,739,000 | 6 | $\frac{6}{6}$ | 10,622,301 | 4 | $\frac{4}{4}$ | 17,640,325 |
| | | Total | | $\frac{40}{40}$ | | | $\frac{8}{8}$ | | | $\frac{6}{6}$ | |
| 0.01 | 1.00 | 1 | 42 | 42 | | 25 | 1 | | 25 | 1 | |
| | | 2 | 1 | 1 | 137,939,900 | 5 | 1 | 1,059,550 | 10 | 4 | 1,153,322 |
| | | 3 | 1 | $\frac{1}{1}$ | 459,739,000 | 14 | $\frac{14}{16}$ | 3,742,377 | 9 | $\frac{9}{14}$ | 5,969,271 |
| | | Total | | $\frac{44}{44}$ | | | $\frac{16}{16}$ | | | $\frac{14}{14}$ | |
| 0.05 | 1.00 | 1 | 42 | 38 | | 32 | 1 | | 26 | 1 | |
| | | 2 | 1 | 1* | 137,939,900 | 6 | 1 | 4,708,409 | 14 | 2 | 1,779,500 |
| | | 3 | 1 | $\frac{1}{1}$ | 459,739,000 | 6 | $\frac{6}{8}$ | 10,622,301 | 4 | $\frac{4}{7}$ | 17,349,902 |
| | | Total | | $\frac{40}{40}$ | | | $\frac{8}{8}$ | | | $\frac{7}{7}$ | |
| 0.10 | 1.00 | 1 | 42 | 30 | | 34 | 1 | | 28 | 1 | |
| | | 2 | 1 | 1* | 137,939,900 | 6 | 1 | 4,848,218 | 13 | 1 | 2,413,800 |
| | | 3 | 1 | $\frac{1}{1}$ | 459,739,000 | 4 | $\frac{4}{6}$ | 16,749,625 | 3 | $\frac{3}{5}$ | 30,091,449 |
| | | Total | | $\frac{32}{32}$ | | | $\frac{6}{6}$ | | | $\frac{5}{5}$ | |

* Requires over allocation to satisfy coefficient of variation.

TABLE 3

**Effect of Increasing the Number of Strata on
Sample Sizes for Two Stratification Methods
p=1, c=0.05**

| Population 1 (N=1221) Stratification Method | Strata | Number of Strata | | | | | | | | |
|--|--------|------------------|-------|-----------|-------|-------|-----------|-------|-------|-----------|
| | | 3 | | | 4 | | | 5 | | |
| | | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ | N_h | n_h | $b(h)$ |
| Mixture | 1 | 1017 | 16 | | 897 | 6 | | 823 | 3 | |
| | 2 | 152 | 11 | 465,180 | 194 | 5 | 311,117 | 194 | 2 | 245,090 |
| | 3 | 52 | 52 | 1,131,961 | 78 | 4 | 641,252 | 101 | 2 | 465,180 |
| | 4 | | | | 52 | 52 | 1,131,961 | 51 | 2 | 751,297 |
| | 5 | | | | | | | 52 | 52 | 1,131,961 |
| | Total | | 79 | | | 67 | | | 61 | |
| Optimum | 1 | 858 | 8 | | 704 | 3 | | 655 | 2 | |
| | 2 | 323 | 16 | 271,920 | 373 | 7 | 173,981 | 358 | 4 | 149,327 |
| | 3 | 40 | 40 | 1,867,254 | 112 | 6 | 604,869 | 163 | 5 | 453,114 |
| | 4 | | | | 32 | 32 | 2,676,449 | 29 | 4 | 1,522,329 |
| | 5 | | | | | | | 16 | 16 | 5,810,487 |
| | Total | | 64 | | | 48 | | | 31 | |
| Population 3 (N=161) | | | | | | | | | | |
| Mixture | 1 | 106 | 6 | | 84 | 2 | | 71 | 1 | |
| | 2 | 39 | 6 | 265,480 | 38 | 2 | 185,320 | 35 | 1 | 155,260 |
| | 3 | 16 | 16 | 553,255 | 23 | 2 | 335,620 | 22 | 1 | 265,480 |
| | 4 | | | | 16 | 16 | 553,255 | 17 | 1 | 385,720 |
| | 5 | | | | | | | 16 | 16 | 553,255 |
| | Total | | 28 | | | 22 | | | 20 | |
| Optimum | 1 | 86 | 4 | | 55 | 1 | | 34 | 1 | |
| | 2 | 65 | 9 | 199,415 | 61 | 3 | 125,572 | 51 | 1 | 83,594 |
| | 3 | 10 | 10 | 680,942 | 39 | 5 | 312,769 | 42 | 2 | 192,215 |
| | 4 | | | | 6 | 6 | 826,942 | 29 | 3 | 382,236 |
| | 5 | | | | | | | 5 | 5 | 906,894 |
| | Total | | 23 | | | 15 | | | 12 | |