



HAL
open science

Allocation constraints in stratification

Marcin Kozak, Pawel Jankowski

► **To cite this version:**

Marcin Kozak, Pawel Jankowski. Allocation constraints in stratification. *Communications in Statistics - Simulation and Computation*, 2008, 37 (09), pp.1763-1775. 10.1080/03610910802278842 . hal-00514328

HAL Id: hal-00514328

<https://hal.science/hal-00514328>

Submitted on 2 Sep 2010

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Allocation constraints in stratification

Journal:	<i>Communications in Statistics - Simulation and Computation</i>
Manuscript ID:	LSSP-2007-0207.R2
Manuscript Type:	Original Paper
Date Submitted by the Author:	13-Jun-2008
Complete List of Authors:	Kozak, Marcin; Warsaw University of Life Sciences, Experimental Design and Bioinformatics Jankowski, Paweł; Warsaw University of Life Sciences, Biometry
Keywords:	constraints, optimization, optimum stratification, sample allocation
Abstract:	When a finite population is to be stratified, one of constraints in stratification is that sample sizes from strata may not be greater than the corresponding strata sizes and may not be smaller than two. There are several ways of treating this allocation constraint, each providing an alternative approach to stratification. In the paper it is shown that a choice of the approach has a bearing on stratification efficiency. Unfortunately, no particular approach out of the four compared is shown to be the most efficient for each population studied. In addition, the approaches are applied to stratify a real population.



Allocation constraints in stratification

Marcin Kozak^{1*}, Paweł Jankowski²

¹ Department of Experimental Design and Bioinformatics, Warsaw University of Life Sciences, Nowoursynowska 159, 02-787, Warsaw, Poland,

² Department of Biometry, Warsaw University of Life Sciences, Nowoursynowska 159, 02-787, Warsaw, Poland,

*Corresponding author, Department of Experimental Design and Bioinformatics, Warsaw University of Life Sciences
Nowoursynowska 159, 02-787, Warsaw, Poland, e-mail:

m.kozak@omega.sggw.waw.pl

Running head: Allocation constraints in stratification

Allocation constraints in stratification

Abstract

When a finite population is to be stratified, one of constraints in stratification is that sample sizes from strata may not be greater than the corresponding strata sizes and may not be smaller than two. There are several ways of treating this allocation constraint, each providing an alternative approach to stratification. In the paper it is shown that a choice of the approach has a bearing on stratification efficiency. Unfortunately, no particular approach out of the four compared is shown to be the most efficient for each population studied. In addition, the approaches are applied to stratify a real population.

Keywords: constraints, optimization, optimum stratification, sample allocation

Mathematics Subject Classification: Primary 62D05, Secondary 62P20, 62P25.

1. Introduction

In this paper we consider an optimization approach to stratification, which has been recently proved to be superior to approximate stratification procedures (see Kozak and Verma (2006) and the citations therein). Suppose we aim to stratify a finite population U based on an auxiliary (stratification) variable X . Let the aim of stratification be minimizing the variance of an estimator of the population total of a study variable Y subject to fixed sample size n . At the design stage of a survey it is usually assumed that a survey variable (Y) and stratification variable (X) be the same and that there be no non-responses. This is, of course, never the case in practice, yet such an approach is common and is not thought of as controversial. Furthermore, let us consider the common practical situation in which the survey and stratification variables are

positively skewed; then, the most efficient approach is to construct a so-called take-all stratum, from which all the elements are taken to the sample (e.g., Hidiroglou, 1986; Lavallée and Hidiroglou, 1988).

The objective function to be minimized in the problem in question is the variance $Var(\hat{t}_X)$ of an estimator \hat{t}_X of the population total of X . Assuming that the L th stratum is the take-all one, the variance, under the take-all stratum approach, takes the form

$$Var(\hat{t}_X) = \sum_{h=1}^{L-1} S_h^2 W_h^2 \left(\frac{1}{n_h} - \frac{1}{N_h} \right) \quad (1)$$

$$\hat{t}_X = \sum_{h=1}^L \frac{N_h}{n_h} \sum_{k=1}^{n_h} X_{hk}; \quad t_X = \sum_{h=1}^L \sum_{k=1}^{N_h} X_{hk}; \quad W_h = N_h / \sum_{h=1}^L N_h;$$

$$S_h^2 = \frac{1}{N_h - 1} \sum_{k=1}^{N_h} (X_{hk} - \bar{X}_h)^2 \text{ for } h = 1, \dots, L-1; \quad \bar{X}_h = \frac{1}{N_h} \sum_{k=1}^{N_h} X_{hk}$$

where \hat{t}_X is the unbiased estimator of the population total t_X of X ; S_h^2 is the population variance of the variable X restricted to the h th stratum; n_h is the sample size from the h th stratum of size N_h ; X_{hk} is the value of X for the k th population element of the h th stratum; and \bar{X}_h is the population mean of X restricted to the stratum h . In equation (1), we have considered the most classical unbiased estimation of the population total under stratified sampling; see, e.g., Särndal et al. (1992).

A vector of stratification points, say $\mathbf{a} = (a_1, \dots, a_{L-1})^T$, which explicitly defines the subdivision of the population U into strata, is a vector of parameters to be searched for (e.g., Lednicki and Wiczorkowski, 2003). It is to be noted that the variance (1) does not directly involve the parameters sought. An objective function $f(\mathbf{a})$ can be written in a general form as

$$f(\mathbf{a}) = Var(\hat{t}_X) \quad (2)$$

where $Var(\hat{t}_X)$ is given in equation (1). The constraints for the function (2) are as follows:

$$2 \leq n_h \leq N_h, \quad h = 1, \dots, L-1 \quad (3)$$

$$N_h \geq 2; \quad h = 1, \dots, L \quad (4)$$

$$N_L + \sum_{h=1}^{L-1} n_h = n \quad (5)$$

Fulfilling the constraints (3) and (4) is required to ensure that the variance (1) can be evaluated. (Note that the constraint $N_L \geq 2$ is not required here, but it is reasonable to use it in order to obtain a take-all stratum comprising at least two elements.) Sample sizes n_h from strata are usually determined through the Neyman optimum sample allocation, which aims at minimizing the variance (1); after adjusting the formula for the take-all stratum approach, the sample sizes are given by

$$n_h = (n - N_L) \frac{W_h S_h}{\sum_{h=1}^{L-1} W_h S_h}, \quad h = 1, \dots, L-1; \quad n_L = N_L \quad (6)$$

There are two possible ways of treating the constraints (3): (i) one does not accept the solution in which any n_h provided by the formula (6) does not fulfil the constraints (3), and changes the stratification points (such an approach was applied, for instance, by Lednicki and Wieczorkowski (2003)); and (ii) one does not accept the allocation and searches for the optimum allocation using numerical optimization. The first approach is obviously easier to implement and provides less time-consuming computation. However, intuition makes us suppose that this approach may give rise to rejecting solutions that either are optimal or may be connections between stratification points considered in a particular step and the optimum points (or a path leading to the optimum points).

The option (i) of treating the constraints can be applied in two manners. First, one can reject points not fulfilling the constraints (3). Second, one can apply the following procedure to adjust the sample sizes for the constraints: determine n_h through the allocation (6) and apply the following formula (for $h = 1, \dots, L-1$)

$$n_h = 2 \text{ if } n_h < 2; n_h = N_h \text{ if } n_h > N_h \quad (7)$$

This procedure makes us accept solutions that would be rejected by the first manner, in which way a set of possible solutions is widened.

In this paper we compare the following approaches to treating the allocation constraints in stratification:

(A) Approach based on not accepting a solution (stratification) in which any n_h

provided by the formula (6) does not fulfil the constraints (3).

(B) Approach based on applying the allocation (6) with the adjustment (7).

(C) Approach based on applying numerical optimization to allocate the sample: if n_h 's

provided by the formula (6) do not fulfil the constraints (3), solve the following

problem to find the allocation. Given a vector of stratification points \mathbf{a} , find such

n_h 's, $h = 1, \dots, L-1$, that minimize the objective function (1), i.e.,

$$f(n_1, \dots, n_{L-1}) = \text{Var}(\hat{f}_X), \text{ under the constraints (3) and (5).}$$

Kozak (2004b) showed that results of stratification determined by numerical optimization depend on stratification points that are used as initial parameters in the optimization. Therefore, we will consider an additional approach as follows:

(D) Approach based on applying approach C with strata boundaries provided by

approach B taken as initial parameters in optimization.

Hence, a question to answer is, does the choice of an approach of treating the allocation constraints in stratification have an influence on stratification efficiency? The aim of the paper is to answer this question through a simulation study.

2. Design of experiment

The following aspects of a population and a stratification variable were considered in the experiment: (i) population size N , viz., $N = \{1000, 2000, 5000, 10000, 15000\}$; (ii) number L of strata to be constructed, viz., $L = \{3, 5, 7, 9\}$; (c) sample fraction $f = n/N$ (n being the assumed sample size), viz., $f = \{0.1, 0.2\}$; and (d) parameter σ of the distribution of the stratification variable (see below), viz., $\sigma = \{0.4, 0.6, 0.8\}$. For the sake of convenience, below the population quantities N , L , f , and σ will be referred to as factors.

Stratification variables were generated based on the following formula:

$$X = [\exp(Z)],$$

where Z is the realization of an $N(10, \sigma^2)$ variable (a normal random variable with mean 10 and standard deviation σ) and the function $[\cdot]$ stands for rounding to integers (to simulate the most often practical situation). As a result of such generation, the variables were positively skewed; the greater the σ value, the greater the skewness was.

For each combination of $N \times L \times f \times \sigma$, 100 independent populations (stratification variables) were generated; thus there were 12000 populations altogether. For an i th population, the four stratification approaches of study were applied and the coefficient of variation $cv_{ki}(\hat{t}_X^i)$ (k referring to the k th approach to stratification, $k = A, B, C, D$) of the estimator \hat{t}_X^i was evaluated using the formula

$$cv_{ki}(\hat{t}_X^i) = (t_X^i)^{-1} \sqrt{Var_k(\hat{t}_X^i)} \quad (8)$$

where $Var_k(\hat{t}_X^i)$ is the variance (1) of the estimator \hat{t}_X^i under the k th approach to stratification, and t_X^i is the total of X in the i th population.

As Lednicki and Wieczorkowski (2003) did, to perform stratification we have applied *optim* function, which implements the algorithm of the simplex method of Nelder and Mead (1965), available in R language and environment (R Development Core Team 2007). Following Kozak's (2004a) results on efficiency of approximate stratification points used as initial parameters in optimization, stratification points determined by Mahalanobis's (1952) procedure were taken as the initial vector of parameters in optimization. However, because this procedure does not take account of the take-all stratum approach, the point defining the last stratum was changed in such a way that this stratum comprised five population elements with the largest values of the stratification variable. (Note that the take-all stratum contained five elements at the initial stage of stratification, but later the number of its elements was not limited to five.) Whenever these stratification points did not fulfil the constraints (3) and/or (4), we used the points provided by (i) the Ekman (1959) procedure; (ii) then, whenever Ekman's points failed, the Dalenius and Hodges (1959) procedure were applied; and finally, (iii) whenever Dalenius and Hodges' points failed, we used the Gunning and Horgan (2004) procedure. Each such procedure was applied with the above-mentioned adjustment for the take-all stratum approach. However, whenever all these procedures failed to fulfil the constraints (3), the initial strata were constructed based on the following procedure. First, the five-element take-all stratum had been constructed, and then the remaining part of the population was subdivided into $L-1$ strata of equal sizes (or nearly equal, if $(N-5)/(L-1)$ was not an integer).

1
2 In the approaches C and D, R's function *optim* was applied to determine the
3 optimum sample allocation; as the initial parameters (sample sizes from strata) in the
4 optimization, the $(L - 1)$ -vector of twos was taken.
5
6
7

8 All the computation was performed in R (R Development Core Team, 2007) using
9 self-implemented functions, which can be obtained from the corresponding author upon
10 request.
11
12
13

14 15 16 **3. Results** 17

18
19
20 Results of the simulation study are presented in Tables 1-5, for $N = 1000, 2000, 5000,$
21 10000 and 15000 , respectively. For a particular population we calculated ranks for the
22 values of the coefficient of variation (cv), given by Eq. (8), of the estimator obtained
23 under the four approaches. Then we determined the number of times in which (i) a
24 particular approach was the best (cv obtained under approach A, C, or D had rank 1; cv
25 obtained under approach B had rank 1.5 provided that cv obtained under approach D
26 had rank 1.5, too); (ii) a particular approach was the worst (cv obtained under approach
27 A or C had rank 4; cv obtained under approach B had rank 3.5, provided that cv
28 obtained under approach D had rank 3.5, or rank 4; cv obtained under approach D had
29 rank 3.5 provided that cv obtained under approach B had rank 3.5, too). Through this
30 analysis of ranks approach B was recognized as the best in a situation when approach D
31 had not improved its results and approaches A and C had been worse (we would not
32 recognize approach D as the best then, since its application did not provide any gain in
33 efficiency in comparison to its initial parameters that had been provided by approach B).
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 the results of approach B, of which the results of approaches A and C had been better
3
4 (in such a situation approach B was also recognized as the worst).
5
6
7

8 [Table 1]
9

10 [Table 2]
11

12 [Table 3]
13

14 [Table 4]
15

16 [Table 5]
17
18
19

20 In addition, Table 6 contains mean ranks determined for the approaches under a
21 particular level of the population quantities studied; from this table it follows that in
22 general approach D appears to be the best (in the sense that it provides the most efficient
23 stratification points); approaches A, B and C seem to provide similar results (inferior to
24 those of approach D), even though a slight tendency of approach C to be better than
25 approaches A and B has been detected. Such a result has been obtained for all the levels
26 of the factors studied except for $L = 3$, in which case all the approaches provided
27 similar mean ranks though approach D appeared to be slightly better than the other
28 approaches. For all the other factor levels, cv obtained under approach D had the
29 smallest mean rank. Usually cv obtained under approach C had a little smaller mean
30 rank than cvs obtained under approaches A and B, but for some combinations cvs
31 obtained under approaches A, B and C had similar mean ranks. Differences between the
32 values of cv obtained by different approaches were sometimes meaningful (e.g., it
33 sometimes happened that $cv_{worst} > 1.5cv_{best}$, the first cv referring to the worst and the
34 second to the best approach).
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 [Table 6]
3
4
5

6 Several interesting situations occurred. Under $\sigma = 0.4$ and 0.6 , the larger L , the more
7 often approach D was the best. However, this situation did not occur under $\sigma = 0.8$. In
8 general, under $L = 5, 7$ and 9 , approaches A and B, and sometimes C, were quite often
9 the worst. In general, approach B was seldom the best; neither was approach C,
10 although under $\sigma = 0.8, f = 0.2$, and $L = 7$ and 9 it was usually the best. The hypothesis
11 that approach C and/or D may always be the best, which was mentioned in Introduction,
12 has not been proven correct by our experiment. Nonetheless, in most situations
13 approach D often appeared the best; in addition, it seldom was the worst. That approach
14 D was usually better than approach C is easy to explain—it resulted from more efficient
15 initial strata points used in the former than those used in the latter. Nonetheless, let us
16 recall the combinations $\sigma = 0.8, f = 0.2$, and $L = 7$ and 9 , in which, without any
17 reasonable and explicable reason, approach D was usually worse than approach C.
18
19

20 Approach A has one important drawback that must be mentioned here. There were
21 many populations, especially under high N, σ , and f values, for which the optimization
22 was unable to perform this stratification based on all the initial values mentioned, for
23 which reason this approach failed to stratify such populations. This situation did not
24 occur for any other approach.
25
26

27 We have not studied the approach in which the results of approach A would be
28 taken as initial parameters to perform optimization in approach C. The main reason was
29 that such an approach would fail in a situation in which approach A did not succeed to
30 provide stratification points fulfilling the constraints (3) and (4).
31
32

33 Based on the results obtained we are not able to choose the best approach explicitly.
34 In practice the best way is to apply all the approaches (even with points provided by the
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2 approach A taken as initial parameters if only they are feasible) and to choose the best
3
4 one. However, if for any reason it is unlikely to be done, approach D should be chosen,
5
6 as the best and the least risky in our experiment. Of course, a simulation study as it is, it
7
8 does consider a somewhat limited range of possible populations and stratification
9
10 approaches (e.g., a population size, variability in a stratification variable, number of
11
12 strata to be constructed, and the like). It is possible that considering other population
13
14 and stratification attributes, we could obtain different results. Nonetheless, this would
15
16 not change our conclusions that we cannot point out any particular approach as the best
17
18 one, and that the choice of an approach does count even though still no explicit
19
20 recommendation can be given.
21

22
23 The most possible reason why approaches C and D were not the best for all the
24
25 populations is that Nelder and Mead's optimization might lead to the local minimum of
26
27 a function optimized. Moreover, initial stratification points have a bearing on the
28
29 optimization results (Kozak, 2004b). We have applied several approximate stratification
30
31 procedures to provide initial stratification points, but there is no certainty that those
32
33 stratification points are really the best ones. Further efforts should be focused on
34
35 determining a procedure that would provide more efficient stratification than Nelder and
36
37 Mead's optimization approach does. Kozak's (2004a) random search algorithm is a very
38
39 promising one, as claimed by Baillargeon and Rivest (2007), but we need to remember
40
41 that as a global optimization method this procedure provides random results. Hence the
42
43 best option is to apply it several times (maybe with various starting points) to ensure
44
45 that the results obtained are indeed globally and not locally optimum. Worth noting is
46
47 that Baillargeon and Rivest (2007) implemented a non-random version of the algorithm,
48
49 which is of course free of the problem of random results.
50
51
52
53
54
55
56
57
58
59
60

Deleted: The main aim of this paper was to show that a choice of treating the constraints in stratification does matter; this has been proven indeed. What is more, we have shown that this choice is very important, since it may cause results of stratification under various treating the constraints be very different. In the next section we will present the application of the four approaches for a real population.

4. Example

Here we apply the four allocations for a real data set SHS available in the package *stratification* (Baillargeon and Rivest, 2007) of R (R Development Core Team, 2007). The set contains data for 16057 units from the 2001 Survey of Household Spending (SHS) Statistics Canada; for stratification we will use one variable, namely “household income before taxes”. The results are given in Table 7.

Apparently approaches A and D were the most and approach C the least efficient for this particular data set. Interestingly, all the approaches provided the same take-all stratum (even though the boundaries for the take-all strata they provided differed, all of them comprised the same eight population units). Note that due to rounding of sample sizes from strata there were some inconsistencies in overall sample sizes as in none of the cases it equalled 1500 (it was either 1499 or 1501), but for so large a sample size these two elements did not make any real difference.

[Table 7]

5. Conclusion

This paper aimed to show that the way of treating the constraint (3) does matter. This has been proven indeed: The choice may cause results of stratification under various treating the constraints be very different in terms of the precision of estimation of a parameter studied. Unfortunately, of four such ways considered in this paper, none was ultimately the best.

1
2 From the results we have concluded that in practical applications the best way is to
3
4 apply all the approaches and to choose the best one for the particular population. If this
5
6 is impossible, approach D should be applied as most often the best and the least risky in
7
8 our experiment.
9

10 11 12 **References**

- 13
14 Baillargeon, S. and L.P. Rivest (2007). stratification: Stratification of Survey
15
16 Populations. R package version 1.0.
17
18 Dalenius, T. and Hodges, J.L. (1959). Minimum Variance Stratification. Journal of the
19
20 American Statistical Association 54:88—101.
21
22 Ekman, G., (1959), An Approximation Useful in Univariate Stratification. Annals of
23
24 Mathematical Statistics 30:219—229.
25
26 Gunning, P. and Horgan, J.M. (2004). A Simple Algorithm for Stratifying Skewed
27
28 Populations. Survey Methodology 30:159—166.
29
30 Hidirolou, M. (1986). The Construction of a Self-Representing Stratum of Large Units
31
32 in Survey Design. The American Statistician 40:27—31.
33
34 Kozak, M. (2004a). Optimal Stratification Using Random Search Method in
35
36 Agricultural Surveys. Statistics in Transition 6 (5):797—806.
37
38 Kozak, M. (2004b). Optimal Stratification with the Auxiliary Variable. Wiadomosci
39
40 Statystyczne (Statistical News) 8:29—34 (in Polish).
41
42 Kozak, M. and Verma, M.R. (2006). Geometric Versus Optimization Approach to
43
44 Stratification: A Comparison of Efficiency. Survey Methodology 32(2):157—163.
45
46 Lavallée, P. and Hidirolou, M. (1988). On the Stratification of Skewed Population.
47
48 Survey Methodology 14:3—43.
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
- Lednicki, B. and Wieczorkowski, R. (2003). Optimal Stratification and Sample Allocation between Subpopulations and Strata. *Statistics in Transition* 6:287—306.
- Mahalanobis, P.C. (1952). Some Aspects of the Design of Sample Surveys. *Sankhya* 12:1—7.
- Nelder, J.A. and Mead, R. (1965). A Simplex Method for Function Minimization. *Computer Journal* 7:308—313.
- R Development Core Team (2007). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria; URL <http://www.R-project.org>.
- Särndal, C.E., Swensson, B. and Wretman, J. (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag.

Table 1. Summary of rank-ordering of efficiency of stratification approaches A, B, C, and D, for $N=1000$ and various combinations of σ, f, L .

σ	f	L	A ^a		A ^b	B ^c		C ^a		D ^e	
			best	worst	failed	best	worst	best	worst	best	worst
0.4	0.1	3	0	0	0	0	0	0	0	7	0
0.4	0.1	5	12	17	0	5	22	13	38	25	3
0.4	0.1	7	20	28	0	9	24	20	42	48	3
0.4	0.1	9	22	19	0	1	24	8	56	69	0
0.4	0.2	3	2	3	0	0	0	0	0	7	0
0.4	0.2	5	18	33	0	3	6	2	6	46	0
0.4	0.2	7	17	39	0	2	20	11	24	66	0
0.4	0.2	9	18	43	3	1	31	12	24	69	0
0.6	0.1	3	1	0	0	0	0	0	1	12	0
0.6	0.1	5	18	21	1	12	16	17	43	33	4
0.6	0.1	7	19	21	1	6	24	19	50	54	3
0.6	0.1	9	14	16	0	2	15	8	68	76	0
0.6	0.2	3	2	6	0	1	0	0	0	9	0
0.6	0.2	5	16	62	0	2	4	0	4	72	0
0.6	0.2	7	17	39	0	0	30	14	22	68	1
0.6	0.2	9	28	37	2	0	31	23	32	49	0
0.8	0.1	3	0	2	0	0	0	1	0	18	0
0.8	0.1	5	12	32	0	10	28	29	25	40	6
0.8	0.1	7	16	21	0	3	16	13	59	66	2
0.8	0.1	9	5	27	0	3	11	8	62	84	0
0.8	0.2	3	2	9	0	0	0	0	0	6	0
0.8	0.2	5	9	58	1	3	3	2	9	73	1
0.8	0.2	7	23	34	7	0	39	19	22	57	1
0.8	0.2	9	13	66	2	1	22	63	12	23	0

^a "A best", "C best", "A worst", and "C worst" indicate number of times in which a particular approach (A or C) had rank 1 (best) or rank 4 (worst); ^b "A failed" indicates number of time in which numerical problems occurred under approach A, so the solution was not be found; ^c "B best" indicates number of times in which this approach had rank 1.5 provided that approach D had rank 1.5, too; ^d "B worst" indicates

1
2
3 number of times in which this approach had rank 4 or 3.5 provided that approach D also had rank 3.5; ^e “D
4 best” indicates number of times in which approach D had rank 1; ^f “D worst” indicates number of times in
5 which this approach had rank 3.5 provided that approach B had rank 1, too.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

Table 2. Summary of rank-ordering of efficiency of stratification approaches A, B, C, and D, for $N=2000$ and various combinations of σ, f, L .

σ	f	L	A ^a		A ^b	B ^c		C ^a		D ^e	
			best	worst	failed	best	worst	best	worst	best	worst
0.4	0.1	3	0	0	0	0	0	0	0	4	0
0.4	0.1	5	19	43	0	1	6	0	2	37	1
0.4	0.1	7	25	35	0	1	31	10	15	62	0
0.4	0.1	9	35	29	0	0	42	11	26	54	0
0.4	0.2	3	2	4	0	0	0	0	0	3	0
0.4	0.2	5	13	42	0	0	4	0	2	62	0
0.4	0.2	7	22	47	0	0	17	6	14	72	0
0.4	0.2	9	16	53	0	0	19	4	20	80	0
0.6	0.1	3	0	1	0	0	0	0	0	5	0
0.6	0.1	5	18	57	0	3	10	1	2	58	0
0.6	0.1	7	31	20	0	0	37	8	33	59	0
0.6	0.1	9	30	19	1	0	40	14	41	56	0
0.6	0.2	3	1	4	0	0	0	0	0	12	0
0.6	0.2	5	10	59	0	0	1	1	3	75	0
0.6	0.2	7	17	39	0	1	18	1	19	80	1
0.6	0.2	9	23	48	1	0	33	28	12	49	0
0.8	0.1	3	0	3	0	0	0	0	0	18	0
0.8	0.1	5	17	54	0	3	7	3	5	62	0
0.8	0.1	7	28	21	0	1	39	20	38	51	1
0.8	0.1	9	31	18	0	1	33	17	49	51	1
0.8	0.2	3	1	10	0	0	0	0	0	19	0
0.8	0.2	5	13	58	4	1	3	2	7	68	1
0.8	0.2	7	10	60	7	2	23	45	7	43	1
0.8	0.2	9	3	82	46	0	18	85	0	12	0

^a "A best", "C best", "A worst", and "C worst" indicate number of times in which a particular approach (A or C) had rank 1 (best) or rank 4 (worst); ^b "A failed" indicates number of time in which numerical problems occurred under approach A, so the solution was not be found; ^c "B best" indicates number of times in which this approach had rank 1.5 provided that approach D had rank 1.5, too; ^d "B worst" indicates

1
2
3 number of times in which this approach had rank 4 or 3.5 provided that approach D also had rank 3.5; ^e “D
4 best” indicates number of times in which approach D had rank 1; ^f “D worst” indicates number of times in
5 which this approach had rank 3.5 provided that approach B had rank 1, too.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

Table 3. Summary of rank-ordering of efficiency of stratification approaches A, B, C, and D, for $N=5000$ and various combinations of σ, f, L .

σ	f	L	A ^a		A ^b	B ^c		C ^a		D ^e	
			best	worst	failed	best	worst	best	worst	best	worst
0.4	0.1	3	0	0	0	0	0	0	0	7	0
0.4	0.1	5	22	35	0	1	3	1	2	41	0
0.4	0.1	7	12	60	0	1	10	2	8	85	0
0.4	0.1	9	8	57	0	0	20	2	12	90	0
0.4	0.2	3	0	0	0	0	0	0	0	8	0
0.4	0.2	5	3	39	0	1	0	2	2	66	0
0.4	0.2	7	17	48	0	0	15	6	14	76	0
0.4	0.2	9	14	49	0	0	18	6	20	80	0
0.6	0.1	3	1	1	0	0	0	0	0	11	0
0.6	0.1	5	10	41	0	1	8	2	6	65	0
0.6	0.1	7	13	47	0	1	20	5	18	81	0
0.6	0.1	9	26	32	0	0	26	2	26	72	0
0.6	0.2	3	0	0	0	0	0	0	0	22	0
0.6	0.2	5	7	47	0	1	3	1	2	73	0
0.6	0.2	7	7	42	0	0	25	9	18	84	0
0.6	0.2	9	4	86	2	0	6	51	8	45	0
0.8	0.1	3	0	6	0	0	0	0	0	11	0
0.8	0.1	5	11	44	0	4	8	1	7	69	0
0.8	0.1	7	22	61	3	0	12	8	10	69	0
0.8	0.1	9	32	17	0	0	37	5	39	63	0
0.8	0.2	3	0	5	0	0	0	0	0	18	0
0.8	0.2	5	14	35	0	1	3	3	7	40	1
0.8	0.2	7	2	76	3	0	23	79	0	19	1
0.8	0.2	9	0	100	85	0	0	96	0	4	0

^a "A best", "C best", "A worst", and "C worst" indicate number of times in which a particular approach (A or C) had rank 1 (best) or rank 4 (worst); ^b "A failed" indicates number of time in which numerical problems occurred under approach A, so the solution was not be found; ^c "B best" indicates number of times in which this approach had rank 1.5 provided that approach D had rank 1.5, too; ^d "B worst" indicates

1
2
3 number of times in which this approach had rank 4 or 3.5 provided that approach D also had rank 3.5; ^e “D
4 best” indicates number of times in which approach D had rank 1; ^f “D worst” indicates number of times in
5 which this approach had rank 3.5 provided that approach B had rank 1, too.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

Table 4. Summary of rank-ordering of efficiency of stratification approaches A, B, C, and D, for $N=10000$ and various combinations of σ, f, L .

σ	f	L	A ^a		A ^b	B ^c		C ^a		D ^e	
			best	worst	failed	best	worst	best	worst	best	worst
0.4	0.1	3	0	1	0	0	0	0	0	12	0
0.4	0.1	5	15	34	0	1	1	3	2	46	0
0.4	0.1	7	18	49	0	0	12	6	10	75	0
0.4	0.1	9	10	54	0	0	13	3	17	87	0
0.4	0.2	3	0	0	0	0	0	0	0	11	0
0.4	0.2	5	6	31	0	0	1	0	1	63	0
0.4	0.2	7	9	59	0	0	5	3	12	88	0
0.4	0.2	9	8	49	0	0	21	6	13	86	0
0.6	0.1	3	0	1	0	0	0	0	0	18	0
0.6	0.1	5	7	32	0	0	6	2	1	67	0
0.6	0.1	7	12	44	0	0	20	8	21	80	0
0.6	0.1	9	9	43	0	0	27	10	18	81	0
0.6	0.2	3	0	0	0	0	0	0	0	29	0
0.6	0.2	5	11	38	0	1	0	2	4	65	0
0.6	0.2	7	15	54	0	0	14	10	12	74	1
0.6	0.2	9	0	99	0	0	1	77	0	23	0
0.8	0.1	3	0	3	0	0	0	0	0	17	0
0.8	0.1	5	13	35	0	0	3	0	3	56	3
0.8	0.1	7	12	62	3	0	13	14	11	72	0
0.8	0.1	9	26	34	0	0	29	7	27	67	0
0.8	0.2	3	0	0	0	0	0	0	0	23	0
0.8	0.2	5	8	26	0	0	1	2	1	55	0
0.8	0.2	7	4	71	1	0	27	89	1	7	1
0.8	0.2	9	0	100	94	0	0	98	0	2	0

^a "A best", "C best", "A worst", and "C worst" indicate number of times in which a particular approach (A or C) had rank 1 (best) or rank 4 (worst); ^b "A failed" indicates number of time in which numerical problems occurred under approach A, so the solution was not be found; ^c "B best" indicates number of times in which this approach had rank 1.5 provided that approach D had rank 1.5, too; ^d "B worst" indicates

1
2
3 number of times in which this approach had rank 4 or 3.5 provided that approach D also had rank 3.5; ^e “D
4 best” indicates number of times in which approach D had rank 1; ^f “D worst” indicates number of times in
5 which this approach had rank 3.5 provided that approach B had rank 1, too.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

Table 5. Summary of rank-ordering of efficiency of stratification approaches A, B, C, and D, for $N=15000$ and various combinations of σ, f, L .

σ	f	L	A ^a		A ^b	B ^c		C ^a		D ^e	
			best	worst	failed	best	worst	best	worst	best	worst
0.4	0.1	3	0	0	0	0	0	0	0	10	0
0.4	0.1	5	8	33	0	0	2	0	0	66	1
0.4	0.1	7	6	67	0	0	11	3	4	89	0
0.4	0.1	9	7	57	0	0	10	5	17	88	0
0.4	0.2	3	0	0	0	0	0	0	0	14	0
0.4	0.2	5	1	30	0	0	1	0	0	75	0
0.4	0.2	7	10	47	0	0	8	10	13	80	0
0.4	0.2	9	10	52	0	0	6	3	14	87	0
0.6	0.1	3	1	0	0	0	0	0	0	25	0
0.6	0.1	5	14	36	0	1	3	1	3	69	1
0.6	0.1	7	12	35	0	0	24	7	21	81	0
0.6	0.1	9	12	41	0	0	20	1	21	87	0
0.6	0.2	3	0	0	0	0	0	0	0	33	0
0.6	0.2	5	6	21	0	1	0	0	0	77	0
0.6	0.2	7	10	45	0	1	19	21	15	67	0
0.6	0.2	9	0	100	0	0	0	86	0	14	0
0.8	0.1	3	0	7	0	0	0	0	0	21	0
0.8	0.1	5	7	20	0	1	2	3	2	58	0
0.8	0.1	7	21	48	3	2	11	7	9	66	0
0.8	0.1	9	25	38	2	0	31	13	21	62	0
0.8	0.2	3	0	2	0	0	0	0	0	27	0
0.8	0.2	5	7	21	0	0	7	5	4	53	1
0.8	0.2	7	2	82	0	0	17	93	1	5	0
0.8	0.2	9	0	100	97	0	0	98	0	2	0

^a "A best", "C best", "A worst", and "C worst" indicate number of times in which a particular approach (A or C) had rank 1 (best) or rank 4 (worst); ^b "A failed" indicates number of time in which numerical problems occurred under approach A, so the solution was not be found; ^c "B best" indicates number of times in which this approach had rank 1.5 provided that approach D had rank 1.5, too; ^d "B worst" indicates

1
2
3 number of times in which this approach had rank 4 or 3.5 provided that approach D also had rank 3.5; ^e “D
4 best” indicates number of times in which approach D had rank 1; ^f “D worst” indicates number of times in
5 which this approach had rank 3.5 provided that approach B had rank 1, too.
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review Only

Table 6. Mean rank for four approaches compared, under different levels of the population quantities studied.

	Items	A	B	C	D
$N=1000$	2400	2.73	2.79	2.77	1.71
$N=2000$	2400	2.74	2.89	2.67	1.70
$N=5000$	2400	2.93	2.84	2.6	1.64
$N=10\ 000$	2400	2.98	2.84	2.54	1.64
$N=15\ 000$	2400	2.99	2.85	2.55	1.61
$\sigma=0.4$	3200	2.81	2.81	2.77	1.61
$\sigma=0.6$	3200	2.85	2.86	2.68	1.61
$\sigma=0.8$	3200	2.95	2.86	2.43	1.76
$L=3$	2400	2.58	2.59	2.59	2.25
$L=5$	2400	2.92	2.81	2.78	1.49
$L=7$	2400	2.94	2.99	2.65	1.42
$L=9$	2400	3.05	2.98	2.49	1.48
$f=0.1$	4800	2.73	2.84	2.81	1.63
$f=0.2$	4800	3.02	2.85	2.45	1.69
Totally	96 000	2.87	2.84	2.63	1.66

Table 7. Stratification of household income before taxes for 16 057 units from the 2001 Survey of Household Spending (SHS) Statistics Canada (source: Baillargeon and Rivest, 2007) for seven strata and sample size $n = 1500$. (Sums of stratum sample sizes are different from 1500 due to rounding to integers.)

Stratum	Approach A				Approach B				Approach C				Approach D	
	k^1	N_h	n_h	k	N_h	n_h	k	N_h	n_h	k	N_h	n_h		
1	0	3717	200	0	3717	178	0	2755	103	0	3717	202		
2	21707.55	3351	168	21862.38	2011	51	17076.72	2750	87	21017.47	3351	170		
3	37177.37	3094	178	30978.86	4434	311	29363.74	2358	76	37365.61	2930	160		
4	55383.14	2722	195	55897.71	2548	149	41727.19	3297	182	54053.91	3570	363		
5	78944.03	2079	228	76158.94	1626	103	62832.64	2595	183	86728.29	1886	281		
6	112557.23	1086	524	99599.61	1713	701	89082.52	2294	860	138399.45	595	315		
7	488693.26	8	8	496513.15	8	8	480472.03	8	8	496177.07	8	8		
CV		0.00374			0.00417			0.00436			0.00376			

¹ k is a stratum boundary; a particular h th stratum comprises units of which the stratification variable's values are within the interval $<k_h, k_{h+1}$)