

Optimization of sampling strata with the SamplingStrata package

Giulio Barcaroli

2018-10-07

Introduction

Let us suppose we need to design a sample survey, having a complete frame containing information on the target population (identifiers plus auxiliary information). If our sample design is a stratified one, we need to choose how to form strata in the population, in order to get the maximum advantage by the available auxiliary information. In other words, we have to decide in which way to combine the values of the auxiliary variables (from now on, the 'X' variables) in order to determine a new variable, called 'stratum'. To do so, we have to take into consideration the target variables of our sample survey (from now on, the 'Y' variables): if, to form strata, we choose the X variables most correlated to the Y's, the efficiency of the samples drawn by the resulting stratified frame may be greatly increased. In order to handle the whole auxiliary information in a homogenous way, we have to reduce continuous data to categorical (by mean of a k-means clustering technique, for example). Then, for every set of candidate auxiliary variables X's, we have to decide (i) what variables to consider as active variables in strata determination, and (ii) for each active variable, what set of values (in general, what aggregation of atomic values) have to be considered. Every combination of values of each active variable determine a particular stratification of the target population, i.e. a possible solution to the problem of 'best' stratification. Here, by best stratification, we mean the stratification that ensures the minimum sample cost, sufficient to satisfy a set of precision constraints, set on the accuracy of the estimates of the survey target variables Y's (constraints expressed as maximum allowable sampling variance on estimates in different domains of interest). When the cost of data collection is uniform over the strata, then the total cost is directly proportional to the overall sample size, and the convenience of a particular stratification can be measured by the associated size of the sample, whose estimates are expected to satisfy given accuracy levels. This minimum size can be determined by applying the Bethel algorithm, with its Chromy variant. In general, the number of possible alternative stratifications for a given population may be very high, depending on the number of variables and on the number of their values, and in these cases it is not possible to enumerate them in order to assess the best one. A very convenient solution to this, is the adoption of the evolutionary approach, consisting in applying a genetic algorithm that may converge towards a near-optimal solution after a finite number of iterations. The methodology is fully described in M. Ballin and Barcaroli (2013), and a complete illustration of the package, together with a comparison with the *stratification* package, is in Barcaroli (2014). Also a complete application in a case of network data is reported in Marco Ballin and Barcaroli (2016). The implementation of the genetic algorithm is based on a modification of the functions in the *genalg* package (see Willighagen 2005). In particular, the crossover operator has been modified on the basis of the indications given by O'Luing, Prestwich, and Tarim (2017).

Procedural steps

The optimization of the sampling design starts by making the sampling frame available, defining the target estimates of the survey and establishing the precision constraints on them. It is then possible to determine the best stratification and the optimal allocation. Finally, we proceed with the selection of the sample. Formalizing, these are the required steps:

- analysis of the frame data: identification of available auxiliary information;
- manipulation of auxiliary information: in case auxiliary variables are of the continuous type, they must be transformed into a categorical form;

- construction of atomic strata: on the basis of the categorical auxiliary variables available in the sampling frame, a set of strata can be constructed by calculating the Cartesian product of the values of all the auxiliary variables;
- characterization of each atomic stratum with the information related to the target variables: in order to optimise both strata and allocation of sampling units in strata, we need information on the distributions of the target variables (means and standard deviations) inside the different strata;
- choice of the precision constraints for each target estimate, possibly differentiated by domain;
- optimization of stratification and determination of required sample size and allocation in order to satisfy precision constraints on target estimates;
- analysis of the resulting optimized strata;
- association of new labels to sampling frame units, each of them indicating the new strata resulting by the optimal aggregation of the atomic strata;
- selection of units from the sampling frame with a selection scheme;
- evaluation of the found optimal solution in terms of expected precision and bias.

In the following, we will illustrate each step starting from a real sampling frame, the one that comes with the R package *sampling* (the dataframe *swissmunicipalities*).

Analysis of the frame data and manipulation of auxiliary information

As a first step, we have to define a frame dataframe containing the following information:

- a unique identifier of the unit (no restriction on the name, may be 'cod');
- the values of m auxiliary variables (named from X1 to Xm);
- the (optional) values of p target variables (named from Y1 to Yp);
- the value of the domain of interest for which we want to produce estimates (named 'domainvalue').

By typing the following statements in the R environment:

```
library(SamplingStrata)
#> Loading required package: memoise
#> Loading required package: doParallel
#> Loading required package: foreach
#> Loading required package: iterators
#> Loading required package: parallel
#> Loading required package: pbapply
require(memoise)
data(swissmunicipalities)
```

we get the *swissmunicipalities* dataframe, that contains 2896 observations (each observation refers to a Swiss municipality). Among the others, there are the following variables (data are referred to 2003):

- REG: Swiss region
- Nom: municipality name
- Surfacesbois: wood area
- Surfacescult: area under cultivation
- Alp: mountain pasture area
- Airbat: area with buildings
- Airind: industrial area
- Pop020: number of men and women aged between 0 and 19
- Pop2040: number of men and women aged between 20 and 39
- Pop4065: number of men and women aged between 40 and 64
- Pop65P: number of men and women aged between 65 and over
- POPTOT: total population

Let us suppose we want to plan a survey whose target estimates are the totals of population by age class in each Swiss region. In this case, our Y variables will be:

- Y1: number of men and women aged between 0 and 19
- Y2: number of men and women aged between 20 and 39
- Y3: number of men and women aged between 40 and 64
- Y4: number of men and women aged between 65 and over

As for the auxiliary variables (X's), we can use all of those characterising the area use (wood, mountain or pasture, cultivated, industrial, with buildings).

Finally, we want to produce estimates not only for the whole country, but also for each one of the seven different regions.

Function *buildFrameDF* permits to organize data in a suitable mode for next processing:

```
swissmunicipalities$id <- c(1:nrow(swissmunicipalities))
swissframe <- buildFrameDF(df = swissmunicipalities,
                           id = "id",
                           X = c("POPTOT",
                                "Surfacesbois",
                                "Surfacescult",
                                "Alp",
                                "Airbat",
                                "Airind"),
                           Y = c("Pop020",
                                "Pop2040",
                                "Pop4065",
                                "Pop65P"),
                           domainvalue = "REG")

str(swissframe)
#> 'data.frame':    2896 obs. of  12 variables:
#> $ id           : int  1 2 3 4 5 6 7 8 9 10 ...
#> $ X1           : int  363273 177964 166558 128634 124914 90483 72626 59496 48655 40377
#> ...
#> $ X2           : int  2326 67 97 1726 1635 2807 1139 408 976 425 ...
#> $ X3           : int  967 31 93 1041 714 1827 1222 183 196 694 ...
#> $ X4           : int  0 0 0 0 0 0 0 0 18 0 ...
#> $ X5           : int  2884 773 1023 1070 856 972 812 524 463 523 ...
#> $ X6           : int  260 60 213 212 64 238 134 27 108 137 ...
#> $ Y1           : int  57324 32429 28161 19399 24291 18942 14337 9533 9127 8128 ...
#> $ Y2           : int  131422 60074 50349 44263 44202 28958 24309 18843 14825 11265 ...
#> $ Y3           : int  108178 57063 53734 39397 35421 27696 21334 18177 15140 13301 ...
#> $ Y4           : int  66349 28398 34314 25575 21000 14887 12646 12943 9563 7683 ...
#> $ domainvalue: int  4 1 3 2 1 4 5 6 2 2 ...
```

As the X variables are of the continuous type, first we have to reduce them in a categorical (ordinal) form.

A suitable way to do so, is to apply a k-means clustering method (see Hartigan and Wong 1979) by using the function *var.bin*:

```
library(SamplingStrata)
swissframe$X1 <- var.bin(swissmunicipalities$POPTOT, bins=18)
swissframe$X2 <- var.bin(swissmunicipalities$Surfacesbois, bins=3)
swissframe$X3 <- var.bin(swissmunicipalities$Surfacescult, bins=3)
swissframe$X4 <- var.bin(swissmunicipalities$Alp, bins=3)
swissframe$X5 <- var.bin(swissmunicipalities$Airbat, bins=3)
swissframe$X6 <- var.bin(swissmunicipalities$Airind, bins=3)
```

Now, we have six different auxiliary variables of the categorical type, the first with 18 different modalities, the others with 3 modalities.

In any case, this dataframe comes with the package *SamplingStrata*: it can be made available by executing:

```
library(SamplingStrata)
data(swissframe)
head(swissframe)
#>   progr REG X1 X2 X3 X4 X5 X6      id   Y1   Y2   Y3   Y4
#> 1     1   4 18  3  2  1  3  3   Zurich 57324 131422 108178 66349
#> 2     2   1 17  1  1  1  3  2   Geneve 32429  60074  57063 28398
#> 3     3   3 17  1  1  1  3  3    Basel 28161  50349  53734 34314
#> 4     4   2 17  2  3  1  3  3    Bern 19399  44263  39397 25575
#> 5     5   1 17  2  2  1  3  2   Lausanne 24291  44202  35421 21000
#> 6     6   4 16  3  3  1  3  3 Winterthur 18942  28958  27696 14887
#>   domainvalue
#> 1           4
#> 2           1
#> 3           3
#> 4           2
#> 5           1
#> 6           4
```

We could also not indicate substantive X variables, if we want that each unit in the sampling frame be considered as an atomic stratum, and let to the optimization step to aggregate them on the basis of the values of the Y variable. In any case, as we have to indicate at least one X variable, we can use to this purpose the unique identifier in the frame:

```
swissmunicipalities$id <- c(1:nrow(swissframe))
newframe <- buildFrameDF(df = swissmunicipalities,
                        id = "id",
                        X = "id",
                        Y = c("Pop020",
                              "Pop2040",
                              "Pop4065",
                              "Pop65P"),
                        domainvalue = "REG")
str(newframe)
#> 'data.frame':   2896 obs. of  7 variables:
#> $ id           : int  1 2 3 4 5 6 7 8 9 10 ...
#> $ X1           : int  1 2 3 4 5 6 7 8 9 10 ...
#> $ Y1           : int  57324 32429 28161 19399 24291 18942 14337 9533 9127 8128 ...
#> $ Y2           : int  131422 60074 50349 44263 44202 28958 24309 18843 14825 11265 ...
#> $ Y3           : int  108178 57063 53734 39397 35421 27696 21334 18177 15140 13301 ...
#> $ Y4           : int  66349 28398 34314 25575 21000 14887 12646 12943 9563 7683 ...
#> $ domainvalue : int  4 1 3 2 1 4 5 6 2 2 ...
```

Construction of atomic strata and association of the information related to target variables

The *strata* dataframe reports information regarding each atomic stratum in the population. There is one row for each stratum. The total number of strata is given by the number of different combinations of X's values in the frame. For each stratum, the following information is required:

- the identifier of the stratum (named 'stratum' or 'strato'), concatenation of the values of the X variables;

- the values of the m auxiliary variables (named from X_1 to X_m) corresponding to those in the frame;
- the total number of units in the population (named 'N');
- a flag (named 'cens') indicating if the stratum is to be censused (=1) or sampled (=0);
- a variable indicating the cost of interviewing per unit in the stratum (named 'cost');
- for each target variable y , its mean and standard deviation, named respectively 'Mi' and 'Si');
- the value of the domain of interest to which the stratum belongs ('DOM1').

For example:

```
data(strata)
head(strata)
#>   stratum      N  X1  X2  X3      M1      M2      S1      S2 cens cost
#> 1      1 2246 x11 x21 x31 148.1598 443.0137 95.41435 202.4569    0    1
#> 2      2 2972 x11 x21 x32 184.2041 513.8995 81.26956 214.9999    0    1
#> 3      3 1905 x11 x22 x31 193.8927 488.8046 79.66667 261.1876    0    1
#> 4      4 3125 x11 x22 x32 181.3437 597.1925 82.77032 226.5086    0    1
#> 5      5 1733 x12 x21 x31 109.9850 418.2234 88.20289 179.1571    0    1
#> 6      6 1060 x12 x21 x32 114.7943 489.8292 52.71574 166.0292    0    1
#>   DOM1
#> 1 tot
#> 2 tot
#> 3 tot
#> 4 tot
#> 5 tot
#> 6 tot
```

If in the *frame* dataframe are also present the values of the target Y variables (from a census, or from administrative data), it is possible to automatically generate the *strata* dataframe by invoking the *buildStrataDF* function. Let us consider again the *swissframe* dataset that we have built in previous steps. To this dataset we can apply the function *buildStrataDF*:

```
swissstrata <- buildStrataDF(swissframe, progress = FALSE)
#>
#> Computations are being done on population data
#>
#> Number of strata: 641
#> ... of which with only one unit: 389
```

This is the structure of the created dataframe:

```
head(swissstrata)
#>   STRATO      N      M1      M2      M3      M4      S1      S2
#> 1 1*1*1*1*1*1 184 48.31522 49.40217 61.44022 28.40761 26.81536 28.49831
#> 2 1*1*1*1*1*2   1 98.00000 106.00000 116.00000 43.00000  0.00000  0.00000
#> 3 1*1*1*2*1*1   2 57.00000 64.00000 70.00000 50.00000  4.00000  0.00000
#> 4 1*1*2*1*1*1  11 77.72727 81.18182 92.36364 47.00000 15.24998 18.69768
#> 5 1*2*1*1*1*1   9 58.22222 61.55556 66.77778 36.22222 25.46360 20.27100
#> 6 1*2*1*2*1*1   8 61.00000 68.00000 84.62500 58.37500 24.56624 19.48076
#>   S3      S4 COST CENS DOM1 X1 X2 X3 X4 X5 X6
#> 1 32.63062 14.63922   1    0   1  1  1  1  1  1  1
#> 2  0.00000  0.00000   1    0   1  1  1  1  1  1  2
#> 3  1.00000 15.00000   1    0   1  1  1  1  2  1  1
#> 4 17.03084 11.12736   1    0   1  1  1  2  1  1  1
#> 5 24.89881 15.49751   1    0   1  1  2  1  1  1  1
#> 6 26.35307 26.55625   1    0   1  1  2  1  2  1  1
```

It is worth while noting that the total number of different atomic strata is 641, lower than the dimension of the Cartesian product of the X's (which is 4374): this is due to the fact that not all combinations of the value of the auxiliary variables are present in the sampling frame. Variables 'cost' and 'cens' are initialised respectively to 1 and 0 for all strata. It is possible to give them different values:

- for variable 'cost', it is possible to differentiate the cost of interviewing per unit by assigning real values;
- for variable 'cens', it is possible to set it equal to 1 for all strata that are of the 'take-all' type (i.e. all units in that strata must be selected).

The *swissstrata* dataframe comes together with *SamplingStrata* package, it can be made available by typing:

```
data(swissstrata)
```

On the contrary, if there is no information in the frame regarding the target variables, it is necessary to build the strata dataframe starting from other sources, for instance a previous round of the same survey, or from other surveys. In this case, we need to read sample data by executing:

```
samp <- read.delim("samplePrev.txt")
```

The only difference is that computed mean and variances of the Y's are sampling estimates, whose reliability should be evaluated by carefully considering their sampling variances. In addition to the naming constraints previously introduced, this case requires that a variable named 'WEIGHT' is present in the samp dataframe. Then we can execute this function in this way:

```
strata <- buildStrataDF(samp)
```

The result is much the same than in the previous case: the function creates a new dataframe, *strata*, and writes out in the working directory the strata file, named 'strata.txt'.

Choice of the precision constraints for each target estimate

The *errors* dataframe contains the accuracy constraints that are set on target estimates. This means to define a maximum coefficient of variation for each variable and for each domain value. Each row of this frame is related to accuracy constraints in a particular subdomain of interest, identified by the *domainvalue* value. In the case of the Swiss municipalities, we have chosen to define the following constraints:

```
data(swisserrors)
swisserrors
#>   DOM CV1 CV2 CV3 CV4 domainvalue
#> 1 DOM1 0.08 0.12 0.08 0.12      1
#> 2 DOM1 0.08 0.12 0.08 0.12      2
#> 3 DOM1 0.08 0.12 0.08 0.12      3
#> 4 DOM1 0.08 0.12 0.08 0.12      4
#> 5 DOM1 0.08 0.12 0.08 0.12      5
#> 6 DOM1 0.08 0.12 0.08 0.12      6
#> 7 DOM1 0.08 0.12 0.08 0.12      7
```

This example reports accuracy constraints on variables Y1, Y2, Y3 and Y4 that are the same for all the 7 different subdomains (Swiss regions) of domain level DOM1. Of course we can differentiate the precision constraints region by region. It is important to underline that the values of 'domainvalue' are the same than those in the *frame* dataframe, and correspond to the values of variable 'DOM1' in the strata

dataframe. Once having defined dataframes containing frame data, strata information and precision constraints, it is worth while to check their internal and reciprocal coherence. It is possible to do that by using the function *checkInput*:

```
checkInput(errors = checkInput(errors = swisserrors,
                                strata = swissstrata,
                                sampframe = swissframe))

#>
#> Input data have been checked and are compliant with requirements
#>
#> No input data indicated
```

For instance, this function controls that the number of auxiliary variables is the same in the *frame* and in the *strata* dataframes; that the number of target variables indicated in the *frame* dataframe is the same than the number of means and standard deviations in the *strata* dataframe, and the same than the number of coefficient of variations indicated in the *errors* dataframe.

If we try to determine the total size of the sample required to satisfy these precision constraints, considering the current stratification of the frame (the 641 atomic strata), we can do it by simply using the function *bethel*. This function requires a slightly different specification of the constraints dataframe:

```
cv <- swisserrors[1,]
cv
#>   DOM  CV1  CV2  CV3  CV4 domainvalue
#> 1 DOM1 0.08 0.12 0.08 0.12          1
```

because the *bethel* function does not permit to differentiate precision constraints by subdomain. In any case, the result of the application of the Bethel algorithm [see bethel:1989] is:

```
allocation <- bethel(swissstrata,cv)
sum(allocation)
#> [1] 893
```

That is, the required amount of units to be selected, with no optimization of sampling strata. In general, after the optimization, this number is sensibly reduced.

Optimization of frame stratification

Once the strata and the constraints dataframes have been prepared, it is possible to apply the function that optimises the stratification of the frame, that is *optimizeStrata*. This function operates on all subdomains, identifying the best solution for each one of them. The fundamental parameters to be passed to *optimizeStrata* are:

- *errors*: the (mandatory) dataframe containing the precision levels expressed in terms of maximum allowable coefficients of variation that regard the estimates on target variables of the survey;
- *strata*: the (mandatory) dataframe containing the information related to 'atomic' strata, i.e. the strata obtained by the Cartesian product of all auxiliary variables *X*'s. Information concerns the identifiability of strata (values of *X*'s) and variability of *Y*'s (for each *Y*, mean and standard deviation in strata);
- *cens*: the (optional) dataframe containing the 'take-all' strata, those strata whose units must be selected in whatever sample. It has same structure than **strata* dataframe;
- *strcens*: flag (TRUE/FALSE) to indicate if 'take-all' strata do exist or not. Default is FALSE;
- *initialStrata*: the initial limit on the number of strata for each solution. Default is NA, and in this case it is set equal to the number of atomic strata in each domain. If the parameter *addStrataFactor* is

equal to zero, then *initialStrata* is equivalent to the maximum number of strata to be obtained in the final solution;

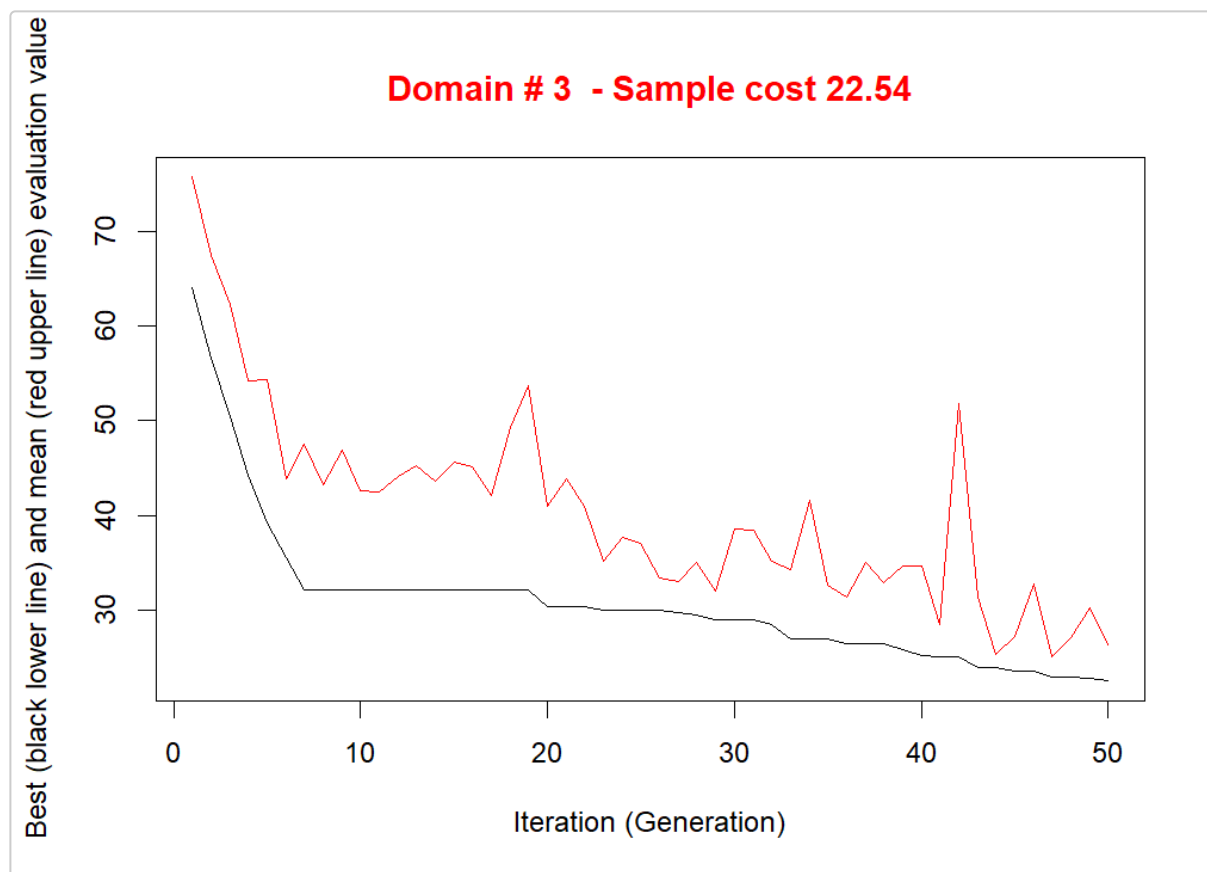
- *addStrataFactor*: indicates the probability that at each mutation the number of strata may increase with respect to the current value. Default is 0.0;
- *minnumstr*: indicates the minimum number of units that must be allocated in each stratum. Default is 2;
- *iter*: indicates the maximum number of iterations (= generations) of the genetic algorithm. Default is 50
- *pops*: dimension of each generations in terms of number of individuals to be generated. Default is 20;
- *mut_chance* (mutation chance): for each new individual, the probability to change each single chromosome, i.e. one bit of the solution vector. High values of this parameter allow a deeper exploration of the solution space, but a slower convergence, while low values permit a faster convergence, but the final solution can be distant from the optimal one. Default is NA, in correspondence of which it is computed as $1/(\text{vars}+1)$ where vars is the length of elements in the solution;
- *elitism_rate*: indicates the rate of better solutions that must be preserved from one generation to another. Default is 0.2;
- *suggestions*: indicates one possible solution (from kmeans clustering or from previous runs) that will be introduced in the initial population. Default is NULL;
- *realAllocation* : if FALSE, the allocation is based on INTEGER values; if TRUE, the allocation is based on REAL values. Default is TRUE;
- *writeFiles* : indicates if the various dataframes and plots produced during the execution have to be written in the working directory /output. Default is FALSE;
- *showPlot* : indicates if the plot showing the trend in the value of the objective function has to be shown or not. In parallel = TRUE, this defaults to FALSE, otherwise default is TRUE.
- *parallel* : Should the analysis be run in parallel. Default is TRUE
- *cores* : if the analysis is run in parallel, how many cores should be used. If not specified n-1 of total available cores are used OR if number of domains < (n-1) cores, then number of cores equal to number of domains are used

In the case of the Swiss municipalities, we can use almost all of default values for parameters with the exception of the errors and strata dataframes, and for the option 'writeFiles':

```
solution1 <- optimizeStrata(  
  errors = swisserrors,  
  strata = swissstrata,  
  parallel = TRUE,  
  writeFiles = FALSE,  
  showPlot = FALSE)  
  
sum(solution1$aggr_strata$SOLUZ)  
#> [1] 192.9996
```

Note that by so doing the *initialStrata* parameter is set equal to the number of atomic strata in each domain . Another possibility is to set a pre-determined value for each domain, for instance equal in each domain, as c(5,5,5,5,5,5,5).

The execution of *optimizeStrata* produces the solution of 7 different optimization problems, one for each domain.



The graphs (we report here the one related to domain 3) illustrate the convergence of the solution to the final one starting from the initial one (i.e. the one related to the atomic strata). Along the x-axis are reported the executed iterations, from 1 to the maximum, while on the y-axis are reported the size of the sample required to satisfy precision constraints. The upper (red) line represent the average sample size for each iteration, while the lower (black) line represents the best solution found until the i-th iteration.

The results of the execution are contained in the list 'solution', composed by two elements:

1. *solution1\$indices*: the vector of the indices that indicates to what aggregated stratum each atomic stratum belongs;
2. *solution1\$aggr_strata*: the dataframe containing information on the optimal aggregated strata.

Initial solution with kmeans clustering of atomic strata

In order to speed up the convergence towards the optimal solution, an initial one can be given as a "suggestion" to 'optimizeStrata' function. The function *KmeansSolution* produces this initial solution by clustering atomic strata considering the values of the means of all the target variables Y.

Also, the optimal number of clusters is determined inside each domain. If the default value for *nstrata* is used, then the number of aggregate strata is optimized by varying the number of cluster from 2 to number of atomic strata in each domain, divided by 2. Otherwise, it is possible to indicate a fixed number of aggregate strata to be obtained.

Other parameters are:

- *minnumstrat*: the minimum number of units to be allocated in each stratum(default is 2);
- *maxcluster*: the maximum number of clusters to be considered in the execution of kmeans algorithm;
- *showPlot*: if TRUE, allows to visualise the optimization.

For any given number of clusters, the correspondent aggregation of atomic strata is considered as input to the function 'bethel'. The number of clusters for which the value of the sample size necessary to fulfil

precision constraints is the minimum one, is retained as the optimal one.

The overall solution is obtained by concatenating optimal clusters obtained in domains. The result is a dataframe with two columns: the first indicates the clusters, the second the domains:

```
solutionKmeans1 <- KmeansSolution(swisstrata,
                                  swisserrors,
                                  nstrata=NA,
                                  minnumstrat=2,
                                  maxclusters=NA,
                                  showPlot=FALSE)

#>
#> -----
#> Kmeans solution
#> -----
#> *** Domain: 1 ***
#> Number of strata: 8
#> Sample size      : 18
#> *** Domain: 2 ***
#> Number of strata: 9
#> Sample size      : 17
#> *** Domain: 3 ***
#> Number of strata: 7
#> Sample size      : 13
#> *** Domain: 4 ***
#> Number of strata: 6
#> Sample size      : 10
#> *** Domain: 5 ***
#> Number of strata: 8
#> Sample size      : 15
#> *** Domain: 6 ***
#> Number of strata: 8
#> Sample size      : 15
#> *** Domain: 7 ***
#> Number of strata: 9
#> Sample size      : 17
head(solutionKmeans1)
#> suggestions domainvalue
#> 1          2          1
#> 2          2          1
#> 3          2          1
#> 4          2          1
#> 5          2          1
#> 6          2          1
```

This solution can be given as argument to the parameter *suggestion* in the *optimizeStrata* function:

```
solution_with_kmeans <- optimizeStrata(
  errors = swisserrors,
  strata = swisstrata,
  suggestions = solutionKmeans1,
  parallel = TRUE,
  writeFiles = FALSE,
  showPlot = FALSE)
```

```
sum(ceiling(solution_with_kmeans$aggr_strata$SOLUZ))
#> [1] 97
```

thus obtaining a much more convenient solution than the one without the kmeans suggestion, with the same number of iterations.

Adjustment of the final sampling size

After the optimization step, the final sample size is the result of the allocation of units in final strata. This allocation is such that the precision constraints are expected to be satisfied. Actually, three possible situations may occur:

- the resulting sample size is acceptable;
- the resulting sample size is too high, it is not affordable with respect to the available budget;
- the resulting sample size is too low, the available budget permits to increase the number of units.

In the first case, no action is required. In the second case, it is necessary to reduce the number of units, by equally applying the same reduction rate in each stratum. In the third case, we could either to set more tight precision constraints, or proceed to increase the sample size by applying the same increase rate in each stratum. This increase/reduction process is iterative, as by applying the same rate we could find that in some strata there are not enough units to increase or to reduce. The function *adjustSize* permits to obtain the desired final sample size. Let us suppose that the obtained sample size is not affordable. We can reduce it by executing the following code:

```
adjustedStrata <- adjustSize(size=150, strata=solution1$aggr_strata, cens=NULL)
#>
#> 162
#> 162
#> Final adjusted size: 162
sum(adjustedStrata$SOLUZ)
#> [1] 162
```

Instead, if we want to increase the size because the budget allows to do this, then this is the code:

```
adjustedStrata <- adjustSize(size=250, strata=solution1$aggr_strata, cens=NULL)
#>
#> 257
#> 257
#> Final adjusted size: 257
sum(adjustedStrata$SOLUZ)
#> [1] 257
```

The difference between the desired sample size and the actual adjusted size depends on the number of strata in the optimized solution. Consider that the adjustment is performed in each stratum by taking into account the relative difference between the current sample size and the desired one: this produces an allocation that is expressed by a real number, that must be rounded. The higher the number of strata, the higher the impact of the rounding in all strata on the final adjusted sample size.

Analysis of results

This function has two purposes:

1. instrumental to the processing of the sampling frame (attribution of the labels of the optimized strata to the population units);

2. analysis of the aggregation of the atomic strata obtained in the optimized solution.

The function *updateStrata* assigns the labels of the new strata to the initial one in the dataframe *strata*, and produces:

- a new dataframe named *newstrata* containing all the information in the *strata* dataframe, plus the labels of the new optimized strata;
- a table, contained in the dataset *strata_aggregation.txt*, showing in which way each optimized stratum aggregates the auxiliary variables X's.

The function is invoked in this way:

```
newstrata <- updateStrata(swissstrata,
                        solution1,
                        writeFiles = TRUE)
```

Now, the atomic strata are associated to the aggregate strata defined in the optimal solution, by means of the variable *LABEL*. If we want to analyse in detail the new structure of the stratification, we can look at the *strata_aggregation.txt* file:

```
strata_aggregation <- read.delim("strata_aggregation.txt")
head(strata_aggregation)
#>   DOM1 AGGR_STRATUM X1 X2 X3 X4 X5 X6
#> 1     1             1 1 1 1 1 1 1
#> 2     1             1 3 2 2 2 1 1
#> 3     1             2 1 1 1 1 1 2
#> 4     1             2 1 2 1 2 1 1
#> 5     1             2 2 1 1 1 1 1
#> 6     1             2 2 1 1 2 1 1
```

In this structure, for each aggregate stratum the values of the X's variables in each contributing atomic stratum are reported. It is then possible to understand the meaning of each aggregate stratum produced by the optimization.

Updating the frame and selecting the sample

Once the optimal stratification has been obtained, to be operational we need to accomplish the following two steps:

- to update the frame units with new stratum labels (combination of the new values of the auxiliary variables X's);
- to select the sample from the frame.

As for the first, we execute the following command:

```
framenew <- updateFrame(swissframe, newstrata, writeFiles=FALSE)
```

The function *updateFrame* receives as arguments the indication of the dataframe in which the frame information is memorised, and of the dataframe produced by the execution of the *updateStrata* function. The execution of this function produces a dataframe *framenew*, and also a file (named *framenew.txt*) with the labels of the new strata produced by the optimisation step. The allocation of units is contained in the *SOLUZ* column of the dataset *outstrata.txt*. At this point it is possible to select the sample from the new version of the frame:

```

sample <- selectSample(framenew, solution1$aggr_strata, writeFiles=FALSE)
#>
#> *** Sample has been drawn successfully ***
#> 191 units have been selected from 71 strata
#>
#> ==> There have been 10 take-all strata
#> from which have been selected 20 units

```

that produces two .csv files:

- *sample.csv* containing the units of the frame that have been selected, together with the weight that has been calculated for each one of them;
- *sample.chk.csv* containing information on the selection: for each stratum, the number of units in the population, the planned sample, the number of selected units, the sum of their weights that must equalise the number of units in the population.

The *selectSample* operates by drawing a simple random sampling in each stratum.

A variant of this function is *selectSampleSystematic*. The only difference is in the method used for selecting units in each strata, that is by executing the following steps:

1. a selection interval is determined by considering the inverse of the sampling rate in the stratum;
2. a starting point is determined by selecting a value in this interval;
3. the selection proceeds by selecting as first unit the one corresponding to the above value, and then selecting all the units individuated by adding the selection interval.

This selection method can be useful if associated to a particular ordering of the selection frame, where the ordering variable(s) can be considered as additional stratum variable(s). For instance, we could decide that it could be important to consider the overall population in municipalities when selecting units. Here is the code:

```

# adding POPTOT to framenew
data("swissmunicipalities")
framenew <- merge(framenew,swissmunicipalities[,c("REG","Nom","POPTOT")],
                  by.x=c("REG","ID"),by.y=c("REG","Nom"))
# selection of sample with systematic method
sample <- selectSampleSystematic(frame=framenew,
                                outstrata=solution1$aggr_strata,
                                sortvariable = c("POPTOT"))

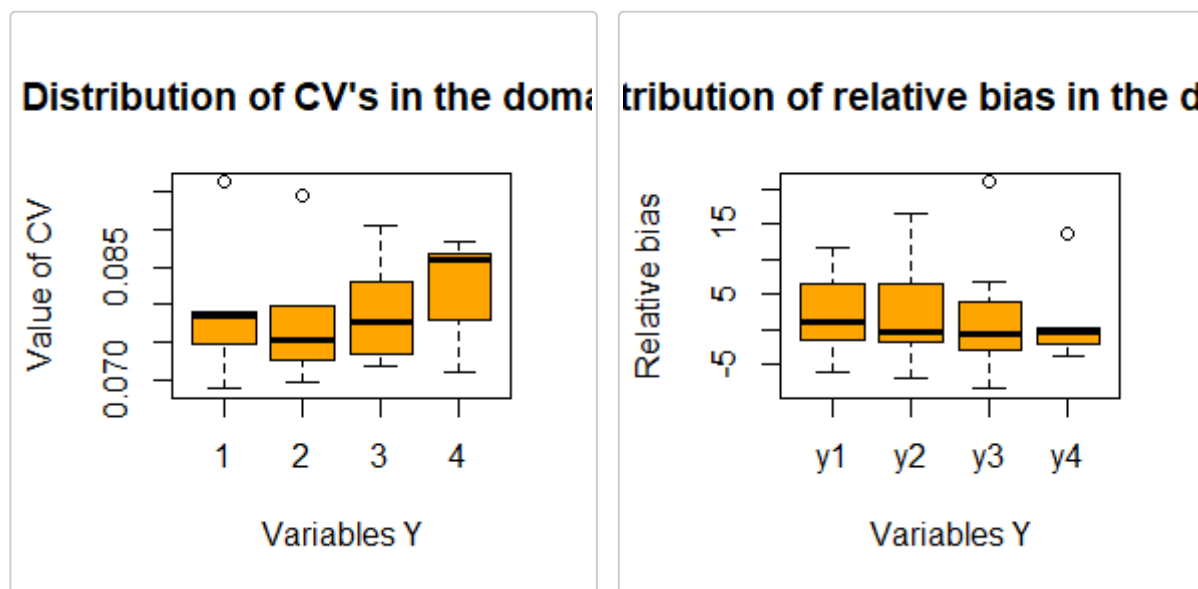
#>
#> *** Sample has been drawn successfully ***
#> 191 units have been selected from 71 strata
#>
#> ==> There have been 10 take-all strata
#> from which have been selected 20 units
head(sample,3)
#>   DOMAINVALUE STRATO REG          ID   STRATUM  PROGR  X1  X2  X3
#> 1           1      1  1 Chesalles-sur-Moudon 1*1*1*1*1  2593  1  1  1
#> 2           1      1  1 Villars-Tiercelin 1*1*1*1*1  2176  1  1  1
#> 3           1     10  1 Concise 2*2*1*1*1*1  1630  2  2  1
#>   X4 X5 X6  Y1  Y2  Y3  Y4 LABEL POPTOT  WEIGHTS      FPC
#> 1  1  1  1  54  40  45  20    1   159 92.50000 0.01081081
#> 2  1  1  1  76  92 136  39    1   343 92.50000 0.01081081
#> 3  1  1  1 155 195 231 103   10   684 25.33333 0.03947368

```

Evaluation of the found solution

In order to be confident about the quality of the found solution, the function `evalSolution` allows to run a simulation, based on the selection of a desired number of samples from the frame to which the stratification, identified as the best, has been applied. The user can invoke this function also indicating the number of samples to be drawn:

```
eval <- evalSolution(framenew,
                     solution1$aggr_strata,
                     nsampl=100,
                     writeFiles=FALSE,
                     progress=FALSE)
```



For each drawn sample, the estimates related to the Y's are calculated. Their mean and standard deviation are also computed, in order to produce the CV related to each variable in every domain. These CV's can be inspected and compared to the constraints:

```
eval$coeff_var
#>      CV1      CV2      CV3      CV4 dom
#> 1 0.07419901 0.07521375 0.07226064 0.08668406 DOM1
#> 2 0.07900759 0.07977563 0.08246782 0.08605109 DOM2
#> 3 0.06870326 0.06975793 0.07778861 0.08843159 DOM3
#> 4 0.07843063 0.07366349 0.07433401 0.07095334 DOM4
#> 5 0.09630069 0.09452957 0.09051767 0.08707496 DOM5
#> 6 0.07873707 0.07964042 0.08364721 0.07916796 DOM6
#> 7 0.07502415 0.07132590 0.07182810 0.07672917 DOM7
swisserrors
#>   DOM  CV1  CV2  CV3  CV4 domainvalue
#> 1 DOM1 0.08 0.12 0.08 0.12          1
#> 2 DOM1 0.08 0.12 0.08 0.12          2
#> 3 DOM1 0.08 0.12 0.08 0.12          3
#> 4 DOM1 0.08 0.12 0.08 0.12          4
#> 5 DOM1 0.08 0.12 0.08 0.12          5
#> 6 DOM1 0.08 0.12 0.08 0.12          6
#> 7 DOM1 0.08 0.12 0.08 0.12          7
```

Handling 'take-all' strata in the optimization step

As input to the optimization step, together with proper sampling strata, it is also possible to provide *take-all* strata. These strata will not be subject to optimisation as the proper strata, but they will contribute to the determination of the best stratification, as their presence in a given domain will permit to satisfy precision constraint with a lower number of units belonging to proper sampling strata.

In order to correctly execute the optimization and further steps, it is necessary to perform a pre-processing of the overall input. The first step to be executed consists in the bi-partition of units to be censused and of units to be sampled, in order to build two different frames:

```
data(swisserrors)
data(swisstrata)
data(swissframe)
#---Selection of units to be censused from the frame
framecens <- swissframe[ (swissframe$domainvalue == 1 |
                        swissframe$domainvalue == 4) &
                        (swissframe$X2 == 1 &
                        swissframe$X3 == 1 &
                        swissframe$X4 == 1 &
                        swissframe$X5 == 1 &
                        swissframe$X6 == 1) , ]
#---Selection of units to be sampled from the frame
# (complement to the previous)
framesamp <- swissframe[!((swissframe$domainvalue == 1 |
                        swissframe$domainvalue == 4) &
                        (swissframe$X2 == 1 &
                        swissframe$X3 == 1 &
                        swissframe$X4 == 1 &
                        swissframe$X5 == 1 &
                        swissframe$X6 == 1)) , ]
```

In this way, we have included all units belonging to ‘take-all’ strata in ‘framecens’, and the remaining in ‘framesamp’. At the end of the process, the sample will be selected from ‘framesamp’, while the units in ‘framecens’ will be simply added to the sample.

We can obtain census strata and sampling strata by applying *buildStrataDF* respectively to *framecens* and *framesamp*:

```
# Build strata to be censused and sampled
cens <- buildStrataDF(framecens,progress = FALSE)
#>
#> Computations are being done on population data
#>
#> Number of strata: 18
#> ... of which with only one unit: 2
sum(cens$N)
#> [1] 405
strata <- buildStrataDF(framesamp,progress = FALSE)
#>
#> Computations are being done on population data
#>
#> Number of strata: 623
#> ... of which with only one unit: 387
sum(strata$N)
#> [1] 2491
```

and

```
sum(cens$N)
#> [1] 405
sum(strata$N)
#> [1] 2491
```

Now we have all required inputs to run 'optimizeStrata' in presence of the 'take-all' strata:

```
solution2 <- optimizeStrata(
  errors = swisserrors,
  strata = strata,
  cens = cens,
  strcens = TRUE,
  parallel = TRUE,
  writeFiles = TRUE,
  showPlot = FALSE
)
```

Once the optimized solution has been produced, the next steps are executed by considering only the sampling part of the frame:

```
newstrata <- updateStrata(strata, solution2)
# updating sampling frame with new strata labels
framenew <- updateFrame(frame=framesamp, newstrata=newstrata)
# selection of sample from sampling strata
sample <- selectSample(frame=framenew, outstrata=solution2$aggr_strata)
#>
#> *** Sample has been drawn successfully ***
#> 209 units have been selected from 75 strata
#>
#> ==> There have been 9 take-all strata
#> from which have been selected 16 units
```

Finally, the units in the 'take-all' strata can be added to sampled ones. First, the census frame needs to be made homogeneous to the sample frame in order to permit the 'rbind' step:

```
# addition of necessary variables to
colnames(framesamp) <- toupper(colnames(framesamp))
colnames(framecens) <- toupper(colnames(framecens))
framecens$WEIGHTS <- rep(1, nrow(framecens))
framecens$FPC <- rep(1, nrow(framecens))
framecens$LABEL <- rep("999999", nrow(framecens))
framecens$STRATUM <- rep("999999", nrow(framecens))
framecens$STRATO <- rep("999999", nrow(framecens))
```

The overall set of units to be surveyed is obtainable in this way:

```
survey <- rbind(sample, framecens)
```

and this is the proportion of sampling and censused units:

```
survey$cens <- ifelse(survey$LABEL == "999999", 1, 0)
table(survey$cens)
#>
```



```
#> 0 1
#> 209 405
```

Handling Anticipated Variance

In the previous sections it has been assumed that, when optimizing the stratification of a sampling frame, values of the target variables Y 's are available for the generality of the units in the frame, or at least for a sample of them by means of which it is possible to estimate means and standard deviation of Y 's in atomic strata. Of course, this assumption is seldom expected to hold. The situation in which some proxy variables are available in the frame is much more likely to happen. In these situations, instead of directly indicating the real target variables, proxy ones are named as Y 's. By so doing, there is no guarantee that the final stratification and allocation can ensure the compliance to the set of precision constraints. In order to take into account this problem, and to limit the risk of overestimating the expected precision levels of the optimized solution, it is possible to carry out the optimization by considering, instead of the expected coefficients of variation related to proxy variables, the anticipated coefficients of variation (ACV) that depend on the model that is possible to fit on couples of real target variables and proxy ones. In the current implementation, only models linking continuous variables can be considered. The definition and the use of these models is the same that has been implemented in the package *stratification* (see Baillargeon and Rivest 2014). In particular, the reference here is to two different models, the linear model with heteroscedasticity:

$$Y = \beta \times X + \epsilon$$

where

$$\epsilon \sim N(0, \sigma^2 X^\gamma)$$

(in case $\gamma = 0$, then the model is homoscedastic)

and the loglinear model:

$$Y = \exp(\beta \times \log(X) + \epsilon)$$

where

$$\epsilon \sim N(0, \sigma^2)$$

In order to make evident the importance of the above, consider the following example, based on the dataset *nations* available in the package.

```
data(nations)
head(nations)
#>      Country TFR contraception infant.mortality GDP region
#> 1  Afghanistan 6.90             63             154 2848 Asia
#> 2   Albania 2.60             47             32  863 Europe
#> 3   Algeria 3.81             52             44 1531 Africa
#> 4 American-Samoa 1.35             71             11 2433 Oceania
#> 5   Andorra 1.61             71              7 19121 Europe
#> 6   Angola 6.69             19             124  355 Africa
#> Continent
#> 1      2
#> 2      1
#> 3      4
#> 4      5
#> 5      1
#> 6      4
```

Let us assume that in the sampling frame only variable *GDP* (Gross Domestic Product) is available for all countries, while *contraception rates* and *infant mortality rates* are available only on a subset of countries (about one third).

```
set.seed(1234)
nations_sample <- nations[sample(c(1:207),70),]
```

In this subset we can fit models between GDP and the two variables that we assume are the target of our survey.

One model for *infant mortality* and *GDP*:

```
mod_logGDP_INF Mort <- lm(log(nations_sample$infant.mortality) ~ log(nations_sample$GDP))
summary(mod_logGDP_INF Mort)
#>
#> Call:
#> lm(formula = log(nations_sample$infant.mortality) ~ log(nations_sample$GDP))
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -1.1292 -0.3765 -0.1455  0.3316  2.6345
#>
#> Coefficients:
#>
#>              Estimate Std. Error t value      Pr(>|t|)
#> (Intercept)      6.86295     0.33620    20.41 < 0.0000000000000002
#> log(nations_sample$GDP) -0.46580     0.04389   -10.61 0.00000000000000452
#>
#> (Intercept)          ***
#> log(nations_sample$GDP) ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 0.6158 on 68 degrees of freedom
#> Multiple R-squared:  0.6236, Adjusted R-squared:  0.6181
#> F-statistic: 112.7 on 1 and 68 DF,  p-value: 0.000000000000004523
```

and one model for *contraception* and *GDP*:

```
mod_logGDP_CONTRA <- lm(log(nations_sample$contraception) ~ log(nations_sample$GDP))
summary(mod_logGDP_CONTRA)
#>
#> Call:
#> lm(formula = log(nations_sample$contraception) ~ log(nations_sample$GDP))
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -1.96139 -0.27360 -0.01435  0.45058  1.25143
#>
#> Coefficients:
#>
#>              Estimate Std. Error t value      Pr(>|t|)
#> (Intercept)      0.98318     0.30538    3.220     0.00197 **
#> log(nations_sample$GDP) 0.34649     0.03986    8.692 0.000000000000122 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 0.5593 on 68 degrees of freedom
```

```
#> Multiple R-squared: 0.5263, Adjusted R-squared: 0.5193
#> F-statistic: 75.55 on 1 and 68 DF, p-value: 0.000000000001217
```

We define the sampling frame in this way:

```
nations$progr <- c(1:nrow(nations))
nations$dom <- 1
frame <- buildFrameDF(nations,
                      id="Country",
                      X="progr",
                      Y=c("GDP", "GDP"),
                      domainvalue = "dom")
```

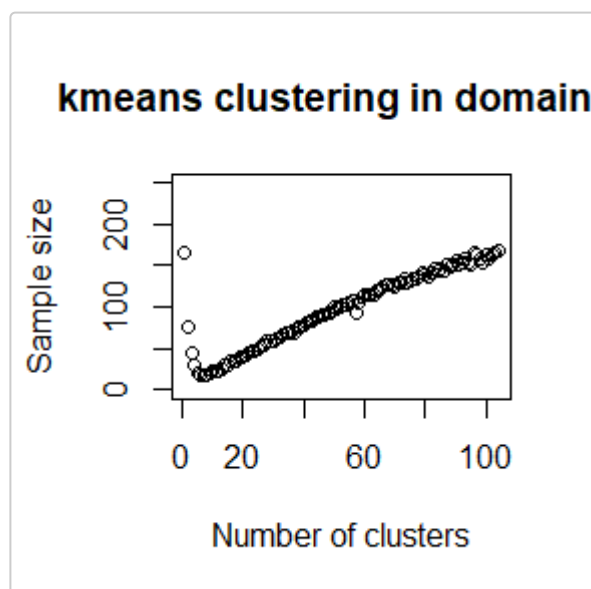
that is, we replicate twice the variable *GDP* because it will be used once for *infant mortality* and once for *contraception*.

We set 10% and 5% precision constraints on these two variables:

```
cv <- as.data.frame(list(DOM=rep("DOM1",1),
                        CV1=rep(0.10,1),
                        CV2=rep(0.05,1),
                        domainvalue=c(1:1)
                        ))
cv
#>   DOM CV1 CV2 domainvalue
#> 1 DOM1 0.1 0.05          1
```

We build the strata without any assumption on the variability of the two target variables, and proceed in the optimization:

```
strata1 <- buildStrataDF(frame, progress = FALSE)
#>
#> Computations are being done on population data
#>
#> Number of strata: 207
#> ... of which with only one unit: 207
solutionKmeans2 <- KmeansSolution(strata1,cv)
```

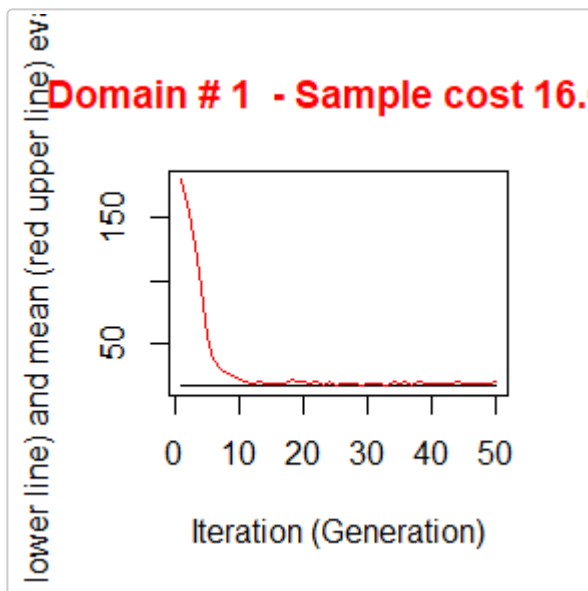


```

#>
#> -----
#> Kmeans solution
#> -----
#> *** Domain: 1 ***
#> Number of strata: 7
#> Sample size      : 17
solution3 <- optimizeStrata(cv,
                             strata1,
                             iter = 50,
                             pops = 20,
                             parallel = FALSE,
                             suggestions = solutionKmeans2,
                             writeFiles = FALSE,
                             showPlot = FALSE)

#>
#> *** Domain : 1 1
#> Number of strata : 207
#> -----
#> Optimal stratification with Genetic Algorithm
#> -----
#> *** Parameters ***
#> -----
#> Domain: 1
#> Maximum number of strata: 207
#> Minimum number of units per stratum: 2
#> Take-all strata (TRUE/FALSE): FALSE
#> number of sampling strata : 207
#> Number of target variables: 2
#> Number of domains: 1
#> Number of GA iterations: 50
#> Dimension of GA population: 20
#> Mutation chance in GA generation: NA
#> Elitism rate in GA generation: 0.2
#> Chance to add strata to maximum: 0
#> Allocation with real numbers instead of integers: TRUE
#> Suggestion: 5 3 3 6 6 7 5 6 7 4 2 7 2 2 2 7 6 6 5 7 6 1 7 6 1 6 5 5 7 6 2 6 7 6 1 5 7 6
2 1 3 3 6 6 6 1 2 2 4 6 4 5 7 7 5 7 6 5 1 1 1 1 7 5 6 4 5 4 7 3 6 4 6 4 6 2 5 1 7 6 6 5 1 6
5 6 7 7 2 2 2 5 6 6 5 6 7 5 7 5 5 6 3 6 6 3 3 2 4 6 7 6 6 5 6 5 5 6 7 6 6 6 6 6 5 5 3 7 6 6
7 6 6 3 7 6 6 4 6 7 6 5 7 4 5 5 7 1 5 6 2 6 3 5 7 3 7 7 7 6 6 5 6 6 5 3 2 3 3 4 6 7 6 2 4 6
1 5 1 3 7 6 6 6 6 4 6 3 5 2 6 7 7 1 3 3 6 3 6 7 2 7 7 6 6 6 1
#> *** Sample cost: 16.62586
#> *** Number of strata: 6

```



```
#>
#> *** Sample size : 16
#> *** Number of strata : 6
#> -----

sum(solution3$aggr_strata$SOLUZ)
#> [1] 16.62586
```

Then, we evaluate the expected CV's on the three variables:

```
newstrata <- updateStrata(strata1,solution3)
framenew1 <- updateFrame(frame,newstrata)
framenew1 <- framenew1[order(framenew1$ID),]
framenew1$Y2 <- nations$infant.mortality
framenew1$Y3 <- nations$contraception
results1 <- evalSolution(framenew1, solution3$aggr_strata, 50, progress = FALSE)
results1$coeff_var
#>          CV1          CV2          CV3  dom
#> 1 0.05332022 0.2326833 0.1431669 DOM1
```

Clearly, the CV's on *infant mortality* and *contraception* are not compliant with the corresponding precision constraints.

We now proceed in building the *strata* dataframe using the models:

```
model <- NULL
model$beta[1] <- mod_logGDP_INFMORT$coefficients[2]
model$sig2[1] <- summary(mod_logGDP_INFMORT)$sigma
model$type[1] <- "loglinear"
model$gamma[1] <- 0
model$beta[2] <- mod_logGDP_CONTRA$coefficients[2]
model$sig2[2] <- summary(mod_logGDP_CONTRA)$sigma
model$type[2] <- "loglinear"
model$gamma[2] <- 0
model <- as.data.frame(model)
model
#>          beta          sig2          type gamma
```

```
#> 1 -0.4658038 0.6157600 Loglinear 0
#> 2 0.3464857 0.5593031 Loglinear 0
```

```
strata2 <- buildStrataDF(frame, model = model, progress = FALSE)
```

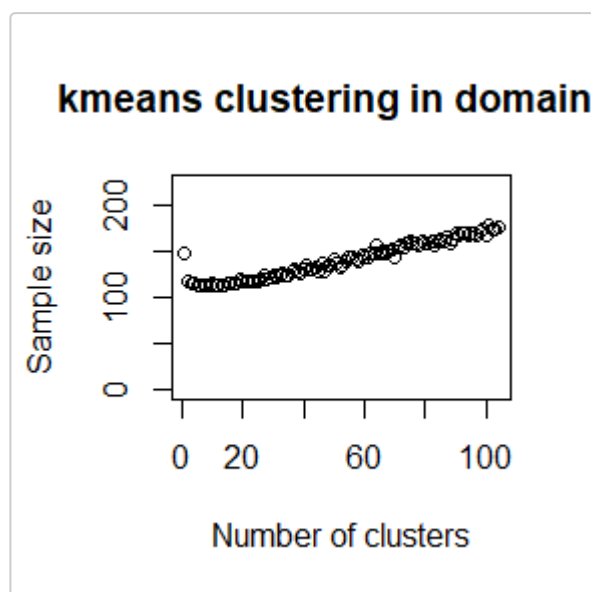
```
#>
#> Computations are being done on population data
#>
#> Number of strata: 207
#> ... of which with only one unit: 207
```

```
head(strata2)
```

```
#>   STRATO N      M1      M2      S1      S2 COST CENS DOM1 X1
#> 1     1 1 0.024595860 15.737965 0.022690436 13.624506    1    0    1    1
#> 2    10 1 0.009910578 30.945032 0.009142812 26.789408    1    0    1   10
#> 3   100 1 0.011086660 28.468488 0.010227784 24.645441    1    0    1  100
#> 4   101 1 0.067027720  7.465936 0.061835128  6.463332    1    0    1  101
#> 5   102 1 0.064539738  7.678981 0.059539888  6.647767    1    0    1  102
#> 6   103 1 0.030744575 13.331076 0.028362814 11.540840    1    0    1  103
```

We proceed with the optimization

```
solutionKmeans3 = KmeansSolution(strata2,cv)
```

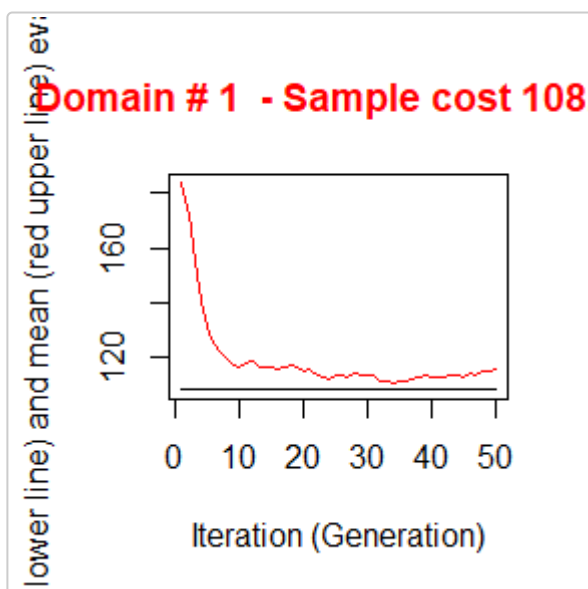


```
#>
#> -----
#> Kmeans solution
#> -----
#> *** Domain: 1 ***
#> Number of strata: 12
#> Sample size      : 113
strata2 <- buildStrataDF(frame, model = model, progress = FALSE)
#>
#> Computations are being done on population data
#>
#> Number of strata: 207
#> ... of which with only one unit: 207
solution4 <-
  optimizeStrata(
    errors = cv ,
```

```

strata = strata2,
iter = 50,
pops = 20,
parallel = FALSE,
suggestions = solutionKmeans3,
showPlot = FALSE,
writeFiles = FALSE)
#>
#> *** Domain : 1 1
#> Number of strata : 207
#> -----
#> Optimal stratification with Genetic Algorithm
#> -----
#> *** Parameters ***
#> -----
#> Domain: 1
#> Maximum number of strata: 207
#> Minimum number of units per stratum: 2
#> Take-all strata (TRUE/FALSE): FALSE
#> number of sampling strata : 207
#> Number of target variables: 2
#> Number of domains: 1
#> Number of GA iterations: 50
#> Dimension of GA population: 20
#> Mutation chance in GA generation: NA
#> Elitism rate in GA generation: 0.2
#> Chance to add strata to maximum: 0
#> Allocation with real numbers instead of integers: TRUE
#> Suggestion: 12 3 3 4 4 1 12 4 9 6 8 1 8 8 8 1 7 11 6 9 7 5 1 4 5 4 12 12 1 4 8 4 9 2 5
12 1 7 8 5 3 3 4 7 10 5 8 8 6 4 6 12 9 1 12 9 7 12 5 5 3 5 1 12 7 6 12 6 9 3 2 6 10 6 7 8
12 5 9 10 11 12 3 10 12 2 9 1 8 8 8 12 11 11 12 4 1 6 1 12 12 4 3 4 10 3 3 8 6 4 9 9 4 12 7
12 12 10 1 4 10 11 9 7 12 6 3 1 11 7 1 11 10 3 9 4 7 6 10 12 4 12 9 6 12 6 1 5 6 11 8 9 3
12 1 3 1 9 1 4 11 12 11 4 12 8 8 8 3 6 4 1 4 8 6 4 5 12 5 3 1 4 11 10 4 5 10 3 6 8 4 9 1 5
3 3 4 3 10 1 8 9 9 4 10 7 5
#> *** Sample cost: 108.3405
#> *** Number of strata: 12

```



```

#>
#> *** Sample size : 108

```

```
#> *** Number of strata : 12
#> -----

sum(solution4$aggr_strata$SOLUZ)
#> [1] 108.3405
```

This time the sample size is much higher.

What about the expected CV's?

```
newstrata <- updateStrata(strata2,solution4)
framenew2 <- updateFrame(frame,newstrata)
framenew2 <- framenew2[order(framenew2$ID),]
framenew2$Y2 <- nations$infant.mortality
framenew2$Y3 <- nations$contraception
results2 <- evalSolution(framenew2, solution4$aggr_strata, 50, progress = FALSE)
results2$coeff_var
#>          CV1          CV2          CV3  dom
#> 1 0.00566401 0.05020429 0.02625483 DOM1
```

This time the expected CV's of all variables are more than compliant with the precision constraints.

Appendix - Methodological approach

In a stratified sampling design with one or more stages, a sample is selected from a frame containing the units of the population of interest, stratified according to the values of one or more auxiliary variables (X) available for all units in the population.

For a given stratification, the overall size of the sample and the allocation in the different strata can be determined on the basis of constraints placed on the expected accuracy of the various estimates regarding the survey variables (Y).

If the target survey variables are more than one the optimization problem is said to be *multivariate*; otherwise it is *univariate*.

For a given stratification, in the univariate case the optimization of the allocation is in general based on the *Neyman allocation*. In the univariate case it is possible to make use of the *Bethel algorithm*.

The criteria according to which stratification is defined are crucial for the efficiency of the sample.

With the same precision constraints, the overall size of the sample required to satisfy them may be significantly affected by the particular stratification chosen for the population of interest.

Given G survey variables, their sampling variance is:

$$Var(\hat{Y}_g) = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{h,g}^2}{n_h} \quad g = 1, \dots, G$$

If we introduce the following cost function:

$$C(n_1, \dots, n_H) = C_0 + \sum_{h=1}^H C_h n_h$$

the optimization problem can be formalized in this way:

$$\min = C_0 + \sum_{h=1}^H C_h n_h$$

under the constraints

$$\begin{cases} CV(\hat{Y}_1) < U_1 \\ CV(\hat{Y}_2) < U_2 \\ \dots \\ CV(\hat{Y}_G) < U_G \end{cases}$$

where

$$CV(\hat{Y}_g) = \frac{\sqrt{Var(\hat{Y}_g)}}{mean(\hat{Y}_g)}$$

Given a population frame with m auxiliary variables X_1, \dots, X_M we define as *atomic stratification* the one that can be obtained considering the cartesian product of the definition domains of the m variables.

$$L = \{(l_1), (l_2), \dots, (l_k)\}$$

Starting from the atomic stratification, it is possible to generate all the different stratifications that belong to the universe of stratifications. For example:

$$\begin{aligned} P_1 &= \{(l_1, l_2, l_3)\} & P_2 &= \{(l_1), (l_2, l_3)\} \\ P_2 &= \{(l_2), (l_1, l_3)\} & P_4 &= \{(l_{31}), (l_1, l_2)\} \\ P_5 &= \{(l_1), (l_2), (l_k)\} \end{aligned}$$

The number of feasible stratifications is exponential with respect to the number of initial atomic strata:

$$B_4 = 15 \quad B_{10} = 115975 \quad B_{100} \approx 4.76 \times 10^{115}$$

In concrete cases, it is therefore impossible to examine all the different possible alternative stratifications. The *Genetic Algorithm* allows to explore the universe of stratification in a very efficient way in order to find the optimal (or close to optimal) solution.

In planning a stratified sampling for a given survey, proceed as follows:

- given the survey variables Y_1, Y_2, \dots, Y_p , set *precision constraints* on their estimates in the different domains, expressed in terms of CVs (coefficients of variation);
- in the available sampling frame build the *atomic stratification*, obtained as cartesian product of the domains of the auxiliary variables X_1, \dots, X_M ;
- in each atomic stratum report the distributional characteristics of the survey variables by calculating their *means* and *standard deviations* calculate the values of the population (directly or by using proxy variables);
- on the basis of these inputs, the Genetic Algorithm determines the *best solution* in terms of both frame *stratification*, *sample size* and *allocation* in optimized strata.

The application of the genetic algorithm is based on the following steps:

- a given stratification is considered as an *individual* in a population (= *generation*) subject to *evolution*;
- each individual is characterized by a *genome* represented by a vector of dimension equal to the number of atomic strata: the position of each element in the vector identifies an atomic stratum;
- each element in the vector is assigned a random value between 1 and K (maximum acceptable number of strata): the vector therefore indicates the way in which the individual atomic strata are aggregated together;
- for each individual (stratification) its *fitness* is calculated by solving the corresponding problem of optimal allocation by means of Bethel's algorithm;
- in passing from one generation to the next, *individuals with higher fitness are favored*;

- at the end of the process of evolution, *the individual with the overall best fitness represents the optimal solution.*

References

- Baillargeon, Sophie, and Louis-Paul Rivest. 2014. *Stratification: Univariate Stratification of Survey Populations*. <https://CRAN.R-project.org/package=stratification>.
- Ballin, M., and G. Barcaroli. 2013. "Joint Determination of Optimal Stratification and Sample Allocation Using Genetic Algorithm." *Survey Methodology* 39: 369–93.
- Ballin, Marco, and Giulio Barcaroli. 2016. "Optimization of Stratified Sampling with the R Package SamplingStrata: Applications to Network Data." *Computational Network Analysis with R: Applications in Biology, Medicine and Chemistry*. John Wiley & Sons.
- Barcaroli, G. 2014. "SamplingStrata: An R Package for the Optimization of Stratified Sampling." *Journal of Statistical Software* 61 (4): 1–24. <https://www.jstatsoft.org/article/view/v061i04>.
- Hartigan, J. A., and M. A. Wong. 1979. "A K-Means Clustering Algorithm." *Applied Statistics* 28: 100–108.
- O'Luing, M., S. Prestwich, and S. Armagan Tarim. 2017. "Constructing Strata to Solve Sample Allocation Problems by Grouping Genetic Algorithm." *arXiv:1709.03076*.
- Willighagen, E. 2005. *Genalg: R Based Genetic Algorithm*. <https://CRAN.R-project.org/package=genalg>.