

# R package **stratification** summary table

Sophie Baillargeon and Louis-Paul Rivest

November 1, 2024

A summary table of the package **stratification** can be found in the appendix of Baillargeon and Rivest (2011). Since the publication of this paper, the package has been updated (see the NEWS file for more details). At the end of this short note you will find an update of this summary table that reflects the changes made to the package. This table aims at providing a quick reference for the R package **stratification**. It lists the five public functions in **stratification** and their arguments. The following notes complete the table

(1) According to the general allocation scheme (Hidirogloou and Srinath, 1993). The stratum sample sizes are proportional to  $N_h^{2q_1} \bar{Y}_h^{2q_2} S_{yh}^{2q_3}$  (see `help(stratification)` for more details).

(2) The elements of the `model.control` argument depend on the model :  
- **loglinear** model with mortality :

$$Y = \begin{cases} \exp(\alpha + \text{beta} \log(X) + \text{epsilon}) & \text{with probability } p_h \\ 0 & \text{with probability } 1 - p_h \end{cases}$$

where `epsilon ~ N(0, sig2)` is independent of `X`. The parameter `p_h` is specified through `ph`, `ptakenone` and `pcertain`.

- heteroscedastic **linear** model :

$$Y = \text{beta}X + \text{epsilon} \quad \text{where } \text{epsilon} \sim N(0, \text{sig2} X^{\text{gamma}})$$

- **random** replacement model:

$$Y = \begin{cases} X & \text{with probability } 1 - \text{epsilon} \\ X_{new} & \text{with probability } \text{epsilon} \end{cases}$$

where `Xnew` is a random variable independent of `X` with the same distribution as `X`.

The following table presents `model.control` default values according to the model.

| model       | beta | sig2 | ph                      | ptakenone | pcertain | gamma | epsilon |
|-------------|------|------|-------------------------|-----------|----------|-------|---------|
| "loglinear" | 1    | 0    | <code>rep(1, Ls)</code> | 1         | 1        | -     | -       |
| "linear"    | 1    | 0    | -                       | -         | -        | 0     | -       |
| "random"    | -    | -    | -                       | -         | -        | -     | 0       |

(3) The default value of `initbh` is the boundaries obtained with the cumulative root frequency method of Dalenius and Hodges (1959) for Kozak's algorithm, and the set of arithmetic starting points of Gunning and Horgan (2007) for Sethi's algorithm. If `takenone=1` and `initbh` is of size  $L_s-1$ , the initial boundary of the take-none stratum is set to the first percentile of  $X$ .

(4) The following table summarize information about elements of `algo.control`. For a complete description of every element see `help(strata.LH)`. Sethi's algorithm has only one customizable parameter, the maximal number of iterations `maxiter`. However, for Kozak's algorithm, every parameter in the table below apply.

| parameter             | description   | format                                     | default  |
|-----------------------|---|--|--|
| <code>maxiter</code>  | maximal number of iterations  | positive integer                           | 500 (Sethi) or<br>10 000 (Kozak)                           |
| <code>minsol</code>   | if the number of solutions is below<br><code>minsol</code> $\Rightarrow$ complete enumeration | integer $\geq 2$ and<br>$\leq 2\ 000\ 000$ | 10 000   |
| <code>idopti</code>   | identification of stratum sample sizes<br>used in optimization criteria calculation           | "nh" or<br>"nhnonint"                      | "nh"   |
| <code>minNh</code>    | minimum size for sampled strata   | integer $\geq 2$                           | 2  |
| <code>maxstep</code>  | maximal step for boundary modification  | integer $\geq 2$                           | * $N_u/10$ , rounded up<br>and truncated to 100            |
| <code>maxstill</code> | maximal number of iterations without<br>boundary modification                                 | positive integer                           | $\text{maxstep} \times 10$ , bounded<br>between 50 and 500 |
| <code>rep</code>      | number of repetition of the algorithm   | integer $\geq 1$                           | 5  |
| <code>trymany</code>  | indicator for trying many initial<br>stratum boundaries                                       | TRUE or FALSE                              | TRUE   |

\* $N_u$  = number of unique values of the stratification variable  $X$  (without considering the units in the certainty stratum)

## References

- Baillargeon, S. and Rivest L.-P. (2011). The construction of stratified designs in R with the package `stratification`. *Survey Methodology*, **37**(1), 53-65. <http://www.statcan.gc.ca/pub/12-001-x/2011001/article/11447-eng.pdf>
- Dalenius, T. and Hodges, J.L., Jr. (1959). Minimum variance stratification, *Journal of the American Statistical Association*, **54**, 88-101.
- Gunning, P. and Horgan, J. M. (2007). Improving the Lavallée and Hidiroglou algorithm for stratification of skewed populations, *Journal of Statistical Computation and Simulation*, **77**, 277-291
- Hidiroglou, M. A. , and Srinath, K. P. (1993). Problems associated with designing subannual business surveys, *Journal of Business and Economic Statistics*, **11**, 397-405

| argument                     | description  | format   | default  |
|------------------------------|--|--|--|
| <code>strata.cumrootf</code> |  |  |  |
| <code>strata.Geo</code>      |  |  |  |
| <code>strata.LH</code>       |  |  |  |
| <code>strata.bh</code>       |  |  |  |
| <code>var.strata</code>      |  |  |  |
| <code>alloc</code>           | ✓ ✓ ✓ ✓ ✓ allocation specification (1)   | list (q1,q2,q3) where $q_i \geq 0$   | Neyman ( $q_1=q_3=0.5, q_2=0$ )  |
| <code>takenone</code>        | ✓ ✓ ✓ number of take-one strata  | 0 or 1   | 0  |
| <code>bias.penalty</code>    | ✓ ✓ penalty for the bias   | numeric $\in [0, 1]$   | 1  |
| <code>takeall.adjust</code>  | ✓ ✓ ✓ number of take-all strata<br>indicator of adjustment for take-all strata           | one of $\{0, 1, \dots, L_s - 1\}$<br>TRUE or FALSE   | 0<br>TRUE (as in the rest of the package)  |
| <code>rh.postcorr</code>     | ✓ ✓ ✓ ✓ anticipated response rates<br>indicator of posterior correction for non-response | numeric (vector or not)<br>TRUE or FALSE   | <code>rep(1, Ls)</code> or <code>rh</code> from <code>strata</code><br>FALSE (no correction) |
| <code>model</code>           | ✓ ✓ ✓ ✓ ✓ model identification   | "none", "loglinear",<br>"linear"*, or "random" * →   | "none"<br>(* unavailable with Sethi's algo)  |
| <code>model.control</code>   | ✓ ✓ ✓ ✓ ✓ model's parameter specification (2)  | list ( <code>beta</code> , <code>sig2</code> , <code>ph.ptlnone</code> ,<br><code>pertain</code> , <code>gamma</code> , <code>epsilon</code> )   | depends on <code>model</code> , but equivalent to <code>model="none"</code>                  |
| <code>nclass</code>          | ✓ number of classes  | integer $\geq L_s$   | $\min(15L_s, Nu)$  |
| <code>initbh</code>          | ✓ initial stratum boundaries   | numeric vector   | depends on <code>algo</code> (3)   |
| <code>algo</code>            | ✓ algorithm identification   | "Kozak" or "Sethi"   | "Kozak"  |
| <code>algo.control</code>    | ✓ algorithm's parameters specification (4)   | list ( <code>maxiter</code> , <code>minsol</code> , <code>idopti</code> , <code>minWh</code> ,<br><code>maxstep</code> , <code>maxstill</code> , <code>rep</code> , <code>trymany</code> ) | depends on <code>algo</code>   |
| <code>strata</code>          | ✓ stratified design  | strata object  | none ( <code>strata</code> is mandatory)   |
| <code>y</code>               | ✓ study variable   | numeric vector   | NULL (model given instead)   |