

# ReGenesees

(R evolved Generalised software for sampling estimates and errors in surveys)

## *Scope*

**Design-Based and Model-Assisted Analysis of Complex Sample Surveys**

## *Main Statistical Functions*

- **Complex Sampling Designs**
  - Multistage, stratified, clustered, sampling designs
  - Sampling with equal or unequal probabilities, with or without replacement
  - “Mixed” sampling designs (i.e. with both Self-Representing and Non-Self-Representing strata)
- **Calibration**
  - Global and partitioned (for factorizable calibration models)
  - Unit-level and cluster-level weights adjustment
  - Homoscedastic and heteroscedastic models
  - Linear, raking and logit distance functions
  - Bounded and unbounded weights adjustment
  - Multi-step calibration
  - Consistent trimming of calibration weights
- **Basic Estimators**
  - Horvitz-Thompson
  - Calibration Estimators
- **Variance Estimation**
  - Multistage formulation (via Bellhouse recursive algorithm)
  - Ultimate Cluster approximation
  - Collapsed strata technique for handling lonely PSUs
  - Taylor linearization of nonlinear smooth estimators
  - Generalized Variance Functions (GVF) method
- **Estimates and Sampling Errors (standard error, variance, coefficient of variation, confidence interval, design effect) for:**
  - Totals
  - Means
  - Absolute and relative frequency distributions (marginal, conditional and joint)
  - Ratios between totals
  - Shares and ratios between shares
  - Multiple regression coefficients
  - Quantiles (variance estimation via the Woodruff method)

- **Estimates and Sampling Errors for Complex Estimators**
  - Handles arbitrary differentiable functions of Horvitz-Thompson or Calibration estimators
  - Complex Estimators can be freely defined by the user
  - Automated Taylor-linearization
  - Design covariance and correlation between Complex Estimators
- **Estimates and Sampling Errors for Subpopulations (Domains)**
  - All the analyses above can be carried out for arbitrary domains

Under development:

- Replication based Variance Estimation for non-analytic estimators, through the Delete-A-group Jackknife (DAGJK) technique: this will integrate the [EVER](#) package with the ReGenesees system.

## ***System Architecture***

ReGenesees is a full-fledged software system entirely developed in R. It has a clear-cut two-layer architecture. The application layer of the system is embedded into an R package named itself **ReGenesees**. A second R package, called **ReGenesees.GUI**, implements the presentation layer of the system. Both packages can be run under Windows, Mac, as well as under most of the Unix-like operating systems. While the **ReGenesees.GUI** package requires the **ReGenesees** package, the latter can be used also without the GUI on its top. This means that the statistical functions of the system will always be accessible by users interacting with R through the traditional command-line interface. On the contrary, less experienced R users will take advantage from the user-friendly mouse-click graphical interface.

## ***Data Input/Output***

The ReGenesees system can import data in a variety of ways. First, it can load R workspace files (with .RData or .rda extensions) storing previously saved data. Second, data can be imported from Text Files (with extensions .txt, .csv, .dat). Third, the system can import data from MS Excel spreadsheets and/or MS Access database tables. Further extensions are possible. Currently, ReGenesees can save output data into R workspace files (.RData, .rda) and/or export them into Text Files (.txt, .csv, .dat). Further extensions are possible.

## ***Development Status***

The current version of the ReGenesees system is **1.9**

## ***Software Documentation***

Both packages composing the system (**ReGenesees** and **ReGenesees.GUI**) come with their own reference manuals, which fulfill R standards for packages' documentation.

## ***Software Distribution***

The ReGenesees system is distributed as Open Source Software, under the EUPL license.

## ***Authors***

Overall Project: Diego Zardetto ([zardetto@istat.it](mailto:zardetto@istat.it))

Application layer (i.e. **ReGenesees** package): Diego Zardetto

Presentation layer (i.e. **ReGenesees.GUI** package): Diego Zardetto, and Raffaella Cianchetta

## ***Download***

The ReGenesees system can be downloaded from:

- The European Commission Repository for Open Source Software (Joinup):  
<https://joinup.ec.europa.eu/software/regenesees/description>
- Istat website
  - English:  
<http://www.istat.it/en/tools/methods-and-it-tools/processing-tools/regenesees>
  - Italian:  
<http://www.istat.it/it/strumenti/metodi-e-strumenti-it/strumenti-di-elaborazione/regenesees>

## Sample GUI Screenshots

**ReGENESEES 1.0 [pkg] - 1.0 [gui]**

**ReGENESEES**

R EVOLVED GENERALISED SOFTWARE  
FOR ESTIMATES AND ERRORS IN SURVEYS

Authors: Diego Zarbetto, Raffaella Diancinetta

Istat

START

**e calibrate**

Population and Survey Data

Select population table: **pop.tbl**

Select a survey design object: **shades**

Variables: **id**, **public**, **emp.num**, **emp.cl**

Formula composer

Formula: **calmodel**

Formula: **calmodel**

Optional Fields

low: **0.1**, high: **10**, aggregate.stage: **NULL**, sigma2: **emp.num**, maxit: **50**, epsilon: **1e-7**, force: **False**

Output Object Name: **shades**

**ReGENESEES 1.0 [pkg] - 1.0 [gui]**

File Data Functions Tools Options Help

**Commands Window**

```
## ReGENESEES session start:
## Sat Mar 31 16:52:01 2012

shades <- e.svydesign(data= shades, id= id, strata= strata, weights= weight, fpc= fpc, seir.rep.st= NULL, check.data= TRUE)

va.area.HT <- svydataTT(design= shades, y= va.imp2, by= area, estimator= "Total",
vartype= "se", conf.int= FALSE, conf.level= 0.95, deff= FALSE, na.rm= FALSE)

totals <- pop.template(data= shades, calmodel= (emp.num + ent):emp.cl - 1, partition= area)

totals.HT <- new.certificate(design= shades, calmodel= (emp.num + ent):emp.cl - 1, partition= area, template= totals)

totals <- fill.template(universe= pop.frame, template= totals, mem.frame= 10)

shades <- e.calibrate(design= shades, df.population= totals, calmodel= (emp.num + ent):emp.cl - 1, partition= area, calvar= "linear", bounds= c(-Inf, Inf), aggregate.stage= NULL, sigma2= emp.num, maxit= 50, epsilon= 1e-07, force= TRUE)
```

**Warnings Window**

**ReGENESEES**

id	public	emp.num	emp.cl	nace5	nace2	area	cena	region	va.cl	va	dona	nace	nace2	dona2
1	1268	0	38	(19,49)	1210	1	32	0	Center	22	5500.0	1	(19,49)	Agriculture
2	1338	0	39	(19,49)	1240	1	32	0	Center	19	1300.0	1	(19,49)	Agriculture
3	13819	0	25	(19,49)	1131	1	41	0	Center	16	400.0	1	(19,49)	Agriculture
4	13749	0	25	(19,49)	1111	1	43	0	Center	1	0.0	1	(19,49)	Agriculture
5	8431	0	29	(19,49)	1121	1	31	0	Center	2	0.5	1	(19,49)	Agriculture
6	7572	0	59	(19,99)	1132	1	41	0	Center	11	400.0	1	(19,99)	Agriculture
7	9701	0	67	(49,99)	1240	1	33	0	Center	23	7000.0	1	(49,99)	Agriculture
8	8461	0	56	(49,99)	1137	1	39	0	Center	14	400.0	1	(49,99)	Agriculture
9	11899	0	52	(49,99)	1131	1	41	0	Center	16	400.0	1	(49,99)	Agriculture
10	15136	0	12	(9,19)	1111	1	43	0	Center	1	0.0	1	(9,19)	Agriculture
11	10980	0	10	(9,19)	1240	1	43	0	Center	18	750.0	1	(9,19)	Agriculture
12	2229	0	143	(99,Inf)	1240	1	33	1	Center	26	30000.0	1	(99,Inf)	Agriculture
13	12258	0	353	(99,Inf)	1113	1	41	1	Center	23	7000.0	1	(99,Inf)	Agriculture
14	5477	0	7	(6,9)	1111	1	43	0	Center	21	3500.0	1	(6,9)	Agriculture
15	3894	0	7	(6,9)	1111	1	31	0	Center	18	750.0	1	(6,9)	Agriculture
16	7640	0	20	(19,49)	1410	14	31	0	Center	16	400.0	14	(19,49)	Industry
17	14165	0	22	(19,49)	1410	14	42	0	Center	16	400.0	14	(19,49)	Industry
18	1186	0	21	(19,49)	1410	14	43	0	Center	19	1500.0	14	(19,49)	Industry
19	1420	0	36	(19,49)	1410	14	31	0	Center	19	1500.0	14	(19,49)	Industry
20	3848	0	90	(49,99)	1410	14	31	0	Center	20	2500.0	14	(49,99)	Industry
21	15162	0	82	(49,99)	14300	14	31	0	Center	31	3500.0	14	(49,99)	Industry
22	12380	0	51	(49,99)	1410	14	41	0	Center	24	9000.0	14	(49,99)	Industry
23	3515	0	13	(6,18)	1410	14	32	0	Center	15	250.0	14	(6,18)	Industry
24	7214	0	13	(9,19)	1410	14	31	0	Center	15	250.0	14	(9,19)	Industry
25	9493	0	172	(99,Inf)	1410	14	42	1	Center	31	3500.0	14	(99,Inf)	Industry
26	15498	0	8	(6,9)	1410	14	41	0	Center	15	250.0	14	(6,9)	Industry
27	8124	0	7	(6,9)	1410	14	32	0	Center	16	400.0	14	(6,9)	Industry
28	5738	0	28	(19,49)	15421	15	32	0	Center	26	30000.0	15	(19,49)	Industry
29	4769	0	21	(19,49)	15811	15	41	0	Center	16	400.0	15	(19,49)	Industry
30	1286	0	22	(19,49)	13980	15	31	0	Center	19	1300.0	15	(19,49)	Industry

## References

- Woodruff, R. S. - (1952)**  
*"Confidence Intervals for Medians and Other Position Measures"*  
 Journal of the American Statistical Association,  
 Vol. 47, n. 260, pp. 635-646.
- Woodruff, R. S. - (1971)**  
*"A Simple Method for Approximating the Variance of a Complicated Estimate"*  
 Journal of the American Statistical Association,  
 Vol. 66, n. 334, pp. 411-414.
- Wilkinson, G.N., Rogers, C.E. - (1973)**  
*"Symbolic Description of Factorial Models for Analysis of Variance"*  
 Journal of the Royal Statistical Society, series C (Applied Statistics),  
 Vol. 22, pp. 181-191.
- Kalton, G. - (1979)**  
*"Ultimate cluster sampling"*  
 Journal of the Royal Statistical Society, series A (General),  
 Vol. 142, Part 2, pp. 210-222.
- Binder, D. A. - (1983)**  
*"On the variances of asymptotically normal estimators from complex surveys"*  
 International Statistical Review,  
 Vol. 51, n. 3, pp. 279-292.

- **Rust, K.** - (1985)  
*"Variance Estimation for Complex Estimators in Sample Surveys"*  
 Journal of Official Statistics,  
 Vol. 1, n. 4, pp. 381-397.
- **Bellhouse, DR.** - (1985)  
*"Computing Methods for Variance Estimation in Complex Surveys"*  
 Journal of Official Statistics,  
 Vol.1, n. 3, pp. 323-329.
- **Rust, K., Kalton, G.** - (1987)  
*"Strategies for Collapsing Strata for Variance Estimation"*  
 Journal of Official Statistics,  
 Vol. 3, n. 1, pp. 69-81.
- **Korn, E.L., Graubard, B.I.** - (1990)  
*"Simultaneous testing of regression coefficients with complex survey data: Use of Bonferroni t statistics"*  
 The American Statistician,  
 Vol. 44, n. 4, pp. 270-276.
- **Särndal, C.E., Swensson, B., Wretman, J.** - (1992)  
*"Model Assisted Survey Sampling"*  
 Springer Verlag.
- **Deville, J.C., Särndal, C.E.** - (1992)  
*"Calibration Estimators in Survey Sampling"*  
 Journal of the American Statistical Association,  
 Vol. 87, n. 418, pp. 376-382.
- **Chambers, J.M., Hastie, T.J.** - (1992)  
*"Statistical Models in S"*  
 Wadsworth & Brooks/Cole.
- **Deville, J.C., Särndal, C.E., Sautory, O.** - (1993)  
*"Generalized Raking Procedures in Survey Sampling"*  
 Journal of the American Statistical Association,  
 Vol. 88, n. 423, pp.1013-1020.
- **Sautory, O.** - (1993)  
*"La macro CALMAR: Redressement d'un Echantillon par Calage sur Marges"*  
 Document de travail de la Direction des Statistiques Demographiques et Sociales,  
 n. F9310.
- **Dorfman, A., Valliant, R.** - (1993)  
*"Quantile variance estimators in complex surveys"*  
 Proceedings of the ASA Survey Research Methods Section,  
 pp. 866-871.
- **Kish, L.** - (1995)  
*"Methods for design effects"*  
 Journal of Official Statistics,  
 Vol. 11, n. 1, pp. 55-77.
- **Estevao, V., Hidirolou, M. A., Särndal, C. E** - (1995)  
*"Methodological principles for a generalized estimation system at Statistics Canada"*  
 Journal of Official Statistics,  
 11, n. 2, pp. 181-204.
- **Singh, A.C., Mohl, C.A.** - (1996)  
*"Understanding calibration estimators in survey sampling"*  
 Survey Methodology,  
 22, pp. 107-115.

- **Rao, J. N. K., Lohr, S. L.** - (1999)  
*"Some Current Trends in Sample Survey Theory and Methods"*  
 Sankhya: The Indian Journal of Statistics, Special issue on Sample Surveys,  
 Vol. 61, Series B, Pt. 1, pp. 1-57.
- **Valliant, R.** - (2000)  
*"Variance estimation for the general regression estimator"*  
 Survey Methodology,  
 28, pp. 103-114.
- **Vanderhoeft, C.** - (2001)  
*"Generalized Calibration at Statistic Belgium"*  
 Statistics Belgium Working Paper n. 3  
[http://statbel.fgov.be/nl/binaries/paper03%5B1%5D\\_tcm325-35412.pdf](http://statbel.fgov.be/nl/binaries/paper03%5B1%5D_tcm325-35412.pdf)
- **Fuller, W.A.** - (2002)  
*"Regression estimation for survey samples"*  
 Survey Methodology,  
 28, pp. 5-23.
- **Rao, J. N. K., Lohr, S. L.** - (2004)  
*"Sample Survey Methods: Recent Developments and Applications"*  
 two-day workshop slides, Joint Statistical Meetings, Toronto.
- **Lumley, T.** - (2006)  
*"survey: analysis of complex survey samples"*  
 R package version 3.6-5.  
<http://cran.at.r-project.org/web/packages/survey/index.html>
- **Wolter, K. M.** - (2007)  
*"Introduction to Variance Estimation"*  
 Second Edition, Springer-Verlag, New York.
- **Scannapieco, M., Zardetto, D., Barcaroli, G.** - (2007)  
*"La Calibrazione dei Dati con R: una Sperimentazione sull'Indagine Forze di Lavoro ed un Confronto con GENESEES/SAS"*  
 Contributi Istat n. 4.  
[http://www3.istat.it/dati/pubbsci/contributi/Contributi/contr\\_2007/2007\\_4.pdf](http://www3.istat.it/dati/pubbsci/contributi/Contributi/contr_2007/2007_4.pdf)
- **Lumley, T.** - (2012)  
*"Complex Surveys: A Guide to Analysis Using R"*  
 John Wiley & Sons, New York.
- **Barcaroli, G., Zardetto, D.** - (2012)  
*"Use of R in Business Surveys at the Italian National Institute of Statistics: Experiences and Perspectives"*  
 Proceedings of the 4<sup>th</sup> International Conference of Establishment Surveys (ICES IV),  
 American Statistical Association.  
<http://www.amstat.org/meetings/ices/2012/papers/302193.pdf>
- **Zardetto, D.** - (2013)  
*"ReGenesees: an Advanced R System for Calibration, Estimation and Sampling Errors Assessment in Complex Sample Surveys"*  
 Proceedings of the 7<sup>th</sup> International Conference on New Techniques and Technologies for Statistics (NTTS 2013), Eurostat.  
[http://www.cros-portal.eu/sites/default/files//NTTS2013fullPaper\\_131-v2.pdf](http://www.cros-portal.eu/sites/default/files//NTTS2013fullPaper_131-v2.pdf)

- **Fallows A., Pope M., Digby-North J., Brown G., Lewis D. - (2015)**  
*"A Comparative Study of Complex Survey Estimation Software in ONS"*  
Romanian Statistical Review,  
n. 3, pp. 46-64.  
<http://www.revistadestatistica.ro/index.php/comparative-study-of-complex-survey-estimation-software-in-ons/>
- **Zardetto, D. - (2015)**  
*"ReGenesees: an Advanced R System for Calibration, Estimation and Sampling Error Assessment in Complex Sample Surveys"*  
Journal of Official Statistics,  
Vol. 31, n. 2, pp. 177-203.  
<http://www.istat.it/it/files/2014/05/Zardetto-jos-2015-0013.pdf>