

Methods and IT tools for statistical production

Content of web area Methods and IT tools for statistical production

To send information and/or updates on methods and IT tools please write to <u>softgen@istat.it</u>

Last update: June 2017

INDEX

Methods and IT tools for statistical production	1
DESIGN phase	
DESIGN phase – METHODS	5
DESIGN phase – TOOLS	6
FS4 (First Stage Stratification and Selection in Sampling)	7
MAUSS-R (Multivariate Allocation of Units in Sampling Surveys)	9
Multiway Sample Allocation	11
SamplingStrata	14
COLLECT phase	17
COLLECT phase – METHODS	20
COLLECT phase – TOOLS	21
Blaise	22
PROCESS phase	24
PROCESS phase – METHODS	32
PROCESS phase – TOOLS	
RELAIS (REcord Linkage At IStat)	
StatMatch	
CIRCE (Comprehensive Istat R Coding Environment)	41
Banff	44
CANCEIS (CANadian Census Edit and Imputation System)	46
CONCORDJava (CONtrollo e CORrezione dei Dati version with Java interface)	47
SeleMix (Selective editing via Mixture models)	49
EVER (Estimation of Variance by Efficient Replication)	51
ReGenesees (R evolved Generalised software for sampling estimates and errors in sur	veys)53
ANALYSE phase	57
ANALYSE phase – METHODS	67
ANALYSE phase – TOOLS	70
COMIC	71
Ranker	73
ARGUS	76

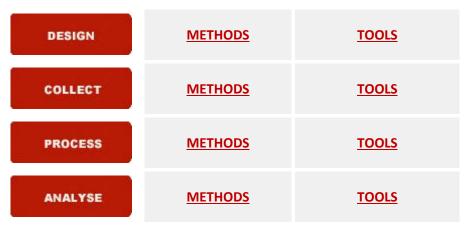
Methods and IT tools for statistical production

The Repository of the methods and IT tools for statistical production contains:

- statistical methods
- generalized IT tools

validated and used at Istat for the production of statistical outputs.

The Repository is organized by **phases of the statistical production process**, in accordance with the <u>Generic Statistical Business Process Model (GSBPM) Version 5.0</u>.



PHASES OF THE PROCESS

The phases of the production process are the "gates" of access to information and specific materials. Only the phases Design, Collect, Process, Analyse, i.e., phases for which methods and tools are currently available, are taken into account.

The PHASES pages contain the most important information about the sub-processes within the selected phase.

Methods and tools are grouped by sub-processes. Methods and tools relevant to more than one sub-process are attributed to the prevailing sub-process.

The METHODS pages contain lists of documents that describe statistical methods. The documents are Istat standards (or European Statistical System standards), or methodological manuals, publications and reports drawn up by Istat researchers (also in collaborations). For a given sub-process, the documents are displayed in reverse chronological order (newest to oldest) and are available in PDF format or by linking to external sites.

The TOOLS pages contain lists of generalized IT tools that implement statistical methods. For a given sub-process, the tools are listed in alphabetical order. For each tool a general description showing the functions implemented is provided together with an information sheet.

Tools developed by Istat can be **downloaded** for free. However, at the first download you are requested to provide a valid e-mail address. How to find the tools not owned by Istat is also shown.

Last updates

- <u>StatMatch</u> Published the release 1.2.5 of the R package for data integration through statistical matching and, as a by-product, the ability to impute missing values in a data set | June 16, 2017
- <u>ReGenesees</u> Published the release 1.9 of the full-fledged R software for design-based and model-assisted analysis of complex sample surveys. New function 'trimcal' trims calibration weights to a bounded interval, while simultaneously preserving all the calibration constraints | May 8, 2017
- <u>COMIC</u> Published version 1.0 of the SAS based software for the construction of composite indices, through various aggregation methods, and evaluation of their robustness | March 15, 2017

DESIGN phase

Design frame and sample

The activities related to the design of the list and the sampling methodology refer to subprocess 2.4 "*Design frame and sample*" of <u>GSBPM</u>. More precisely:

- Construction of the selection frame of the target population, containing, for each unit of the population, all the identifying information needed for the contact, any auxiliary variables used to define the sample design (stratification variables, identifying variables of any selection stage);
- Design of the sampling design which, on the basis of the objectives specified in Phase 1 "Specify Needs" and of the operating and cost constrains, allows obtaining estimates as accurate as possible.

Design of the selection list

The characteristics of the sample list are essential for the correct definition of the sampling design.

The list should meet quality criteria in terms of refreshing, coverage and accuracy of the information contained in it. From a theoretical point of view the selection list should ideally have the following requirements:

- it is composed only of the units belonging to the population of interest at the reference time of the survey;
- it includes all units of the population only once;
- it contains the most updated data for identifying variables (name and address) and for any descriptive information (other relevant structural data) of the units.

Possible situations of departure from the ideal list are:

- *under-coverage*, which occur when some elements of the population are not contained in the list and cannot, therefore, be included in the sample;
- over-coverage, when some elements of the list are non-existent and / or do not belong to the population of interest;
- *duplication of some units,* when some elements of the population are enumerated more than once in the list;
- *clusters of units,* when some elements of the list contain clusters of elements of the population.

Planning and implementation of the sample design

The planning of the sample design consists firstly of the following activities:

• definition of the sampling scheme, performed on the basis of the cost related to the chosen data collection technique and the information contained in the selection list (multistage sampling selection, stratified sampling). The choice of a multistage design

generally derives from the need to concentrate the sample locally in order to limit the cost of interview in case of survey using a direct mode of administration of the questionnaire (face to face interview). The choice of a stratified sample design has the purpose of improving the precision of the estimates and guaranteeing planned sample size for the domains of estimate. The division of the units of the population in strata is carried out using auxiliary variables contained in the list and related to the variables under investigation.

On the basis of the adopted sampling scheme, the following steps may be undertaken:

- **choice of stratification criteria** (choice of variables, choice of the number of strata, definition of the criterion of strata construction);
- choice of the probabilistic method for the selection of the sample units (selection with equal probabilities, selection with unequal probabilities). For designs in two or more stages the selection of the primary sampling units (units of the first stage) is generally done with probability proportional to a measure of size suppository correlated with the variables under investigation.
- determination of sample size for the different stages of selection and allocation of the sample between the strata on the basis of the sampling error admitted for the main estimates, in relation to the reference domains and subclasses of the population. Since surveys are usually designed for the production of a variety of estimates for different domains of interest, it is necessary to use approaches that address the problem in a global perspective of determining the optimal sample size in the presence of a high number of objectives and constraints.

DESIGN phase – METHODS

• DESIGN FRAME AND SAMPLE

Generalized Framework for Defining the Optimal Inclusion Probabilities of One-Stage Sampling Designs for Multivariate and Multi-domain Surveys 2015 Survey Methodology, 41.

MEMOBUST – Handbook on Methodology of Modern Business Statistics: <u>Statistical Registers and Frames</u> <u>Sample Selection</u> 2014 MEETS ESSnet MEMOBUST

References

Istat. 2008. Strategia di campionamento e precisione delle stime. In "<u>L'indagine europea sui</u> <u>redditi e le condizioni di vita delle famiglie (Eu-silc)</u>", Collana Metodi e Norme, n. 37, Istat.

Istat. 2006. Il disegno campionario della nuova indagine e la fase di estrazione. In "<u>La rilevazione</u> <u>sulle forze di lavoro: contenuti, metodologie, organizzazione</u>", Collana Metodi e Norme, n. 32, Istat.

Istat. 2006. Il piano di campionamento. In "<u>Il sistema di indagini sociali multiscopo. Contenuti e</u> <u>metodologia delle indagini</u>", Collana Metodi e Norme, n. 31, Istat.

Istat. 2006. Strategia di campionamento e livello di precisione delle stime. In "*L'indaqine campionaria sulle nascite: obiettivi, metodologia e organizzazione*", Collana Metodi e Norme, n. 28, Istat.

Cicchitelli G., A. Herzel, G.E. Montanari. 1992. Il campionamento statistico. Il Mulino, Bologna.

Särndal C.E., B. Swensson, J. Wretman. 1992. *Model Assisted Survey Sampling*. Springer-Verlag, New York.

Bethel J. 1989. <u>Sample Allocation in Multivariate Surveys</u>. Survey Methodology 15:47-57.

Chromy J. 1987. <u>Design Optimization with Multiple Objectives</u>. Proceedings of the Survey Research Methods Section, *American Statistical Association*, 194-199.

Cochran W.G. 1977. Sampling Techniques. J. Wiley & Sons, New York.

DESIGN phase – TOOLS

• DESIGN FRAME AND SAMPLE

FS4 (First Stage Stratification and Selection in Sampling). FS4 is a generalized software for first stage stratification and selection in sampling related to two or more stages, implemented completely in R and with a GUI (Graphical User Interface).

<u>MAUSS-R</u> (Multivariate Allocation of Units in Sampling Surveys – version R with Java interface). Software for the determination of the sample allocation in the multivariate case and for multiple domains of estimate for single selection stage samples.

<u>Multiway Sample Allocation</u> R-package implementing the sample allocation for one stage multiway stratified (simple or with varying inclusion probabilities) sampling and for incomplete stratified sampling designs. The allocation allows to control the expected sampling errors of the estimates of many parameters and several reference subpopulations.

<u>SamplingStrata</u>. Optimal stratification of sampling frames for multipurpose sampling surveys.

FS4 (First Stage Stratification and Selection in Sampling)

Description

FS4 is an open source generalized software for first stage stratification and selection in sampling related to two or more stages, implemented completely in R language and with a GUI (Graphical User Interface).

It merges two data frames (whatever PSU – Primary Sampling Unit – population and allocation data frames) and then, computes stratification and selection of a fixed number of sample PSUs per stratum using the Sampford's method (unequal probabilities, without replacement, fixed sample size), implemented by the UPsampford function of R package "sampling".

As for stratification, the function carries out, for each estimation domain, the computation of a size threshold for a given size PSU. PSUs with measure of size exceeding a calculated threshold are identified as SR (Self Representative) and each constitutes a stratum on its own. The remaining NSR (Non Self Representative) PSUs are ordered by measure of size and divided into strata having sizes approximately constant to the corrected threshold and with PSUs having sizes as homogeneous as possible. For each PSU, FS4 determines the size of final sample units on the basis of the size planned in input.

The main features of FS4 are the following:

- it is a flexible software allowing the user to:
 - choose the measure of size to define stratification and inclusion probability for PPS (Probability Proportional to Size);
 - insert as input a single or plural (variable between domains) planned size of final sample units to observe in each PSU;
 - launch the procedure in two separate steps, so to observe an initial output and on the basis of this modify the input parameters;
- it is an open source package released under EUPL licence;
- it has a user friendly GUI that also allows to non-expert R users to use it, implemented in Tcl/Tk by means the R "tcltk" package, although for expert R users the StratSel function from R command line can be used directly.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	2.4 Design frame and sample4.1 Create frame and select sample
Programming language:	R
Language of the GUI:	EN
Keywords:	stratification, selection, sampling, PSU
Contact:	name: Raffaella Cianchetta email: <u>cianchet@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

R (version \geq 3.0.1)

R packages: tcltk2, svMisc, plyr and sampling

COPYRIGHT

Copyright 2013 Istat

Licensed under the European Union Public Licence (EUPL), version 1.0 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 1.0 Windows binaries
- Version 1.0 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – FS4 v. 1.0

User guide – FS4 v. 1.0

MAUSS-R (Multivariate Allocation of Units in Sampling Surveys)

Description

MAUSS-R is an open source generalized software for the determination of the sample allocation, developed using the R language and provided with a Java interface to facilitate the selection of data and parameters of interest.

The reference methodology is an extension of the Neyman allocation method to the case of several variables and adopts as a method of solving a generalization of the Bethel's proposal (1989).

The generalized Bethel algorithm allows to calculate the sample allocation for stratified sampling designs: the overall sample size and the allocation among the different strata is determined starting from the constraints imposed in the survey on the accuracy of the estimates of interest.

The software in the current version allows the determination of the sample allocation in the multivariate case and for multiple domains of estimate for single selection stage samples.

The interface allows:

- changing the default values of the processing parameters (such as the minimum number of units per stratum);
- the management of the file containing the constraints of precision and the storing of different versions of the constraints and the obtained results;
- the execution of the calculation of the optimal solution;
- the comparison with the proportional allocation and equal;
- viewing reports on population, the results and the comparison of the results obtained changing the constraints.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.0</u>
GSBPM code:	2.4 Design frame and sample
Programming language:	R, Java
Language of the GUI:	IT, EN
Keywords:	sampling design, optimal allocation, sample size, Bethel
Contact:	name: Maria Teresa Buglielli email: <u>bugliell@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

R (version $\geq 2.7.0$)

Only for graphical interface: Java 2 Runtime Environment (version \geq 6.0)

COPYRIGHT

Copyright 2013 Istat

Licensed under the European Union Public Licence (EUPL), version 1.0 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/en/document/7330.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

MAUSS-R version 1.1

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

User manual – MAUSS-R v. 1.1

Multiway Sample Allocation

Description

R-package Multiway Sample Allocation allows to allocate the sample of multi-way stratified sampling design.

The multi-way stratification is achieved by combining the variables defining the marginal cells of a multi-way contingency table. In the context of the sampling theory the marginal cells are identified by the categories of the variables defining the interest domain partitions being independent of each other. For example, the stratification obtained by combining the categories of two partitions into domains of interest defines a two-way stratification. A partition is denoted as independent or marginal when the domains cannot be got as aggregation of other domains defined in other partitions. The allocation of the sample in the dependent domains is done by aggregating the sample sizes of the marginal domains. In case of multi-way stratified designs, the sample size is fixed for each stratum. In case of incomplete stratified sampling design the sample size is fixed for each domain of interest.

The Multiway.Sample.Allocation package defines the allocation of the sample in accordance with the precision constraints of the estimates for:

- different parameters of interest (total) (multivariate problem) at sub-population or reference domain level (multi-domain issue);
- one stage stratified sampling designs simple or with varying inclusion probabilities variables within the multi-way strata, in which the sample size is fixed at the stratum level; incomplete stratified sampling designs in which the sample size is fixed at the domain level (except for rounding to the above/below integer);
- varying inclusion probability designs where the sample size is fixed (approximately) at the domain level (the performance of the allocation process subject to computational constraints associated with the population size and the number of the estimates to be considered).

The algorithm performing the allocation (Falorsi and Righi, 2015) is an extension of the Chromy (1987) and Bethel (1989) algorithm. The algorithm solves an optimization problem formalized according to a general expression of the variance of the estimates (Falorsi and Righi, 2015) depending on:

- the superpopulation model used to define input parameters (mean and variance of the variable defining the parameters of interest);
- the implemented sampling design.

The main output of the package is the inclusion probability of the population units.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	2.4 Design frame and sample
Programming language:	R
Keywords:	sample allocation, multi-way stratification, incomplete stratified sampling design
Contact:	name: Paolo Righi email: <mark>parighi@istat.it</mark>

Software and documentation

SOFTWARE DEPENDENCIES

R (version \geq 3.1.1), R-package <u>MASS</u>.

COPYRIGHT

Copyright 2016 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 1.0 Windows binaries
- Version 1.0 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual - Multiway Sample Allocation v. 1.0

OTHER DOCUMENTATION

De Vitiis C., P. Righi, M. D. Terribili. 2016. *Optimal sample allocation for the Incomplete Stratified Sampling design*. *Rivista di Statistica Ufficiale*, in press.

Falorsi P. D., P. Righi. 2016. A flexible tool for defining optimal sampling designs. <u>The Survey</u> <u>Statistician</u>, 73:21-31.

Falorsi P. D., P. Righi. 2015. <u>Generalized Framework for Defining the Optimal Inclusion</u> <u>Probabilities of One-Stage Sampling Designs for Multivariate and Multi-domain Surveys</u>. *Survey Methodology*, 41.

Falorsi P. D., P. Righi. 2008. <u>A Balanced Sampling Approach for Multi-way Stratification Designs</u> for Small Area Estimation. *Survey Methodology*, 34(2):223-234.

SamplingStrata

Description

When designing a sampling survey, usually constraints are set on the desired precision levels regarding one or more target estimates (the Y's). If a sampling frame is available, containing auxiliary information related to each unit (the X's), it is possible to adopt a stratified sample design. For any given stratification of the frame, in the multivariate case it is possible to solve the problem of the best allocation of units in strata, by minimizing a cost function subject to precision constraints (or, conversely, by maximizing the precision of the estimates under a given budget). The problem is to determine the best stratification in the frame, i.e., the one that ensures the overall minimal cost of the sample necessary to satisfy precision constraints. The X's can be categorical or continuous; continuous ones can be transformed into categorical ones. The most detailed stratification is given by the Cartesian product of the X's (the atomic strata). A way to determine the best stratification is to explore exhaustively the set of all possible partitions derivable by the set of atomic strata, evaluating each one by calculating the corresponding cost in terms of the sample required to satisfy precision constraints. This is unaffordable in practical situations, where the dimension of the space of the partitions can be very high. Another possible way is to explore the space of partitions with an algorithm that is particularly suitable in such situations: the genetic algorithm. The R package SamplingStrata, based on the use of a genetic algorithm, allows to determine the best stratification for a population frame, i.e., the one that ensures the minimum sample cost necessary to satisfy precision constraints, in a multivariate and multi-domain case.

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:

- 1. analysis of the frame data: identification of available auxiliary information;
- 2. manipulation of auxiliary information: in case auxiliary variables are of the continuous type, they must be transformed into a categorical form;
- 3. 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);
- 5. choice of the precision constraints for each target estimate, possibly differentiated by domain;
- 6. optimization of stratification and determination of required sample size and allocation in order to satisfy precision constraints on target estimates;
- 7. analysis of the resulting optimized strata;
- 8. association of new labels to sampling frame units, each of them indicating the new strata resulting by the optimal aggregation of the atomic strata;

- 9. selection of units from the sampling frame with a stratified random sample selection scheme;
- 10. evaluation of the found optimal solution in terms of expected precision and bias.

Information

Status:	validated
Author:	Istat
Licence:	GPL-2 GPL-3
GSBPM code:	2.4 Design frame and sample4.1 Create frame and select sample
Programming language:	R
Keywords:	optimal stratification, sample design, sample allocation, genetic algorithm
Contact:	name: Giulio Barcaroli email: <u>barcarol@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

R (version \geq 2.15.0)

COPYRIGHT

Copyright 2016 Istat

Licensed under the GNU General Public License (GPL), version 2 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://www.gnu.org/licenses/licenses.en.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 1.1 Windows binaries
- Version 1.1 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – SamplingStrata v. 1.1

Vignettes – SamplingStrata v. 1.1

OTHER DOCUMENTATION

Barcaroli G. 2014. <u>SamplingStrata: An R Package for the Optimization of Stratified Sampling</u>. *Journal of Statistical Software*, 61(4):1-24.

Ballin M., Barcaroli G. 2013. <u>Joint determination of optimal stratification and sample allocation</u> <u>using genetic algorithm</u>. *Survey Methodology*, 39(2):369-393.

COLLECT phase

Create frame and select sample

The activities related to the creation of the list and the selection of the sample refer to subprocess 4.1 *"Create frame and select sample"* of <u>GSBPM</u>.

The creation of the frame consists in the construction of the archive of units belonging to the target population. The selection of the sample consists in the identification of the sample units on the basis of the sampling scheme.

For a given iteration of the collection, the creation of the frame and the selection of the sample are made on the basis of the specifications defined in the sub-process 2.4 "*Design frame and sample*".

Data collection

Collection of data comes after a quite complex set of activities aimed at designing the survey questionnaire used to observe the different aspects of the phenomena to be measured.

What above stated is represented in the <u>GSBPM</u> by a set of sub-processes that belong to different phases of the model, from phase 1 to phase 4, as described here after.

Phase 1 "Specify needs". Sub-processes 1.1 to 1.5 are involved for this phase. They allow to identify the survey needs and to translate them into concepts that have to be comprehensible to respondents. Besides, concepts should be easily measurable and therefore convertible in statistical variables that will be designed in the following Phase 2. At this stage, it is very important to check whether current data sources (for example administrative data) could meet survey needs in order to reduce the number of variables to collect. In this way both cost and response burden can be reduced.

Phase 2: "Design". Once Phase 1 is over, activities described in sub-processes 2.1, 2.2 and 2.3 can be carried on. Survey variables, specified in Phase 1, are used to design the tabulation plan (also useful for the Dissemination phase) that will also allow to define derived variables as well as any statistical classifications that will be used in the collection phase. At this point, the creation of the questionnaire can start: statistical variables can be "translated" into survey questions, relations among variables and variable's characteristics into questionnaire's paths and/or checking rules. Institutional metadata system should be taken into account when designing survey variables, in order to use existing definitions or to update the system with new ones. This will foster national and international standardisation processes and the re-use of any "element" from previous or similar surveys.

Variables design (sub-process 2.2) should be run in parallel with sub-process 2.3 "*Design Collection*" that determines the most appropriate collection method(s) and instrument(s). This is because variables design and, therefore, questions' structure and wording, strictly depend on the mode used to collect data. The same is true for the design of checking rules in case computer assisted techniques are used (example: CATI, CAPI, CAWI, described in Phase 3). With these modes, data can be validated while they dare collected. It is therefore advisable to design checking rules in such a way to solve the greatest number of inconsistencies (major or more frequent inconsistencies), paying attention, at the same time, to the fluency of the interview. In other words, a balance between data quality and respondent burden should be respected: a too

high number of checking rules can increase response burden and negatively affect total or partial non response rate.

Phase 3 "Build": The sub-process 3.1 "Build collection instruments" describes the activities to build the collection instruments to be used during the next Phase 4 "Collect". The collection instrument is generated or built based on the design specifications created during the "Design" phase. Data collection may use one or more modes to receive the data, e.g., personal or telephone interviews, paper, electronic or web questionnaires. Collection instruments may also be data extraction routines used to gather data from existing data sources. In this last case instruments of data exchange like EDI (Electronic Data Interchange) or XBRL (eXtensible Business Reporting Language) can be used for information exchange between Istat and reporting units (enterprises or public organizations that provide data in the form of administrative data sets or registers). This sub-process also includes preparing and testing the contents and functioning of that instrument, e.g., testing the questions in a questionnaire. It is recommended to consider the direct connection of collection instruments to the statistical metadata system, so that metadata can be more easily captured in the collection phase. Connection of metadata and data at the point of capture can save work in later phases, like for instance the Dissemination one. Capturing the metrics of data collection (paradata) is also an important consideration in this sub-process.

Computer assisted techniques play an increasingly important role in the data collection phase. The main features of these modes, **CADI** (*Computer Assisted Data Imputing*), **CAPI** (*Computer Assisted Personal Interviewing*), **CATI** (*Computer Assisted Telephone Interviewing*), **CAWI** (*Computer Assisted Web Interviewing*) is represented by the possibility of performing the editing phase during the data collection, allowing the collection of only valid data. **CADI** differs from the others techniques because checking rules are only used to limit keying errors or to support the revision during the data entry of paper questionnaires.

Another distinctive feature of all computer assisted techniques, except **CADI**, is the customisation of the electronic questionnaire: questions wording can be personalised according to respondent's characteristics (name, gender) or to answers to previously asked questions or to already available information (previous wave of the survey). In this way, the questionnaire appears more friendly and the respondents cooperation is enhanced. How many and what type of checking rules can be implemented in the electronic questionnaire depends on the technique. Main differences are between CATI/CAPI and CAWI:

- **CATI** and **CAPI** are interviewer administered and, therefore, it is possible to use a greater number of checks than for CAWI. Besides, relying on the fact that the interviewer is well trained on both technical aspects and survey content, it is also possible to use a greater number of blocking checking rules, that require to solve the inconsistency before proceeding with the interview;
- As to CAWI, which is instead self-administered, a smaller number of checks should be used. Besides, checking rules should be more like warnings than blocking checks, because they should advise respondents about possible inconsistencies among data, letting them free to solve or not the "error" before proceeding with new questions. This is to avoid the abandoning of the interview before its very end.

It is important to remind here, that the set of checking rules should be such to guarantee a balance between data quality and response burden, whatever technique is used.

Besides, at this point, it is also useful to remind, that the questionnaire design phase has to take into account the technique chosen for collection of data, in order to design the electronic questionnaire in the most suitable way. The electronic questionnaire should be deeply tested to check for its compliance with technical specifications and for its usability and fluency.

Phase 4 "*Collect*": Once collection instruments have been built and tested, the collection of data can start. Activities described in sub-processes 4.2 "*Set up collection*", 4.3 "*Run collection*" and 4.4 "*Finalise collection*" are involved: from interviewers training to the storage of collected data into suitable electronic environment for further processing (Phase 5 "*Process*").

COLLECT phase – METHODS

• CREATE FRAME AND SELECT SAMPLE

MEMOBUST – Handbook on Methodology of Modern Business Statistics: <u>Statistical Registers and Frames</u> <u>Sample Selection</u> 2014 MEETS ESSnet MEMOBUST

• DATA COLLECTION

La modernizzazione delle indagini via Web sulle imprese - Pratiche Raccomandate per il disegno dei questionari 2015 Gdl Armonizzazione dei questionari di impresa, Istat MEMOBUST – Handbook on Methodology of Modern Business Statistics:

Questionnaire Design Data Collection 2014 MEETS ESSnet MEMOBUST

Handbook of Recommended Practices for Questionnaire Development and Testing in the European Statistical System 2006 European project QDet

References

Murgia M., A. Nunnari. 2012. <u>Improving the quality of data collection: minimum requirements</u> for generalised software independent from the mode. *Seminar on New Frontiers for Statistical Data Collection*, UNECE, Ginevra, 31 October-2 November.

Istat. 2008. *La progettazione e lo sviluppo informatico del sistema Capi sulle forze di lavoro*. Collana Metodi e Norme, n. 36, Istat.

Istat. 2006. <u>L'indagine campionaria sulle nascite: obiettivi, metodologia e organizzazione</u>. Collana Metodi e Norme, n. 28, Istat.

COLLECT phase – TOOLS

• DATA COLLECTION

<u>Blaise</u>

Computer Assisted Data Collection Software.

<u>Blaise</u>

Description

Blaise is a software tool for computer assisted data collection developed by Statistics Netherlands. It consists of several modules to perform data collections functions:

- the scheduling of telephone contacts for CATI surveys (CATI Call Management System)
- the questionnaire administration (Data Entry)
- the assisted coding of free text answers (*Coding*)

The calls scheduling is managed by an ad hoc module called *CATI Call Management System*. Generally speaking, it selects the telephone numbers (sample units) to be called in each survey day according to a call priority that depends on previous contact results and on the values of the scheduling parameters settled by the survey manager.

Scheduling parameters indicate how many times each phone number has to be tried before abandoning it. In particular, they indicate how many times not definitive contact results (busy, no one answer, appointment) have to be tried during the survey period before they get a definitive contact results (interviews, refusal, definitive interruptions).

The *Data Entry* module is an user interface for the questionnaire administration suitable for any kind of computer assisted data collection techniques. It allows to manage the screen layout, any question format and to organise the questionnaire in sections. Besides, it is possible to use interactive editing to manage consistency, coherency and routing checks. Error messages can be customised thus helping the interviewer and/or the respondent in filling in the questionnaire.

The assisted *Coding* module, based on sub-strings algorithms, supports three kinds of coding:

- Hierarchical coding: it looks for the code following the tree structure of the classification, from the highest levels (branches) to the lowest ones (leaves). It can be used only for hierarchical classifications.
- Alphabetical or trigram coding: it looks for the code directly using the input text, making partial matches based on sub-strings of texts called bigram or trigrams or based on the alphabetical order.
- Mixed coding: it looks for code making an alphabetical or trigram search inside the selected classification level.

Information

Status:	validated
Author:	Statistics Netherlands
GSBPM code:	3.1 Build collection instrument4.2 Set up collection4.3 Run collection4.4 Finalise collection5.2 Classify and code
Keywords:	CAPI, CATI, CADI, scheduler, assisted coding

Software and documentation

To get the Blaise software as well as the methodological and technical documentation, please contact <u>Statistics Netherlands</u>.

OTHER DOCUMENTATION

Degortes M., S. Macchia, M. Murgia. 2003. <u>The usage of BLAISE in an integrated application for</u> <u>the Births Sample Survey in CATI mode</u>. In *Proceedings of the 8th International BLAISE Users Conference (IBUC 2003)*, International Blaise Users Group, Copenhagen, 21-23 May 2003.

PROCESS phase

Data integration

Record linkage

Record linkage is an important process for the integration of data coming from different sources. The purpose of record linkage is to identify the same real world entity that can be differently represented in data sources, even if unique identifiers are not available or are affected by errors.

In official statistics, record linkage is needed for several applications: for instance,

- to enrich the information stored in different data sets;
- to create, update and de-duplicate a frame;
- to improve the data quality of a source;
- to measure a population amount by capture-recapture method;
- to check the confidentiality of public-use microdata.

The complexity of the whole linking process relies on several aspects. The lack of unique identifiers requires sophisticated statistical procedures, the huge amount of data to process involves complex IT solutions, constraints related to a specific application may require the solution of difficult linear programming problems.

Statistical matching

The goal of statistical matching (sometimes named as data fusion) is the integration of two or more data sources referring to the same population with the aim of exploring the relationships between variables that are not jointly observed in the same data source. The sources to be integrated are composed of different non-overlapping units; as usually happens when data from several sample surveys are integrated. The typical situation of statistical matching is the one in which there are two data sources A and B; variables X and Y are available in A, variables X and Z are observed in B; the objective is to study the relationship between Y and Z by exploiting the common information in X. The objective of statistical matching can be macro or micro; in the first case the interest is in one or more parameters that summarize the relationship between Y and Z (correlation coefficient, regression coefficient, contingency table, *etc.*); in the second case the result of integration is a synthetic data set in which all the variables of interest, X, Y and Z are present.

The objectives of matching can be achieved by means of a parametric or nonparametric approach, or a mixture of them (mixed methods).

The parametric approach requires the specification of a model and the estimation of the related parameters. In absence of auxiliary information, it is generally assumed the conditional independence of Y and Z given the common variables X. This assumption is rather strong and unfortunately in the typical situation of the matching it is not testable.

Nonparametric methods are usually applied when you have a micro objective. In this case hotdeck imputation methods are frequently used. They aims at imputing missing variables in the data set chosen as recipient (for instance A) by using the observed values in the data set (B) chosen as donor. The donor unit for a given unit in A is the most similar observation in B in terms of the values of the common variables X. The mixed approach is composed of two steps: 1) a parametric model is assumed, parameter estimation is performed, and imputation is carried out; 2) a hot-deck imputation procedure is applied, it makes use of imputed values for choosing the donor observation.

It is worthwhile noting that an alternative approach based on the quantification of the uncertainty inherent the estimation of a particular parameter can be used. This approach does not require the conditional independence assumption nor of auxiliary information on non-identifiable parameters, i.e., those related to relations between Y and Z. The study of uncertainty does not lead to a point estimate but to a set of plausible estimates. The set of parameter estimates referring to Y and Z is composed of those consistent with the estimates obtainable from the data at hand, i.e., parameters concerning the couples (Y, X) and (Z, X).

The application of the matching data from complex surveys poses additional problems. In such circumstances it is required to take into account the sampling design as well as the other methodologies used to deal with nonsampling errors (coverage and non-response).

Coding of textual answers

The coding activity is generally performed in case the survey questionnaire contains textual variables that refer to official classifications that allow for national and/or international data comparability. Example of this kind of variables are Economic Activity (NACE), Occupation, Education, Places (of birth, of residence, *etc.*).

Coding means to assign a unique code to a textual answer according to a classification scheme. The level of detail of the matched code depends on the survey aims and/or the dissemination needs. Coding can be performed manually or trough automated systems. Manual coding can be performed only at the end of the data collection phase, while if automated systems are used, it can be run during or after collection of data: in the first case it is called assisted coding (on-line coding) while in second case automated coding (batch coding).

With reference to <u>GSBPM</u>, coding belongs to the sub-process 5.2 "*Classify and code*" of the Phase 5 "*Process*" that includes those activities that are necessary to make data ready for the analysis (Phase 6 "*Analyse*"). Obviously, in case of assisted coding some of the activities of sub-process 5.2 can start before Phase 4 "Collect" ends, improving the timeliness of data delivery.

Coding is, in general, a very hard activity of the survey process. Besides, if it is manually performed it is also difficult to standardise, because coding results strictly depend on coders. Despite coders are well trained about criteria and principles of each official classification, coding is influenced by the cognitive process of each coder that might lead to different (subjective) interpretations and, therefore, different coding of the same textual answer.

The use of specialised coding software can produce a considerable saving of time and resources and will also guarantee a higher standardisation level of the coding process, increasing the expected quality of the coding results.

As already said, computer assisted coding can be distinguished in "automatic coding" and in "assisted coding". They differs in terms of aims and coding process:

• Automatic coding: the coding software analyses and codes, on the basis of a reference dictionary, a data file containing all the textual answers collected during the collection

phase (batch coding). The aim is to look for and to assign a single code to each textual answer according to quality thresholds;

• Assisted coding: the coding software is an interactive instrument, that aids the coder/respondent in coding the textual answer. The aim is to offer the user a wider set of possible matching codes among which to choose the correct one.

The key point of any coding system, automated or assisted, is the implementation of the informative basis that represents the reference dictionary containing codes and texts of the official classification and enriched with textual answers collected by Istat surveys (and correctly coded). In order to be processed by a software, the reference dictionary has to undergo a number a standardisation operations aimed at producing analytic, synthetic and not ambiguous descriptions. Besides, in general, the richest the dictionary the higher the coding rate.

Generally speaking, coding systems varying according to the algorithm used to match the textual answers with the dictionary descriptions. They can be classified as follows:

- dictionary algorithms: they look for exact matches on the bases of key words (or groups of key words);
- weighting algorithms: they look for partial or exact matches on the basis of similarity functions among texts that assign weights to each word according to its informative content;
- sub-strings algorithms: they look for partial or exact matches processing portions of texts (bigrams or trigrams).

Besides, for what concern assisted coding, there are three possible methods to consult (to navigate) the reference dictionary:

- tree search: it navigates inside the classification hierarchical structure, from the higher branch to the lowest one (leave) that represents the most detailed code (highest number of digits) that can be assigned to a textual answer;
- alphabetic search: it navigates inside the entire dictionary looking for the definition which is equal or the most similar to the textual answer to be coded;
- mixed mode search: it makes an alphabetic search inside the selected classification branch.

Data collection technique highly influences the choice of the searching method. A special distinction is among interviewer administered and self-administered modes. For the latter, where respondents are not trained on classifications and coding like interviewers are, it is extremely important to provide a coding system that is user friendly and guarantees high quality results.

The quality of coding activity is highly influenced by the update of both the dictionary content and the matching rules (training phase). It is advisable to perform the training phase periodically, in general after the coding of textual answers collected by a survey. To this aim, after a coding application, it is important to:

- verify the quality of the coded cases;
- use the not coded cases to update the coding application (dictionary and checking rules);
- highlight eventual lacks of the classification used.

Per la valutazione della qualità delle due modalità di codifica, è possibile utilizzare i seguenti indicatori:

Indicators for assisted and automated coding can be used to evaluate the performance of the coding phase:

Automated coding indicators:

- efficacy/coding rate: ratio of "number of coded texts" to "total number of texts to be coded";
- accuracy: ratio of "number of correctly coded texts" to "number of coded texts";
- efficiency: unitary coding time.

Assisted coding indicators:

- average time to assign a single code;
- coherence among each collected textual description and the assigned code.

Detection and treatment of measurement errors and imputation of partial nonresponses

Partial non-responses (PNR) and measurement errors are specific types of non-sampling errors which are generally treated in the editing and imputation phase. In the present context, a measurement error is defined as a discrepancy between "true value" and observed value of a variable in a statistical unit. Discrepancy could be originated in any phase of the measurement process (data collection, coding, storing, *etc.*). Since NRPs and measurement errors may seriously compromise the accuracy of the target estimates, they should be prevented by means of suitable strategies. However, even though efforts are made to limit the impact of non-sampling errors, a proportion of collected data are typically affected by partial non-responses and measurement errors. Thus, use of editing and imputation methods are necessary.

In the context of <u>GSBPM</u>, the two sub-processes 5.3 "Review and validate" e 5.4 "Edit and impute" concern input data validation and editing and imputation respectively. However, often in the real contexts, these two sub-processes are hardly distinguishable.

Detection of measurement errors

Methods for error detection can be classified according to the different error typologies. A first important distinction is between **systematic errors** and **random errors**.

An error is called **systematic** if it depends on structural problems in the statistical production process such as defects in the questionnaire design or in the storing system. They usually determine deviation from the true values "in the same direction" for one or more variables of interest. Systematic errors are generally treated by means of deterministic rules base on the knowledge of the error mechanism. A typical systematic error for quantitative variables is the unity measure error.

An error is called **random** if it depends on stochastic elements that are not identifiable. In contrast with the case of systematic errors, a deterministic approach is generally not appropriate for random errors.

An important typology of errors involves errors causing "**out of range**" values, that is, values that do not belong to a known set of acceptable values. A similar typology involves errors determining **inconsistencies** between different observed variables. These errors are detected by applying **coherence rules (edit)** on data. Consistency errors that are believed to be "not influent" are generally detected via automatic methods based on some "general" principles. In this context, a popular approach is the **minimum change principle**. According to this principle, erroneous values are detected in such a way that for each record the minimum number of items is to be changed in order for the record to pass all the edits. The **Fellegi-Holt** methodology implements the minim change principle. It has been originally developed for categorical variables and later extended to numerical variables.

Another important class of errors includes errors which determine **outliers** in data. Outliers are observations whose behavior deviates from the one typical for most observations. Usually, technics for outlier identifications assume, at least implicitly, a model for the "typical" data and try to identify deviations from the model. Methods for outlier detection are often used to identify **influential errors**. However, the concept of influential error is to be distinguished from the one of outlier. In fact the latter is related to a model assumed for the data, while the former depends on the population parameter to be estimated. Specifically, outlier may not be caused by influential errors, and influential errors may not cause outliers.

Correction of errors and imputation of partial non-responses

Once errors have been detected, incorrect values have to be replaced (imputed) by correct or nearly correct values. Imputation is also used for partial non-responses. Imputation is commonly used for a number of technical and practical reasons. First, released data must be complete and coherent at micro level. Furthermore, imputation allows the users using standard methods and software on the final dataset.

A lot of imputation methods are available in many statistical packages both proprietary and free. A frequent classification of imputation methods is in terms of **parametric methods** (e.g., **regression** type methods), relying on some explicit assumptions on the probability distribution generating the data, and **non-parametric methods**, generally based on weaker not explicit assumptions (e.g., **hot-deck, nearest neighbor donor**).

Sometimes imputation methods are called **deterministic** when different applications on the same dataset provide same outputs, and stochastic, if the outcomes are characterized by a certain level of variability.

The choice of the imputation method depends on the analyses that have to be performed on data. For instance, If the quantity of interest is a linear quantity such as a mean or a total, a deterministic method may be appropriate, while, if also are distributional characteristics (such as the ones involving the second moment of the data distribution), a stochastic imputation method is generally preferable.

Weighting, estimation and sampling error evaluation

The activities concerning the production of target estimates and the evaluation of sampling errors refer to sub-process 5.6 "*Calculate weights*" and 5.7 "*Calculate aggregates*" of <u>GSBPM</u>.

Production of the estimates of interest

Each estimation method is based on the principle that the subset of the units of the population included in the sample must also represent the complementary subset consisting of the remaining units of the population. This principle is generally achieved by assigning to each unit included in the sample a weight that can be seen as the number of population elements represented by that unit.

The sample surveys carried out by Istat are large-scale surveys that have the purpose of providing a large number of estimates of population parameters such as counts, totals, proportions, averages, etc.

The estimation of the parameters of the population can be made using two different approaches:

- Methods based on the direct approach using values of the variable of interest observed on the sample units belonging to the domain of interest. They are the standard methods used by Istat and by all the main National Institutes of Statistics to produce survey estimates.
- Methods based on the indirect approach that make use values of the variable of interest observed on sample units belonging to a wider domain containing the domain of interest and/or other survey occasions. They are, usually, used for particular estimation problems, such as those associated with the generation of estimates referring to domains in which the sample size is too small for the production of estimates using direct methods.

Direct methods

In general, for the estimation a total the following two operations should be performed:

- 1. computation of the weight to be assigned to each unit included in the sample;
- 2. calculation of the estimates as weighted sums of the values of the target variables using with weights determined in step 1.

The weight given to each unit is obtained according to a procedure divided in several steps:

- 1. the *starting weight* of each sample unit, named *direct weight*, is calculated according to the sampling design, as the reciprocal of the inclusion probability;
- 2. the starting weight is adjusted in order to account for non-response, obtaining the *base weight*;
- correction factors of the base weight basis are computed to take into account equality constraints between some known parameters of the population and the corresponding sample estimates;
- 4. the *final weight* is obtained as the product between the base weight and the correction factors.

The class of estimators corresponding to the operations described above is known as *calibration estimators*, since both the adjustment to correct for non-response and the weight correction to achieve consistency with known population parameters is obtained by solving a constrained minimization problem. In details we want to minimize the distance between the weight before and after the calibration phase.

The main problem for the choice of the estimation method is to find an estimator satisfying:

- *efficiency* criteria in terms of sample variance and bias due to the presence of *total and partial non-responses*, frame under-coverage;
- external and internal coherence. The external consistency of the estimates arises
 whenever known totals are available from external sources. Estimates of the total
 produced by the survey should generally match or not deviate too much from the known
 values of these totals. The internal consistency of the estimates is achieved when all the
 estimates of the same aggregate coincide with each other. This result can be obtained
 using a unique system of weights.

Calibration estimators meet the above criteria since:

- they yield, generally, more efficient estimates than those obtained using direct estimators; the higher the correlation between the auxiliary variables and the target variables the greater the efficiency;
- they are approximately design unbiased;
- they produce estimates of totals that coincide with the known values of these totals;
- they mitigate the non-response bias effect;
- they reduce the bias due to the under-coverage of the frame from which the sample is selected.

Calibration estimators are used for calculating the final weights for most social and business surveys carried out by Istat.

Indirect methods

Indirect estimation methods are used by Istat to give clear responses to the growing need of local governments for accurate information for small geographical areas, or more generally domains, called small areas. Sample surveys conducted by Istat are, however, designed to provide reliable information for the main aggregates of interest for planned domains defined at design stage and may not be able to respond appropriately to the production of estimates for larger level of detail.

The solution adopted in the past by Istat to obtain estimates at unplanned domain level, was to increase the sample size without changing the sampling strategy, i.e. without modifying neither the sampling design nor the estimator. However over-sampling, besides rising collection costs and increasing the difficulty of organizational issues, implies the increase of non-sampling errors due to the difficulty to handle with too large sample sizes. In addition, the over-sampling is a partial solution to the problem of small area estimation, since not being able to increase the sample size over a certain threshold makes it possible to provide reliable estimates only for a subset of the small areas of interest.

For these reasons, Istat makes use of estimation methods based on:

- the use of auxiliary information, related to the phenomena under study, known at small areas level;
- the adoption (implicit or explicit) of statistical models linking the variable of interest observed in the small area with the values of the same variable related to a larger area containing the small area of interest and/or related to other survey occasions.

An important problem related with these methods is that they are based on models and, therefore, the properties of the results depend on the validity of the model assumptions. Since models are never expected to perfectly match reality, these estimators introduce unmeasurable bias which may arise serious concerns about their use.

Assessment of sampling errors

For the evaluation of the sampling errors of the estimates, Istat usually uses approximated variance estimation methods. In fact, for most of the estimation procedure an analytical expression of the estimator variance is not available, since:

- ISTAT surveys are carried out using complex sampling designs, generally based on multiple selection stages, on the stratification of the units, and on without repetition selection scheme with varying selection probabilities among units;
- estimates are determined through the use of calibration estimators which are non-linear functions of the sample information.

The estimation methods of the sample variance generally used in Istat are based on the method of linearization of Woodruff (1971) which provides estimates of the sample variance in the case where the estimators used are non-linear functions of the sample data.

This variance estimation methods are implemented by Istat in the generalized software GENESEES and ReGenesees, which feature a user friendly interface and are currently used to estimate the sampling errors of the estimates produced by Istat surveys.

In addition, by means of these software, it is possible to compute important analysis statistics which provide useful tools to evaluate the adopted sampling design. In particular it is possible to evaluate:

- the overall efficiency of the sampling design;
- the impact on the efficiency of the estimates due to stratification, the number and type of selection stages and weighting effect.

It is important to point out that a generalized variance errors representation is also provided. It is a summary obtained making use of regression models that relate the estimates with the corresponding sampling errors. These models provide important summary information on sampling errors and are disseminated together with the tables reporting the estimate values.

PROCESS phase – METHODS

• DATA INTEGRATION

MEMOBUST – Handbook on Methodology of Modern Business Statistics <u>Micro-Fusion</u> 2014 MEETS ESSnet MEMOBUST State of the art on statistical methodologies for data integration

Methodological developments 2011 ESSnet on Data Integration

Old and new approaches in statistical matching when samples are drawn with complex survey designs 2010 in Proceedings of the SIS Conference, Padua

ISAD Work packages and executive summary

2008

ESSnet Statistical Methodology – Area ISAD (Integration of Survey and Administrative Data)

Metodi statistici per il record linkage

2005 Collana Metodi e Norme, n. 16, Istat

References

D'Orazio M., M. Di Zio, M. Scanu. 2006. <u>Statistical Matching for Categorical Data: Displaying</u> <u>Uncertainty and Using Logical Constraints</u>. *JOS*, 22(1):137-157.

D'Orazio M., M. Di Zio, M. Scanu. 2006. *Statistical Matching: Theory and Practice*. J. Wiley & Sons, Chichester.

CODING OF TEXTUAL ANSWERS

Metodi e software per la codifica automatica e assistita dei dati 2007 Collana Tecniche e strumenti, n.4, Istat

References

Istat. 2008. <u>L'ambiente di codifica automatica dell'ATECO 2007. Esperienze effettuate e</u> <u>prospettive</u>. Collana Metodi e Norme, n. 41, Istat.

Istat. 2005. <u>Una soluzione per la rilevazione e codifica della Professione nelle indagini CATI</u>. Collana Contributi Istat, n. 11, Istat.

Macchia S., M. D'Orazio. 2001. <u>A system to monitor the quality of automated coding of textual</u> <u>answers to open questions</u>. *Research in Official Statistics*, 4(2).

DETECTION AND TREATMENT OF MEASUREMENT ERRORS AND IMPUTATION OF PARTIAL NON-RESPONSES

MEMOBUST – Handbook on Methodology of Modern Business Statistics Statistical Data Editing Imputation 2014 MEETS ESSnet MEMOBUST A Contamination Model for Selective Editing 2013 Journal of Official Statistics. Volume 29, Issue 4, Pages 539-555.

Recommended Practices for Editing and Imputation in Cross-Sectional Business Surveys 2007 European project EDIMBUS Results of the EUREDIT project

Euredit – The Development and Evaluation of New Methods for Editing and Imputation 2003 European project EUREDIT

References

De Wall T., J. Pannekoek and S. Sholtus. 2011. *Handbook of Statistical Data Editing and Imputation*. J. Wiley & Sons, Hoboken, N.J.

Di Zio M., U. Guarnera and R. Rocci. 2007. <u>A mixture of mixture models for a classification</u> problem: the unity measure error. *Computational Statistics and Data Analysis*, 51:2573-2585.

Di Zio M., U. Guarnera and O. Luzi. 2005. <u>Editing Systematic Unity Measure Errors Through</u> <u>Mixture Modelling</u>. *Survey Methodology*, 31(1):53-63.

Pannekoek J. and T. De Waal. 2005. <u>Automatic Edit and Imputation for Business Surveys: The</u> <u>Dutch Contribution to the EUREDIT Project</u>. *Journal of Official Statistics*, 21(2):257-286.

Särndal C. E. and S. Lundström. 2005. *Estimation in Surveys with Nonresponse*. J. Wiley & Sons, New York.

Chambers R., A. Hentges and X. Zhao. 2004. <u>Robust automatic methods for outlier and error</u> <u>detection</u>. *Journal of the Royal Statistical Society*, Series A, 167(2):323-339.

Little J. and D. Rubin. 2002. Statistical Analysis with Missing Data. J. Wiley & Sons, New York.

Chen J. and J. Shao. 2000. <u>Nearest Neighbor Imputation for Survey Data</u>. *Journal of Official Statistics*, 16(2):113-131.

Schafer J. L. 2000. Analysis of Incomplete Multivariate Data. Chapmann and Hall/CRC, New York.

Latouche M. and J.M. Berthelot. 1992. <u>Use of a Score Function to Prioritize and Limit Recontacts</u> <u>in Editing Business Surveys</u>. *Journal of Official Statistics*, 8(3):389-400.

Little R.J.A. 1988. Missing-data adjustments in large surveys. *Journal of Business & Economic Statistics*, 6(3):287-296.

Rubin D. 1987. *Multiple Imputation for Nonresponse in Surveys*. J. Wiley & Sons, New York.

Hidiroglou M. A. and J.M. Berthelot. 1986. Statistical editing and imputation for periodic business surveys. *Survey Methodology*, 12(1):73-83.

Kalton G. and D. Kasprzyk. 1982. <u>Imputing for missing survey responses</u>. In *Proceedings of the section on Survey Research Methods*, American Statistical Association.

Fellegi P. I. and D. Holt. 1976. A systematic approach to automatic edit and imputation. *Journal of the American Statistical Association*, Applications Section, 71:17-35.

• WEIGHTING, ESTIMATION AND SAMPLING ERROR EVALUATION

MEMOBUST – Handbook on Methodology of Modern Business Statistics Weighting and Estimation 2014 MEETS ESSnet MEMOBUST Results of the SAE project SAE – Small Area Estimation 2012 ESSnet SAE Riponderazione 2005 Note metodologiche, Istat

<u>Stime ed Errori</u> 2005 Note metodologiche, Istat

References

Istat. 2008. Strategia di campionamento e precisione delle stime. In "<u>L'indagine europea sui</u> <u>redditi e le condizioni di vita delle famiglie (Eu-Silc)</u>", Collana Metodi e Norme, n. 37, Istat.

Istat. 2006. La procedura di stima e la valutazione degli errori campionari. In "<u>Il sistema di</u> <u>indagini sociali multiscopo. Contenuti e metodologia delle indagini</u>", Collana Metodi e Norme, n. 31, Istat.

Istat. 2006. Strategia di campionamento e livello di precisione delle stime. in "<u>L'indagine</u> <u>campionaria sulle nascite: obiettivi, metodologia e organizzazione</u>", Collana Metodi e Norme, n. 28, Istat.

Rao, J. N. K. 2003. Small Area Estimation. J. Wiley & Sons, New York.

Cicchitelli G., A. Herzel, G.E. Montanari. 1992. Il campionamento statistico. Il Mulino, Bologna.

Deville J.C., C.E. Särndal. 1992. Calibration Estimators in Survey Sampling. *Journal of the American Statistical Association*, 87:376-382.

Särndal C.E., B. Swensson, J. Wretman. 1992. *Model Assisted Survey Sampling*. Springer-Verlag, New York.

Cochran W.G. 1977. *Sampling Techniques*. J. Wiley & Sons, New York.

PROCESS phase – TOOLS

• DATA INTEGRATION

<u>RELAIS</u> (REcord Linkage At IStat)

Toolkit for dealing with record linkage projects.

StatMatch

R package for data integration through statistical matching and, as a by-product, the ability to impute missing values in a data set.

• CODING OF TEXTUAL ANSWERS

<u>CIRCE</u> (Comprehensive Istat R Coding Environment)

Automated coding system for textual answers.

DETECTION AND TREATMENT OF MEASUREMENT ERRORS AND IMPUTATION OF PARTIAL NON-RESPONSES

<u>Banff</u>

Edit and imputation system for numeric and continuous variables.

CANCEIS (CANadian Census Edit and Imputation System)

Edit and imputation system based on the Nearest-neighbour Imputation Methodology (NIM). The NIM allows the simultaneous hot-deck imputation of numeric and categorical variables based on a single donor. The basic entity that is processed (the unit) can be composed of one or more sub-units of lower hierarchical level.

<u>CONCORDJava</u> (CONtrollo e CORrezione dei Dati version with Java interface) Integrated system to check and correct the data (imputation). One of the modules (SCIA) implements the Fellegi-Holt Methodology to treat the inconsistencies between the values of qualitative variables.

<u>SeleMix</u> (Selective editing via Mixture models)

R package to treat quantitative data, which aims to identify a set of units affected by errors which potentially influence the interest estimates (selective editing).

• WEIGHTING, ESTIMATION AND SAMPLING ERROR EVALUATION

EVER (Estimation of Variance by Efficient Replication)

R package for calibration, estimation and sampling error assessment in complex sample surveys, based on replication methods .

<u>ReGenesees</u> (R evolved Generalised software for sampling estimates and errors in surveys) R-based software system for design-based and model-assisted analysis of complex sample surveys.

RELAIS (REcord Linkage At IStat)

Description

RELAIS (REcord Linkage At IStat) is a toolkit providing a set of techniques for dealing with record linkage projects.

The purpose of record linkage is to identify the same real world entity that can be differently represented in data sources, even if unique identifiers are not available or are affected by errors. In statistics, record linkage is needed for several applications, including: enriching the information stored in different data-sets; de-duplicating data-sets; improving the data quality of a source; measuring a population amount by capture-recapture method; checking the confidentiality of public-use micro data. In fact, record linkage areas; moreover, several different techniques can be adopted for each phase. We believe that the choice of the most appropriate technique not only depends on the practitioner's skill but, most of all, it is application specific.

Moreover, in some applications, there is no evidence to prefer a given method to others or of the fact that different choices, at some linkage stage, could bring to the same results. This is why it is reasonable to dynamically select the most appropriate technique for each phase and to combine the selected techniques for building a record linkage work-flow of a given application. RELAIS is a toolkit relying on these ideas.

The principal features of RELAIS are:

- It is designed and developed to allow the combination of different techniques for each of the record linkage phases, so that the resulting work-flow is actually built on the basis of application and data specific requirements.
- It has been developed as an open source project, so several solutions already available for record linkage in the scientific community can be easily re-used. It is released under the EUPL license (European Union Public License).
- It has been implemented by using two languages based on different paradigms: Java, an object-oriented language, and R, a functional language. This choice depends on our belief that a record linkage process is composed of techniques for manipulating data, for which Java is more appropriate, and of calculation-oriented techniques for which R is a preferable choice. The choice of Java and R is also in line with the open source philosophy of the RELAIS project.
- It has been implemented using a relational database architecture, in particular it is based on a MySQL environment that is also in line with the open source philosophy of the RELAIS project.

The RELAIS project aims to provide record linkage techniques easily accessible to not-expert users. Indeed, the developed system has a GUI (Graphical User Interface) that on the one hand permits to build record linkage work-flows with a good flexibility. On the other hand it checks the execution order among the different provided techniques whereas precedence rules must be controlled.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	5.1 Integrate data
Programming language:	R, Java
Language of the GUI:	EN
Keywords:	data integration, probabilistic record linkage, string comparators, blocking/sorting/indexing, deduplication, open source software
Contact:	name: Luca Valentino email: <u>luvalent@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

Java 2 Runtime Environment (version \geq 6.0)

R (≥ 2.5.1)

R packages: <u>lpSolve</u> (\geq 5.5), RODBC

MySQL Server

MySQL ODBC ($\geq 5.x$)

COPYRIGHT

Copyright 2015 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

RELAIS version 3.0

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

User manual – RELAIS v. 3.0

OTHER DOCUMENTATION

Cibella N., G.L. Fernandez, M. Guigò, F. Hernandez, M. Scannapieco, L. Tosco, T. Tuoto. 2009. <u>Sharing Solutions for Record Linkage: the RELAIS Software and the Italian and Spanish</u> <u>Experiences</u>. In *Proceedings of New Techniques and Technologies for Statistics (NTTS) Conference*, Eurostat, Brussels, 18-20 February 2009.

Eurostat. 2009. Theory and practice of developing a record linkage software. In <u>"Insights on</u> <u>Data Integration Methodologies. ESSnet-ISAD workshop, Vienna, 29-30 May 2008</u>". Methodologies and working papers, Eurostat.

Cibella N., M. Fortini, M. Scannapieco, L. Tosco, T. Tuoto. 2007. <u>RELAIS: Don't Get Lost in a</u> <u>Record Linkage Project</u>. In *Proceedings of the FCSM 2007 Conference*, Federal Committee on Statistical Methodology, Arlington, 5-7 November 2007.

Fortini M., P.D. Falorsi, C. Vaccari, N. Cibella, T. Tuoto, M. Scannapieco, L. Tosco. 2006. <u>Towards</u> <u>an Open Source Toolkit for Building Record Linkage Workflows</u>. In *Proceedings of the International Workshop on Information Quality in Information Systems (IQIS),* Chicago, 30 June 2006.

<u>StatMatch</u>

Description

StatMatch is an add-on package for the R environment that functions for data integration through statistical matching and, as a by-product, the opportunity to impute the missing values in the data set.

The package contains functions that implement methods of matching and functions to support matching (calculation of the distances, etc.). There are three functions dedicated to the application of nonparametric methods of matching at the micro level:

- NND.hotdeck: mimimum distance donor selection; implements many distance functions; allows the definition of donation classes. It is possible to introduce a constraint for avoiding the selection of a donor more than once (constrained matching).
- RNDwNND.hotdeck: random selection of the donor in fixed or "moving" classes. In the
 latter case, it is possible to select at random a donor among the k closest ones; random
 choice of a donor among those with a distance less than a given threshold, etc. The
 selection of the donor can be done with probability proportional to a weighing variable.
 Different distance functions are implemented.
- rankNND.hotdeck: select the donor closer based on the distance between the
 percentiles of the empirical cumulative distribution of a continuous variable available in
 both the data set. In estimating the empirical cumulative distribution it is possible to
 account for different weights assigned to units. The empirical cumulative distribution
 may be estimated in given classes of units.

These functions can be used to impute the missing values in a data set through the corresponding hot-deck methods.

There is only one function mixed.mtc that allows to implement methods of matching in parametric macro or mixed framework. The function assumes that the variables *X*, *Y* and *Z* are distributed according to a multivariate normal distribution. The estimation of the parameters can be done via the maximum likelihood method or by applying a method based on sample estimates of the parameters of interest (means and variances).

Two functions implement methods for matching at the macro level in the presence of data from complex surveys. These functions are based on a series of calibrations of the weights associated to the units in the source data sets according to the methods suggested by Renssen (1998). At this time the application of these methods is limited to the case where Y and Z are both categorical and the goal consists in estimating the contingency table Y vs. Z. In particular, the function harmonize.x harmonizes the marginal/joint distribution of the chosen X variables so to make it coherent between the two initial data sets; then the function comb.samples estimates the table Y vs. Z through the methodologies proposed by Renssen, with or without an additional data source C where Y and Z or X, Y and Z are jointly observed.

Finally the functions Frechet.bounds.cat and Fbwidths.by.x allow the exploration of uncertainty when all the variables (X, Y, and Z) are categorical. The first function estimates the ranges of uncertainty bounds for all the cells in the contingency table Y vs. Z. In the presence of many common variables X, the function Fbwidths.by.x permits to identify the variables most related to Y and Z, and that allow a reduction of the width of the uncertainty bounds.

Information

Status:	validated
Author:	Istat
Licence:	<u>GPL-2 GPL-3</u>
GSBPM code:	5.1 Integrate data 5.4 Edit and impute
Programming language:	R
Keywords:	statistical matching, data fusion, hot deck imputation, uncertainty analysis
Contact:	name: Marcello D'Orazio email: <u>madorazi@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

R (≥ 2.7.0)

R packages: proxy, <u>lpSolve</u>, <u>survey</u>.

In some cases, it may be necessary to have the additional R packages: <u>optmatch</u>, <u>SDaA</u>, <u>simPopulation</u>, <u>MASS</u>. It should be noted that the use of the package <u>optmatch</u> is subject to some limitations.

COPYRIGHT

Copyright 2016 Marcello D'Orazio

Licensed under the GNU General Public License (GPL), version 2 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://www.gnu.org/licenses/licenses.en.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 1.2.5 Windows binaries
- Version 1.2.5 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – StatMatch v. 1.2.5

Vignettes – StatMatch v. 1.2.5

CIRCE (Comprehensive Istat R Coding Environment)

Description

CIRCE is a software package based on R aimed at automatically coding textual variables according to official classifications. It is a generalised software system with respect to the language and the classification used. CIRCE replaces Actr v3, that has been adopted by Istat since 1998, because it was no more commercialised and maintained from its producers and because it was not compatible with the software platform used in Istat (Windows 7, Windows Server 2008).

To prevent lower quality results, the same matching algorithms of ACTR v3 has been developed for CIRCE.

Being an R package it is portable to different platform with no need of compilation. This made it possible to have one single package running on both Windows and Linux operating systems. It can be used on Windows environment through an User Graphical Interface and on web through a "call" to a web service.

CIRCE belongs to the weighting algorithms category and manages three types of coding procedures:

- automated coding, for set of records (batch coding);
- interactive coding, to analyse coding results of single record (a GUI is provided to coders);
- web coding, a web service for single record coding. In this case is currently available a web service dedicated to the identification of the activity code (in Italian language) accessible through the page: http://www.istat.it/it/strumenti/definizioni-e-classificazioni/ateco-2007.

Notwithstanding the type of procedure, the coding phase is performed in two consecutive steps:

1) standardization of texts, called parsing;

2) matching of parsed texts.

The parsing step is a quite sophisticated phase of text standardisation totally customisable, that provides (till now) 14 different functions such as characters mapping, deletion of trivial words, definition of synonymous, suffixes removal, etc.. The parsing aims at removing grammatical or syntactical differences in order to make equal two different descriptions but with the same semantic content.

The second step is the matching phase. The parsed response is compared with the parsed descriptions of the informative base. If this search returns a perfect match or direct match, then a unique code is assigned, otherwise the software uses an algorithm to find the best partial matches, providing an indirect match.

CIRCE is developed by Istat. This will make it easier adding or changing its functionalities with respect to standardization parings and/or matching steps.

Please note: for the moment, both CIRCE user guide (Manuale Utente.pdf) and its GUI are in Italian. English versions will eventually be provided in the future.

Information

Status:validatedAuthor:lstatLicence:EUPL-1.1GSBPM code:5.2. Classify and codeProgramming language:R, VB.NETLanguage of the GUI:ITKeywords:automated coding, weighting coding algorithmsContact:name: Laura Capparuc@istat.it

Software and documentation

SOFTWARE DEPENDENCIES

- R (version \geq 3.1.1).
- Windows (version ≥7).
- Microsoft Framework .net 4 (only for the graphical user interface).

COPYRIGHT

Copyright 2016 Istat

Licensed under the European Union Public Licence (EUPL), version 1.0 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

CIRCE version 1.0

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

User manual – CIRCE v. 1.0

OTHER DOCUMENTATION

Istat. 2007. <u>Metodi e software per la codifica automatica dei dati</u>. Collana Tecniche e strumenti, n. 4, Istat.

Istat. 2005. *La codifica delle variabili testuali nel 14° Censimento Generale della Popolazione*. Collana Documenti Istat, n. 1, Istat.

Macchia S., M. D'Orazio. 2001. <u>A system to monitor the quality of automated coding of textual</u> <u>answers to open questions</u>. *Research in Official Statistics*, 4(2).

De Angelis R., S. Macchia, L. Mazza. 2000. Applicazioni sperimentali della codifica automatica: analisi di qualità e confronto con la codifica manuale. <u>Quaderni di ricerca – Rivista di statistica</u> <u>Ufficiale</u>, 1.

<u>Banff</u>

Description

Banff is a generalized system for editing and imputing survey data based on the SAS architecture. Banff was developed by Statistics Canada to process numeric and continuous variables. It uses consistency rules (edit) that must be expressed in linear form.

Banff has a modular structure: each module corresponds to a particular sub-function of the general structure of an edit and imputation process of quantitative variables:

- edit specification;
- check edits for consistency and redundancy;
- error localization;
- detection of outlier values;
- imputation.

The error localization module uses the *Chernikova algorithm* based on the *minimum change principle* or Fellegi-Holt paradigm. For each record that fails at least one edit, the algorithm identifies the minimum number of fields to change (impute) so that the record passes all the rules. In general, the Fellegi-Holt paradigm is considered appropriate to treat stochastic errors.

Banff implements several imputation methods:

• Deterministic imputation

It checks if there is one and only one value that, once assigned to the field to impute, allow the record to pass all the edits.

• Donor imputation

The nearest neighbour record (according to a specific distance function) to the current failed record is chosen among the potential donors, i.e. units that pass all the edits. All required fields are imputed by transferring the corresponding values from the nearest neighbour record.

It is important to note that a potential donor will be actually chosen as the donor, if the imputed values are such that the imputed record pass the user-specified post-imputation edits.

• Estimator Imputation

Values to be imputed are obtained through modeling or observed data. Examples are mean imputation and regression imputation.

Banff also provides outputs that allow the user to analyze the impact of the editing process on the data (for example, the list of redundant rules or the failure frequency of the erroneous record).

Information

Status:	decommissioned
Author:	Statistics Canada
GSBPM code:	5.3 Review and validate 5.4 Edit and impute
Keywords:	editing for numerical variables, error localization, minimum change principle, nearest neighbour donor

Software and documentation

To get the Banff software as well as the methodological and technical documentation, please contact <u>Statistics Canada</u>.

Only for Istat staff: contact Francesco Dell'Orco.

CANCEIS (CANadian Census Edit and Imputation System)

Description

CANCEIS is an edit and imputation system developed by Statistics Canada. It was designed based upon the Nearest-neighbour Imputation Methodology (NIM) and was first used to perform edit and imputation on census data. CANCEIS performs minimum change nearest-neighbour imputation and deterministic imputation. The NIM allows the simultaneous hot-deck imputation of numeric and categorical variables based on a single donor. The NIM identifies donors for the entire household, not only for the individuals. For each failed household, the NIM identifies a set of potential donors (nearest neighbours) which are as similar as possible to the failed household to be imputed. For each nearest neighbour, the smallest subsets of variables which, if imputed, allow the imputed record to pass the edits, are identified. One of those imputation actions is randomly selected by giving a better chance to those imputation actions that are simultaneously closer to the failed household and the potential donor.

Information

Status:	validated
Author:	Statistics Canada
GSBPM code:	5.3 Review and validate 5.4 Edit and impute
Keywords:	editing, nearest-neighbour donor imputation, numeric and categorical variables

Software and documentation

To get the CANCEIS software as well as the methodological and technical documentation, please contact <u>Statistics Canada</u>.

OTHER DOCUMENTATION

Istat. 2007. <u>Indagine sulle Cause di Morte: Nuova procedura automatica per il controllo e la</u> <u>correzione delle variabili demo-sociali</u>. Collana Documenti Istat, n. 6, Istat.

Manzari A., A. Reale. 2001. <u>Towards a new system for edit and imputation of the 2001 Italian</u> <u>Population Census data: A comparison with the Canadian Nearest-neighbour Imputation</u> <u>Methodology</u>. In *ISI World Statistics Congress Proceedings 53rd Session*, International Statistical Institute, Seoul, 2001.

CONCORDJava (CONtrollo e CORrezione dei Dati version with Java interface)

Description

CONCORDJava is an *open source* software for editing and imputing data. The application integrates software previously developed and used in Istat: SCIA, RIDA and GRANADA.

The application, currently released in beta edition for the part about the deterministic corrections, is available for download in Italian and English.

The different methods in the software are implemented in separate modules:

• SCIA (Sistema di Controllo e Imputazione Automatici):

Editing and imputation of qualitative variables by applying the Fellegi-Holt methodology. First the system identifies the minimum number of fields to impute for each failed record and then it imputes ensuring the imputed record will pass all edits.

• **RIDA** (Ricostruzione dell'Informazione con Donazione Automatica):

It performs imputation of both categorical and continuous variables through donor of minimum distance. Preliminary operations are:

- the classification of units in passed and failed;
- their recording into two separate files;
- $\circ~$ the identification of values to be imputed using a preset character (error character).
- GRANADA (Gestione delle Regole per l'ANAlisi dei DAti):

It imputes qualitative and quantitative variables according to the deterministic approach, in other words applying the rules of the type IF [error condition] THEN [corrective action]. Using this module you can also check data according to edits which admit logical and arithmetic operators (and therefore valid for qualitative and quantitative variables).

Preparatory to the various steps is the phase of definition of the variables, i.e. of the fields of the record to be checked, and the edits.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	5.3 Review and validate 5.4 Edit and impute
Programming language:	Fortran, Java
Language of the GUI:	EN, IT
Keywords:	localization, imputation, nearest neighbour donor, Fellegi-Holt
Contact:	name: Maria Teresa Buglielli email: <u>bugliell@istat.it</u>

Software and documentation

Minimum size hardware

256 Mb RAM

30 Mb on C:/

SOFTWARE DEPENDENCIES

Java 2 Runtime Environment (version \geq 6.0)

COPYRIGHT

Copyright 2014 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

CONCORDJAVA version 2.2

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

User manual – CONCORDJava v. 2.2

SeleMix (Selective editing via Mixture models)

Description

Selemix is an R package to treat quantitative data, which aims to identify a set of units affected by errors which potentially influence the estimates of interest (selective editing).

The underlying methodology is based on particular latent class models known in the literature as contamination models. Specifically, it is assumed that "true" data (that is not affected by errors), possibly in logarithmic scale, are independent realizations of a multivariate Gaussian distribution with the mean vector which can in turn be expressed as a linear combination of a set of covariates not contaminated. The "intermittent" nature of the error mechanism is captured by Bernoulli variables that have the role of indicators for the occurrence of error on each unit. Moreover, the error is assumed additive and associated with a Gaussian vector with zero mean and variance covariance matrix proportional to the variance covariance matrix that characterizes the distribution of data without errors. The explicit modeling of the not contaminated data distribution and the error mechanism allow to obtain the distribution of actual data conditionally on the observed data. On the basis of this distribution estimates of the true values, and therefore of errors, are obtained. For each unit, a score is calculated in terms of the difference (possibly weighted with sample weight) between the predicted value and the observed value. So all units are sorted (in descending order) according to their score. Assuming that the parameter of interest is an average or total population, selection of observation to be reviewed as interactive, is made considering the estimated error that remains in the data after interactive editing. The number of units selected according to this criterion relies on a userspecified threshold that is related to the accuracy of the estimate of interest.

In the following, the main functions of the package SeleMix are described:

- ml.est: this function performs the maximum likelihood estimates of the parameters of a contamination model by ECM algorithm and it provides the expected values of the "true" data for all units that were used for the estimation. Also it returns, to each unit, the posterior probability of occurrence of the error and the flags of classification as outliers no outlier calculated on the basis of a threshold for the probability of error specified by the user. It requires the specification of the type of model assumed for the true data (normal or lognormal) and some technical parameters for the algorithm ECM. It requires the specification of the type of model assumed for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the true data (normal or lognormal) and some technical parameters for the algorithm ECM.
- **pred.y**: on the basis of a set of contamination model parameters, and a set of observed data, it calculates the expected values of the corresponding real data. Missing values for the variables response as well as are allowed, but not for covariates.
- sel.edit: it performs Selective Editing. On the basis of a set of observed data and the corresponding predictions for the true data, it selects the units required for interactive editing. It requires in input the wanted accuracy threshold and, if present, the sample weights associated with the units. It provides the score for each unit and the corresponding rank.

Given the opportunity to use the features of the package also in the presence of incomplete data, it can also be used as a tool for robust imputation of multivariate Gaussian data.

Information

Status:validatedAuthor:IstatLicence:EUPL-1.1GSBPM code:5.3 Review and validate
5.4 Edit and imputeProgramming language:RKeywords:latent class models, selective editing,
influential errorContact:name: Maria Teresa Buglielli@istat.it

Software and documentation

COPYRIGHT

Copyright 2013 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 0.9.1 Windows binaries
- Version 0.9.1 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – SeleMix v. 0.9.1

Vignettes – SeleMix v. 0.9.1

OTHER DOCUMENTATION

Barcaroli G., D. Zardetto. 2012. <u>Use of R in Business Surveys at the Italian National Institute of</u> <u>Statistics: Experiences and Perspectives</u>. In *Proceedings of the 4th International Conference of Establishment Surveys (ICES IV)*, American Statistical Association Montréal, 11-14 June 2012.

EVER (Estimation of Variance by Efficient Replication)

Description

The EVER package (the acronym stands for: Estimation of Variance by Efficient Replication) is mainly intended for calculating estimates and standard errors in complex surveys. Variance estimation is based on the extended DAGJK (Delete-A-group Jackknife) technique proposed by P. S. Kott.

The advantage of the DAGJK method over the traditional jackknife is that, unlike the latter, it remains computationally manageable even when dealing with "complex and big" surveys (tens of thousands of PSUs arranged in a large number of strata with widely varying sizes). In fact, the DAGJK method is known to provide, for a broad range of sampling designs and estimators, (near) unbiased standard error estimates even with a "small" number (e.g. a few tens) of replicate weights.

Besides its peculiar computational efficiency, the DAGJK method takes advantage of the strong points it shares with the most common replication methods. As a remarkable example, EVER is designed to fully exploit DAGJK's versatility: the package provides the user with a user-friendly tool for calculating estimates, standard errors and confidence intervals for estimators defined by the user themselves (even non-analytic). This functionality makes EVER especially appealing whenever variance estimation by Taylor linearisation can be applied only at the price of crude approximations (e.g. poverty estimates).

Main Statistical Functions

- Delete-A-Group Jackknife replication
- Calibration of replicate weights
- Estimates and Sampling Errors (standard error, variance, coefficient of variation, confidence interval, design effect) for:
 - Totals
 - Means
 - Absolute and relative frequency distributions
 - Ratios between totals
 - Multiple regression coefficients
 - o Quantiles
- Estimates and Sampling Errors for user-defined Complex Estimators (even non-analytic)
- Estimates and Sampling Errors for Subpopulations (Domains)
 - All the analyses above can be carried out for arbitrary domains

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	5.6 Calculate weights 5.7 Calculate aggregates
Programming language:	R
Keywords:	calibration, estimation, replication variance, complex surveys, complex estimators, Delete-A-Group Jackknife, R
Contact:	name: Diego Zardetto email: <u>zardetto@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

R (version \geq 2.5.1)

COPYRIGHT

Copyright 2013 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

- Version 1.2 Windows binaries
- Version 1.2 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – EVER v. 1.2

<u>ReGenesees (R evolved Generalised software for sampling estimates and errors</u> <u>in surveys)</u>

Description

ReGenesees (R Evolved Generalized Software for Sampling Estimates and Errors in Surveys) is a full-fledged R software for design-based and model-assisted analysis of complex sample surveys.

This system is the outcome of a long-term research and development project, aimed at defining a new Istat standard for calibration, estimation and sampling error assessment in large-scale sample surveys.

Main Statistical Functions

- Complex Sampling Designs
 - Multistage, stratified, clustered, sampling designs
 - o 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
 - o Ultimate Cluster approximation
 - Collapsed strata technique for handling lonely PSUs
 - Taylor-linearization of nonlinear "smooth" estimators
 - o Generalized Variance Functions method
- Estimates and Sampling Errors (standard error, variance, coefficient of variation, confidence interval, design effect) for:
 - o Totals
 - o Means
 - Absolute and relative frequency distributions (marginal, conditional and joint)
 - Ratios between totals

- Shares and ratios between shares
- Multiple regression coefficients
- Quantiles
- 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)

System Architecture

The ReGenesees system has a clear-cut two-layer architecture. The application layer of the system is embedded into an R package named **ReGenesees**. A second R package, called **ReGenesees.GUI**, implements the presentation layer of the system (namely a Tcl/Tk GUI). Both packages can be run under Windows as well as under Mac, Linux and 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. On the contrary, less experienced R users will take advantage from the user-friendly mouse-click graphical interface.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	5.6 Calculate weights 5.7 Calculate aggregates
Programming language:	R
Language of the GUI:	EN
Keywords:	calibration, estimation, variance estimation, complex surveys, complex estimators, automated linearization, R
Contact:	name: Diego Zardetto email: <u>zardetto@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES for package ReGenesees

R (≥ 2.14.0)

SOFTWARE DEPENDENCIES for package ReGenesees.GUI

R (≥ 2.14.0)

R packages: ReGenesees, <u>tcltk2</u>, <u>RODBC</u> and <u>svMisc</u>

COPYRIGHT

Copyright 2015 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

Package ReGenesees (statistical engine)

- Version 1.9 Windows binaries
- Version 1.9 Package source

Package ReGenesees.GUI (graphical user interface)

- Version 1.9 Windows binaries
- Version 1.9 Package source

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Reference manual – ReGenesees v. 1.9

Reference manual – ReGeneseesGUI v. 1.9

OTHER DOCUMENTATION

Fallows A., Pope M., Digby-North J., Brown G., Lewis D. 2015. <u>A Comparative Study of Complex</u> <u>Survey Estimation Software in ONS</u>. Romanian Statistical Review, 3:46-64.

Zardetto D. 2015. <u>ReGenesees: an Advanced R System for Calibration, Estimation and Sampling</u> <u>Error Assessment in Complex Sample Surveys</u>, (extended version). Journal of Official Statistics, 31(2):177-203.

Zardetto D. 2013. <u>ReGenesees: an Advanced R System for Calibration, Estimation and Sampling</u> <u>Errors Assessment in Complex Sample Surveys</u>. In *Proceedings of the* 7th International Conferences on New Techniques and Technologies for Statistics (NTTS), Eurostat, Brussels, 5-7 March 2013. Barcaroli G., D. Zardetto. 2012. <u>Use of R in Business Surveys at the Italian National Institute of</u> <u>Statistics: Experiences and Perspectives</u>. In *Proceedings of the 4th International Conference of Establishment Surveys (ICES IV)*, American Statistical Association, Montréal, 11-14 June 2012.

ANALYSE phase

Prepare draft outputs

This sub-process is where the data are transformed into statistical outputs. It includes the production of additional measurements such as indices, trends or seasonally adjusted series, as well as the recording of quality characteristics.

A. Composite indices

A composite index is a mathematical combination (or aggregation as it is termed) of a set of indicators that represent the different dimensions of a phenomenon to be measured.

Constructing a composite index is a complex task. Its phases involve several alternatives and possibilities that affect the quality and reliability of the results. The main problems, in this approach, concern the choice of theoretical framework, the availability of the data, the selection of the more representative indicators and their treatment in order to compare and aggregate them.

In particular, we can summarize the procedure in the following main steps:

- 1. Defining the phenomenon to be measured. The definition of the concept should give a clear sense of what is being measured by the composite index. It should refer to a theoretical framework, linking various sub-groups and underlying indicators. Also the model of measurement must be defined, in order to specify the relationship between the phenomenon to be measured (concept) and its measures (individual indicators). If causality is from the concept to the indicators we have a reflective model indicators are interchangeable and correlations between indicators are explained by the model; if causality is from the indicators to the concept we have a formative model indicators are not interchangeable and correlations between indicators are not explained by the model.
- 2. Selecting a group of individual indicators. The selection is generally based on theory, empirical analysis, pragmatism or intuitive appeal. Ideally, indicators should be selected according to their relevance, analytical soundness, timeliness, accessibility and so on. The selection step is the result of a trade-off between possible redundancies caused by overlapping information and the risk of losing information. However, the selection process also depends on the measurement model used: in a reflective model, all the individual indicators must be intercorrelated; whereas in a formative model they can show negative or zero correlations.
- 3. Normalizing the individual indicators. This step aims to make the indicators comparable. Normalization is required before any data aggregation as the indicators in a data set often have different measurement units. Therefore, it is necessary to bring the indicators to the same standard, by transforming them into pure, dimensionless, numbers. Another motivation for the normalization is the fact that some indicators may be positively correlated with the phenomenon to be measured (positive polarity), whereas others may be negatively correlated with it (negative polarity). We want to normalize the indicators so that an increase in the normalized indicators corresponds to increase in the

composite index. There are various methods of normalization, such as re-scaling (or Min-Max), standardization (or z-scores) and 'distance' from a reference (or index numbers).

- 4. Aggregating the normalized indicators. It is the combination of all the components to form one or more composite indices (mathematical functions). This step requires the definition of the importance of each individual indicator (weighting system) and the identification of the technique (compensatory or non-compensatory) for summarizing the individual indicator values into a single number. Different aggregation methods can be used, such as additive methods (compensatory approach) or multiplicative methods and unbalance-adjusted functions (non-compensatory or partially compensatory approach).
- 5. Validating the composite index. Validation step aims to assess the robustness of the composite index, in terms of capacity to produce correct and stable measure, and its discriminant capacity (Influence Analysis and Robustness Analysis).

B. Seasonal adjustment of time series

Seasonality can be defined as the systematic intra-year movement caused by various factors, e.g. weather changes, calendar, vacation or holidays and usually consists of periodic, repetitive and generally regular and predictable patterns in the level of a time series. Seasonality can be influenced also by production and consumption decisions made by economic agents taking into account several factors like endowments, their own expectations, as well as preferences and the production techniques available in the economy.

Differently, cyclic pattern presents non fixed rises and falls and its fluctuation length is usually not shorter than 2 years.

The overlap of this two kinds of fluctuations (seasonal and cyclic) in a time series could provide some problems for short term (monthly or quarterly) variation interpretation, above all when the seasonal component is highly represented in the observed data. For this reason, in order to measure cyclical changes, short term variations are computed from seasonal adjusted series. In turn, seasonal adjustment is the process of seasonal and calendar effects removal from a time series. This process is performed by means of analytical techniques that break down the series into components with different dynamic features. These components are unobserved and have to be identified from the observed data based on an ex-ante assumptions on their expected behavior. Broadly speaking, seasonal adjustment includes the removal of both within-a-year seasonal movements and the influence of calendar effects (such as the different number of working days, or Easter and moving holidays).

Notice that calendar effects are not constant among different countries or economic sectors, so that, time series which include them are not comparable each other. For this reason, generally, calendar effects are removed together with seasonal component in the seasonal adjusted series, so that, it is possible to better catch the yearly variation (computed with respect to the same period of the previous year), as well as, the mean yearly variation. Moreover, together with the seasonal adjusted series, can be also produced time series corrected only for calendar effects.

Once removed the repeated impact of these effects, seasonally adjusted data highlight the underlying long-term trend and short-run innovations in the series.

Seasonal adjustment approaches

All the seasonal adjustment methods are based on the assumption that each time series, Yt (with a time index t = 1,2,...T), can be decomposed into three different unobserved components:

- A trend-cycle (CT_t) component representing long-run movement of the series (like those associated to business cycles). It generally depends on structural conditions like institutional situations, technological and demographic trends or patterns of civil and social organization.
- A seasonal component (S_t) representing the intra-year (monthly, quarterly) fluctuations.
- An **irregular component (It)** representing the short term fluctuations that are not systematic and, to a certain extent unpredictable, e.g. uncharacteristic weather patterns.

Although the series may be decomposed in different ways, generally two main approaches consistent with the European guideline (Eurostat 2015), are considered:

- 1. Arima Model Based (AMB) approach, developed among the others by Burman (1980), Box, Hillmer and Tiao (1978) and Hillmer and Tiao (1982), based on the assumption that there exists a statistical parametric model (ARIMA) representing the probabilistic structure of the stochastic process connected to the observed time series. Time series assumed to be a finite part of a particular realization of a stochastic process. The linear filters used in this approach depend, consequently, on the features of the time series considered. This kind of approach is adopted in the TRAMO-SEATS (Time series regression with ARIMA noise, missing observations and outliers and Signal Extraction in ARIMA time series - TS) procedure developed by Gómez and Maravall (1996).
- 2. Filter Based Approach (FLB), a non-parametric or semiparametric approach, which, differently from AMB approach, does not require to hypothesize a statistical model representing the series. Indeed, it is based on an iterative application of several linear filters on the series based on central moving averages. These procedures are referred to as ad hoc because the filters are chosen according to empirical rules, not taking into account the probabilistic structure of the stochastic process generating the series. To this approach belong the classical methods of the X-11 (X11) family: from the first X11 and X11-Arima (X-11A) to the more recent X-12-ARIMA (X-12A) (Findley et al. 1998) and X-13-ARIMA-SEATS (X-13AS) (Findley, 2005) which include several improvements over the previous versions; among which, the most remarkable is the use of reg-Arima models aimed at pre-treating the data and at improving the forecasting performances of the series, that in turn translates into an improvement of the estimated seasonal factor.

In both cases data are pre-treated for selecting the decomposition scheme to be applied to the time series (additive, multiplicative, log-additive, etc.). Moreover, some deterministic effects like outliers or calendar effects are removed. This pre-treated series is the input of the following step whose output is the seasonal adjusted series (SA). Once the seasonal adjusted series is obtained there is a last step in which some elements identified in the pre-treatment phase and

related to the trend-cycle components (like level shift) or to the irregular component (like additive outlier or temporary changes) are included back; while are taken out from the final series the calendar effects and the seasonal outliers.

Apply disclosure control

The primary function of a public statistical system is to produce official statistics for its own country. In fact, the Legislative Decree no.322 of September 6, 1989, constituting the National Statistical System (Sistan), cites: "The official statistical information is provided to the country and to international organizations through the National Statistical System" (Article 1, paragraph 2), and "the data processed in the framework of statistical surveys included in the national statistics program are the heritage of the community and are distributed for purposes of study and research to those who require them under the rules of the present decree, subject to the prohibitions contained in art. 9 "concerning statistical confidentiality (art. 10 paragraph 1).The Legislative Decree no.322 / 1989, also states that "the data collected in the context of statistical surveys included in the national Statistical Programme cannot be communicated or disseminated to any external, public or private entity, nor to any office of public administration except in aggregate form so as to be unable to derive any reference to identifiable individuals." In any case, the data cannot be used in order to re-identify the parties concerned.

Further principles concerning the protection of the data confidentiality, are established by the Code of Professional Ethics and Good Conduct for the Treatment of personal data for statistical purposes and scientific research purposes performed within the framework of the National Statistical System (Legislative Decree no. 196, June 30, 2003). In particular, the Code defines the concept of identifiability of a statistical unit, in terms of possibilities, through the use of reasonable means, to establish a significantly likely relationship between the combination of the mode of the variables related to the statistics unit and its data identification. Moreover, the means used for the identification of the person concerned, for example, the economic resources, time, the possibility of crosschecking with name records or other sources, etc. are specified.

The translation of the concepts laid down by the law in operational rules from a statistical point of view requires a preliminary identification of the statistical units subject to identification risk and thus a precise definition of what constitutes a breach of confidentiality. The subsequent quantification of the probability of breaching the confidentiality shall define the most suitable techniques to ensure data protection.

The definition of a breach of confidentiality adopted by the National Statistical Institutes is based on the concept of identifiability of a unit of the population observed (the respondent). By defining as *intruder* a person who has an interest in breaching the confidentiality of the released data, the intrusion occurs when the intruder is able to match, with a certain degree of certainty, the information released to the respondent. The release of statistical information with confidential data in no case involves the so-called *direct identifiers* (i.e. the variables that uniquely identify the person such as tax code, name or company name, address, etc.). The problem arises for so-called *indirect identifiers* (or key variables). These are the variables that do not identify the person directly but allow to circumscribe the belonging population and which the intruder will use for his/her own purposes. Indirect identification could be determined, for example, by the combined use of territorial variables, economic activity and size class. The mechanism by which identification can happen may be immediate (e.g. direct recognition) or assigned to more or less complex algorithms of information combination (*record linkage, statistical-matching*, etc.).

To limit the risk of identification the National Statistical Institutes may modify data (for example by using disturbance techniques), or have an effect on the indirect identifiers by removing it in whole or in part, or reducing its details (i.e. by deciding to not release such detail as *municipality*

and leaving in its place the variable *district, or region*). The application of the protection techniques, both for the dissemination of tables as for the communication of elementary data, leads to a reduction or a change in the information content of the data released (loss of information).

A. The breach of confidentiality in the disseminated tables

The table represents the tool mostly used by the national statistical institutes for the dissemination of aggregate data, or grouped together in cells, defined by the intersection of the classification variables. The concept of a breach of confidentiality does not depend upon the type of product used for the dissemination. Consistently in line with the previous section, also in the case of aggregate data, a breach occurs when it is possible to draw information that allows the identification of the individual. The definition of "confidential" information also includes sensitive data and judicial data (as defined in the Legislative Decree no. 196, June 30, 2003, art. 4), while public variables are not considered confidential (the character or combination of characters, qualitative or quantitative, subject to a statistical survey that refers to information present in public registers, lists, records, documents or sources accessible to anyone – definition contained in the Code of Conduct). When a table is to be released an initial assessment concerns the content information on data to be published: if it is not confidential than there is no necessity to implement statistical security procedures of data, otherwise it is necessary to apply the rules of protection of privacy. The evaluation of the risk of a breach of confidentiality of data in the table is carried out for each cell individually: when the value inside one of the cells refers to (with a certain degree of certainty) the subject to which relates the data itself (sensitive cell), then the table does not respect the rules on the protection of confidentiality.

The process of aggregated data protection comprises several phases. The first stage defines the area in which one is working, which tables are to be processed and what their characteristics are. Then we define the risk rule or the criterion according to which to determine whether or not a cell is at risk of a breach of confidentiality. The final phase concerns the implementation of procedures for the protection of confidentiality. These depend on the type of tables that one intends to release and on any restrictions of publication, but also on the type of reserved variables, on the underlying complexity to each processing and on data availability.

Although some of the principles described below, with particular reference to the threshold rule, are also used for frequency tables, the rules listed mainly refer to the magnitude tables. In the case of frequency tables, cells at risk are identified as a result of an evaluation done on a case by case basis and not by resorting to general rules as is the case of the magnitude tables.

Magnitude tables and risk rules

The risk rules used for magnitude tables are those based on the size of the cells (threshold or frequency rule), and those based on concentration measurements (dominance rule and ratio rule). The threshold rule is widely used by Istat according to which a cell is sensitive if the number of units contained therein is less than an *n* value (threshold) fixed a priori. In order to apply this rule to magnitude tables you need the relative frequency table. The protection depends upon the value of *n* that is applied to the table: the higher the threshold value the higher the level of protection applied. There is no univocal criterion for identifying the threshold value that will depend on the assumed scenario of intrusion and data processed. The minimum value of the threshold shall be three (as provided by the Code of Conduct).

According to the dominance rule [(n, k) -dominance] a cell is at risk if the first *n* contributors hold a proportion of its total value which is higher than the threshold k% fixed a priori. The level of protection that should be applied to the table depends from the two values of *n* and of *k*. There are no univocal criteria for fixing the two parameters. Based on the statistical units involved and the desired levels of security one can define the parameters by identifying a maximum allowable concentration.

The ratio rule (p-rule) relies on the accuracy with which the value of the first contributor can be estimated, assuming that the second contributor attempts the breach. The cell is considered at risk if the relative error is less than a *p* threshold fixed a priori.

In case of tables with possible data of opposite sign, the risk rules based on concentration measurements become meaningless. However, their application is possible by using the absolute values of the contributors.

Operating a breach of confidentiality in a context with possible negative data is much more complex. The general recommendation is to assign parameters to the risk functions with less stringent values in respect to solely positive data.

In the case of sample tables that are obtained by detecting data on a subset of the reference population, the evaluation of the breach of confidentiality risk has to take into account the sampling plan used. The value listed in the cells is an estimate made by extending a partial value (detected in the sample) to the reference population. The units observed are unknown and also the true value of the population is not recognized. The risk of violation is contained for cells containing data with a survey's weight greater than the unit. In this context suggestion of a breach of confidentiality appears unlikely. However, especially for tables of economic data, a careful assessment of the breach of confidentiality risk is necessary even in the case of sample tables. In fact, in some cases the most representative units (dominant) are included in the sample with certain probability. Furthermore, in the case of stratified samples, some cells are sampled at 100% and then the detected value matches (unless missing responses) with the value of the population.

Except for special cases where the sample design and the number of sampled units allow considering a safe table in terms of confidentiality, rules of confidentiality must be applied to the sample tables too.

The criterion used by Istat considers the implementation of the risk rules on the obtained estimated cell values using survey's weights. That requires that the sampled units are "similar" to those present in the population.

Frequency tables and risk rule

Frequency tables are used primarily to represent social phenomena and census data. The only criterion to determine whether or not a cell is at risk for this type of tables is based on the size of the cells. The risk rules cannot in fact be applied based on measures of concentration. There are no univocal rules for determining whether a frequency table is at risk of a breach of confidentiality or not. In fact, a cell with a low frequency (for example equal to 1) not always indicates a cell at risk, and vice versa a cell that contains a high number of units not always can be considered safe in terms of statistical confidentiality.

As a general rule, the frequency tables that have one of the cases listed below are considered at risk of confidentiality breach:

- marginal with less than three contributors;
- all units belong to a single category (group disclosure) or the sole contributor of a cell (self recognition) acquires confidential information on all other units (all concentrated in another cell).

Statistical protection of tables

After the cells at risk are identified, one must modify the table in an appropriate manner ensuring the anonymity of information contained therein. There are many techniques of data protection and they range from a unification of adjacent modes, to methods based on the original data modification, to the introduction of missing values (suppressions). The methods used by Istat are: the amendment of classification variables mode and the introduction of missing values.

A method of protection of the tables that is not based on the change of the values in the cells is the definition of a different modes combination. After identifying the risk rule, the method consists in determining the said modes in such a way that the distribution of the characters and / or units in the cells must be such that it does not present any sensitive cell.

By changing the modality accordingly it is possible, for example, to obtain a table which has a minimum size (for example greater than or equal to three) in each cell, or a table with a default maximum concentration of the character in each cell.

The change in the modes of classifying variables can be considered a practical solution only when the nature of the classifying variables is transferable, and if the tables to be released need not satisfy strict rules dictated by regulations that constrain the details of the classifying variables.

The technique that involves insertion of a missing value stipulates that the value of the cells at risk is deleted (obscured). The suppression operated on cells at risk is also called primary suppression. With the introduction of missing values within the sensitive cells, the process of protecting the table does not exhaust. In the first place, it is necessary to evaluate that the deleted cells cannot be calculated from the issued data, for example, by difference from the marginal values. Suppressions are to be distributed among the cells of the table to ensure that the table is properly protected according to the set criteria. When this does not occur it is necessary to introduce additional missing values between not-at-risk cells: the secondary suppressions. The literature proposes several algorithms for the determination of the secondary suppressions. Currently the most widely used by Istat is the HiTas algorithm available in some generalized software such as Tau-ARGUS.

Linked tables

The tables are defined as linked when they contain data on the same response variable and have at least the same classifying variable. The most frequent case of linked tables is represented by tables with common cells, with particular reference to the marginal values. The connection between the statistical data may also form part of a wider context. Sometimes indeed different surveys show the same aggregates.

The application of the confidentiality rules to linked tables implies that common information (cells) have assigned the same deployability status.

To optimize the process of protection it would be appropriate, where possible, to simultaneously operate the protection of all the linked tables.

B. The breach of confidentiality in the release of elementary data

Elementary data can be defined as the end product of a statistical survey after the phases of design, implementation, monitoring and correction. The elementary data in the dissemination phase is an archive of records each containing all the validated information (usually a subset of those recorded) relating to an individual statistical unit. These variables, as well as in the case of aggregate data disseminated through tables, can be classified as key variables, as indirect identifiers, or as reserved variables.

With respect to tables' release, there will be a substantial change both in the set of key variables, that in general will be more numerous, and in the content of a possible violation, as the variables reserved in the elementary data are shown all together. By contrast, the release of micro-data only covers sample collections and file access is much more controlled (for research purposes only and behind the signing of a form / contract). However, there is no doubt that the release of elementary data is a most sensitive issue with respect to disseminated tables. For this reason, specific models for measuring the risk of identification often based on probability models have been developed. The methods of data elementary protection fall into three categories:

- recoding of variables (global recoding) consists in reduction of release detail of some variables (e.g. the age in the five-year classes rather than annual classes);
- suppression of information (local suppression): to eliminate features that make some records more easily identifiable;
- disturbance of the published data: different methods but with the same purpose of the tables.

Among the initiatives involved in "protected" release of elementary data are included the so called Microdata File for Research (MFR), the public use files (mIcro.STAT) and the Laboratory for the Analysis of the ELEmentary data (ADELE). MFR files are produced for statistical surveys regarding both individuals and families and businesses and are made specifically for the needs of scientific research. The release of these files is subject to the fulfilment of certain requirements relating both to the organization and to the characteristics of the research project for which the file is required. MIcro.STAT files are public use files, obtained starting from the respective MFR file, properly treated in terms of protection of privacy and downloaded directly from Istat website.

The ADELE Laboratory, active since 1999, is a so-called Research Data Centre (RDC) or a "safe" place which can be accessed by the researchers and academics to make their own statistical

analysis of elementary data produced by the National Statistics Institution in compliance with the confidentiality rules. Main objective of the ADELE Laboratory is to provide external "expert" users the ability to analyze the basic data of the Istat surveys, by shifting the phase of verification of the protection of confidentiality of statistical analysis on output rather than input (as in the case of files for research purposes and public use files). The protection of confidentiality for processing carried out at the ADELE Laboratory is ensured in several ways:

- legally: the user subscribes to a form in which he/she undertakes to comply with specific rules of conduct;
- physically: through the control of the working environment. The Laboratory is located at the head office of Istat with workers in charge of the control room; input and output operations and access to external network are disabled to users;
- statistically: by controlling the user analysis results prior to the release.

ANALYSE phase – METHODS

• PREPARE DRAFT OUTPUTS

Computation and evaluation of composite indices

Methods for constructing composite indices: one for all or all for one? 2013 Rivista Italiana di Economia Demografia e Statistica, Volume LXVII n. 2.

Seasonal adjustment of time series

Relazione del team tecnico incaricato della definizione di metodi standard per la destagionalizzazione di serie storiche con metodi implementati in diversi strumenti IT (TS, X12-Arima, X13-Arima-Seats, JDemetra) 2015 Gdl per la definizione di standard per l'Istat, Istat Destagionalizzazione di serie storiche con metodologia Arima model based (AMB) implementata nel software JDemetra+ 2015 Gdl per la definizione di standard per l'Istat, Istat

ESS guidelines on seasonal adjustment (2015 Edition)

2015 Manuals and Guidelines, Eurostat

References

Computation and evaluation of composite indices

Massoli P., Mazziotta M., Pareto A., Rinaldelli C. 2014. <u>Indici compositi per il BES</u>. *Giornate della Ricerca, Istat*, 10-11 Novembre.

Massoli P., Mazziotta M., Pareto A., Rinaldelli C. 2013. <u>Metodologie di sintesi sperimentali per</u> <u>i domini del BES</u>. XXXIV Conferenza Italiana di Scienze Regionali (AISRE), Palermo, 2-3 Settembre.

Mazziotta M., A. Pareto 2011. <u>Un indice sintetico non compensativo per la misura della</u> <u>dotazione infrastrutturale: un'applicazione in ambito sanitario</u>. *Rivista di Statistica Ufficiale*, 1:63-79.

Mazziotta C., Mazziotta M., Pareto A., Vidoli F. 2010. <u>La sintesi di indicatori territoriali di</u> <u>dotazione infrastrutturale: metodi di costruzione e procedure di ponderazione a confronto</u>. *Rivista di Economia e Statistica del Territorio*, 1:7-33.

OECD. 2008. <u>Handbook on Constructing Composite Indicators. Methodology and user guide</u>. OECD Publications.

Aureli Cutillo E. 1996. Lezioni di statistica sociale. Parte II. Sintesi e graduatorie. CISU, Roma.

Delvecchio F. 1995. Scale di misura e indicatori sociali. Cacucci, Bari.

Silvio-Pomenta J. F. 1973. <u>Typological study using the Wroclaw Taxonomic Method (A study of</u> regional disparities in Venezuela). Social science project on human resources indicators, n. 28, UNESCO.

Harbison F. H., J. Maruhnic, J.R. Resnick. 1970. *Quantitative Analyses of Modernization and Development*. Princeton University Press, New Jersey.

Seasonal adjustment of time series

Eurostat. 2015. *ESS Guidelines on Seasonal Adjustment (2015 Edition)*. Manuals and Guidelines, Eurostat.

Grudkowska S. 2015. JDemetra+ Reference Manual, v. 0.1. Narodowy Bank Polski.

Eurostat. 2014. <u>ESS Handbook for Quality Reports (2014 Edition)</u>. Manuals and guidelines, Eurostat.

Findley D. F. 2005. <u>Some Recent Developments and Directions in Seasonal Adjustment</u>. *Journal of Official Statistics*, 21(2):343-365.

Giovannini E., D. Piccolo. 2000. Seasonal adjustment procedures – experiences and perspectives. Istat.

Findley D. F., B. C. Monsell, W. R. Bell, M. C. Otto, B. C. Chen. 1998. <u>New Capabilities and</u> <u>Methods of the X-12-ARIMA Seasonal-Adjustment Program</u>. Journal of Business & Economic Statistics, 16(2):127-152.

Gómez V., A. Maravall. 1996. Programs TRAMO and SEATS, Instruction for User (Beta Version: September 1996). Banco de España.

Chen C., L. Liu. 1993. Joint Estimation of Model Parameters and Outlier Effects in Time Series. *Journal of the American Statistical Association*, 88(421):284-297.

Hillmer S. C., G. C. Tiao. 1982. <u>An ARIMA-Model-Based Approach to Seasonal Adjustment</u>. *Journal of the American Statistical Association* 77(377):63-70.

Burman J. P. 1980. <u>Seasonal Adjustment by Signal Extraction</u>. *Journal of the Royal Statistical Society*, 143(3): 321-337.

Box G. E. P., S. C. Hillmer, G. C. Tiao. 1978. <u>Analysis and Modeling of Seasonal Time Series</u>. *Seasonal Analysis of Economic Time Series*, 309-334.

Box G. E. P., G. M. Jenkins, G. C. Reinsel. 1970. *Time Series Analysis: Forecasting and Control.* Wiley, Holden Day, San Francisco.

APPLY DISCLOSURE CONTROL

Disclosure protection of non-nested linked tables in Business Statistics

2009

Supporting paper, Joint UNECE/Eurostat work session on statistical data confidentiality, Bilbao, Spain

Metodologie e tecniche di tutela della riservatezza nel rilascio dell'informazione statistica 2004

Collana Metodi e Norme, n. 20, Istat

References

Hundepool A., J. Domingo-Ferrer, L. Franconi, S. Giessing, E. Sculte Nordholt, K. Spicer, P.P. De Wolf. 2012. *Statistical Disclosure Control*. Wiley & Sons, London.

ANALYSE phase – TOOLS

• PREPARE DRAFT OUTPUTS

<u>COMIC</u> (COMposite Indices Creator)

The COMIC software computes composite indices, by various aggregation methods, and evaluates their robustness

Ranker

Software system based on Visual Basic language for the analysis and benchmarking of results produced by different methods of statistical summary on composite indicators available in the literature.

APPLY DISCLOSURE CONTROL

<u>ARGUS</u>

Generalized software for applying disclosure control.

<u>COMIC</u>

Description

The COMIC software computes composite indices, by various aggregation methods, and evaluates their robustness. COMIC calculates composite indices relevant to different periods of time.

COMIC deals with:

- different format acquisition (.csv or .xls or .txt or SAS) for different individual indicators available for each unit, as specified by users;
- standardization/normalization of the individual indicators;
- different aggregation methods;
- output display in both tables and graphics;
- comparison of the results;
- sensitivity and robustness analysis to evaluate the composite indices.

COMIC has implemented the following aggregation methods:

- Mean indices 0-1. Arithmetic mean of individual indicators normalized by min-max method;
- Mean z-scores. Arithmetic mean of standardized individual indicators;
- Jevons Index. Geometric mean of individual indicators computed as index numbers;
- Mazziotta-Pareto Index (MPI). Arithmetic mean of standardized individual indicators; a
 penalty is subtracted from the mean to take into account the variability of the individual
 indicators (MPI with negative penalty);
- Adjusted Mazziotta-Pareto Index (AMPI). Arithmetic mean of the individual indicators normalized by the min-max method; a penalty is subtracted from the mean to take into account the variability of the individual indicators (AMPI with negative penalty);
- Geometric Mean Index (IMG). Geometric mean of the individual indicators normalized by the min-max method.

Information

Status:	validated
Author:	Istat
Licence:	EUPL-1.1
GSBPM code:	6.1 Prepare draft outputs
Programming language:	SAS
Language of the GUI:	
Keywords:	composite indices, normalization, aggregation, ranking, sensitivity, robustness
Contact:	name: Pierpaolo Massoli email: <u>pimassol@istat.it</u>

Software and documentation

SOFTWARE DEPENDENCIES

SAS for Windows (SAS base)

COPYRIGHT

Copyright 2013 Istat

Licensed under the European Union Public Licence (EUPL), version 1.1 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

DOWNLOAD

COMIC version 1.0

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

Manuale utente - COMIC v. 1.0

OTHER DOCUMENTATION

Massoli P., Mazziotta M., Pareto A., Rinaldelli C. 2015. <u>COMIC: a tool for composite indices</u> <u>evaluation</u>. *Conference Dealing with complexity in society: from plurality of data to synthetic indicators*, Padova, 17-18 Settembre.

Ranker

Description

Ranker is a software system, for the analysis and benchmarking of results produced by different composite indicators.

Ranker is a prototype of a generalized tool, which allows to:

- acquire in standard format (.xls) values of the different individual indicators available for each unit (eg. geographical areas);
- compute, for each unit, one or more implemented methods;
- display the values and rankings of each method both in tables and graphics;
- compare the rankings of the different methods.

Ranker has implemented eight methods:

- Mazziotta-Pareto Index (MPI) (De Muro et al. 2010);
- Wroclaw (Wroclaw);
- Arithmetic Mean of the z-scores (M1Z);
- Ranking (Grad.RNK);
- Relative indices (IR);
- Arithmetic Mean of index numbers (ANIM);
- Geometric Mean of index numbers (GNIM);
- Squared Mean of index numbers (QNIM).

The data analysis is based on five different phases:

- The organization and the reading of data. The acquired data must be organized in tabular form, in a matrix that describes each territorial unit (in line) according to the value of the individual indicators (in column). For each indicator considered, it is necessary to provide, in addition, the polarity, distinguishing those that describe an effect "positive" with respect to the dynamics of the phenomenon and those that, on the contrary, are related in the reverse direction and to which corresponds a decreasing ranking of the territorial units;
- 2. The normalization. This step is intended to obtain indicators purified by the specific units of measurement, which are of equal width (eg. between 0 and 100);
- The aggregation. After loading the matrix with individual indicators and after their standardization, the application makes the processing of all the methods mentioned above in order to calculate the relative composite indicators and describe each territorial unit;
- 4. The visualization. Through the application you can view both the values obtained and the rankings derived from them. Besides the description of the numerical values in tabular form, Ranker enables the production of maps that graphically represent the results obtained by each methods;
- 5. The evaluation of the methods. The software enables you to evaluate comparatively the results produced by different methods: in particular, the impact of the choice of the

method on the final result (the ranking). A first instrument is given by the matrix of rank correlation obtained comparing different methods. A second instrument used is the matrix plot, which can be produced on the values of the indicators, or better, on the rankings obtained. A final tool is given by the scatter plot calculated on pairs of selected methods.

Istat offers the availability on choosing the online version (<u>i.ranker</u>) or the desktop one. The online version implements the first five methods above descripted and offers the ability to view the values and the ranking resulting from the application of each method both in tabular and graphical layout. The desktop version implements all methods but only allows the tabular display of results.

Information

Status:	validated
Author:	Istat
Licence:	<u>EUPL-1.1</u>
GSBPM code:	6.1 Prepare draft outputs
Programming language:	Visual Basic
Language of the GUI:	IT
•	composite indicators, normalization, aggregation, ranking
Contact:	name: Marco Broccoli email: <u>broccoli@istat.it</u>

Software and documentation

COPYRIGHT

Copyright 2014 Istat

Licensed under the European Union Public Licence (EUPL), version 1.0 or subsequent. You may not use this work except in compliance with the Licence. You may obtain a copy of the Licence at: <u>http://ec.europa.eu/idabc/eupl.html</u>. Unless required by applicable law or agreed to in writing, software distributed under the Licence is distributed on an "AS IS" basis, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the Licence for the specific language governing permissions and limitations under the Licence.

DISCLAIMER

Istat assumes no responsibility for the results arising from use of the instrument that is inconsistent with the methodological guidance contained in the documentation available.

ONLINE SYSTEM

i.ranker

DOWNLOAD

RANKER version 1.2

TECHNICAL AND METHODOLOGICAL DOCUMENTATION

<u>User manual – Ranker v. 1.2</u>

ARGUS

Description

Istat participated in the European project, CASC (Computational Aspects of Statistical Confidentiality) whose goal was the development of the software ARGUS for the protection of the data confidentiality in the release of statistical information.

The ARGUS software consists of two modules Mu -ARGUS, for elementary data, and Tau-ARGUS, for tables, and contains many of the methods of protection proposed in the literature.

Information

Status:	validated
Author:	CASC
GSBPM code:	6.4 Apply disclosure control
Keywords:	confidentiality, risk of identification, value at risk, suppression of value, risk rule

Software and documentation

To get the ARGUS software as well as the methodological and technical documentation, please contact <u>CASC</u>.