

TEIDE: A new software for data editing

Juan José Salazar González
Universidad de La Laguna
Tenerife - Spain

-

Work partially funded by

"Instituto Canario de Estadística" & "Instituto Andaluz de Estadística"

UNECE, Bonn, September 26, 2006

Content

- 1 Introduction
 - Preface
 - Previous experience
 - What is TEIDE?
 - Implemented methodologies
 - What is TEIDE 1?
 - What is TEIDE 2?
 - A few details
- 2 On-line demonstration
 - Statistical agencies using TEIDE
 - Computational experiments
- 3 Conclusions

You are here. . .

- 1 Introduction
 - Preface
 - Previous experience
 - What is TEIDE?
 - Implemented methodologies
 - What is TEIDE 1?
 - What is TEIDE 2?
 - A few details
- 2 On-line demonstration
 - Statistical agencies using TEIDE
 - Computational experiments
- 3 Conclusions

Limitations

Aiming to replace humans by an automatic guided software is impossible when editing statistical data.

However, using hardware and software easily available today, we can reduce part of the effort, specially when the same operations (checks, computations,...) must be done on hundreds or thousands of records.

Then, the saved effort can be devoted to improve the more complex and challenging parts of the survey.

As noticed by Tom de Waal ("Processing of Erroneous and Unsafe Data", 2003), there is a narrow relationship between Statistical Disclosure Control and Statistical Data Editing:

- Finding suppressions in SDC is like finding errors in SDE.
- Finding the expected attacker value of a suppression is like finding the imputation of an error.

University of La Laguna (Tenerife) was involved in SDC:

- European project, SDC 1996-1998
- European project, CASC 2001-2003
- U.K. project funded by ONS, 2004

What is TEIDE?

TEIDE is the name of the highest mountain in Spain, which is a vulcan in Tenerife with 3718 meters.



What is TEIDE?

In our context, it is a software for editing and imputation:

- 1 **T**écnicas de **E**dición e **I**mputación de **D**atos **E**stadísticos comes from "Técnicas de Edición e Imputación de Datos Estadísticos", which means "Techniques for Editing and Imputation of Statistical Data".
- 2 It works on continuous, categorical, and mixed data.
 - Discrete
 - list ($x_i \in \{v_{i1}, v_{i2}, \dots, v_{il_i}\}$)
 - range ($x_i \in \{v_i^{min}, \dots, v_i^{max}\}$)
 - Continuous ($x_i \in [v_i^{min}, v_i^{max}]$)
- 3 The edits for categorical data are logical constraints, like:
 - Arithmetic ($x_1 = 2 * x_4 + x_2$)
 - Logical (*if* ($x_1 = a$) *then* ($x_3 \leq b$))
 - Mixed (*if* ($x_4 * x_2 = 6$) *then* ($x_1 = 1$))

Assumptions in TEIDE

- 1 The edits for numerical data must be linear equalities or inequalities.
- 2 On categorical data, unfeasible edits are automatically detected. On numerical data, no yet!
- 3 Both on categorical data and numerical data, the syntaxis of all edits are checked. Invalid edits are detected, so the user can redefine them. Redundant edits are not checked (no yet!)
- 4 Data are automatically checked according to the edits, but data cannot be modified by the users inside TEIDE.

Methodologies under TEIDE

- 1 The syntaxis of all edits are checked.
- 2 Validation: Each record is validated on each edit, individually.
- 3 Error localization: the minimum number of modifications is computed through optimal or near-optimal approaches – Fellegi y Holt (JASA 1976)
- 4 Imputation: Donor and regression are optional to all (categorical or numerical) variables

How "TEIDE 1" is being implemented?

- It is written in C++ Programming Language
- Compiled and linked in Borland C++ Builder
- Running under Microsoft Windows
- It is a stand-alone executable, so it can be copied into a new computer and executed.
- Read and Write in Microsoft Access database.
- It is free software.
- Download a copy from <http://www.goma.u11.es>

How "TEIDE 2" is being implemented?

- It is written in C++ Programming Language
- It is open source
- It can be compiled on any C++ compiler, and therefore, it runs on any computer under any operating system (including Windows, Linux, Mac, Sun, HP, ...)
- no virtual machine is required (no JAVA, no .NET,...)
- Read and Write in XML format.
- For the moment we have a beta version in progress.

Finishing this version depends on the economical support.

How are displayed the edits

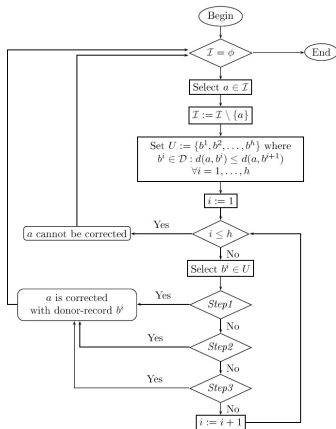
Edits are located in three groups:

- 1 type: checking absolute values
- 2 filters: checking values relative to other values
- 3 general: standard edits in general form

How the process works

- 1 The input data are loaded
- 2 The edits are evaluated
- 3 Registers are classified into two groups: correct and incorrect
- 4 Individual data of an invalid record are classified in different status: wrong values, values in wrong edits, values in edits with variables in wrong edits,...
- 5 Different weights are assigned to individual values depending on the status, the importance of the values, etc.
- 6 Weights will be used to define the distance between two records (typically between an incorrect and a correct record) when using the donor technique, and to be used in a potential regression analysis.

Flowchart



You are here. . .

- 1 Introduction
 - Preface
 - Previous experience
 - What is TEIDE?
 - Implemented methodologies
 - What is TEIDE 1?
 - What is TEIDE 2?
 - A few details
- 2 On-line demonstration
 - Statistical agencies using TEIDE
 - Computational experiments
- 3 Conclusions

Real-world surveys:

- TEIDE was developed and improved by working on real-world surveys at “Instituto Canario de Estadística” (ISTAC) and at “Instituto Andaluz de Estadística” (IAE).
- ISTAC: “Encuesta de Condiciones de Vida de Hogares e Individuos Canarios”, “Encuesta de Salud”, “Encuesta de Turismo”, “Cesta de la Compra”, “Encuesta de Tecnología en Hogares”.
- IAE: “Encuesta Mundial de Valores”, “Encuesta Social”

We thank these organizations for supporting TEIDE

Some benchmarks:

Block	Information	Survey		
		HEALTH	HOUSE	INDIVIDUAL
Description	<i>#variables</i>	375	176	234
	<i>#records</i>	5633	7797	22584
	<i>#microdata</i>	2112375	1372272	5284656
	<i>#real microdata</i>	1129379	1108119	2035238
	<i>#continuous variables</i>	2	12	16
	<i>#discrete variables</i>	373	164	218
	<i>#imputable variables</i>	335	128	210
	<i>#filter edits</i>	283	42	207
	<i>#general edits</i>	33	7	45
	<i>Loading-Checking time</i>	105.29	69.43	251.79
Imputation	<i>Total time</i>	3835.47	258.14	5573.30
	<i>#donor records</i>	4138	5049	11915
	<i>#correctable records</i>	1495	2748	10669
	<i>#non-corrected records</i>	5	0	0
	<i>#corrected records</i>	1490	2748	10669
	<i>#warning records</i>	30	19	49
	<i>affected var. average</i>	7.92	7.16	12.08
	<i>imputation var. average</i>	5.22	6.44	11.13
	<i>wrong-range var. average</i>	1.47	5.93	9.75
	<i>wrong-edits var. average</i>	6.73	1.69	4.34
	<i>connected var. average</i>	225.13	5.50	195.63
	<i>donor record distance average</i>	0.98	1.10	1.35
Overall time		3965.94	341.16	5899.28

You are here. . .

- 1 Introduction
 - Preface
 - Previous experience
 - What is TEIDE?
 - Implemented methodologies
 - What is TEIDE 1?
 - What is TEIDE 2?
 - A few details
- 2 On-line demonstration
 - Statistical agencies using TEIDE
 - Computational experiments
- 3 Conclusions

Finally:

- Download TEIDE 1 from <http://www.goma.u11.es>
- University of La Laguna (Tenerife) needs economical support to keep working on TEIDE 2, the open-source version.
- If you need help when using TEIDE:
Contact me to jjosalaza@ull.es, or convince your boss to visit me in Tenerife.