

**UNITED NATIONS STATISTICAL COMMISSION and  
ECONOMIC COMMISSION FOR EUROPE**

**CONFERENCE OF EUROPEAN STATISTICIANS**

**Work Session on Statistical Data Editing**  
(Neuchâtel, Switzerland, 5-7 October 2009)

Topic (i): Automated editing and imputation and software applications

**TEIDE 2: A modern tool for Editing and Imputation of Statistical Data**

Supporting Paper

Prepared by Juan José Salazar González and M<sup>a</sup> Salomé Hernández García, University of La Laguna,  
Tenerife, Spain

**I. INTRODUCTION**

1. Techniques for editing and imputation of data are used to detect and correct "errors" or "lost data" contained in surveys. This is a complex process that needs sophisticated and effective procedures to guarantee good solution quality. Therefore it is essential to take advantage of modern technology, including computational tools (computers, programming languages, databases, GUI, etc.).

2. The software TEIDE ("Técnicas de Edición e Imputación de Datos Estadísticos", which means "Techniques for Editing and Imputation of Statistical Data") was born in order to make faster, reliable and more comfortable several tasks related to Statistical Data Editing. The software was designed and implemented by a research team at University of La Laguna (Tenerife, Spain) and as a result it:

- is multi-platform (32-bits and 64 bits running under Linux, Windows, Mac,...).
- Makes use only of free software (C++, Qt library).
- Reads and writes in standard formats including XML, Microsoft Access and ORACLE.

3. Since an early version TEIDE1 started in 2003, the software has been considerably improved thanks to the use of it by several statistical agencies, including ISTAC (Statistical Office of Canary Islands) and IEA (Statistical Office of Andalucía). TEIDE has been used to debug and clean different real-world surveys, and has learned a lot from these experimentations. The current version is TEIDE2. It works on surveys with both quantitative and qualitative variables, allowing a very wide variety of consistency rules. For the imputation TEIDE2 combines the standard donor record procedure with a regression analysis, both based on a subset of correct records selected from the whole survey. Many parameters allow a user to control and debug the whole process through a set of friendly windows.

**II. DEVELOPMENT**

4. TEIDE2 was used to debug and clean the "Survey on Income and Living Conditions of Households in Canary Islands", "Survey on the Implementation of Information and Communication Technologies in Households in Canary Islands", "Survey on the Implementation and Use of Technology in Information and Communication in Business in Canary Islands", "Tourist Expenditure Survey", and others.

5. TEIDE 2 is a new version, totally re-implemented after our experiences using the first version on several real-world surveys. The old version TEIDE1 was exclusive for computers running Windows and Microsoft Access, and compiled and linked using the commercial software Borland C/C++.. The new version TEIDE2 has been re-implemented in standard C++ and linked only using the free Qt library.

5.1. C++ is a general-purpose programming language. It is regarded as a middle-level language, as it comprises a combination of both high-level and low-level language features. It was developed by Bjarne Stroustrup starting in 1979 at Bell Labs as an enhancement to the C programming language and originally named "*C with Classes*". It was renamed to C++ in 1983.

C++ is widely used in the software industry. Some of its application domains include systems software, application software, device drivers, embedded software, high-performance server and client applications, and entertainment software such as video games. Several groups provide both free and proprietary C++ compiler software, including the GNU Project, Microsoft, Intel, Borland and others.

5.2. Qt is a cross-platform application development framework, widely used for the development of GUI programs (in which case it is known as a *widget toolkit*), and also used for developing non-GUI programs such as console tools and servers. Qt is most notably used in KDE, Google Earth, Skype, Qt Extended, Adobe Photoshop Album, VirtualBox and OPIE. It is produced by Nokia's Qt Software division, which came into being after Nokia's acquisition of the Norwegian company Trolltech, the original producer of Qt, on June 17, 2008.

Qt comes under the following licenses:

<b>Qt Commercial Version</b>	The Qt Commercial version is the appropriate version to use for the development of proprietary and/or commercial software. This version is for developers who do not want to share the source code with others or otherwise comply with the terms of the GNU Lesser General Public License version 2.1 or GNU GPL version 3.0.
<b>Qt GNU LGPL v. 2.1</b>	This version of Qt is appropriate for the development of Qt applications (proprietary or open source) provided you can comply with the terms and conditions contained in the GNU LGPL version 2.1.
<b>Qt GNU GPL v. 3.0</b>	This version of Qt is appropriate for the development of Qt applications where you wish to use such applications in combination with software subject to the terms of the GNU General Public License version 3.0 or where you are otherwise willing to comply with the terms of the GNU General Public License version 3.0.

A clear advantage of working with TEIDE2 is that it is all open and free software. Then it is the ideal framework to implement and test new methods, and compare how these new methods works respect to other previously implemented.

### III. WAY OF USE

6. The application TEIDE2 reads a metafile which should contain the name of the database Oracle, Microsoft Office Access (.mdb) or in format XML. These databases contain the necessary tables:

*Variables:* The variables that we have to study, with the necessary information (ID: identifying of the variable; NAME: name of the variable; INFO\_VARIABLE: information about the meaning of the variable; TIPO: it indicates if it is discreet in range, discreet in list or constant; RANGO: values that it can take (of a list or of a range of values); FILTRO: necessary filter to indicate when the variable must be NO\_PROCEDE; INFO\_FILTRO: information about the meaning of the filter; SENTIDO\_FILTRO: be able to take the values a, b or c "a) if (non-filter) then (valor=NP)" "b) if (filter) then (value! =NP)" "c) a and b"; IMPUTABLE: it indicates if variable debit or not to be imputed; NO\_PROCEDE: it indicates if it can or not to take the value NO\_PROCEDE; NO\_SABE: it indicates if it can or not to take the value NO\_SABE; NO\_CONTESTA: it indicates if it can or not to take the value NO\_CONTESTA; NS\_NC: it indicates if it can or not take the value NO\_SABE\_O\_NO\_COTESTA; PESO: it indicates the weight that this variable takes with regard to other; MAPPING: the table indicates the mapping and contains the meaning of the values that this variable may take; IMP\_NUM: type of imputation, by default it is the record donor).

*Microdata:* it contains all the values that take the variables in every record.

*Edits:* they are the rules that have to meet the variables. The table contains ID: identifying of the edit; CONDICION: the edit; DESCRIPCION: information about the meaning of the edit.

*Missing*: it is the table that contains the values missing (lost) that takes the variables.

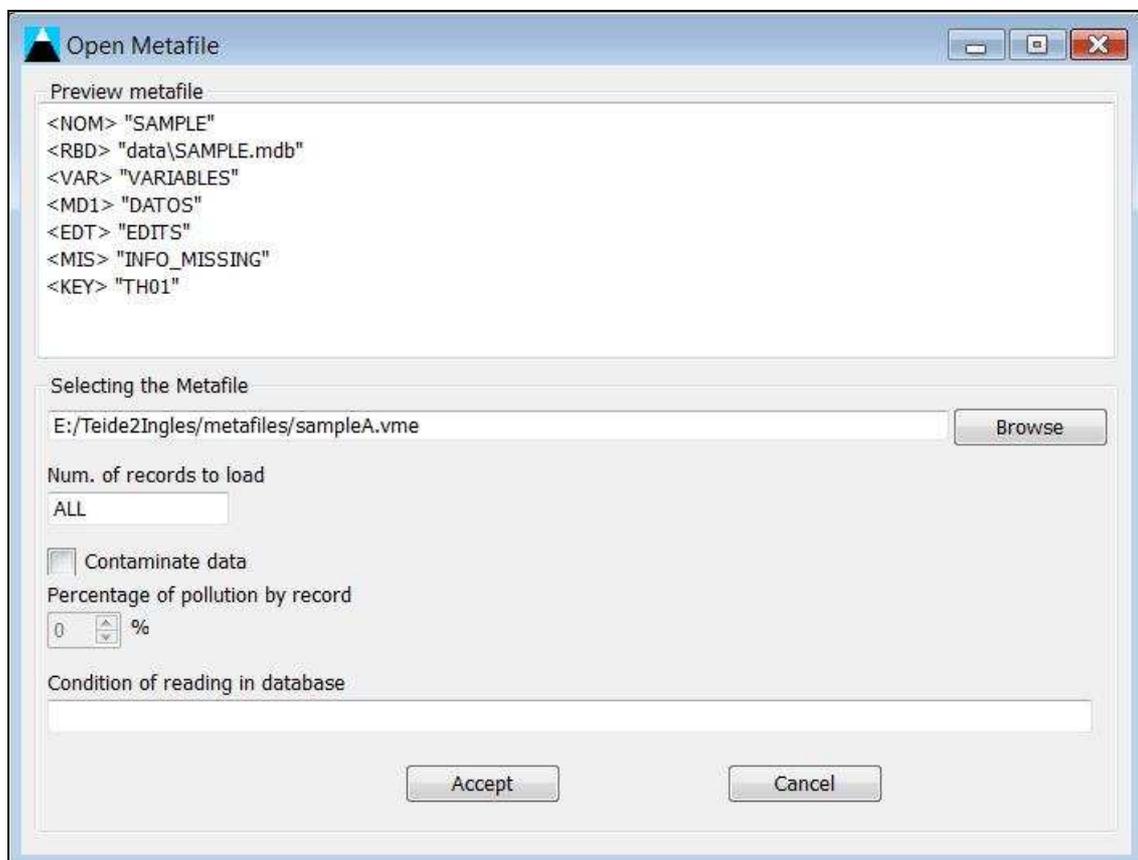
*Mapping*: they are the tables with the meaning of every value that takes this variable.

7. As soon as the application is executed the user has to choose an option of the menu that appears later:



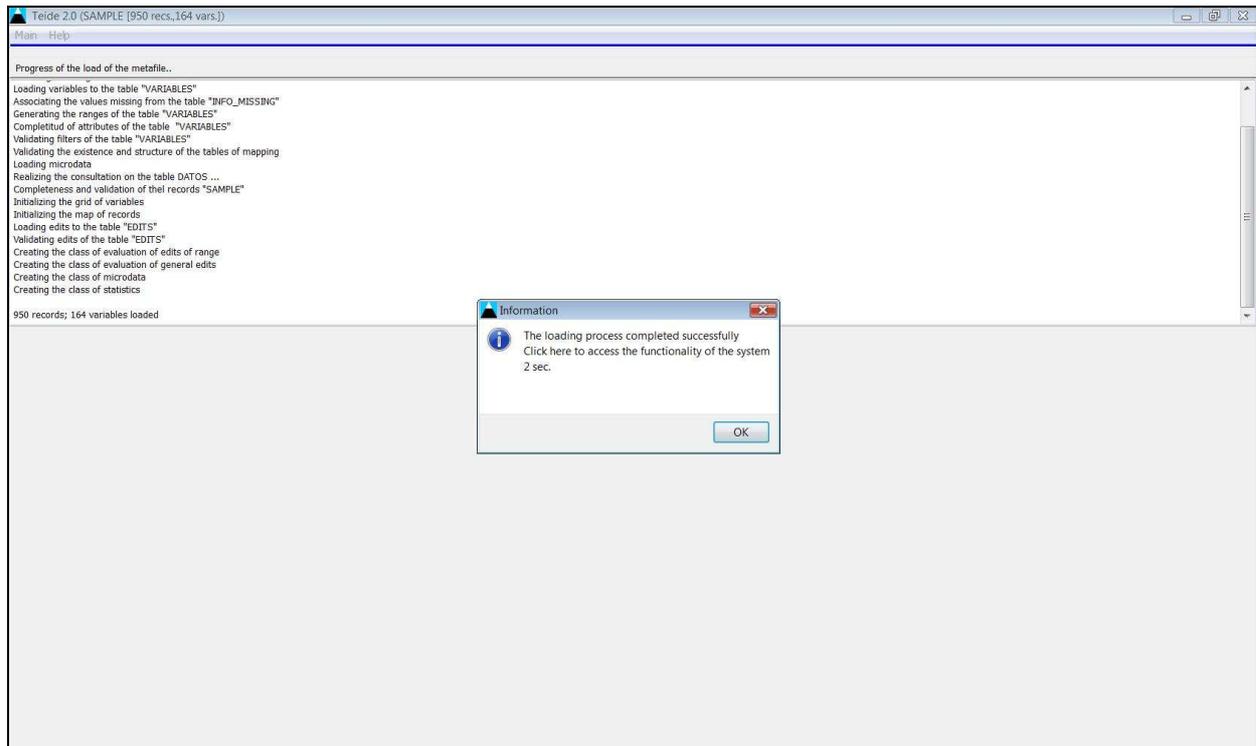
Menu.

8. The user has to open a metafile created in order to load data.



Open metafile.

9. When the metafile is selected, it begins the process of load in which the program will read all the tables and will show the information in the tab: Variables, Microdata and Edits.



Progress of the load.

	NOMBRE	INFO_VARIABLE	RANGO	FILTRO	INFO_FILTRO	SENTIDO_FILTRO	IMPUTABLE
2	TC001A		1:20				False D
18	TUH001A		1:5			a	True D
19	TUH002A		1:6	TUH001A > 1		a	True
20	TUH002B		1:6	TUH001A > 1		a	True
21	TUH002C		1:6	TUH001A > 1		a	True
22	TUH002D		1:6	TUH001A > 1		a	True
23	TUH002E		1:6	TUH001A > 1		a	True
24	TUH002F		1:6	TUH001A > 1		a	True
25	TVV001A		1890:2005			a	True D
26	TVV002A		1:6			a	True D
27	TVV003A		1:3	TVV002A > 1		a	True D
28	TVC001A		1:7			a	True D
29	TVC002A		1:20			a	True D
30	TVC003A		30:600			a	True D
31	TVC003B		30:350			a	True D
32	TVC004AA		1:6			a	True
33	TVC004AB		1:6			a	True
34	TVC004AC		1:6			a	True
35	TVC004AD		1:6			a	True
36	TVC004AE		1:6			a	True
37	TVC004AF		1:6			a	True
38	TVC004AG		1:6			a	True
39	TVC004AH		1:6			a	True
40	TVC004AI		1:6			a	True
41	TVC004AL		1:6			a	True
42	TVC004AM		1:6			a	True
43	TVC004AN		1:6			a	True
44	TVC004BG		1:10	TVC004AG = 1		a	True D
45	TVC004BM		1:10	TVC004AM = 1		a	True D
46	TVC005A		1:6			a	True D
47	TVC005B		1:6			a	True D
48	TVC006A		1:6			a	True
49	TVC006B		1:6			a	True
50	TVC006C		1:6			a	True
51	TVC006D		1:6			a	True
52	TVC006E		1:6			a	True
53	TVC006F		1:6			a	True
54	TVC006G		1:6			a	True
55	TVC006H		1:6			a	True
56	TVC006I		1:6			a	True
57	TVC006J		1:6			a	True

Tab Variables.

10. In the tab of variables, there appears the information of every variable: name, range, filter (in case it has it), if it is attributable, etc.

11. The tab of microdata shows all the records in the survey and the values that take the variables in each of them. Also it shows information about the selected variable.

The screenshot shows the 'Microdata' tab in the Teide 2.0 application. It displays a grid with columns for variables (TC001A, TUH001A, TUH002A, TUH002B, TUH002C, TUH002D, TUH002E, TUH002F, TVV001A, TVV002A, TVV003A, TVC001A, TVC002A, TVC003A, TVC003B, TVC004A) and rows for individual records (AW00001 to AW00037). The values are mostly -1, 1, or 2, with some missing values (indicated by '...').

Tab Microdata.

12. The tab of edits shows all the consistency rules. It indicates if they are correct or not. It allows to modify, erase or insert a new edit in an easy form.

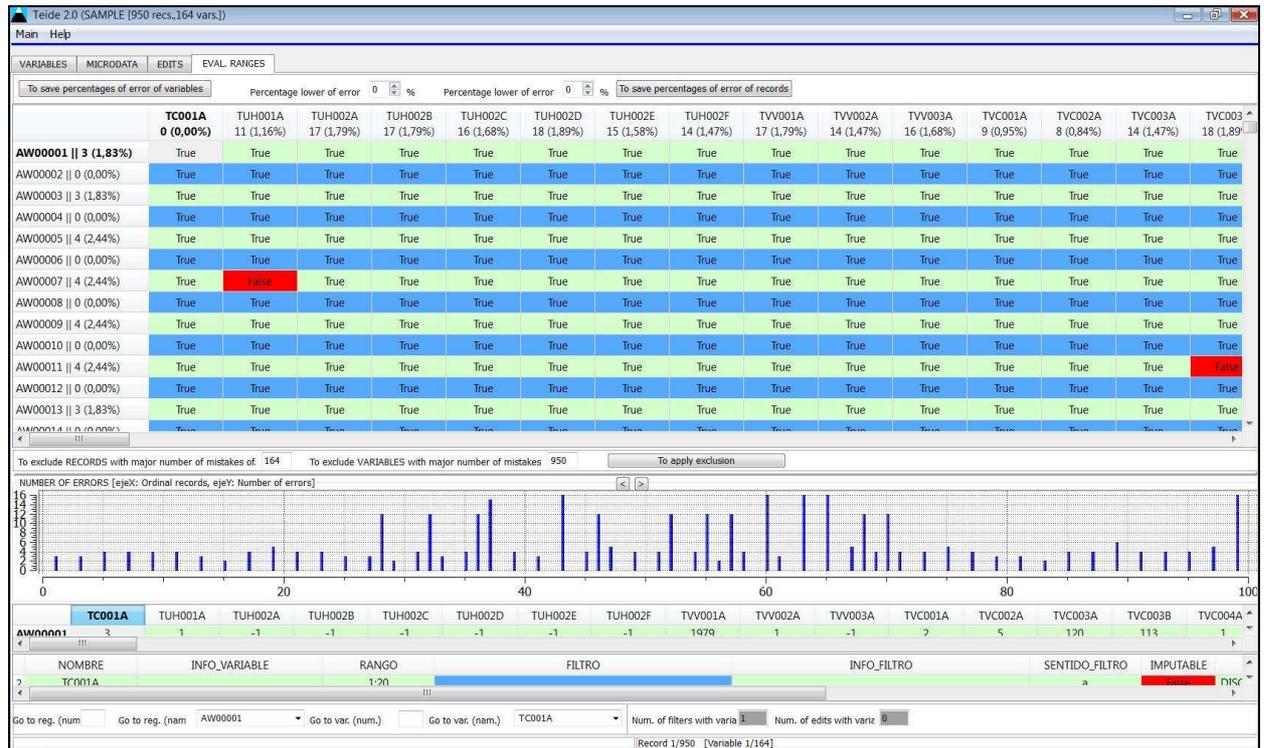
13. Once one has seen the information showed after the process of load, the application allows that it should follow with the process of edition and imputation. It is possible to do it in one over two ways: stepwise or complete form. If the chosen option is stepwise, the user will be able to see each of the tabs that are created. If the option is complete form, the user will not be able to stop neither the process of edition nor that of imputation, and the application will show it all the tabs once it has finished the whole data-editing process.

The screenshot shows the 'Edits' tab in the Teide 2.0 application. It displays a list of consistency rules under the heading 'CONDICION' and their corresponding descriptions under 'DESCRIPCION'. The rules are numbered 1 through 7.

CONDICION	DESCRIPCION
1 IF (TVC003B > 0 AND TVC002A > 0) THEN ((TVC003B >= 3*TVC002A) AND (TVC003B <= 60*TVC002A))	Descripción edit
2 IF (TVC003B > 0 AND TVC003A > 0) THEN TVC003B <= TVC003A	"
3 IF TVC004AM = 1 THEN TVC004AB = 1	"
4 IF (TVE001CA = 1 OR TVE001DA = 1 OR TVE001EA = 1) THEN TVE001AA = 1	"
5 IF TVE001NA = 1 THEN TVE001MA = 1	"
6 IF TVE001NB = 2 THEN (TVE001KA = 1 OR TVE001LA = 1)	"
7 IF TVE002A = 1 THEN (TVE002B >= 1 OR TVE002B = -3 OR TVE002B = -7 OR TVE002B = -9)	"

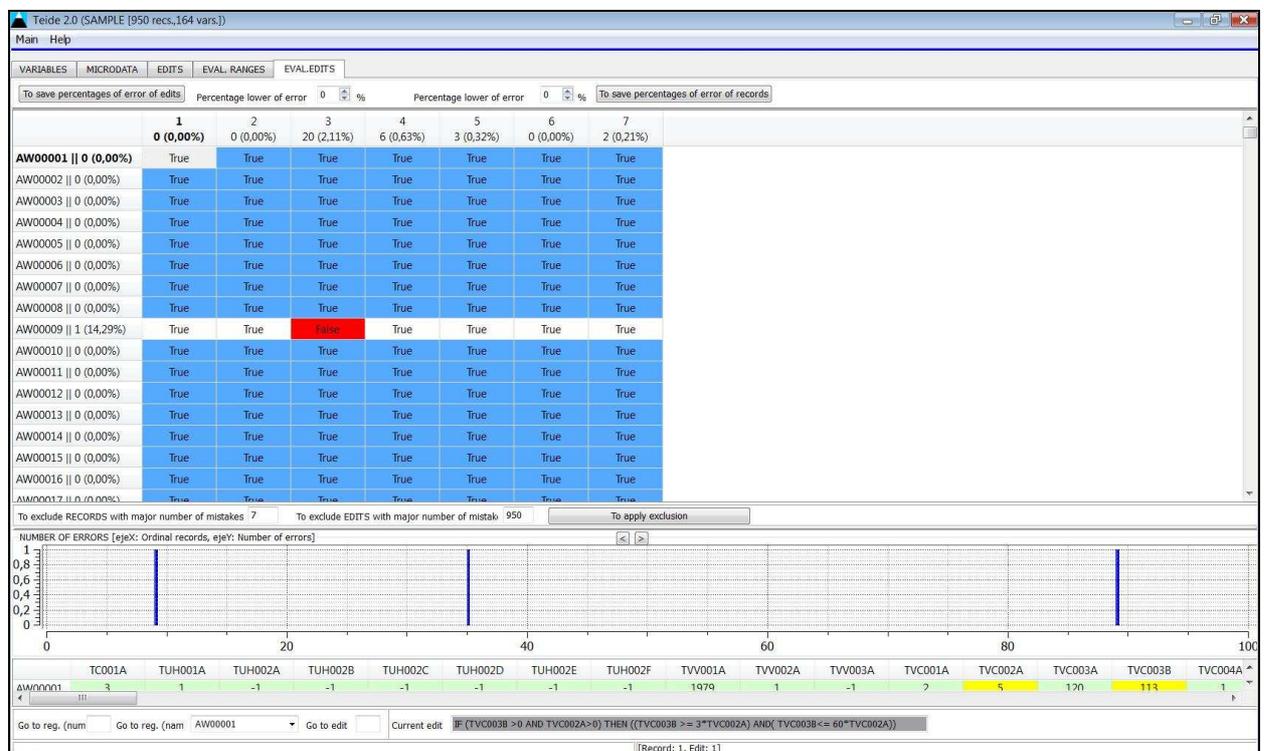
Tab Edits.

14. The following tab shows the evaluation of variable ranges, in which all the records will appear with variables in colors, indicating if the range or the filter are correct or not.



Tab Evaluation of Ranges.

15. When the ranges and filters have been studied, the following step is the study of the edits. The tab of evaluation of edits shows what records with it variables fulfill or do not fulfill the proposed edits.



Tab Evaluation of Edits.

16. Once finished the processes of evaluation, the following step is that of imputation. The tab of imputation shows for all the records and it variables the following information:

- Records donors, in blue color.
- Incorrect records (they could not have corrected), in red color.
- Values of an incorrect record that TEIDE could not have corrected but that can that he is the culprit of whom this record is incorrect, in color magenta.
- Excluded records, in orange color.
- Incorrect records corrected in white color.
- Values of an incorrect record that TEIDE has modified, in yellow color.

Teide 2.0 (SAMPLE [950 recs,164 vars.])

Man Help

VARIABLES MICRODATA EDITS EVAL\_RANGES EVAL\_EDITES IMPUTATION

	TC001A	TUH001A	TUH002A	TUH002B	TUH002C	TUH002D	TUH002E	TUH002F	TVV001A	TVV002A	TVV003A	TVC001A	TVC002A	TVC003A	TVC003B	TVC004A
AW00001	3	1	-1	-1	-1	-1	-1	-1	1979	1	-1	2	5	120	113	1
AW00002	4	1	-1	-1	-1	-1	-1	-1	1991	1	-1	2	2	77	72	1
AW00003	2	3	6	1	1	1	1	1	1984	2	2	2	3	92	86	1
AW00004	3	1	-1	-1	-1	-1	-1	-1	1963	1	-1	2	3	92	86	1
AW00005	3	1	-1	-1	-1	-1	-1	-1	1969	2	2	2	8	164	155	1
AW00006	4	1	-1	-1	-1	-1	-1	-1	1984	1	-1	2	4	107	100	1
AW00007	1	1	-1	-1	-1	-1	-1	-1	1954	2	2	2	3	92	86	1
AW00008	5	1	-1	-1	-1	-1	-1	-1	1999	3	3	2	5	120	113	1
AW00009	1	1	-1	-1	-1	-1	-1	-1	1996	2	2	2	4	80	78	1
AW00010	2	1	-1	-1	-1	-1	-1	-1	1981	2	2	2	4	107	100	1
AW00011	1	1	-1	-1	-1	-1	-1	-1	1951	2	2	2	2	77	72	1
AW00012	4	1	-1	-1	-1	-1	-1	-1	1982	2	2	2	4	107	100	1
AW00013	4	1	-1	-1	-1	-1	-1	-1	1972	1	-1	2	4	107	100	1
AW00014	5	1	-1	-1	-1	-1	-1	-1	1959	1	-1	2	4	107	100	1
AW00015	2	1	-1	-1	-1	-1	-1	-1	2004	3	3	2	4	107	100	1
AW00016	5	1	-1	-1	-1	-1	-1	-1	1962	1	-1	2	3	92	86	1
AW00017	2	1	-1	-1	-1	-1	-1	-1	1971	1	-1	2	4	107	100	1
AW00018	2	1	-1	-1	-1	-1	-1	-1	1954	1	-1	2	5	120	113	1
AW00019	2	1	-1	-1	-1	-1	-1	-1	1953	1	-1	2	6	135	127	1
AW00020	2	1	-1	-1	-1	-1	-1	-1	1996	3	3	2	3	92	86	1
AW00021	2	1	-1	-1	-1	-1	-1	-1	1949	1	-1	2	4	150	100	1

Legend: Original data (white), Atributable data (yellow), Rec. correct (blue), Rec. exclude (orange), Rec. incorrect (red), Possible incorrect data (magenta)

AW00021: 2, 1, -1, -1, -1, -1, -1, -1, -1, 1949, 1, -1, 2, 4, 150, 100, 1

Browser rec. donors: Distance to the current rec: 1,01 CHANGES DONOR RECORD: # Rec. donor: 1/436

NOMBRE INFO\_VARIABLE RANGO FILTRO INFO\_FILTRO SENTIDO\_FILTRO IMPUTABLE

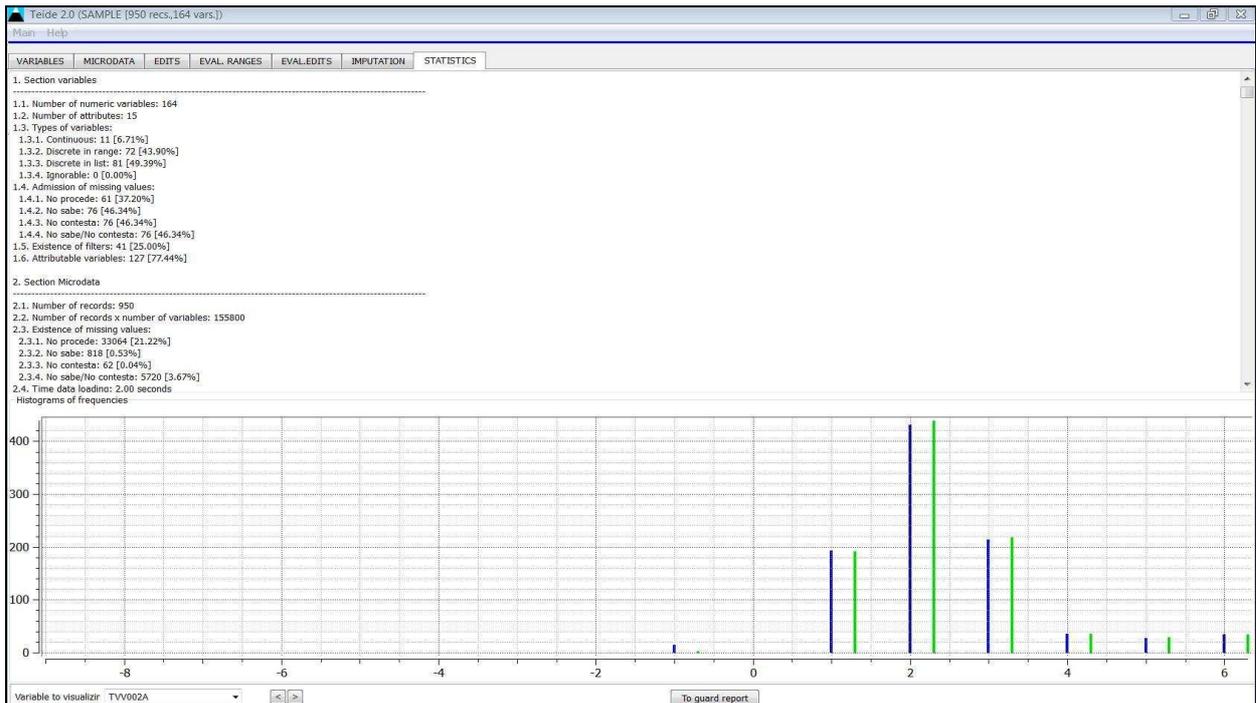
TC001A TUH001A TUH002A TUH002B TUH002C TUH002D TUH002E TUH002F TVV001A TVV002A TVV003A TVC001A TVC002A TVC003A TVC003B TVC004A

Go to reg. (num): AW00001 Go to var. (num.): TC001A #Var. Bas: 5 #Var. Ext: 7 #Var. Imp: 3 Original data: 3

950 Records, 164 Variables Record 1/950 [Variable = 1]

Tab Imputation.

17. When the edition and imputation phases have finished TEIDE shows the tab of statistics (a report) which shows all the information relative to the obtained results. As well as a graph in which it is possible to see for every variable the values that had before and that have now.



Tab Statistics.

18. The application displays the results and computational statistics, i.e. the errors found and corrections made. At this stage we analyze the performance of executive tasks and data, to decide whether the results have an acceptable quality level. We analyze and evaluate the coverage obtained, the level of non-response, accuracy of responses, errors occurring during the recording, the level and distribution of errors made and the imputations.

#### IV. TESTS OF COMPILATION AND EXECUTION

19. To give a proper idea on how TEIDE2 works on a real-world survey, here we illustrate technical details using TEIDE2 on one survey. Due to confidentiality issue we cannot give the survey for further experiments. The database consists of 289 variables, 2789 records and 66 edits. Most of the variables are qualitative and the rest are quantitative. We conducted tests on two machines with the following characteristics:

Intel Core 2 Duo 3 GHz. Memory (RAM) 4 GB Operating Systems Ubuntu 9.04 64 bits Windows Vista Business 64 bits
Intel Core 2 Duo 3 GHz Memory (RAM) 3GB Operating Systems Ubuntu 9.04 32 bits Windows Vista Home 32 bits

As for C/C++ compilers we have used:

In Windows: Qt includes MinGW compiler (g++). We have also used Microsoft Visual Studio 2008 to obtain a different executable on the same computer. The MinGW which provides Qt is 32 bits and they do not have a 64-bit version. However, with Visual Studio we have a version of 32 and 64 bits.

In Linux we have been used g++ and icpc (Intel compiler c++). Both are 32 and 64 bits.

Executables have been compiled on two different machines, with a 32-bit OS for 32 bit executable and 64-bit OS for 64 bit executable. All tests were performed in 64-bit OS.

##### A. WINDOWS

20. Time of load of data

MinGW 32 bits	Visual Studio 32 bits	Visual Studio 64 bits
13 seconds	8 seconds	8 seconds

Process time in the section of ranges

MinGW 32 bits	Visual Studio 32 bits	Visual Studio 64 bits
7 seconds	5 seconds	4 seconds

Process time in the section of test

MinGW 32 bits	Visual Studio 32 bits	Visual Studio 64 bits
7 seconds	6 seconds	5 seconds

Process time in the section of imputation

<b>MinGW 32 bits</b>	<b>Visual Studio 32 bits</b>	<b>Visual Studio 64 bits</b>
1741 seconds	1356 seconds	1152 seconds

## B. LINUX

### 21. Time of load of data

<b>G++ 32 bits</b>	<b>ICPC 32 bits</b>
10 seconds	10 seconds

<b>G++ 64 bits</b>	<b>ICPC 64 bits</b>
9 seconds	8 seconds

Process time in the section of ranges

<b>G++ 32 bits</b>	<b>ICPC 32 bits</b>
10 seconds	10 seconds

<b>G++ 64 bits</b>	<b>ICPC 64 bits</b>
22 seconds	22 seconds

Process time in the section of test

<b>G++ 32 bits</b>	<b>ICPC 32 bits</b>
8 seconds	8 seconds

<b>G++ 64 bits</b>	<b>ICPC 64 bits</b>
20 seconds	20 seconds

Process time in the section of imputation

<b>G++ 32 bits</b>	<b>ICPC 32 bits</b>
1251 seconds	1165 seconds

<b>G++ 64 bits</b>	<b>ICPC 64 bits</b>
1053 seconds	1026 seconds

## C. COMMENTS

22. According to this in Windows it seems that the best version is compiled in Visual Studio 2008 to 64 bits, resulting in the imputation phase a difference of nearly 600 seconds with respect to MinGW compiled to 32-bit.

Linux is the best compiled with icpc to 64 bits.

## V. RESULTS OF THE IMPUTATION

23. The database is read in XML as it is the only format that currently supports all versions. The results were as follows:

<b>Donor records</b>	2023
<b>Records to correct</b>	766
<b>Correct records</b>	671
<b>Incorrect records</b>	95

<b>Records warning</b>	33
------------------------	----

<b>Average imputed variables per record (total)</b>	1.94
<b>Average imputed variables per record (no reg. warning)</b>	1.56
<b>Average of errors in range per record</b>	0.00
<b>Average of variables involved in edits incorrect for record</b>	3.18
<b>Average of variables involved in total mistakes for record</b>	3.18
<b>Average of variables involved in comp. connected with mistake for record</b>	37.51
<b>Average distance to donor records</b>	0.06

Over the total (2789 records) 72.53% are correct records, thus 27.47% records have been modified.

## VI. CONCLUSIONS AND FUTURE

24. TEIDE2 can impute great variety of surveys, with quantitative and qualitative variables. Since the beginning of the data collection, TEIDE2 is a nice tool to start checking that the collecting is being done as expected, or to detect any kind of unexpected misunderstanding on the query or other type of fundamental errors. The whole editing and imputation process is carried out step by step so that the statistician can see the results of each study or modification. It is also possible to make a complaint about a given fact. There are different parameters to drive the imputation process as desired. TEIDE2 will never remove the human in the data editing process, but will help him/her to make the best of his/her time. TEIDE2 is also a great simulator to measure the impact of different consistent rules on a data survey.

Any application can be improved and in this case the future tasks we intend to address are:

- Reduce processing times. This will require thinking about improvements in the algorithms of evaluation of expressions and data structures.
- Improve gradually the imputation procedures that are being used.
- Exploit the parallelization that it is now available in my computers to make our data-editing algorithm faster.
- Add new formats to better integrate TEIDE2 with other software.
- Test and compare TEIDE2 with other software on the same data surveys.

25. We invite statistical agencies to download and check our software TEIDE2 from our webserver [www.goma.ull.es](http://www.goma.ull.es), and to help us to improve it with new experiments on different surveys. Since TEIDE2 is a self-contained software, fully implemented inside a research team in a public university, it is the best framework to also implement and test different imputation techniques. Therefore, please, feel free to suggest the authors all type of suggestions that would make TEIDE2 more useful for the data editing processes in your statistical agency. Contact [jjisalaza@ull.es](mailto:jjisalaza@ull.es) for further details.

26. As a consequence of the economical crisis in all Europe, we are currently finding troubles to find economical support to continue improving TEIDE2. We would be gratefully to join any kind of proposal to collaborate in any National or European funded project or contract.

27. In summary, our objective is to show and share with you an application (TEIDE2) which is now faster, more flexible and more reliable to automate as far as possible the process of purification of a survey. In this way, the statistical institute can devote their resources to explore more deeply sensitive and more specific tasks.