

Business Intelligence Tools for Statistical Data Analysis

Savoska Snezana, Manevska Violeta

Faculty for Administration and Management of Information systems, University "St.Kliment Ohridski" Bitola, Partizanska bb, 7000 Bitola

snezana.savovska@uklo.edu.mk, violeta.manevka@uklo.edu.mk

Abstract. Dealing with the up flow of information is one of the most significant tasks of the information systems nowadays, especially when they are official information from state institutions. Getting information from all state sectors is a very complicated task, dedicated to the State Statistical Office (SSO). This task cannot function without using highly sophisticated Business Intelligence (BI) tool for Data warehousing and information visualization. The objective of this paper is to elaborate the study on utilization of BI tools for complex statistical calculations in the SSO in the Republic of Macedonia. This tool has been implemented in SAS software modules, sophisticated statistical software for data extracting, cleansing, transforming and loading in Data Warehouse (DW).

Keywords. Business intelligence, Data Warehouse, Statistical data analysis.

1. Introduction

The SSO DW is the most professional statistical usage of DW for BI in the Republic of Macedonia. It is installed on the SAS DW software environment with precisely defined procedures for extraction, cleansing, transformation and data loading. The physical model of SSO DW is shown on Figure 1.

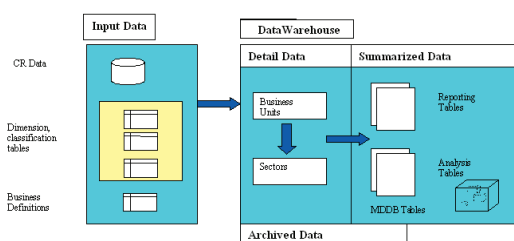


Figure 1: Physical model of State Statistical office DW

Data which is an object of data processing in the DW is taken from different administrative sources, areas from where data is obtained with statistical researches made in SSO and from Statistical Registry. They are extracted from many sources, prepared with manual or automated software procedures, transformed in specific code pages and loaded in temporary DW files.

2. The Data Warehouse content

The administrative data sources are: The State central registry (the final account report of subjects), The office of public finance (the taxes-information database), ministry of finance (public revenue – the treasury accounts), National bank (monetary statistics), Sanitary insurance fund (contribution of sanitary insurance) as Retirement and invalid insurance fund (contribution for retirement and invalid insurance).

The statistical researches are made in the areas of: Inside trade market, foreign trade market, industry, building and construction sector, tourism, agriculture, poll for household spending, prize statistic and poll for labour.

The data which are loaded in DW are taken from different data sources, working in different operating systems and platforms. First, they are processed with Microsoft Visual Studio.NET (Visual basic or ACCESS) or other preparation procedures. They use SAS Data integration studio for data import (Figure 2), selection and transformation. It is a very powerful tool, which can provide integration of diverse databases laid in different data servers or from external sources. Also, it can provide creating of multidimensional cubes, modelling of software processes and creating calculating methods for summarized (aggregated) tables and DW. These procedures are used for finding irregular or missing data, transformation of these data, conversion of Macedonian code page in ACSII cods, prepared

The screenshot shows the SAS Data Integration Studio 3.6.1 interface. The 'New Object Wizard' dialog is open, titled 'Select the type of object to create'. The 'Table' category is selected, and the 'Table' option is chosen from the list. The background shows the 'Project Explorer' on the left and the 'SAS Data Integration Studio' title bar at the top.

3. Data preparation for Data Warehouse

```
graph BT; keep_variable[keep variable] --> keep_var[keep variable-analitticki_S11]; keep_var --> SAS_Sort[SAS Sort]; SAS_Sort --> Analitticki_varjabli[Analitticki varjabli]; DT_DRU[DT_DRU] --> ST_S11[ST_S11]; DT_SMA[DT_SMA] --> ST_S11; DT_BUK[DT_BUK] --> ST_S11; Analitticki_varjabli --> ST_S11; ST_S11 --> ST_S11_YEAROB[ST_S11_YEAROB]
```

The flowchart illustrates the data processing workflow for ST_S11 data. It begins with a 'keep variable' step (blue square), which leads to a 'keep variable-analitticki_S11' step (green circle). This is followed by a 'SAS Sort' step (blue square), then an 'Analitticki varjabli' step (orange circle). The 'Analitticki varjabli' step feeds into the 'ST_S11' step (blue square). Additionally, three data sources—'DT_DRU' (orange circle), 'DT_SMA' (orange circle), and 'DT_BUK' (orange circle)—also feed into the 'ST_S11' step. Finally, the 'ST_S11' step leads to the 'ST_S11_YEAROB' step (orange circle).

These defined processes and procedures made on the data in the DW perform gaining of the desired outputs. They are defined through SAS process designer, visual tools that depict all defined and needed processes. Depending on the

4. Metadata of the Data Warehouse

[illegible]

With all these transformations, we calculate the new variables, called category for calculation of GDP. The final outputs are settled in many different summary tables, grouped depends on institutional sector. This data are disseminated as public data and placed in the SSO web site.

For updating data of the State central registry - the final account report of subjects in the DW, data is merging with dimensional tables of each subject of DW being made. The input data from

the State central registry, obtained in ACCESS tables and saved in the temporary folders, is cleansed, recoded and transformed. Then, the data is loaded with SAS in a DW Integration studio. Cyrillic data is recoded, merged with identified data in the subject's data of the State central registry and some additional variables are added in the dimensional tables. This additional data is prepared for the next phase of data preparation. These dimensional tables have the territorial classification, valid in the country, the institutional sector which possesses this data, National classification of activities (NCA) and other.

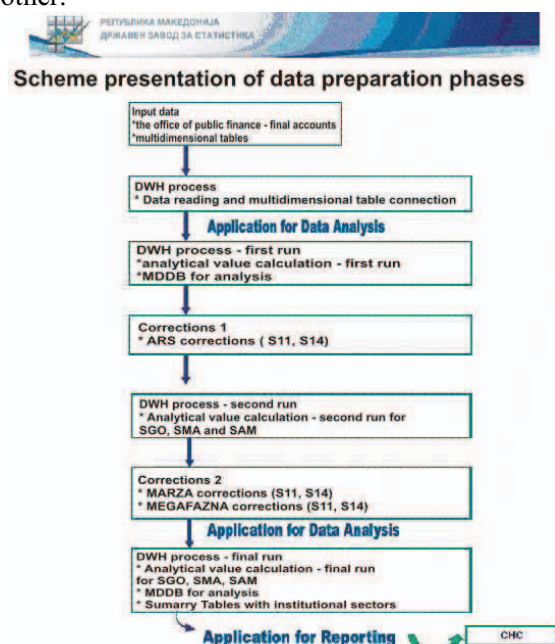


Figure 5: Schematic presentation of data preparations phases

A very important subject who gives additional information for business subjects and institutional sectors is the National accounts sector (NAS) has the responsibility for all category data calculation, keeping the formula up to date and data control. This data is defined as metadata for the transformation steps and is saved in "Kategorii" Excel file. This is the most important file for making all data calculations and must be updated and also, the oldest "Kategorii" file must be saved for historical analysis. This file is the base for the intermediate file creation and for the next calculation steps. This phase is called First run and the result is aggregated tables with all the data, but unlashd of grey economy (illegal market). The second run programs enrich these calculations with mythological solutions for

embedding some influence of the grey economy. The data from these two phases is compared and ultimately made some corrections. Then, the Final run phase can be started. With this calculation phase, the ESA codes for sector accounts are prepared. This phase is called also ESA run. The visual representation of data transformation and preparation for final reports are shown on Figure 5.

5. BI Tools for Statistical Data Analysis

A tool for data analysis and reporting used in SSO is SAS Adjustment analysis tool, Corrections and Reporting tools which are the parts of SAS EIS software tool. It is powerful and user-friendly software with many possibilities provides many-phrasal data and much table analysis (Figure 6).

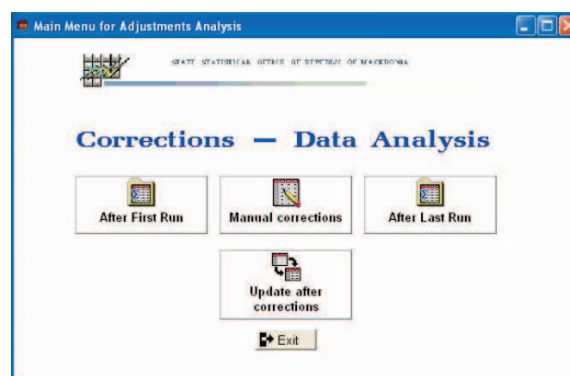


Figure 6 SAS Corrections – Data Analysis screenshot

After previous phases of data preparation, the analyst of each institutional sector in the National Accounts sectors makes analysis of the classification tables (Figure 7) with Adjustment analysis - Corrections tools. There are possibilities for analysis after the First run phase, making some manual corrections, after last run analysis and some other updates and corrections. These analyses are made with the selection of desired data from dimensional tables of DW for couple of classification codes. The experienced analysts from the NAS can explore data in each phase and can find some irregularities, if there is something like this. After this analysis, data must be saved in the final fact and dimensional tables and can be used for the next step – data dissemination. The final data can have some correction, but they must be strongly controlled from the NAS analysts.

A17-en	Amortizacija iz curenj		Amortizacija iz lasti		Bruto dodana iz curenj		Bruto dodana iz lasti		Bruto iz lasti		Bruto iz curenj		Mogufasna iz curenj		Mogufasna iz lasti		Sredstva iz curenj		Sum
	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	Sum	
A	114051.10	101522.44	402034.30	412034.30	117048.62	159719.62	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	100407.30	291303.04
B	4726.75	5495.25	47011.32	30388.83	100001.24	171431.89	84203.57	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	74722.81	20893.18
C	59270.79	43044.62	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	520047.67	143627.39
D	9621476.34	1080000.9	4795472.4	5637895.2	2041305.4	2184048.9	1703576.9	1474205.9	2504209.4	22	1703576.9	1474205.9	2504209.4	22	1703576.9	1474205.9	2504209.4	22	1703576.9
E	1593912.52	790702.03	1136886.3	860044.46	262033.3	438336.1	3171709.8	2016448.5	625709.08	58	3171709.8	2016448.5	625709.08	58	3171709.8	2016448.5	625709.08	58	3171709.8
F	1517380.6	102702.65	128619.3	895584.4	3420501.1	4261348.4	2970180.1	2489889.7	643450.26	95	2970180.1	2489889.7	643450.26	95	2970180.1	2489889.7	643450.26	95	2970180.1
G	440504.60	401246.79	3634701.4	3088081.1	5800856.7	7180256.9	3474516.5	2701857.6	1523459.3	13	3474516.5	2701857.6	1523459.3	13	3474516.5	2701857.6	1523459.3	13	3474516.5
H	835476.96	807011.15	420715.29	3023417.1	808886.12	1001086.2	5817516.6	477237.95	2347039.69	20	5817516.6	477237.95	2347039.69	20	5817516.6	477237.95	2347039.69	20	5817516.6
I	570157.87	860175.06	263845.3	2437709.3	6114007.1	6895319.5	435844.8	3672027.7	936379.79	79	6895319.5	435844.8	3672027.7	79	6895319.5	435844.8	3672027.7	79	6895319.5
J	1301145.95	110326.43	126637.7	1218433.9	182807.4	197208.4	710476.09	543071.50	546776.45	46	197208.4	710476.09	543071.50	46	197208.4	710476.09	543071.50	46	197208.4
K	188160.86	163871.07	148546.3	1014489.1	228823.35	364152.5	1545687.2	1188333.9	746773.63	55	364152.5	1545687.2	1188333.9	55	364152.5	1545687.2	1188333.9	55	364152.5
L	454142.32	401607.37	357732.7	227258.3	384708.7	467434.3	203708.6	1572134.8	2211342.18	18	467434.3	203708.6	1572134.8	18	467434.3	203708.6	1572134.8	18	467434.3
M	120531.69	183746.95	124305.7	113306.1	1551483.7	196509.5	602003.76	401779.57	1952010.4	97	196509.5	602003.76	401779.57	97	196509.5	602003.76	401779.57	97	196509.5
N	247083.73	244231.39	1348410.9	1219705.2	2004448.2	2263642.1	947521.27	784743.00	1051429.4	92	2263642.1	947521.27	784743.00	92	2263642.1	947521.27	784743.00	92	2263642.1
O	270568.57	102744.41	1088274.5	6827481.54	1195420.5	1718121.5	678889.94	538780.84	403546.83	28	1718121.5	678889.94	538780.84	28	1718121.5	678889.94	538780.84	28	1718121.5
A17-en	211052.87	142420.14	2106211.1	6127020.82	6249230.4	524242.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9	423230.9

Figure 7 – Tables compared after first and final run phases

Data dissemination is made with another Adjustment analysis – Reports tool (Figure 8 and 9). This tool provides many types of reports, depending of analyst's need and desired data. The analyst can provide analysis of integrated accounts, selected accounts or other Additional analysis depending of its demands. Reports can be enlarged if there are some demands from NAS.

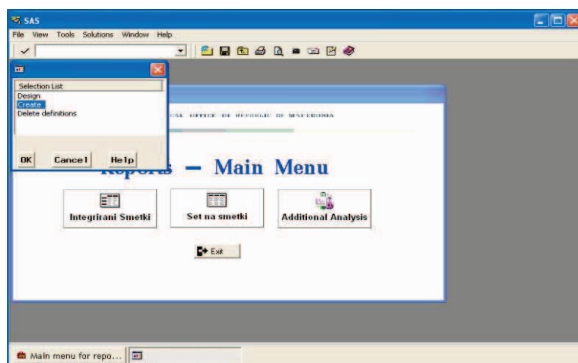


Figure 8 – Adjustment analysis - Reports tool

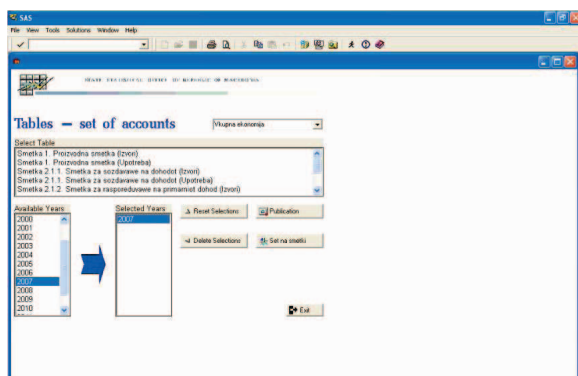


Figure 9 – Selection of parameters with Reports tool

This module can be enriched with SAS Graph tool, but it isn't implemented yet. In this moment, the data dissemination is made only with the table's reports. But, this is an excellent base for visual data analysis with other software tools

6. Tools for Visual Data Dissemination

However, the information dissemination isn't DW administrator task, but it is specific task which is dedicated to the SSP sector for data publishing. Some of the annual editions cover printed data dissemination. Although printed editions, SSO possess web site with embedded databases and possibility to browse desired data from this web site. The web site is made in the Microsoft ASP.NET platform and has a possibility to combine and present selected data from the DW. The user has a spectre of tables from which one can pick desired information and parameters for statistical analysis. The human-computer interface is easy to use and permits data visualization in the form of business graphics or maps. Some examples of use of these HCI are shown in Figures 10 (a, b, c and d). This example presents some statistic selections of data from the population in the Republic of Macedonia.

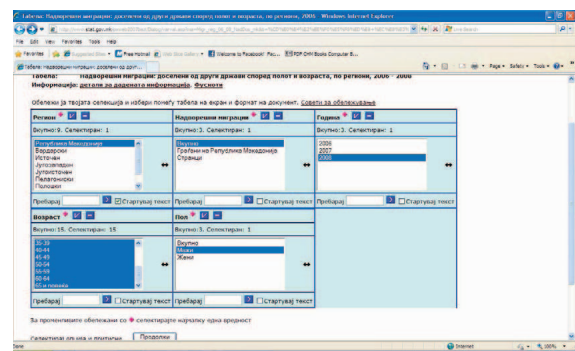


Figure 10 a – Web selection screen from SSA database

The web site is compatible with EU standards for statistical data representation. Data dissemination can be reached via web browser or PC-Axes for Windows. This easy HCI interface provides a wide range of data combination from the SSO Data Warehouse. The data can be viewed as tables or visualizations like business graphs or mapped data in the Macedonian maps. Also, they can be saved as Excel files in the user's desktops, prepared with other software tools or import in the company's databases. The first selection is

population, foreign trade, population State inventory, agriculture and regional statistics. Depending on the selected thematic area, the application selects data from multidimensional tables of DW, the selected category and data is collected in temporary tables and can be presented in the selected visual format, graphics or maps with graphical data presentation (Figure 11 a and b).

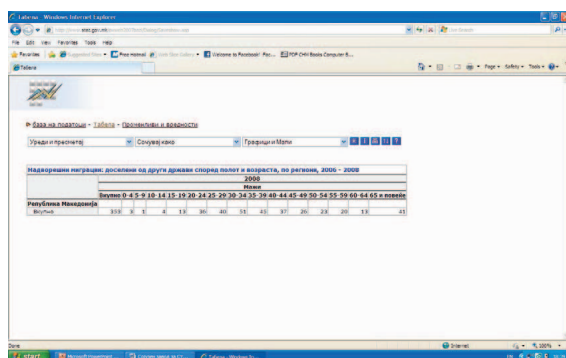


Figure 10 b – Selected data, presented in table format

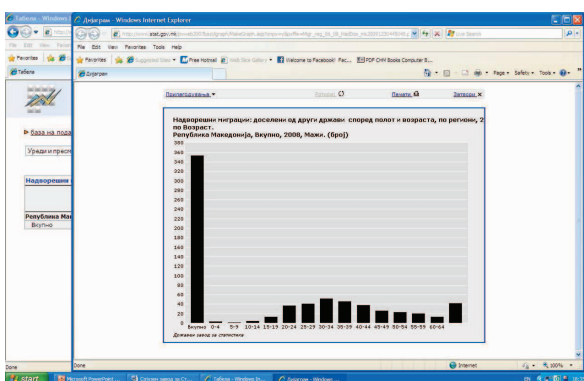


Figure 10 c – Data from 10 b presented in a visual format

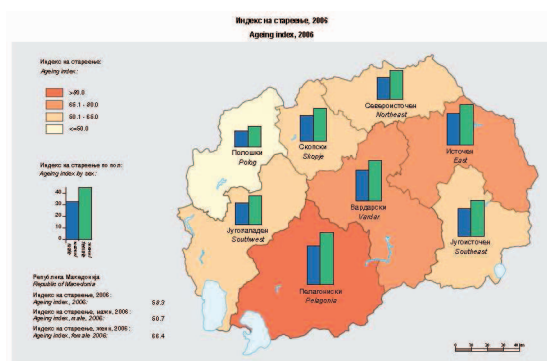


Figure 10 d – Combination of Map and BI graphics (coloured bars)

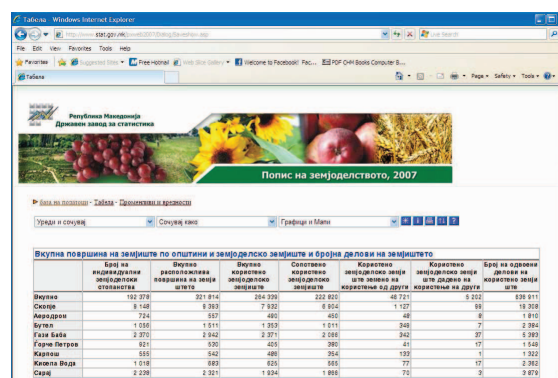


Figure 11 a - Agriculture Data analysis for selected data in a web table format

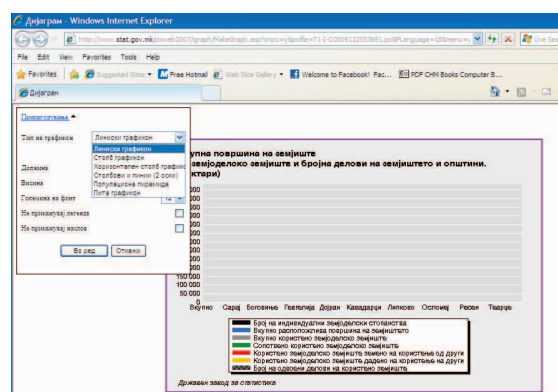


Figure 11 b Agriculture Data analysis and selection of desired visual format for data representation on the web

7. Conclusion

SAS DW is a base which is used for gaining very complex statistical calculations for the need of SSO and official statistical data, produced in the Republic of Macedonia. They use data from different sources, under different platforms and systems. The data must be very carefully cleansed, transformed and loaded in the target destinations in multidimensional tables in DW. After that, the data is the base of calculation of the sector's data in aggregated tables, depending by the statistical logic, nestled in the Excel file. These files are the fundamental databases for data dissemination for public usage via SSO web file and user friendly web software tools for data selection and data visualization. It is a very successful story for gaining BI information via statistical calculations from integrated state data. The data is disseminated via SSO web. It is a very user-friendly web site with a wide range of possibilities for statistical data representation.

7. References

- [1] Kozielski S., Wrembel R., New Trends in Data Warehousing and Data Analysis, Springer, 2009, [Pages 93-113];
- [2] Ponniah P, Data Warehousing Fundamentals, A Comprehensive guide for IT Professionals, A Wiley - Interscience Publication, 2001, [Pages 111-130];
- [3] SAS Institute Inc., SAS Data Integration server, USA, 2009; www.sas.com/offices, [2009/2010];
- [4] Cartier J., Painless Graphics: The Click, Drag and Drop Approach of Graph-N-Go, SAS Institute Inc, 2009
- [5] Karp A.N.H., Working with SAS Date and Time Functions, Sierra Information Services Inc, 2008;
- [6] <http://www.sas.com/rnd/webgraphs/>, [2008/2009]
- [7] <http://www.stat.gov.mk/>[2009/2010]