

# Web scrapping techniques for price statistics - the Romanian experience

July 1, 2019

## Abstract

Internet has been widely recognized as a new data source that can be used to compile new statistics or the enhance the traditional ones in several fields of official statistics. Considering that online commerce has a growing share in the overall household's consumption, price statistics is one of the areas of official statistics that can have important benefits from this new data source. There have been several projects around European countries exploring the potential of Web scrapping techniques to enhance the production of the official statisticals, including the classical consumer price index (CPI). In this paper we will describe the experience of NSI Romania regarding the collection of prices from Internet and using these data to compile an experimental consumer price index. The aim of our pilot project was mainly to investigate whether the web-scrapping method of data collection for prices can be introduced in the production of official statistics in the near future and which are the methodological challenges that we have to deal with. We developed a chain of software tools that automates the whole process, starting with data collection, transforming the semi-structured data into structured data, going to a data validation procedure and finally to a computation procedure that outputs a price index. We started from the traditional methodology used for CPI and added some new features such as a clustering technique and a distance-based method for matching similar products to take advantage of the specificity of the web-scrapping collection method. The whole process was represented in terms of GSBPM.

## 1 Introduction

Whether we use Internet for doing business, social networking, shopping or education, the quantity of data that we produce in our daily activities has recorded and exponential growth. Together with data produced by machines and sensors, these new and almost real-time large volumes of data that are generated today are commonly called Big Data. The first signs that Big Data sources can generate value and useful insights were given by private companies during the late 90s' but nowadays the European and national statistical systems

has also witnessed major transformations because of the challenges raised by the big data sources.

The incorporation of Big Data sources in the official statistical production does not aim to entirely replace the traditional methodologies but it is rather an iterative and incremental approach in which certain components of the traditional statistical production process are augmented by the Big Data sources inputs and the related processing algorithms [1], [2]. Alternatively, big data sources can contribute to the reduction of the response burden or they can be used only to study some economic or social phenomena before designing a statistical survey which may be expensive.

Speaking in other words, the incorporation of big data sources into the official statistics means maintaining a net competitive advantage and relevance of the official statistics products compared to those provided by a plethora of commercial players, with reference to large corporations that are active in the field of information technology [3].

One of the main big data sources is the Web system that can be considered an immense reservoir of information and this source cannot be neglected by official statistics institutes. In order to take advantage of the data publicly available on Web sites some automatic procedure for data collection should be designed first. These procedures are referred to under the term of web-scraping.

Automatic data collection and its use to derive statistical indicators was pioneered by MIT [19] where the prices collected from online shops were used to build a consumer price index for some South-American countries. Since this first experiment, several statistical offices throughout the world started to collect data from online retailers and study how these data sets can be used for consumer price index calculation. We can mention here Statistics Netherlands [20], ISTAT [22] or Destatis [23] as some of the first statistical offices in Europe that experimented the web-scraping technique for online prices, although they didn't follow the classical big data approach of MIT and only monitored some prices or tried to collect prices only for the products included in the traditional collection method. The web-scraping technique was used to collect data in other areas of statistics too, for example to improve some statistical registers [24] or for job vacancies [25]. No matter how it was used, for bulk scraping of all prices, or for only specific prices in certain areas [21] the web-scraping technique proved to be a very useful method in the hand of statisticians.

Under these auspices, the overall objectives of our experimental project were to streamline the statistical production process by lowering the overall production costs, reduce the response burden and the dissemination term. Such projects, through the incorporation of modern computing technologies, could create the premises for developing a framework for testing and piloting new methodologies and technologies in a systematic and rigorous manner [4].

Our project experimented how web-scraping collection method can be used to compute a new/experimental consumer price index (CPI) or to improve the classical CPI computation [5]. We started our work by identifying and selecting online channels that have significant weights in the process of trading goods and services for household consumption. This is not an easy task given that there is

no information on the volume of online transactions made by firms, issue found in other projects too [6]. The eloquent example is given by retailers in the hypermarket category, which although they have a physical trading correspondent with very high trading volumes, the volume of online transactions is unknown. The criteria used to select the online trading channels included in our study was to have a physical correspondent and record significant sales volumes at national level. Next, we proceeded with the task of identifying the appropriate means to implement the automated price collection process from e-commerce sites. The criteria used to identify the optimal solutions are expressed in terms of flexibility, ease of use, scalability and cost. An essential task to achieve this goal was to explore other approaches and test the existing solutions.

Another objective of our project was to carry out the automatic price collection process over a relevant period: 6 months - 2 years. Achieving a maturity level specific to official price statistics that are currently published will require a much wider period of rigorous and systematic testing of the collection process and the results obtained. The resources available for running the data collection, technology and skills are critical and a continuity plan should be devised if some data sources become unavailable, legislative changes occur during this period, or the technology and skills are outdated by the evolution of the Web architecture.

The next objective of our project was to compute an elementary price index at article/varietal and assortment level and compare it with those obtained using the traditional data collection method in order to emphasize the issues related to the difficulties of applying and/or adapting the traditional consumer price index (CPI) methodology [12] to the new data sources. A compromise to ensure a certain degree of comparability is the use of traditional CPI methodology [13], [14] to estimate price indices, although traditional methodology may be incompatible from some points of view with the new data source.

Last but not least, we intended to identify the legally sensitive aspects regarding the reconciliation between National Statistical Law, the European Statistics Code of Practice, other regulations on official statistics and legislation on access to online data [15].

The paper is structured as follows. In section 2 we present details of the data collection process, in section 3 we provide a description of the methodological approach, in section 4 we present our first results and section 5 concludes our paper

## 2 Data collection

Some of the official statistics bureaus that have run similar projects have opted to outsource this component to companies specialized in collecting, processing and storing the data instead of acquiring the data directly. We explored several existing software solutions: Robot framework [7], Scrapy [8], [9] Apache Nutch [10], RSelenium [17], [18] and [11]. Based on an analysis of the advantages and disadvantages of each solution we chose to work with the Robot framework.

The observation unit was the web site of the retail companies. In this case, the assumption from which we started was that the companies cover the entire national territory through their site. Sites selection was based on establishing a sales-turnover relationship, sorting by decreasing order the sales figures reported by the firms that own the sites. At the present moment, there are certain barriers, for example the most important player in terms of turnover on the hypermarket segment in Romania, does not have a section dedicated to online transactions. We selected 4 sites for food, 5 sites for clothing and 5 sites for footwear products. However, moves made at European level by firms that have physical stores on this segment suggest that market forces will require online migration of the most important players in the field.

The main variable collected from these sites was the price with VAT. The automatic collection method allows us to also record the prices for the goods and services affected by discounts, promotions, or other forms of attracting customers through prices, so we can record the old price, and the discount shown as a percentage alongside with the displayed price. This facilitates, for example, the easy identification of seasonality factors that affect the price variation for certain categories of goods and services. Prices are recorded in .csv files that contain the following variables: the article/varietal name (the name under which the article is marketed), current retail price, old price and/or retail discount if displayed, composition for clothes/footwear, a short description (manufacturer and technical specifications), the date collection and the website address. The selection of the products whose prices are kept under observation is based on the CPI national standard classification. We collected about 50,000 to 70,000 records every month.

Data collection took place through the Robot Framework and RSelenium software solution. It is worth mentioning that Robot Framework has a high degree of configurability through the possibility of introducing specific procedures regarding the technology behind the sites. This software solution proved to be a scalable web-scraping tool that can fulfill the requirements of a large organization. The automatic collection of prices observed on the sites included in the sample was made during the same period as for the traditional CPI survey. Due to the complexity of the data extracted through the web-scraping process, i.e. of the semi-structured data gathered from the sites, the decomposition at the core components of the CPI classification is required first.

The structure of the data collected from the retailers' sites for the food group of products contains the product name, the manufacturer, the quantity, certain technical-quality details, the price per unit or the price per piece, the article/varietal and assortment type, and the category according to the structure of the site. From the point of view of the classification of products in a given product category, these data may appear at a first glance as inputs for a manual or automatic classification procedure, but the very large number of products and the fact that the description is not standardized for all sites targeted by the collection process makes this stage to be considered as the most difficult one.

A trivial observation about the form of data is that they cannot be directly used in the process of classifying and estimating price indices. To address this

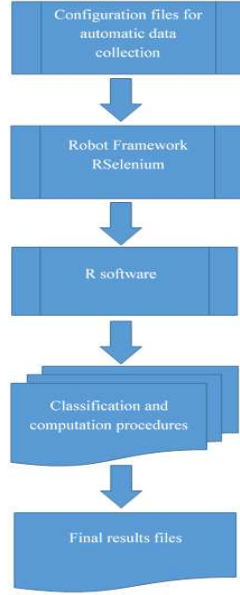


Figure 1: Data collection and processing session

issue, we have developed a series of R scripts that transform the data in a way that allows flexible handling. The CPI computation steps are sequentially deployed, the data input for each stage depending on the output of the previous stage, except for the first step whose input depends on the result of the automatic data collection.

In the following, the activities carried out at each stage will be detailed, noting that we attempted to keep the traditional CPI methodology as much as possible intact. A graphical representation of the data collection and processing is shown in 1.

The first activity was the data cleaning. We started with the web scrapped files and performed some basic operations checking for missing data and other basic validation operations. In case there are missing items among the data sets, the web-scraping process resumes, after checking the online accessibility of the site and the log files of the web-scraping application. Some possible error sources could be:

- sites were unavailable or have undergone changes;
- the web-scraping application encountered web content elements that cannot be directly processed;
- web server identified the web-scraping application as a malicious software and imposed an access restriction to the site at the IP address level.

Next, all the files obtained from the data collection process for a certain month are joined automatically. The resulting file is read by an R script and transformed into a data structure suitable for an automatic processing procedure. Some basic transformations are again performed using an automated R script before classifying and linking the products according to the CPI classification. We started with the manual product linking and classification according to the standard CPI classification which implies identifying the observations which contain a description similar to the one provided in the classical CPI classification. This activity can generate errors whose propagation can significantly influence the quality of the results. The principle that we used in the absence of a previous experience in working with methodological aspects of selection of the articles was to assume that the consumer will choose a product or products substitutable to the one present in the standard classification within a reasonable price limit ( $\leq 150\%$  of the price of an article from the standard classification).

Thus, we chose to select several articles for one assortment within the same observation point. To reinforce the strict tracking rule of the same articles found in the standard CPI methodology, we performed join operations between the data structures for all decades and observed months. The join operation between two or more tables was based on the "name" variable containing the product description by matching strings in a 1 to 1 ratio. After performing this activity, from an initial number of about 10,000 of articles, they were restricted to 545 articles, 216 assortments, and 52 expenditure groups, identified as constant during the 6 months of observation, assuming that the description given in the observations made for the variable "name" represents a guarantor for the invariance of the technical and qualitative characteristics of the articles. This technique was used to encode the entire sample.

Several attempts were made to develop an automatic encoding procedure with encouraging results. However, their use would involve deviations from the established methodological standard, manifested by the appearance and disappearance of the articles in the sample with a high frequency. We tried several machine learning and distance-based algorithms for this procedure and the results are presented in table 1. The best results, as it can be observed in 1, were obtained using the Levenstein distance.

| Algorithm                           | Accuracy |
|-------------------------------------|----------|
| Boosting                            | 0.56     |
| Support Vector Machines             | 0.34     |
| Random Forests                      | 0.41     |
| Scaled linear discriminant analysis | 0.17     |
| Bagging                             | 0.28     |
| Regex                               | 0.70     |
| Levenstein Distance                 | 0.80     |

Table 1: A summary of the methods used for automatic classification

### 3 Some methodological aspects

The scope of the project was to assess if online observed prices can be successfully used as a substitute data set for computing, either the traditional CPI or a similar experimental statistics, e.g. online observed CPI. Therefore, in order to retain the results, as much as possible, comparable with the traditional CPI, the collection periods within a month, along with the goods and services included in the CPI national classification were preserved. Due to practical limitations regarding the allocated resources for this project, the data collection process was focused on food and beverages and items covering clothing and footwear categories, as these types of goods hold the biggest share in household's consumption expenses, e.g. food accounts for nearly 40% of total expenses[26].

CPI is computed by aggregating at different stages indices a product/variety, assortment and group/category level, the entire process being graphically described by 3. After data pre-processing, prices are aggregated into an arithmetic monthly average for each item, given that data was collected 3 times per month for food and beverages and once per month for clothes and shoes.

$$\bar{p}_v = \frac{\sum p_{v_n}}{n}, \quad (1)$$

,

$\bar{p}_v$  = monthly price average for any given variety,  
 $n$  = number of times price data was collected in a given month.

The average is used to compute elementary price indices at product/variety level, by dividing the current monthly average for an item to its respective base period monthly average.

$$i_{p_v} = \frac{\bar{p}_{v_{current}}}{\bar{p}_{v_{base}}}, \quad (2)$$

$i_{p_v}$  = elementary prices index for any given variety,  
 $\bar{p}_{v_{current}}$  = current monthly price average,  
 $\bar{p}_{v_{base}}$  = base monthly price average.

To ensure that results are comparable and capture only pure price change, ideally would be to collect price data for the same products indefinitely[13]. In reality, this is impractical due to different reasons. Therefore, price data collectors are equipped with a list of strict rules when products or services are no longer available and substitutes are needed. These rules may target product description(producer, weight, composition, etc.) to ensure that qualitative differences between products no longer available on the market and substitute new products are minimal, store placement and local or national market share, or a combination of these[12]. According to different studies conducted at National Statistical Offices [].The collection of price data from online sources restrict the number of articles within the same observation point and the same assortment for the exact application of the classical CPI methodology. We used a clustering procedure, meaning that we computed a geometric mean to aggregate the

results in the form of a generic article specific to that observation point. For example, if one observation point encounters 3 articles within the same assortment, in another observation point 2 articles within the assortment encountered at the previous point, we apply a geometric mean for the prices of the three articles for the first observation point, the result being a generic article for the first observation point and the same procedure for the second observation point according to the formula:

$$i_{vg} = \sqrt[n]{\prod_{i=1}^n i_{v_n}} \quad (3)$$

where  $i_{vg}$  is the elementary index of the generic article at observation unit level,  $n$  is the number of articles within the same assortment and  $i_{v_n}$  is the elementary price index at the article level. The subsequent calculation steps follow roughly the steps from the classical CPI methodology, noting that in order to obtain the price index at the expenditure group level we assigned to each assortment a weight equal to  $\frac{1}{n}$ , where  $n$  is the number of assortments identified as belonging to that group and we applied a weighting recalibration procedure to aggregate price indices at expenditure group level using the following formula:

$$coef_{rp} = \frac{\sum_1^{pn} pg_{pn}}{\sum_1^p pg_p} \quad (4)$$

where  $coef_{rp}$  is the recalibration coefficient,  $pg$  is the weight of the expenditure group in the total of initial weights,  $pn$  the number of expenditure groups identified after the classification and data cleaning and  $p$  the initial number of expenditure groups from the classical CPI methodology.

The recalibration coefficient was applied to each weight from the classical CPI expenditure groups according to the formula:

$$pr = pg \times coef_{rp} \quad (5)$$

This intermediate stage is necessary due to the absence of certain articles from the offer present on retailers' web sites, subsequently transferred to assortments and/or expenditure groups. The process diagram is shown in 2 while in 3 we build a process diagram in terms of GSBPM.

## 4 Results and discussion

Using August 2017 as the basis for computation of the monthly price index, we obtained the aggregated indices at the groups of food, clothing and footwear presented in figures 4, 5, and 6.

From the evolution of the two price indices considered, it can be noticed that the online collection method implies a different trajectory due to the different samples used and the use of equal weights at assortment and expenditure group level. Another possible explanation can be found in the non-probabilistic



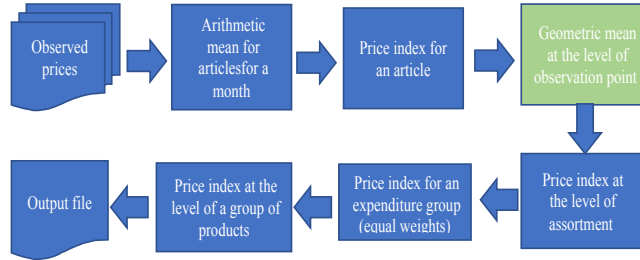


Figure 2: The process diagram for computing online price index - the green box adds a new phase to the traditional price index methodology

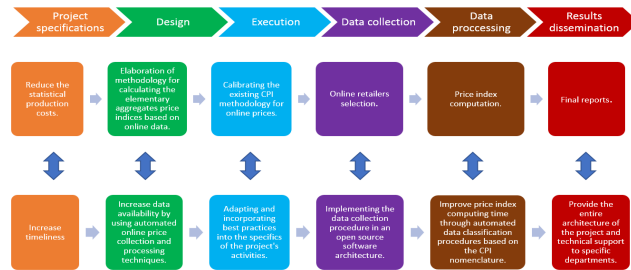


Figure 3: The GSBPM diagram



Figure 4: The comparative evolution of the price indices for food



Figure 5: The comparative evolution of the price indices for clothes



Figure 6: The comparative evolution of the price indices for footwear

sampling process through which online stores are selected ignoring the representativeness at national level due to the lack of specific information. Selected food, clothing and footwear stores can serve large cities and neighboring areas, having complex pricing policies which are different from small shops serving small city areas and rural communities.

## 5 Conclusions and future developments

This project was the first experiment that implemented a web-scraping technique for data collection in our NSI. While we gained experience with the software tools involved in such a project we also identified some limitations for our specific study of online price collection which are briefly described below:

- Generalization hypothesis of online transactions. The number of households purchasing an online product is relatively small, and generally depends on several factors such as the geographical position, income level, education level, etc.
- Not all businesses with a significant volume of transactions included in the list of observation units for traditional consumer price index has a website;
- The IT technology can have a significant impact on price variation. An example of this may be the discrimination based on the geographic position of a user when displaying prices on a particular site;
- The components of the classical consumer basket and the weights used at the level of the expenditure groups do not entirely reflect the consumption habits and the budget restrictions of the segment of the population addressed by the online stores.

Based on the results obtained and the potential of the web-scraping collection method we intend to implement it to other official statistics areas and we will continue to develop a specific online price index [13], by extending the current collection procedures to the entire products and services nomenclature and by developing a new methodology based on online prices. Secondly, a separate product and service nomenclature may be developed specifically for online observations based on measurements such as the longevity of certain products and services in the online offer, and a series of metadata related to those products and services, for example, analysis of online interaction based on reviews of buyers with the respective brands and the online store.

## References

- [1] Robert Griffioen, R., Olav ten Bosch, and Els Hoogteijling, Challenges and solutions to the use of internet data in the Dutch CPI, Workshop on Statistical Data Collection, The Hague, The Netherlands, 3-5 October 2016.

- [2] Griffioen, Robert and Ten Bosch, Olav On the use of Internet data for the Dutch CPI, Conference of European Statisticians, Geneva, 2-4 May, 2016.
- [3] European Commission, Internet as a data source, Luxembourg, Publications Office of the European Union, 2012.
- [4] Himanshi Bhardwaj, Tanya Flower, Philip Lee and Matthew Mayhew, Research indices using web scraped price data: August 2017 update, ONS, UK, 2017.
- [5] Josef Auer and Ingolf Boettcher, From price collection to price data analytics. How new large data sources require price statisticians to re-think their index compilation procedures. Experiences from web-scraped and scanner data, Ottawa Group - International Working Group on Price Indices, 2017, url=<http://www.ottawagroup.org/>.
- [6] Leon Willenborg, Elementary price indices for internet data, Discussion Paper no. 8, CBS, The Hague, Netherlands, 2017.
- [7] CBS, Robot Framework, 2018, url=<http://research.cbs.nl/Projects/RobotFramework/index.html>.
- [8] Dimitrios Kouzis-Loukas Learning Scrapy, Packt, 2016.
- [9] Daniel Myers and James W. McGuffee, Choosing Scrapy, Journal of Computing Sciences in Colleges, Volume 31, Issue 1, October 2015, Pages 83-89.
- [10] The Apache Software Foundation, Nutch, A highly extensible, highly scalable Web crawler, 2018, url=<http://nutch.apache.org/>.
- [11] Hadley Wickham, rvest: Easily Harvest (Scrape) Web Pages, 2016, R package version 0.3.2, url = <https://CRAN.R-project.org/package=rvest>.
- [12] INS Romania, Ancheta statistică a prețurilor de consum al populației (in Romanian), 2017.
- [13] ILO/IMF/OECD/UNECE/Eurostat/The World Bank, Consumer price index manual: Theory and practice, Geneva, International Labour Office, 2004.
- [14] UNECE/ILO/IMF/OECD/EUROSTAT/The World Bank/ONS, Practical Guide to Producing Consumer Price Indices, United Nations, New York and Geneva, 2009
- [15] Nigel Swier, How should web scraping be organised for official statistics? 61st ISI World Statistics Congress, Marrakech, 2017.
- [16] Leon Willenborg, Transitivity of elementary price indices for internet data using the cycle method, CBS Discussion Paper, September, 2017,
- [17] John Harrison, RSelenium: R Bindings for 'Selenium Web-Driver'. R package version 1.7.5, 2019. url=<https://CRAN.R-project.org/package=RSelenium>

- [18] John Harrison, RSelenium: R Bindings for 'Selenium WebDriver', 2019  
url=<https://CRAN.R-project.org/package=RSelenium>
- [19] Cavallo, Alberto, Scraped Data and Sticky Prices, MIT Sloan, December 28, 2010.
- [20] Hoekstra, R., ten Bosch, O. and Harteveld, F., Automated data collection from web sources for official statistics: First experiences. Statistical Journal of the IAOS: Journal of the International Association for Official Statistics 28 (3-4). pp. 99-111 2012.
- [21] Olav ten Bosch, Dick Windmeijer, Arnout van Delden and Guido van den Heuvel, Web scrapping meets survey design: combining forces, BigSurv Conference, October 25 - 27, Barcelona, Spain, 2018.
- [22] Polidoro F., Giannini R., Lo Conte R., Mosca S., Rossetti F. Web scraping techniques to collect data on consumer electronics and airfares for Italian HICP compilation, Statistical Journal of the IAOS 31 pp. 165–176 2015.
- [23] Brunner, K., Automated price collection via the internet, DESTATIS, 2014.
- [24] Barcaroli G., Scannapieco M., Scarno M., Summa D. Using Internet as a Data Source for Official Statistics: a Comparative Analysis of Web Scraping Technologies, New Techniques and Technologies for Statistics, Bruxelles, 2015.
- [25] Swier, Nigel, Webscraping for Job Vacancy Statistics, Eurostat Conference on Social Statistics: Towards more agile social statistics, Luxembourg, 2016
- [26] INS Romania, Coordonate ale nivelului de trai în România (in Romanian), 2019, page 52.