

# Using Web scrapping techniques for price statistics - the Romanian experience

Dr. Bogdan OANCEA

*Director, Department of Innovative Tools in Statistics, NSI Romania*

**Abstract.** Internet has been widely recognized as a new data source that can be used to compile new statistics or to 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 classical consumer price index. In this paper, we will describe the experience of NSI Romania regarding the collection of prices from Internet and compiling a consumer price index. The aim of our pilot project was 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 what are the methodological challenges that we have to deal with. We developed a chain of 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 but we 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

The European and the national statistical system has witnessed major transformations because of the challenges raised by the new massive real-time data generation, commonly called the Big Data revolution, whether we refer to the data generated by individuals, processes or machines. The incorporation of Big Data sources in the 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 [5]. 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 the information technology [4].

Under these auspices, the overall objectives of our experimental project are to streamline the statistical production process by lowering the overall production costs, reducing 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 [2].

Our project experimented how new techniques of data collection such as web scrapping can be used to compute a new/experimental consumer price index or to improve the classical CPI computation [1]. We started 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 [14]. 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. 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 [3], Scrapy [8], Apache Nutch [10], and rvest [12]. Based on an analysis of the advantages and disadvantages of each solution we chose to work with the Robot framework.

Another objective of our project was to carry out the automatic price collection process over a relevant period: 6 months - 1 year. 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 the elementary price indices at article/variety and assortment level and compare them with those obtained with 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 [7] to the new data sources. A compromise to ensure a certain degree of comparability is the use of traditional CPI methodology [6, 11] to estimate price indices, although traditional methodology may be incompatible from some points of view with the new data source. Finally yet importantly, 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 [9].

## **2. Some methodological aspects**

To keep the results as much as possible comparable with the traditional CPI, the observation periods within a month, along with the goods and services included in the consumer price index national classification are retained, starting the data collection process with the group of food and the items covering clothing and footwear, these goods having the biggest share in households' consumption.

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/variety 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 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 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. The first activity was the data cleaning. We started with the web-scraped 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, the web server identified the web-scraping application as a malicious software and imposed an access restriction to the site at the IP address level, etc.

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 that 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 table 1, were obtained using the Levenstein distance.

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

ALGORITM	BOOSTING	SVM	RF	SLDA	BAGGING	REGEX	LEVENSTEIN DISTANCE
ACCURACY	0.56	0.34	0.41	0.17	0.28	0.70	0.80

The next phase of calculating the elementary price indices at the level of article, assortment, and expenditure group is the last stage at which the final results were obtained. In the first part we computed the arithmetic averages for the articles for each month and observation point. This average was used to compute the elementary price indices at the article level. In order to calculate the indices at the assortment level we have decided to 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:

$$i_{vg} = \sqrt[n]{\prod_1^n i_{v_n}} \quad (1)$$

where  $i_{vg}$  is the elementary index of the generic article at the observation point 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  $1/n$ , 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:

$$coef_{rp} = \frac{\sum_1^{pn} pg_{pn}}{\sum_1^p pg_p} \quad (2)$$

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:

$$pr = pg * coef_{rp} \quad (3)$$

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 figure 1 while in figure 2 we build a process diagram in terms of GSBPM.

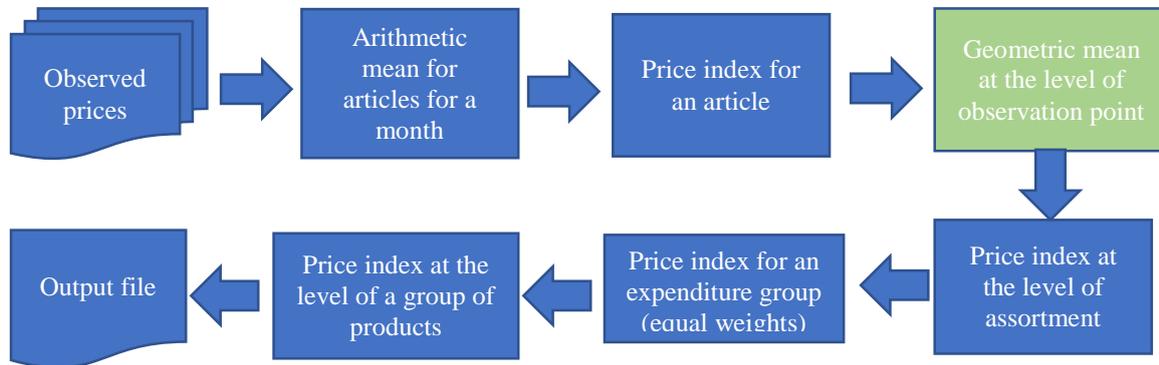


Figure 1. The process diagram for computing online price index (the green box adds a new phase to the traditional price index methodology)

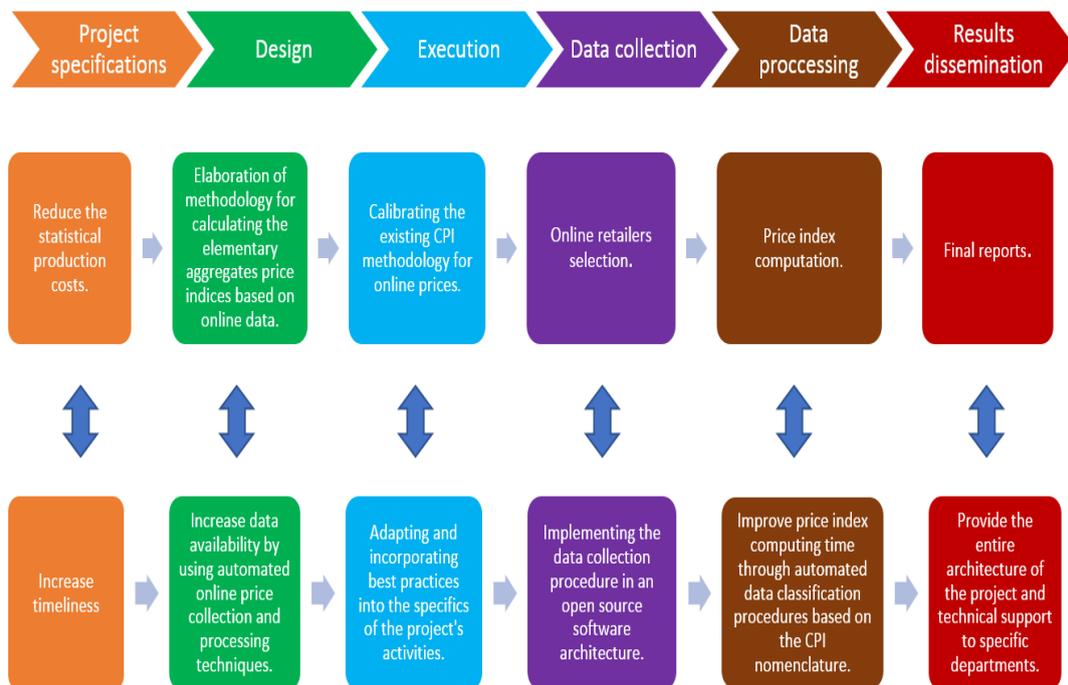


Figure 2. The GSBPM diagram

### 3. 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 3, 4, and 5.

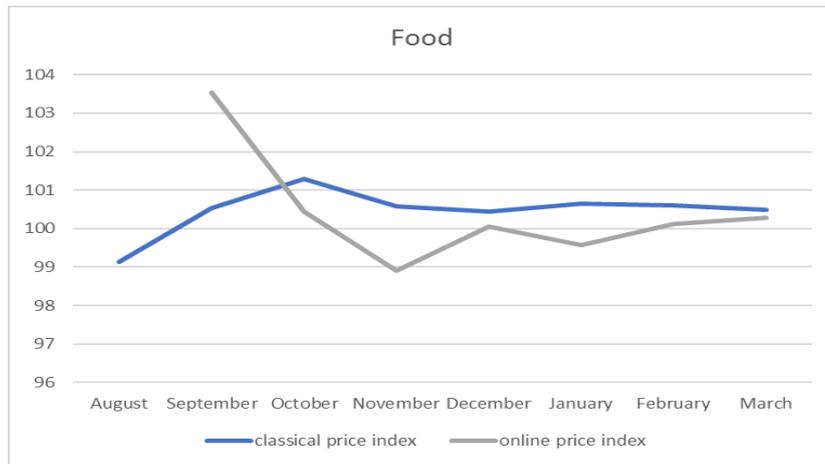


Figure 3. The comparative evolution of the price indices for food



Figure 4. The comparative evolution of the price indices for clothes



Figure 5. The comparative evolution of the price indices for footwear

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 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 that are different from small shops serving small city areas and rural communities.

#### **4. 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 that 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.

## 5. References

- [1] Auer, J., and Boettcher, I., (2017), From price collection to price data analytics, Ottawa Group - International Working Group on Price Indices, <http://www.ottawagroup.org/>.
- [2] Bhardwaj, H., Flower, T., Lee, P., and Mayhew, M., (2017), Research indices using web scraped price data: August 2017 update, ONS, UK.
- [3] CBS, (2018), Robot Framework, <http://research.cbs.nl/Projects/RobotFramework/index.html>
- [4] European Commission, (2012), Internet as data source, Luxembourg, Publications Office of the European Union.
- [5] Griffioen, R., ten Bosch, O., and Hoogteijling, E., (2016), Challenges and solutions to the use of internet data in the Dutch CPI, Conference of European Statisticians, Workshop on Statistical Data Collection, „Visions on Future Surveying”, The Hague, Netherlands.
- [6] ILO/IMF/OECD/UNECE/Eurostat/The World Bank, (2004), Consumer price index manual: Theory and practice, Geneva, International Labour Office.
- [7] INS, (2017), Ancheta statistică a prețurilor de consum al populației (IPC), INS Romania.
- [8] Scrapy developers, (2018), Scrapy - An open source and collaborative framework for extracting the data you need from websites in a fast, simple, yet extensible way, <https://scrapy.org/>.
- [9] Swier, N., (2017), How should web scraping be organised for official statistics? 61st ISI World Statistics Congress, Marrakech.
- [10] The Apache Software Foundation (2018), Nutch, A highly extensible, highly scalable Web crawler, <http://nutch.apache.org/>.
- [11] UNECE/ILO/IMF/OECD/EUROSTAT/The World Bank/ONS, (2012), Practical Guide to Producing Consumer Price Indices, United Nations, New York and Geneva, 2009.
- [12] Wickham, H., (2018), rvest - Easily Harvest (Scrape) Web Pages, <https://github.com/hadley/rvest>.
- [13] Willenborg, L., (2017), Transitivity of elementary price indices for internet data using the cycle method, Technical Report, CBS, DOI: 10.13140/RG.2.2.14338.27846.
- [14] Willenborg, L., (2017), Elementary price indices for internet data, Discussion Paper no 8., CBS, The Hague, Netherlands.