17th Ottawa group meeting Rome, 7 – 10 June 2022

The use of different web scraping techniques and API to estimate Italian price inflation (draft)

Antonietta D'Amore (Istat - Italy, antonietta.damore@istat.it) Rolando Duma (Istat - Italy, rolando.duma@istat.it) Cristina Dormi (Istat - Italy, cristina.dormi@istat.it) Rosabel Ricci (Istat - Italy, rosabel.ricci@istat.it) Antonella Simone (Istat - Italy, antonella.simone@istat.it)

Abstract

The Italian National Institute of Statistics (Istat) for the centralized data collection of consumer prices currently adopts a system that gathers in high percentage price information on the Internet for several goods and services. After the introduction of web scraping techniques in 2014 for consumer electronics (currently no longer used) and the preliminary tests related to airfares, in 2018 there was the very first case of application of the use of web scraping to collect the prices of transport by train. One year later the same technique has been used for Electricity on the free market, and in 2019 the web scraping collection has been extended to other two products: town gas and food delivery. The improvement obtained in terms of efficiency of the process and in terms of quality of the data has driven the choices of the last period to answer to the growing requests of statistical information in the field of consumer prices.

The paper is aimed at showing the use of different techniques to scrape online data in order to improve the indicators produced by Istat to measure inflation. In particular, the paper focuses the attention on one side on the use of a screen collection (imitating the behavior of consumers visiting websites and buying products or services) and, on the other side, on the data capture of information available in the API (Application Programming Interfaces) that allows direct access to structured data.

The software used for screen collection is a browser based application for macro recording, editing and playback for web automation and testing. It is provided as a standalone application and as an extension for the most common web browser that adds the recording and playback functions of repetitive activities, such as sending and receiving files and e-mails, refreshing, saving and printing web pages, filling out the forms and remember passwords.

The screen collection is used to detect prices of electricity on the free market (since January 2019). This approach is extended in 2022 to cover the prices of gas on the free market that will complement the rates applied in the protected market, to give the right picture of the energy market in our country, thanks also to many changes in the regulatory environment that have favored liberalization. The free gas market was introduced in the 2022 basket for the Italian consumer price indices. The paper illustrates the use of the automatic data collection procedures (web scraping), the analysis of the data acquisition status, the selection and stratification phases of the sample and the trends of the indices with a comparison between the free market and the protected one.

The collection of data through the API's is the other main topic of the paper in the framework of the research on new survey techniques that Istat is testing. In this context, some important results of the experimentations carried out for the detection of air transport prices will be analyzed.

Then pros and cons of the two web-scraping techniques are discussed.

Finally, it is analyzed the use of big data to support the current consumer price survey as an alternative to the traditional techniques.

Index

1 Introduction					
2	The collection on the gas and electricity free market for the estimation of inflation				
	2.1	Description of the reference market	. 5		
	2.2	Online Portal of offers	. 5		
	2.3	The sampling definition and data collection	. 6		
	2.4	The trend of Electricity and Town gas and natural gas prices	. 7		
3	API'	s, the new price collection project	10		
	3.1	IT aspects and data flow	10		
	3.2	Some results	11		
4	14 Concluding remarks				

1 Introduction

The data that contribute to the compilation of the monthly indices of consumer prices are collected through the use of a plurality of sources and techniques: territorial survey, carried out by the Municipal Statistics Offices (UCS); centralized survey, conducted directly by Istat or through collaboration with large data providers; data scanners; administrative data sources.

In 2022, the centralized survey accounts for 22.3% of the basket in terms of weights and currently adopts a system that gathers in high percentage price information collected on the Internet for several goods and services.

Alternative data sources, namely scanner and web-scraped data, and methods to utilize these data sources, are being introduced systematically into the production of IT consumer price statistics since 2018.

The collection conducted manually at the very beginning has taken advantage of the use of alternative data acquisition like the web-scraping that is made of a set of techniques aimed at finding information from the web automatically through intelligent software that interpret, select and collect the information content of a web-site. Data are then structured, stored and analyzed in a local database in order to be used for the indices calculation.

The gain in efficiency for the process and in quality for the data has driven the choices of the last period to answer to the growing requests of statistical information.

In 2018 after few years of experiments, there was the very first case of application of the use of web scraping to collect the prices of transport by train. Since then, twice a month the first and second week with a fixed calendar over the year, it is carried out a simulation of purchase of a train ticket on the web sites of the two main Italian companies of railway transportation. Prices refer to several types of services and to a total of 167 travel routes selected on the base of actual consumer purchasing behavior.

In 2019 the same technique has been used for Electricity on the free market. Once a month, within the first two weeks of the month, prices are scraped on the Regulatory Authority's public website <u>www.ilportaleofferte.it</u>, referred to tariffs applied to the domestic customers of the free market divided by region, by type of contract (resident or non-resident; variable or blocked price) and by type of time slot (single-hour or multi-hourly).

In 2020 the web scraping collection has been extended to other two products: town gas and food delivery.

For town gas and natural gas, on the regulated market, once a month, by the first two weeks of the month prices are scraped. Data are used in combination with the traditional price collection realized by the municipalities participating in the survey, to have the full coverage of the national territory.

To collect data on food delivery, once a month the website of the delivery company more widespread on the national territory, is scraped to collect the prices referred to an order of a standard meal with drink, delivered to the customer's home. The cities included in the sample are the 12 municipalities in the center of the metropolitan areas.

The last one, in chronological order, is the collection for Gas on the free market; as for the Electricity on the free market, prices are scraped on the Regulatory Authority's public website. To this aim two

profiles are defined according to the type of contract: fixed or variable rate. A sample of detection units identified with a cut-off criterion is selected among the main suppliers, based on the regional market shares derived from the regulated sector survey, carried out annually by Regulatory Authority's (ARERA).

Over the past year, an innovative project has taken its preliminary steps. Istat has launched a new experimentation in the field of web scraping investigating the potentiality of the APIs (Application Programming Interfaces) to collect airfares. Istat has obtained access to the APIs of an international flight search engine which allows users to navigate between the prices and destinations of flight offers and to compare the different offers available.

The test of this data collection started in September 2021 and as 13 consecutive months of data will be available, it will be carried out the comparison with the traditional manual internet data collection and the check of the impact of this innovation.

The results presented in this paper show the magnitude of the information that is becoming available through this new data source to improve the Consumer Prices Index quality.

In the following paragraphs we will focus on the web collection on the gas and electricity free market and the one on airfares through the API's, illustrating the use of the automatic data collection procedures (web scraping) and some of the results achieved.

2 The collection on the gas and electricity free market for the estimation of inflation

In recent years, an important innovation has been introduced for the estimation of inflation, starting with the survey on energy prices on the free market in 2019 for the electricity, and in January 2022 for the gas network, which complements the survey already carried out by Istat on the protected market.

After the start of liberalization, the law established the coexistence of the free market regime and the, so-called, restricted one (Servizio di Maggior Tutela), where the authority for energy the Italian Regulatory Authority for Energy, Networks and the Environment (Autorità di Regolazione per Energia Reti e Ambiente, ARERA) every three months establishes contractual and economic conditions (price).

The free market nowadays represents a substantial part of the electricity and gas market in Italy, thanks also to the changes in the regulatory environment that have encouraged liberalization. The number of families in Italy who have chosen the free market is growing. Domestic customers reached the 59.7% in the electricity sector and 62% in the gas sector (71% for condominiums), confirming a growing trend in recent years¹.

¹ Report on the half-yearly monitoring of retail markets published by ARERA - January 2022.

2.1 Description of the reference market

In the free market, the economic and contractual conditions are defined by the companies that offer different contractual solutions in the context of free competition whereas in the protected market the economic and contractual conditions are regulated by ARERA which updates the tariffs every quarter based on the fluctuations of the value on the raw material market.

All customers, at any time, can choose amongst the various offers available on the free market the one that best suits their needs by stipulating a new supply contract.

The withdrawal, or the choice of a free market supplier, can be exercised at any time and it does not involve additional charges and takes place without interruption of the supply in progress.

The progressive exit of final customers from the protected market in the electricity and gas sector, is continuing at a constant pace, albeit with a certain lack of homogeneity in the national territory (with regard to the consistency of the free market, in the great majority of Italian regions and provinces, more than half of customers have left the system of protection by choosing a free market contract, for both sectors).

2.2 Online Portal of offers

The online Portal of Offers, provided by the Law on Competition, resolution 51/2018/R/com implementing the Italian Law No. 124 of 2017, is the public website where all the contractual offers of the energy market, electricity and gas, must be published. Therefore, families and small businesses can compare and choose immediately, clearly and freely, all offers of electricity and natural gas. It provides a user-friendly search engine and provides information on the operation and expected evolutions of the electricity and natural gas markets.

With reference to the calculation of the estimated annual expenditure associated with each offer, the methods of estimating the user's annual consumption (and its distribution over time) are identified and, in the case of offers at variable price, these estimates are made with reference to the forward values of the price / index indicated in the contract, to take into account the trend in commodity prices.

The online Portal of Offers allows comparison between fixed and variable offers of the free market, as well as the PLACET offers, the offers that all the suppliers will have to give, with pre-established and uniform contractual conditions for all the users. The design and implementation of the Portal of Offers is focused on ensuring the easy consultation by the end user. To this scope, an analysis of usability of the Portal of Offers was carried out. The monitoring of accesses shows that overall in 2021, the site had a total of 992,732 visits (+49.4% compared to 2020). The number of users who use the Portal of Offers has increased both in absolute terms and in percentage terms compared to the total visits. On average, almost 72,000 users visited the Portal in 2021, with a peak in December '21 of over 109,000 users.

At the end of December 2021 the available offers for online consultation were more than two thousands for the domestic users, basically splitted in half between the electricity and the natural gas sector.

2.3 The sampling definition and data collection

For the electricity the products in the basket are:

- ✓ the *Electricity regulated market* is in the basket since 1954; prices collected directly by ISTAT (via internet and other sources). Electricity prices for households in the 'protected market' are updated quarterly by the Regulatory Authority for Electricity and Gas (January, April, July and October); in the prices updating mechanism, the Regulatory Authority takes into account different price components: the price change of raw material oil which is the most important component; cost changes of supply and retail trade for electricity; general system charges which include incentives targeted to renewable and assimilated energy resources, special tariffs for nuclear safety measures; contributes to the research and to the promotion of energy efficiency;
- ✓ the *Electricity free market* (introduced in the basket since 2019); prices are directly collected by Istat through web scraping procedures: to this aim the main regional providers are considered. Eight profiles are defined according to the following three variables: Type of household: residential/non-residential; Type of rate: mono/dual hourly time slot; type of contract: fixed/variable rate.

In order to start the web scraping collection eight profiles have been defined according to the following three variables:

- type of household: residential/non-residential; some Italian regions (Valle d'Aosta, Liguria, Abruzzo, Molise, Calabria, Trentino-Alto Adige, Sardinia and Sicily) have a greater tourist vocation consequently these are the regions with the largest share of non-residents;

- type of time slot: mono and multi-hourly rate; in 2020 there has been a constant and substantial preference for the mono-hourly price, which was chosen by 62.1% of customers (equivalent to 60.7% of volumes); 29.5% of customers chose the two-hourly mode and 8.4% the multi-hourly in reduction of 9.2% compared with the previous year (the last two were grouped in "multi_hourly");

- type of contract: fixed/variable rate; with regard to the type of preferred price, it was found that 84% of domestic customers signed a fixed price contract in the free market (i.e. with the price not changing for at least one year from the moment of signing), while only 16% chose a variable price contract (i.e. with the price changing at a time and in a manner determined by the contract itself).

The prices collected refer to the main regional providers. The largest suppliers have been selected using cutoff sampling according to the market share, resulting in seventy-eight suppliers (collection unit) and six hundred and twenty-four strata (product by region by collection units). Every month around 60,000 prices are scraped, and around 4,000 prices are used to contribute to estimate Italian inflation.

Gas prices are quarterly updated by the Regulatory Authority for Electricity and Gas (January, April, July and October) taking into account the following price components: raw material cost; transport and storage distribution; retail trade, wholesale and additional costs; taxes; price collection is carried out by 92 municipalities out of 107. For the remaining 15, prices are collected through targeted webscraping, and used in combination with the classical price collection.

For Town gas and natural gas - free market, the prices are directly collected by Istat through web scraping procedures, in the portal offers, public website created and managed by "Acquirente Unico", in accordance with the procedures established by the Regulatory Authority for Energy Networks and Environment (ARERA).

Two profiles are defined according to the following variable:

- type of contract: fixed/variable rate; 73.9% of domestic customers have signed a fixed-price contract, while 26.1% chose a variable price contract.

For each product, all the expenses relating to the offers that can be selected by customers with a specific consumption profile (corresponding to our product) are scraped. For each regional capital, which is used as estimator of the region, an average annual consumption is considered, estimated using data coming from the Istat survey on household energy consumption.

At the beginning of the year, a sample consisting of the main sellers is selected, with the cut-off criterion based on the market shares derived from the "Survey on regulated sectors", carried out annually by ARERA (based on law no.481/1995), with a regional detail; the selection considers coverage ranging from 60% to 80% of the market shares for each region.

The prices scraped every month are approximately 10,000, a total of 54 collection units (the collection unit is given by the combination of region code and supplier. For example, the same operator can be repeated n times if it is present in the different territorial areas). One hundred and height strata layers given by the combination of product, regions (19, with the exclusion of Sardinia), and collection unit, from which about four hundred prices are selected and used in the calculation.

For each stratum (defined by the different categories of the variables: product / region / collection unit), an elementary index is calculated. Starting from this, regional index is obtained using the weighted arithmetic mean (Laspeyres index), with weights proportional to the domestic users' consumption, in each region, on the free market.

2.4 The trend of Electricity and Town gas and natural gas prices

The analysis of the trend of the consumer price indices for the Whole Nation (HICP) for *Electricity* and *Town gas and natural gas* shows how in the first months of 2022 prices have increased reaching levels never observed since the series exist for these types of products, that is since 1996.

FIGURE1 shows the time series for two consumption segments of *Electricity* from January 2019 (the year in which the new product *Electricity free market* was introduced) to March 2022.

For the *Electricity regulated market*, after the decrease of prices on monthly basis (the largest in April 2019 and in April 2020), in the following months there were increases, the greatest growth was +47.9% in January 2022, month in which ARERA published the rates for the first quarter of 2022.



Figure 1. Italian harmonized Consumer Price Index (HICP), Electricity. January 2019 - March 2022, monthly price changes (index, 2015=100)

The prices of *Electricity free market* show a dynamics that is lower if compared with the tariffs on the regulated market; from July 2021 the prices of the free market recorded opposite trends of the tariffs of the regulated market in the same month (July 2021 and in October 2021 down by -9.2% and -10.2%, against the +6.4% and +24.8% of the regulated market), whereas in the following month prices increased on monthly basis (the greatest growth was +34.3% in November 2021) showing a one-month lagging behavior with respect to the prices on the protected market.

The FIGURE 2 shows the economic dynamics on annual basis from January 2019 to March 2022.

Tariffs of *Electricity regulated market* between October 2019 and December 2020 were the wide decreases in April 2020 (-13.6%), an accelerating trend begins in January 2021, with a marked trend change in January 2022, equal to +103.4%, a variation that has never been recorded since 1996.

The prices of *Electricity free market* decreased from February 2020 to February 2021 (the largest in May 2020 and equal to -5.9% on annual basis), less than the regulated market; after a positive trend from March 2021, with a peak in June 2021 (+9.1%), in October the prices of the free market recorded a decrease (-7.9%), opposite to the dynamics of the regulated market (in strong acceleration). From October 2021 prices of *Electricity free market* sped up, with the huge increase to +65.5% in March 2022.



Figure 2. Italian harmonized Consumer Price Index (HICP), Electricity. January 2019 - March 2022, annual changes (index, 2015=100)

The **FIGURE 3** shows the economic dynamics on monthly basis from January 2019 to March 2022 for the consumption segments *Town gas and natural gas*.

The tariff of *Town gas and natural gas regulated market* decreased on monthly basis in 2020 (the largest in April 2020 and equal +12.3%), in the following months there were increases, the greatest growth was +39.0% in January 2022.

The *Town gas and natural gas free market* is a new product included in the basket from January 2022, so it is only possible to analyze the dynamics on monthly basis of the first three months of the year: prices increased to +10.7% in January 2022, less than the price of regulated market.



Figure 3. Italian harmonized Consumer Price Index (HICP), Town gas and natural gas; January 2019 - March 2022, monthly price changes (index, 2015=100)

The **FIGURE 4** analyzes the economic dynamics on annual basis from January 2019 to March 2022 for the product *Town gas and natural gas regulated* (as *Town gas and natural gas free market* a new product, it does not have the trends on annual basis).

After a decrease on annual basis of price of *Town gas and natural gas regulated* from July 2019 to March 2021, the tariff increased until they peak in March 2022 (+86.5%)

Figure 4. Italian harmonized Consumer Price Index (HICP),, Town gas and natural gas; January 2019 - March 2022 annual changes (index, 2015=100)



3 API's, the new price collection project

The web-scraped data for the production of consumer price statistics introduced since 2018 is basically a screen data collection, imitating the behavior of consumers visiting websites and buying products or services. As such proved to be a good choice although presents some critical issues and for this reason it was decided to explore new data collection techniques via web.

In the last year an alternative web scraping collection has taken its preliminary steps, involving the use of APIs for the transport by air. Like web scraping this technique allows to gather data from web sites but having direct access to some structured data.

To do so we are using APIs, free of charge, of a web-agency that allowed us to scrape their data on air fares and to have information, as much as possible, to take the needed decisions to allow the use of this data in the production of HICP.

3.1 IT aspects and data flow

For accessing APIs (**FIGURE 5**), an application (named APIstotele) has been created using JHipster, a free and open source application builder used to rapidly develop modern web applications.

APIstotele queries a famous metasearch engine's API to gather flight prices; the raw response, in json format, is saved in an Oracle database.

The API data are grabbed on a daily basis; once the process has finished, an Oracle job starts, parses all the stored json files and saves them in a database schema for further processing.



To observe, analyze and exploit the data a data warehouse used by BI (Business Intelligence) instruments is used.

In order to prepare the queries, the following flights features have been defined: one traveler, in economy travel class, using all the Carrier except Charter. To start, the routes considered are those currently used for the manual collections: 42 national, 61 European, and 65 intercontinental. More over different lags between observation date and departure date have been considered, and the return date was assumed between two and fourteen days.

The APIs requests run from midnight for about 7 hours, due to the search engine API restrictions for free access everything has been set in order to have at most 100 requests per minute. In general we have been collecting more than 8 million quotations per day, in terms of bytes this means 300 bytes per quotation and more than 73 GB per month.

3.2 Some results

TABLE 1 shows data captured since August 2021 for national, European, and intercontinental destinations for departures in October, November, and December 2021. The number of prices scraped for the European destinations is increasing in the period considered and is the highest in December exceeding 400,000 quotations. These huge numbers are due to the fact that for each route, date and time of departure, carrier, agent, the number of price quotes available depends on the distance in terms of days and time of the back travel.

Month	Туре	No. of travels	No. of prices	
	National	12,429,296	37,068,356	
October	European	57,888,087	172,008,016	
	Intercontinental	24,801,328	57,201,724	
	National	11,210,247	36,640,611	
November	European	113,244,065	343,102,698	
	Intercontinental	25,114,469	57,190,815	
	National	12,203,525	36,923,383	
December	European	151,192,155	451,802,783	
	Intercontinental	30,919,323	69,072,548	

Table 1 Number of observations for type of destination and month of departure (Year 2021). Absolute values

TABLE 2 shows data collected for flights departing in the same period as those for specific routes, one for each group. Also in this case the greater number of prices are collected for the European destination.

Table 2. Number of observations for single destinations and month of departure (Year 2021). Absolute values

Month	Туре	Departure	Arrive	No. of travels	No. of prices
	National	Milan Linate (LIN)	Rome Fiumicino (FCO)	923,937	4,129,568
October	European	Rome Fiumicino (FCO)	Barcellona El Prat (BCN)	2,767,653	9,038,992
	Intercontinental	Milano Malpensa (MXP)	New York JFK (JFK)	800,193	1,578,074
	National	Milan Linate (LIN)	Rome Fiumicino (FCO)	1,042,642	5,328,989
November	European	Rome Fiumicino (FCO)	Barcellona El Prat (BCN)	2,594,148	8,631,414
	Intercontinental	Milano Malpensa (MXP)	New York JFK (JFK)	776,235	1,618,229
	National	Milan Linate (LIN)	Rome Fiumicino (FCO)	1,087,397	4,776,265
December	European	Rome Fiumicino (FCO)	Barcellona El Prat (BCN)	2,970,285	10,858,359
	Intercontinental	Milano Malpensa (MXP)	New York JFK (JFK)	849,750	1,800,686







Figure 7. Number of observations per route and month of departure (Year 2021). Absolute and percentage values

In the following **FIGURES** (8-10) the average prices for the three routes are showed, in all cases it has been considered the same time slot for departure and return and the data have been grouped by number of days in advance of purchase.

FIGURE 8 shows that the level of the average prices depend on the lag between the purchase simulation and the day of departure. The higher level of prices on the weekdays might reveal that this route is used mainly for work reasons.



Figure 8. Departure/return average price by advance class for route Milan Linate - Rome Fiumicino (Year 2021)

The European route showed in **FIGURE 9** also shows lower level of the average prices when the lag between the purchase simulation and the day of departure is higher in all the three months



Figure 9. Departure/return average price by advance class for route Rome Fiumicino - Barcelona El Prat (Year 2021)

FIGURE 10 shows the average prices for the international route Milan Linate-New York JFK. The higher prices in December might be the effect of the relaxation of travel-related control measures to contain the COVID-19 pandemic.



Figure 10. Departure/return average price by advance class and time slot for route Milan Linate -New York JFK (Year 2021)

4 Concluding remarks

In Italian CPI/HICP the use of web scraping collection has been increasing in recent times to complement or replace traditional survey sources, to reduce burden on the respondents and costs of the data collection and to improve the accuracy of the CPI and it has allowed to deal with the issues deriving from COVID 19 crisis.

The work done until now is very promising. Our plans foresee going on with the massive data collection strengthening IT infrastructures; deepening the analysis of the data collected by several variables and extending experimentation to other groups of product like clothing and hotel facilities.

The different techniques considered till now are not perfectly interchangeable, they have different pros and cons, and the choice of one technique rather than another is determined by the type of information to be collected and by the structure of the website being queried.

Technical issues including frequent changes in the website structure, as well as the presence of automatic blockage of high frequency web scraping may lead to a preference for API techniques. Instead, the automatic data collection procedures might be preferable when the query of a site and the choice of different options lead to a calculation obtained on the basis of hypotheses that may change over time as for the energy tariffs.

References

ARERA (1 February 2022) "Monitoraggio sull'evoluzione dei mercati di vendita al dettaglio dell'energia elettrica e del gas - Rapporto di aggiornamento di gennaio 2022" (RAPPORTO 37/2022/I/COM)

ARERA (1 February 2022) "Annual report to the international agency for the cooperation of energy regulators and the European commission on the regulatory activities and fulfilment of duties of the Italian regulatory authority for energy, networks and environment (Report 344/2021/I)

ARERA (1 January 2022). Report on the half-yearly monitoring of retail markets

ARERA "Relazione annuale stato dei servizi 2020- Volume 1 e Volume 2"

Istat (4 February 2019) "Gli indici dei prezzi al consumo: aggiornamenti del paniere, della struttura di ponderazione e dell'indagine. Anno 2019", *Nota Informativa* (https://www.istat.it/en/archivio/226765), Rome;

Istat (2 February 2022) "Gli indici dei prezzi al consumo: aggiornamenti del paniere, della struttura di ponderazione e dell'indagine. Anno 2022", *Nota Informativa* (https://www.istat.it/it/archivio/265952), Roma;

Istat (15 April 2022) "Consumer prices. Final Data- March 2022", *Press release* (<u>https://www.istat.it/en/archivio/269478</u>), Roma.