Study of new data sources and techniques to improve CPI compilation in Brazil: first steps

High-level seminar on the The Digital Economy: A policy and Statistical Perspective Beijing, Cina, November 15-17, 2018

Vladimir Gonçalves Miranda

vladimir.Miranda@ibge.gov.br

Instituto Brasileiro de Geografia e Estatística IBGE



E3 IBGE Instituto Brasileiro de Geografia e Estatística IBGE

Structure of the presentation

- i) Brief introduction to price indices at IBGE and some details of our IPCA.
- ii) Introduction and motivation of the studies.
- iii) Study of web scraping as support for hedonics.
- iv) Study of web scraping techniques to collect airfares.
- v) Final remarks: other projects and next steps.

Brief introduction of Indices at IBGE

Brief introduction of Indices at IBGE

Indices currently produced: PPI (at industry coordination), CPI and construction (at prices indices coordination - COINP).

RPPI under development at COINP.

Recent change in the structure of the prices indexes coordination.

We have many projects focused in the improvement of work routines and the accuracy and methodologies of our indices.

Among such projects is the study of new data sources and scraping techniques to improve the CPI. This is the focus of this presentation.

Brief summary of our CPI

Brazil' s most important CPI measure is the IPCA.

The index covers families with an income between 1 and 40 minimum wages.

Geographically, 16 states are covered. Which covers about 90% of the total population according to our HBS, adopting the income as the weighting criterium.

Central office is situated at Rio de Janeiro.

There are local units in each state.

Over 480.000 prices collected each month for approximately 380 subitems (sub-classes in COICOP).



MBGE

Introduction and motivation

The ONS challenges in the digital era

Provide more data information in a timely, attractive manner and with robust scientific methods.

(Of course this should be performed with a reduced budged and staff).

Deal with the "breakdown" of the monopoly of providing information for some fields where private companies used to have no or limited data access (CPI for instance) due to large costs to get them.

In this scenario, the use of new data sources and new techniques to improve the surveys and methods is mandatory.

Big data sources for price indices: a not so recent story

New Technologies

New data sources and commerce practices.

Main data sources: scanner data, administrative records and web sites.



Big data sources for price indices: a not so recent story

Internet, new source of commerce and data: e-commerce "borns" with the internet but intensifies after the mid 90' s (amazon and ebay birth).



Use of Big data sources for price indices: a recent story

Though the data sources are not so recent, its massive use in price indices is.

Initial proposal of use of scanner data in CPIs in 1994 (work presented at the Ottawa group meeting by Saglio). However, until now the number of countries that have implemented it in their CPIs is small (around 10), though the interest and number of adepts is growing.

Massive use of scanner data date' s back the early 2000' s. Use of web data for prices indices has approximately a decade, triggered by MIT' s Billion Price Project (Cavallo and Rigobon, Journal of Economic Perspectives, 2016).

Why the delay?

Access to such data: necessity of legislation or negotiation to access the data. Critical for scanner data and administrative records.

IT infrastructure and techniques to deal with such data.

Main Uses of new data sources for CPIs

- i) Improve the traditional ways of data collection and work routines.
- ii) Capture new forms of commerce and improve sample representativeness.
- iii) Development and improvement of methodologies in price indexes.
- iv) More frequent update of weighting structure.
- v) Extension of the number of goods in the CPI basket.

Our initial steps and choices: how and where can we use web data?

We also want to join the game.

We still do not have access to scanner data and adminitrative registers on prices transactions (dream whishlist).

But we can have access to web data.

Inital ideas on using such source:

Use of web data to improve collection methods (specially in sectors where prices are already obtained via web sites).

Use of web to improve index compilation: implementation of hedonics for quality adjustment.

Use of web prices is a good choice to start dealing with such big

(Traditional) CPI based on a fixed basket of goods and services. "Same" products should be compared between months Matched model method.

Same product, in the same outlet defined in the reference period 0, should be collected in subsequent dates.



Market dynamics implies that products have a finite lifetime. Hence, oftently items need to be replaced in the basket.

Replacement goods/services may be of different quality respective the old ones.

Some examples of quality change:





Airfares that used to contain meal and luggage in the price and now are bought apart.

This process bias due lac adjustment.	s leads to k of quality	$r_{\{t+3,t+1\}}$	$-2\} = \frac{p_{\{n\}}^{\{t+3\}}}{p_{\{m\}}^{\{t+2\}}}$	} - }	
Item/period	t	<i>t</i> +1	t+2	<i>t</i> +3	<i>t</i> +4
1	p_l^t	p_l^{t+1}	p_l^{t+2}	p_{l}^{t+3}	p_l^{t+4}
т	p_m^t	p_m^{t+1}	p_m^{t+2}		
п				p_{n}^{t+3}	p_n^{t+4}

Problem is more serious with high tech products, products with some sort of depreciation and those with high turnover.



(Van Loon and Roels, UNECE CPI meeting 2018)

Most celebrated method to deal with these examples is hedonic modelling.

This essentially states that each good is composed by a bundle of attributes and each has a marginal contribution for the good's final price.

The problem is that markets usually do not reveal the prices of the attibutes and this need to be estimated. That is what the modelling does.

Main message: to deal with this, we need to have data on prices of products and products most importante characteristics.

This data is used to build hedonics "patching" (low rate of substitutions) or hedonic indexes for products with high turnover or depreciation (used-cars example).

Formulation for patching: multivariate regression based on the item characteristics "Z":

 $Price = \beta_0 \,\beta_1^{z_1} \,\beta_2^{z_2} \,\beta_3^{z_3} \,... \beta_n^{z_n} \,\varepsilon$

 $\ln \text{Price} = \ln \beta_0 + z_1 \ln \beta_1 + z_2 \ln \beta_2 + z_3 \ln \beta_3 + \dots + z_n \ln \beta_n + \ln \varepsilon$

Model allows the estimation of new item in the previous period based on data on t+2 and account for quality adjustment.



Problem is that to get information on products characteristics via field collection is very costfull, most demanding for the collector and increases respondant burden.

It is also more difficult to control the process, for instance, guarantee that the correct attributes are being collected etc.

Using the web data we can get such data in a cheap, controlled and efficient manner.

Our pilot: again build a home-made scraper using R to extract the products characteristics and prices of selected goods from the most important retailers of household appliances.

We started with refrigerators and extract data on:

Dimensions and weight;

Capacity;

Kind of door (oneside, two side or inverse);

Presence of water and ice dispenser or both;

Power consumption;

Coating material;

Make.

CARACTERÍSTICAS E FUNÇÕES DO CONGELADOR / FREEZER				
Recursos do Congelador/Freezer		<u>Compartimento de Congelamento</u> Rápido Fábrica de Gelo		
Tipo de Degelo	O que é isso?	Frost Free		
EFICIÊNCIA ENERGÉTICA				
Eficiência Energética / Faixa Selo Procel	? O que é isso?	Selo Procel A (Mais Eficiente)		
RECOMENDAÇÕES PARA UTILIZAÇÃO E SEGURANÇA				
Antes de Utilizar o produto recomenda-se		Consultar manual de instruções. Verificar se o produto possui selo de Certificação do Inmetro.		
ESPECIFICAÇÕES TÉCNICAS				
Altura		<u>184 cm</u>		
Largura		<u>70.5 cm</u>		
Profundidade		76,7 cm		
Peso		73,5 kg		
Tensão / Voltagem		127V 220V		
Consumo		55 kWh		

Output of the regression model (default method used for patching) using step-wise approach:

Call: lm(formula = Preco ~ Acabamento.Externo.c Dispenser.Externo + Marca + Tipo.de.F	da.Porta + Capacidade.Total + Porta, data = dataframefinal3)			
Residuals: Min 1Q Median 3Q -1822.31 -277.27 -70.94 147.40 2678	Max 3.50			
Coefficients:				
Estima	ate Std. Error t value Pr(> t)			
(Intercept) -4.12	216 541.5261 -0.008 0.993940			
Acabamento.Externo.da.PortaInox 285.18	361 122.7746 2.323 0.021830 *			
Acabamento.Externo.da.PortaVidro 997.64	492 706.6592 1.412 0.160539			
Capacidade.Total 6.08	339 0.9389 6.480 1.99e-09 **			
Dispenser.ExternoÁgua e Gelo 4183.44	19 595.7203 7.022 1.30e-10 **			
Dispenser.ExternoNenhum -115.14	427 393.1984 -0.293 0.770141			
MarcaConsul -471.78	383 177.8074 -2.653 0.009021 **			
MarcaElectrolux -463.25	516 153.9039 -3.010 0.003170 **			
MarcaPanasonic -460.60	027 232.3431 -1.982 0.049661 *			
MarcaSamsung 687.89	J29 244.8490 2.809 0.005775 **			
Tipo.de.PortaDuplex 132.10	579 252.8819 0.523 0.602160			
Tipo.de.PortaFrench Door Inverse 1647.01	L82 430.7828 3.823 0.000208 **			
Tipo.de.PortaInverse 890.59	J24 327.2040 2.722 0.007435 **			
Tipo.de.PortaSide by Side 1134.98	332 598.6418 1.896 0.060315 .			
Signif. codes: 0 '***' 0.001 '**' 0.01	**' 0.05'.' 0.1'' 1			
Residual standard error: 656.1 on 123 degrees of freedom				
Multiple R-squared: 0.9158, Adjusted R-squared: 0.907				
F-statistic: 103 on 13 and 123 DF, p-value: < 2.2e-16				

Price = -4 + 285*stainless_steel_coat + 997*glass_coat+ 6*total_capacity + 4183*water_ice_Dispenser - 115*no_Dispenser + ...

Process seems promising to adopt for hedonic patching. But we still need to perform some additional studies on other products and test the adequacy of the estimated models.

Also necessary to study the use of such technique for hedonic indices.

This approach is also useful to identify absence of products and introduction of new products in the Market.

Important open issues:

- i) Can the web prices be used to calibrate the models? In principle they can provide much larger samples at cheaper costs that can lead to more precise models parameters.
- ii) Should they be used at least for small sample products where the models obtained are not reliable?

Improving collection techniques: airfares

I.

Improving collection techniques: airfares

Tradional price collection: each of the 16 local units collects manually the prices of selected routes in the web sites of main airline companies.

Previous Week Also change my inbound Next Week Also					
Sat Sun 08 Sep 09 Sep from from \$595 \$723	Mon 10 Sep from \$510 Tue 11 Sep from \$550	Wed Thu Fri 12 Sep 13 Sep 14 Sep from from from \$550 \$530 \$550		Travel Classes	
Departs	Arrives	Flight Operator	Economy	Premium Economy	Business Class
	🕉 Sorry	, there are no direct flights for th	is route, flights with connec	tions are below. 🛛 🐳	
Outbound cor	nnecting flight opt	ions			
21:50	18:00	British Airways GIG-LHR	©\$550	○\$1450	○\$1425
GIG	GVA	British Airways LHR-GVA	Lowest	Lowest	Lowest
Ø 1 Connection	Show journey details +		.i	🕒 Total journey t	ime:15 hours 10 minutes
21:50 20:15	British Airways GIG-LHR	◎\$550	●\$1450	○\$1425	
GIG	GVA BA CityElyer LCY-SVA		Lowest	Lowest	Lowest
1 Connection Show journey details + Show journey details +					
21:50 20:25 11 Sap 12 Sap GIG GVA	British Airways GIG-LHR	◎\$550	○\$1450	○\$1425	
	British Airways LHR-GVA	Lowest	Lowest	Lowest	
1 Connection Show journey details +					
21:50 22:25 11 Sep 12 Sep GIG GVA	22:25	British Airways GIG-LHR	◎\$550	●\$1450	○\$1425
	GVA	British Airways LHR-GVA	Lowest	Lowest	Lowest
1 Connection Show journey details + O Total journey time: 19 hours 35 minutes				ime:19 hours 35 minutes	

Collection for all flights for predetermined routes, for tickects bought 2 months previous to the departure date and some conditions for the departure and arrival dates and days.

Data collected for different tickets categories.

Collection for different airline companies.

Data collected once a week.

Improving collection techniques: airfares

- I. Process is time demanding: approx. 4h for each collection. About 16h a month for each area.
- II. More subjective to errors. Demands extra analysis time by central office team. 1-2 hours montly.
- III. Since last quarter of 2017, extra collection was necessary due the implementation of the continuous ICP program of CEPAL:
 - a) new routes and companies added.
 - b) Collection 3 weeks per month for 2 areas.
 - c) Extra 6h of manual collection per month.

Main question: Is it possible to improve this process using

We need a Web scraper: a program that extracts the data of the page.

Html behind the web site of an airline company: need to arrange the data into a structured format for use.

Reexibir todo

OS VOOS

Organizar por MAX PLUS LIGHT decolagem mais Image: The 2b bagagens gratuitas Image: The 2b bagagens gratuitas Image: The 3b bagagen gratuita Image: The 3b bagagen gratuita <th>1 opções de voos</th> <th>Filtro +</th>	1 opções de voos	Filtro +		
e mais vantagens + e mais vantagens + R\$ 733,17 R\$ 643,17	Organizar por decolagem mais	MAX □ 1° e 2° bagagens gratuitas ✓ R\$ 1 = 4 milhas Smiles ✓ Assento GOL+ Conforto gratuito ✓ Antecipação gratuita	PLUS ☐ 1º bagagem gratuita ✓ R\$ 1 = 3 milhas Smiles ✓ Marcação de assento gratuita ✓ Antecipação gratuita	LIGHT gratuita ✓ R\$ 1 = 2 milhas Smiles ✓ Marcação de assento gratuita no período de check-in
		<u>e mais vantagens +</u> R\$ 733,17	<u>e mais vantagens +</u> R\$ 643,17	e mais vantagens + MENOR PREÇO DO DIA

```
<div>
       <div class="lessPriceBox">
        <div class="lessPrice">Menor Preco do Dia</div>
       </div>
       <div class="taxaSelected">
        <div class="checkTaxaSelected"></div>
      </div><span class="smlesAndMoneyValue"></span>label class="textIdentFareValue" for="
                           /</span><span>tarifa Light</span></label><input id="ControlGr
                     R$ 598.
     </div>
   <div id="market2_journey1" class="infoGrid bgGrid popupANAC"><span class="anacInformatio"</pre>
                 </span></div>
```

um 1 1 a 11

at Constant to Constant

Extra demand: it may be necessary to emulate an user navigating in the page.

×	Só voo	As minhas pesquisas recentes De Brasil	Choice of origin
•	Ferramentas de viagem	Rio de Janeiro ▼ ● viagem de ida e volta ● Só ida Várias cidades/escalas Para Digite 3 letras	and destination.
		11/09/18 Image: 28/09/18 Classe de voo ❶ Económica Tipo de bilhete ❶ Bilhete flext(vol)	Choice of travel dates.
		Image: New Procurar voos	Choice of ticket kind.

Pilot home-maid web scraper built combining R and selenium tools. Project developed along with the Methods and Quality Coordination.

Initial steps, focus on scrape data for airfares to reproduce the manual collection processes.

Results:

Robots take about 30min to perform the collection of an area, 8 times faster than the manual process. A single desktop machine was used, so this time can be decreased by using more machines and parallelising the process.

Results (continuation)

In a controlled test ("sincronized" robots and manual collection), first comparison between manual and automatic process showed a high agreement between prices collected as expected. Divergences essentially due to small diferences in the time of collection.

Scrapers were also built for the ICP program and are also used in this Project.



Some challenges for the implementation

Legal issues: anti-robots politics.

Please show you're not a robot

I'm not a robot	2
	reCAPTCHA Privacy - Terms

Each site has its own "design" .

Instabilities: sites change without a previous warning.



Cookies: according to a certain profile different prices might be offered.

How to get the support of the respondent and be safe against prices manipulation?

Current status:

Implementation of this for the CPI is more demanding since it is a continuous process that doesn't allow failures.

IT team developed a tool to implement the automatic process (combining C# and Selenium and based in the COMEQ' s code) in the production routine and is running more tests in order to compare manual and automatic results and reporting the challenges to put this in production.

Error control system is being developed and is already under test in order to guarantee that the robots are collecting the correct data from the pages. This involves an automatic print of the screen from which data was collected which allows

Current status:

We had a meeting with representatives of the air companies, mediated by ANAC (national regulatory agency of the air sector), to ask for collaboration.

We introduced the methodology and pointed the challenges to implement this process in production.

Some issues treated: robots block, API access, adequacy of our model, sites instability, massive data collection.

Some important lessons:

Robots access: relatively simple, the company just needs to include the ONS IP in a 'whitelist', though requires some negotiation.

Data collection: the site of the air companies usually are maintained by a third party and they pay for each attempt made to get a price quotation. So if you perform massive access you are increasing their costs and that's why they want to block you.

Sites instability: is it possible to know in advance when a site change will occur?

API access: possible, but requires more negotiation with commercial sector since they sell this product.

Some hints on measuring inflation for digital products: Fitting into the matched model method

"Another version of matching methodology is to select some elements of a tariff as "representative items" and re-price them in subsequent periods. For instance, for air fares this could be for each airline carrier, a non-refundable and non-changeable airline fare from one pre-specified location to another, with predetermined outbound and inbound dates chosen by time of day and day of week, including all surcharges."

"Once chosen in the base period, the route, the departure and arrival times and the ticket type and class of travel should remain the same throughout the year"

"Transport ticket prices should feed into the index at the time of travel, not at the time they were booked. So, a ticket price for a December flight should feed into the December index even if the price relates to a purchase in October." CPI manual, Chapter 11.

We have a HBS running and in the new basket we will probably have products like Uber.

We will be faced with the question of how to fit it into the MMM. Looking at the airfares case, we have a hint on how to act. However, some important issues arise:

- i) The number of routes here is much larger;
- ii) How to distinguish users that fall into the CPI reference population?
- iii) How to get access to such data? Here we can' t have access even to offer prices. We would need to negotiate some special user access which can see the offer prices or get the transacted prices for some trips.

Final remarks

Conclusions, other projects and next steps

The results obtained by web collection implemented for airfares and quality adjustment seem very promising and we are performing some extra tests before implement the results in the CPI.

Robots collection already running for the airfares of the CEPAL's ICP program.

We are also working together with CEPAL in a pilot Project to improve price collection techniques for the ICP.

We are initially developing collection for housing rents. CEPAL has already built a scraper and has some data collected. We are at the moment giving technical support on the kind of data they should get and interpret the results. We are also going to build a scraper to compare the results obtained and check if such data can be used to improve our CPI.

Other projects and next steps

Keep negotiating with air companies to guarantee the data access. (Though we are in a digital era, people still role the world.)

Check all the itens whose collection have potential to be replaced by a robot collection. Some items like airfares are essentially data collected from the web.

Implement web stores in our sample to improve its representativity, according the results of our most recent HBS.

Develop internet indices to acquire skills in this area and get prepared to use scanner data.

E2IBGE

Thank you for your attention.