

MEASURING VARIATIONS OF CPI PRICES THROUGH WEB SCRAPING



Strategy

/DANEColombia

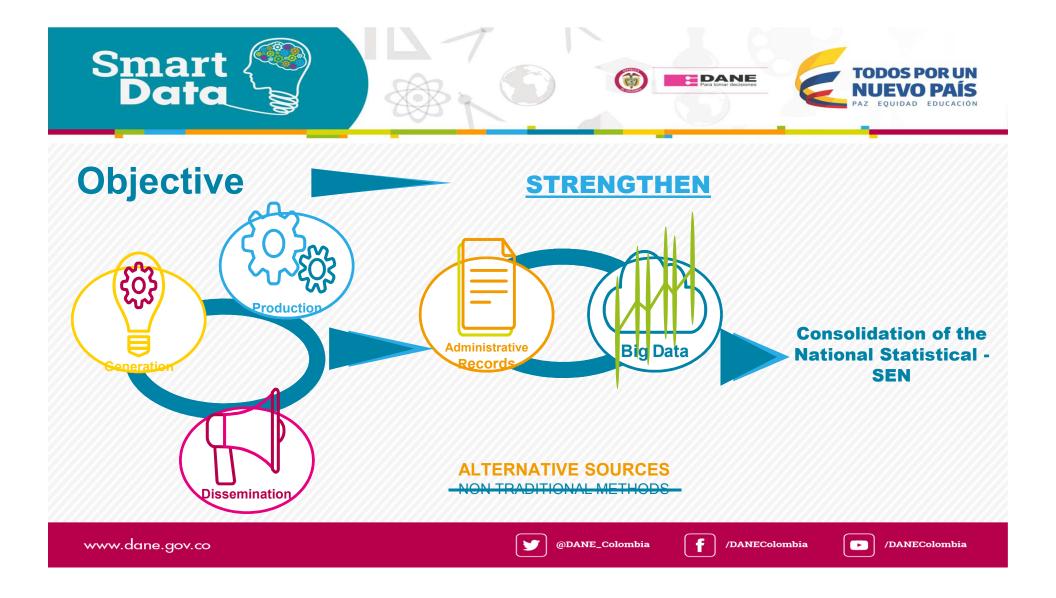
f

August 2016

www.dane.gov.co

@DANE_Colombia

/DANEColombia





Innovation Plan 2015 - 2018





MEASURING VARIATIONS OF CPI PRICES THROUGH WEB SCRAPING

www.dane.gov.co



@DANE_Colombia

/DANEColombia

f

/DANEColombia





To propose an alternative method of gathering information for the CPI, in which the use of alternative sources (e. G. Websites supermarkets) is done through web scraping techniques.





/DANEColombia

/DANEColombia



Specific objetives

- Reduce costs and time in the process of collecting prices for the CPI.
- Show the relevance of applying alternative sources of information for statistical production.
- Exploring alternative uses of prices and products obtained through alternative sources of information.
- Develop algorithms to automate the process of extracting data from websites selected supermarkets.



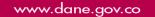


Colombian CPI

. Groups: Households, Foods, clothing, health care, education, recreation, transport, communications, and other expenditure.

- Sub groups: 34
- Classes: 88
- basic expenditures: 181

Sources: 58600 monthly **Quotes:** 211.000 prices monthly



@DANE_Colombia

/DANEColombia

f

/DANEColombia

Design of the pilot and expected outcomes



- Target: supermarkets websites
- Coverage: Bogota prices
- Variable of analysis are the variations of prices
- Frecuency: 15 days

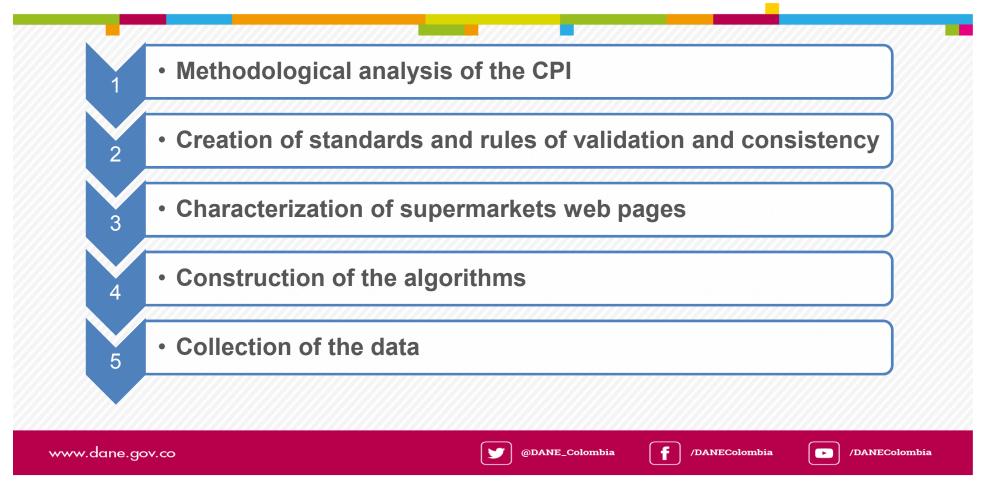
OUTCOMES:

- 1. The collecction program (code)
- 2. Data base: code of products, date, specifications including unit, price, quantity.
- 3. Comparable tables between the downloaded data and the collected by the traditional collection method.





Phases of execution



Methodological analysis of the CPI



Taking to account:

- 1. Basket
- 2. Structure : levels, weights
- 3. Collection frequency
- 4. Products specifications: brand, variety, unit, primary and secondary characteristics
- 5. Sources: supermarkets
- 6. Application of technical novelties





Creation of standards and rules of validation and consistency

- 1. The algorithm must ensure the scraped data contains the information required: specifications. Since there are several specifications the algorithm must select those with a higher frequency of occurrence and are relevant to the product in question.
- 2. As part of the preparation of the data, there must be a data cleansing of duplicates, missing values etc. There is also a filter to identify products with complete information.





/DANEColombia

/DANEColombia



Characterization of supermarkets web pages

- Identification of CPI supermarkets and revision of its web sites:
 12 sites.
- 2. There is no standardization in the structure of the web sites or in the way the information about the products is presented.
- 3. 6 web pages with favorable conditions for the scraping.
- 4. It was fount that the construction of specific algorithms for each web page was required.





Construction of the algorithms

- Exploratory analysis: web scraping through Wget command + data cleansing with Stata. Wget is a command created in C language for UNIX systems (particularly used in Linux) to download information from any website.
- 2. web scraping using R: has the advantage that the algorithms is designed to only download the information required, excluding pictures or other unnecessary objects for price analysis.
- 3. The algorithms were built to search a common pattern in the presentation of the data for products.





Collection of the data



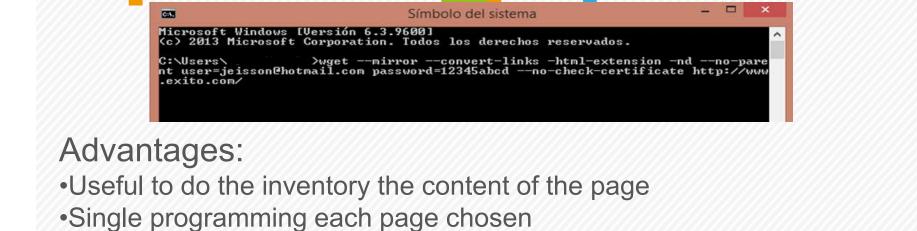


Supermarket Website example

🛄 Aplicaciones ★ Bookmarks 🛐 Facebool	k 🚺 YouTube - Broadcas 📕 T	Construction for the second states and the second states are set of t			confusi » ն	Otros marcadores
		Puntos Éxito y	más <mark>Servicios</mark> para ti, <mark>Clic ac</mark>			î ()
exito	m para servirte				Iniciar sesión Registrarse	
		Todos los productos 🧿 Lo qu	ie deseas, esta en exito.com	Buscar 🤇	Comprar 0 productos	
	tul Bassa Maste	Denotes				
Mercado Salue belle		Tecnologia Deportes entretenimie			A Ofertas Mundo éxito	
TV v video Computado	res Tablets Teléfonos y	Consolas Audio	Cámaras Lanza			
e impresora	as rabiets celulares	y videojuegos y reproducto	ores y videocámaras Canza	Atrás Reenviar	Alt+Flecha izquierda Alt+Flecha derecha	
tecnología > teléfonos y celulares > Celulares y accesorios x Volver a cargar					Ctrl+R	
Mi ciudad Bogotá		PUBLICID	AD	Guardar como	Ctrl+S	
	Mostrando 1 - 20 de 144 resultados			Imprimir	Ctrl+P 8 *	
Cambiar ubicación	Ordenar por: Relevancia >	×	Ver como: Lista Cuadrícula	Traducir a español		
Tecnologia	Relevancia >			Ver codigo fuente de pagina	Ctrl+U	
Celulares y accesorios				Ver información de <mark>l</mark> a página		
Celulares (108) Accesorios para celular (35)				O AdBlock		
	Lexar	Lexar	08:08	Inspeccionar elemento	Ctrl+Mayús+I	
Marcas	4GB Miss	8GB Misse	and the second se	MOBILE HIROSOMC		
MOTOROLA (6)	© E	@ E E				
SAMSUNG (32) APPLE (26)						
HUAWEI (15)				- And and a second s		
🔲 LG (1)	Memoria Micro Sd De	Memoria Micro Sd De	Huawei Ascend Y330	Memoria Para	Huawei Ascend G730	
📄 L G (20)	4gb Para C -	8gb Para C -	White - Y330 WHITE	Celular+adap Usb -	Black - SIN REF	
NOKIA (3)	LSDMI4GBAC	LSDMI8GBAC	HUAWEI	LSDMI8GBBC	HUAWEI	T
	🗩 💌 💌				- 记 🖿 😫 📶 🌵 🗆	ESP 8:01 p. m. 26/05/2015



Command wget



Disadvantages:

- •Download all the information even "junk" (.txt , .png , .jpg) .
- •Inefficient in terms of download time.





Results and lessons learned

- 1. The websites of selected as sources turned out to be perfectly viable to using web scraping technique.
- 2. The method is valid to obtain prices and specifications of several products of the CPI mostly food products, that represent the 28,1% of the CPI.
- 3. The method is viable to complement the collection process of the CPI.
- 4. The use of free software as R is a good alternative to explore new methods in a cost effective way.





