

CHALLENGES AND SOLUTIONS TO THE USE OF INTERNET DATA IN THE DUTCH CPI

Robert Griffioen, Olav ten Bosch and Els Hoogteijling
Workshop on Statistical Data Collection, The Hague, The Netherlands, 3-
5 October 2016



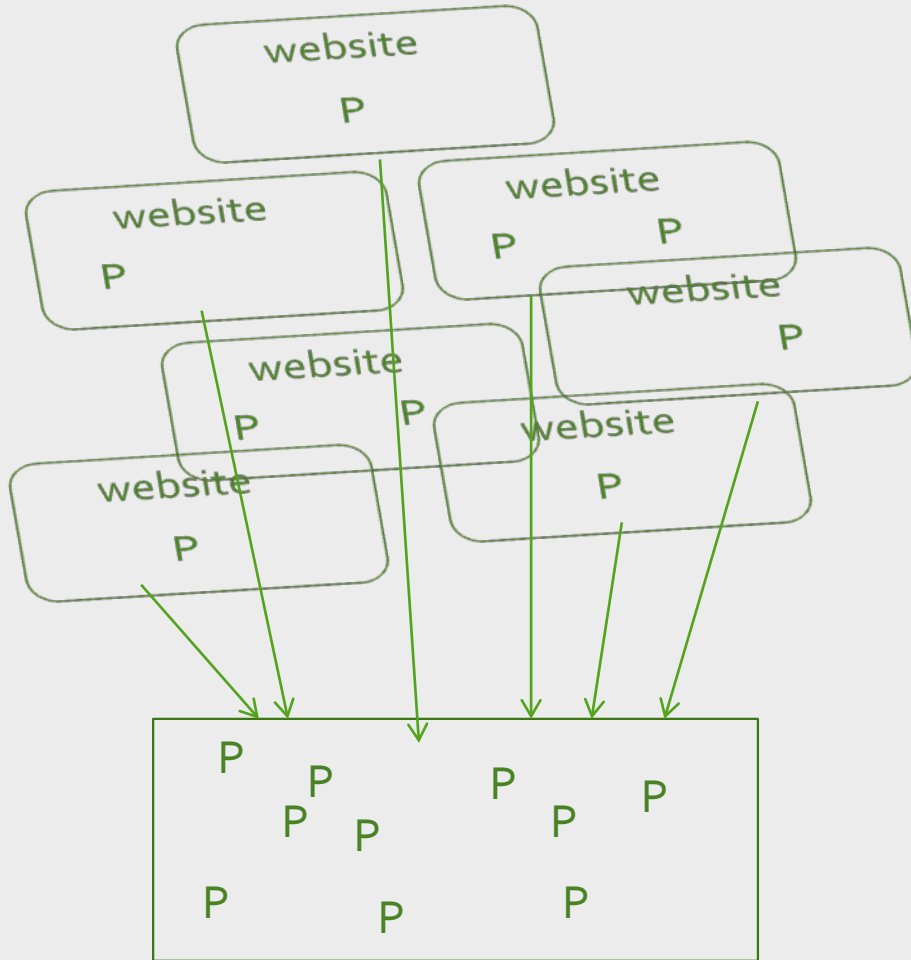
Statistics
Netherlands

Overview

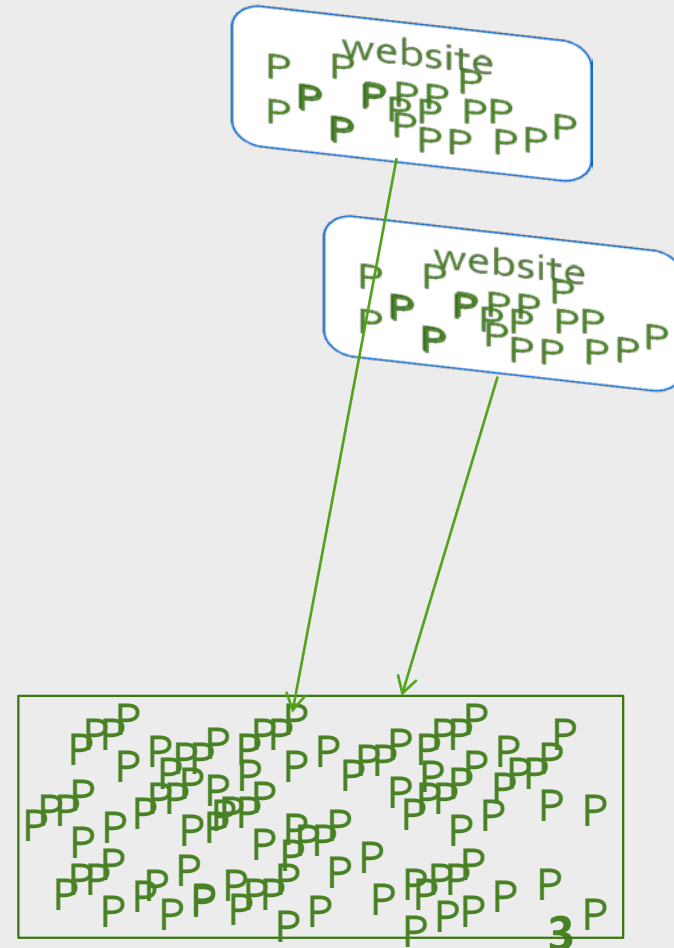
- Introduction CPI and internet data collection
- Core value of using internet data in the CPI
- Technical, organizational and legal challenges of using internet data in the CPI
- Other challenges
- Conclusions and discussions

Introduction: a two tier approach of web scraping

Robot tool



Bulk-scraping

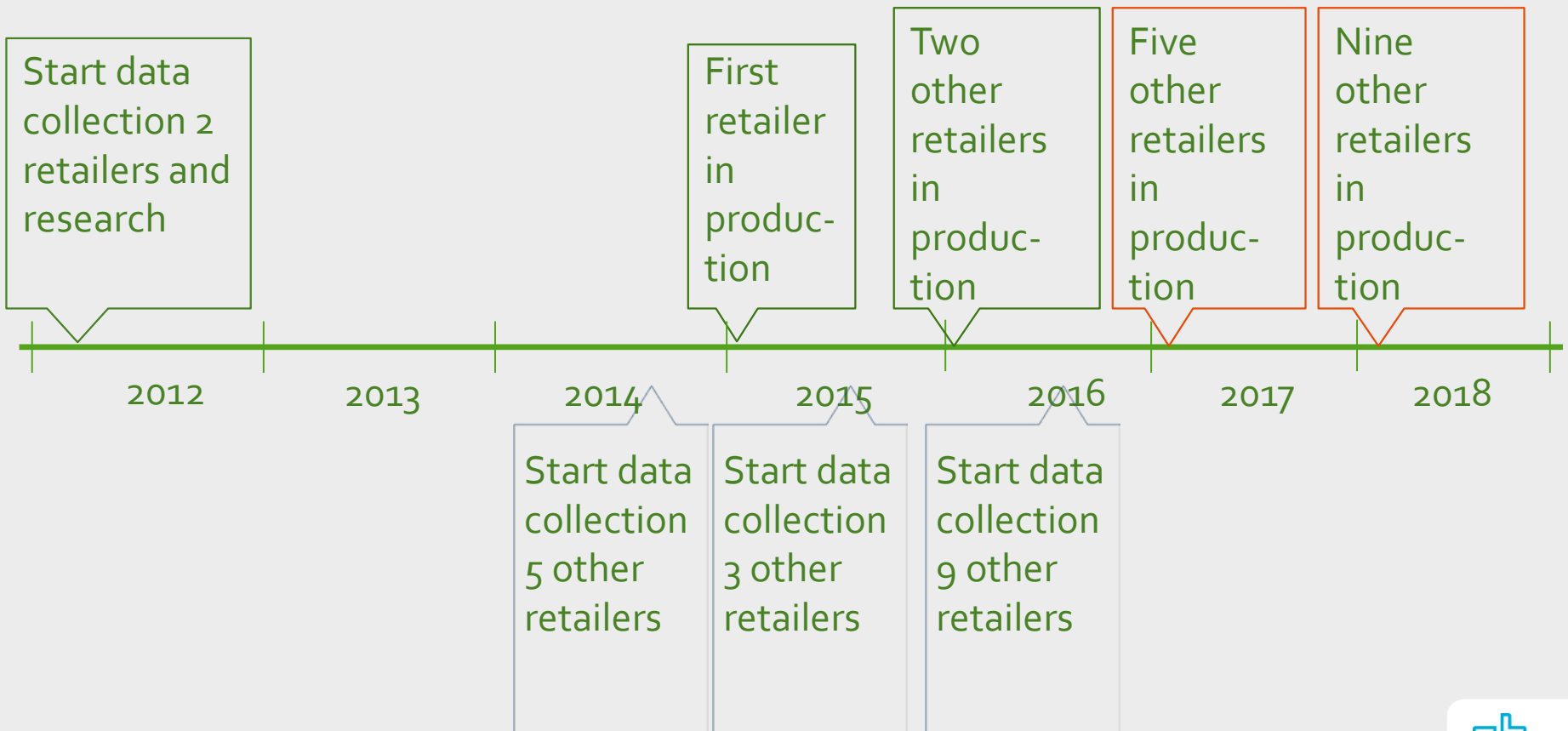


Introduction: history and plans robot tool

Robot tool: started at 2012, 2014 in production

ECOICOP	
31410	Laundry and dry cleaner's (service)
62390	Paramedical services (like acupuncturist, chiropractors and physiotherapist)
72300	Services for maintenance and repair of private cars
72430	Driver license lessons, driver licence examinations and cost of a driver license.
73220	Cost of a taxi
94120	Cost of participation in recreation and sports
94210	Cinema, theatre and concerts
94220	Museums, libraries and zoos
94250	Photography services (for instance printing)
111110	Catering services
111120	Accommodations

Introduction: history and plans bulk scraping



Core value internet data collection

- Internet is a very important medium in our daily lifes:
 - Communication
 - Information/consulting
 - Entertainment
 - Education
 - Shopping
 - Et cetera ... internet is becoming a natural part of our daily lifes
 - → If official statistics does not use it, we may develop a bias towards traditional channels
- Cost reduction



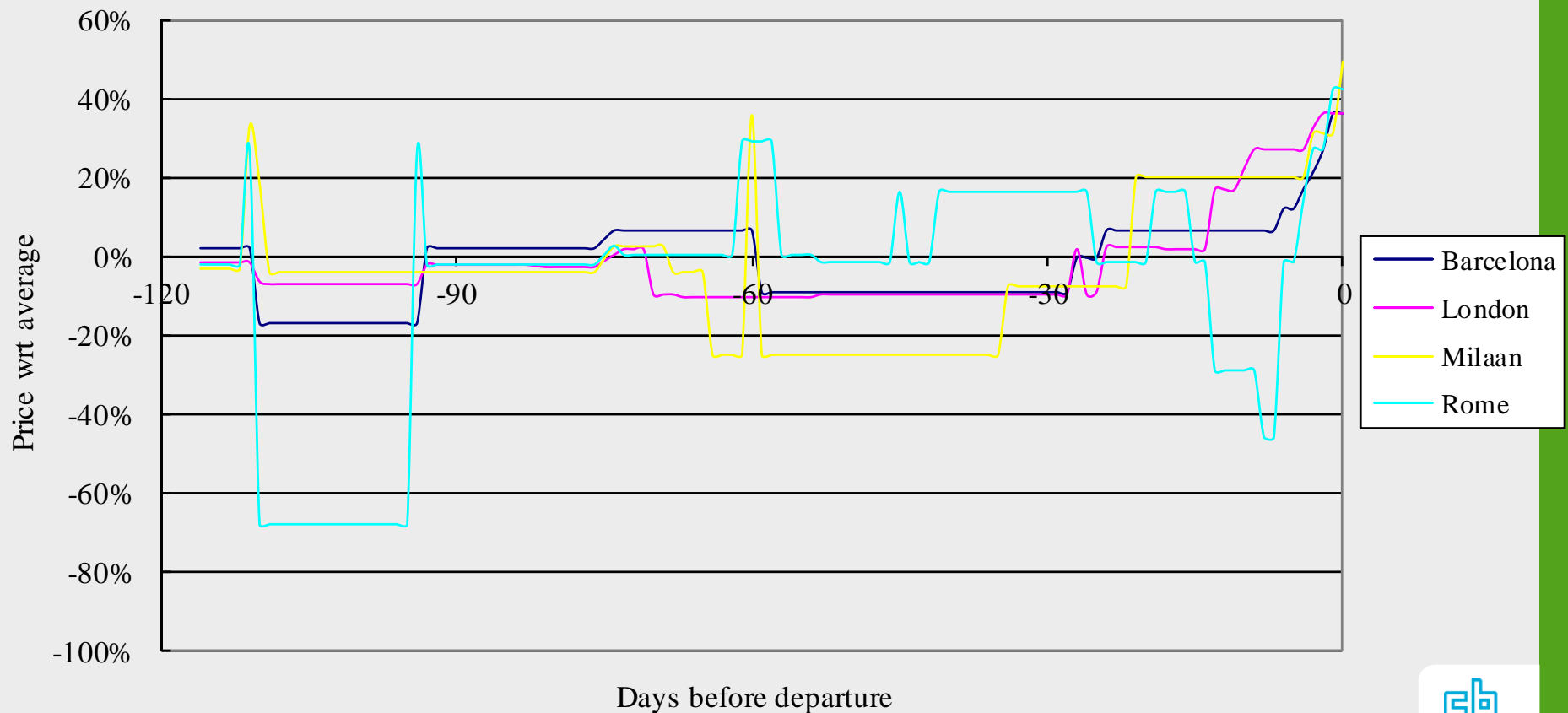
Core value internet data collection: good quality

The screenshot shows a navigation menu with 'WOMEN' highlighted. Below it, a breadcrumb trail reads 'Women / Clothing / Tops'. A filter section titled 'Select type of product' is open, showing 'Tops' as the selected category. To the right, there are checkboxes for 'Showing all brands' and other brand filters. Below the filters, two product images are shown: a 'Laybba Shirt' by Malene Birger for EUR 369.95 and a 'Cleopatra Tee' by Selected Femme for EUR 34.95. The 'WOMEN' tab in the navigation menu is circled in blue, and a blue line connects it to the 'Text classification' list on the right.

Text classification

- L-Coicop
 - Women
 - Men
 - Kids
- Clothing type
 - Jeans
 - Trousers
 - Top
 - Sweaters
 - ...
- Classification rule:
“sweater” and not “jack”

Core value internet data collection: better insight in data



Challenges & solutions: organization

Team

Robot team

- Research
- IT

CPI team

Process

- Web scraper
Development
and
Maintenance
- Data
Monitoring

Data
monitoring /
validation

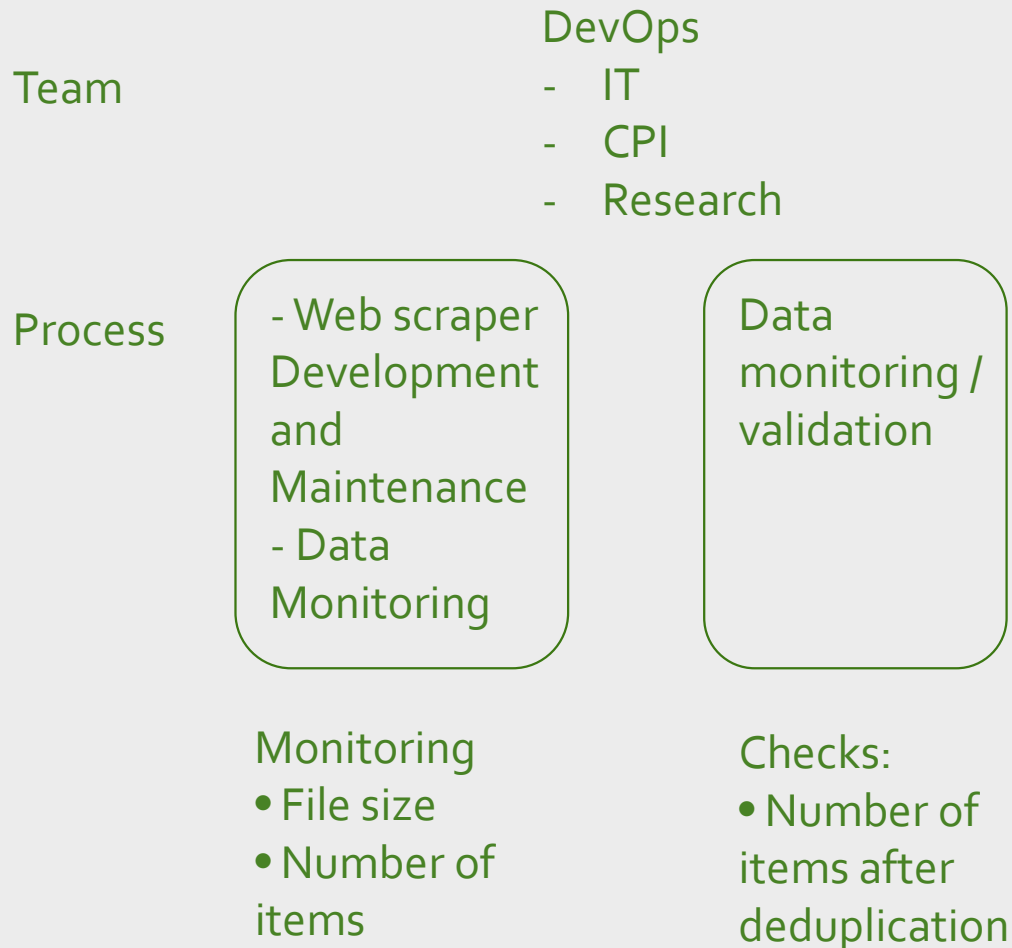
Monitoring

- File size
- Number of items

Checks:

- Number of items after deduplication

Challenges & solutions: organisation



Challenges and solution: new skills and knowledge

- IT-skills for web scraping and data-analyses
- Machine learning skills and knowledge to process big data
- New employees
- Acquisition of new knowledge and skills
 - External courses
 - Internal knowledge transfer



Challenges & solutions: other

- Mental: shared vision, a common urgency, enthusiasm and courage
- Legal issue



Conclusions & discussion

- Internet is (becoming) a natural part of our daily life →
Not using internet as a data source in official statistics might not only cause a bias in data sources but also in statistical publications.
- The business case for internet data will probably grow stronger over the years, as experience, knowledge and tools can also apply to other statistics.
- New opportunity's

