EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS 2024 ESTORIL - PORTUGAL

# Statistical scraping: informed plough begets finer crops

Olav ten Bosch, Statistics Netherland

Alexander Kowarik, Statistics Austria

Sónia Quaresma, Statistics Portugal

David Salgado, Statistics Spain

Arnout van Delden, Statistics Netherlands

# Contents

- The use of web data

- Statistical scraping: high level view

- Examples

- Definition, consequences

- Wrap-up

# 15 years of web data at Statistics Netherlands

2008-2010
Fuel prices
Real estate
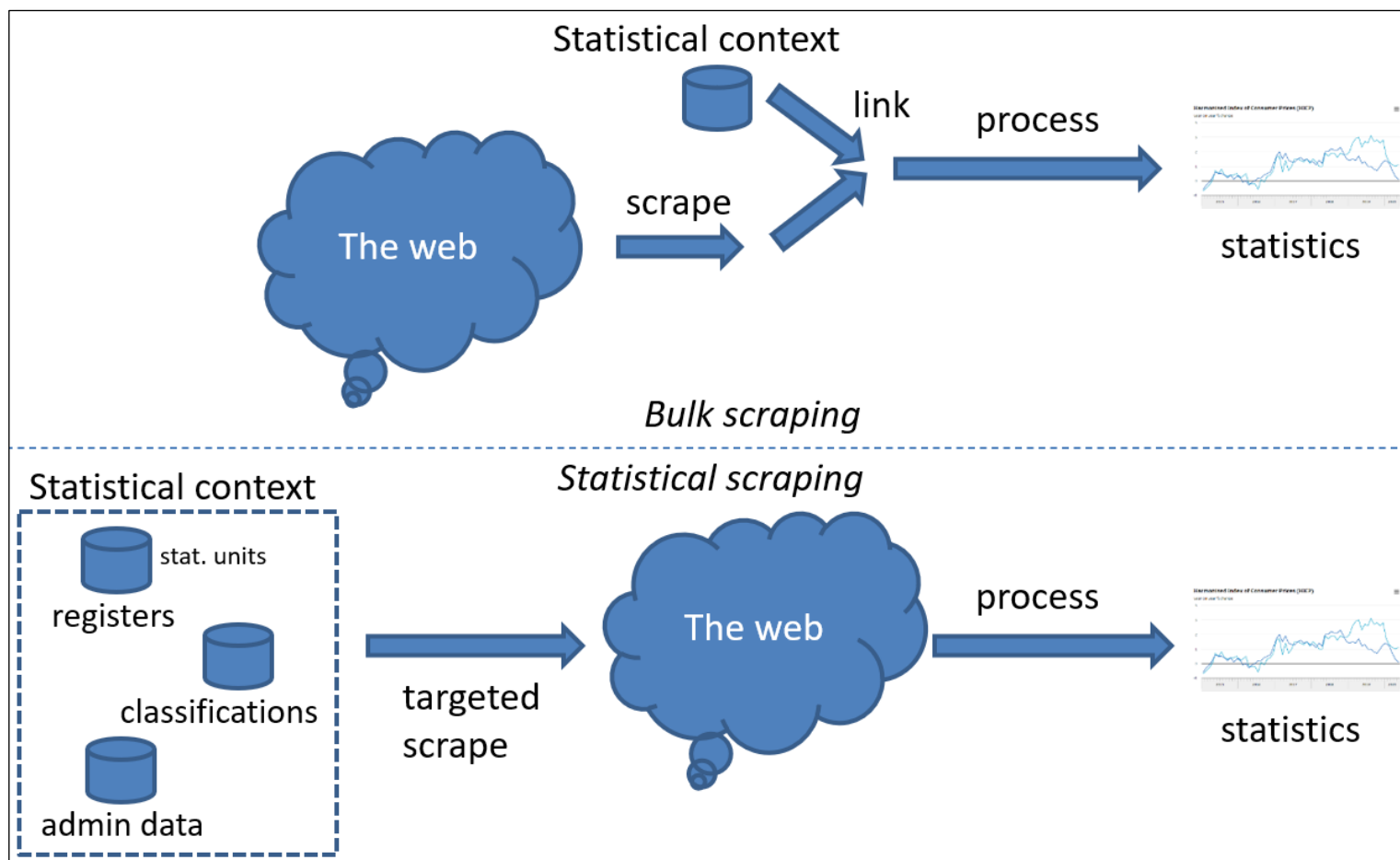Airlines

2011-2019
Experimenting
towards offstats

- *Webshops*: CPI (inflation): prices (clothing), books, travel, consumer electronics
- *Enterprise websites*: ecommerce activities, webshop detection, social media use, NACE, innovative companies, family businesses, jobs, drone companies
- *Annual reports*: financial and institutional data
- *Social media*: social tension indicator, (social) networks, community statistics
- *Property portals*: housing market dynamics
- *Job portals*: trends on job market (Textkernel), skills
- *Hotels / holiday homes portals*: tourism
- *Wikipedia*: community data, i.e. on international enterprises, network topology of train tracks, ..
- *DNS*: domain dynamics / relation with organisations
- *Municipality portals*: environmental permits
- *School portals*: courses offered; education
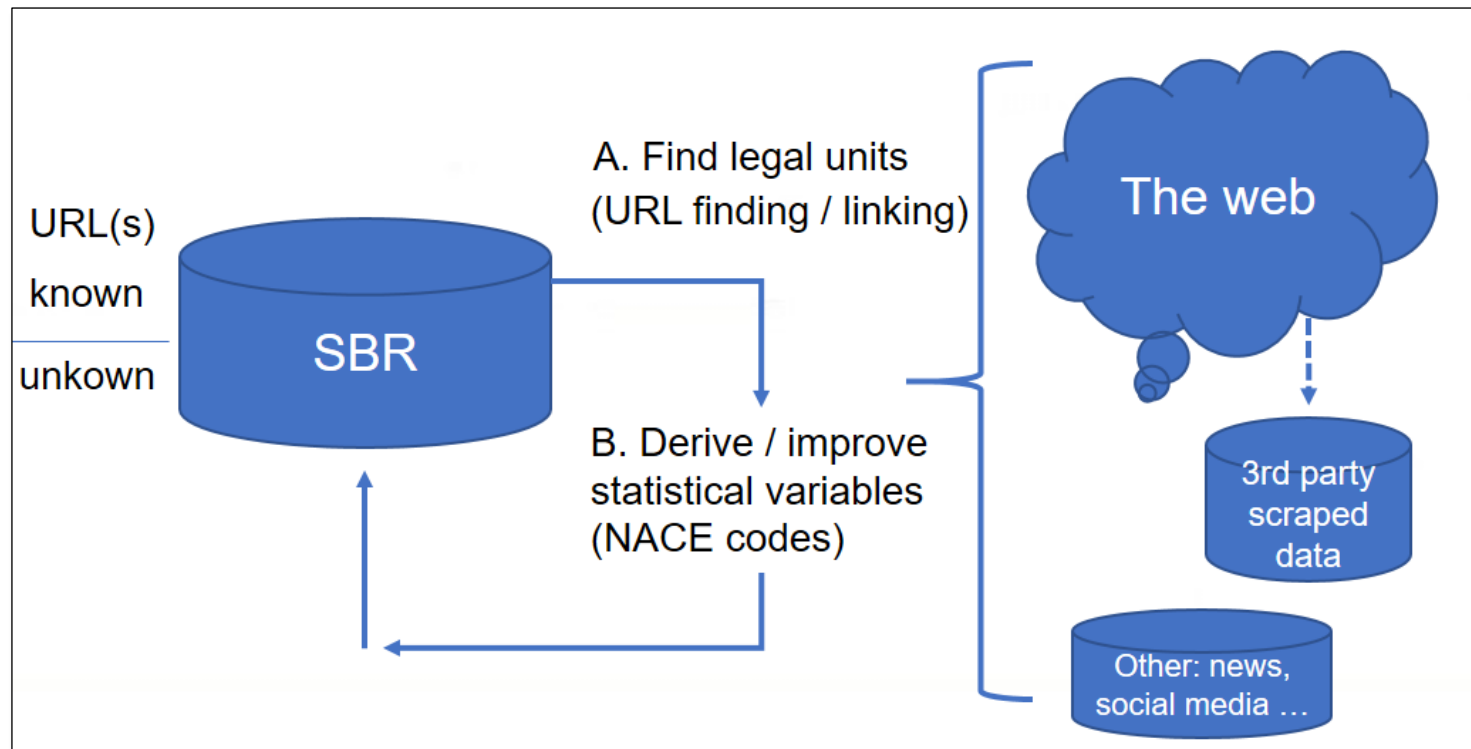- *Search engines*: enterprises, products

Manual retail price observations discontinued
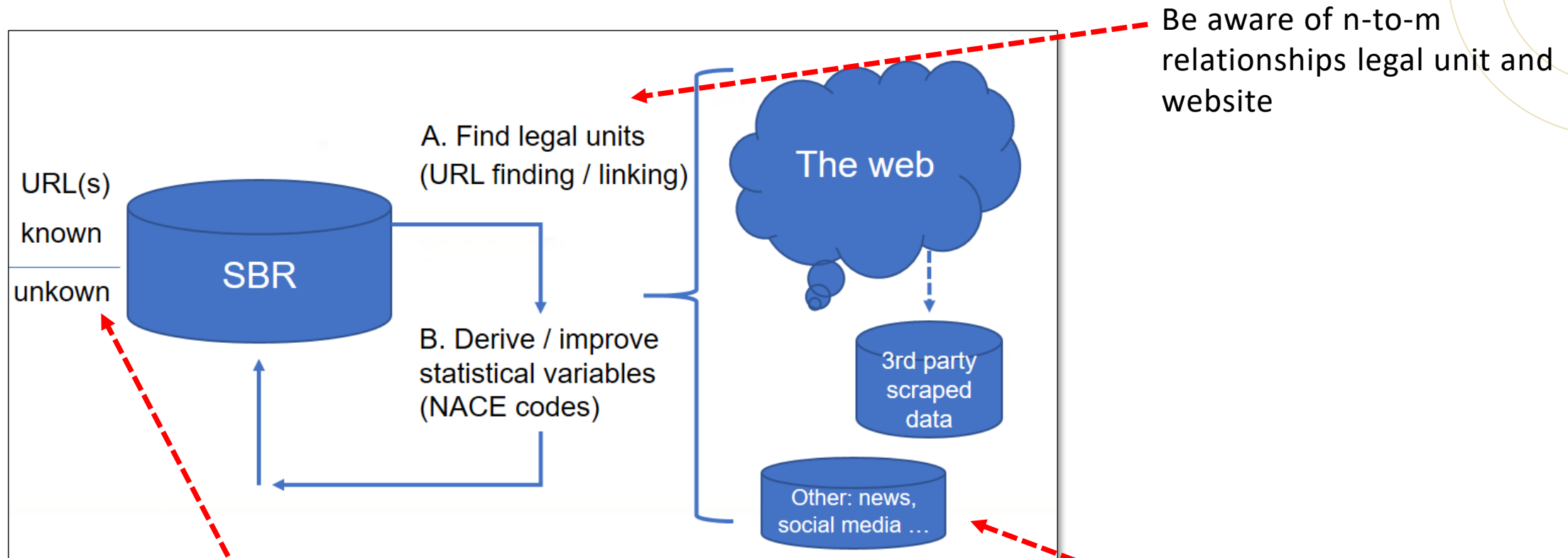
13/01/2020 14:00

2020

# High level view

# Example: Business register enhancements

# Example: Business register enhancements



Be aware of n-to-m relationships legal unit and website

A. Find legal units (URL finding / linking)

The web

URL(s) known unkown

SBR

B. Derive / improve statistical variables (NACE codes)

3rd party scraped data

Other: news, social media …

The ratio known/unknown is country- specific

Web data is more than websites only

# Example: Business register enhancements



using **search engine**(s): query on enterprise name / details

**focused** scraping: focus on interesting parts of web pages, like 'about us page' etc.

Both are target scrapes: It starts from what we already know in the business register
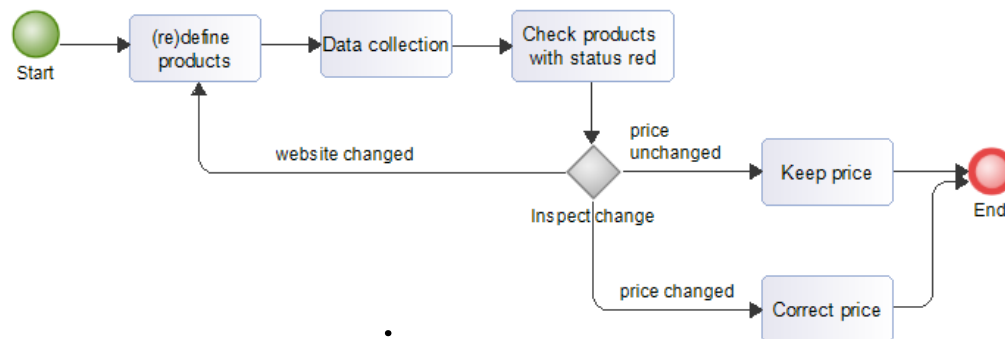
# Example: price statistics

- Early approaches used **bulk scraping** (Cavallo, 2009)

- Approaches evolved into a more **targeted** scrape of subset of products relating to COICOP categories (statistical context)

- NL: **Robot-assisted data collection**: basket-based price collection on to the web. For each product a number of measuring spots on the web are defined and regularly visited. In production for over 10 years.

- Even more advanced: use site-specific or general **search engines** to discover well-defined products (targeted scrape)

# Robot-assisted price collection

- ***Robottool***: (2012->ongoing) 8 users; 2850 price observations/month

- ***Check*** products with infrequent price changes ***easily***:
  - Examples: Cinema tickets, drivers lessons, car / bike repair, music instruments, farmacy, snackbars, dentists, sports, museum

- Price specialists define ***path*** to price and product to be checked

Green: nothing changed -> last price saved
Red: needs attention

Open source version:
https://github.com/SNStatComp/RobotTool

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is partly
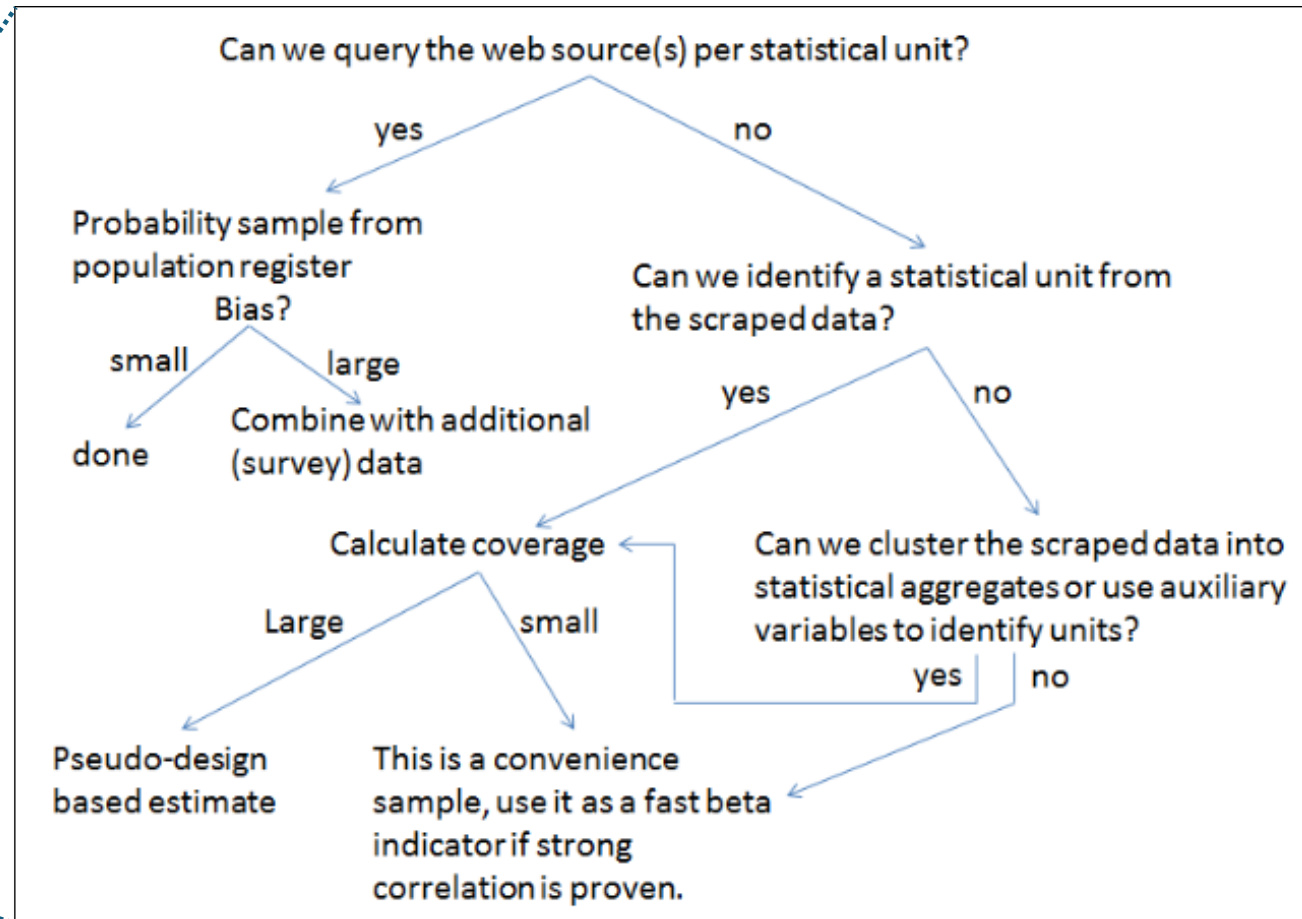financed by the European Union

# Other examples

- ***Tourism***: for each hotel star rating category a number of hotels are selected from the national ***star rating overview site*** (sample design). This sample is then scraped regularly at ***low frequency***.

- ***Education***: visit school web sites starting from a register of schools to draw a sample of teachers. Collecting data on a ***representative subset*** of doctorate holders.

- ***Labour market***: is statistical scraping applicable? Start from business register and visit enterprise websites and possibly job portals, keeping the relationship with the statistical unit?

EUROPEAN CONFERENCE ON
QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is partly
financed by the European Union

# Generic workflow for webdata



BigSurv paper, 2018

https://www.researchgate.net/publication/327385487_Web_scraping_meets_survey_design_combining_forces

# Definition

Def 1.1: *Statistical scraping is the use of online data starting from a-priori information in the respective statistical domain keeping a clear relation with the statistical context.*

# Consequences

Methodological:

- In general, statistical scraping **helps** cope with different types of **representation errors**

- If applied on unit level it becomes possible to calculated proven survey methodology **quality indicators**

Other:

- A targeted scrape leads to **smaller**, **more manageable** data streams

- Web queries may need **possibly sensitive** statistical **input data**, which should be handled with care

- Search (engine) **leakage**: needs attention bot **manageable**

EUROPEAN CONFERENCE ON QUALITY IN OFFICIAL STATISTICS
2024 ESTORIL - PORTUGAL

INSTITUTO NACIONAL DE ESTATÍSTICA
Statistics Portugal

eurostat

The conference is partly
financed by the European Union

# Wrap-up

- **_Online data_** have shown to be valuable for official statistics, either as primary or auxiliary data source

- In addition to bulk scraping a **_statistical scraping_** approach has been presented where a-priori information is used for a **_targeted scrape_** resulting in selected samples with high statistical value

- Examples can be found in in **_business register enhancement_** and **_price_**, **_tourism_** and **_education_** statistics. Applicability in other domains is to be explored.

- Statistical scraping has methodological advantages, to cope with representation errors and if applied on unit level proven **_quality indicators_** can be calculated

- Statistical scraping leads to **_smaller_** and more manageable data streams

- Attention should be paid to treatment of possible **_sensitive_** scraping input data

- The concept may **_complement_** or in some cases **_replace_** bulk scraping methods
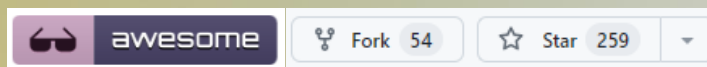
# Thanks

Questions, comments, ideas:

Olav ten Bosch

o.tenbosch@cbs.nl

And please like the awesome list of official statistics software:

awesomeofficialstatistics.org