


Webscraping in FinEco Research: Risks and Opportunities

Cynthia A. Huang

Econometrics & Business Statistics, Monash University

Introduction

About Me

-  PhD Candidate in Monash EBS, affiliated with
 - NUMBATS (EBS), SoDa Labs (ECON), EmVis (Faculty of IT)
 - Monash Data Futures Institute
- Researching principles and methods for "alternative" data
 - conceptual and practical data provenance tools for harmonised multi-source datasets
 - ***adapting web-scraped retail product & price data for public health research***
 - Statistical properties of alternative data, grammar of graphics
- Open and reproducible research tools, research software design
 - Quarto, git, replication packages
 - NUMBATS Hacky hour, community building

Objectives


Research opportunities & risks

- Big data, novel data, alternative data
- Operational vs. scientific risk
- Web technologies and web-scraping methods
- Ethical & legal risk

 How to code a web-scrapers

 How to code big-data analysis

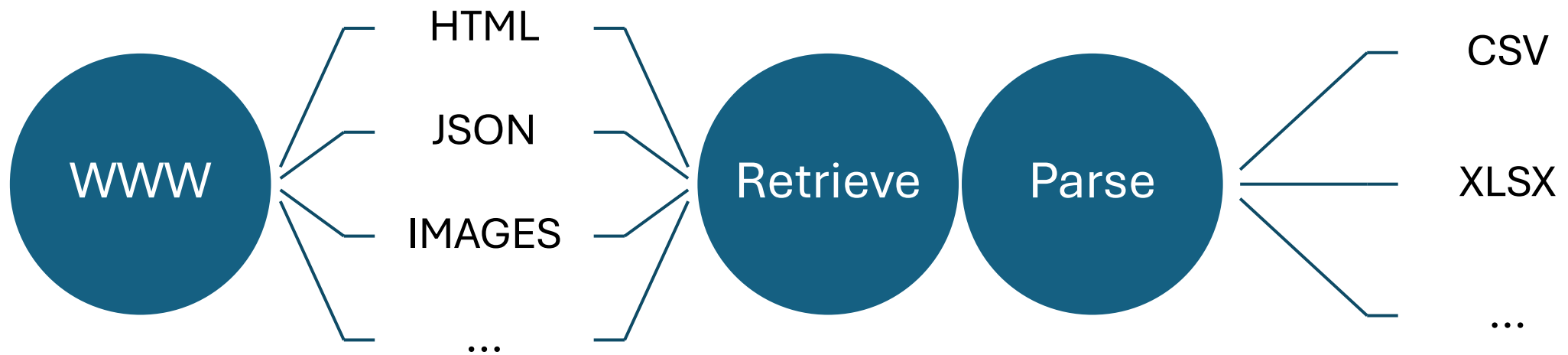
 Project level considerations

 Resourcing web-scraping projects (skill sets, RA hires etc.)

Web scraping—*verb*.

Cambridge dictionary:

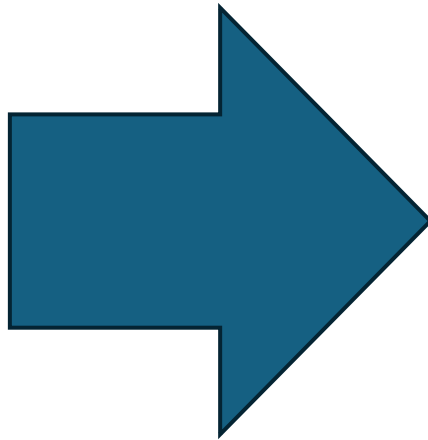
*the activity of **taking** information from a website or computer screen and **putting** it into an ordered document on a computer*



Web scraped data—*noun*.

Cambridge dictionary:

*the activity of taking **information from a website or computer screen** and putting it into an ordered document on a computer*



Common properties of web scraped data



- Some degree of messiness
 - Collection, preparation
 - Analysis, statistical properties
- Non-trivial parsing from raw web extraction to analysis-ready data
- Murky data dimensions/quality
 - How many records?
 - How many observational units?
 - Complete/missing observations?

Research Opportunities & Risks

Big, novel, alternative (data) opportunities

- The internet is a rapidly expanding universe of **(collectable*) data**
 - Multiple modalities: text, images, video...
 - Multi-source harmonization: shopping aggregators, knowledge graphs
 - Novel signal sources: social media sentiment, digital economy indicators
- **Lower* cost** than traditionally available data collection methods
 - Digital vs. Analogue access and delivery
 - Automated data extraction vs. Manual data entry
- Volume, Velocity, Variety... **Veracity?**
 - Noise vs. signal for data mining, prediction, casual inference
 - Collecting data from a Complex System vs. Library

Research Risk: Traditional Data

Operational Risk

- Availability and Access
- Cost
- Legal and ethical

Scientific Risk

- Suitable analysis methods
- Signal-to-noise ratio (in context)
- Data quality

Research Risk: Web Scraped Data

Operational Risk

- **Cost = Time**
- Availability and access
 - Sampling bias
 - Unstable web interfaces
- Legal and ethical
 - Usage terms

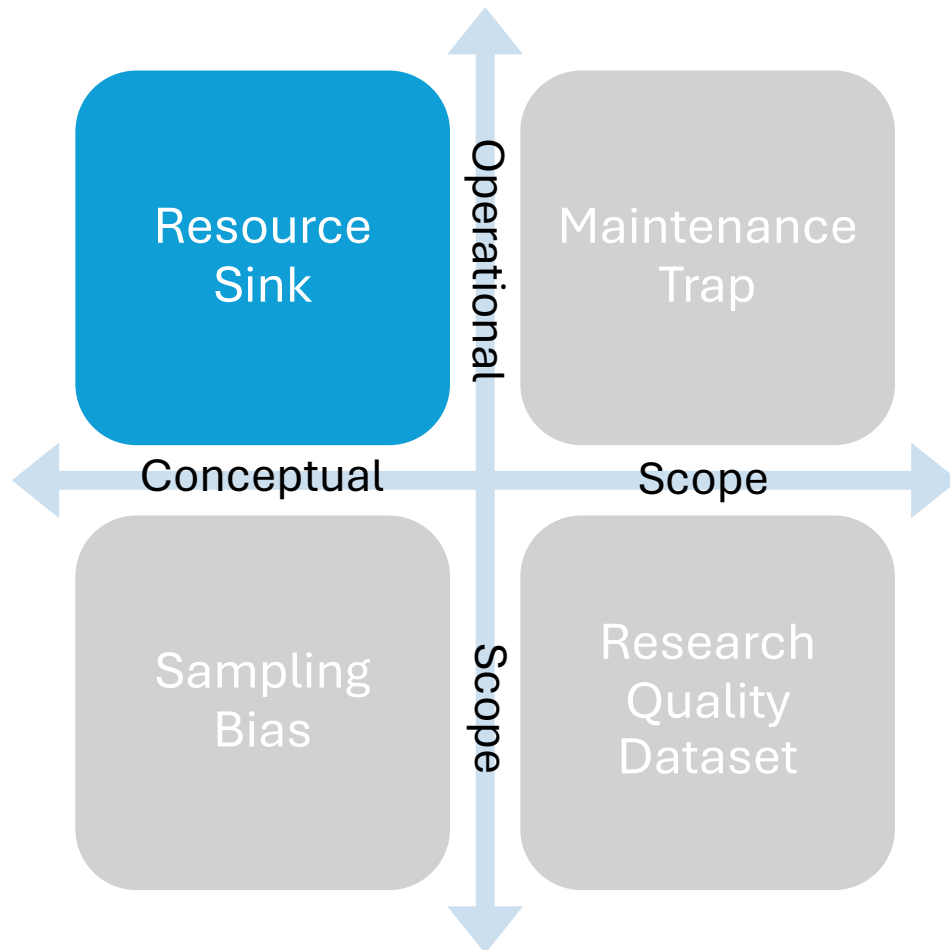
Scientific Risk

- **Data quality**
 - Completeness
 - Accuracy
- **Signal-to-noise ratio**
 - Big Data Paradox (Bradley et al., 2021)
- Suitable analysis methods

Data Collection Scope

Conceptual Scope

Web Scraped Data Risk



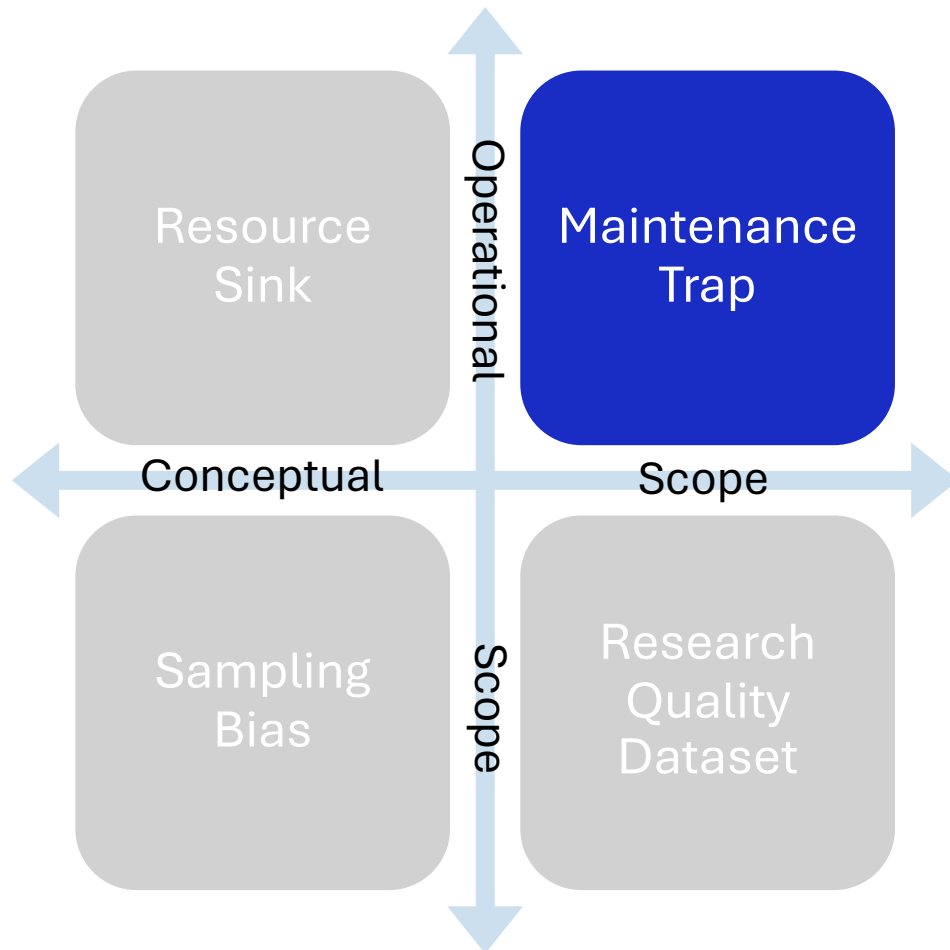
Resource Sink

Example dataset description:

Daily stock prices

- Conceptually and operationally open-ended
- Unclear dataset dimensions
 - Which firms?
 - What time period?

Web Scraped Data Risk



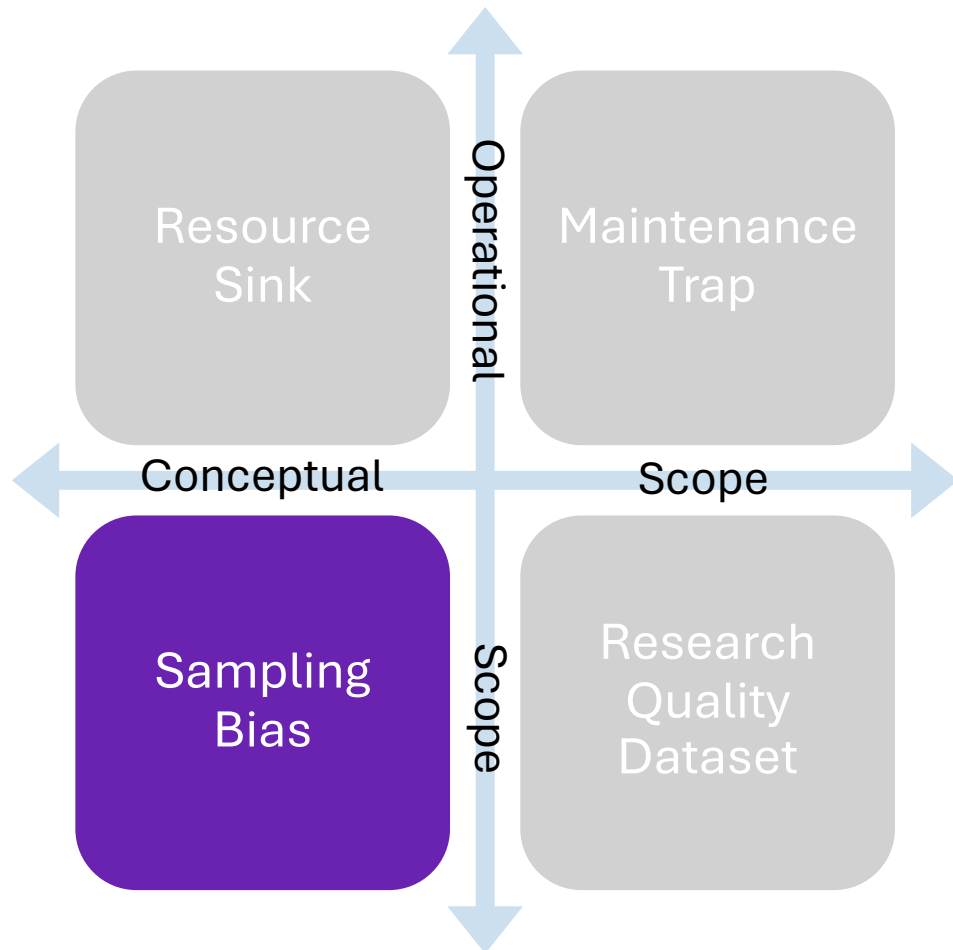
Maintenance Trap

Example dataset description:

Daily prices for alcoholic beverages sold online by retailers with physical stores in Australia

- Conceptually limited, but operationally open-ended
- implies **ongoing** collection:
 - Repair “broken” scrapers
 - Upgrade storage & compute
 - (Cavallo & Rigobon, 2016)

Web Scraped Data Risk



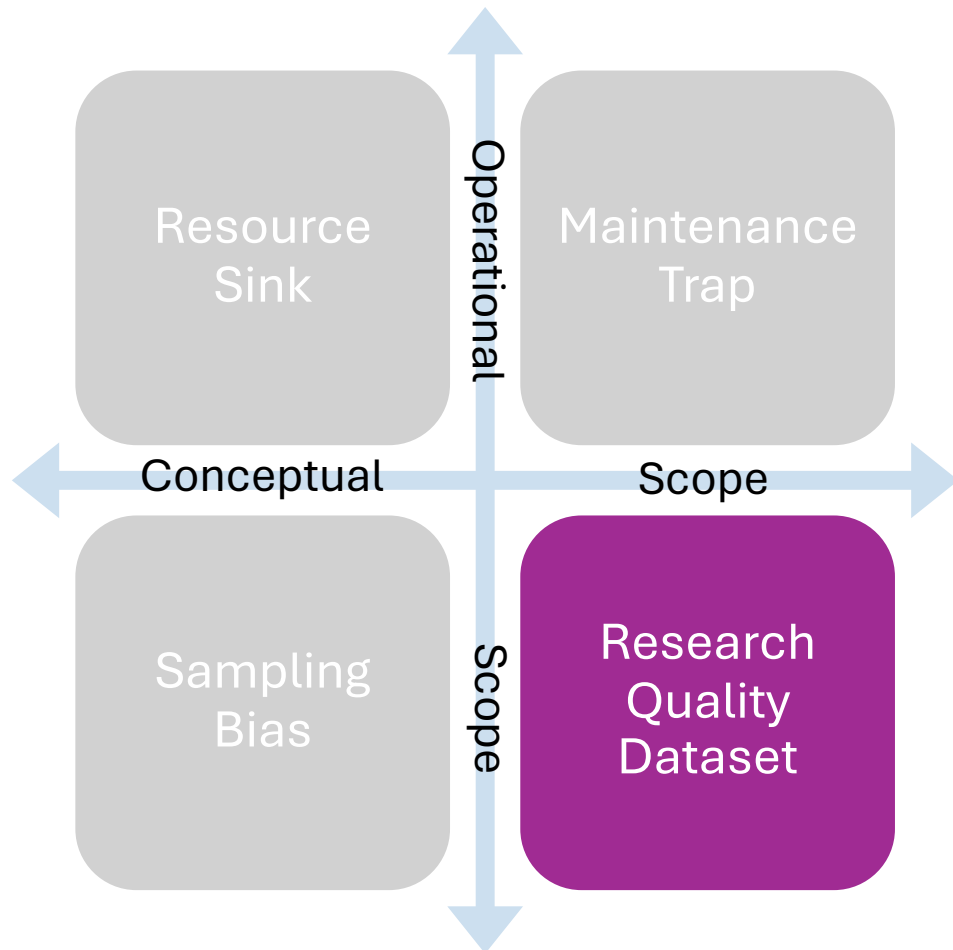
Sampling bias

Example dataset description:

Daily prices for alcoholic beverages available online between Jan-Dec 2023

- Operationally limited, but conceptually open-ended
- What population is being sampled from?
 - Index coverage (Foerderer, 2023)

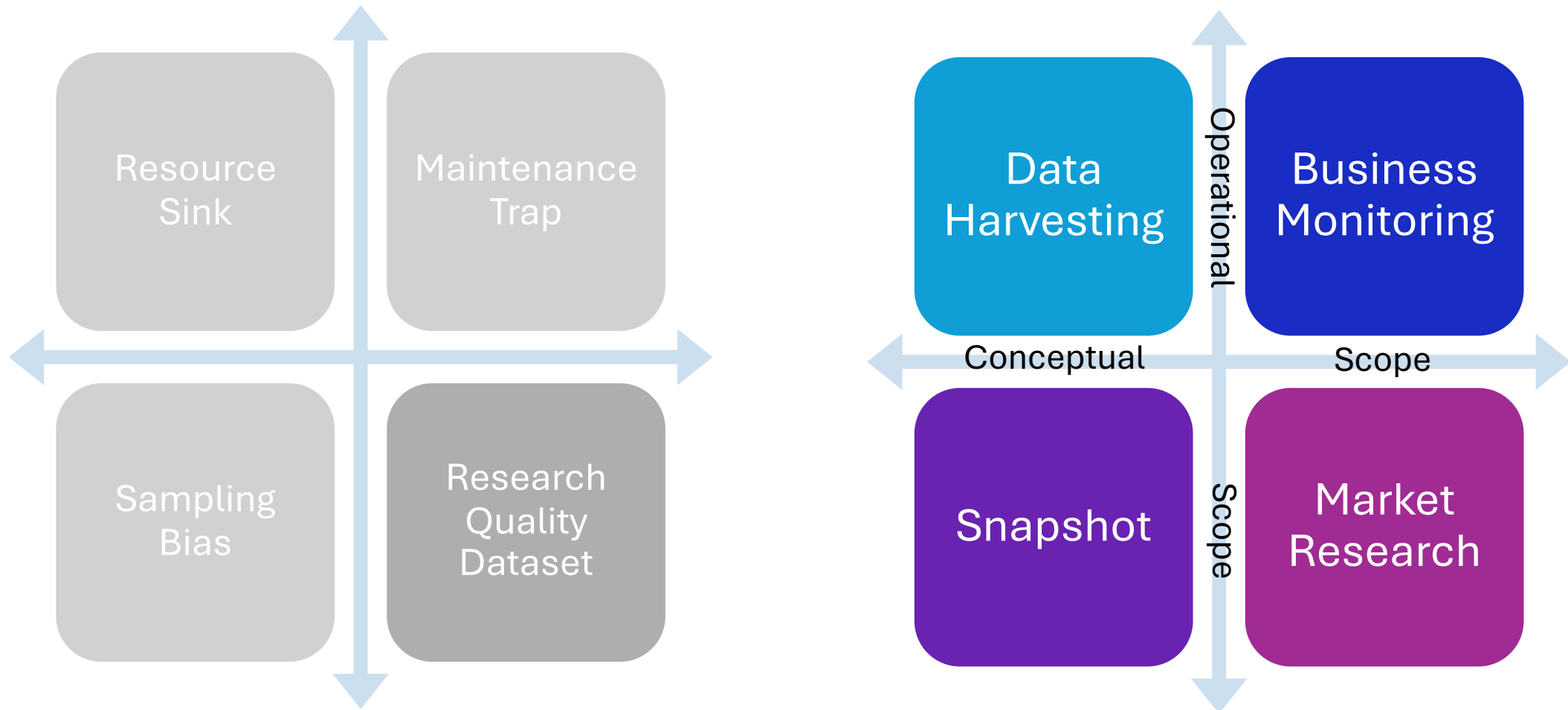
Web Scraped Data Risk



Research Quality Dataset(s)

- Conceptually and operationally limited
- Quality is **context** dependent!
- Requires deliberate design decisions
 - Filtering criteria (e.g. top products, largest retailers)
 - Conversion between observation and analysis resolution (e.g. from retailer to brand)

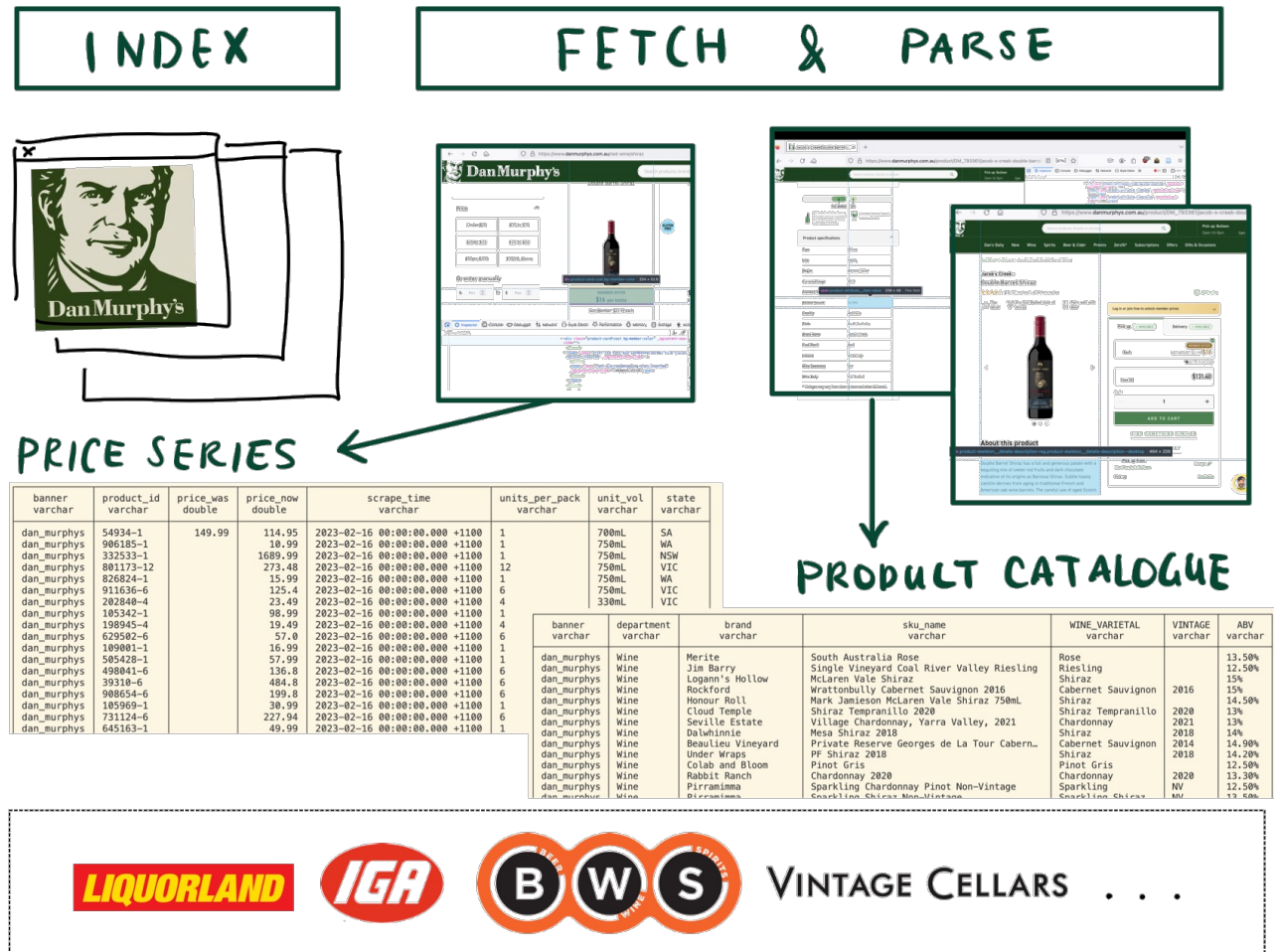
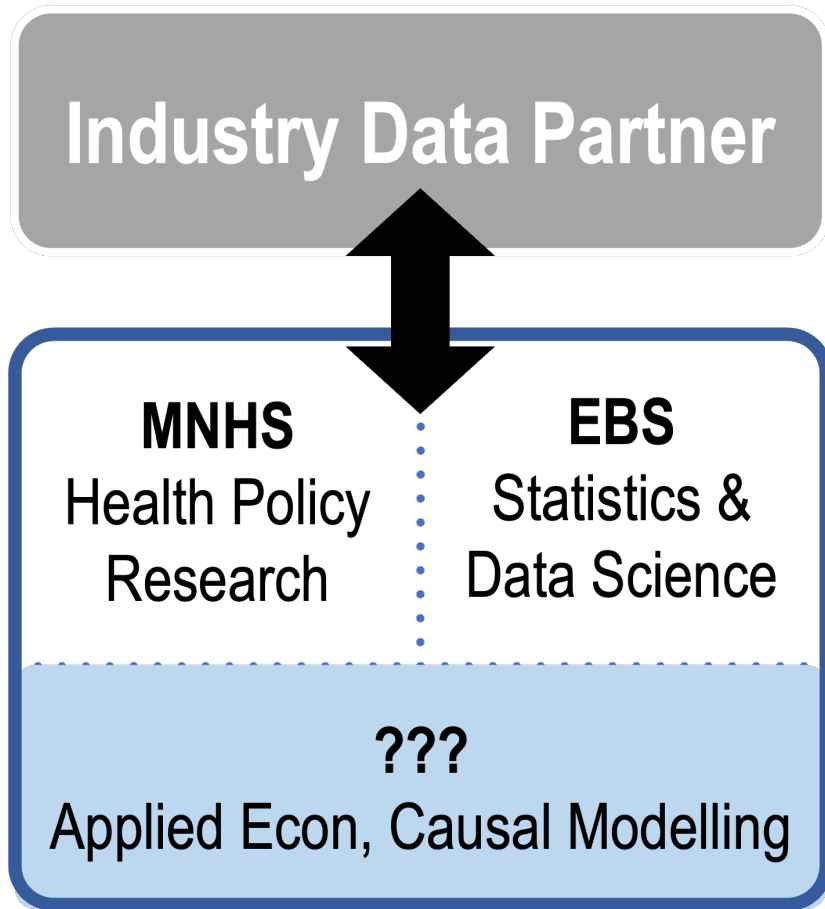
Business vs. Research Needs



Case Study

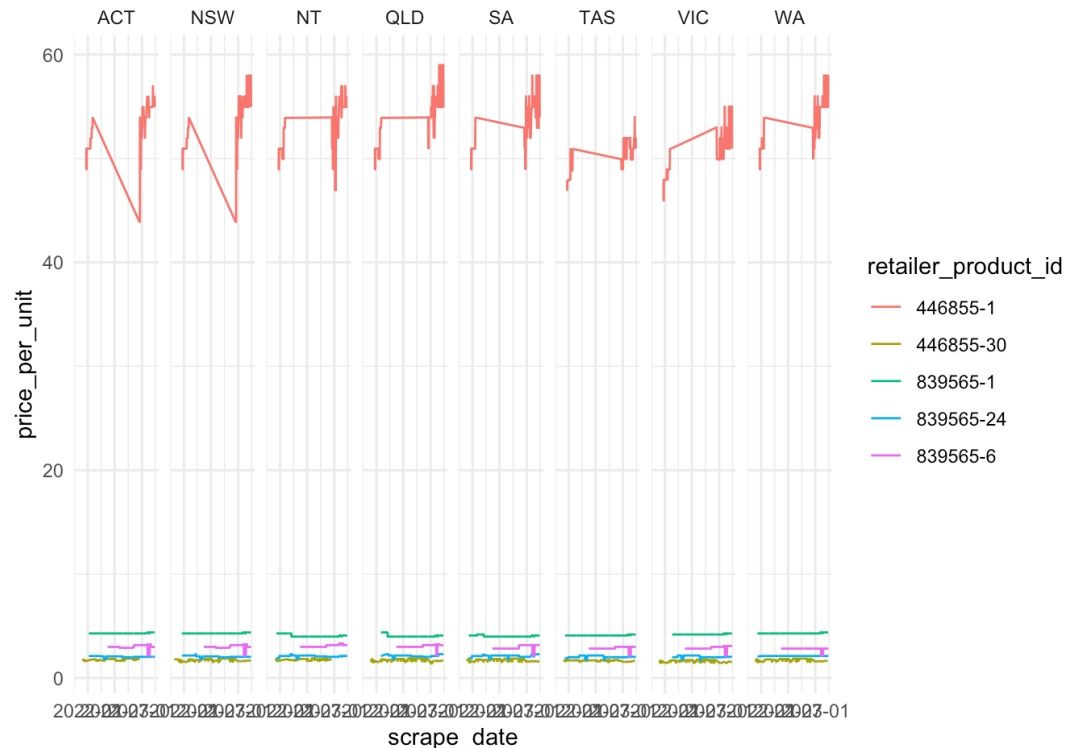
Adapting web scraped retail price data for public health research

Outsourcing web scraping



Data refinements & augmentation

Anomalies in cross-retailer harmonisation



Augmenting product attributes

- Product name
- Brand name
- Manufacturer
- Product category
- Volume per unit
- Alcohol by Volume
- Standard drinks per unit

Web technologies, access interfaces, and collection methods

Web scraping: Expectations vs. Reality



Web scraping: Expectations vs. Reality



Live Demo

- HTML elements
- API responses

The screenshot shows the Yahoo Finance website with the Chrome DevTools Inspector open. The Inspector displays the HTML structure of a stock quote for DJT. The selected element is an anchor tag with the following attributes: `class="Fw(b) Ell D(b) C($linkColor) Pos(r) Z(2)" href="/quote/DJT" title="Trump Media & Technology Group Corp."`. The 'Rules' panel shows a CSS rule for visibility: hidden.

Symbol	Price	Change	% Change
DJT	48.66	-13.30	-21.47%
UNH	489.70	-5.00	-1.01%
CXAI	6.15	+3.70	+151.02%
GOOG	156.50	+4.24	+2.78%
MU	124.30	+6.41	+5.44%

Web technologies

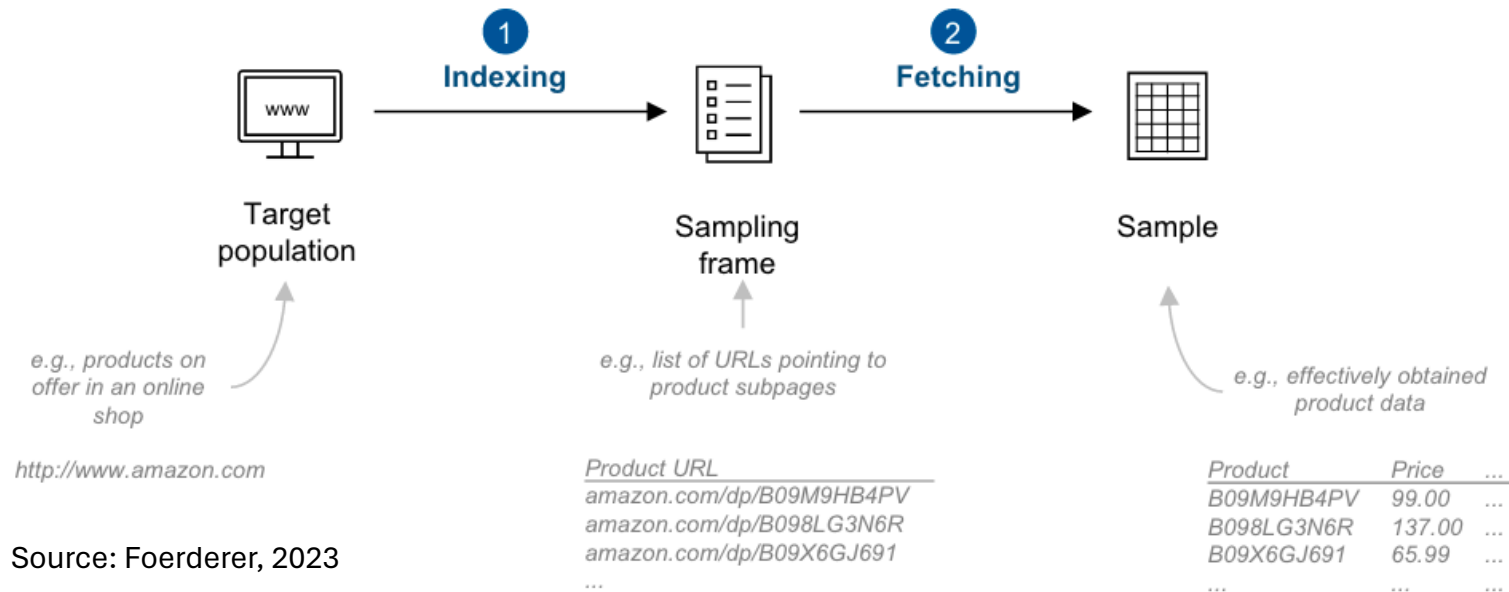
- Webpages and websites
 - Styling: CSS
 - Static content: HTML/XML
 - Dynamic content: Javascript, API queries
- Web Protocols (HTTP/HTTPS)
- Application Programming Interface (APIs)
 - Query parameters
 - Response formats: JSON, CSV...



Web scraping methods: Indexing & Crawling

Figure 1: Web Scraping Process: From the Target Population to the Sample.

Web scraping entails two steps, *indexing* and *fetching*. In indexing, the target population is systematically registered. Indexing yields the frame in terms of a register of all units in the population, together with the URLs pointing to each unit. Fetching automatically visits each URL listed in the frame and downloads the resource at which it points, typically an HTML document.



Web scraping methods: Retrieval

Manual Copy-Paste

- Concurrent* collection & parsing
- Human errors and judgement

Interactive collection

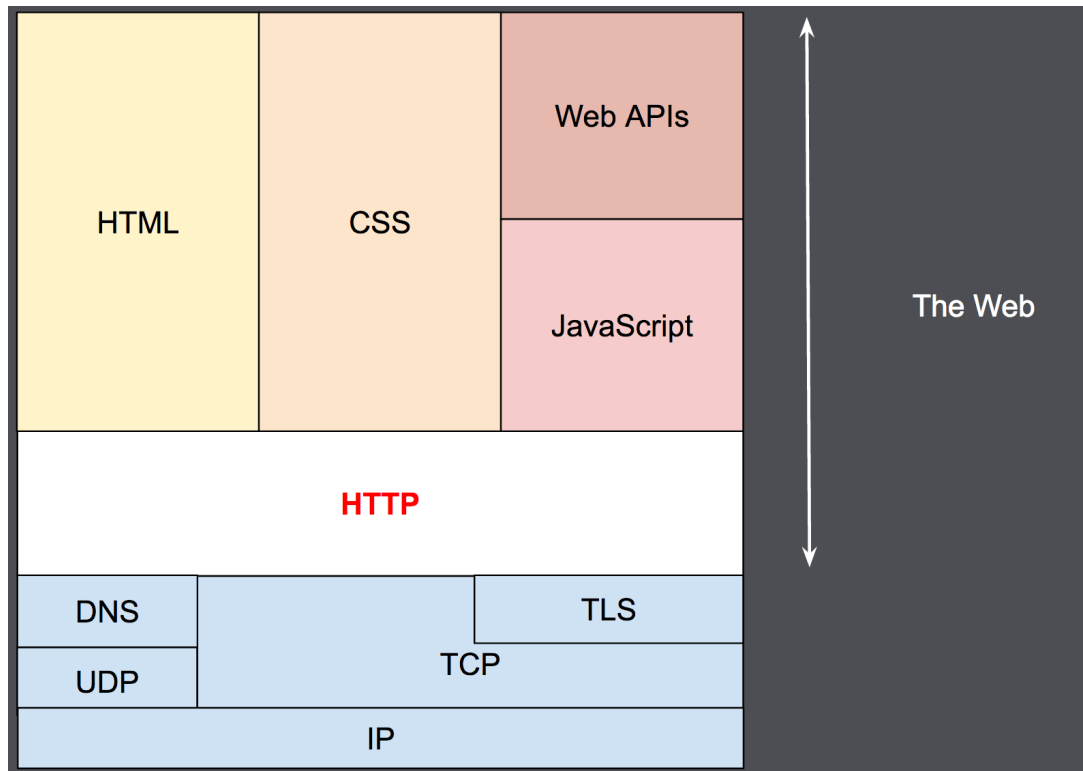
- Automating mouse/keyboard inputs
- Via headless browsers
- Selenium, chromium etc.

Request-based extraction

- Request resources using machine readable methods
- HTTP CONNECT/GET
- API Queries



Web scraping methods: Parsing



An overview of HTTP - HTTP | MDN. (2023, December 16).

<https://developer.mozilla.org/en-US/docs/Web/HTTP/Overview>

- HTML element extraction:
 - Tables, if you're lucky
 - Often requires detangling style information from metadata and data
- API response parsing:
 - A way for computers to talk to computers without human readability/presentation layer
 - Deciphering JSON keys can be tricky without access to dictionary

Collection Risk by Method / Format

Method	Access Stability	Access Difficulty	Parsing Complexity	Common raw formats
Manual	M	L	L-H???	...
Interactive sessions	L	H	M-H	HTML elements
HTTP page requests	M	L	M-H	HTML elements
API queries	H	L-M	L	JSON, CSV...

Live Demo

- HTML elements
- API responses

The screenshot shows a web browser at the URL `https://www.danmurphys.com.au/red-wine/pinot-noir`. The page displays two wine products: "Devil's Corner Pinot Noir" for \$19.99 each and "Jacob's Creek Classic Pinot Noir" for \$25.99 each. The Chrome DevTools Network tab is open, showing a list of 317 requests. The selected request is a POST to `api.da... Browse` with a response size of 41.21 kB. The response is a JSON object containing product information.

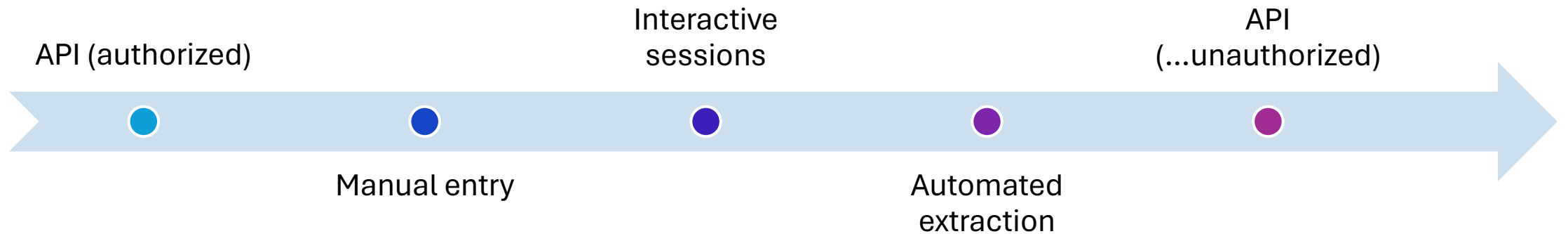
Sta...	Method	Domain	File	Initiator	T...	Transferred	Si...
200	OPTIO...	api.ed...	SddPrompt	fetch	plain	739 B	0 B
200	GET	api.da...	Keywords?path=/red-wine/pinot-noi	polyfills.67...	json	1.18 kB	8...
200	GET	api.da...	Keywords?path=/red-wine/pinot-noi	polyfills.67...	json	1.18 kB	8...
200	POST	api.da...	SponsoredAds	polyfills.67...	json	8.95 kB	3...
200	GET	api.da...	Banners?Location=/red-wine/pinot-i	polyfills.67...	json	946 B	3...
200	POST	api.da...	Browse	polyfills.67...	json	41.21 kB	2...
200	GET	api.da...	Banners?Location=/&Channel=Web	polyfills.67...	json	946 B	3...
200	GET	aem.d...	pinot-noir.model.json	polyfills.67...	json	4.80 kB	11...
200	GET	api.da...	221057,322002	polyfills.67...	json	4.44 kB	1...
200	POST	api.da...	new	polyfills.67...	json	15.77 kB	8...
200	POST	api.da...	member offers	polyfills.67...	json	31.64 kB	1...
200	GET	api.ed...	murphy	/:1 (fetch)	json	1.05 kB	2...
200	GET	api.ed...	SddPrompt	/:1 (fetch)	json	946 B	6...
200	POST	api.da...	SeoMetatags	polyfills.67...	json	1.55 kB	9...
200	GET	aem.d...	master.model.json	polyfills.67...	json	3.69 kB	1...
200	GET	auth.d...	openid-configuration	polyfills.67...	json	2.46 kB	1...
200	GET	api.da...	Banners?Location=/&Channel=Web	polyfills.67...	json	946 B	3...
200	GET	api.da...	meganav	polyfills.67...	json	52.30 kB	2...
200	GET	api.da...	Categories?Level=0&Location=meg	polyfills.67...	json	233.86 kB	1...
200	GET	api.da...	Trolley?summary=true&IncludeOrde	polyfills.67...	json	1.42 kB	1...
200	GET	api.da...	Preferences?IsCheckoutV2=true	polyfills.67...	json	1.94 kB	2...
200	GET	aem.d...	home.model.json	home:184 ...	json	15.86 kB	1...
200	GET	www.w...	15-h8-968898-62hh54424 ie	runtime:00	ie	11.49 kB	4...

The response for the selected request is a JSON object with the following structure:

```
{  "InspirationCards": Object { "Cards": {...} },  "SeoMetaTags": Object { "Title": "Buy Wine, Beer, Spirits online at Dan Murphy's | Alcohol Delivery & Bottle Shop", "MetaDescription": "Now with contactless delivery, shop online to get drinks delivered to your door or pick up in-store in 30 minutes. Lowest Liquor Price Guarantee. Biggest Range.", "NoIndex": false, ... },  "Bundles": [ {...}, {...}, {...}, {...}, {...}, {...}, {...}, {...}, {...}, {...}, ... ],  "0": Object { "Name": "Little Giant Pinot Noir", "PackDefaultStockCode": "140587", "PackParentStockCode": "140587", ... },  "Products": [ {...} ],  "IsInDefaultList": false,  "IsPersonalised": false,  "1": Object { "Name": "Yarra View Yarra Valley Pinot<br>Noir", "PackDefaultStockCode": "79190", "PackParentStockCode": "79190", ... },  "2": Object { "Name": "Devil's Corner Pinot Noir", "PackDefaultStockCode": "902552", "PackParentStockCode": "902552", ... },  "3": Object { "Name": "Jacob's Creek Classic Pinot<br>Noir", "PackDefaultStockCode": "328206", "PackParentStockCode": "328206", ... },  "4": Object { "Name": "Mud House Claim 431 Vineyard<br>Pinot Noir", "PackDefaultStockCode": "917721", "PackParentStockCode": "917721", ... },  "5": Object { "Name": "Cold Snap Cool Climate Pinot<br>Noir", "PackDefaultStockCode": "171492", "PackParentStockCode": "171492", ... },  "6": Object { "Name": "Elephant In The Room Pinot<br>Noir", "PackDefaultStockCode": "608178", "PackParentStockCode": "608178", ... },  "7": Object { "Name": "Oyster Bay Pinot Noir", "PackDefaultStockCode": "904243", ... }
```

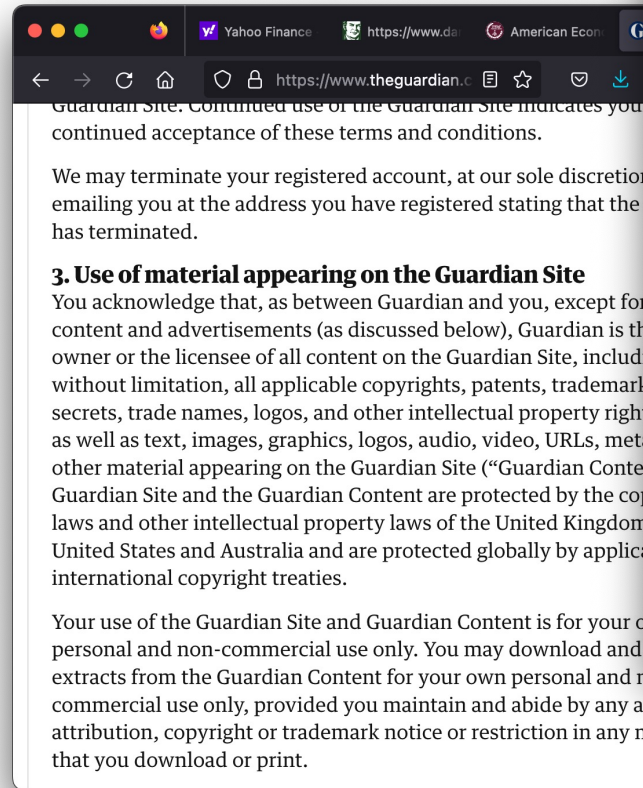
Final considerations

Ethical & Legal Risk: Retrieval



Ethical & Legal Risk: Data Usage

- It depends, check...
 - Website T&Cs
 - Journal policies
 - robots.txt
 - ...



Starting a Web Scraped Data Project?

Assess Feasibility & Quality:

- Inspect webpages
- Assess consistency
 - How many different webpage architectures?
- Consider data collection vs. analysis unit/resolution
 - Define "research quality" in your context
- Consider website and data owner motives
 - Are there any incentives to obfuscate data?

Resource appropriately:

- For in-house hires, consider experience & knowledge of:
 - Retrieval methods,
 - Scraping etiquette,
 - Parsing methods,
 - Data wrangling
- Consider industry partners if available
- Paid services can be suitable for limited data collection

Appendix

Resources

Web-scraping in R

- Extended tutorial by Hadley Wickham:
 - <https://github.com/hadley/web-scraping>
- Scraping tools:
 - <https://rvest.tidyverse.org>
- Polite sessions:
 - <https://github.com/dmi3kno/polite>

Case Study: Resource Sink

- “Although gathering this massive amount of prices was cheaper online than with traditional methods, it required funding that could not be sustained through grants. Thus, in 2011 we started a company called **PriceStats** that now collects the data and produces high-frequency indexes for central banks and financial-sector customers.” (Cavallo and Rigobon, 2016, p. 153)

AEA Data Legality Policy

- “A particular concern expressed by many researchers is the treatment of computer-assisted acquisition of data (“scraped data”) when such acquisition contravenes terms of use of the data owner. While scraping may contravene terms of use, it may not be illegal. Its legality is not settled under current US law (as of January 2023). Editors will treat papers with scraped data as legally acquired as long as this issue is unsettled when the paper was submitted. Should it in the future become settled law that scraped data is illegal, the AEA will communicate how scraped data will be treated under the Policy.”

<https://www.aeaweb.org/journals/data/data-legality-policy>

References

- **Slide 10:** Sagiroglu, S., & Sinanc, D. (2013). Big data: A review. *2013 International Conference on Collaboration Technologies and Systems (CTS)*, 42–47. <https://doi.org/10.1109/CTS.2013.6567202>
- **Slide 12:** Bradley, V. C., Kuriwaki, S., Isakov, M., Sejdinovic, D., Meng, X.-L., & Flaxman, S. (2021). Unrepresentative big surveys significantly overestimated US vaccine uptake. *Nature*, 600(7890), 695–700. <https://doi.org/10.1038/s41586-021-04198-4>
- **Slide 14:** Cavallo, A., & Rigobon, R. (2016). The Billion Prices Project: Using Online Prices for Measurement and Research. *Journal of Economic Perspectives*, 30(2), 151–178. <https://doi.org/10.1257/jep.30.2.151>
- **Slide 15,26:** Foerderer, J. (2023). Should we trust web-scraped data? (arXiv:2308.02231). arXiv. <http://arxiv.org/abs/2308.02231>
- **Slide 24:** Nigam, H., & Biswas, P. (2021). Web Scraping: From Tools to Related Legislation and Implementation Using Python. In J. S. Raj, A. M. Ilyasu, R. Bestak, & Z. A. Baig (Eds.), *Innovative Data Communication Technologies and Application* (Vol. 59, pp. 149–164). Springer Singapore. https://doi.org/10.1007/978-981-15-9651-3_13