

Webscraping in FinEco Research: Risks and Opportunities

Cynthia Huang

Econometrics & Business Statistics

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-scraper



How to code big-data analysis



Project level considerations

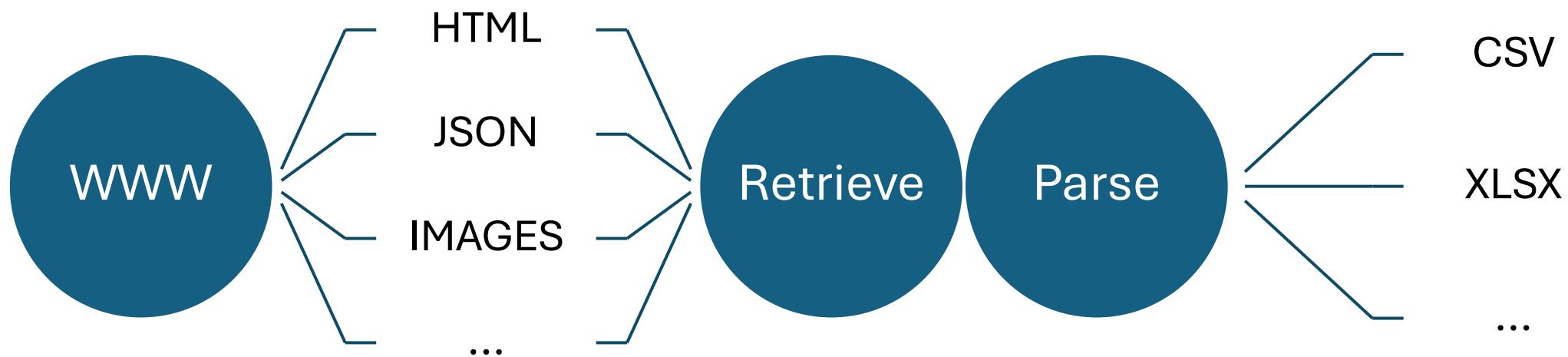


Resourcing web-scraping projects (skill sets, RA hires etc.)

Web scraping—verb.

Cambridge dictionary:

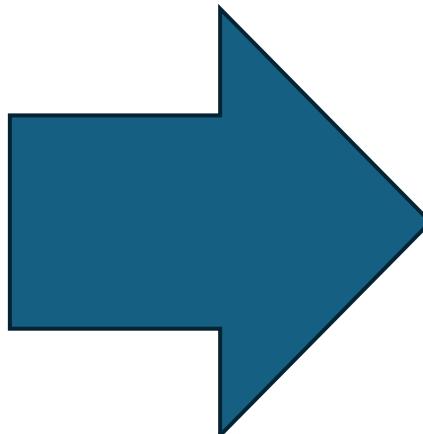
*the activity of **taking** information from a website or computer screen and **putting** it into an ordered document on a computer*



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, causal 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 Scrapped Data

Operational Risk

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

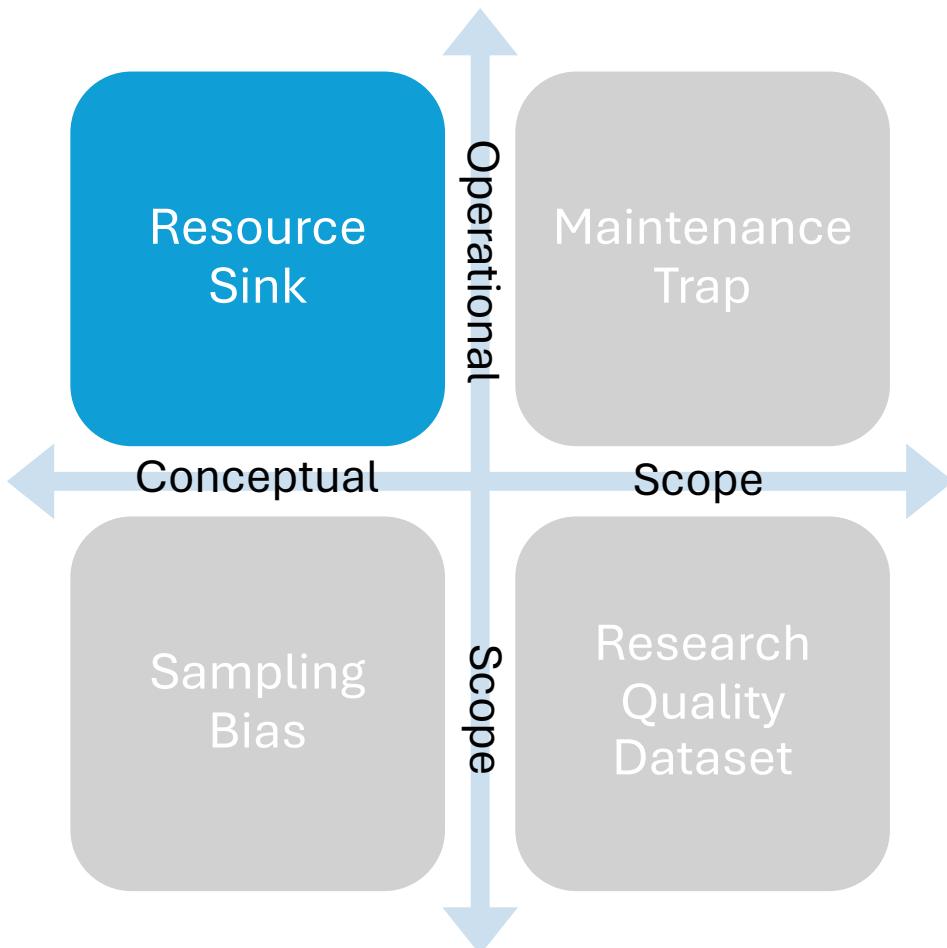
Scientific Risk

- **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



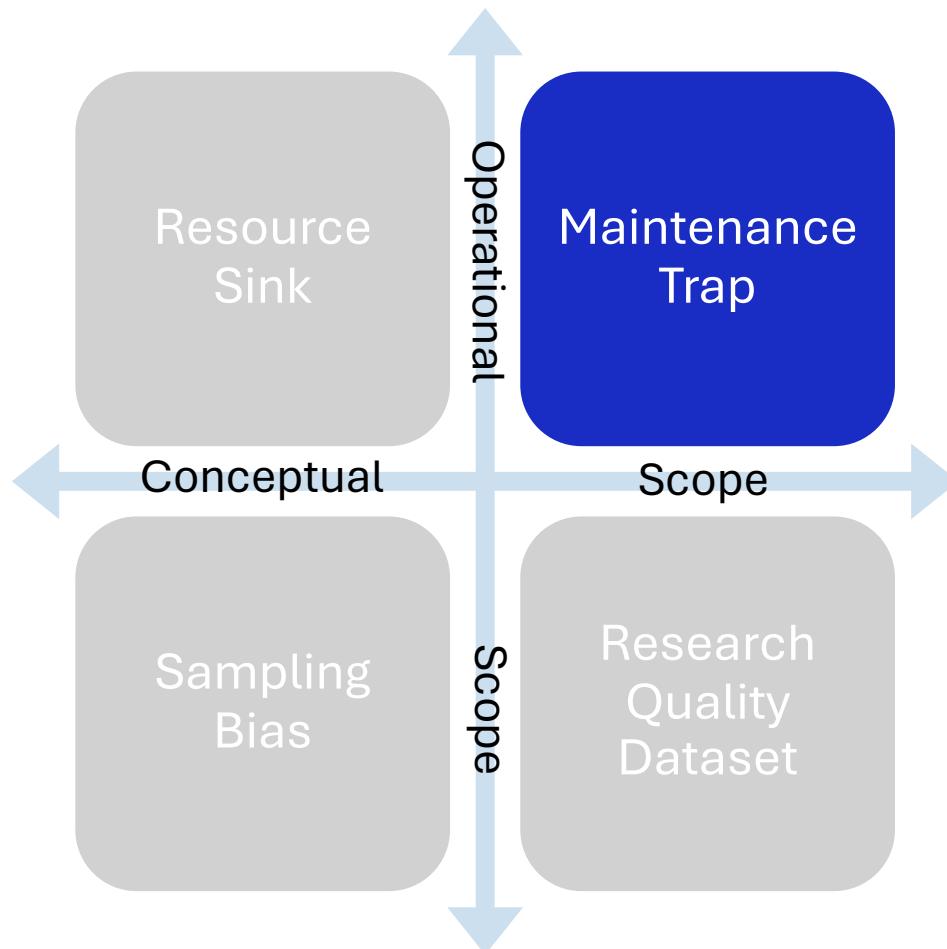
Resource Sink

Example dataset description:

Daily stock prices

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

Web Scraped Data Risk



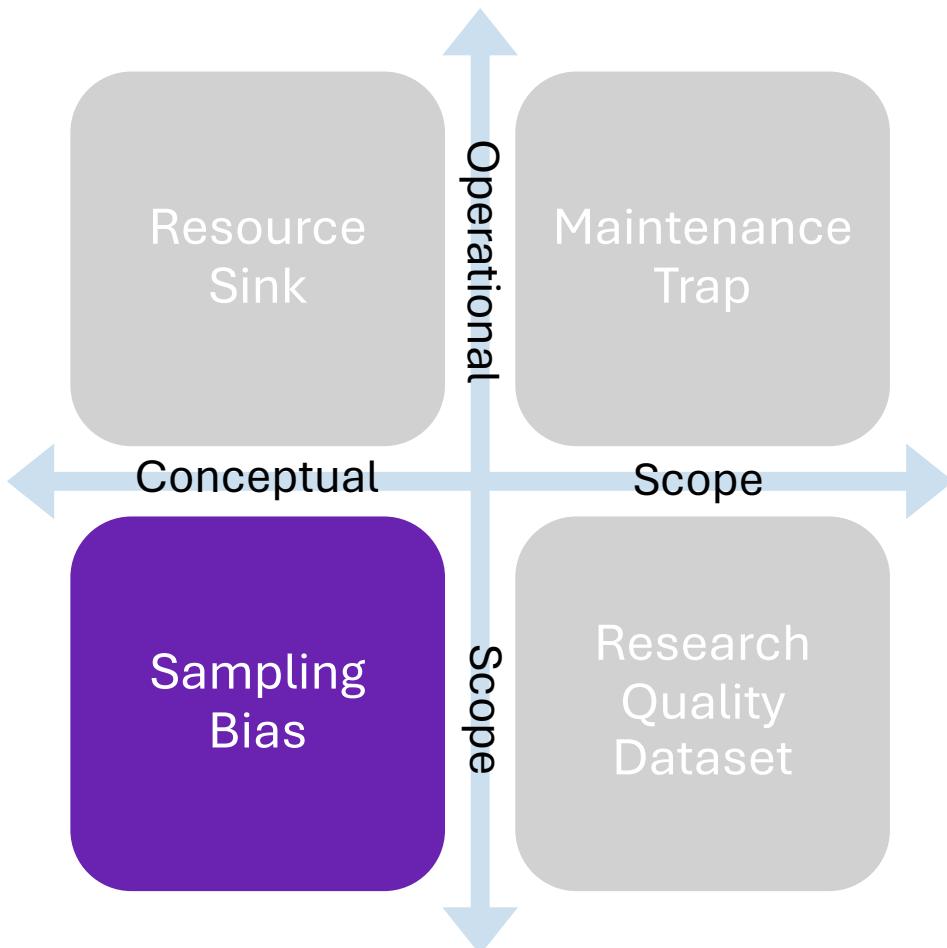
Maintenance Trap

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 data storage

Web Scrapped Data Risk



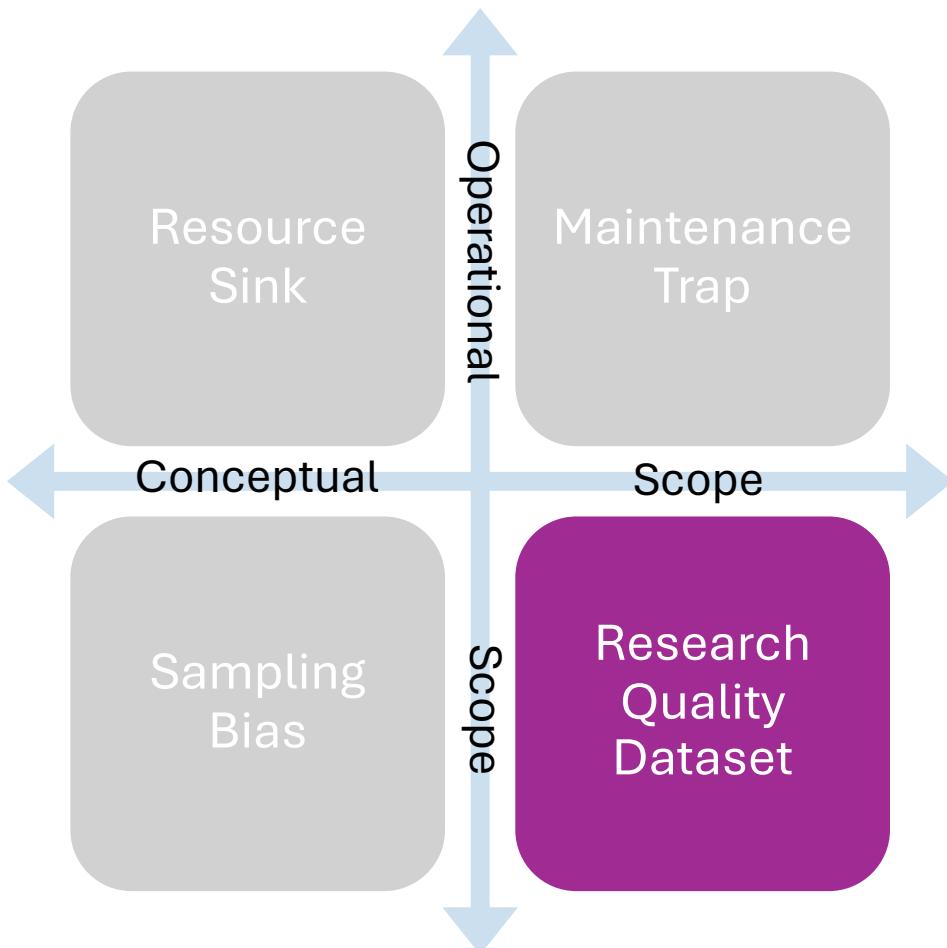
Sampling bias

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

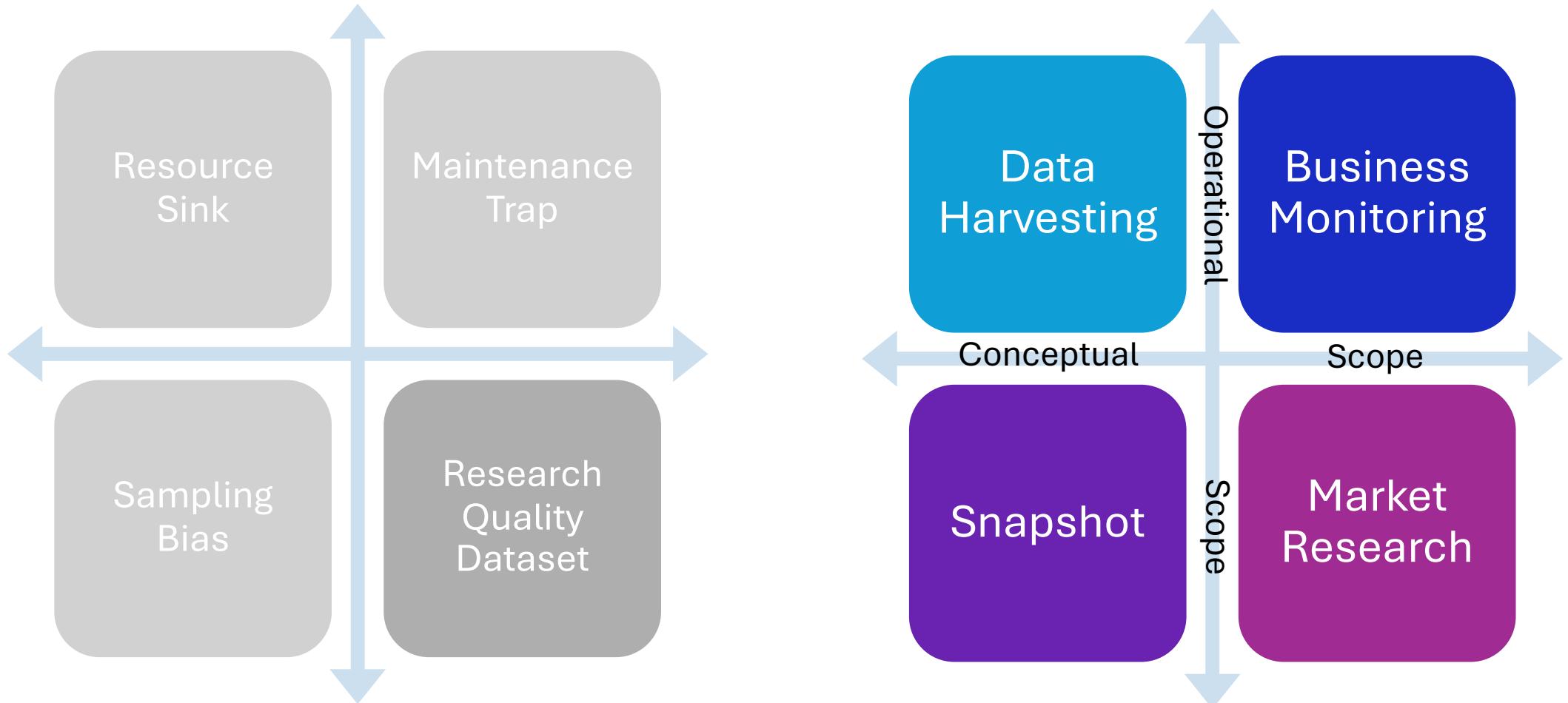
Web Scraped Data Risk



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)

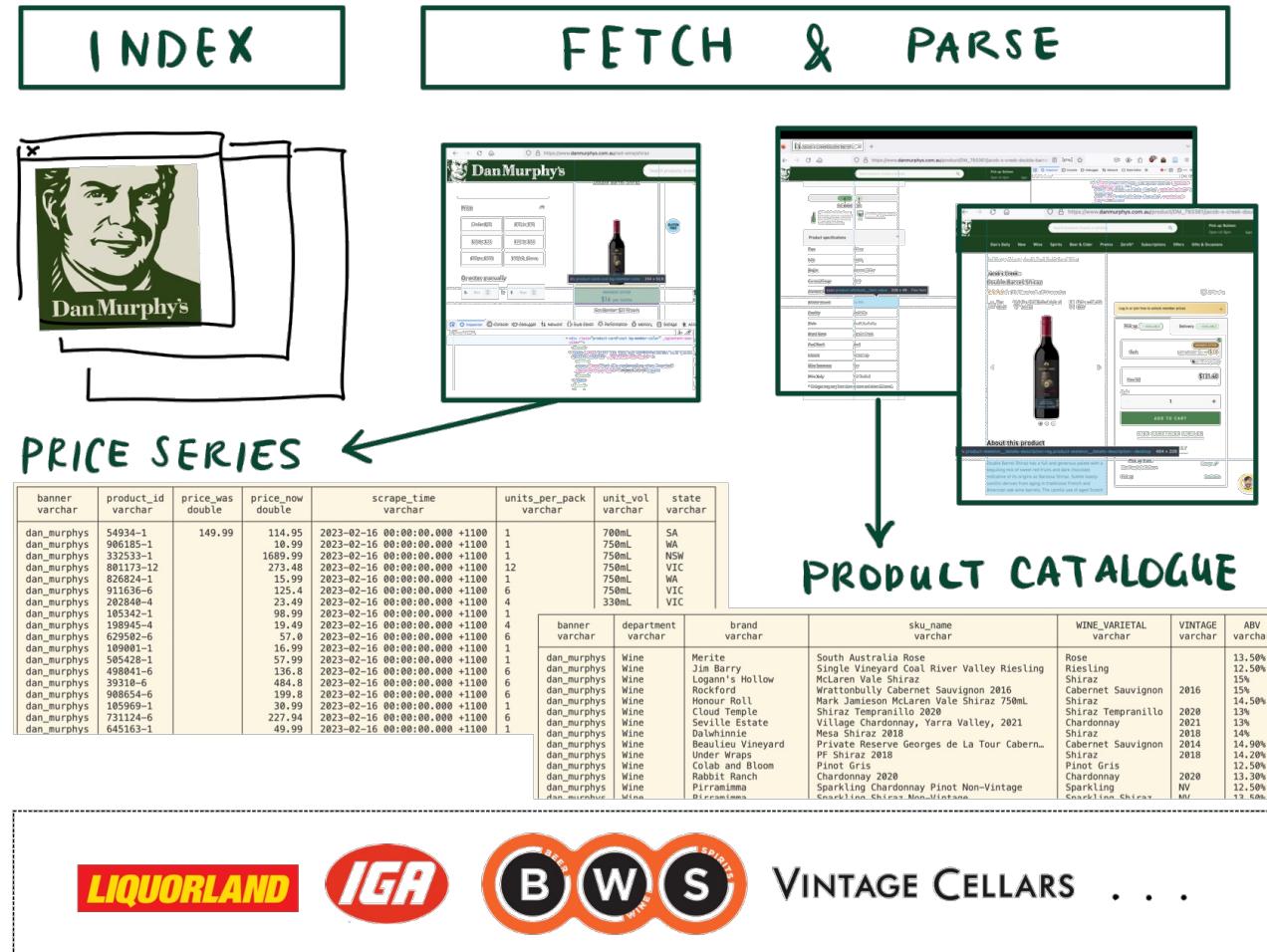
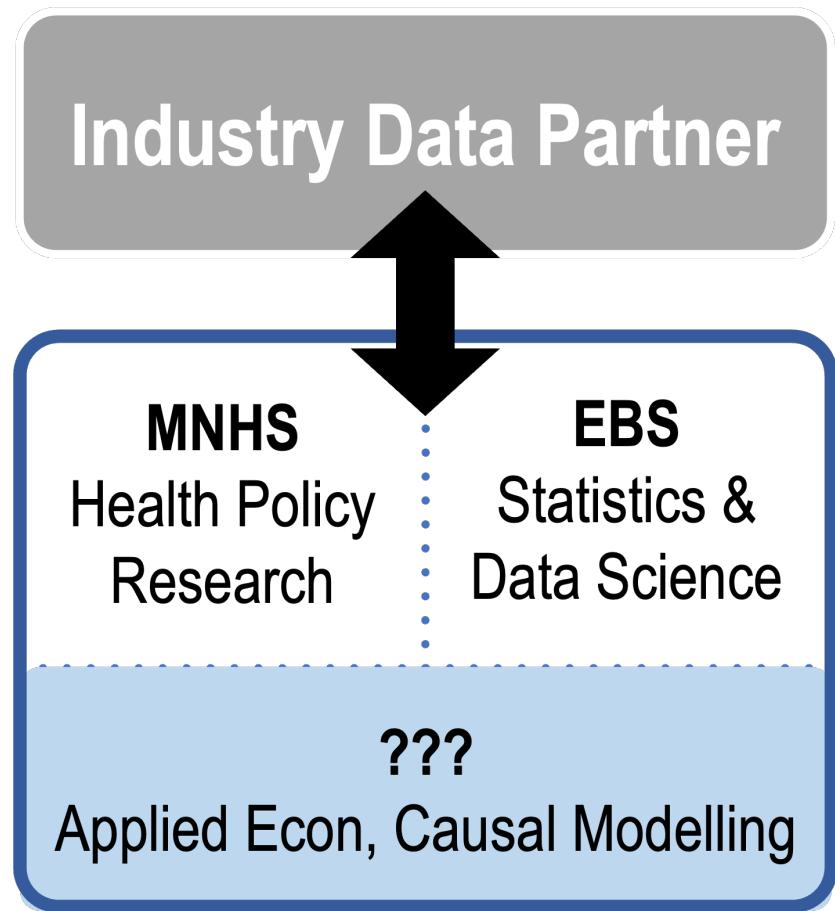
Business vs. Research Needs



Case Study

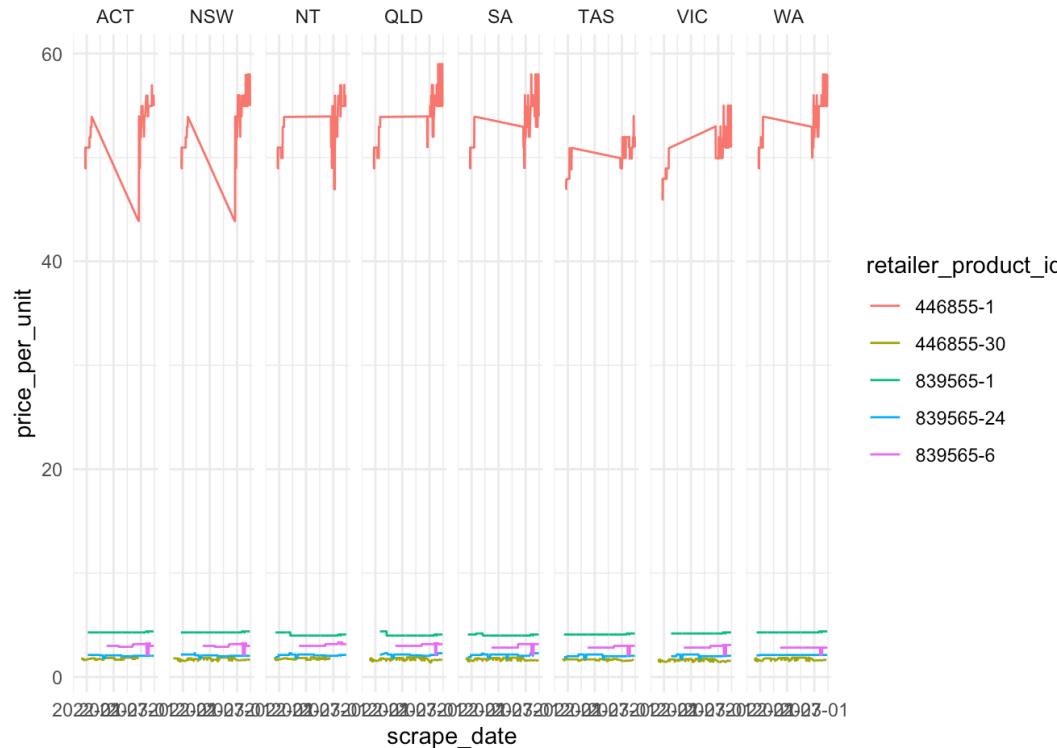
Adapting web scraped retail price data for
public health research

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 a Firefox browser window with the Yahoo Finance homepage loaded at <https://finance.yahoo.com>. The page features a navigation bar with links like Watchlists, My Portfolio, Markets, News, Videos, Screeners, Personal Finance, Crypto, Sectors, and Contact Us. Below the navigation is a search bar and a news ticker. The main content area displays several news stories and a stock market summary table. The table lists stocks such as DJT, UNH, CXAI, GOOG, and MU, along with their current price, change, and percentage change. The developer tools' Inspector tab is open over the table, highlighting specific HTML elements. The highlighted element is a table cell (`<td>`) containing a link to the company's quote page, with its class attribute set to `a.Fw(b).Ell.D(b).C($linkColor).Pos(r).Z(2)`.

Web technologies

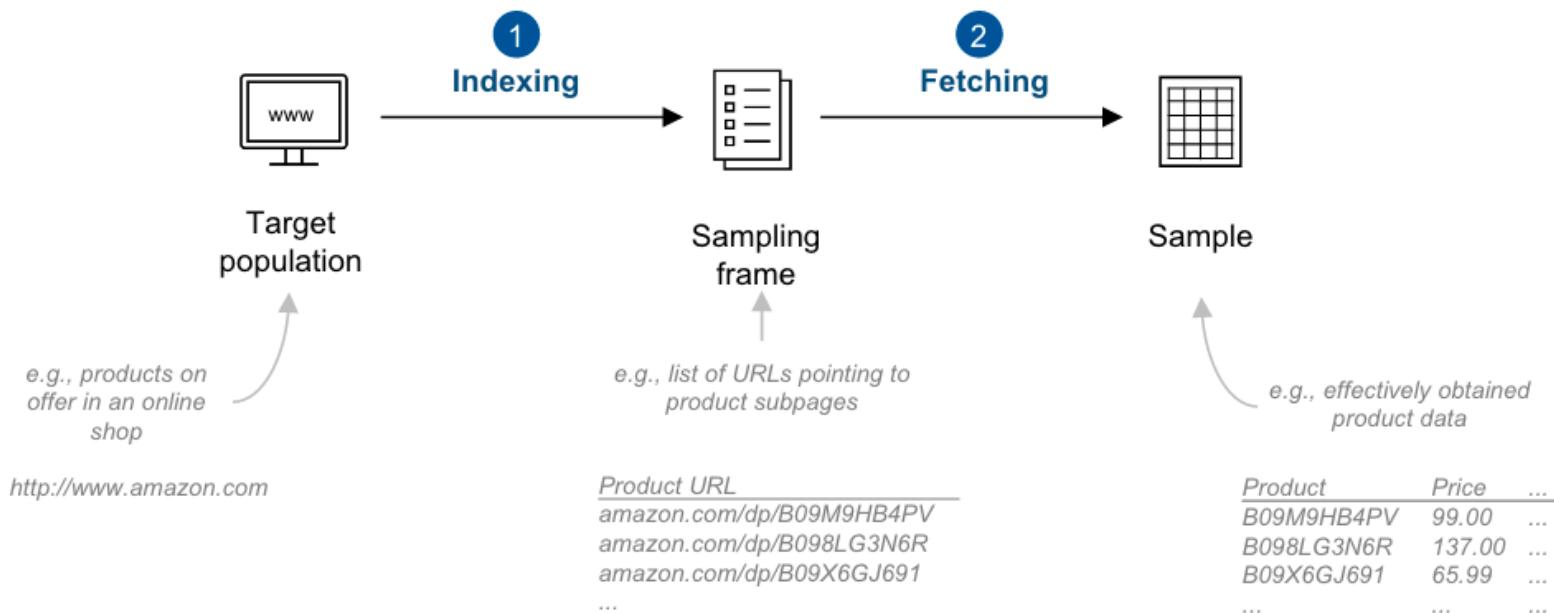
- Webpages and websites
 - Styling: CSS
 - Static content: HTML/XML
 - Dynamic content: Javascript, API queries
- Web Protocols (HTTP/HTTPS)
- Application Programming Interface (APIs)
 - Query parameters
 - Response formats: JSON, CSV...



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

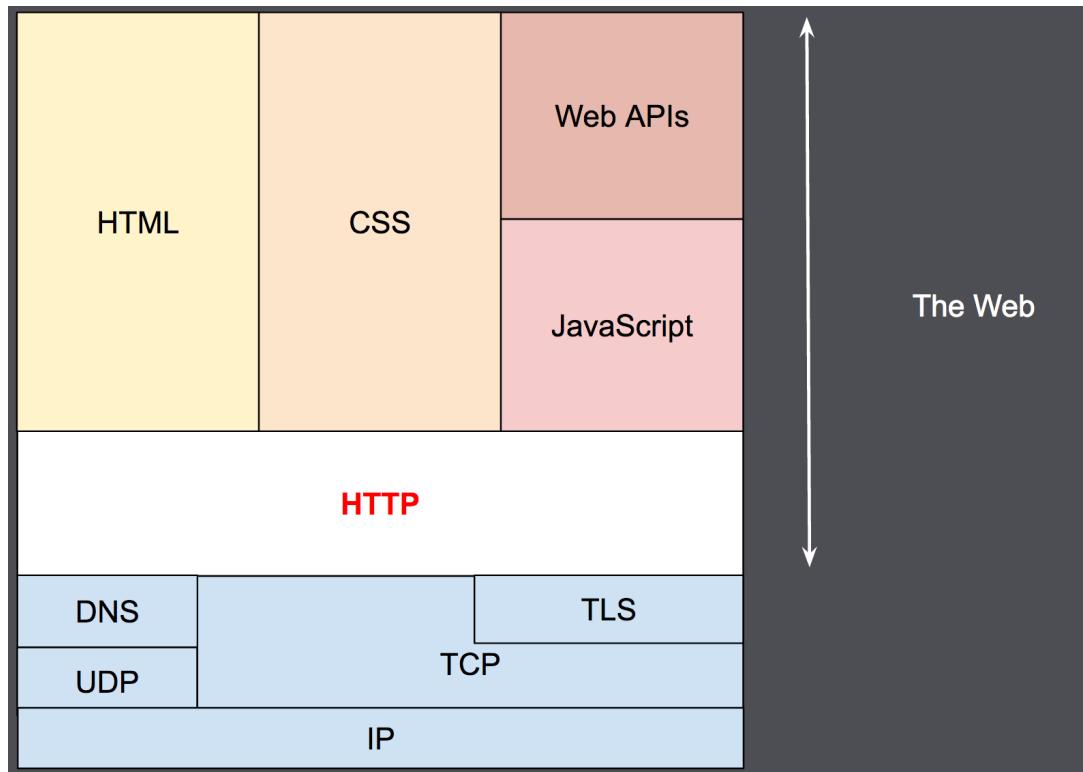
- Automating mouse/keyboard inputs
- Via headless browsers
- Selenium, chromium etc.

Request-based extraction

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



Web scraping methods: Parsing



- **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

An overview of HTTP - HTTP | MDN. (2023, December 16).
<https://developer.mozilla.org/en-US/docs/Web/HTTP/Overview>

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

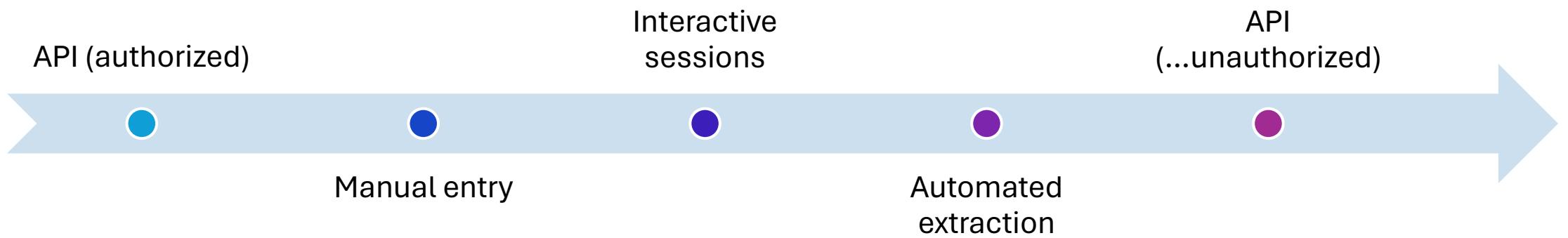
- HTML elements
- API responses

The screenshot shows a Firefox browser window with two tabs open: "Yahoo Finance - Stock Market" and "Buy Pinot Noir Wine Online [Untitled]". The main content area displays a search results page for "red-wine/pinot-noir" on the website [danmurphys.com.au](https://www.danmurphys.com.au). The results show two wine products: "Devil's Corner Pinot Noir" and "Jacob's Creek Classic Pinot Noir", both labeled as "GLUTEN FREE". The "Jacob's Creek" product has a price of "\$24.95 each" and a rating of "★★★★★ (0)". The "Devil's Corner" product has a price of "\$25.99 each" and a rating of "★★★★★ (48)". The "Jacob's Creek" product is currently selected, with its details highlighted.

The right side of the browser window shows the Network tab of the developer tools. It lists 317 requests made to the server. The requests are primarily 200 OK responses from the "api.da..." domain, with some 200 OK responses from "aem.d...", "auth.d...", and "aem:184...". The "Response" column shows the JSON data for the selected "Jacob's Creek" product. The JSON object contains information such as the product name ("Jacob's Creek Classic Pinot Noir"), pack default stock code ("140587"), and parent stock code ("140587"). Other objects listed include "Yarra View Yarra Valley Pinot Noir", "Devil's Corner Pinot Noir", "Jacob's Creek Classic Pinot Noir", "Mud House Claim 431 Vineyard Pinot Noir", "Cold Snap Cool Climate Pinot Noir", "Elephant In The Room Pinot Noir", and "Oyster Bay Pinot Noir". The "Timings" and "Stack Trace" sections are also visible at the bottom of the Network tab.

Final considerations

Ethical & Legal Risk: Retrieval



Ethical & Legal Risk: Data Usage

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

The image displays two side-by-side screenshots of web browser windows. The left window shows a terms of service page for 'Guardian Site' from the Guardian website. It includes sections about continued acceptance of terms, account termination, and copyright notice. The right window shows the 'AEA Data Legality Policy and Explanations' page from the American Economic Association website, featuring a header with the AEA logo and navigation links.

Guardian Site. Continued use of the Guardian Site indicates your continued acceptance of these terms and conditions.

We may terminate your registered account, at our sole discretion, by emailing you at the address you have registered stating that the account has terminated.

3. Use of material appearing on the Guardian Site

You acknowledge that, as between Guardian and you, except for content and advertisements (as discussed below), Guardian is the owner or the licensee of all content on the Guardian Site, including without limitation, all applicable copyrights, patents, trademarks, secrets, trade names, logos, and other intellectual property rights as well as text, images, graphics, logos, audio, video, URLs, meta-data and other material appearing on the Guardian Site ("Guardian Content"). Guardian Site and the Guardian Content are protected by the copyright laws and other intellectual property laws of the United Kingdom, United States and Australia and are protected globally by applicable international copyright treaties.

Your use of the Guardian Site and Guardian Content is for your own personal and non-commercial use only. You may download and print extracts from the Guardian Content for your own personal and non-commercial use only, provided you maintain and abide by any author attribution, copyright or trademark notice or restriction in any material that you download or print.

AMERICAN ECONOMIC ASSOCIATION

Home > Journals > AEA Data and Code Policies and Guidance > AEA Data Legality Policy and Explanations

AEA Data Legality Policy and Explanations

Advice to Authors

A more extensive explanation of the Data Legality Policy's rationale is given below, but a key rationale is that under the Policy the discussion of data legality occurs early in the editorial process so that there are no surprises at the end.

AEA Journal Editors or the AEA Data Editor would be happy to answer questions about data legality before you submit your paper. If you use data that you believe might not be considered legally acquired, please flag this in the online submission portal when submitting your paper.

Starting a Web Scrapped 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.”