

JOURNAL of

# Marketing



Co-Host: Christine Moorman  
Editor-in-Chief  
Duke University

June 23, 2022

AM> | AMERICAN MARKETING  
ASSOCIATION

Webinar

---

JOURNAL of

# Marketing

The *Journal of Marketing* develops and disseminates knowledge about real-world marketing questions relevant to scholars, educators, managers, policy makers, consumers, and other societal stakeholders around the world.



# **FIELDS OF GOLD: SCRAPING WEB DATA FOR MARKETING INSIGHTS**



Johannes Boegershausen

Erasmus University,  
The Netherlands



Hannes Datta

Tilburg University,  
The Netherlands



Abhishek Borah

INSEAD,  
France



Andrew Stephen

Oxford University,  
The UK

# ► Agenda

- **Introduction**
  - Web data in academic marketing research & how to extract it
  - Pathways for creating new marketing knowledge
- **Methodological framework**
  - Managing the idiosyncratic legal, technical and validity challenges of web data
  - Focus on three key stages: source selection, design, extraction
- **Food for thought & conclusion**
  - Key insights
  - Exploiting new fields of gold
- **Q&A**

Web data in academic marketing research

# INTRODUCTION

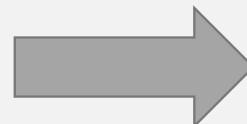
## ► Enormous & diverse data for marketing research

**7:11**  
hours

**85%**

time spent online  
per day by the  
average American  
consumer

proportion of US  
consumers that  
use the Internet  
every single day



~ 244m reviews



> 1b reviews & opinions



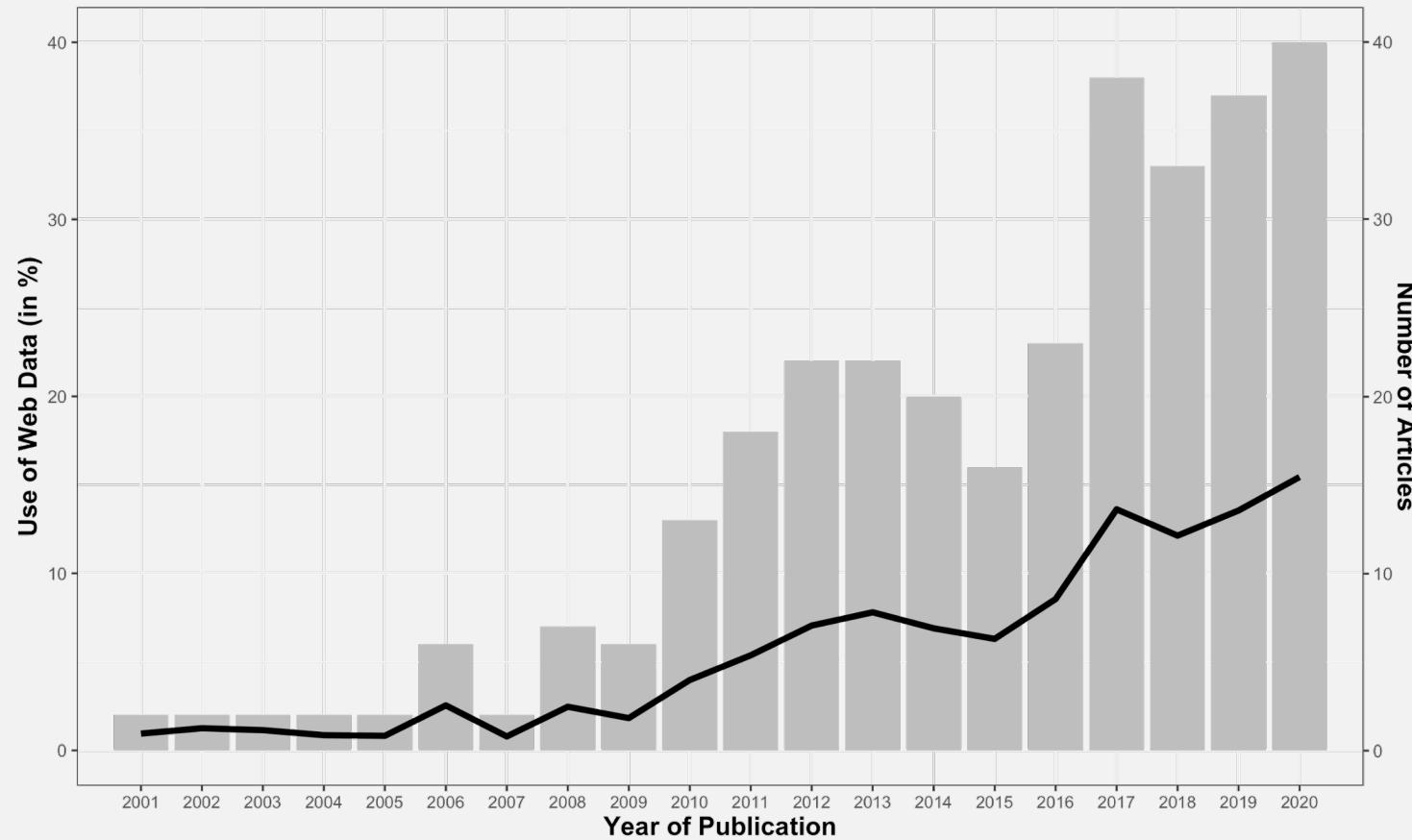
500m/day



556K projects

based on available company and market research statistics in May 2022

## ► Increasing usage of web data in marketing research



# ► Extracting web data at scale via...



## Web Scraping

... the process of developing software to automatically collect information **displayed in a web browser**

EXAMPLE SOURCES

The block contains three logos side-by-side: the red Yelp logo, the black and orange Amazon logo, and the yellow and green TripAdvisor logo.

Example articles:  
Chevalier & Mayzlin (2006); Ludwig et al. (2013)



Application Programming Interface

... allow programmatic access to the **internal databases or algorithms of data providers**

EXAMPLE SOURCES

The block contains four logos: the Spotify for Developers logo (white circle with three horizontal lines), the Google Cloud Vision API logo (hexagonal geometric pattern), the Amazon Product Advertising API logo (black Amazon logo with "Product Advertising API" text), and the Twitter API logo (blue bird icon).

Example articles:  
Tellis et al. (2019); Toubia & Stephen (2013)

# ► Highly versatile data collection technique

Pathway ①

Studying new phenomena



e.g., Zervas et al. (2017); Datta et al. (2018)

Pathway ③

Facilitating methodological advancement



e.g., Netzer et al. (2012); Liu et al. (2020)

Pathway ②

Boosting ecological value



e.g., Du et al. (2015); Ludwig et al. (2013)

Pathway ④

Improving measurement



e.g., Li et al. (2017); Datta et al. (2022)

Source: Boegershausen, Datta, Borah, and Stephen (2022)

## ► Collecting web data can be challenging

- Data generation process is often opaque
- Highly dynamic and unstable environment
- Mostly poorly or undocumented measures
- Cannot be “downloaded” → needs to be generated through automated browsing
- Numerous *idiosyncratic pitfalls* (*more on that later*)

Single vs.  
multisource?

Algorithmic  
biases?

Extraction frequency?

How to sample?

Which software  
to use?

Keep the raw  
(HTML, JSON)  
data?



Single vs.  
multisource?

Algorithmic  
biases?

Extraction frequency?

## **Existing guidance limited**

- Focus on technicalities (and not validity)
- Unclear how to deal with or mitigate legal concerns
- Scope of methodological guidance rather narrow

How to sample?

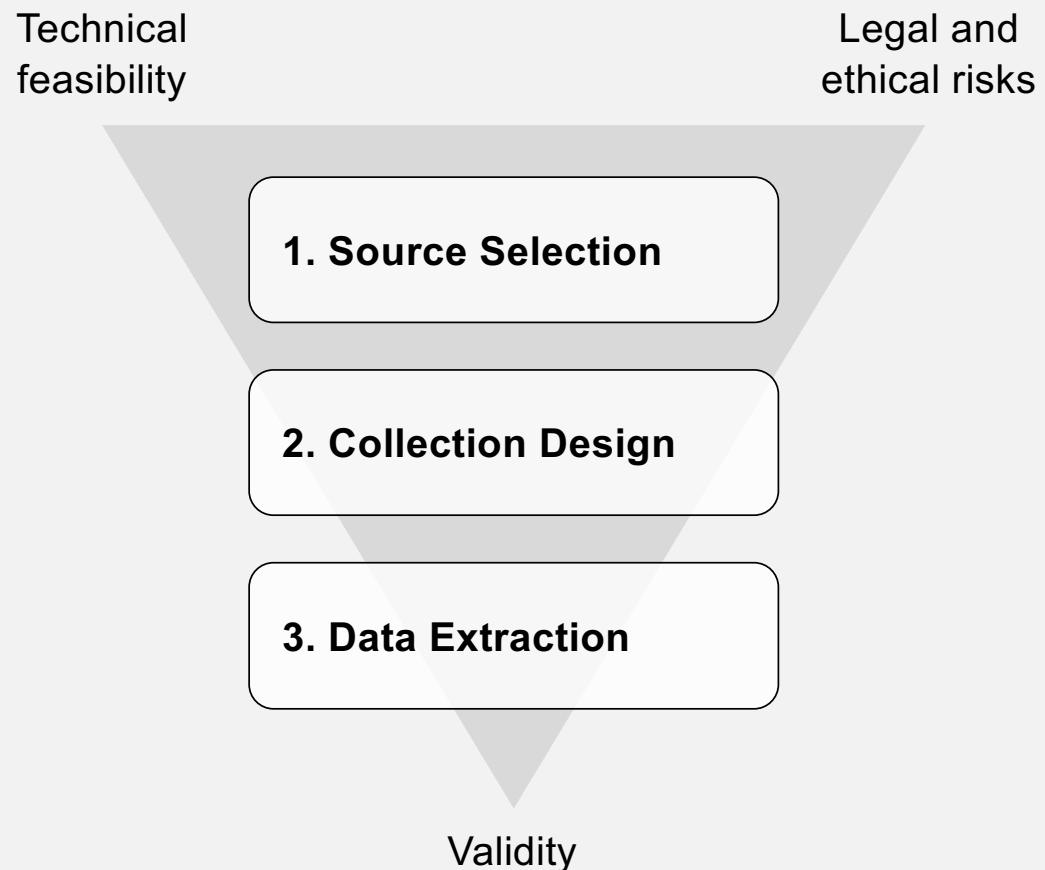
Which software  
to use?

Keep the raw  
(HTML, JSON)  
data?

Managing the idiosyncratic legal, technical and validity challenges of web data

## **METHODOLOGICAL FRAMEWORK**

# ► Methodological framework



Source: Boegershausen, Datta, Borah, and Stephen (2022)

## ► **Source selection:** challenges

- Access to near-to infinite number of potential sources without traditional gatekeepers. Different forms of access.
- But sources vary vastly in terms of quality, stability, and retrievability.  
→ Might prompt researchers to primarily consider dominant or familiar platforms only.





## ► **Source selection:** recommendations I

- **Explore the universe** of potential web sources
  - Broaden geographic search criteria (e.g., non-Western)
  - Identify adjacent data sources (e.g., Google Trend's "related search queries")
  - Expand search to non-primary data providers (e.g., aggregators like SocialBlade)



## ► Source selection: example





## ► Source selection: example



tripadvisor



How America finds a doctor.\*



## ► Source selection: example





## ► Source selection: recommendations II

- Explore the universe of potential web sources
  - Broaden geographic search criteria (e.g., non-Western)
  - Identify adjacent data sources (e.g., Google Trend's "related search queries")
  - Expand search to non-primary data providers (i.e., aggregators, databases)
- Consider alternatives to web scraping
  - Expand search by explicitly including terms such as "API" or "dataset download"
  - APIs? How does the data compare to data that could be scraped?

**Recommender Systems and  
Personalization Datasets**

Julian McAuley, UCSD

yelp\* Dataset

kaggle

# ► Source selection: recommendations III



- Explore the universe of potential web sources
  - Broaden geographic search criteria (e.g., non-Western)
  - Identify adjacent data sources (e.g., Google Trend's "related search queries")
  - Expand search to non-primary data providers (i.e., aggregators, databases)
- Consider alternatives to web scraping
  - Expand search by explicitly including terms such as "API" or "dataset download"
  - APIs? How does the data compare to data that could be scraped?
- **Map the data context**
  - Screen blogs, press releases, a source's software "changelogs," ...
  - Understand changes to the data-generating process (e.g., archive.org)
  - Algorithms present? Visit source using different devices/times, inspect source code

# ► Designing the data collection

Technical  
feasibility

Legal and  
ethical risks

**1. Source Selection**

**2. Collection Design**

**3. Data Extraction**

Validity

# ► Which information to extract? Example

ASTRO Gaming A20 Wireless Headset Gen 2 for Xbox Series X | S, Xbox One, PC & Mac - White /Grey

◀ Back to results



Gaming Headset with Microphone,  
Gaming Headphones Stereo 7.1  
Surround Sound PS4 Headset 50mm  
Drivers, 3.5mm Audio Jack Over Ear  
Headphones Wired for PC Switch  
Playstation Xbox PS5 Laptop

VISIT the FEIYING Store

★★★★★ 1,215 Ratings | 35 answered questions

Amazon's Choice for "gaming headsets"

List Price: \$49.97 Details  
With Deal: **\$17.31**  
You Save: **\$32.66 (65%)**

No Import Fees Deposit & \$11.60 Shipping to Netherlands

Details ▾  
Coupon:  Save an extra 7% when you apply this coupon.  
Terms ▾

Roll over image to zoom in



# ► Which information to extract? Example

Sponsored ⓘ

**Customer reviews**

★★★★★ 4.5 out of 5

1,215 global ratings

Rating	Percentage
5 star	75%
4 star	10%
3 star	6%
2 star	3%
1 star	6%

How customer reviews and ratings work

**By feature**

Feature	Rating
Value for money	★★★★★ 4.6
Comfort	★★★★★ 4.6
For gaming	★★★★★ 4.5

See more

**Review this product**

Share your thoughts with other customers

**Reviews with images**

Sponsored ⓘ

See all customer images

Read reviews that mention

sound quality noise cancellation son loves highly recommend  
gaming headset noise cancelling definitely recommend really good  
high quality great price comfortable to wear listening to music

Top reviews from the United States

Zane

★★★★★ Very Nice Gaming Headset with Microphone

Reviewed in the United States on January 15, 2022

Color: A Camo Gray | Verified Purchase

## ► Which information to extract? Example

Zane

Have you checked out your Public Profile yet? Make sure it's up to date!

[View your public profile](#)

**About**

**Impact**

24

**Community activity**

1 Reviews    0 Idea Lists

Zane reviewed a product · Jan 15, 2022

★★★★★ Verified Purchase  
Very Nice Gaming Headset with Microphone

# ► Which information to extract? Example



## Validity implications

- Is information subject to algorithmic biases or missing data?  
**Delete cookies & check?**
- Are there significant changes to the data-generating process?  
**Archive.org**
- Is meta data required to make sense of variables?  
**Save timestamps/IP addresses**

## Legal/ethical risks

- Publicly accessible vs. login? Consent to ToS? Implicit or explicit?  
**Focus on public pages**
- Personal or sensitive information?  
**Anonymize while collecting**
- Overlap original intent of posting & research question / scientific justification?  
**Formulate scientific justification**

## Technical feasibility

- All information extractable?  
**Build prototype**
- Limits to iterating through pages?  
**Check last page, try a few in-between**

# ► How to sample? Challenges & considerations



- How to capture the entire population (or a sample) of...?
  - Internal pages (e.g., bestseller, category, search page)
  - Externally available lists?
- Sampling frames (might) create different datasets or even induce systematic biases
- Which sample size is technically feasible?

## ► At what frequency to extract data? Challenges

- Validate “data” assumptions early on
  - Configuration (e.g., “data is historically available”)
  - Data-generating process (e.g., “website hasn’t changed”)
  - Characteristics (e.g., measurement is clear; use of interpolation)
- Examples
  - Archival versus “live” data → discover fake reviews
  - Gains from capturing information more than once? → build longitudinal data set
  - Balance sample size and extraction frequency → sufficient power?



## ► How to process data during the extraction?

- Web data is “messy”
  - BUT “on-the-fly” processing can create significant threats to validity
- **Keep the raw data whenever possible**

# ► How to process data during the extraction?



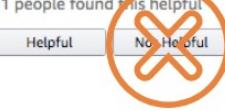
- Web data is “messy”
- BUT “on-the-fly” processing can create significant threats to validity  
→ Keep the raw data whenever possible
- **Opportunity:** “stumbling” into natural experiments

★★★★★ Well worth its cost.  
October 5, 2017  
Style: W/ CR123A Batteries | Package Type: Plastic Clamshell Pa

Without a doubt, a top notch light instrument for everyday carry, never leaves my possession. I've kept it clipped into a back pocket. Furthermore, the lumen power is plenty powerful enough to more than hold its own. I've exposed it to free flowing water... to extended day and overnight shifts. You won't be disappointed... especially if you also purchase the 18650 Button Top AC Li-Ion 120V which is also found here on Amazon.

11 people found this helpful

 3 people found this helpful  
 35 people found this helpful

 Helpful | No Helpful | Comment | Report abuse

NEWS & EVENTS

## An update to dislikes on YouTube

By The YouTube Team  
Nov. 10. 2021

# ► Data extraction

Technical  
feasibility

Legal and  
ethical risks

**1. Source Selection**

**2. Collection Design**

**3. Data Extraction**

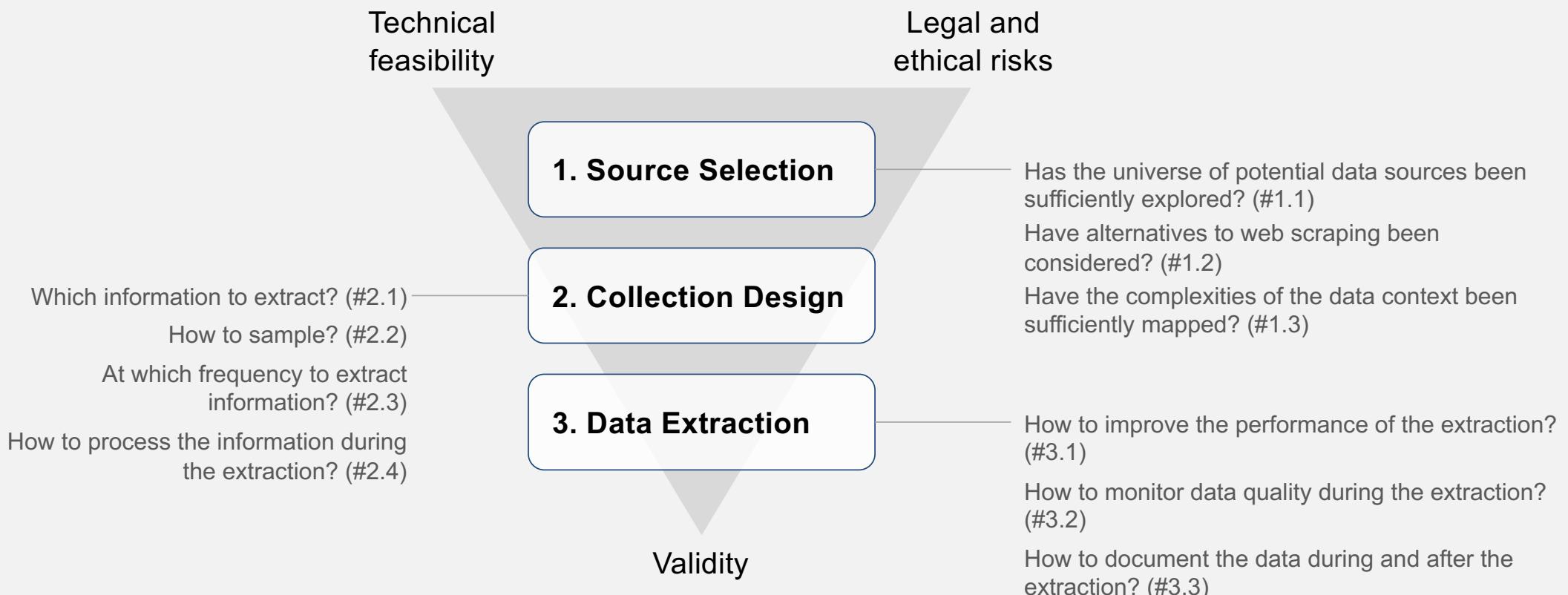
Validity

# ► Data extraction



- How to **improve** the performance of the data extraction?
  - Keep the collection running for some time – does it continue to work?
  - Log the (timestamped) URLs of scraped pages and visualize performance over an extended period.
- How to **monitor** data quality during the extraction?
  - Collect and report metadata to diagnose issues in real-time
- How to **document** the data during and after the extraction?
  - Nobody, except you, knows how the data was generated!
  - Start early! Logbook. Collect information around the focal source(s).

# ► Methodological framework: summary



Source: Boegershausen, Datta, Borah, and Stephen (2022)

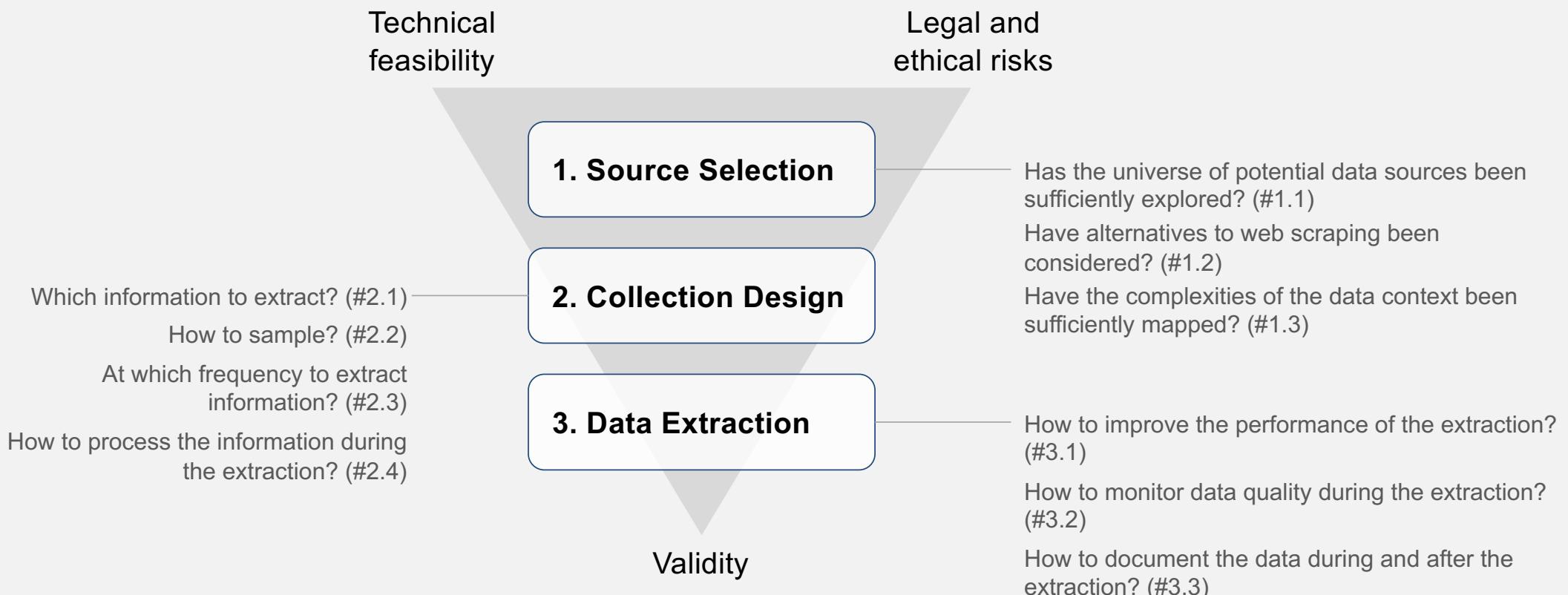
# ► Our paper helps reasoning through design challenges...

TABLE 3 CHALLENGES AND SOLUTIONS IN DESIGNING WEB DATA COLLECTIONS			
Challenge #2.1: Which information to extract from which pages?			
Validity Challenges [V]	Legal/Ethical Challenges [L]	Technical Challenges [T]	Solutions and Best Practices
<ul style="list-style-type: none"> <li>• Which information is necessary to justify construct operationalization and allow analysis?</li> <li>• Which metadata might enhance internal and external validity?</li> <li>• Is information subject to algorithmic biases or missing data?</li> <li>• Are there significant changes to the data-generating process?</li> </ul>	<ul style="list-style-type: none"> <li>• Is all of the required information publicly accessible, or is a login required?</li> <li>• Does the data contain personal or sensitive information, and can subjects be identified?</li> <li>• Is there a sufficient scientific justification for using the data?</li> <li>• How large is the overlap between the research objective and the original intent of subjects disclosing the data?</li> </ul>	<ul style="list-style-type: none"> <li>• Explore different types of pages to determine vs. identical information [V]</li> <li>• Explore whether alternative ways to browse/interact with the site (e.g., URLs, clicking, scrolling, logging in to the site) provides different or reveals new information [T]</li> <li>• Explore how extraction methods (e.g., “headless” HTTP requests vs. simulated browsing, different user agents, screen width, login status, use of different packages) affect information displays [V, T]</li> <li>• Assess the accuracy of timestamps (e.g., time zones) [V]</li> <li>• Explore (temporarily available) information in the source code of a website using the browser’s “inspect” tools [V]</li> <li>• Assess the presence of technical roadblocks (e.g., captchas) [T]</li> <li>• Assess how data was generated historically at the source (e.g., via archive.org) [V]</li> <li>• Explore limits to iterating through pages [T]</li> <li>• Obtain information from various sources to reduce dependency on data provider [L]</li> <li>• If possible, opt-out of firm-administered experiments or block cookies; alternatively, identify relevant metadata that can be used to control for the presence of algorithms [V]</li> </ul>	
Challenge #2.2: How to sample?			
Validity Challenges [V]	Legal/Ethical Challenges [L]	Technical Challenges [T]	Solutions and Best Practices
<ul style="list-style-type: none"> <li>• Is the sample size sufficient to effectively inform the research question?</li> <li>• To which population does the sample generalize?</li> <li>• Is the sampling frame corresponding to the research objective (e.g., randomness)?</li> <li>• How prevalent is panel attrition?</li> </ul>	<ul style="list-style-type: none"> <li>• Does the data represent an excessive portion relative to all data available?</li> <li>• Can external information (e.g., IDs) be consistently matched to the data?</li> <li>• Are some of the sampled units (potentially) vulnerable?</li> </ul>	<ul style="list-style-type: none"> <li>• Is the required sample size technically feasible?</li> <li>• Assess characteristics of the population (e.g., using secondary sources) [V]</li> <li>• Explore options to sample directly from the source (e.g., from different pages, randomization through filtering/searching, obtaining usernames from forums; see also Neudorf 2017 and Humphreys and Wang 2018) [V]</li> <li>• Choose lists or pages that are not affected by algorithmic influence [V]</li> <li>• Refresh sample (or use multiple types of sampled units) to assess the stability of sample and counterbalance panel attrition [V]</li> <li>• Discard units from the sample to prevent data collection from subjects falling under prohibitive national and supranational legislation (e.g., GDPR) [L]</li> <li>• Explore external sources to inform the sampling frame (e.g., or facilitate linkage [T]</li> <li>• Assess the efficiency of different navigation paths and their impact on sample size [T]</li> <li>• Pseudo-anonymize or discard sensitive or personal information [L]</li> <li>• Ensure no excessive amount of data (e.g., data on all users) is collected (absolute volume, relative volume) [L]</li> <li>• Re-examine alternative sources to improve justification of data extraction [L]</li> </ul>	

TABLE 3 CHALLENGES AND SOLUTIONS IN DESIGNING WEB DATA COLLECTIONS [CONTINUED]			
Challenge #2.3: At which frequency to extract the data?			
Validity Challenges [V]	Legal/Ethical Challenges [L]	Technical Challenges [T]	Solutions and Best Practices
<ul style="list-style-type: none"> <li>• Is the extraction frequency in sync with the studied phenomena?</li> <li>• Is the refresh rate of the source sufficient?</li> <li>• Is the data (thought to be archival) really archival?</li> <li>• Is the information consistently available across all periods of interest?</li> <li>• Does the order and frequency in which information is retrieved induce bias?</li> </ul>	<ul style="list-style-type: none"> <li>• Does the extraction frequency pose an excessive load on the source?</li> <li>• Does collecting more data at higher frequencies make the data more sensitive?</li> <li>• Does the desired frequency pose new technical hurdles?</li> <li>• How can the stability of data collection be guaranteed, and different collection batches be distinguished?</li> </ul>	<ul style="list-style-type: none"> <li>• Explore the gains in collecting data multiple times rather than once (e.g., in a “live” data collection) [V]</li> <li>• Adhere to best practices in setting the extracting frequency (e.g., 5 requests per second for APIs, 1 request per 2 seconds for web scraping) [L, T]</li> <li>• Experiment with technical parameters (e.g., number of computers) to balance technically feasible sample size and desired frequency of data extraction [T]</li> <li>• Formulate test, and refine data source theory (Landers et al. 2016) [V]</li> <li>• Reinspect the robots.txt file to avoid exceeding retrieval limits for selected pages [T]</li> <li>• Consider randomizing extraction order for sampled units over time [V]</li> <li>• Consider (cost) implications for storage and computation time [T]</li> <li>• Consider getting in touch with the data provider if the targeted data set is infeasible to extract via web scraping or APIs [T, L]</li> <li>• If possible, opt-out of firm-administered experiments or block cookies; alternatively, identify relevant metadata that can be used to control for the presence of algorithms [V]</li> </ul>	
Challenge #2.4: How to process the data during the collection?			
Validity Challenges [V]	Legal/Ethical Challenges [L]	Technical Challenges [T]	Solutions and Best Practices
<ul style="list-style-type: none"> <li>• Could erroneous processing lead to unexpected data loss?</li> <li>• Could there be any significant scientific value in retaining the raw data?</li> </ul>	<ul style="list-style-type: none"> <li>• Is the collected data in conflict with prohibitive laws (e.g., GDPR)?</li> <li>• Is the collected data sufficiently secured from unauthorized access?</li> <li>• Is normalization or pseudonymization required?</li> </ul>	<ul style="list-style-type: none"> <li>• Which storage facilities to use to accommodate the expected data (size, location, format, encoding)</li> <li>• Is normalization necessary?</li> </ul>	<ul style="list-style-type: none"> <li>• Retain raw data (e.g., HTML pages, JSON responses) whenever possible [V, T]</li> <li>• Always parse some minimal amount of data (e.g., timestamps) to facilitate monitoring checks in real-time [V, T]</li> <li>• Remove sensitive and personal information on the fly; if personal or sensitive information is strictly required to meet the research objective, consider pseudo-anonymizing (potentially via third parties) [L]</li> <li>• Verify data storage during collection meets legal requirements for potentially sensitive or personal data [L]</li> <li>• Ensure proper encoding of (non-English) characters, retain correct digit separators and correct data format</li> </ul>

IMPORTANT: trade-offs are (almost) inevitable  
**MAKE TRADE-OFFS EXPLICIT IN THE MANUSCRIPT**

# ► Questions

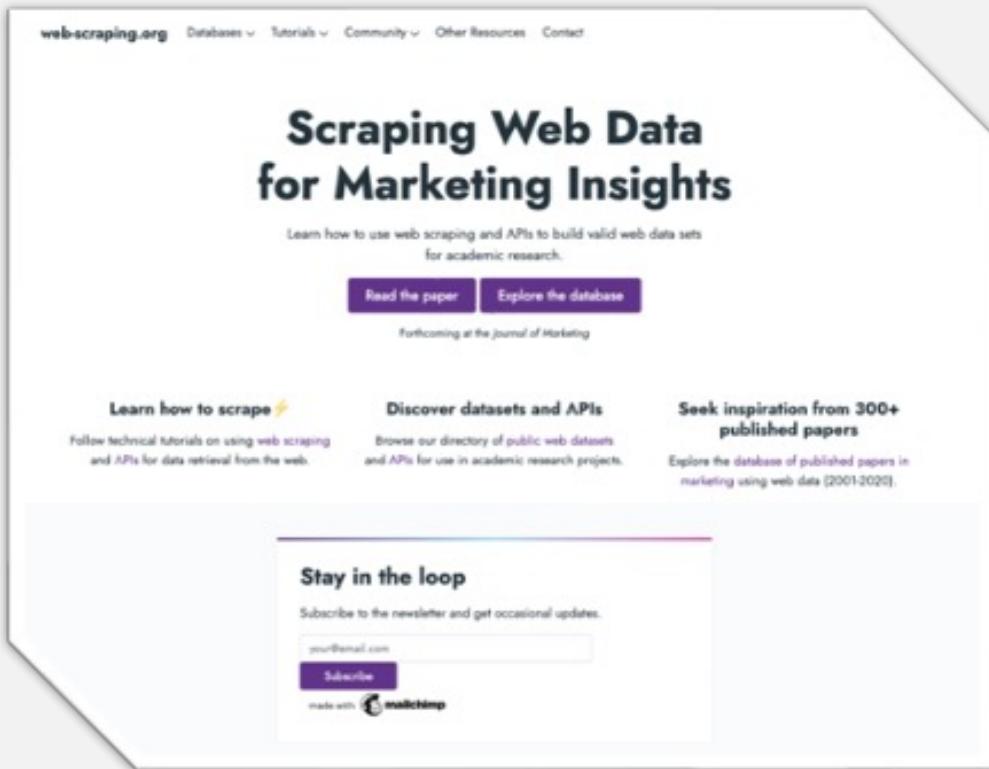


Source: Boegershausen, Datta, Borah, and Stephen (2022)

Key insights & exploiting new fields of gold

## **FOOD FOR THOUGHT**

# ► Our framework & companion website



- Explore our **database with 300+ published marketing articles** using web data.
- Discover **web datasets & APIs** for your research projects.
- **Tutorials and example code** for collecting web data using web scraping & APIs.

# ► ... for teaching

**Consumer Behavior, Ph.D. seminar, 2021-2022 | Rotterdam School of Management**

**Session 3:** [Exploring Marketplaces, People Perception, and Morality with Web Data](#)

**Faculty:** Johannes Boegershausen

**Readings\*:**

[1] Goodwin, Geoffrey P. (2015), "Moral Character in Person Perception," *Current Directions in Psychological Science*, 24 (1), 38-44.

[2] Hoover, Joseph, Morteza Dehghani, Kate Johnson, Rumen Iliev, and Jesse Graham (2018), "Into the Wild: Big Data Analytics in Moral Psychology," in *The Atlas of Moral Psychology*, Jesse Graham and Kurt Gray, eds. New York: Guilford Press, 525-36.

[3] Boegershausen, Johannes, Abhishek Borah, Hannes Datta, and Andrew T. Stephen (2021), "Fields of Gold: Generating Relevant and Credible Insights Via Web Scraping and APIs", *working paper*, <https://dx.doi.org/10.2139/ssrn.3820666>.

[4] Howe, Lauren C. and Benoit Monin (2017), "Healthier Than Thou? "Practicing What You Preach" Backfires by Increasing Anticipated Devaluation," *Journal of Personality and Social Psychology*, 112 (5), 718-35.

[5] Kirmani, Amna, Rebecca W. Hamilton, Debora V. Thompson, and Shannon Lantzy (2017), "Doing Well Versus Doing Good: The Differential Effect of Underdog Positioning for Moral and Competent Service Providers," *Journal of Marketing*, 81 (1), 103-17.

## Ph.D. seminars

 **UNIVERSITY**  
Open Education

**Online Data Collections (oDCM)**

Search

Course

Modules

Tutorials

Team Project

Final Exam

About

**Online Data Collection and Management (oDCM)**

Instructor: dr. Hannes Datta  Follow @hannesdatta

Course codes: 328060-M3 (fall, block 1) and 328061-M3 (spring, block 3)

This edition: August - October 2022 | Next edition: February - April 2023

**Learn how to mine the web**

Welcome to the course website of oDCM.

This course teaches you the nuts and bolts about collecting web data using web scraping and APIs. Unlike most other courses on this topic, this course not only teaches you the *technicalities* of using web scraping and Application Protocol Interfaces (APIs), but also introduces a *comprehensive framework* that helps you to *think* about scraping - specifically with regard to its application in empirical marketing research.

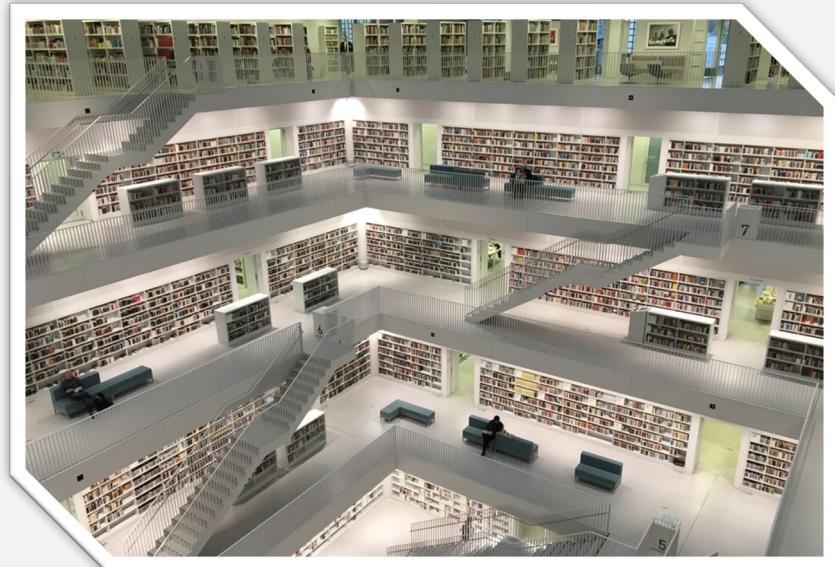
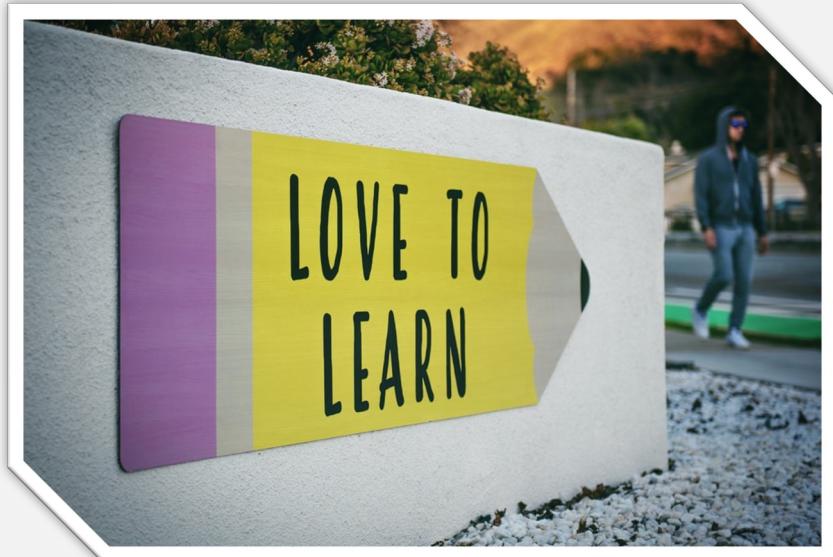
[Check out the course syllabus & learning objectives](#)

[View the schedule](#) [Enroll now!](#)

## method courses

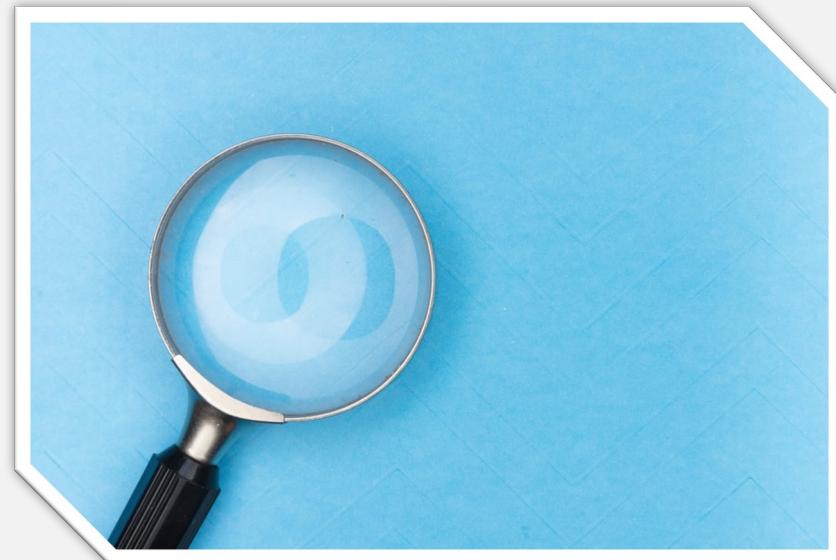
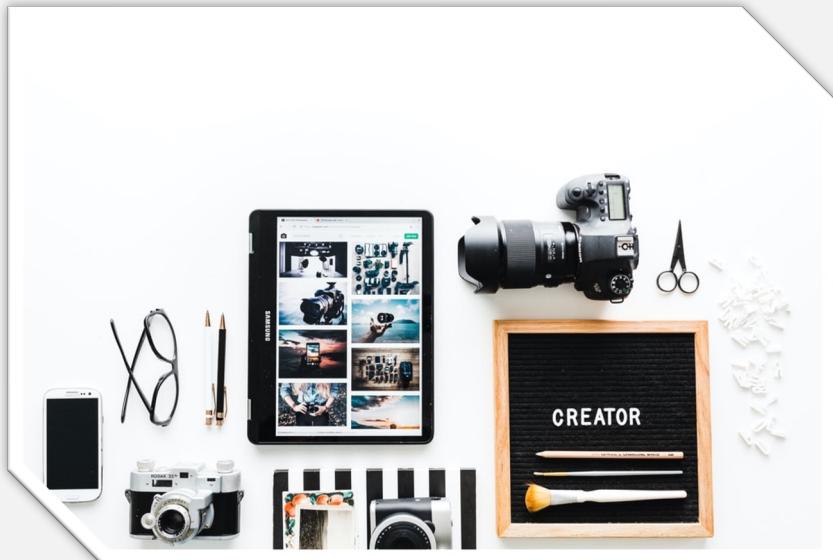
## ► ... as an inspiration?

- If you want to learn scraping...



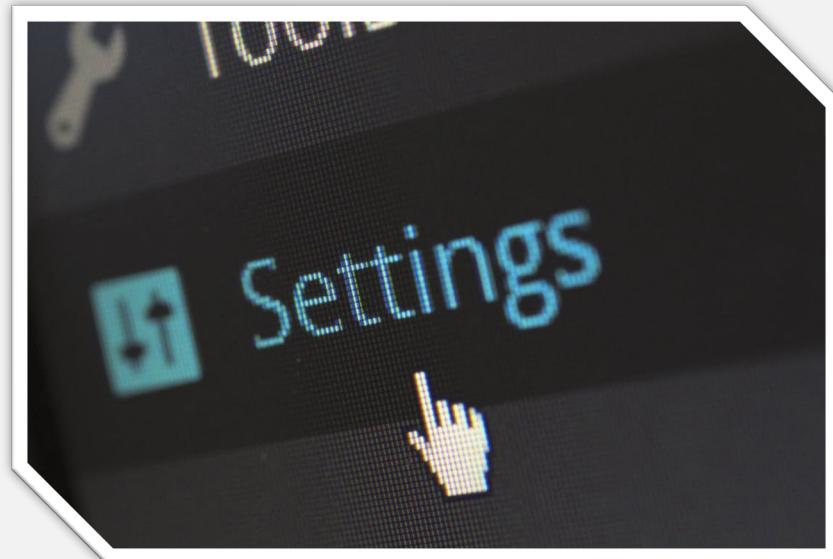
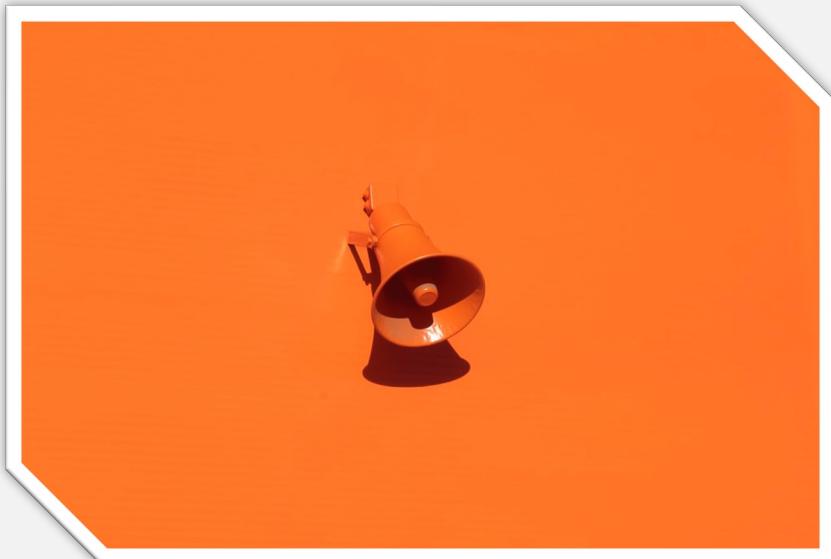
## ► ... as an inspiration?

- If you already work on a scraping/API project...



## ► ... as an inspiration?

- If you are an expert in working with web data ...



## ► Conclusion

- Webinar = primer
- Web (data) is here to stay (and grow)
- Important methodological tool for (early-career) researchers regardless of substantive or methodological focus.
- Many unexploited fields of gold! Potential for innovation!

# THANK YOU & TIME FOR MORE QUESTIONS

SLIDES AVAILABLE AT  
[HTTPS://WEB-SCRAPING.ORG](https://web-scraping.org)

BOEGERSHAUSEN@RSM.NL | @JOBOEGERSHAUSEN  
H.DATTA@TILBURGUNIVERSITY.EDU | @HANNESDATTAA  
ABHISHEK.BORAH@INSEAD.EDU  
ANDREW.STEPHEN@SBS.OX.AC.UK | @ANDREWSTEPHEN

<https://www.ama.org/journal-of-marketing/>

# All forthcoming articles



## JM Scholarly Insights:

Using Spatial Distance  
Strategically with Luxury and  
Popular Product Displays



## JM Webinars:

- For Marketers
- For Scholars
- Archive



## Special Issue:

Read the “Better Marketing  
for a Better World”  
Special Issue



## Editorial Teams

- Editors
- Associate Editors
- Editorial Review Board
- Advisory Board
- Meet JM at a Conference

## JM Articles

- Current Issue
- Articles in Advance
- Accepted Manuscripts
- Go to the Journal Archive
- AMA Member Access
- Awards

## Authors

- Editorial Mission
- Submission Guidelines
- AMA Journal Policies
- Manuscript Central
- Appeal Policy



New editors began handling new manuscripts on February 1 and will assume full operational control on July 1.

---

JOURNAL of

# Marketing



Hari Sridhar  
Editor-in-Chief Designate  
Texas A&M University



Cait Lamberton  
Editor Designate  
University of Pennsylvania



Detelina Marinova  
Editor Designate  
University of Missouri



Vanitha Swaminathan  
Editor Designate  
University of Pittsburgh