



Part 1: The art of web scraping



WHAT IS WEB SCRAPING?

Using a script to pretend to be a user and automatically collect data from websites.

Transform unstructured web data, typically HTML or JavaScript, into structured data that can be stored and analyzed in a format like CSV, JSON, or a pandas dataframe.

EXAMPLE USES OF WEB SCRAPING

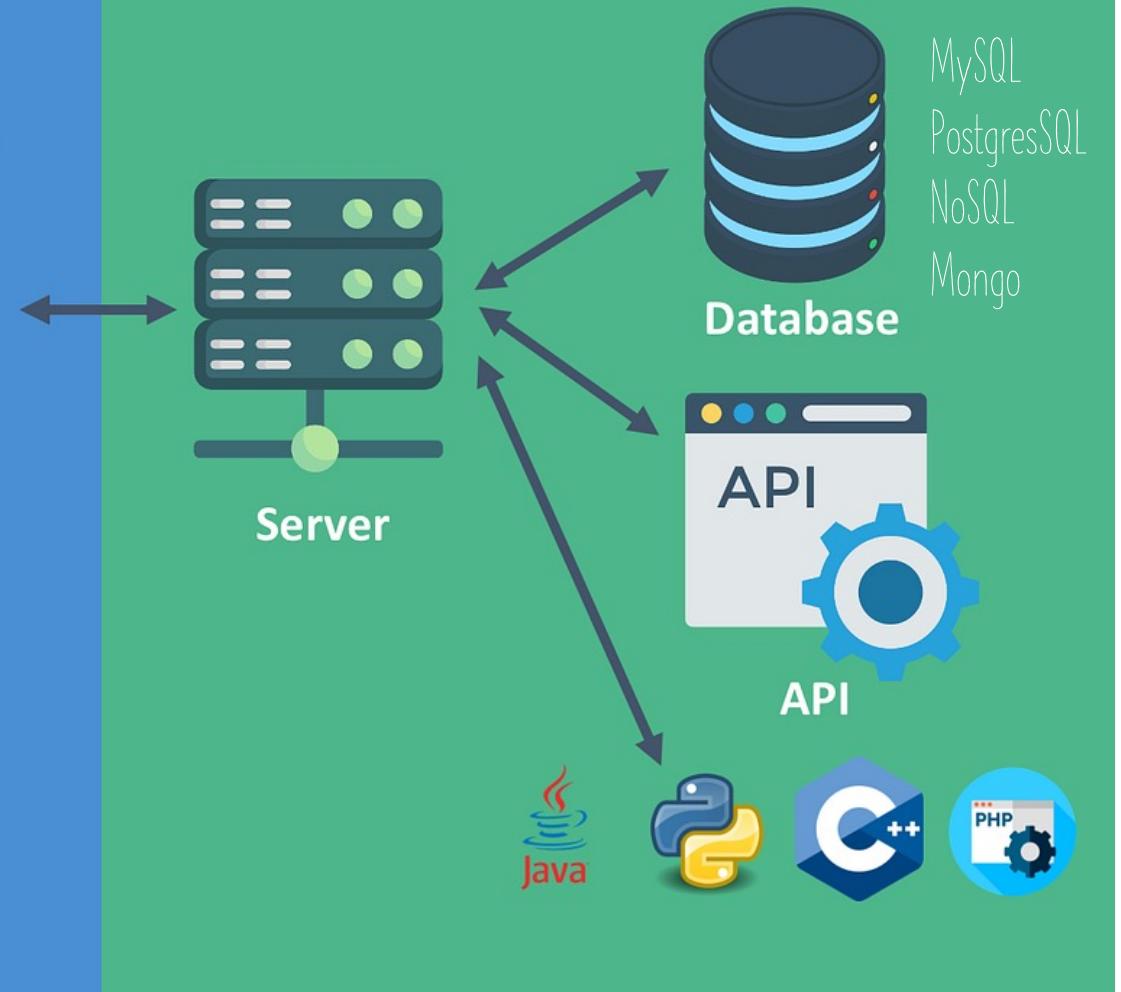
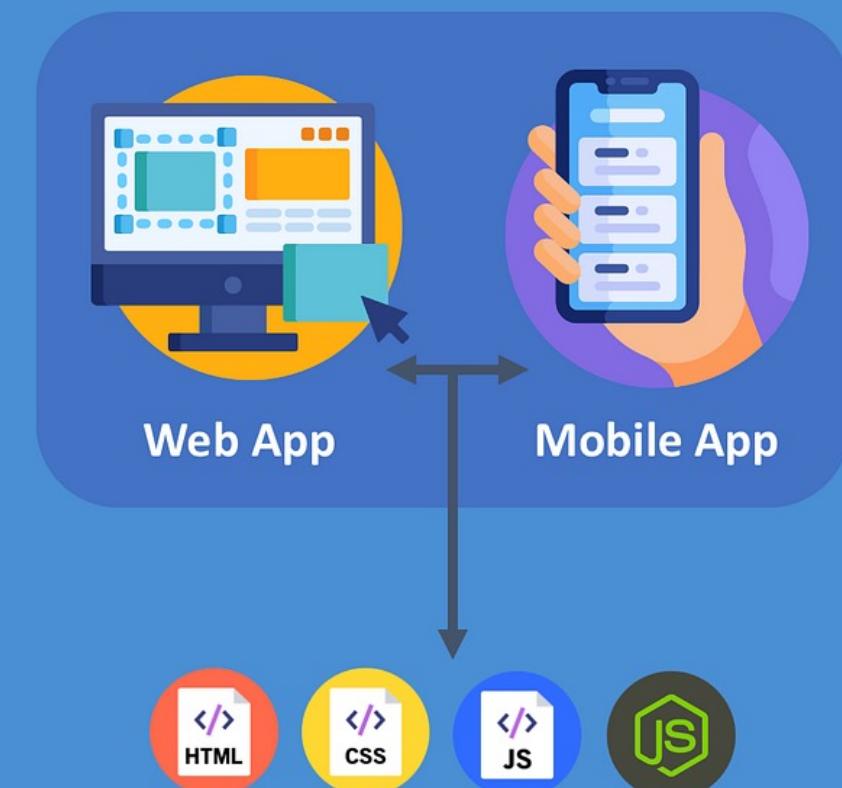
Observational data from websites and platform are among the most important data sources for CSS research:

- Collect social media posts from a platform
- Collecting listings and reviews from Airbnb
- Getting job listings from a job portal
- Getting news paper articles
- Get websites and the links between them
- ... *other examples?*

FRONT-END

Anatomy of a website

BACK-END



FRONT-END LANGUAGES

HTML (Hypertext Markup Language): HTML provides the basic structure of the website, using a system of tags and attributes to denote text, links, images, and other content.

CSS (Cascading Style Sheets): CSS is used for styling and visually formatting the HTML elements. It controls layout, colors, fonts, and even some animations.

JavaScript: JavaScript adds interactivity to web pages. It's a scripting language that enables dynamic content, control multimedia, animate images, and much more.

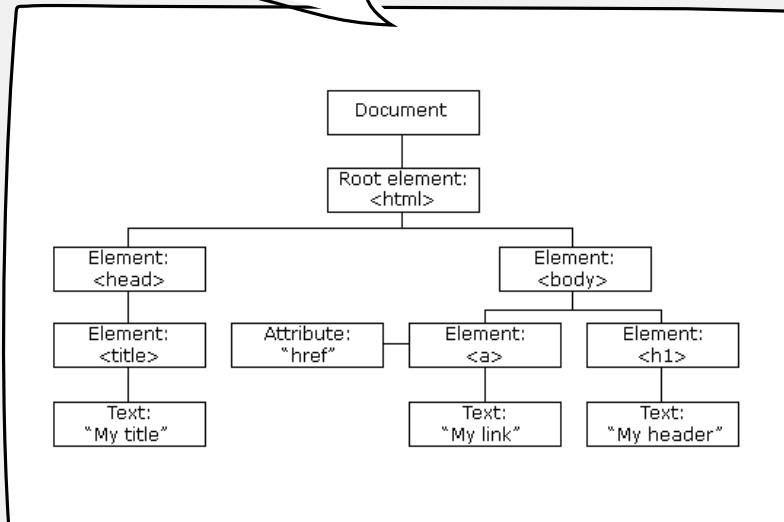


HTML USUALLY HAS OUR CONTENT

HTML is a series of elements or tags that structure content, allowing browsers to display text, links, images, and other resources.

- <p>This is a paragraph.</p>
- This is a link.
- Thisis a list
- <div class='box'>This is a object of the class 'box'</div>

THE DOCUMENT OBJECT MODEL (DOM)



- The DOM is a programming interface for web documents.
- Each element in an HTML document is represented as a node in the DOM tree, making it possible to traverse the hierarchy and access elements programmatically.
- When you use web scraping tools, you're essentially navigating the DOM to access specific parts of the page.

The screenshot shows the University of Amsterdam's Social Sciences landing page. At the top, there is a navigation bar with the university logo, a search bar, and language selection (EN). Below the header is a large banner featuring a group of people walking, with the text "Social Sciences" overlaid. To the left of the banner is a sidebar containing accessibility information and a "Vergelijk" button. The main content area includes a section about the Research Master's Social Sciences (RMSS) program, details about study mode and location, and a timeline for events. A testimonial from a postgraduate student, Kyriaki, is also present. At the bottom, there is a call-to-action for chat with students.

The screenshot shows the developer tools (Elements tab) with the page source code. The code is heavily annotated with class names and component identifiers, such as "p.lead", "c-programmepageheader", and "c-sectionmenu". The code is written in HTML and includes CSS styles and JavaScript snippets. The right side of the developer tools interface shows the corresponding CSS styles for the selected elements.

ROBOTS.TXT

A document that tells scrapers what the website owner is ok with you scraping or not.

Do you have to respect it?

google.com/robots.txt:

```
User-agent: *
Disallow: /search
Allow: /search/about
Allow: /search/static
Allow: /search/howsearchworks
Disallow: /sdch
Disallow: /groups
Disallow: /index.html?
Disallow: /?
Allow: /?hl=
Disallow: /?hl=*&
Allow: /?hl=*&gws_rd=ssl$
Disallow: /?hl=*&gws_rd=ssl
Allow: /?gws_rd=ssl$
Allow: /?pt1=true$
```

LIBRARIES FOR SCRAPING

We use programming libraries to scrape:

Requests: simple protocol to communicate with the website and get the HTML

Simple, fast, lightweight and best choice for static pages.

Selenium: runs an entire web browser and pretends to be a user clicking around

Slow and heavy, but very capable and works with dynamic websites

BeautifulSoup: to parse the HTML

WEB SCRAPING CHALLENGES

Web scraping can be difficult as websites try to prevent scraping.

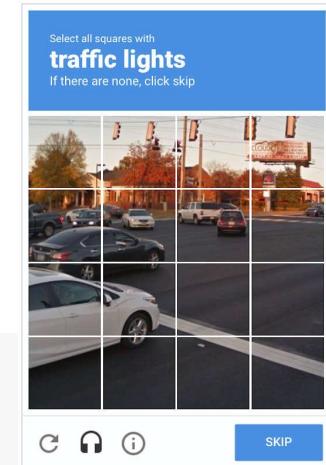
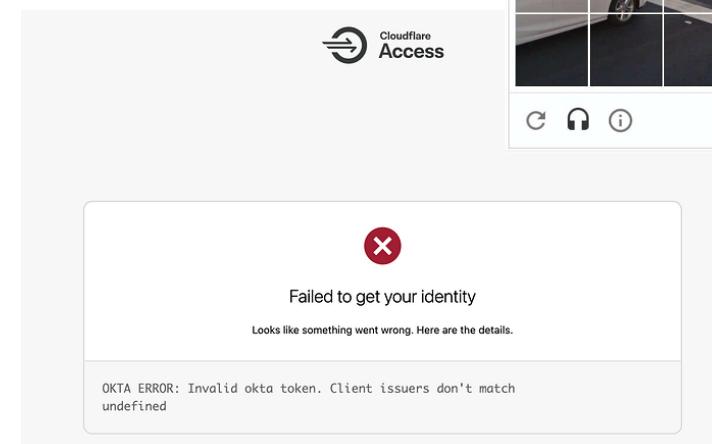
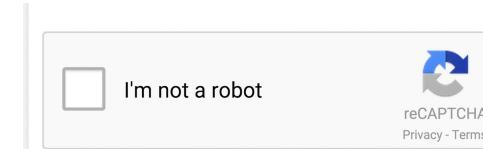
User-agent checks: You have to fake your identity

IP rate limit: Blocks too many calls

Cloudflare: US company. key infrastructure of the internet.
Prevents script access to websites.

Captcha: Designed to prevent scraping.

Handling cookies and sessions: Websites can be complex



SCRAPED DATA MUST BE CLEANED

- Remove duplicates
- Remove rows with missing data
- Fill in missing data
- Remove outliers
- Harmonize data
- Fix data types



The screenshot shows a web browser displaying the Stormfront.org website. The header features a circular logo with "WORLD WIDE" and "STORMFRONT.ORG" text, along with a "Stormfront" banner. The main navigation menu includes links for "Stormfront > General > History & Revisionism", "The Hitler We Loved And Why", "Donate", "Register", and "Blogs". A sidebar on the left shows user information for "Gott Mit Uns" (Alter Kämpfer, "Friend of Stormfront", Sustaining Member) with a profile picture of a classical statue and stats: Join Date: Jan 2023, Posts: 892. The main content area displays a post titled "The Hitler We Loved And Why" by Eric Thomson and Christof. It includes a PDF link (<https://der-fuehrer.org/bucher/engli...0and%20Why.pdf>) and a video link ([The Hitler We Loved and Why by Eric Thompson](#)). Below the post is a red book cover for "The Hitler We Loved". The right side of the page has a sidebar with "Every month is White history month" text, user login fields, and a search bar.

EXAMPLE: STORMFRONT.ORG

- Long-standing large Nazi community
- How do we understand 'echo chambers'? What actually takes place within them?
- Scrapped all comments and discussions

ROUTLEDGE, FOCUS
INTIMATE COMMUNITIES OF HATE
Why Social Media Fuels Far-Right Extremism
Anton Törnberg and Petter Törnberg

kaggle

Search

Sign In Register

PETTERTORNBERG · UPDATED 2 YEARS AGO

1 New Notebook Download (2 GB) ...

Stormfront.org - White Power forum posts 2001-2020

Hate speech from White Supremacists: 10M posts over nearly 20 years

Data Card Code (0) Discussion (0) Suggestions (0)

About Dataset

Stormfront.org is a white supremacist forum that has been active for 20 years, making it one of the longest-running forums on the Internet. This dataset contains the time-stamped text from all 10M posts made on Stormfront in the 2001-2020 period. The data can for instance be used to track the evolution of language as users engage with a far-right community. Each post is tagged by language (automatically identified). Usernames are anonymized, but a unique id is provided.

The file can be read from Python by:

```
df = pd.read_csv('stormfrontposts.csv.gz', compression='gzip')
```

More information about the dataset can be found in the following article:

Törnberg, P. Törnberg, A. 2022. Inside a White Power Echo Chamber: Why Fringe Digital Spaces are Polarizing Politics. *New Media & Society*.

If you use any of the provided material in your work, please cite us as follows:

Usability 7.06

License CC BY-SA 4.0

Expected update frequency Never

Tags Online Communities Text People and Society



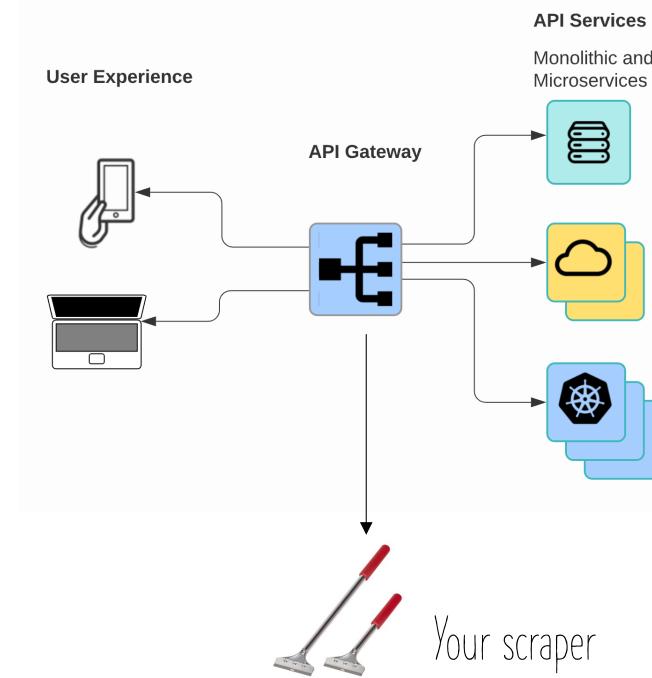
API (APPLICATION PROGRAMMING INTERFACE) DATA

- Many platforms provide API access, which allows you to get data.
- APIs send structured data, such as JSON, rather than HTML.
- This data is very easy to parse programmatically
- Comes with manuals! APIs often have libraries for making it even easier! E.g. Tweepy for Twitter.

```
{  
  "longitude": 47.60,  
  "latitude": 122.33,  
  "forecasts": [  
    {  
      "date": "2015-09-01",  
      "description": "sunny",  
      "maxTemp": 22,  
      "minTemp": 20,  
      "windSpeed": 12,  
      "danger": false  
    },  
    {  
      "date": "2015-09-02",  
      "description": "overcast",  
      "maxTemp": 21,  
      "minTemp": 17,  
      "windSpeed": 15,  
      "danger": false  
    },  
    {  
      "date": "2015-09-03",  
      "description": "raining",  
      "maxTemp": 20,  
      "minTemp": 18,  
      "windSpeed": 13,  
      "danger": false  
    }  
  ]  
}
```

GETTING DATA FROM APIs

- APIs are made to facilitate data exchange between applications.
- Most websites and apps run on APIs
- Can also be a strategy for scraping: direct access to internal APIs by reverse engineering
- Research data from Twitter, Reddit, Facebook, etc., generally come from their APIs



APIS BUILD INFRASTRUCTURAL POWER

Anne Helmond argues that platforms offer APIs to extend their ecosystems beyond their own websites, effectively turning third-party websites and apps into extensions of the platform.

This dependency gives platforms a significant amount of control over third-party developers and the wider platform ecosystem.



Helmond, A. (2015). The platformization of the web: Making web data platform ready. *Social media+ society*, 1(2)



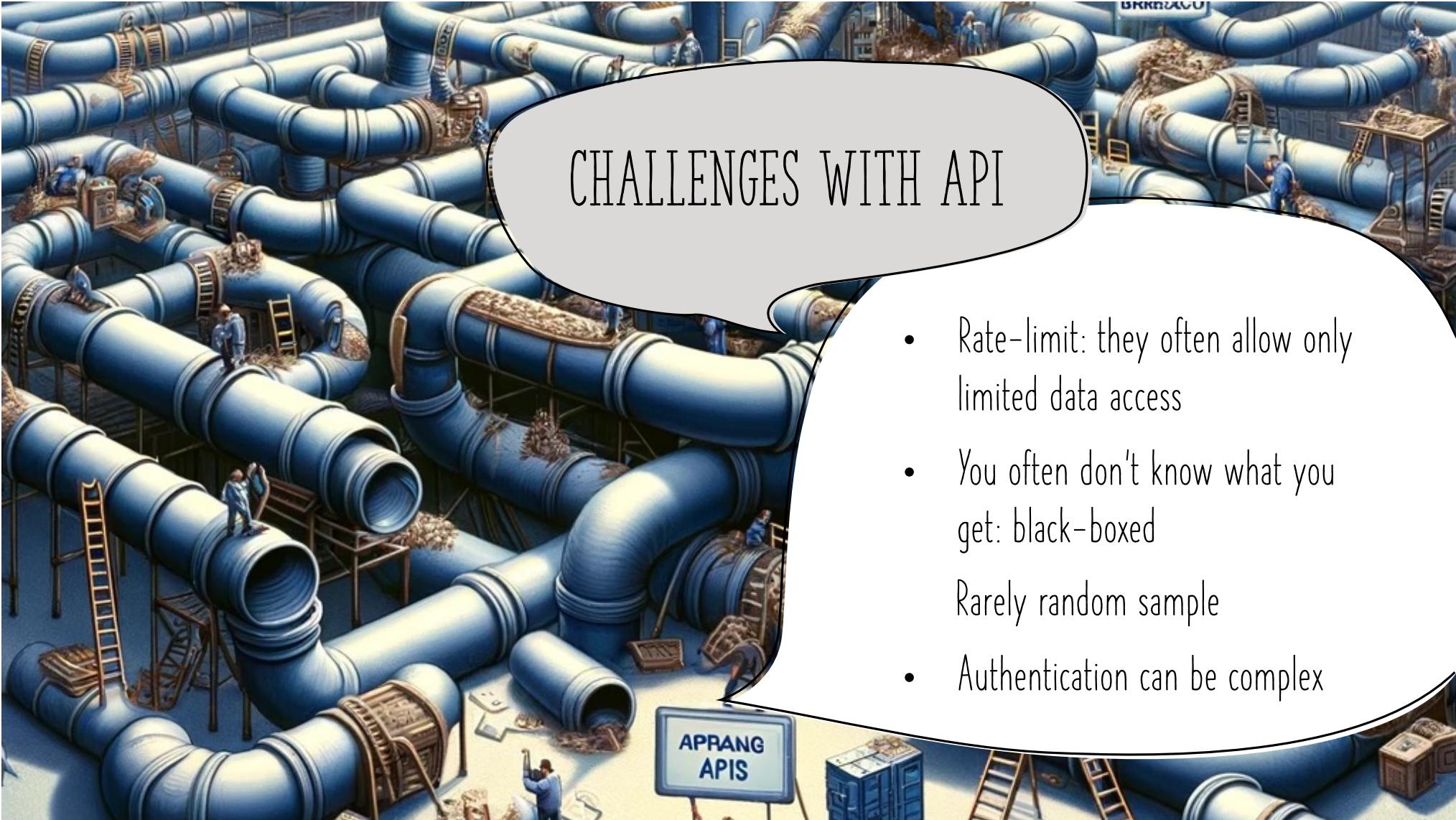
THE APICALYPSE

- Many APIs have been shutting down
- Instagram, Reddit, Twitter, Facebook...
- In part due to costs, privacy concerns, and in part to LLMs
- But DSA might come to the rescue

SOCIAL SCIENCE ANALYSIS AS A SERVICES

APIs also provide access to services: they allow code to interact with an online service.

- MTurk is based on an API: humans-as-a-service
- OpenAI ChatGPT is available through API.
- Many AI, NLP and image analysis tools are available through APIs
- We will learn using Perspective Toxicity API to analyze how toxic a comment is



CHALLENGES WITH API

- Rate-limit: they often allow only limited data access
- You often don't know what you get: black-boxed
Rarely random sample
- Authentication can be complex

RESPONSIBLE DATA COLLECTION PRACTICES

- When scraping, or using APIs, don't collect too quickly since it can overburden the server.
Pause between requests e.g., `time.sleep(0.5)`
- Don't collect data on people that you don't need (data minimization)
- Consider anonymizing the data before storing (anonymization)

```
import requests

url = "https://www.uva.nl/en/programmes/bachelors/computational-social-science/computational-social-science.html"

response = requests.get(url)

if response.status_code == 200:
    print("Here is the result:")
    print(f"{response.text[:300]} ...")
else:
    print(f"Failed to retrieve the webpage. Status code: {response.status_code}")
```

The screenshot shows the official website for the Bachelor's Computational Social Science programme at the University of Amsterdam. The page features a large banner image of a modern building complex. The main title 'Computational Social Science' is prominently displayed. Below the title, there is a brief description: 'Would you like to learn how to make the world a better place by using digital technology? Are you eager to learn all about data science? Do you have a hands on attitude and excellent team spirit?' There are also sections for 'Bachelor' and 'Compare programme'. At the bottom, there is a footer with links for 'Education', 'Research', 'News & Events', 'About the UvA', and 'Library'.

```
<!doctype html>
<html class="no-js" lang="en">
<head>
  <meta charset="utf-8"/>

  <title>Bachelor's Computational Social Science – University of Amsterdam</title>
  <link rel="canonical" href="https://www.uva.nl/en/programmes/bachelors/computational-social-science/computational-so ..."
```

This week's workshop!
Scraping and APIs



PART 2:
LAW & ETHICS
IN CSS



Cambridge
Analytica

facebook



ETHICS != LAW

- There are legal things that are not ethical.
- There are ethical things that are not legal.
- But we usually want to be both ethical and legal.

LAW

1. Terms of Service
2. Copyright Infringement
3. EU Database Directive in the EU
4. GDPR



TERMS OF SERVICE & SCRAPING

Can websites contractually limit scraping in their terms of use? Yes, they can.

But are those provisions enforceable? The legal theory behind contract enforceability is rather complex, but when talking about web scraping, the number one thing to check is the way how the contract was created.

- Browsewrap agreement: contracts that were concluded simply by visiting a website. Legal theory generally does not accept agreements of this type as valid.
- Clickwrap agreements: require an action by the user. "By continuing you agree to our Terms and Conditions". Clickwrap agreements are perfectly fine and fair contracts and the courts will readily enforce them.

Basically: If you need to log in or click "I agree", you have signed a ToS and need to follow it.

If not, *go wild!*

APIS AND TERMS OF SERVICE

- APIs often have ToS that you need to follow
- (In)famously, Twitter does not allow you to share any data to outsiders.
- Be mindful of the ToS when using APIs!

COPYRIGHT INFRINGEMENT

Scraping and reproducing content from a website could infringe on copyright laws.

In the EU, the scraping of copyrighted content is permitted by Article 3 and 4 of the Directive 2019/790 on copyright and related rights in the Digital Single Market (DSM Directive). The DSM Directive permits text and data mining, which means:

- "any automated analytical technique aimed at analysing text and data in digital form in order to generate information which includes, but is not limited to, patterns, trends and correlations";

Scraping of copyrighted content is only permitted for the purposes of generating information.

For example, you can scrape a webpage to extract prices from it, or books for natural language analysis, but you cannot scrape news articles and then republish them on your own website.

For scientific research, you can freely scrape almost anything!

Do not publish original works

EU DATABASE DIRECTIVE (DIRECTIVE 96/9/EC)

- Databases are protected - which can impact web scraping.
- The maker of a database has the right to prevent extraction of a substantial part of the contents of a database. For researchers, this means that scraping data from databases could potentially infringe on the rights of the database maker if a substantial part of the database is extracted without permission.
- The directive offers exemptions for teaching and scientific research, provided that you cite your source and that the use is for non-commercial purposes.

GDPR: GENERAL DATA PROTECTION REGULATION

- This is the big one.
- The GDPR requires that data collection practices comply with principles like lawfulness, fairness, and transparency.

Researchers must ensure explicit consent for the collection and use of personal data, provide clear information about data usage, and implement robust data protection measures. The GDPR also grants individuals the right to access, correct, and delete their data, affecting how researchers store and manage data. Additionally, the regulation's emphasis on data minimization and purpose limitation necessitates careful planning to justify the data collected and its use in research. Procedures and paperwork!

WHAT IS PERSONAL DATA ANYWAY?

GDPR only applies to personal data: "Personal data means any information relating to an identified or identifiable natural person"

If the data is not possible to link to a person, it does not fall under GDPR.

Direct identifying personal data refers to data which contains for example the name of the data subject, identity numbers, telephone numbers, email addresses, postal addresses or bank account numbers, etc.

Indirect identifying personal data refers to data that can be traced back to an individual when combined with other information. Since these data can be traced back to an individual (even though not in a direct way), the data are personal data and should be treated as such.

E.g.. *Name, surname, date of birth, address, social security number, passport number, national ID number, employment information, Contact details, phone number, email address, IP address, Facebook, Twitter, and other network handles, location either by address or GPS, shopping preferences, behavioral data, Video + audio recordings of people and biometric data. Special categories of personal data: sex, gender, and sexual orientation, racial or ethnic origin, religious beliefs, political opinions, medical records*

DE-IDENTIFICATION BY ANONYMIZATION / PSEUDONYMIZATION

If research data are anonymised/pseudonymised, to what extent do they still fall under the GDPR?

- Anonymisation is the process of removing all the direct and indirect information that can link the data to an individual.

After anonymisation, nobody can link the data to an individual anymore, i.e. no key files, pseudonyms, etc. are used. When data is fully anonymised, the data isn't personal data anymore and it doesn't fall in the scope of the GDPR.

- When the data is pseudonymised it is still possible to identify the person. Generally, a key file is used so that at least one person can link the data to an individual.

Even though most people cannot trace the data back the participant, the data is still personal data and therefore falls in the scope of the GDPR.

Re-identification is a constant risk!

PRINCIPLES OF GDPR

Data minimization: only data necessary should be processed, short storage period, limited accessibility

Privacy by design: any actions involving the processing of personal data is done with data protection and privacy in mind.

Privacy by default: all technical and organisational measures are taken to process data with the highest privacy protection.

As a researcher, you need to be proactive about data protection, and evidence (document) the steps you take to meet your obligations and protect people's rights.



GDPR PAPERWORK

1. Data management plan (DMP)
2. Privacy Protection Review (PPR)
3. Data Protection Impact Assessment (DPIA)
4. Processor agreement (PA)
5. Data exchange agreement (DEA)

DATA MANAGEMENT PLAN (DMP)

A Data Management Plan (DMP) is always required under GDPR.

- What data are you going to collect? What type of data or what file formats? How many?
- Where and how will you store your data? How will you provide back-ups?
- How are you going to organise and describe your data?
- Who gets access to your data? When? How are you going to manage access?
- What data will be archived when the project is finished? Where and for how long?
- Will the archived data be made available to others? When? Under what licence?
- Who owns the data? Who is responsible for the management of your data?
- Is there any funding to cover the costs of the implementation of the plan?

PRIVACY PROTECTION REVIEW (PPR)

PPRs are generally not mandated by law but are adopted by organizations as a best practice to ensure comprehensive privacy governance.

UvA requires this for any research projects.

- What type of data are you collecting?
- How will you store it? Anonymize? Share? Etc.

1. During your research AND the pre-research phase (i.e. screening, selection of participants), what (special) personal data relating to subjects are processed? *

For more information see infobox below.

- In this research no personal data are processed
- A. Name
- A. Contact details (e-mail, telephone number, or home address)
- A. IP-addresses
- A. Student number
- A. Citizen service number (BSN)
- A. Facial images/video and/or voice recordings
- A. Financial details (account number or creditcard number)
- A. Copies of passport or other identity documents
- A. Location data (GPS)
- A. Username
- B. Gender
- B. Age in years
- B. Date of Birth
- B. Nationality
- B. City or area of residence/postal code
- B. Unique identifier (e.g. number that can be used to re-identify a research participant)
- B. Number plates and device numbers (e.g. car, mobile phone IMSI)
- B. Language Background
- C. Personal data concerning convictions, criminal offenses or relevant safety precautions
- C. Personal data that point to ethnic background
- C. Personal data that point to political views
- C. Personal data that point to religious or philosophical beliefs
- C. Physical or mental health details
- C. Data on sexual behaviour or orientation
- C. Genetic details
- C. Biometric data for the purpose of uniquely identifying a person
- C. Trade union
- C. (increased risk of) Dyslexia and/or language disorders
- Other (please specify)

DATA PROTECTION IMPACT ASSESSMENT (DPIA)

Required by GDPR for high-risk projects:

- If you're using new technologies (e.g., AI)
- If you're tracking people's location or behavior
- If you're systematically monitoring a publicly accessible place on a large scale
- If you're processing personal data related to "racial or ethnic origin, political opinions, religious or philosophical beliefs, or trade union membership, and the processing of genetic data, biometric data for the purpose of uniquely identifying a natural person, data concerning health or data concerning a natural person's sex life or sexual orientation"
- If your data processing is used to make automated decisions about people that could have legal (or similarly significant) effects
- If you're processing children's data
- If the data you're processing could result in physical harm to the data subjects if it is leaked

Processor agreement (PA)

You draw up a processor agreement when you ask a third party to process personal data for your research, for example to transcribe audio files or to deliver an instrument for online surveys. This third party solely carries out your requests and does not take any decisions about the data themselves.

Data exchange agreement (DEA)

Sometimes you collaborate with other organisations or universities on a research project and you collectively decide how personal data will be collected, for what purpose and how they will be processed. In that case a data exchange agreement is required to establish who is responsible for what. If you process personal data, it needs to be clear whom a participant can contact with a remark or complaint.

EXAMPLE PROJECT: "TWITTER MOOD PREDICTS THE STOCK MARKET"

by Johan Bollen, Huina Mao, and Xiao-Jun Zeng in 2011.

This study utilized Twitter data to analyze the collective mood states of users and examined how these mood states could predict the rise and fall of the Dow Jones Industrial Average.

Data Analysis: The researchers applied text analysis techniques to Twitter feeds to extract mood states using a psychometric instrument known as the Profile of Mood States (POMS). They then correlated these mood states with stock market movements.

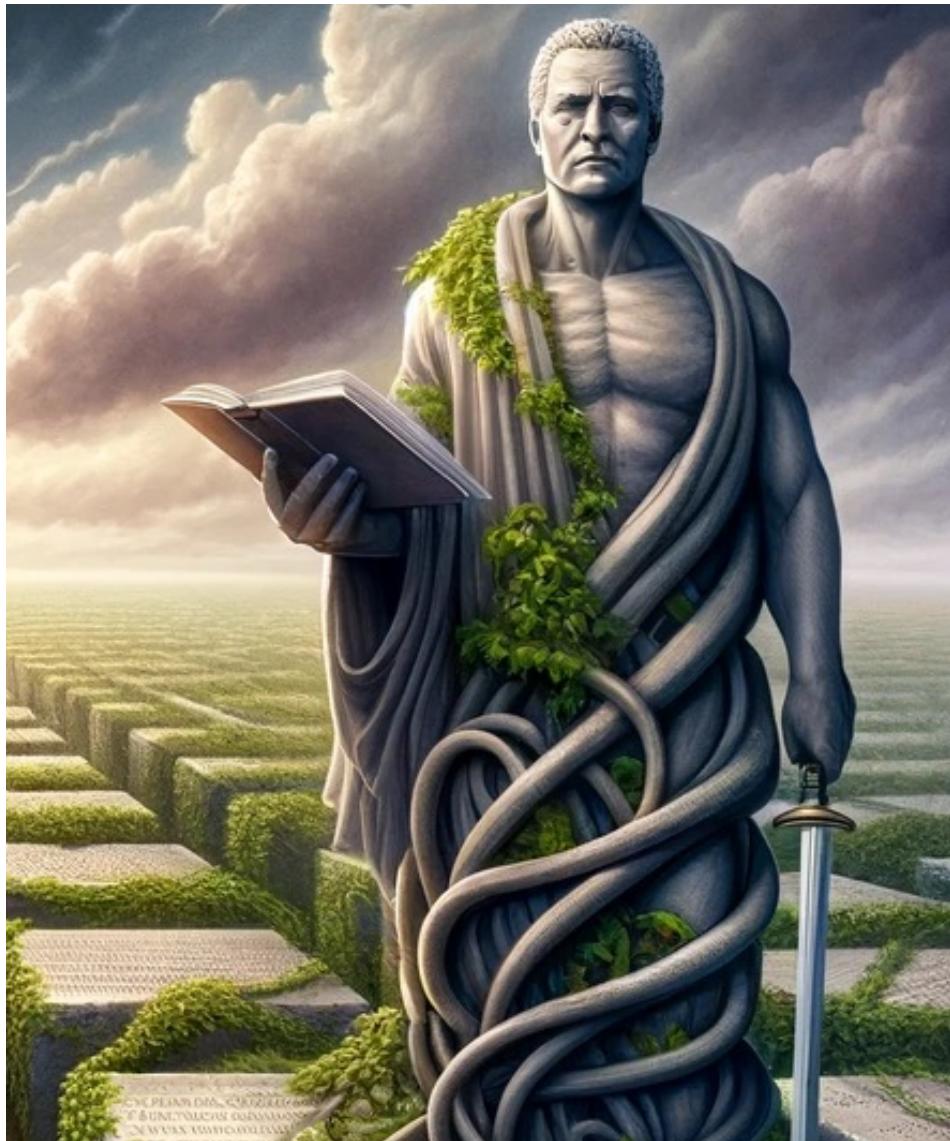
Findings: The study found a significant correlation between certain mood states captured through Twitter and the Dow Jones Industrial Average. Leading to the conclusion that social media could forecast economic indicators.

TASK: DISCUSS THROUGH LEGAL CONSIDERATIONS OF THE MOOD STUDY

What do you need to do to fulfill legal requirements for the study?

What data should you collect? How store?

What paperwork do you need for GDPR compliance?



ETHICS

1. Ethical frameworks
2. Ethical pluralism
3. AoIR guidelines
4. ERB application
5. Ethics section in paper

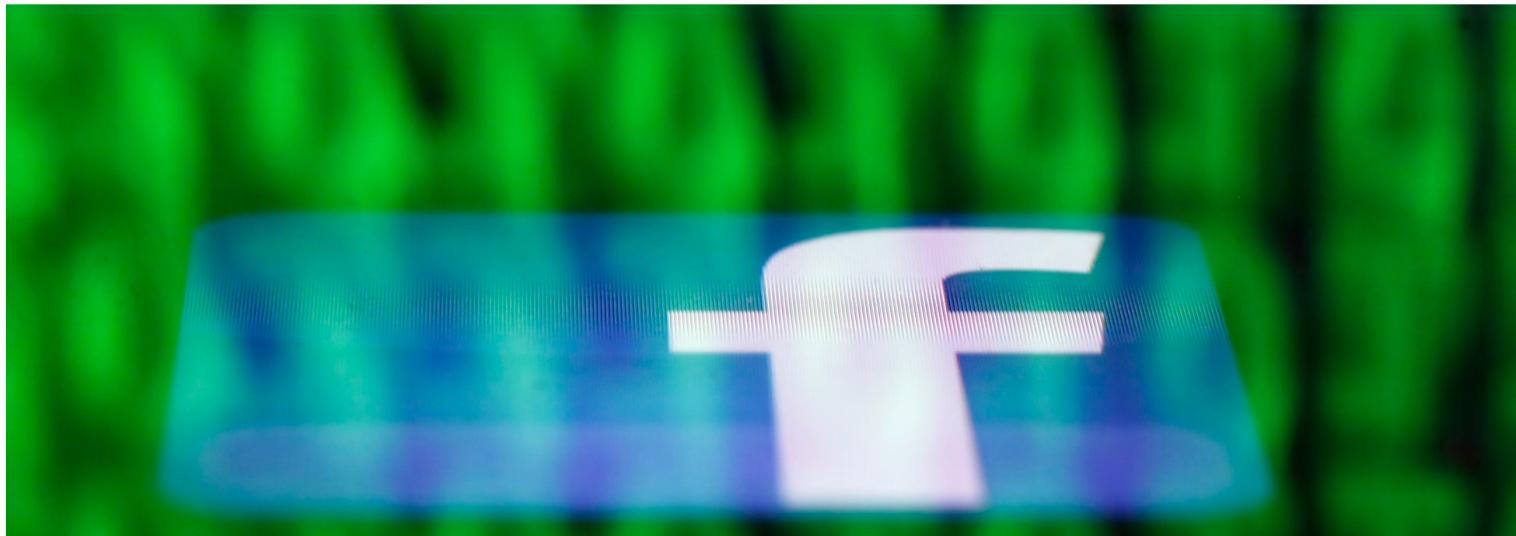
TECHNOLOGY

Everything We Know About Facebook's Secret Mood-Manipulation Experiment

It was probably legal. But was it ethical?

By Robinson Meyer

- 700,000 users were exposed either to negative posts or positive posts
- Measured effects on their own posts
- People got sadder
- Published in PNAS



RESEARCH ETHICS IN CSS

- Social scientists have long had IRB: come from the tradition of experiments.
 Stem from highly unethical experiments in the 1960s (e.g., Stanford Prison Experiment)
- Computer science has little traditions of ethics.

CSS is in the process of debating ethics. Still not well-established.

Varies depending on institution and region.

WHAT IS ETHICAL ANYWAY?

We can analyze our ethical dilemmas and challenges through a range of ethical frameworks:

Deontological: the rightness or wrongness of actions does not depend on their consequences but rather on whether they fulfill our duty or follow certain moral rules. Rule-based. Strong in Europe and Scandinavia.

Consequentialism: The outcomes and consequences of actions are what matter. (UK and US).

Utilitarianism: calculate the greater good for the collective and society in general.

Virtue-ethics: Actions are right if they are what a virtuous person would do in the same context.

Feminist ethics: Critically examine and address ethical theories and practices that perpetuate gender inequalities

Ethics of care: Prioritizes maintaining and fostering relationships, empathy, and caring for others, particularly those dependent on or vulnerable to one's actions.

So which one do we use?

ETHICAL PLURALISM AND CROSS-CULTURAL AWARENESS

Ethical pluralism: see the issue from multiple frameworks and discuss their competing views.

Ambiguity, uncertainty, and disagreement are inevitable.

Cross-cultural awareness: What are the expectations of the person being studied? What ethical framework do they follow?

If you're studying a distant culture, you must understand their views and expectations on ethics.

Research ethics is rarely black and white: there are competing interests that must be weighed

ASSOCIATION OF INTERNET RESEARCHERS (AOIR) ETHICS GUIDELINES

Key Principles

- Respect for Persons: Emphasizes the importance of respecting the autonomy, privacy, and dignity of individuals and communities involved in or affected by the research.
- Beneficence: Encourages actions that contribute to the welfare of research participants, aiming to maximize benefits and minimize harm.
- Justice: Ensures fair treatment and equitable distribution of the benefits and burdens of research across different groups, avoiding exploitation.

AIOR QUESTIONS TO GUIDE ETHICS

Discussion point: use these questions to think through the Twitter mood project!

1. Have participants been informed and consented? For publicly available data, have you considered the context in which the data was shared and whether the subjects anticipated their data being used for research? What are their expectations of privacy?
2. How does the project ensure the anonymity and privacy of individuals whose data is being analyzed? Are there measures in place to anonymize the data, especially when working with potentially sensitive information? Have you assessed the potential risks of re-identification of anonymized data?
3. Does the research aim to do good, and has it been designed to minimize potential harm to participants and communities involved?
4. How will the results be used, and could the publication of the research findings potentially harm individuals or communities? Have you considered the broader social implications of the research findings and how they contribute to the public good?
5. What measures are in place to secure the data during and after the research process?
6. Does the research involve groups or communities that might be considered vulnerable? If so, what additional protections are in place to safeguard these participants?
7. Are the research methods and ethical considerations clearly documented and accessible for review?

WRITING ETHICS REVIEW APPLICATION - ERB/IRB

- Ethics reviews are often mandatory for CSS research projects
- Offers opportunity to think through the ethics of your research
- The ERB is a floor, not a ceiling.

Task: discuss with your neighbor and go through the ethics review form, thinking of the Mood Stock Market project

WRITING AN ETHICS SECTION IN A PAPER

- Your paper should have an ethics section! Usually under "Methods & data"
- Describe your ethics process. Show that you've considered ethical implications.
- *Make an argument* for your work being ethical based on an existing framework or guidelines.
- *Very brief!*

We follow the ethical guidelines for Internet research provided by the Association of Internet Researchers ([Franzke et al., 2020](#)) and by the [British Sociological Association \(BSA, 2017\)](#). To ensure anonymity for the members of the community, usernames were deleted when collecting the corpus. Since the analytical focus is on broader discursive patterns, it is not possible to identify individual users from the results here presented. The project has been assessed and approved by the Regional Ethical Review Board.

DISCUSSION QUESTION: HARMING NAZIS



We generally want to avoid harming our research subjects.

But when studying Stormfront, the question was how fast to collect the data:
fast scraping means getting data in weeks instead of months, but may slow down their website.

But their website fuels neo-Nazi violence and has been linked to 100s of murders.

If we're harming self-proclaimed Nazis, does it make it more okay?

What would different ethical frameworks say?

BEFORE NEXT WORKSHOP

Sign up for Perspective API:

<https://perspectiveapi.com/>

Get API key for Perspective API and YouTube API!