# *Introduction to Data Science*

# Data Collection Methods & Ethics (Lecture 1)

Dr. Oul Han

Special thanks to:
Juhi Kulshrestha  (http://www.juhikulshrestha.com/)
Aniko Hannak (http://ancsahannak.me)

# What is data? One example of MANY

**Qualitative and small**

- Observation of individual cases
- Explaining, characterizing
- *Example*:
  - How many friends do overseas students have?
  - How to they interact with them in the holidays?

# What is data? Another example of MANY

**Quantitative and large**

- Observation of categories

- Counting, sorting by features

- *Example*:

  - How many people in Koblenz are registered as self-employed?

  - How many of them are merchants, manufacturers, or service providers?

# Data from online platforms

- Data is everywhere!
- The only limitation is your imagination
  - *And the terms of service*



(iconimage / stock.adobe.com)

# Data from online platforms - Pros & Cons

- Larger & cheaper than surveys or field experiments
- Examine human interactions in their natural environments
- Immediate feedback after external events

# Data from online platforms - Pros & Cons

- Larger & cheaper than surveys or field experiments
- Examine human interactions in their natural environments
- Immediate feedback after external events

- Big data isn't more representative or of better quality
- Data-driven analysis, or over-simplified representation?
- Your conclusions may be wrong about the world (external validity)

# Example: WeST Facebook group

**Representative?**

- "Does the WeST Facebook group represent the opinion of all WeST students?"

**Wrong about the world?**

- „A high number of likes in the WeST Facebook group means that all WeST students are highly satisfied"

# In this class

Collecting data:
Ethical & legal
considerations

Overview of data
collection tools:
APIs
Web scraping
Browser automation

Sharing data:
Ethical & legal
considerations

Image source: https://www.flaticon.com/

Am I harming the users?

Am I harming the platform?

# Collecting Data: Ethical and Legal

# Am I harming the users?

- Interference through experiment
- Manipulating the user's behavior without consent



Facebook reveals news feed experiment to control emotions

Protests over secret study involving 689,000 users in whic[h] were moved to influence moods



I Agree ☐

# Am I harming the users?

- Collecting personal or sensitive information



Image reused from CESSDA ERIC Expert Tour Guide on Data Management which is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License

# Am I harming the users?

- Personal information
  - Any data that can be used to identify living (or deceased) individuals
- Sensitive information
  - name or date of birth
  - person's origin,
  - political opinion,
  - religious beliefs,
  - health,
  - trade union membership,
  - sexual orientation …
- If sample size is small
  - number of children a person has
  - shoe size …



Image reused from CESSDA ERIC Expert Tour Guide on Data Management which is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License

# Am I harming the users?

- Collect anonymously
- Use consent forms
- Anonymize data
- Securely store, access, transfer

# EU General Data Protection Regulation (GDPR)[1]

- **Transparency**
  - processing personal data "lawfully, fairly and in a transparent manner"

- **Data Minimization**
  - data use shall be limited to the purpose of the respective research

- **Accuracy**
  - inaccurate data must be "erased or rectified without delay"

- **Integrity and Confidentiality**
  - data must be protected by appropriate security measures (technical and organizational)

1 https://www.fosteropenscience.eu/learning/data-protection-and-ethics/#/id/5ace27ca8ee5d6920ab94c13

# Am I harming the platform?

- Interference with algorithms on the site by clicking or searching
  - Bots
  - Search words hijacking



THE STRATEGIST

ASPI
AUSTRALIAN
STRATEGIC
POLICY
INSTITUTE

SUBM

Bots foment political polarisation through social media

7 Aug 2020 | Elise Thomas

SHARE

# Am I harming the platform?

- Click fraud
  - Advertisers pay for every click or impression[1]

AdFisher may have cost advertisers a small sum of money. AdFisher never clicked on any ads to avoid per click fees, which can run over $4 [34]. Its experiments may have caused per-impression fees, which run about $0.00069 [35]. In the billion dollar ad industry, its total effect was about $400.

1 Datta, A., Tschantz, M. C., & Datta, A. (2015). Automated experiments on ad privacy settings: A tale of opacity, choice, and discrimination. *Proceedings on privacy enhancing technologies, 2015*(1), 92-112.

# Am I harming the platform?

- Load balance
  - Check if an API exists or if data are available for download
  - Do not overload the servers such that their service is affected!

**Web Science**

# Terms of Service



1. You will not provide any false personal information on Facebook, or create an account for anyone other than yourself without permission.
2. You will not create more than one personal account.
7. If you collect information from users, you will: obtain their consent, make it clear you (and not Facebook) are the one collecting their information, and post a privacy policy explaining what information you collect and how you will use it.

**Twitter Terms of Service**                    Download: Twitter User Agreement

You may not do any of the following while accessing or using the Services: (i) access, tamper with, or use non-public areas of the Services, Twitter's computer systems, or the technical delivery systems of Twitter's providers; (iii) access or search or attempt to access or search the Services by any means (automated or otherwise) other than through our currently available, published interfaces that are provided by Twitter (and only pursuant to the applicable terms and conditions), unless you have been specifically allowed to do so in a separate agreement with Twitter (NOTE: crawling the Services is permissible if done in accordance with the provisions of the robots.txt file, however, scraping the Services without the prior consent of Twitter is expressly prohibited);

# Terms of Service



1. You will not provide any false personal information on Facebook, or create an account for anyone other than yourself without permission.
2. You will not create more than one personal account.
7. If you collect information from users, you will: obtain their consent, make it clear you (and not Facebook) are the one collecting their information, and post a privacy policy explaining what information you collect and how you will use it.



You may not do any of the following while accessing or using the Services: (i) access, tamper with, or use non-public areas of the Services, Twitter's computer systems, or the technical delivery systems of Twitter's providers; (iii) access or search or attempt to access or search the Services by any means (automated or otherwise) other than through our currently available, published interfaces that are provided by Twitter (and only pursuant to the applicable terms and conditions), unless you have been specifically allowed to do                                                                                          h the provisions of tl

Robots.txt

Rate limits

```
https://varvy.com/robots.txt

User-agent: *
Disallow: /folder/
Disallow: /file.html
Disallow: /image.png
```

# Exceptions for research

- Computer Fraud and Abuse Act (CFAA) (US)
  - Crawling is considered illegal for business practices such as discrimination

1 https://www.fosteropenscience.eu/learning/data-protection-and-ethics/#/id/5ace27ca8ee5d6920ab94c13

# Exceptions for research

- Computer Fraud and Abuse Act (CFAA) (US)
  - Crawling is considered illegal for business practices such as discrimination
- GDPR research exemptions (EU)[1]
  - If for "the public interest, scientific or historical research purposes or statistical purposes" (Art. 5.1 2016/679/EU)
  - If "the data subject has given consent to the processing of his or her personal data for one or more specific purposes" (Art. 6.1 2016/679/EU)

1 https://www.fosteropenscience.eu/learning/data-protection-and-ethics/#/id/5ace27ca8ee5d6920ab94c13

# Data collection methods

- APIs
- Web scraping
- Browser automation
- Personalized data collection

# APIs

# APIs

- API - Application Programming Interface
  - set of http requests that returns structured data (JSON, XML)
- Two types
  - Restful APIs
    - queries for static information at current moment
      - user profiles, posts, …
  - Streaming APIs
    - changes in users' data in real time
      - new tweets, weather alerts, …

# Twitter API example

# Twitter API example



```
{
  "created_at": "Wed Oct 10 20:19:24 +0000 2018",
  "id": 1050118621198921728,
  "id_str": "1050118621198921728",
  "text": "To make room for more expression, we will
  now count all emojis as equal—including those with
  gender and skin t… https://t.co/MkGjXf9aXm",
  "user": {}, "entities": {}
}
```

# Twitter API example

# GET followers/ids





```
{
  "created_at": "Wed Oct 10 20:19:24 +0000 2018",
  "id": 1050118621198921728,
  "id_str": "1050118621198921728",
  "text": "To make room for more expression, we will
  now count all emojis as equal—including those with
  gender and skin t… https://t.co/MkGjXf9aXm",
  "user": {}, "entities": {}
}
```

# APIs - Pros

- ToS compliant
- Easy to use, replicate
- Well-formatted data

# APIs - Cons (I)

- ToS compliant
- Easy to use, replicate
- Well-formatted data

- Authentication
  - IP based
  - User tokens based
- Rate limits
  - Restricted number of API call per user/IP address
  - x% of data

# APIs - Cons (II)

- Possible incompleteness of data
  - Missing info, such as images

"profile_background_image_url"  :   "http://abs.twimg.com/images/themes/theme7/bg.gif"  ,
"profile_background_image_url_https"  :   "https://abs.twimg.com/images/themes/theme7/bg.gif"  ,
"profile_background_tile"  :   false  ,
"profile_image_url"  :   "http://pbs.twimg.com/profile_images/448483168580947968/pL4ejHy4_normal.jpeg"  ,
"profile_image_url_https"  :   "https://pbs.twimg.com/profile_images/448483168580947968/pL4ejHy4_normal.jpeg"  ,
"profile_banner_url"  :   "https://pbs.twimg.com/profile_banners/12/1347981542"  ,

# APIs - Cons (III)

- Possible unknown biases in data
  - Unclear what the platform provider may be giving



Morstatter et al, 2013, ICWSM

# Web scraping

# Web scraping

- Extracting data from source code of website

```
BASH: curl https://en.wikipedia.org/wiki/Data_science > wiki_ds.html
Python Requests:
requests.get("https://en.wikipedia.org/wiki/Data_science")
```

# Web scraping

- Extracting data from source code of website

```
BASH: curl https://en.wikipedia.org/wiki/Data_science > wiki_ds.html
Python Requests:
requests.get("https://en.wikipedia.org/wiki/Data_science")
```

# Web scraping

- Tool to parse HTML code: beautifulsoup
- More in exercise hour!

# Web scraping - Pros & cons

- Easy to set up
- Easy to parallelize

# Web scraping - Pros & cons

- Easy to set up
- Easy to parallelize

- Not ToS compliant
- No ajax, no javascript, no images
- Parsing needs to be updated every time the platform makes a change

# Browser automation

# Browser automation

- Automate browser to scrape dynamically rendered webpages
- Could be used to:
  - fill forms
  - enter text
  - scroll
  - click on buttons

# Browser automation - Pros

- Mimics human browsing

- Loads ajax, images, etc.

- Pops open a browser so you can check if it works - easy to debug

- Can design the flow of events : log-in, search, etc.

# Browser automation - Cons

- Not ToS compliant
- Need to parse content
- Difficult to scale
- Unpredictable bugs
- Platform may throw a pop up (e.g. a captcha) if you collect too much

# Headless browser

- Tool - PhantomJS
- Headless => does not pop open a browser  window

# Headless browser - Pros

- Mimics human browsing

- Loads ajax, images etc

- Can design the flow of events: log-in, search etc.

- Easy to parallelize

- Less memory intensive

# Headless browser - Cons

- Not ToS compliant
- Need to parse content
- Unpredictable bugs
- Messy code
- Hard to debug - no browser window

# Other options for collecting data?

# Summary

| | Sample tools | Pros | Cons |
|---|---|---|---|
| **API** | | ToS compliant, easy to use | possible bias, incompleteness Auth and rate limits |
| **scraping static pages** | Curl, python requests | easy to use, parallelizable | no ajax, no images, no javascript you have to parse data |
| **Automated Browser** | Selenium | mimics real humans, possible to log-in, design flow of events | not possible to parallelize, unpredictable bugs (pop-ups, ads) you have to parse data |
| **Headless Browser Implementation** | phantomJS, selenium | fast, parallelizable | hard to debug since there is no physical browser window you have to parse data |

# How to decide?



Image source: http://pablobarbera.com

# Case study:  Collecting Personalized data

# Collecting personalized data (I)

- Volunteered data - crowdsourced
- Browser plugin

# Collecting personalized data (I)

- Volunteered data - crowdsourced
- Browser plugin

| | |
|---|---|
| • Easier to scale<br>• Leverage real accounts<br>• Does not break ToS | • Must recruit volunteers<br>• Sampling bias<br>• Less control |

# Collecting personalized data (II)

- Controlled experiments
    - Create test accounts with preferred characteristics
    - Collect the personalized recs or search results for these accounts.

# Collecting personalized data (II)

- Controlled experiments
    - Create test accounts with preferred characteristics
    - Collect the personalized recs or search results for these accounts.

| | |
|---|---|
| • Easier to measure impact of features<br>• Clean data<br>• Can mitigate biases | • Synthetic data<br>• Breaks ToS<br>• Harder to scale |

Can I anonymize the data?

Am I violating Terms of Service?

# Sharing data:  Ethical & legal considerations

# Sharing data publicly

- Anonymize the users
- Even then, users can be fingerprinted if:
  - Sample size is small
  - there are outliers or minorities
  - it can be merged with other available data sets
  - etc.

# Sharing data publicly

- Anonymize the users
- Even then, users can be fingerprinted if:
  - Sample size is small
  - there are outliers or minorities
  - it can be merged with other available data sets
  - etc.

- K-anonymization: Data is said to have k-anonymity if the information for each person contained in the dataset cannot be distinguished from at least k-1 individuals whose information also appears in the release

# Sharing data publicly



A 3-diverse patient table

| Zipcode | Age | Salary | Disease |
|---|---|---|---|
| 476** | 2* | 20K | Gastric Ulcer |
| 476** | 2* | 30K | Gastritis |
| 476** | 2* | 40K | Stomach Cancer |
| 4790* | ≥40 | 50K | Gastritis |
| 4790* | ≥40 | 100K | Flu |
| 4790* | ≥40 | 70K | Bronchitis |
| 476** | 3* | 60K | Bronchitis |
| 476** | 3* | 80K | Pneumonia |
| 476** | 3* | 90K | Stomach Cancer |

| Bob | |
|---|---|
| **Zip** | **Age** |
| 47678 | 27 |

Source: https://elf11.github.io/2017/04/22/kanonymity.html

- K-anonymization: Data is said to have k-anonymity if the information for each person contained in the dataset cannot be distinguished from at least k-1 individuals whose information also appears in the release

# Sharing data publicly

- Anonymize the users
- Securely store, control access & transfer your data

# Sharing data publicly

- Anonymize the users
- Securely store, control access & transfer your data
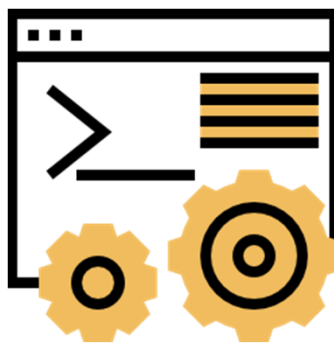- Read ToS carefully

# Sharing data publicly

- Anonymize
  - protect against fingerprinting - k-anonymity
- Securely store, control access & transfer your data
- Read TOS carefully
  - allowed to share tweet ids, and a sample of tweet objects for non commercial use
- Do not share copyrighted content
  - Just because you can download it, doesn't mean you can share someone else's intellectual property

# In this class… we learnt

**Collecting data:
Ethical & legal
considerations**

**Overview of data
collection tools:
APIs
Web scraping
Browser automation**

**Sharing data:
Ethical & legal
considerations**

Image source: https://www.flaticon.com/

# How would you collect data for…

- Which method would you use?
  - Do you need more information to decide?
- What are the method's pros and cons for this scenario?
- What are the ethical and legal issues of this method?

# How would you collect data for…

- What politicians post on Twitter

# How would you collect data for…

- Wikipedia data about scientists from different countries

# How would you collect data for…

- News articles from nytimes.com

# How would you collect data for…

- Multiple pages of search results on Twitter/ Google

# How would you collect data for…

- YouTube recommendations

End of Lecture 1

# Questions:
## Tutorials 1. TA Office hours (book slot)