# DIGITAL DEMOGRAPHY: ANALYZING WEB AND SOCIAL MEDIA DATA

## DAY 1 – INTRODUCTION, INTERNET AND SCRAPING
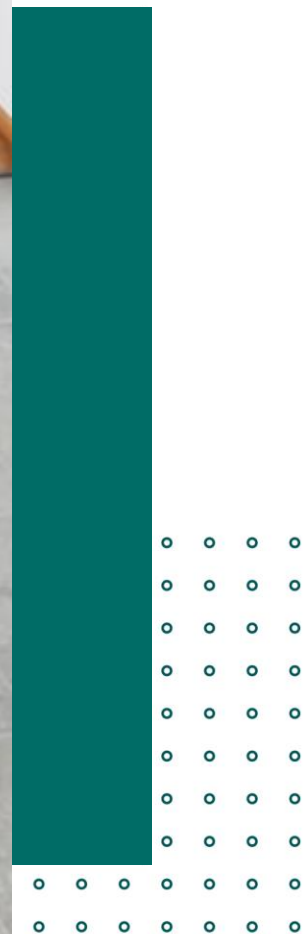
### EDSD NOVEMBER 2022

Tom Theile

Lab of digital and computational demography,

Max-Planck-Institute for Demographic Research, Rostock

# LAB OF DIGITAL AND COMPUTATIONAL DEMOGRAPHY

# AGENDA AND SCHEDULE OF THE WEEK

- 9:30 – 13:00 Monday to Thursday

- Friday: assignment

- Block 1: Internet, Webscraping, APIs

- Block 2: Digital Demography

- Block 3: Social Media Data

- Extra: Digital Datasets

# FINAL ASSIGNMENT

- A small number of short tasks

- To be done and submit until Friday afternoon 6pm

- Everyone has to submit the answers by email

- Email has to contain your full name!

- No group work

- More information on Thursday

# WHO ARE YOU?

- 4 Things in common - related to migration, mortality, fertility and digital data

## AGENDA

- What is the Internet? How do websites work?

- How to scrape data from websites?

- What to do with such data?

# WHAT HAPPENS WHEN WE VISIT A WEBPAGE?

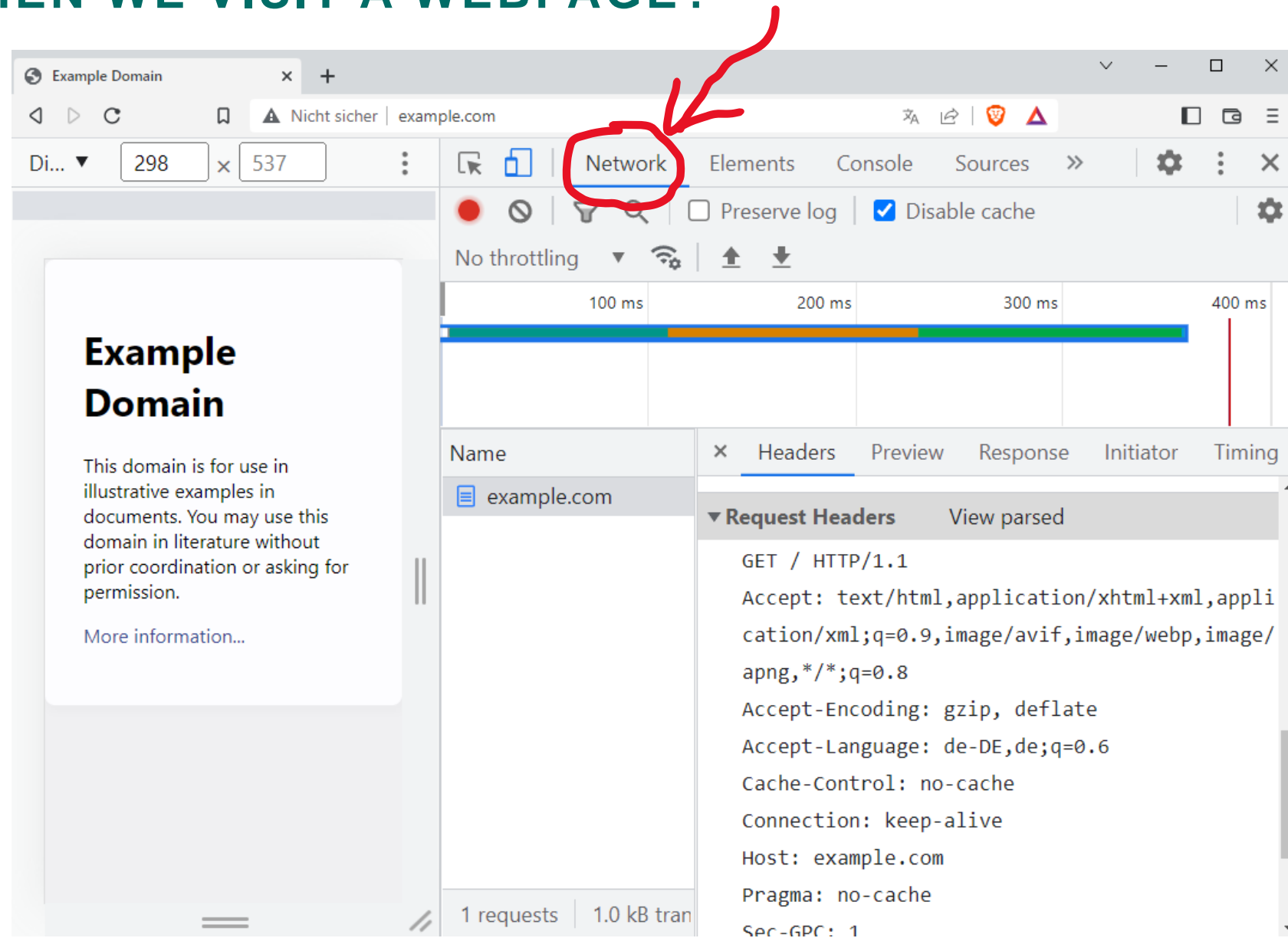Open the DevTools in

your Browser:


Windows or Linux:

F12 or

"CTRL + shift + i"


Mac:

"Fn + F12" or

"Cmd + Option + I"

Open the Network panel!

# EXAMPLE.COM HTML

- If you are not familiar with the basics of HTML, this is a good resource to learn the basics or refresh your knowledge:

- https://developer.mozilla.org/en-US/docs/Learn/Getting_started_with_the_web/HTML_basics

```html
1   <!doctype html>
2   <html>
3   <head>
4       <title>Example Domain</title>
5
6       <meta charset="utf-8" />
7       <meta http-equiv="Content-type" content="text/html; charset=utf-8" />
8       <meta name="viewport" content="width=device-width, initial-scale=1" />
9       <style type="text/css">
10      body {
11          background-color: #f0f0f2;
12          margin: 0;
13          padding: 0;
14          font-family: -apple-system, system-ui, BlinkMacSystemFont, "Segoe UI", "Open Sa
15
16      }
17      div {
18          width: 600px;

34      }
35      </style>
36  </head>
37
38  <body>
39  <div>
40      <h1>Example Domain</h1>
41      <p>This domain is for use in illustrative examples in documents. You may use this
42      domain in literature without prior coordination or asking for permission.</p>
43      <p><a href="https://www.iana.org/domains/example">More information...</a></p>
44  </div>
45  </body>
46  </html>
47
```

# CSS - CASCADING STYLE SHEETS

# CSS - CASCADING STYLE SHEETS

| Selector | Example | Example description |
|---|---|---|
| *#id* | #firstname | Selects the element with id="firstname" |
| *.class* | .intro | Selects all elements with class="intro" |
| *element.class* | p.intro | Selects only <p> elements with class="intro" |
| *\** | * | Selects all elements |
| *Element* | p | Selects all <p> elements |
| *element,element,..* | div, p | Selects all <div> elements and all <p> elements |
| *element element* | div p | Selects all <p> elements inside <div> elements |
| *element>element* | div > p | Selects all <p> elements where the parent is a <div> element |

- **https://www.w3schools.com/cssref/css_selectors.php**

- **https://developer.mozilla.org/en-US/docs/Web/CSS/CSS_Selectors**

# A MORE COMPLICATED WEBSITE

open https://eaps.nl/edsd

**Mpidr.de** 195.37.34.73

**Example.com** 195.37.34.73

Phone 192.168.2.18

Router 94.179.18.10

Laptop 192.168.2.13

Router 192.168.2.1

**Mpidr.de** 195.37.34.73

**Example.com** 195.37.34.73

Phone 192.168.2.18

Router 94.179.18.10

Laptop 192.168.2.13

Router 192.168.2.1

Source: https://www.youtube.com/watch?v=vv4y_uOneC

# TRACEROUTE

- You can list all the routers on this journey with traceroute (tracert on Windows)

- `tracert example.com`

- Visualize it with https://stefansundin.github.io/traceroute-mapper/

# PART 2 – SURFING THE WEB WITH R - WEBSCRAPING

Install Selector Gadget from https://selectorgadget.com/

- Show/enable the bookmark-toolbar in your browser

- Drag link to the bookmark toolbar

Alternatively you can install the chrome extension

Open a website, open SelectorGadget and click on the text you want to select.

# PART 2 – SURFING THE WEB WITH R - WEBSCRAPING

Please open the script 01_webscraping.R with Rstudio

https://nextcloud.demogr.mpg.de/s/iDQaEGNDqsd5XtQ

https://www.tidyverse.org/

https://rvest.tidyverse.org/

**Discussion: Is web scraping legal?**

For most owners of webpages, it is fine to scrape their site, as long as you "behave"

"Behaving" means:

- You don't induce interruptions or unreasonable costs to their service by scraping too fast or too much (good rule of thumb: only scrape one domain every 2 seconds)

- You don't use the data to the disadvantage of the scraped site. Since content is often copy-righted, you are mostly not allowed to share your scraped data publicly. (Exemption: you alter or aggregate the data enough)

- You respect their robots.txt

Some companies think it is not okay to scrape their webpage - but you are probably still legally allowed to scrape it!

- A US court ruled, that LinkedIn has to remove technical measures that prevented a startup from scraping public profile information of LinkedIn users source

- Sometimes you have to have an account and accept terms and contents which may forbid scraping.

You still have to respect local laws like the GDPR in Europe, which limits the collection of PII

# PART 3 - WEB-APIS

- Websites are finicky and prone to change

- CSS-selectors might not work tomorrow

- Webpages are made to be read by humans, not R-scripts

# WEB-APIS

APIs!

→ Application Programming Interfaces

→ Interfaces that are meant to be used by machines

There are different classes of APIs:

- Operating System APIs (e.g.: write a file to a folder)

- R-package APIs (e.g.: read a dataframe from a csv-file with read.csv())

- Web-API

- == webpage for machines

# WEB-APIS

Simple APIs, test them in your browser:

- https://dog-facts-api.herokuapp.com/api/v1/resources/dogs?number=3

- https://docs.openalex.org/about-the-data/work https://api.openalex.org/authors?filter=display_name.search:tom+theile

- https://api.openalex.org/works?filter=author.id:A773951722 as webpage

## WEB-APIS

Please open and follow the second script:

`02_using_web_APIs.R`

# WEB-APIS

Please open and follow the second script:

`02_using_web_APIs.R`

# DISCUSSION

You have read the paper "Do Anti-Immigrant Laws Shape Public Sentiment? A Study of Arizona's SB 1070Using Twitter Data" for yesterdays course

Now that you know how to scrape headlines: Imagine that you have dataset with a lot of headlines from various sources of the same timeframe that was handled in the Flores paper. What would be better or worse compared to twitter data?

# THANK YOU FOR YOUR ATTENTION!

**Tom Theile**

Research Software Engineer

theile@demogr.mpg.de