Digital media lab

Digital data in academia

DiMeLab Digital Media Lab · Roskilde University



Why the Digital Media Lab?

- Research current technology and data practice
- Developing new methods from new data sources and infrastructure
- Understand the limitations and possibilities of digital data for academic and commercial purposes
- Enhancing critical reflection through technical insight





Key services

- Homepage: <u>https://digitalmedialab.ruc.dk/</u>
- Tools
 - Twitter TCAT is running
 - Facebook tools via DMI
 - Instagram in test
 - YouTube via DMI
- Office hours every other Wednesday from 13-15 (Sander)
- Special events such as workshops and research seminars
- Supporting student projects at IKH





Other digital labs (INT/DK)

- DMI, Amsterdam
 - wiki.digitalmethods.net/
- TANT, Aalborg (Kbh)
 - <u>www.tantlab.aau.dk/</u>
- ITU, København
 - <u>ethos.itu.dk/</u>
- AAU, Aarhus
 - <u>www.digitalfootprints.dk/</u>
- RUC, Roskilde
 - https://digitalmedialab.ruc.dk/



| | | • • • • • • • • • • • • • • • • • • • | | ##### ################################ | | 말이 %~ 제 말 이 야 한 해요. 이 한 것 같 것 같 것 같 것 같 것 같 것 같 것 같 것 같 것 같 것 |
|--|---|---------------------------------------|--|---|---------------------|--|
| ************************************** | (1) 四位へを) (1) 四位人を) (1) 四位 | 448 | и «Гели и ва с с с с с с с с с с с с с с с с с с | (1) · · · · · · · · · · · · · · · · · · · | | 2001 #1.01.800 41.00 40. |



Media revolution





Data revolution

Regulating the internet giants The world's most valuable resource is no longer oil, but data

The data economy demands a new approach to antitrust rules



https://www.economist.com/news/leaders/21721656-dataeconomy-demands-new-approach-antitrust-rules-worlds-mostvaluable-resource

- "Data is just like crude. It's valuable, but if unrefined it cannot really be used. It has to be changed into gas, plastic, chemicals, etc., to create a valuable entity that drives profitable activity; so must data be broken down, analyzed for it to have value." (Palmer, 2006)
- http://ana.blogs.com/maestros/ 2006/11/data_is_the_new.html

Big data and 4V

- Volume
- Velocity
- Variety
- Veracity
- (Value?)



DiMeLab



http://www.ibmbigdatahub.com/infographic/extracting-business-value-4-vs-big-data



Datafication

- To datafy a phenomenon is to put it in a quantified format so it can be tabulated and analyzed
- **Digitalization** is the process of converting analog information into the zeroes and ones of binary code so computers can handle it













Cambridge Analytica





Digital methods

"Digital methods is a term that seeks to capture a recent development in Internet-related research, summarized <u>as approaches to the web as data set</u>. Joining a larger computational turn in the social sciences and the digital humanities, it asks a series of questions about the quality of web data, the productivity of online collection and analytical methods, and ultimately the prospects of having the web serve as a site for grounding findings. <u>When may</u> <u>the web become the baseline for</u> <u>findings about social change?</u>" (Rogers, 2015)

Rogers, R. (2015). Digital Methods for Web Research. In Robert A. Scott & S. M. Kosslyn (Eds.), Emerging Trends in the Behavioral and Social Sciences. Hoboken, NJ: Wiley. doi:10.1002/9781118900772



Rogers, R. (2013). Digital methods. Cambridge, Massachusetts: The MIT Press.



Digital Methods Initiative

- "The Digital Methods Initiative is a contribution to doing research into the "natively digital". [...], the focus is on how methods may change, however slightly or wholesale, owing to the technical specificities of new media."
 - <u>https://wiki.digitalmethods.net/D</u> <u>mi/WebHome</u>





Digital collection or analysis

• Digital collection

- 1. Scraping
- 2. Crawling (Google spider)



• Digital analysis

- Traditional methods
 - Observation
 - Content analysis
 - Discourse analysis
 - ...
- Digital methods
 - Automatic
 - Software defined
 - Quantitative or mixed-methods



Digital methods in Danish (2017)









API according to text

 "An API is an interface provided by an application that lets users interact with or respond to data or service requests from another program, other applications, or Web sites. APIs facilitate data exchange between applications, allow the creation of new applications, and form the foundation for the "Web as a platform" concept." (Murugesan, 2007 in Helmond, 2015)

Helmond, A. (2015). The Platformization of the Web: Making Web Data Platform Ready. Social Media+ Society, 1(2), 2056305115603080.



API according to Wikipedia

Technical

- Application Programming Interface
- "In computer programming, an application programming interface (API) is a set of subroutine definitions, protocols, and tools for building application software. In general terms, it is a set of clearly defined methods of communication between various software components." (Wikipedia)
- <u>https://en.wikipedia.org/wiki/Applicat</u> <u>ion_programming_interface</u>

Less technical

- APIs are protocols that define how different types of software can communicate with each other
- APIs allow requests and give structured response
- Metaphorically APIs are like waiters that walks between the customer (client) and the kitchen (server)
- <u>https://youtu.be/s7wmiS2mSXY</u>



APIs and academia

 "Some social media companies make their data banks on users and usage patterns available through their APIs. Hence, the API is also an interface for researchers to collect data from a given social media service. Through small software scripts, researchers can access the API to retrieve, store, and manipulate digital traces left by the users of a service for further empirical analysis." (Lomborg og Bechmann, 2014)



Social media API interfaces/collection tools

| | Text | Metrics | Comments | Visual | Network |
|---------------------------------|------|---------|----------|--------|---------|
| Facebook (Netvizz) | Х | Х | Х | | |
| Twitter (TCAT) | Х | | | | Х |
| Instagram (Instaload) | Х | | | Х | |
| YouTube (Data tool) | Х | Х | Х | | Х |



Netvizz

Netvizz v1.6

Netvizz is a tool that helps you analyze different sections of the Facebook platform – mainly pages – for research purposes. For questions, please consult the FAQ and privacy sections. Please reference this paper in academic work. Netvizz is being updated regularly. If you encounter a problem, please check the FAQ for how to report it. Since this application has not passed Facebook's app review for the "Page Public Content Access" permission, it may stop working in the near future. More details here.

Developing and hosting netvizz costs time and money. If the tool is useful for you, please consider to Donate

The following modules are currently available:

page like network - analyze networks of pages connected through the likes between them page posts - analyze user activity around posts on pages page timeline Images - analyze images from the "Timeline Photos" album on pages search - interface Facebook's search function link stats - generate statistics for links shared on Facebook

Big pages can take some time to process (minutes or hours). Be patient and try not to reload!

Version History

| 1.6 | 7.8.2018 | Added interactive visualizations to the page like, page posts, and timeline images modules |
|-----|-------------|---|
| 1.5 | 28.7.2018 | Getting ready for app review: removed group module, moved to API v3.1, some cleanup. |
| 1.4 | 5 21.2.2018 | following the February 5 2018 API changes, a number of features have been removed from the page module; reaction counts were added to the basic statistics for the page module and to the image module |
| 1.4 | 4 5.6.2017 | added attachment retrieval for comments, both for page (top 200 & full stats) and group modules |
| 1.4 | 3 25.5.2017 | link stat module updated to API version 2.9, which adds valid comment and reaction counts |
| 1.4 | | solided the ability for excitate result to see like exception. Instead hereaftering and anti-hereitering |

- 1.42 3.3.2017 added the ability for multiple seeds to page like network, internal housekeeping and optimizations
- 1.41 3.12.2016 added caching to page like networks to reduce API strain: page information is now only retrieved once per 24 hours
- 1.4 5.11.2016 added page timeline images module and reduced the number of posts to calculate page activity in page like networks to 50 to reduce strain

- Pros
 - Popular platform
 - Engagement metrics
- Cons
 - Restrictive access (public pages)
 - Unreliable API



TCAT

(Part of) URL: (Part of) media URL: Startidate:

Number of tweets: 623.570 Number of distinct users: 295.687

Enddate

2018-11-01

2018-11-12 623.570

| Data selection | | | | |
|----------------------------|--|--|--|--|
| Select the dataset | | | | |
| US_midterms_election_201 | 8 640.942 tweets from 2018-10-30 22:18:00 to 2018-11-13 09:42:1 | 6.418.221 tweets archived so far (and counting) | | |
| Select parameters: | | | | |
| Query: | | (empty: containing any text*) | | |
| Exclude: | | (empty: exclude nothing*) | | |
| From user: | | (empty: from any user") | | |
| Exclude user: | | (empty: exclude no users*) (empty: anything in biography*) (empty: any singulage*) (empty: from any clent*) | | |
| User blo: | | | | |
| User language: | | | | |
| Twitter client URL/descr: | | | | |
| (Part of) URL: | | (empty: any or all URLs*) | | |
| (Part of) media URL: | | (empty: any or all media URLs*) | | |
| Startdate (UTC): | 2018-11-01 | (YYYY-MM-DD or YYYY-MM-DD HH:MM:SS) | | |
| Enddate (UTC): | 2018-11-12 | (YYYY-MM-DD or YYYY-MM-DD HH MM:SS) | | |
| update overview | | | | |
| * You can also do AND or 0 | R queries, although you cannot mix AND and OR in the same query. | | | |
| | | | | |
| Overview of your select | on | | | |
| Dataset: | US_midterms_election_2018 (#midterms2018) | | | |
| Search query: | | Torets containing lines | | |
| Comments: | For alle dem der vil undersøge midtvejsvalget i USA | • Tedets | | |
| Exclude: | | oonlaking ta Sinka | | |
| From user: | | | | |
| Exclude from user: | | | | |
| the liter statistics and | | | | |

• Pros

- Reliable data access
- Public by default platform
- Cons
 - Streaming API (not historical data)
 - 1% limitation sample
 - Metrics
 - Mostly media and politics



Instagram

Instaloader v4.1.1 Site Contents -

>

Useful Links +

Instaloader

Instaloader is a tool to download pictures (or videos) along with their captions and other metadata from Instagram.

With Python installed, do:

\$ pip3 install instaloader

\$ instaloader profile [profile ...]

See Install Instaloader for more options on how to install Instaloader.

Instaloader

- · downloads public and private profiles, hashtags, user stories, feeds and saved media,
- · downloads comments, geotags and captions of each post,
- automatically detects profile name changes and renames the target directory accordingly,
- allows fine-grained customization of filters and where to store downloaded media,
- is free open source software written in Python.

instaloader [--comments] [--geotags] [--stories] [--highlights] [--tagged] [--login YOUR-USERNAME] [--fast-update] profile | "#hashtag" | :stories | :feed | :saved

See Download Pictures from Instagram for a detailed introduction on how to use Instaloader to download pictures from Instagram.

• Pros

- Popular platform with the youth
- Public by default
- Cons
 - Not public API but scraping
 - Only collection from ind. profiles
 - No metrics
 - Messy output (can be fixed)



YouTube

| YouTube Data Tools | blog software research DMI about |
|---|--|
| name) (Channal Life) (Channel Network) (Related Channel Network) (Marci U.P. (Marci Network) (Marci 1965) (1965) | |
| This is a collection of simple tools for extracting data from the YouTube platform via the COMPLETERS. For some context and a small introduction, please check out this COMPLETERS. Each of the modules has a basic description of how it works; there is a COMPLETERS section with additional information, and an COMPLETERS | 00 ¥ 40000. |
| Modules | |
| Channel Infe | |
| This module retrieves different kinds of information for a channel from a specified channel id. | |
| Chonnel Network | |
| This module crawls a network of channels connected via the "featured channels" (and via subscriptions) tab from a list of seeds. So | reds can be channels retrieved from a search or via manual input of channel ids |
| Video List | |
| This module creates a list of video infos and statistics from one of four sources: the videos uploaded to a specified channel, a play a list of ids. | ist, the videos retrieved by a particular search query, or the videos specified by |
| | |
| Video Network | |

- Pros
 - API seems reliable (so far)
 - Popular platform
 - Networks
- Cons
 - Video centric



$\mathbf{D} - \mathbf{I} - \mathbf{Y}$



What is the workshop about?

- 1. Collecting social media data via API
- 2. Getting to know your data: first steps
 - Getting from csv to Excel
 - Sorting data in Excel
 - Doing a time series analysis
- 3. Thinking about next steps: quant and qual
- i. <u>Not</u> included in this workshop (but potential future workshops)
 - Inferential statistics and predictive analysis
 - NLP, AI and machine learning methods
 - Network analysis
 - Image recognition
- <u>Please go to: digitalmedialab.ruc.dk/launch</u>



Output in CSV/TSV/TAB to Excel

- Delimiter-separated values
 - CSV: comma-separated values
 - TSV: Tab-separated values (.tab)
- Two-dimensional data divided into rows and columns in Excel
- Each row is a row but columns are defined through commas, tabs etc.



Pivot tables in Excel



| Data | Vindue | Hjælp | |
|--------|--------------|--------------|-------------|
| Sort | er | | 企 器R |
| Filtre | ér | | |
| Ryd | filtre | | |
| Avar | nceret filte | er | |
| Sub | totaler | | |
| Valio | dering | | |
| Data | atabel | | |
| Teks | st til kolon | ner | |
| Kons | solider | | |
| Grup | opér og di | sposition | • |
| Redi | iger links. | •• | |
| Ops | ummer m | ed pivottabe | el. |
| Tabe | elværktøje | er | • |
| Hen | t eksterne | data | • |
| Opd | ater data | | |