# 5.
# Tools and methods for data collection, secondary use of collected datasets

# Universially useful skills and tools

- R
- Python
- (Excel)
- Tableau (http://www.tableau.com/)

# Data Collection

# Data access options

| | | |
|---|---|---|
| APIs | Official resellers | „manual" forms of collection |
| Re-using published datasets | Third party tools | (Crowdsourcing) |
| | (via social media companies) | |

# Inequality in data access possibilities

- Data haves and data have nots
  - Financial reasons
  - Connections to companies
  - Different skills
  - …

# Data access options

**Via APIs (Application Programming Interface)**

- Not developed for scholarly uses – but still frequently used.

- With several limitations

- May change their terms of services

# Wikipedia API for Python

# Wikipedia



- Editable articles about present and past events
- Collective intelligence?

# Wikipedia: Edit Wars

## TABLE I

NUMBER OF REVERTS DETECTED. THE UPPER PART CORRESPONDS TO A GROUP OF PAGES WITH SEVERE CONFLICTS (EXCEPT THOSE IN *italics*); BELOW THE HORIZONTAL LINE THERE ARE PEACEFUL PAGES (EXCEPT THOSE IN *italics*).

| Both txt and cmt | Only in txt | Only in cmt | Article titlee |
|---|---|---|---|
| 4103 | 930 | 328 | Global warming |
| 2375 | 478 | 142 | Homosexuality |
| 1847 | 617 | 201 | Abortion |
| 1494 | 260 | 35 | *Benjamin Franklin* |
| 1425 | 437 | 130 | Elvis Presley |
| 1396 | 233 | 67 | Nuclear power |
| 1298 | 536 | 104 | Nicolaus Copernicus |
| 1071 | 211 | 51 | Tiger |
| 1036 | 248 | 58 | Euthanasia |
| 937 | 204 | 58 | Alzheimer's disease |
| 870 | 192 | 50 | Gun politics |
| 836 | 172 | 23 | *Sherlock Holmes* |
| 689 | 213 | 49 | Arab-Israeli conflict |
| 659 | 496 | 138 | Israel and the apartheid analogy |
| 652 | 387 | 88 | Liancourt Rocks |
| 642 | 250 | 39 | Schizophrenia |
| 516 | 164 | 472 | Gaza war |
| 431 | 186 | 30 | 1948 Arab-Israeli war |
| 416 | 73 | 9 | *Pumpkin* |
| 380 | 284 | 58 | Gdańsk |
| 318 | 158 | 20 | SQL |

Yasseri, T., Sumi, R., Rung, A., Kornai, A., & Kertész, J. (2012). Dynamics of conflicts in Wikipedia. PLOS ONE 7(6): e38869.

# Wikipedia: Cultural Disparities?

## Top 10 [Bearbeiten]

Die zehn Wikipedias mit den meisten Artikeln (Stand: 2. November 2014) sind:

| Platz ⇕ | Sprache ⇕ | Hauptseite ⇕ | Artikel ⇕ | Artikelanzahl ⇕ | Bearbeitungen ⇕ | aktive Benutzer ⇕ |
|---|---|---|---|---|---|---|
| 1. | Englisch | en.wikipedia.org ⮷ | englischsprachige Wikipedia | 4.636.938 | 741.415.918 | 133.327 |
| 2. | Schwedisch | sv.wikipedia.org ⮷ | schwedischsprachige Wikipedia | 1.946.828 | 28.362.530 | 2.884 |
| 3. | Niederländisch | nl.wikipedia.org ⮷ | niederländischsprachige Wikipedia | 1.794.646 | 43.378.140 | 4.136 |
| 4. | Deutsch | de.wikipedia.org | deutschsprachige Wikipedia | 1.771.852 | 141.065.828 | 19.583 |
| 5. | Französisch | fr.wikipedia.org ⮷ | französischsprachige Wikipedia | 1.557.616 | 111.055.268 | 15.797 |
| 6. | Wáray-Wáray | war.wikipedia.org ⮷ | Wáray-Wáray-Wikipedia | 1.258.797 | 6.061.557 | 91 |
| 7. | Cebuano | ceb.wikipedia.org ⮷ | Cebuanosprachige Wikipedia | 1.208.461 | 5.648.796 | 58 |
| 8. | Russisch | ru.wikipedia.org ⮷ | russischsprachige Wikipedia | 1.159.533 | 78.861.652 | 10.604 |
| 9. | Italienisch | it.wikipedia.org ⮷ | italienischsprachige Wikipedia | 1.153.610 | 74.240.982 | 7.969 |
| 10. | Spanisch | es.wikipedia.org ⮷ | spanischsprachige Wikipedia | 1.135.390 | 83.351.817 | 17.215 |

# Wikipedia dumps

- Also available via the Internet Archive.

- Capture Wikipedia at one point in time.

- E.g. German Wikipedia is about 3 GB plus 15 GB for images.

https://dumps.wikimedia.org/
https://en.wikipedia.org/wiki/Wikipedia:Database_download

# What can you do with Wikipedia data?

- Look at edit history.

- Use as text corpus.

- Compare text length, language use etc. for example across categories.

- Extract links.

- Compare language versions.

- …

# Twitter APIs

STREAMING API

- Push-based, Live-Stream

- Public stream vs. User stream

- Researchers need tools to keep connected to the stream.

REST API

- Limited number of requests per hour

- Social graph data (who is following whom)

- Trending topics

- And many more

https://dev.twitter.com/docs/using-search

# How To Search Twitter

Twitter searches have a wide range of options.

| Example | How It Works |
|---|---|
| twitter search | containing both "twitter" and "search". This is the default operator. |
| "happy hour" | containing the exact phrase "happy hour". |
| love OR hate | containing either "love" or "hate" (or both). |
| beer -root | containing "beer" but not "root". |
| #haiku | containing the hashtag "haiku". |
| from:alexiskold | sent from person "alexiskold". |
| to:techcrunch | sent to person "techcrunch". |
| @mashable | referencing person "mashable". |
| "happy hour" near:"san francisco" | containing the exact phrase "happy hour" and sent near "san francisco". |
| superhero since:2010-12-27 | containing "superhero" and sent since date "2010-12-27" (year-month-day). |
| ftw until:2010-12-27 | containing "ftw" and sent up to date "2010-12-27". |
| movie -scary :) | containing "movie", but not "scary", and with a positive attitude. |
| flight :( | containing "flight" and with a negative attitude. |
| traffic ? | containing "traffic" and asking a question. |
| hilarious filter:links | containing "hilarious" and linking to URLs. |
| news source:twitterfeed | containing "news" and entered via TwitterFeed |

**Twitter Politik** @twitter_politik — 23 Sep

Deutschland hat gewählt. Insgesamt gab es am Wahltag mehr als 350.000 Tweets zur #btw13. blog.twitter.com/de/2013/deutsc... pic.twitter.com/Ki9Bj523hQ
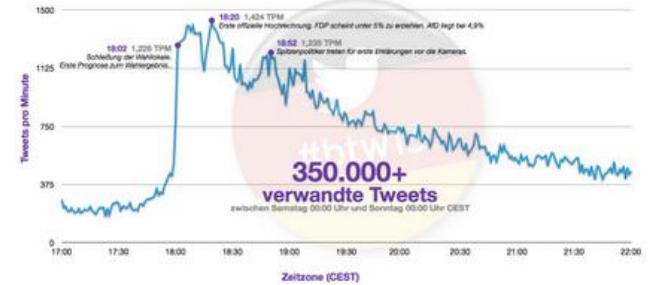
Retweeted by Katrin Weller

Hide photo     ← Reply   ⇄ Retweeted   ★ Favorite   ••• More

**#BTW13: Der Wahltag**
Tweets mit #Btw13, Spitzenkandidaten, Parteien und relevanten Schlagwörtern — 22. September 2013
Folgt @twitter_politik für weitere Updates und Informationen.

350.000+
verwandte Tweets
zwischen Samstag 00:00 Uhr und Sonntag 00:00 Uhr CEST

Zeitzone (CEST)

45 RETWEETS     10 FAVORITES

```
1  {
2       "created_at": "Mon Sep 23 09:28:38 +0000 2013",
3       "id": 382074015386329100,
4       "id_str": "382074015386329088",
5       "text": "RT @twitter_politik: Deutschland hat gewählt. Insgesamt gab es am Wahltag mehr als 350.000 Tweets zur #btw13. https://t.co/TBg8b8f50q http:…",
6       "source": "<a href="http://www.tweetdeck.com" rel="nofollow">TweetDeck</a>",
7       "truncated": false,
8       "in_reply_to_status_id": null,
9       "in_reply_to_status_id_str": null,
10      "in_reply_to_user_id": null,
11      "in_reply_to_user_id_str": null,
12      "in_reply_to_screen_name": null,
13      "user": {
14        "id": 78688499,
15        "id_str": "78688499",
16        "name": "Katrin Weller",
17        "screen_name": "kwelle",
18        "location": "Düsseldorf / Köln",
19        "description": "Information scientist. Interests in Web Science, Social Semantic Web, Social Media, Twitter analytics and more... \r\nworking at @gesis_org",
20        "url": "http://t.co/SJp02wO6",
21        "entities": {
22          "url": {
23            "urls": [
24              {
25                "url": "http://t.co/SJp02wO6",
26                "expanded_url": "http://katrinweller.net",
27                "display_url": "katrinweller.net",
28                "indices": [
29                  0,
30                  20
31                ]
32              }
33
```

Find your own user ID:
http://www.idfromuser.com/

https://dev.twitter.com/console

# Sina Weibo

## Practical Guide for using Sina Weibo's API

Hi! We are Wang Ling and Guang Xiang from the Language Technologies Institute in Carnegie Mellon University.

This is a practical guide for programmers and researchers who intend to use Sina Weibo's Open API as part of their work. We found during our research (description found here) that the current official documentation is unclear and inconsistent in some extent, and in several occasions we had to learn using a trial and error process, which took some time and effort. Thus, we are writing this guide to share what we know and facilitate the usage of this API to new users.

This guide is intended for both Mandarin and non-Mandarin speakers.

In this guide you will learn to:

1. Create a Sina Weibo Account
2. Set up Open Weibo Account and Register an Application
3. Using the Sina Weibo API

## 1 - Creating a Sina Weibo Account

The first step is to create an Weibo account, this should be relatively straight-forward, but might be hard if you cannot read Mandarin. If tha your case follow the steps below.

1-1 - Go to http://www.weibo.com/ and you should see the interface below

http://www.cs.cmu.edu/~lingwang/weiboguide/

https://builtvisible.com/gathering-data-sina-weibo-api/

18

# Sina Weibo

# Data access options

**Via offical resellers**

- Rather costly fees

- Contracts may need to be made with the head of an institute rather than with individual researchers.

- Some technical understanding is required

- May prohibit sharing the data

- Main players: GNIP and Datasift

# GNIP

# DataSift

# Data access options

**Via third party tools**

- Black box effect – lack of transparancy

- Fees may apply

- Changes are possible

- Additional layer of legal restrictions (e.g. limiting what might be published)

# Data access options

**„Manual" forms of collection**

- E.g. copy-and-paste, screenshots
- Time-consuming, error-prone
- Little exchange about best practices, non standardized
- Good in addition to other methods, e.g. to capture images and as„proof"

# Data access options

**Reuse of existing datasets**

- Lack of available datasets
- Lack of standards for documentation and formats
- Lack of persistent identifyers
- Legal and ethical restrictions

# Data access options

**Crowdsourcing**

- Little guidance / best practices
- Via Amazon Mechanical Turk and others: ethical challenges
- Via Apps (e.g. on Facebook)

# Data access options

**Via social media companies**

- Usually in form of priviledged access (special agreements)
- High-quality data
- Limitations on usage and sharing

# Research groups that offer different tools or training

# Digital Methods Initiative

- https://wiki.digitalmethods.net/Dmi/ToolDatabase

# Digital Methods Initiative

Includes

- Tools for extracting and resolving URLs

- Tool for scraping google images

- Wikipedia Edits Scraper

- …

They also offer Summer / Winter Schools.

- https://wiki.digitalmethods.net/Dmi/ToolDatabase

# Statistical Cybermetrics Research Group Wolverhampton



**Statistical Cybermetrics** Research Group

> Home   > Meetings   > People   > Projects   > Publications   > Awards   > Resources   > Services   > Database   > About

## Free Link Crawler, Hyperlink Database and Resources

- The sentiment strength detection program SentiStrength is available free for non-commercial use.

- A brief desrciption of information-centred research (ICR), a paradigm used extensively by the research group.

- A MySpace ID map is available, especially for researchers using the CPAN Perl module WWW::Myspace.

- A summary of instructions for webometric searches in the main search engines.

- Mozdeh - Twitter time series and sentiment analysis is online with instructions.

- Webometric Analyst link analysis software for lists of links or URLs - together with an interface for automatically fetching results from Bing and for various other webometric tasks.

- SocSciBot free Web link crawler with free link analysis and text analysis tools.

- Free database of web link files for the UK, Australian, New Zealand, Spanish, Mainland Chinese and Taiwanese academic Webs.

- SCIT Cybermetrics Web intelligence consulting - our commercial consultancy venture offering services to businesses and non-commercial organisations.

- The web site for the book:Link analysis: An information science approach. Academic Press.

http://cybermetrics.wlv.ac.uk/resources.htm

# Statistical Cybermetrics Research Group Wolverhampton

Includes for example:

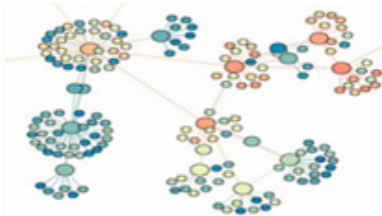- Text analysis, Sentiment Analyse

- Webometrics: „Web citation data"

http://cybermetrics.wlv.ac.uk/resources.htm

# Open Educational Ressources

FREE ONLINE COURSE

# Social Media Analytics: Using Data to Understand Public Conversations

How can you tap into social media conversations, and discover what's being said about the things that matter to you?

**Join now – started 18 Jul**

FREE online course     Duration: 3 weeks     3 hours pw     Certificates available

#imhooked     #sogood
#winning     #boycott
#bigfan     #yawn

**CREATED BY**

QUT

View transcript     Download video: standard or HD

**SHARE**

https://www.futurelearn.com/courses/social-media-analytics

# Data collection tools

# Data collection tools

- Run on your own server

- Install on your own laptop

- Web-based

# YourTwapperkeeper



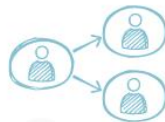| Archive ID | Keyword / Hashtag | Description | Tags | Screen Name | Count | Create Time |
|---|---|---|---|---|---|---|
| 5 | opco12 | Open Course 2012 e-learning | e-learning | kwelle | 2600 | Wed, 01 Feb 2012 16:35:14 +0100 |
| 4 | irhhu | Tweets zur Vorlesunr IR 2011/12 | | terliesner | 1161 | Fri, 20 Jan 2012 12:03:07 +0100 |
| 6 | iwhhu | InformationsWissenschaftHeinrichHeineUni | | kwelle | 67 | Wed, 01 Feb 2012 16:36:37 +0100 |
| 7 | Düsseldorf | just for fun | | kwelle | 241217 | Wed, 01 Feb 2012 16:46:22 +0100 |
| 8 | dgi2012 | DGI-Konferenz 2012 | conference | kwelle | 794 | Wed, 01 Feb 2012 16:55:20 +0100 |
| 9 | cosci12 | Conference on Science and the Internet 2012 | conference | kwelle | 11 | Wed, 01 Feb 2012 16:55:48 +0100 |
| 10 | diata12 | Twitter Workshop DIATA12 | conference | kwelle | 83 | Wed, 01 Feb 2012 16:56:04 +0100 |
| 11 | #grammys | grammy verleihung 2012 | event | kwelle | 244131 | Fri, 10 Feb 2012 |

Bruns, A., & Liang, Y. E. (2012). Tools and methods for capturing Twitter data during natural disasters. First Monday, 17(4). http://doi.org/10.5210/fm.v17i4.3937

http://mappingonlinepublics.net/2011/06/21/switching-from-twapperkeeper-to-yourtwapperkeeper/

# Digital Footprints (for Facebook)

# COSMOS



**gesis** Leibniz Institute for the Social Sciences

SOCIAL DATA SCIENCE LAB.

Scalable real-time social data analytics for research, policy & practice

Home    Projects    Publications    People    Events    **Software**    News    Ethics

The Social Data Science Lab maintains and distributes the COSMOS software. COSMOS is available at no cost to academic institutions and not-for-profit organisations. Please email us to request a download link using your institutional email address. NOTE: COSMOS works best on Mac OS X and Linux Ubuntu. Only install on Windows as a last resort.

Your Name (required)

Your Email (required)

Organisation (required)

**P 8 R 4**

Send

Credits:

- COSMOS user interface, data analysis and visualisation by Cardiff University
- COSMOS data collection by University of St. Andrews
- Sentiment Analysis via SentiStrength by University of Wolverhampton

# Social Feed Manager

**gesis** Leibniz Institute for the Social Sciences

http://www.tubekit.org/

**TubeKit**
A Youtube Crawling Toolkit

Home    Download    Tools    Toolbar    Documentation

Note: due to some recent changes in YouTube APIs, many of the TubeKit's processes may not function. We are working hard to fix this issue.

## What is it?

TubeKit is a toolkit for creating YouTube crawlers. It allows one to build one's own crawler that can crawl YouTube based on a set of seed queries and collect up to 16 different attributes.

TubeKit assists in all the phases of this process starting database creation to finally giving access to the collected data with browsing and searching interfaces. In addition to creating crawlers, TubeKit also provides several tools to collect a variety of data from YouTube, including video details and user profiles.

**⬇ Download TubeKit**

🐙 Fork on GitHub

### Collecting data from YouTube without crawling:

TubeKit includes several tools, which lets one

- extract YouTube video links from any webpage
- extract YouTube video data (author, keywords, genre, number of views, ratings, comments, etc.)
- collect text comments for YouTube videos
- extract a YouTube users' profile data

## Create a Youtube Crawler in 5 EASY STEPS!

1 Provide basic information
(project name, directory to store the crawler, etc.)

2 Set up the database

3 Select up to 16 different attributes to collect for

We have also developed a Firefox toolbar for TubeKit, which lets one accomplish YouTube data collection with just a click!

**Manual data collection**
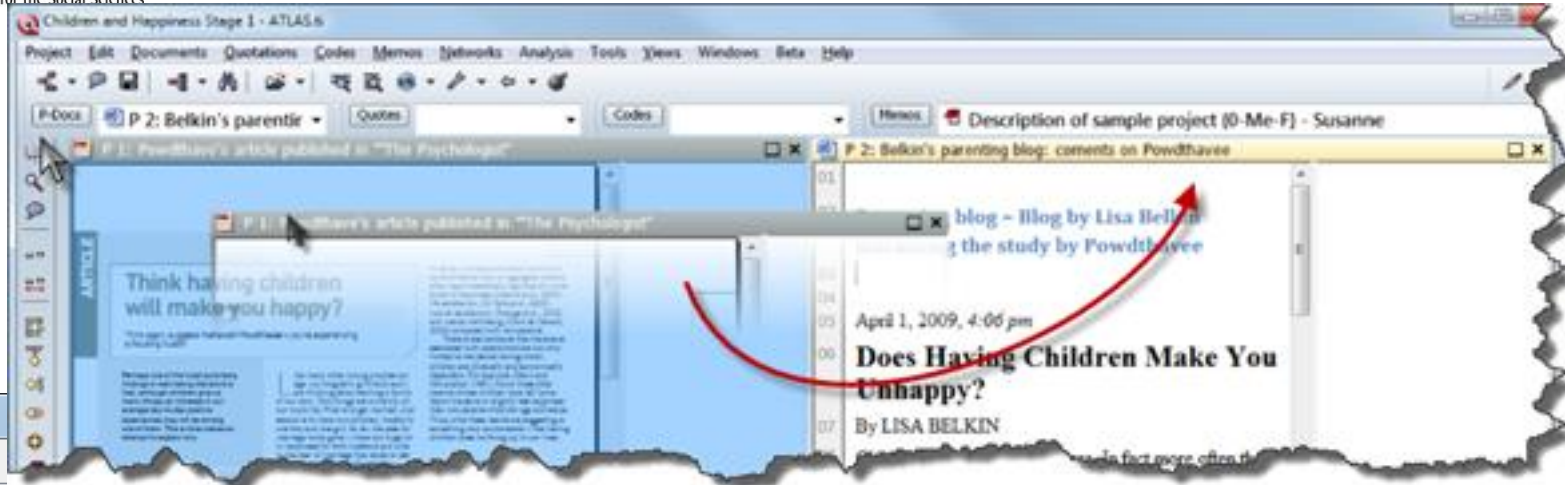**Handling multimedia data?**

# Explore options for capturing and managing data

- Evernote

- Zotero

- Content Analysis Software.

# Content Analysis

CAQDAS Computer-Assisted Qualitative Data AnalysiS

- ▶ (e.g. MAXQDA, QDAMiner, ATLAS.ti, Qualrus, Nvivo)
- ▶ Speech Act Analyse
- ▶ Statistics

- ▶ Some of them now come with import options for social media data!

**ATLAS.ti**

1 Families:
2 term-stack pane
3
4 term-stack pane
5 feedback pane
6 result pane

Create Super Code

output (view, save or print reports)

# NVIVO

- Also able to import tweets.

http://www.qsrinternational.com/

# Data analysis

# Data analysis

Network analysis

Sentiment analysis

Visualization

Content analysis

(Crowdsourcing)

# Network Analysis

# NodeXL



Network Analysis

Based on Excel

Can collect data from (amongst others):

- Facebook

- Twitter

- YouTube

- Flickr

- Wikipedia

- If you can get data into excel, you can get it into NodeXL…

# NodeXL

# NodeXL

# NodeXL: SocialNetImporter

http://socialnetimporter.codeplex.com/

# NodeXL

- NodeXLGraphGallery.org

# Visualization

# Gephi

# Gephi

# Sentiment analysis

# NLP / Sentiment Analysis

- Discover text:  http://discovertext.com/

- Twitter NLP:
  http://www.ark.cs.cmu.edu/TweetNLP/

- Internet Community Text Analyzer (ICTA)
  http://textanalytics.net

- Twitter NLP and Part-of-Speech Tagging:
  http://www.ark.cs.cmu.edu/TweetNLP/

- SentiStrength http://sentistrength.wlv.ac.uk/

# Sentiment Analysis

# Others

# RIP: TOPSY



## Topsy, The Popular Social Analytics Service Bought By Apple, Closes Down

Posted Dec 15, 2015 by *Jon Russell* (*@jonrussell*)

Topsy, the Twitter and general social media analytics service quietly bought by Apple for more than $200 million two years ago, closed its doors today.

That's according to the company itself, which announced the shutdown in what is ironically its first tweet since being bought by the iPhone-maker.

**Topsy** ✔
@Topsy

**Follow**

We've searched our last tweet.

### CrunchBase

**Topsy Labs**

**FOUNDED**
2007

**OVERVIEW**
Topsy is a social analytics company that gives you instant answers to critical business questions through real-time analysis of public conversations. Operating the world's only index of the public social web, Topsy is used by some of the largest global marketing, news, entertainment, and financial organizations to get real

# Oxford Internet Institute

Measuring impact online:

http://microsites.oii.ox.ac.uk/tidsr/welcome

# More

- Mapbox (add data to maps): https://www.mapbox.com/

- classify websites into categories? http://diffbot.com/products/automatic/classifier/

# More?

- Keep looking for the following keywords:
  - ▶ Web Science
  - ▶ Data Science
  - ▶ Computational Social Science
  - ▶ Digital Humanities
  - ▶ Digital Methods
  - ▶ Internet Research
  - ▶ Social Network Analysis

  Open Educational Ressources:
  - ▶ Coursera
  - ▶ EdX
  - ▶ …

# Crowdsourcing
# – for data collection and analysis

# Get your own Facebook data



**Help Centre**                                        ⬆ Go to News Feed

Hi Katrin, how can we help?                                    🔍

🏠 **Desktop Help** ▸ **Manage Your Account**                    English (UK)

Names on Facebook

Account Settings                                              ›

Warnings & Blocks                                            ›

Resetting Your Password                                      ›

Deactivating & Deleting Accounts

Memorialized Accounts

**Downloading Your Info**

Disabled Accounts

About Facebook Ads

Interacting with Ads

## Downloading Your Info

How can I download my information from Facebook?                 ⌄

Can I pick and choose which information I would like to download?  ⌄

What security measures are in place to make sure someone else doesn't  ⌄
download a copy of my information?

Learn more about what's included in your download.

# Crowdsourcing

distribution of number of friends for our users

friend ages for people of different ages

population vs. age

71

# Amazon Mechanical Turk

https://www.mturk.com/mturk/welcome

http://www.wearedynamo.org/forum/actions/topics/call-for-comments-turker-authored-guidelines-for-research-on-amt?page=1

# Research ethics

# Legal Framework?

Terms of Services

Ethics

# Privacy and protecting subjects

- Anonymisation?

- Informed consent?

- Guidelines for handling sensitive user data?

# Data Sharing

# One of the Challenges: Data Sharing

*"But you can't make your data available for others to look at, which means both your study can't really be replicated and it can't be tested for review. But also it just means your data can't be made available for other people to say, Ah you have done this with it, I'll see what I can do with it, (…) There is no open data."*

Weller, Katrin, and Katharina E. Kinder-Kurlanda. 2015. "Uncovering the Challenges in Collection, Sharing and Documentation: The Hidden Data of Social Media Research?." In Standards and Practices in Large-Scale Social Media Research: Papers from the 2015 ICWSM Workshop. Proceedings Ninth International AAAI Conference on Web and Social Media Oxford University, May 26, 2015 – May 29, 2015, 28-37. Ann Arbor, MI: AAAI Press.

# Reusing datasets

*"I actually only use [other researcher's datasets] where I'm very sure about where it comes from and how it was processed and analyzed. There is too much uncertainty in it."*

Weller, Katrin, and Katharina E. Kinder-Kurlanda. 2015. "Uncovering the Challenges in Collection, Sharing and Documentation: The Hidden Data of Social Media Research?." In Standards and Practices in Large-Scale Social Media Research: Papers from the 2015 ICWSM Workshop. Proceedings Ninth International AAAI Conference on Web and Social Media Oxford University, May 26, 2015 – May 29, 2015, 28-37. Ann Arbor, MI: AAAI Press.

# How much should I share?

**Most reproducibility**



**What is being shared?**

- whole dataset plus additional research information (e.g. scripts)
- whole dataset
- whole dataset, but without direct identifiers (pseudonymization)
- parts of the dataset removed (anonymization)
- changed dataset (e.g. only tweet IDs)

**Most privacy**

# Examples of current practices

| Approach | Example | Retrievability | Documentation + standards | Long-term perspective |
|---|---|---|---|---|
| „Grey market" | Data shared with colleagues (often upon request) | **Low**, via personal connections | **Low**, no standardized documentation or data format | **Low**, no guaranteed long term availability |
| Researchers' personal/ professional websites | | **Medium**, URLs may be referenced | **Low**, no standardized documentation or data format | **Low**, no guaranteed long term availability |
| Social media providers | e.g. Wikipedia dumps | **High**, if directly provided from social media platform | **Medium/high**, depending on the producer | **Medium/high**, depending on the producer |
| Project-based or thematic collections | e.g. KONECT, CrisisLex | **Medium**, URLs may be referenced | **Medium**, if same principles are applied within entire collection | **Low**, no guaranteed long term availability |
| Conferences & journals | e.g. ICWSM datasets with conference papers | **High**, usually related to accepted publications | **Medium/high**, depending on the publisher | **Medium/high**, depending on the publisher |
| Professional archives | e.g. datasets at GESIS data archive | **Medium/high**, datasets may be referenceable with DOIs. | **Medium/high**, depending on the archive's requirements | **High**, guaranteed availability for different time spans |

- For additional examples see our paper: http://dx.doi.org/10.1145/2908131.2908172
- See also: Thomson, S.D. 2016. Preserving Social Media. DPC Technology Watch Report. Retrieved from http://dpconline.org/publications/technology-watch-reports

80

# Available datasets

- From individual researchers/groups (sometimes „grey market").

- From conferences: e.g. ICWSM

- Archival institutions: e.g. GESIS (*doi:10.4232/1.12319*)

# Reddit dataset

This is an archived post. You won't be able to vote or comment.

**855**

**I have every publicly available Reddit comment for research. ~ 1.7 billion comments @ 250 GB compressed. Any interest in this?** (self.datasets)

eingereicht vor 1 Jahr * von Stuck_In_the_Matrix

I am currently doing a massive analysis of Reddit's entire publicly available comment dataset. The dataset is ~1.7 billion JSON objects complete with the comment, score, author, subreddit, position in comment tree and other fields that are available through Reddit's API.

I'm currently doing NLP analysis and also putting the entire dataset into a large searchable database using Sphinxsearch (also testing ElasticSearch).

This dataset is over 1 terabyte uncompressed, so this would be best for larger research projects. If you're interested in a sample month of comments, that can be arranged as well. I am trying to find a place to host this large dataset -- I'm reaching out to Amazon since they have open data initiatives.

**EDIT:** ~~I'm putting up a Digital Ocean box with 2 TB of bandwidth and will throw an entire months worth of comments up (~ 5 gigs compressed)~~ **It's now a torrent. This will give you guys an opportunity to examine the data. The file is structured with JSON blocks delimited by new lines (\n).**

_____

One month of comments is now available here:

**Download Link:** Torrent

**Direct Magnet File:** magnet:?xt=urn:btih:32916ad30ce4c90ee4c47a95bd0075e44ac15dd2&
dn=RC%5F2015-01.bz2&tr=udp%3A%2F%2Ftracker.openbittorrent.com%3A80&tr=udp%3A%2F
%2Fopen.demonii.com%3A1337&tr=udp%3A%2F%2Ftracker.coppersurfer.tk%3A6969&tr=udp%3A%2F
%2Ftracker.leechers-paradise.org%3A6969

**Tracker:** udp://tracker.openbittorrent.com:80

**Total Comments:** 53,851,542

https://www.reddit.com/r/datasets/comments/3bxlg7/i_have_every_publicly_available_reddit_comment/

# CrisisLex

# ICWSM-16

Cologne, Germany
May 17-20, 2016

The 10th International
Conference on Web
and Social Media

HOME    SUBMITTING    ATTENDING    PROGRAM    **DATASETS**    ORGANIZATION

CONTACT / ABOUT

DATASETS

G+1 0    Tweet    Share    Like 0

## DATASETS

### ICWSM DATASET SHARING SERVICE

As part of the ICWSM Data Sharing Initiative, ICWSM provides a hosting service for new datasets used by papers published in the proceedings of the annual ICWSM conference. All datasets are released as openly available community resources. Please see the instructions on the registration process in order to gain access to the datasets.

### AVAILABLE DATASETS

ICWSM-16 is the fifth year of this data sharing initiative.
Here is the list of available datasets:

- **iFeel 2.0: A Multilingual Benchmarking System for Sentence-Level Sentiment Analysis**
  Matheus Araújo, João P. Diniz, Lucas Bastos, Elias Soares, Manoel Júnior, Miller Ferreira, Filipe Ribeiro, Fabricio Benevenuto

- **Analyzing Personality through Social Media Profile Picture**

**WHEN, WHERE, WHO**

January 5, 2016   Abstract Submission

January 11, 2016 Full Paper Submission

February 25, 2016
Paper & Poster Notification

March 8, 2016
Camera Ready Version Due

May 17, 2016   Tutorials / Workshops

May 18-20, 2016 Main Conference

**SPONSORS**

**GOLD SPONSOR**

Microsoft
Research

**SILVER SPONSORS**

84

gesis — Leibniz Institute for the Social Sciences

gesis 🏠

∴ datorium

🏠 | **View Item**

## General Description

**Back**

cite this data

| | |
|---|---|
| **Title** | Elite actors in the U.S. political Twittersphere |
| **URI** | http://dx.doi.org/10.7802/1178 |
| **Primary Researcher** | Sebastian Stier; GESIS Leibniz Institute for the Social Sciences |
| **Publication Year** | 2016 |
| **Availability** | Free Access (without Registration) |
| **Subject Area** | Politikwissenschaft Kommunikationswissenschaften |
| **Topic Classification** | Kommunikation, öffentliche Meinung, Medien |
| **Abstract** | Lists of Twitter acccounts of U.S. political elites from government, news media and political parties. Metadata on number of followers (as of December 2015), date of account registration and Twitter ID. List of news media actors contains major newspapers and national TV stations as well as politically influential media persons and individual journalists. List of politicians contains all members of U.S. congress, governors, presidential candidates and main accounts of the Republican and Democratic parties. Metadata on politicians contains political office and party affiliation. List of government actors contains federal agencies, secretaries, the President and affiliated accounts. |
| **Geographical Area** | Vereinigte Staaten / US |
| **Survey Period** | 2015-12; |

### Files in this submission

☐ I accept the terms of use of datorium

Mitglied der
Leibniz-Gemeinschaft

Search

Add data

About datorium ▶

# Unvailable datasets

**TECHNOLOGY**

## Library of Congress' Twitter archive is a huge #FAIL

More than five years on, the library's Twitter archive project is in limbo — with no end in sight.

By **NANCY SCOLA** | 7/11/15 5:09 PM EDT



http://www.politico.com/story/2015/07/library-of-congress-twitter-archive-119698.html

**INSIDE HIGHER ED**

| View Exclusive AAUP Compensation Survey Data | SUBSCR

NEWS | VIEWS | CAREER ADVICE | BLOG U | SURVEYS | WEBINAR

ADMISSIONS | BOOKS | TECHNOLOGY | COMMUNITY COLLEGES | DIVERSITY | TEACHING & LEARN

Views

## The Archive Is Closed

June 3, 2015
By Scott McLemee

**Five years ago**, this column looked into scholarly potential of the Twitter archive the Library of Congress had recently acquired. That potential was by no means self-evident. The incensed "my tax dollars are being used for this?" comments practically wrote themselves, even without the help of Twitter bots.

For what -- after all -- is the value of a dead tweet? Why would anyone study 140-character messages, for the most part concerning mundane and hyperephemeral topics, with many of them written as if to document the lowest possible levels of functional literacy?

As I wrote at the time, papers by those actually doing the research treated Twitter as one more form of human communication and interaction. The focus was not on the content of any specific message, but on the patterns that emerged when they were analyzed in the

THE COUNCIL OF INDEPENDENT COLLEGES

Woodrow Wilson Visiting Fellows

https://www.insidehighered.com/views/2015/06/03/article-difficulties-social-media-research

# Background: LoC Twitter archive

- 2006: Twitter was made public
- 14.04.2010: Twitter donates archive to the Library of Congress: *"Recently, the Library of Congress signaled to us that the public tweets we have all been creating over the years are important and worthy of preservation."* https://blog.twitter.com/2010/tweet-preservation
- 2010-2012: nothing new on the "Twitter archive"…
- 04.01.2013: Update from the Library of Congress: *"We now have an archive of approximately 170 billion tweets and growing. The volume of tweets the Library receives each day has grown from 140 million beginning in February 2011 to nearly half a billion tweets each day as of October 2012."* http://blogs.loc.gov/loc/2013/01/update-on-the-twitter-archive-at-the-library-of-congress/
→ But still no access to the Twitter archive at the LoC and no information. *"Although the Library has been building and stabilizing the archive and has not yet offered researchers access, we have nevertheless received approximately 400 inquiries from researchers all over the world."*

# Conclusions 5

# Lessons learned

- There is no perfect tool. But you may find solutions that are sufficient for different usage scenarios.

- Tools come and go, you need to keep track of what is available.

# If you have time to read 3 papers…

- Borra, E., & Rieder, D. (2014). Programmed method: developing a toolset for capturing and analyzing tweets, Aslib Journal of Information Management, 66(3), 262 – 278. DOI: http://dx.doi.org/10.1108/AJIM-09-2013-0094
- Bruns, A., & Liang, Y. E. (2012). Tools and methods for capturing Twitter data during natural disasters. First Monday, 17(4). doi:10.5210/fm.v17i4.3937
- Gaffney, D., & Puschmann, C. (2014). Data collection on Twitter. In Weller, A. Bruns, J. Burgess., M. Mahrt and C. Puschmann (Ed.), Twitter and Society (pp. 55–68). New York: Peter Lang.