



Things to know when working with reddit data

Show & Tell – Social Media-Daten in
der Forschungspraxis II

Dr. Carolina Haensch





Thank you to Ashley Amaya et al., authors of the article „New Data Sources in Social Science Research: Things to Know Before Working With Reddit Data” on which I base part of my talk





LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

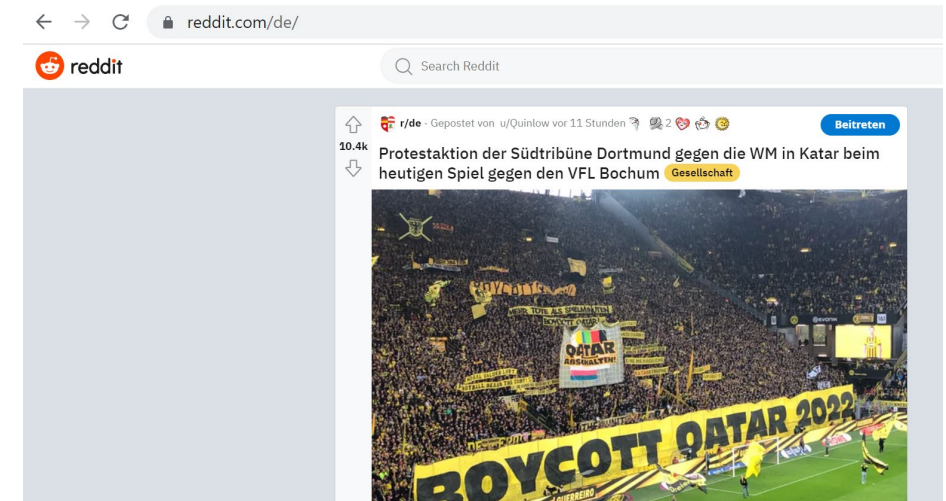
1. Reddit overview





Reddit as platform

- Popular social media site
 - Worldwide Nr. 20
 - Germany Nr. 37
 - US Nr. 9 (all numbers from 2021)
- Self-proclaimed “front page of the Internet”
- Founded in 2005





LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

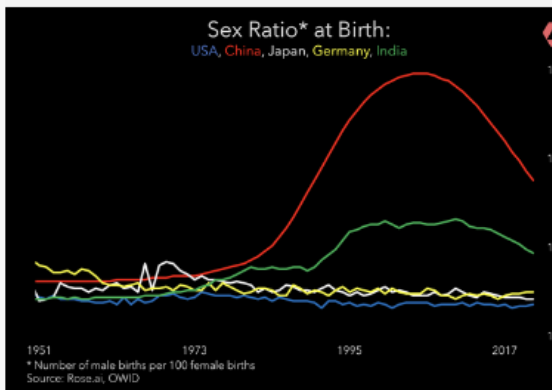
 **DataIsBeautiful** [Beitreten](#)
r/dataisbeautiful

[Beiträge](#) [Posting Rules](#) [Top OC of the Week](#)

Heiß Neu Top Heute

Gepostet von [u/rosetechnology](#) OC: 22 vor 21 Stunden

14.3k OC **Sex Ratio at Birth: USA, China, Japan, Germany, & India [OC]**



Sex Ratio* at Birth:
USA, China, Japan, Germany, India

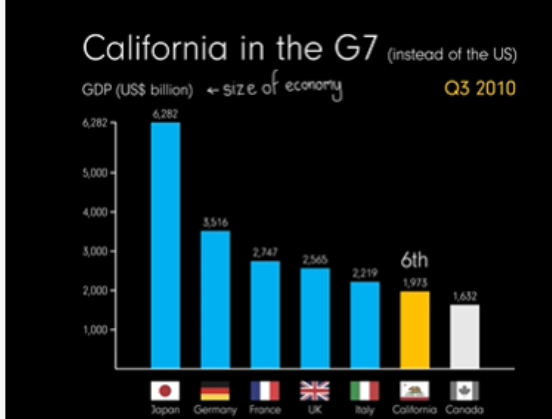
1951 1973 1995 2017

* Number of male births per 100 female births
Source: Roser, OWID

1.3k Kommentare Teilen Speichern

Gepostet von [u/jceagle](#) OC: 97 vor 20 Stunden

6.3k OC **[OC] The GDP of California in the G7 instead of the United States**



California in the G7 (instead of the US)
GDP (US\$ billion) ← size of economy Q3 2010

6,282 3,516 2,747 2,565 2,219 1,975 1,632

Japan Germany France UK Italy California Canada

Source: Federal Reserve Bank of St. Louis

EEAGLI

Über diese Community

DataIsBeautiful is for visualizations that effectively convey information. Aesthetics are an important part of information visualization, but pretty pictures are not the sole aim of this subreddit.

Erstellt am 14. Feb. 2012

18.8m Mitglieder 7.4k Online

Nach Flair filtern

Discussion OC

- r/dataisbeautiful Regeln**
1. A post must be (or contain) a qualifying data visualization.
 2. Directly link to the original source article of the visualization.
 3. [OC] posts must state the data source(s) and tool(s) used.
 4. DO NOT claim "[OC]" for visualizations that are not yours.
 5. All diagrams must have at least one computer generated element.
 6. No reposts of popular posts within 1 month.
 7. Post titles must describe the data plainly without using sensationalized headlines.
 8. Posts involving U.S. politics are allowed only on Thursday (ET)
 9. Posts regarding Personal Data are permissible only on Mondays (ET).
 10. Hate Speech or Dogwhistling
 11. Comment is unconstructive
 12. This is plagiarism. I'll explain in modmail.





Reddit basics

- Registered users can submit content
 - Content includes links, text posts, images, videos, GIFs
 - Upvoting and downvoting possible, influences position on pages
- Content is divided into boards/subpages called subreddits or communities
 - Example: `r/dataisbeautiful` or `r/dataisugly`
- Reddit administrators and subreddit-specific members moderate the subreddits



reddit r/brexit Subreddit Search r/brexit

Posted by 2 years ago

6
↑
↓

can anyone enlighten me of that the UK leaving means and or benefits? (in avg joe terms please)

i see all the fuss but i don't know what is all about?

32 Comments Give Award Share Save Hide Report 88% Upvoted

3 points · 2 years ago

↑
↓

The UK doesn't have to help support failed countries or follow laws forced on them by other countries any more. They get more freedom and the ability to choose what they want to do globally more readily now.

Give Award Share Report Save

0 points · 2 years ago

↑
↓

disingenuous to call controlling a countries national border = racism

Give Award Share Report Save

-3 points · 2 years ago

↑
↓

Are you sure? Stopping immigrants at your border because of where they come from seems to be pretty cut and dry racism. But justify it however you want. I'm all for the UK getting out of the EU, but closing your borders to Arab immigrants is a very very bad idea and will do a lot of harm to the economy that is already aimed at the gutter with your banking practices.

Give Award Share Report Save

1 point · 2 years ago

↑
↓

Never mind the language used here, this is what will happen and they have said so all along.

Give Award Share Report Save

0 points · 2 years ago

↑
↓

British Exit from the EU. They regain their independence and national sovereignty. They will control their own country again instead of unelected outsiders! It is a great day for the people of Britain!

Give Award Share Report Save

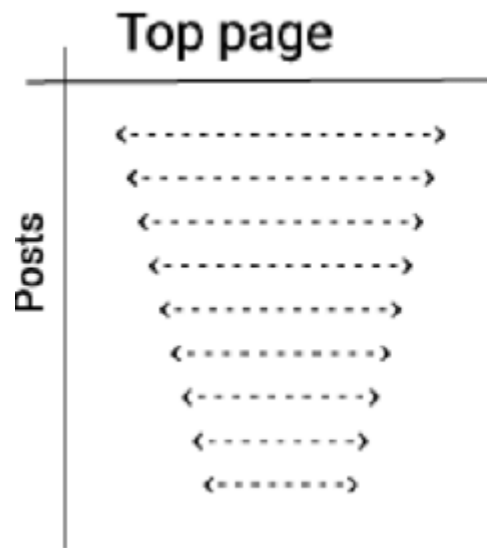
Thread root

First level Comment

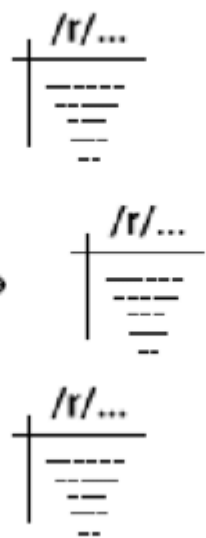
Third level Comment

https://www.researchgate.net/publication/348757045_Linking_User_Opinion_Dynamics_and_Online_Discussions

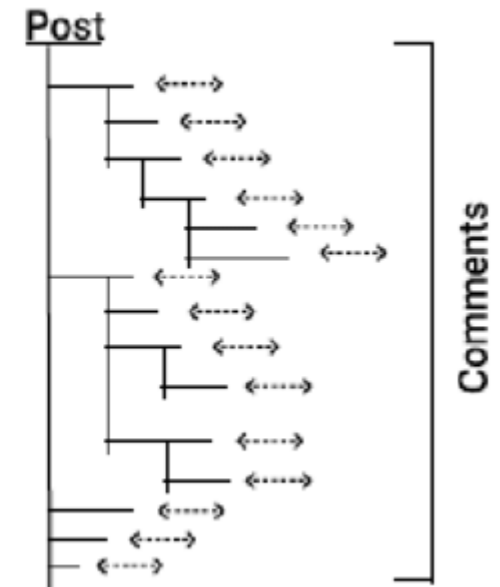
Reddit



Subreddits



Discussion tree



<https://sciendo.com/article/10.2478/subbs-2020-0005>



Reddit as a data source for research

- Amount/Size: Due to its popularity, there is A LOT of data available
- Availability: Most subreddits (exact statistics are unavailable) are public
 - Public subreddits can be viewed by all
 - Posts may be commented and voted on by all registered users (free and anonymous)
 - The data may be downloaded for free by any registered user
- Identification of rare subpopulations: 130,000 active subreddits
 - Researchers of rare populations can identify and study eligible participants



FRANKEN



Frankenkultur

Beitreten

r/Frankenkultur



Beitrag erstellen



Heiß



Neu



Top



VON MODS ANGEHEFTET

3

Gepostet von u/NicolaiWM **Franke** vor 2 Jahren



r/Frankenkultur Lounge

LIVE jetzt beitreten

0 Nachrichten

Auszeichnen

Teilen

Speichern



Gepostet von u/NicolaiWM **Franke** vor 8 Monaten

5

Bamberger Bierkeller-Empfehlungen? **Dorsch**

Über diese Community

Ein Subreddit für alle, die Franken lieben & im schönen Frankenland leben. Hier könnt ihr euch über Kultur, schöne Orte, gutes Essen und sonstige Interessen und Neuigkeiten in Franken austauschen! Teilt gerne eure Tipps.

Erstellt am 15. Dez. 2020

227

Mitglieder

4

Online

Oberste 50%

Sortiert nach Größe

Beitrag erstellen

Content created on Reddit in 2021



Posts
429,952,687

Comments
2,717,216,237

PMs
713,505,544

Chats
1,912,574,822

Total
5,773,249,290



Examples for data analyses with Reddit data

- Gender roles (Ammari, Schoenebeck, & Romero, 2018)
 - Variation across parenting subreddits (mommit and daddit). Similarities in topics across the boards, such as sleep training, as well as differences, such as fathers talking about custody cases and Halloween
- News consumption (Wasike, 2011)
 - Human interest frame was the most common generic frame followed by technology. Science and technology was the most common issue-specific frame.
- Sexual identity (En, En, & Griffiths, 2013)
 - Evaluate through subreddit self-descriptions which aspects of identity are seen – and thereby made – to be markers of personal and collective sexual identities
- Mental health (Choudhury & De, 2014)
 - Self-disclosure in mental illness communities and characterization of mental health social support



Examples

- Right-wing accounts (Gaudette et al., 2021)
 - Group building dynamics in subreddits through calls for violence against minorities and political opponents of right-wingers
- Corona-pandemic and news consumption (Chipidza et al., 2022)
 - News sources in liberal and conservative subreddits



LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

2. A closer look from the social science perspective





LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

2.1 The reddit population





Reddit population

- Reddit is anonymous + anyone can view most subreddits
- Posting, voting, or commenting on Reddit requires registration
- **Demographics of the 330 million active users are unknown**
 - Users are more likely to be male and younger
 - Pew:
 - 6% of all U.S. adults used Reddit
 - Men twice as likely as women to use
 - Reddit:
 - 79% of worldwide users are 18– 34 years of age



Reddit population

- First half of 2022, the **United States** accounted for 47.13 percent of traffic
- The **United Kingdom** was ranked second, accounting for 7.48 percent of web visits to the social media platform
- **Canada, Australia and Germany** are the next biggest countries



Reddit may be an appropriate source of data if a researcher is attempting to target a *rare population* but less ideal for *general population research*.

The definition of Redditors (e.g., Reddit posters, viewers/consumers, commenters) influences your findings



LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

2.2 Getting reddit data





Getting reddit data

- Two ways to access Reddit content
 - Via the Reddit API
 - PRAW for python
 - RedditExtractoR for R
 - Pushshift.io (maintained by redditor Jason Baumgartner, Stuck_In_The_Matrix)
 - Download via <http://files.pushshift.io/reddit/>
 - PSAW for python
 - pushshiftR for R



Getting reddit data via the Reddit API (python)

- API rules: <https://github.com/reddit-archive/reddit/wiki/API#rules>
- Any registered user may use the Reddit API, but they must fill out a special request and be authenticated. (This will take less than an hour.)
- The guide of the PRAW (Python Reddit API Wrapper) package provides documentation to start the process in Python
- Allows data collection from subreddits, threads, users
- Depending on the search being conducted and unit of analysis, there are limits to the number of units the system will return



Submission

```
class praw.models.Submission(reddit: praw.Reddit, id: Optional[str] = None, url: Optional[str] = None, _data: Optional[Dict[str, Any]] = None)
```

A class for submissions to Reddit.

Typical Attributes

Note

This table describes attributes that typically belong to objects of this class. PRAW dynamically provides the attributes that Reddit returns via the API. Since those attributes are subject to change on Reddit's end, PRAW makes no effort to document any new/removed/changed attributes, other than to instruct you on how to discover what is available. As a result, this table of attributes may not be complete. See [Determine Available Attributes of an Object](#) for detailed information.

If you would like to add an attribute to this table, feel free to open a [pull request](#).

Attribute	Description
<code>author</code>	Provides an instance of <code>Redditor</code> .
<code>clicked</code>	Whether or not the submission has been clicked by the client.
<code>comments</code>	Provides an instance of <code>CommentForest</code> .
<code>created_utc</code>	Time the submission was created, represented in Unix Time .
<code>distinguished</code>	Whether or not the submission is distinguished.
<code>edited</code>	Whether or not the submission has been edited.

<code>clicked</code>	Whether or not the submission has been clicked by the client.
<code>comments</code>	Provides an instance of <code>CommentForest</code> .
<code>created_utc</code>	Time the submission was created, represented in Unix Time .
<code>distinguished</code>	Whether or not the submission is distinguished.
<code>edited</code>	Whether or not the submission has been edited.
<code>id</code>	ID of the submission.
<code>is_original_content</code>	Whether or not the submission has been set as original content.
<code>is_self</code>	Whether or not the submission is a selfpost (text-only).
<code>link_flair_template_id</code>	The link flair's ID.
<code>link_flair_text</code>	The link flair's text content, or <code>None</code> if not flaired.
<code>locked</code>	Whether or not the submission has been locked.
<code>name</code>	Fullname of the submission.
<code>num_comments</code>	The number of comments on the submission.
<code>over_18</code>	Whether or not the submission has been marked as NSFW.
<code>permalink</code>	A permalink for the submission.
<code>poll_data</code>	A <code>PollData</code> object representing the data of this submission, if it is a poll submission.
<code>saved</code>	Whether or not the submission is saved.
<code>score</code>	The number of upvotes for the submission.
<code>selftext</code>	The submissions' selftext - an empty string if a link post.
<code>spoiler</code>	Whether or not the submission has been marked as a spoiler.





Getting reddit data via Reddit API

- Overview over options in R from APIs for Social Scientists (Bauer et al)
- https://bookdown.org/paul/apis_for_social_scientists/reddit-api.html
- Simple API calls via subreddit webaddresses *without authentication*
 - `https://www.reddit.com/r/cats/.json`
 - Possible to extract threads, upvotes and so on
- RedditExtractoR: Reddit Data Extraction Toolkit
 - Limited functionalities

```
find_subreddits . . . . .
find_thread_urls . . . . .
get_thread_content . . . . .
get_user_content . . . . .
```

Getting reddit data via pushshift download

- Contains all Reddit content
 - Data sets are stored as large compressed json files organized by month and by type of content (e.g., post or comment)
 - Data is copied into Pushshift at the time it is posted to reddit
- The data set structure and its contents changes over time as changes were implemented in Reddit
- Contains variables and content that may have since been removed or changed (not available on Reddit API)
- Current scores not available via pushshift

Getting reddit data via pushshift download

Directory Contents

Please consider making a donation (<https://pushshift.io/donations>) if you download a lot of data. This helps offset the costs of my time collecting data and providing bandwidth to make these files available to the public. Thank you!

If you have any questions about the data formats of the files or any other questions, please feel free to contact me at jason@pushshift.io.

Filename	Type	Size (bytes)	Date Modified
69M_reddit_accounts.csv.gz	69M_REDDIT_ACCOUNTS.CSV.GZ File	1,051,903,601	Sep 8 2021 6:10 AM
RA_2018-09.gz	RA_2018-09.GZ File	1,104,136,829	Sep 8 2021 6:11 AM
RA_2020-06-28.ndjson.zst	RA_2020-06-28.NDJSON.ZST File	1,933,724,618	Sep 8 2021 6:12 AM
RS_2019-09-01.gz	RS_2019-09-01.GZ File	256,537,945	Sep 8 2021 6:12 AM
authors	<Directory>	<Directory>	Jun 23 2022 2:09 AM
authors.dat.zst	AUTHORS.DAT.ZST File	1,444,191,053	Sep 8 2021 6:14 AM
comments	<Directory>	<Directory>	Oct 11 2022 11:47 AM
daily	<Directory>	<Directory>	Sep 8 2021 5:54 AM
ioel_data.ndjson.zst	JOEL_DATA.NDJSON.ZST File	105,415,750	Oct 19 2021 11:30 AM



Getting reddit data via pushshift API

- Pushshift also offer an API and there are API wrappers for R and python
 - Python package psaw
 - <https://psaw.readthedocs.io/en>
 - R Package pushshiftR
 - <https://github.com/nathancunn/pushshiftR>
- The Python package offers more options and is better documented.



Comparison of pushshift and Reddit API

Table 2. Differences Between Two Reddit Data Sources.

Downloadable Data Set (from http://files.pushshift.io/reddit/)	Reddit API
<ul style="list-style-type: none"> • Contains all Reddit content, making it larger than one may want or be able to process • Maybe some missingness (due to download error), but the rate of missing data has declined over time • The data set structure and its contents changes over time as changes were implemented in Reddit • Contains variables and content that may have since been removed or changed (not available on Reddit API) 	<ul style="list-style-type: none"> • Caps ability to pull data to 60 items per minute. An item would be a single post or a comment along with the metadata associated with that submission • More variables at more levels of analysis (e.g., thread, user, subreddit) are available through the API than the downloadable data set • Will be more consistent if downloaded all at once • Depending on the search being conducted and unit of analysis, there are limits to the number of units (i.e., data rows) the system will return (either 100, 1,000, or unlimited)

Amaya et al.
2019

Note. API = Application Programming Interface.



Combination of pushshift and Reddit API

- Search for content through Pushshift, save IDs of the content
- Forward IDs to the Reddit API
- Get updated (meta-) data through the Reddit API



Summary

- Reddit is large and will likely have sufficient data for several types of research objectives and questions.
- The sociodemographic characteristics of redditors are not comparable to the general population but insights into rare subpopulations might be possible.
- Pushshift and the Reddit API have different advantages and disadvantages for collecting Reddit data, but can also be combined.



LUDWIG-
MAXIMILIANS-
UNIVERSITÄT
MÜNCHEN

Thank you for your attention!

anna-carolina.haensch@stat.uni-muenchen.de

