# Collecting and Analyzing Reddit Data Best Practices

Christine Sowa
csowa@andrew.cmu.edu

**Carnegie Mellon**

**Center for Computational Analysis of Social and Organizational Systems**
http://www.casos.cs.cmu.edu/

---

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Agenda

- Overview of Reddit
- How to Get Data
- Importing into ORA

**11 June 2020**          **Christine Sowa**          **2**

# What is Reddit?

- Reddit is the 6th most popular website in the USA with users averaging 11 minutes and 28 seconds on the site every day.
- Globally it's the 20th most visited site in the world.
- Users are 71% male, and 59% are between the ages of 18 and 29.
- Users are highly reliant on the platform for news.
  - 45% of all Reddit users reported "learning something about the presidential campaign or candidates on the site in a given week"

CASOS

11 June 2020          Christine Sowa          3

# How do users interact with Reddit?

- Over a million distinct subcommunities, called subreddits, exist.
- Community members can 'upvote' or 'downvote' new content.
- 'Karma' is a sum of a user's post and comment scores.
- Posts can be 'gilded' by users for money.
- A post or comment's 'score' is the number of upvotes it receives minus its downvotes.

CASOS

11 June 2020          Christine Sowa          4

## Slide 5

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# What makes Reddit unique?

- Moderation
  - Each subreddit has moderators that enforce community standards for posts

CASOS

## Slide 6

**Carnegie Mellon**
ISr institute for SOFTWARE RESEARCH

# Example Interactions



CASOS

# The Reddit API

- First must read the terms and register to use the API
- API data format comes out as a JSON
  - One JSON per post or comment
- Can use wrappers (like praw or PushShift for Python).



```
import praw

reddit = praw.Reddit(client_id="my client id",
                     client_secret="my client secret",
                     user_agent="my user agent")
```

**11 June 2020** **Christine Sowa** **7**

# Type of Data to Pull

- Get all of the posts (Submissions) from a given subreddit from the past 30 days
  - Get post title, score, id, url, number of comments, author, score
- Get all posts from a given Redditor
- Obtain all comments to a set of posts
  - Get comment author, time, score, text

**11 June 2020** **Christine Sowa** **8**

## Reddit Networks

- User x Subreddit
- User x Post
- User x User
- …

## Walking through API using PushShift

```python
1   import pandas as pd
2   import requests
3
4   def get_pushshift_data(data_type, **kwargs):
5       base_url = f"https://api.pushshift.io/reddit/search/{data_type}/"
6       payload = kwargs
7       request = requests.get(base_url, params=payload)
8       return request.json()
9
10  data_type="submission"
11  query="coronaravirus|coronavirus|wuhan virus|wuhanvirus|2019nCoV|NCoV|NCoV2019|covid-19|covid19|covid 19"
12  size=1000
```

## Pulling Data with Pushshift

```
14    count = 0
15    all_data = []
16    while count <48:
17        count += 1
18        print(f"now printing hour in the past {count}")
19        hour = f"{count}h"
20        data = get_pushshift_data(data_type=data_type, before=hour,size=size, q=query)
21
22        all_data.extend(data["data"])
23
24        df = pd.DataFrame(all_data)
25        print(df)
26
27        pd.DataFrame(df).to_csv(f"corona_posts_{count}.csv", encoding="utf-8")
28
```

**11 June 2020**       **Christine Sowa**       **11**

## Uploading Data into Ora

▼ Import other data formats
- JSON data
- Twitter data
- Blogtrackers data
- GitHub data
- YouTube data
- Talkwalker data
- Pulse data
- NexaIntelligence data
- VK data
- Survey Monkey data
- Shapefile data
- Bibliography & Citations data
- TAVI data
- THINK data
- Reddit data

**11 June 2020**       **Christine Sowa**       **12**