

## Archiving the Web

Sofia Huang  
due 10/11/2022

### 1. Collect URIs from Tweets

- Write a Python program that collects English-language tweets that contain links.
- Write a Python program that extracts the links shared in tweets.
- Resolve all URIs to their final target URI (i.e., the one that responds with a 200).
- Save only unique final URIs (no repeats).
- if after this step, you don't have 1000 unique URIs, go back and gather more until you are able to get at least 1000 unique URIs
- Save this collection of 1000 unique links in a file and upload it to your repo in GitHub

I used the function from `get_tweets.py` from a previous homework to collect 12,500 English tweets using the keywords: 'Queen', 'hurricane', 'Ukraine', 'NASA', and 'Biden'. Then, I read the JSONL file, one line at a time, to obtain any links from each tweet. I also used regex to check if the link led to a Twitter page or any video/audio only site and did not store those links. I also used the requests library to resolve links to their final URI and stored them all in a .txt file. I ended up with 16,004 resolved links.

```
1 import sys
2 import json
3 from twarc import Twarc2, expansions
4 from configparser import ConfigParser
5 import re
6 from subprocess import call
7 import requests
8 import os
9 import time
10 import csv
11
12 TWARC_CONFIG_FILE = "/Users/sofiahuang/Library/Application Support/
    twarc/config"
13 OUTPUT_FILE = "/Users/sofiahuang/Documents/WM/FALL2022/DATA440/tweets.
    jsonl" # line-oriented JSON
14 MAX_TWEETS = 100
15
16 def get_tweets(input_search_term):
```

```
17     # read Twitter API keys from twarc config file, setup twarc2 object
18     config = ConfigParser(interpolation=None)
19     with open(TWARC_CONFIG_FILE) as twarc_config:
20         config.read_string("[TWARC]\n" + twarc_config.read())
21     bearer_token = config['TWARC']['bearer_token'].strip('\n')
22     t = Twarc2(bearer_token=bearer_token)
23
24     search_term = input_search_term
25
26     # limit search results to English language, with links, and no
    retweets
27     query = search_term + " lang:en has:links -is:retweet"
28
29     search_results = t.search_recent(query=query, max_results=
MAX_TWEETS)
30     # max_results is the number of results per page
31
32     num_tweets = 0
33     for page in search_results:
34         # use expansions.flatten to get all the information in a single
    JSON
35         result = expansions.flatten(page)
36
37         # open the file and write one JSON object per new line (jsonl
    format)
38         with open(OUTPUT_FILE, 'a+') as filehandle: # if you want to
    append, change 'w' to 'a+'
39             for tweet in result:
40                 print(tweet)
41                 filehandle.write('%s\n' % json.dumps(tweet))
42                 num_tweets = num_tweets + 1
43                 if num_tweets == MAX_TWEETS:
44                     # must include this to stop after a certain # of tweets
45                     search_results.close()
46
47     print (num_tweets, "tweets written to " + OUTPUT_FILE + " for query
    \"\" + query + "\"\n");
48
49 def get_links():
50     f = open(OUTPUT_FILE)
51     # read in all the lines
52     lines = f.readlines()
53     links = []
54     resolved_links = []
55     # each line is a json
56     for line in lines:
57         tweet_data = json.loads(line)
```

```

58     # collect links
59     if 'urls' in tweet_data['entities']:
60         for link in tweet_data['entities']['urls']:
61             # check if link leads to a Twitter or video/audio only
page using regex
62             pattern = re.compile('(https:\\/\\/)(www\\.\\.\\.*) (twitter|
youtube|youtu|tiktok|twitch|soundcloud) (\\.com|\\.be) (\\.*)')
63             m = pattern.match(link['expanded_url'])
64             if (m == None):
65                 # add to list if link leads to an acceptable page (
not Twitter or video/audio)
66                 links.append(link['expanded_url'])
67                 #print('original: ' + link['expanded_url'])
68                 try:
69                     # resolve links to their final URI
70                     resolved_link = requests.head(link['
expanded_url'], allow_redirects=True, timeout=5).url
71                     resolved_links.append(resolved_link)
72                     #print('final: ' + resolved_link)
73                 except:
74                     continue
75             print('Originally ' + str(len(links)) + ' links.')
76             print('Resolved ' + str(len(resolved_links)) + ' links.')
77             os.chdir("/Users/sofiahuang/Documents/WM/FALL2022/DATA440")
78             # create txt file to store the resolved links
79             resolved_uri_file = os.path.join(os.getcwd(), 'resolved_uris.txt')
80             for link in resolved_links:
81                 try:
82                     with open(resolved_uri_file, 'a+') as f:
83                         f.write("\n%s" % link)
84                         print(link)
85                 except Exception as e:
86                     print(e)
87             return resolved_links

```

**Listing 1:** Collecting tweets and extracting links

```

1 if __name__ == "__main__":
2     search_terms = ['Queen', 'hurricane', 'Ukraine', 'NASA', 'Biden']
3     for term in search_terms:
4         for i in range(25):
5             get_tweets(term)
6
7     final_links = get_links()

```

**Listing 2:** Running get\_tweets() and get\_links()

Then, I used the Unix tools, sort and uniq to save only the unique URIs from resolved\_uris.txt to a separate .txt file. I ended up with 2,146 unique URIs stored.

```
1 (base) sofiahuang@Sofias-MacBook-Pro DATA440 % sort resolved_uris.txt |
   uniq > unique_uris.txt
2 (base) sofiahuang@Sofias-MacBook-Pro DATA440 % wc -l unique_uris.txt
3      2146 unique_uris.txt
```

**Listing 3:** Keeping only unique URIs

## 2. Get TimeMaps for Each URI

**Obtain the TimeMaps for each of the unique URIs from Q1 using the ODU Memento Aggregator, MemGator.**

I read the .txt file with the unique URIs into a list and created a function to run MemGator, locally, on each of the URIs. I stored the timemaps I obtained in folder and ended up with 1,072 timemaps. I was unable to get timemaps for all of the unique URIs I had stored and I believe this has something to do with the MemGator server and maybe it didn't work for all of them.

```
1 def read_unique_links_file(txt_filename):
2     # check if we are in the correct directory
3     os.chdir("/Users/sofiahuang/Documents/WM/FALL2022/DATA440")
4     summary_filename = os.path.join(os.getcwd(), txt_filename)
5     # read the resolved links txt file into a list to be returned
6     try:
7         summary_file = open(summary_filename, 'r')
8         url_data = summary_file.read()
9         unique_list = url_data.split("\n")
10        summary_file.close()
11        print(str(len(unique_list)))
12        return unique_list
13    except Exception as e:
14        print(e)
15        return []
16
17 def get_timemaps(list, directory):
18     # to keep track of how many URIs timemaps are obtained for
19     # and to name the timemaps json file returned from MemGator
20     link_index = 1
21     for uri in list:
22         filename = str(link_index) + ".json"
23         # create a folder to store the timemaps
24         if os.path.exists(os.path.join(directory, filename)):
25             # file has already been created, just increment index and
26             # go to the next file
27             print("Skipping url " + uri + ", file already exists")
28             link_index += 1
29             continue
```

```

29         else:
30             # run MemGator on each URI in the unique links list
31             print('running memgator on: ' + uri)
32             time.sleep(1)
33             os.system('/Users/sofiahuang/Downloads/memgator-darwin-
amd64 -F 2 -f JSON ' + uri + ' > /Users/sofiahuang/Documents/WM/
FALL2022/DATA440/timemaps/' + filename)
34             time.sleep(2)

```

**Listing 4:** MemGator to get timemaps for all unique URIs

```

1 if __name__ == "__main__":
2     # ...more function calls before, shown in previous question
3     links = read_unique_links_file('unique_uris.txt')
4     os.chdir("/Users/sofiahuang/Documents/WM/FALL2022/DATA440")
5     timemap_directory = os.path.join(os.getcwd(), 'timemaps')
6     if not os.path.exists(timemap_directory):
7         os.mkdir(timemap_directory)
8     get_timemaps(links, timemap_directory)

```

**Listing 5:** Running functions to get timemaps for the URIs

### 3. Analyze Mementos Per URI-R

**Use the TimeMaps you saved in Q2 to analyze how well the URIs you collected in Q1 are archived. Create a table showing how many URI-Rs have certain number of mementos.**

To store the number of mementos for each URI, I created a dictionary. I iterated through all of the timemap json files in the directory and got the number of mementos archived for each URI. If the json file was zero bytes, it had 0 mementos and I stored that information in the dictionary but did not try to read the file as this would cause an error. I then saved this dictionary as a .csv file so I could easily see the data obtained.

```

1 def read_json():
2     # create dictionary and add entry for URI-Rs with 0 mementos
3     mem_dict = {0:0}
4     os.chdir("/Users/sofiahuang/Documents/WM/FALL2022/DATA440/timemaps"
)
5     file_count = 0
6     # iterate over json files
7     for json_file in os.listdir("/Users/sofiahuang/Documents/WM/
FALL2022/DATA440/timemaps"):
8         # open json file
9         f = open(json_file, 'r')
10        if (json_file == '.DS_Store'):
11            continue

```

```
12     # if file is 0 bytes, don't read, just add data to dictionary
13     if (os.stat(json_file).st_size == 0):
14         current_val = mem_dict.get(0)
15         mem_dict.update({0: current_val+1})
16         file_count+=1
17         continue
18     # read json file
19     data = json.loads(f.read())
20     # get number of mementos and add to dictionary,
21     # checking if key is already in the dictionary
22     mementos = data['mementos']['list']
23     mem_count = len(mementos)
24     if mem_count in mem_dict:
25         current_val = mem_dict.get(mem_count)
26         mem_dict.update({mem_count: current_val+1})
27     else:
28         mem_dict.update({mem_count: 1})
29     f.close()
30     file_count+=1
31     print(file_count)
32
33     # open file for writing dictionary to csv file
34     w = csv.writer(open("/Users/sofiahuang/Documents/WM/FALL2022/
DATA440/memento_counts.csv", "w"))
35     # loop over dictionary keys and values
36     for key, val in mem_dict.items():
37         # write every key and value to file
38         w.writerow([key, val])
```

**Listing 6:** Find how many mementos obtained from each URI and save to CSV file

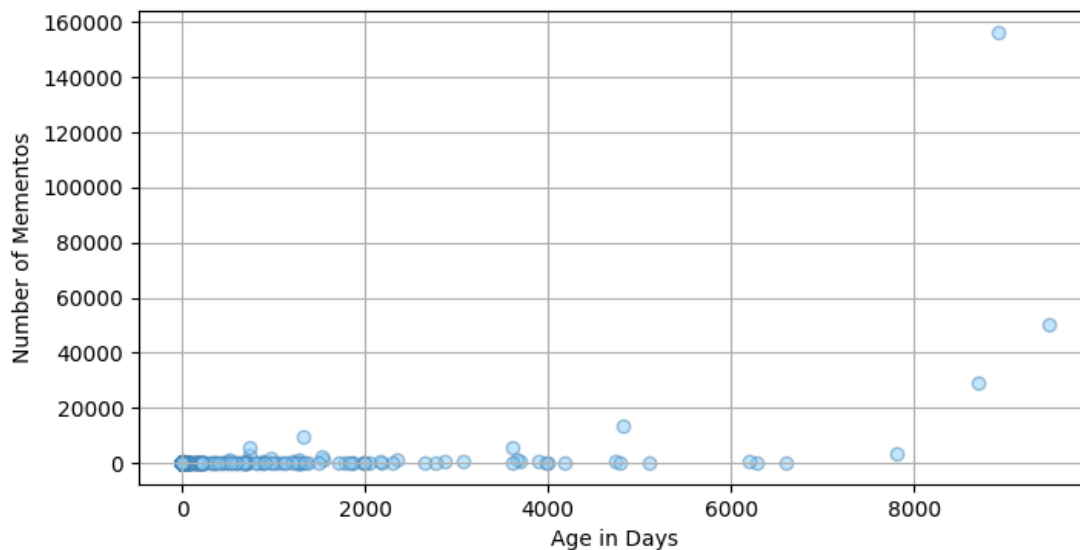
I created a table to show how many URIs had a certain number of mementos. We can see that the majority had 100 or less mementos archived, with 449 URIs having 0. The number drastically decreases as the number of mementos increases. With less than 50 URIs having from 101-1000 mementos. I had 12 URIs that had over 1000 mementos archived, with memento counts ranging from 1,273 to 156,369. I checked to see what webpage had the highest number of mementos and it was a [dailymail.co.uk](http://dailymail.co.uk) webpage from 1998.

Mementos	URI-Rs
0	449
1-100	568
101-200	17
201-300	11
301-400	7
401-500	0
501-600	2
601-700	2
701-800	0
801-900	1
901-1000	3
> 1000	12

## 4. Analyze Datetimes of Mementos

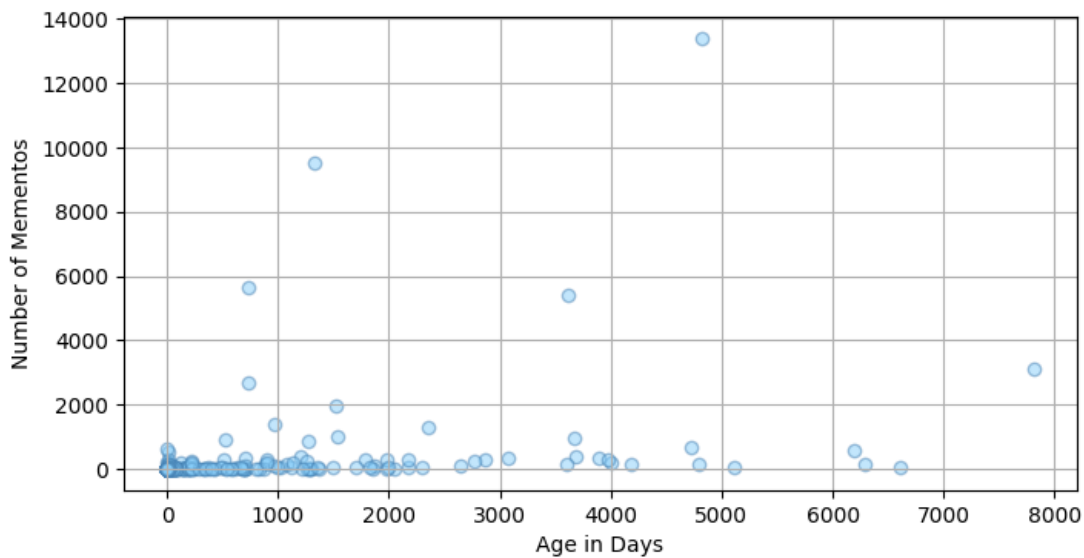
For each of the URI-Rs from Q3 that had  $> 0$  mementos, create a scatterplot with the age of each URI-R (today - earliest memento datetime) on the x-axis and number of mementos for that URI-R on the y-axis.

Figure 1 shows the age and number of mementos for each URI-R.



**Figure 1:** Age vs Memento Counts for URI-Rs

Figure 2 shows the age and number of mementos for each URI-R without outliers.



**Figure 2:** Age vs Memento Counts for URI-Rs \*excluding URI-Rs with more than 20,000 mementos\*

I iterated through the memento json files and stored the age (using the datetime library), memento count, and URI-R in lists for each URI-R. I created a dataframe out of the 3 lists and then stored the dataframe as a .csv file. Then, I wrote another function to create a scatterplot (using Matplotlib). After I looked at the resulting scatterplot, I decided to create another one without the outliers or the URI-Rs that had greater than 20,000 mementos so that the relationship between age and memento count could be seen more clearly for the lower values.

**Q: What can you say about the relationship between the age of a URI-R and the number of its mementos?**

Typically, the older URI-Rs have more mementos than newer ones, however the correlation isn't very strong and there are certainly a few outliers. Most of the URI-Rs were only a couple of years old and had less than 1,000 mementos. I included a scatterplot of the same data just excluding the URI-Rs that had over 20,000 mementos so you could see the majority of the data better.

**Q: What URI-R had the oldest memento? Did that surprise you?**

The URI-R with the oldest memento was '<http://archive.md/19961101130030/http://www.redcross.org/>' and it was 9,471 days old or almost 26 years old. The Red Cross has been around for a long time and is a well-known organization so it doesn't surprise me that their webpage has been archived well.

**Q: How many URI-Rs had an age of  $< 1$  week, meaning that their first memento was captured the same week you collected the data?**

175 URI-Rs had an age of less than a week. I looked at the .csv file to figure this out.



```
1 def get_memento_age():
2     # create dictionary and add entry for URI-Rs with 0 mementos
3     mem_dict = {0:0}
4     os.chdir("/Users/sofiahuang/Documents/WM/FALL2022/DATA440/timemaps"
5     )
6     file_count = 0
7     memento_counts = []
8     memento_ages = []
9     uris = []
10    # iterate over json files
11    for json_file in os.listdir("/Users/sofiahuang/Documents/WM/
12    FALL2022/DATA440/timemaps"):
13        # open json file
14        f = open(json_file, 'r')
15        # if file is 0 bytes, don't read
16        if (os.stat(json_file).st_size == 0):
17            continue
18        # read json file
19        data = json.loads(f.read())
20        # get earliest date
21        memento_earliest_date_str = data['mementos']['list'][0]['
22        datetime']
23        memento_earliest_date = datetime.strptime(
24        memento_earliest_date_str[0:10], '%Y-%m-%d')
25        memento_age = (datetime.today() - memento_earliest_date).days
26        # add data to lists
27        memento_ages.append(memento_age)
28        memento_counts.append(len(data['mementos']['list']))
29        uris.append(data['mementos']['list'][0]['uri'])
30        f.close()
31        file_count+=1
32    # zip lists in dataframe and store as .csv file
33    df = pd.DataFrame(list(zip(uris, memento_ages, memento_counts)))
34    df = df.rename(columns={"0": "uri", "1": "age", "2": "num_memento"
35    })
36    df.to_csv("/Users/sofiahuang/Documents/WM/FALL2022/DATA440/
37    memento_ages_counts.csv", index=False)
38
39 def memento_age_scatterplot(df, filename):
40     # rename columns
41     df.rename(columns={"0": "uri", "1": "age", "2": "num_memento"},
42     inplace=True)
43     plt.figure(figsize=(8, 4))
44     # create scatterplot, add grid, axis labels, and save plot
45     plt.scatter(df['age'], df['num_memento'], c="lightskyblue", alpha
46     =0.5, edgecolors='steelblue')
47     plt.grid()
```

```
40 plt.xlabel('Age in Days')
41 plt.ylabel('Number of Mementos')
42 plt.savefig(filename)
```

**Listing 7:** Create scatterplot to show age of mementos

## References

- StackOverflow - Pretty Print JSON <https://stackoverflow.com/questions/20265439/how-can-i-pretty-print-a-json-file-from-the-command-line>
- MemGator API <https://memgator.cs.odu.edu/api.html>
- Python parse JSON file <https://www.freecodecamp.org/news/python-parse-json-how-to-read-a-json-file/>
- Python Save Dictionary to File <https://pythonspot.com/save-a-dictionary-to-a-file/>
- Check if Key Exists in Dictionary <https://www.geeksforgeeks.org/python-check-whether-given-key-already-exists-in-a-dictionary/>