# HW2 - Web Archiving

Brenden Lewis
Due: 10/11/2020 11:59PM

# Q1

When observing Twitter posts through the API, I decided to opt for posts that contained the string 'Youtube' in an attempt to keep my scope to YouTube video links. The Tweepy and Request libraries were utilized to extract the links from each Twitter post. I was not able to figure out how to check for duplicate URIs in an efficient manner, but I was able to extract about 1,200 total links related to YouTube in some capacity. Each twitter post streamed in was parsed to JSON format for observation.

I did not keep an exact count of the number of shortened links that redirected elsewhere, but observing all the links I was collecting both in testing and in my final collection were largely comprised of redirecting links. I created a recursive function called *def unshorten_url(short _url)* to "unshorten" these links and return the final destination after however many redirects are needed.

Due to this method being called for every link to ensure it is the final location, the connection to the API was often breaking due to the large amount of processing happening during the stream. In order to alleviate this issue, I had to collect the links in batches to ensure a stable connection; after a good deal of testing, I had the most consistent success with a maximum of 70 links collected in a single run. With each run of the program, I added the collected links to a text file *TwitterURLs.txt*.

```
1  import tweepy
2  from tweepy.streaming import StreamListener
3  from tweepy import OAuthHandler
4  from tweepy import Stream
5  import sys
6  import re
7  import json
8  import datetime
9  import csv
10 import urllib.request
11 import requests
12 import urllib.parse
13 import http.client
14 from http.client import IncompleteRead
15
16 # Keys ommitted
17 consumer_key="***"
18 consumer_secret="***"
19 access_token="***"
20 access_secret="***"
```

```python
21
22 # Handles authorization with Twitter
23 auth = OAuthHandler(consumer_key,consumer_secret)
24 auth.set_access_token(access_token,access_secret)
25 api = tweepy.API(auth)
26
27 link_contain = set() # Collection for uri links
28
29 class StdOutListener(StreamListener):
30
31     def __init__(self):
32         super().__init__()
33         self.max_tweets = 70
34         self.tweet_count = 0
35
36     def on_data(self,data):
37                 if self.tweet_count == self.max_tweets:
38                     print ("Finished")
39                     export_data(link_contain)
40                     return False
41                 else:
42                     try:
43                         tweet = json.loads(data)
44                         for url in tweet["entities"]["urls"]:
45                             tweet_url = str(url["expanded_url"])
46                             r = requests.head(tweet_url,
    allow_redirects=True).url
47                             full_url = unshorten_url(r)
48
49                             if "https://twitter" in full_url:
50                                 print("----Twitter URL ommited----")
51                             else:
52                                 self.tweet_count += 1
53                                 print ("URL"+str(self.tweet_count)+": "
    + full_url)
54                                 link_contain.add(full_url)
55
56                                 return True
57                     except Exception:
58                         print("Some issue occurred; Ending process")
59                         export_data(link_contain)
60                         return False
61
62
63     def on_error(self,status):
64         print (status)
65
```

```python
66  # If given a shorten url, expand to final url
67  def unshorten_url(short_url):
68      parsed = urllib.parse.urlparse(short_url)
69
70      if(parsed.scheme == 'https'):
71          h = http.client.HTTPSConnection(parsed.netloc)
72      else:
73          h = http.client.HTTPConnection(parsed.netloc)
74
75      resource = parsed.path
76      if parsed.query != "":
77          resource += "?" + parsed.query
78      h.request('HEAD', resource )
79      response = h.getresponse()
80      if response.status/100 == 3 and response.getheader('Location'):
81          return unshorten_url(response.getheader('Location')) #
        Recurisive call if not final location
82      else:
83          return short_url
84
85  def export_data(collection):
86      url_storage=open("TwitterURLs.txt","a")
87      for url in collection:
88          url_storage.write(url+"\n")
89      url_storage.close()
90
91
92
93
94  if __name__ == '__main__':
95      l = StdOutListener()
96      stream = Stream(auth,l)
97
98      stream.filter(track=['Youtube']) # Filter Twitter posts to thosse
        that contain the string 'YouTube'
```

**Listing 1:** Python code for Q1

# Q2

I opted to create a Python script to run MemGator from the command line for each link I had collected and return their TimeMaps in JSON format. To make creating the graphs easier after processing, I created a *uri_dict* dictionary for each entry containing the URI, number of mementos for that URI, and the age of its first memento. I capped my amount of processing at 1000 URIs worth of data. The program had to run continuously for about 10 hours to collect the

TimeMap for each link; each link took an average of about 30s - 1:10 mins to process. After all
URIs are processed, the program immediately creates several graphs from the data provided by the
TimeMaps. Listing 3 shows the terminal output of the program which prints each URI dictionary
and the total number of links processed and the number of links that had at least 1 memento.

```python
1  import os
2  import subprocess
3  import json
4  import matplotlib.pyplot as plt
5  import numpy as np
6  import pandas as pd
7  import datetime
8  from datetime import date
9
10 timeMaps = []    # Holds TimeMaps for each URI
11 non_zero_data = []   # Holds TimeMaps for each URI with more than 0
       mementos
12 tot_url_count = 0    # Total number of URIs
13 tot_mem_count = 0    # Total number of URIs with mementos
14
15 def collect_Timemaps():
16     with open("testUrls.txt") as f:
17         for line in f:
18             global tot_url_count
19             print("#"+str(tot_url_count+1)+": Collecting from: " + line
       )
20             result = subprocess.run(['memgator', '-f', 'JSON', '--
       hdrtimeout=1m30s', line], shell=True, capture_output=True)
21             output = result.stdout
22             global tot_mem_count
23
24             if output == None or output == b'': # URI has no mementos
       and memgator call returns nothing
25                 memento_count = 0
26                 age_days = 0
27                 tot_url_count += 1
28                 print("Done")
29             else:
30                 output_json = json.loads(output)    #Parses memgator
       output to a JSON object
31                 write_timeMaps_to_file(output_json)
32                 memento_count = count_mementos(output_json)
33                 age_days = calculate_age(output_json)
34                 tot_url_count += 1
35                 tot_mem_count +=1
36                 print("Done")
37
```

```python
38            uri_dict = {"uri":line, "mementos": memento_count, "age":
     age_days}  #Dictionary for URIs
39            if memento_count > 0:
40                non_zero_data.append(uri_dict)
41            else:
42                timeMaps.append(uri_dict)
43
44
45
46
47 def write_timeMaps_to_file(output):
48     outfile = open("timeMaps.txt","a")
49     outfile.write(json.dumps(output, indent=2, sort_keys=True))
50     outfile.close()
51
52 # Counts the number of mementos a link has
53 def count_mementos(output):
54     mementos = len(output['mementos']['list'])
55     return mementos
56
57 # Calculate the age of the first memento a link contains in days
58 def calculate_age(output):
59     first_mem = output['mementos']['first']['datetime'] # 2020-10-01T03
     :07:29Z
60
61     collection = datetime.date.today().strftime('%Y-%m-%d')
62     datetime_coll = datetime.datetime.strptime(collection,'%Y-%m-%d')
63
64     first_date = datetime.datetime.strptime(first_mem,'%Y-%m-%dT%H:%M:%
     SZ').strftime('%Y-%m-%d')
65     datetime_first = datetime.datetime.strptime(first_date,'%Y-%m-%d')
66
67     age = datetime_coll - datetime_first
68     return age.days
69
70 def print_timeMaps():
71     global tot_url_count
72     global tot_mem_count
73     for item in timeMaps:
74         print(item)
75         # tot_url_count+=1
76
77     print("\nTotal Urls: "+str(tot_url_count))
78     print("\nTotal Urls w/ Mementos: "+str(tot_mem_count))
79
80
81 def create_histogram():
```

```
82      output_data = []
83      for diction in timeMaps:
84          output_data.append(diction['mementos'])
85
86      non_zero_output = []
87      for item in non_zero_data:
88          non_zero_output.append(item['mementos'])
89
90      range = (0,1000)
91      bins = 20
92
93      plt.hist(output_data,bins,range,color="red",histtype='bar', rwidth
        =0.8)
94      plt.xlabel("# of Mementos")
95      plt.ylabel("Frequency")
96      plt.title("Q2: Total URIs")
97      plt.show()
98      plt.savefig('HistogramAll.png')
99
100     # Excludes URIs with zero mementos
101     plt.hist(non_zero_output,bins,range,color="blue",histtype='bar',
        rwidth=0.8)
102     plt.xlabel("# of Mementos")
103     plt.ylabel("Frequency")
104     plt.title("Q2: Total URIs with > 0 mementos")
105     plt.show()
106     plt.savefig('HistogramNonZero.png')
107
108 def create_scatter():
109     x = []
110     y = []
111     for diction in non_zero_data:
112         x.append(diction['mementos'])
113         y.append(diction['age'])
114
115     plt.xlabel('# of mementos')
116     plt.ylabel('Age in days')
117     plt.title("Q3 Total URIs")
118     plt.scatter(x,y, c="blue")
119     plt.show()
120     plt.savefig('Scatterplot.png')
121
122 if __name__ =='__main__':
123     collect_Timemaps()
124     create_histogram()
125     create_scatter()
```

```
126    print_timeMaps()
```

**Listing 2:** Python code for Q2 and Q3

```
1  {'uri': 'https://www.youtube.com/watch?v=5Zk_Vz-wz4g&feature=youtu.be\n
      ', 'mementos': 0, 'age': 0}
2  {'uri': 'https://www.youtube.com/watch?v=bFUCiVcGK9g&feature=youtu.be\n
      ', 'mementos': 1, 'age': 7}
3  {'uri': 'https://www.youtube.com/watch?v=3ej_bkYiQnE&feature=youtu.be\n
      ', 'mementos': 0, 'age': 0}
4  {'uri': 'https://www.youtube.com/watch?v=Vm7l3YxEF8g&feature=youtu.be\n
      ', 'mementos': 1, 'age': 2}
5  {'uri': 'https://www.youtube.com/watch?v=0zdS8hNM-sc&feature=youtu.be\n
      ', 'mementos': 0, 'age': 0}
6  {'uri': 'https://exe.io/LgNUIAt\n', 'mementos': 0, 'age': 0}
7  {'uri': 'https://www.youtube.com/watch?v=5hfYJsQAhl0', 'mementos':
      1361, 'age': 3871}
8
9  Total Urls: 1000
10
11 Total Urls w/ Mementos: 496
```

**Listing 3:** Terminal output after program execution

Listing 4 below shows the JSON output format for each link; this example in particular is from the very first link in the url text file. The program counts each memento listed under *list* (if available) and the datetime of the first memento to calculate the age. The output for each link is stored collectively in a seperate textfile labeled *timeMaps.txt*. The final size of this file came out to about 33MB, with many of the URIs having several thousands of mementos.
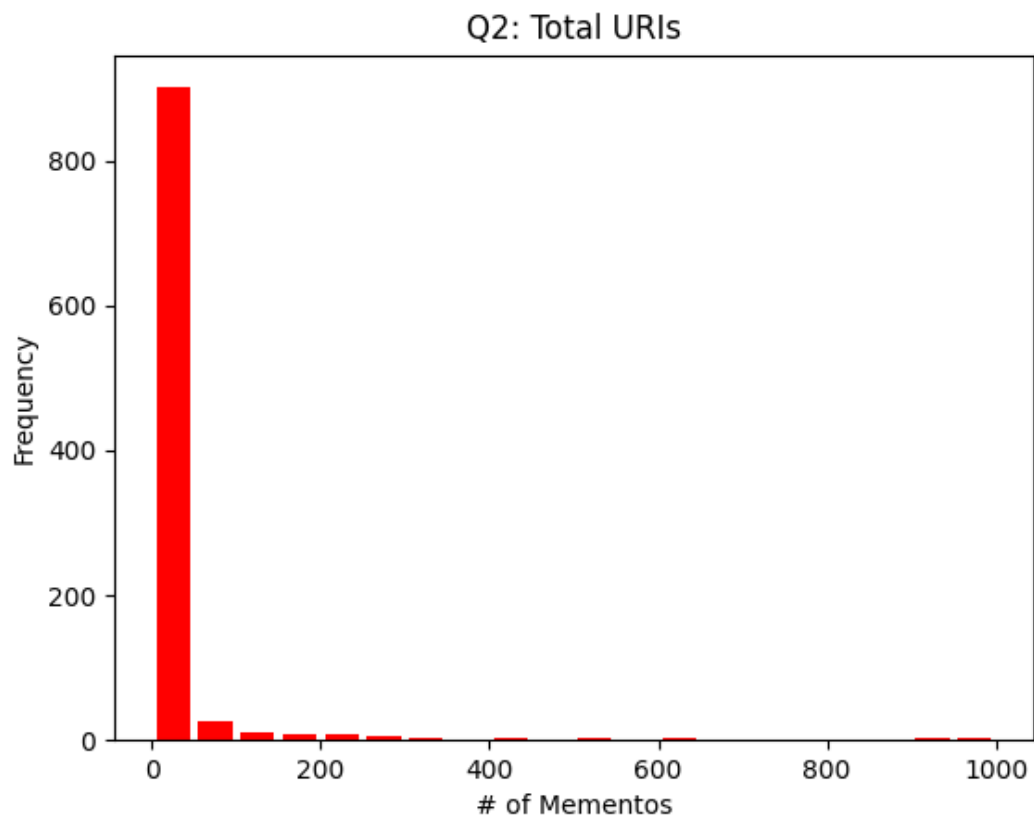
```
1  {
2    "mementos": {
3      "first": {
4        "datetime": "2020-10-03T13:47:32Z",
5        "uri": "https://web.archive.org/web/20201003134732/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ&gl=US&hl=en"
6      },
7      "last": {
8        "datetime": "2020-10-03T13:47:32Z",
9        "uri": "https://web.archive.org/web/20201003134732/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ&gl=US&hl=en"
10     },
11     "list": [
12       {
13         "datetime": "2020-10-03T13:47:32Z",
14         "uri": "https://web.archive.org/web/20201003134732/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ&gl=US&hl=en"
15       }
16     ]
```

```
17    },
18    "original_uri": "https://www.youtube.com/watch?v=4yAO5hbnfQQ",
19    "self": "http://localhost:1208/timemap/json/https://www.youtube.com/
      watch?v=4yAO5hbnfQQ",
20    "timegate_uri": "http://localhost:1208/timegate/https://www.youtube.
      com/watch?v=4yAO5hbnfQQ",
21    "timemap_uri": {
22      "cdxj_format": "http://localhost:1208/timemap/cdxj/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ",
23      "json_format": "http://localhost:1208/timemap/json/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ",
24      "link_format": "http://localhost:1208/timemap/link/https://www.
      youtube.com/watch?v=4yAO5hbnfQQ"
25    }
26 }
```
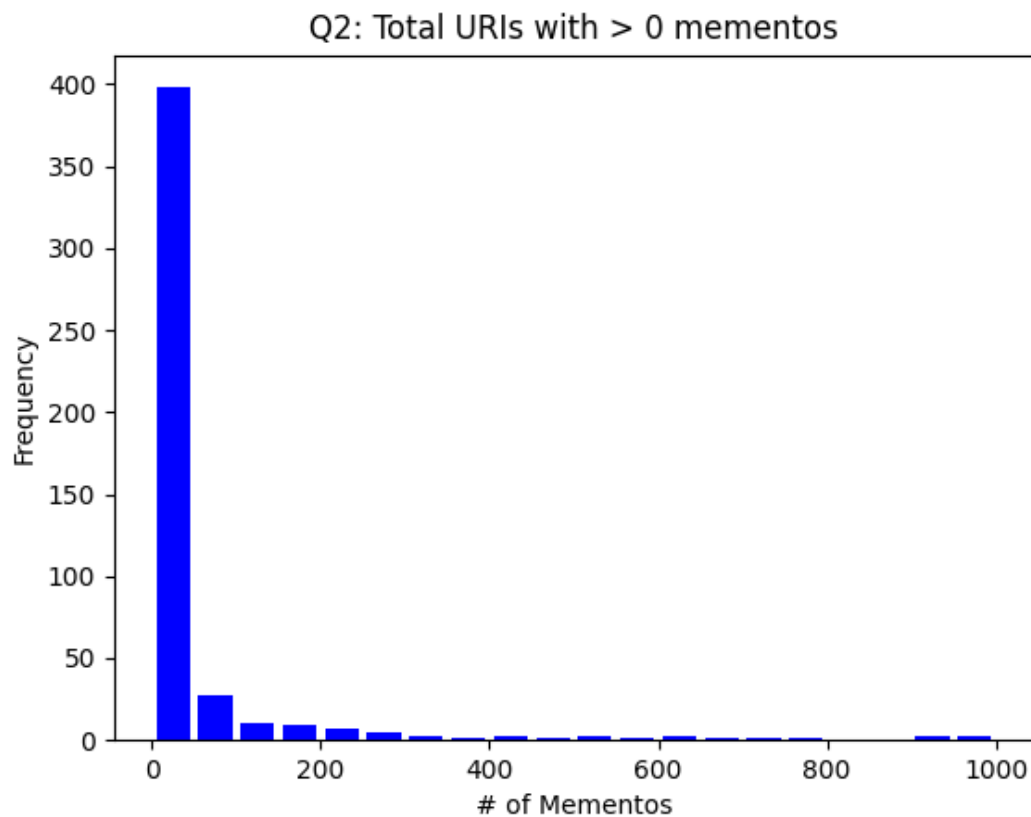
**Listing 4:** TimeMap of first link

Figures 1 and 2 below show histograms of the number of mementos from each URI and their frequency. Figure 1 shows the data from all URIs while Figure 2 attempts to show the same data excluding any URIs with 0 mementos. As also indicated in the program output, almost exactly half of the 1000 URIs observed had 0 mementos. Of the 496 links remaining, the large majority had very few mementos ranging between 1-100 mementos. Outside of a few outliers with thousands of mementos, most of the observed links are not well archived.
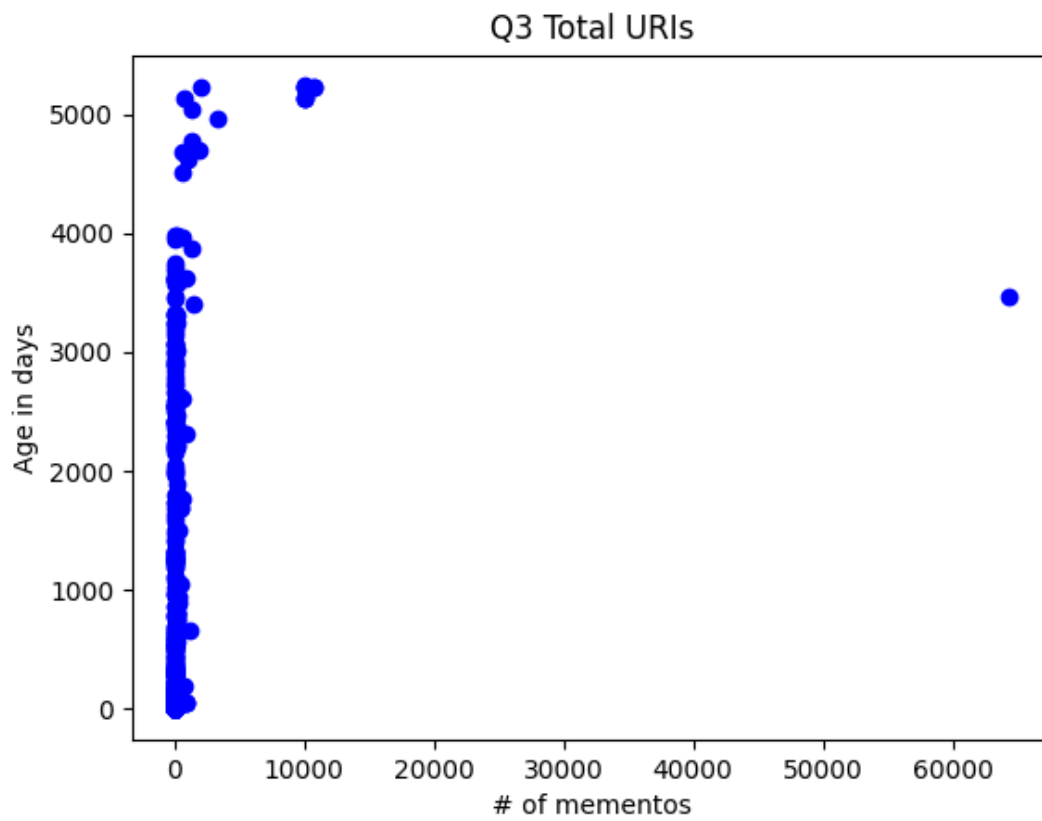
**Figure 1:** Histogram showing the number of URIs vs. the number of mementos

**Figure 2:** Histogram showing the data from Figure 1 excluding URIs with no mementos

# Q3

For each URI with non-zero memento counts, the age (in days) of the first memento is calculated through the *def calculate _age(output)* function that utilizes the *datetime* library and takes the TimeMap of a URI in JSON as an argument. Figure 3 below is a scatterplot of all the URIs with non-zero mementos. Each URI is plotted by the number of mementos its TimeMap contains and the age of its first memento. Unfortunantly, due to an unidentified error, there is an outlier plotted that heavily skews the graph to the right. Observing the rest of the plotted data provides convincing evidence to the hypothesis that, for YouTube links, the number of mementos a URI contains directly correlates with its age.

**Figure 3:** Scatterplot with the number of mementos vs. their age