

Assignment 2

CS 595: Introduction to Web Science

Fall 2013

Shawn M. Jones

Finished on September 26, 2013

1

Question

Write a Python program that extracts 1000 unique links from Twitter. You might want to take a look at:

<http://thomassileo.com/blog/2013/01/25/using-twitter-rest-api-v1-dot-1-with-python/>

But there are many other similar resources available on the web. Note that only Twitter API 1.1 is currently available; version 1 code will no longer work.

Also note that you need to verify that the final target URI (i.e., the one that responds with a 200) is unique. You could have different shortened URIs for www.cnn.com.

You might want to use the search feature to find URIs, or you can pull them from the feed of someone famous (e.g., Tim O'Reilly).

Hold on to this collection -- we'll use it later throughout the semester.

Answer

Extracting 1000 unique links from Twitter was a more complex endeavor than expected.

It required the use of 4 distinct scripts run in the following order:

1. `gimmetweets.py` <starting tweet id><tweet count per call><# of api calls ><screen name>
 - Starts at a given tweet ID and works backward, printing out all tweets for the given screen name, making only a certain number of api calls to ensure Twitter does not lock out the application.
 - The intention is to run this once for each screen name for which we want to gather tweets, then pass the output to a file via redirection ('>')
2. `extractURIsFromTweets.py` <filename #1 containing tweets><filename #2... >
 - Runs through all of the tweets, extracting URIs that match a simple URI pattern of "`http://[\S]*`"
 - Each URI is dereferenced, passing through all redirects. If a status code of 200 comes back, the URI is saved for the next step.
3. `combineLists.py` <filename #1 containing uris><filename #2... >
 - Combine all of the files gathered together in the last step, sort the entries, and eliminate duplicates.
4. `extractRandom1000Links.py` <input file containing links>
 - Because more than 1000 links were acquired, this script randomly generates a representative sample of 1000 links from the file generated in the last step, preventing needless hours processing more than 1000 links in the rest of this assignment.

The 4 scripts allowed resumption of processing if an error state was encountered at any point. For example, `gimmetweets.py` could be cut off because it reached Twitter's hourly request limit. Because it was time consuming (in minutes, sometimes hours) to complete some of these tasks, splitting the scripts up in this way allowed work to be saved across runs. Even

if there was some overlap, `combineLists.py` took care of any extra URIs that were processed repeatedly.

Early testing of the 'expanded_url' field provided by the Twitter API showed that even though a URI was present in the tweet text, sometimes it did not appear as a value in 'expanded_url', so it was inconsistent. To account for this, all tweets were extracted by `gimmetweets.py` and later processed with `extractURIsFromTweets.py`.

The following screen names were followed for data:

- BarackObama - President of the United States
- BillGates - Former Microsoft CEO, now attempting to rid the world of Malaria, in addition to other feats
- CNN - original 24 hour news channel
- GeorgeTakei - Former helmsman on Star Trek, now Internet celebrity
- MaddowBlog - left leaning Rhodes Scholar political pundit Rachel Maddow
- SethMacFarlane - writer and comedian, best known as the creator of the TV series *Family Guy*
- dailykos - faaaaar left leaning news source
- nprnews - National Public Radio

In retrospect, I could have saved a lot of time with `extractURIsFromTweets.py` by using parallel execution. The program could have broken a given list of tweets up and then executed against each sublist, combining the whole once all were done.

Parallel execution would have saved little for `gimmetweets.py` because of the number of requests limited by Twitter.

The following sections show the source code of each of these scripts.

gimmetweets.py

```
1 #!/usr/local/bin/python3
2
3 # -*- encoding: utf-8 -*-
4 from __future__ import unicode_literals
5 import requests
6 from requests_oauthlib import OAuth1
7 from urllib.parse import parse_qs
8 import json
9 import time
10 import sys
11
12 # ugly, but necessary, globals; saw no need to change this
13 # strategy from the example
14 REQUEST_TOKEN_URL = "https://api.twitter.com/oauth/request_token"
15
16 AUTHORIZE_URL = "https://api.twitter.com/oauth/authorize?"
17     oauth_token="
18 ACCESS_TOKEN_URL = "https://api.twitter.com/oauth/access_token"
19
20 CONSUMER_KEY = "n7jt1uMTwGCcIzDvey8g0A"
21 CONSUMER_SECRET = "0r6HUrVD36W4MULgWETKMxrQsCICNy1OFFNc2iW4o"
22
23 OAUTH_TOKEN = "528649269-
24     SffJ0Rei5PzLYd2NSJPnnm28dP5nlAnt7E1gRGwo"
25 OAUTH_TOKEN_SECRET = "htrwXF09pS8tP8cMzFrxmMryavdPXd0zPiJHRnLs"
26
27 class APIError(Exception):
28     """
29     If something goes wrong with the API, throw one of these
30     .
31     (avoids sys.exit in the middle of the program)
32     """
33
34     def __init__(self, value):
35         self.value = value
36
37     def __str__(self):
38         return repr(self.value)
39
40 def setup_oauth():
41     """
42     Authorize your app via identifier.
43     Code inspired by:
44     http://thomassileo.com/blog/2013/01/25/using-twitter-
45     rest-api-v1-dot-1-with-python/
46     """
```

```

42
43     # Request token
44     oauth = OAuth1(CONSUMER_KEY, client_secret=CONSUMER_SECRET)
45     r = requests.post(url=REQUEST_TOKEN_URL, auth=oauth)
46
47     credentials = parse_qs(r.content)
48
49     resource_owner_key = credentials[b'oauth_token'][0].decode(
50         encoding='UTF-8')
51     resource_owner_secret = credentials[b'oauth_token_secret']
52     ] [0].decode(encoding='UTF-8')
53
54     # Authorize
55     authorize_url = AUTHORIZE_URL + resource_owner_key
56     print('Please go here and authorize: ' + authorize_url)
57
58     verifier = input('Please input the verifier: ')
59     oauth = OAuth1(CONSUMER_KEY,
60                     client_secret=CONSUMER_SECRET,
61                     resource_owner_key=resource_owner_key,
62                     resource_owner_secret=resource_owner_secret,
63                     verifier=verifier)
64
65     # Finally, Obtain the Access Token
66     r = requests.post(url=ACCESS_TOKEN_URL, auth=oauth)
67     credentials = parse_qs(r.content)
68     token = credentials[b'oauth_token'][0].decode(encoding='UTF
69     -8')
70     secret = credentials[b'oauth_token_secret'][0].decode(
71         encoding='UTF-8')
72
73     return token, secret
74
75 def get_oauth():
76     """
77     Code inspired by:
78     http://thomassileo.com/blog/2013/01/25/using-twitter-
79     rest-api-v1-dot-1-with-python/
80     """
81     oauth = OAuth1(CONSUMER_KEY,
82                     client_secret=CONSUMER_SECRET,
83                     resource_owner_key=OAUTH_TOKEN,
84                     resource_owner_secret=OAUTH_TOKEN_SECRET)
85     return oauth
86
87 def call_api(ident, oauth, count, screenName):
88     #url = \

```

```

86 # "https://api.twitter.com/1.1/statuses/home_timeline.
    json?max_id=" + \
87 # str(ident) + "&count=" + str(count)
88 url = \
89 "https://api.twitter.com/1.1/statuses/user_timeline.json
    ?screen_name=" + \
90 screenName + "&count=" + str(count) + "&max_id=" + str(
    ident)
91
92 r = requests.get( url, auth=oauth )
93
94 if 'errors' in r:
95     raise APIError(
96         json.dumps(
97             r.json(), sort_keys=True, indent=4, separators=(
98                 ', ', ': ' ) )
99         )
100
101 return r
102
103 def print_tweets_with_ids(
104     startingid, oauth, tweetCountPerCall, numberOfCalls,
105     screenName):
106     """
107     Make the API call, print the tweets returned.
108     """
109
110     ident = startingid
111
112     for i in range(0, numberOfCalls):
113
114         response = call_api(ident, oauth, tweetCountPerCall,
115                             screenName)
116
117         tweetListSize = len(response.json())
118
119         print(tweetListSize)
120
121         if tweetListSize == 1:
122             raise APIError("Ran out of tweets?")
123
124         for i in range(0, tweetListSize):
125             tweet = response.json()[i]
126             ident = str(tweet['id'])
127             text = tweet['text']
128             print( str(ident) + ":" + str(text) )

```

```

129 def usage():
130
131     print("Usage: " + sys.argv[0] +
132           " <startingid> <tweetCountPerCall> <apiCalls> <
              screenName>")
133
134
135 if __name__ == "__main__":
136
137     #startingid = "400000000000000000"
138     try:
139         startingId = sys.argv[1]
140         tweetCountPerCall = int(sys.argv[2])
141         apiCalls = int(sys.argv[3])
142         screenName = sys.argv[4]
143     except IndexError as e:
144         usage()
145         sys.exit(1)
146
147     if not OAUTH_TOKEN:
148         token, secret = setup_oauth()
149         print("OAUTH_TOKEN: " + token)
150         print("OAUTH_TOKEN_SECRET: " + secret)
151         print()
152     else:
153         oauth = get_oauth()
154
155         try:
156             print_tweets_with_ids(
157                 startingId, oauth, tweetCountPerCall, apiCalls,
158                 screenName)
159         except APIError as e:
160             sys.stderr.write(e.value)
161             sys.exit(254)

```

Listing 1: Python program for acquiring tweets for a given screen name

extractURIsFromTweets.py

```
1  #!/usr/local/bin/python3
2
3  import sys
4  import re
5  import urllib.request
6
7  # for IRI support, thank you evitan for your 0-point answer that
   was helpful:
8  # http://stackoverflow.com/questions/4389572/how-to-fetch-a-non-
   ascii-url-with-python-urlopen
9  import httpplib2
10
11  URLPATTERN = re.compile("http://[\\S]*")
12
13  def extractURIsFromLine(line):
14      """
15          Attempts to extract the URIs from the string given.
16
17          Unfortunately, it requires the global URLPATTERN for
   performance.
18      """
19
20      return URLPATTERN.findall(line)
21
22
23
24  def extractRealSupportedURI(uri):
25      """
26          Returns "real" URI if it survives redirects and returns
   a 200.
27
28          Returns None otherwise.
29      """
30
31      realURI = None
32
33      try:
34          # this function follows the URI, resolving all redirects
   ,
35          # and detects redirect loops
36          # iri2uri is needed for IRIs
37          request = urllib.request.urlopen(httpplib2.iri2uri(uri))
38
39          if request.getcode() == 200:
40              realURI = request.geturl()
41
```

```

42     except urllib.error.HTTPError as e:
43         # something went wrong, we don't care what
44         realURI = None
45
46     except urllib.error.URLError as e:
47         # something went wrong, we don't care what
48         realURI = None
49
50     except UnicodeError as e:
51         # something went very wrong with the IRI decoding
52         realURI = None
53
54     return realURI
55
56 if __name__ == "__main__":
57
58     filenames = sys.argv[1:]
59
60     for filename in filenames:
61         f = open(filename)
62
63         for line in f:
64             sys.stderr.write("Working on: " + line + '\n')
65             uris = extractURIsFromLine(line)
66
67             for uri in uris:
68                 sys.stderr.write("Trying URI: [" + uri + ']\n')
69                 goodURI = extractRealSupportedURI(uri)
70
71                 if goodURI:
72                     print(goodURI)
73
74         f.close()

```

Listing 2: Python program for extracting URIs from a file of tweets

combineLists.py

```
1 #!/usr/local/bin/python3
2
3 import sys
4
5 # this file performs the equivalent of cat <filenames> | sort |
6   uniq
7
8 inputFilenames = sys.argv[1:]
9
10 urilist = []
11
12 # combine all of the entries from the given files together
13 for filename in inputFilenames:
14     f = open(filename)
15     urilist.extend(f.readlines())
16     f.close()
17
18 # sort the list
19 urilist.sort()
20
21 # expensive, but agreed upon simplest way to "uniq" a list
22 urilist = list(set(urilist))
23
24 # sort the list again
25 urilist.sort()
26
27 # output!
28 for entry in urilist:
29     print(entry.strip())
```

Listing 3: Python program for combining, sorting, and uniquing several files

extractRandom1000Links.py

```
1 #!/usr/local/bin/python3
2
3 # I realized that I didn't want to spend time waiting for more
4   than 1000 links,
5 # so this script gives a random sample
6 import random
7 import sys
8
9 inputfile = sys.argv[1]
10
11 f = open(inputfile)
12 links = f.readlines()
13 f.close()
14
15 count = 1000
16
17 selectedLinks = []
18
19 while count > 0:
20     # grab a random entry from the list
21     index = random.randint(0, len(links) - 1)
22     newlink = links[index].strip()
23
24     if newlink not in selectedLinks:
25         selectedLinks.append(newlink)
26         count -= 1
27
28 for link in sorted(selectedLinks):
29     print(link)
```

Listing 4: Python program for randomly extracting 1000 lines from a given file

2

Question

Original:

Download the TimeMaps for each of the target URIs. We'll use the ODU Memento Aggregator, so for example:

URI-R = `http://www.cs.odu.edu/`

URI-T = `http://mementoproxy.cs.odu.edu/aggr/timemap/link/http://www.cs.odu.edu/`

Create a histogram of URIs vs. number of Mementos (as computed from the TimeMaps). For example, 100 URIs with 0 Mementos, 300 URIs with 1 Memento, 400 URIs with 2 Mementos, etc.

Answer

The Data Collection

The script `countMementos.py` acquires the TimeMap for each link from a file containing many links, then parses that TimeMap and counts the number of mementos present. Requests that generate a 404 are recorded as having 0 mementos. The mementos are counted using a regular expression.

Because the first run of this script took almost 10 hours to complete, I used the Unix `split` command to divide the 1000 links into 5 files of 200 links each. I then ran the script against each file in turn. Running these in parallel drastically reduced the run time. In fact, the original 10 hour run was on Sunday afternoon. The second run was on Monday night during dinner, which took 1 hour.

The source code for doing the dereferencing and processing of each Time Map is shown in Listing 5.

Once the 5 files were acquired, they were combined together using the Unix `cat` command.

The Results

The histogram in Figure 1 on page 14 is difficult to read. This histogram was generated in R using bins of size 1. This was done to keep the y-axis low, because there are 375 URIs with 0 mementos and 244 URIs with 1 memento, totaling more than half of all data recorded.

Outliers also skew the results. Two outliers have 22,157 and 10543 mementos. Stripping out these outliers gets us to Figure 2 on page 15, which shows a slightly better visualization, but there is still such a jump between the low number of mementos to those at 3000 that it is still difficult to see a pattern.

If we focus on the highest number of records, mementos less than 10, the plot actually begins to look like a histogram in Figure 3 on page 16, and we can see the precipitous drop between 0 and 2.

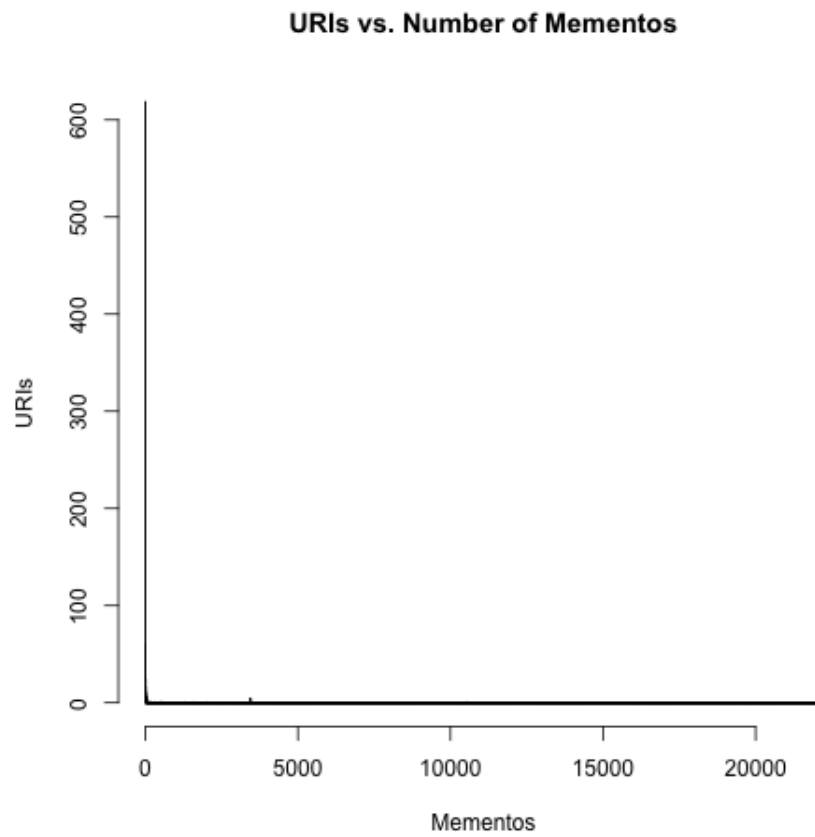


Figure 1: Histogram of URIs vs. number of Mementos for the entire dataset

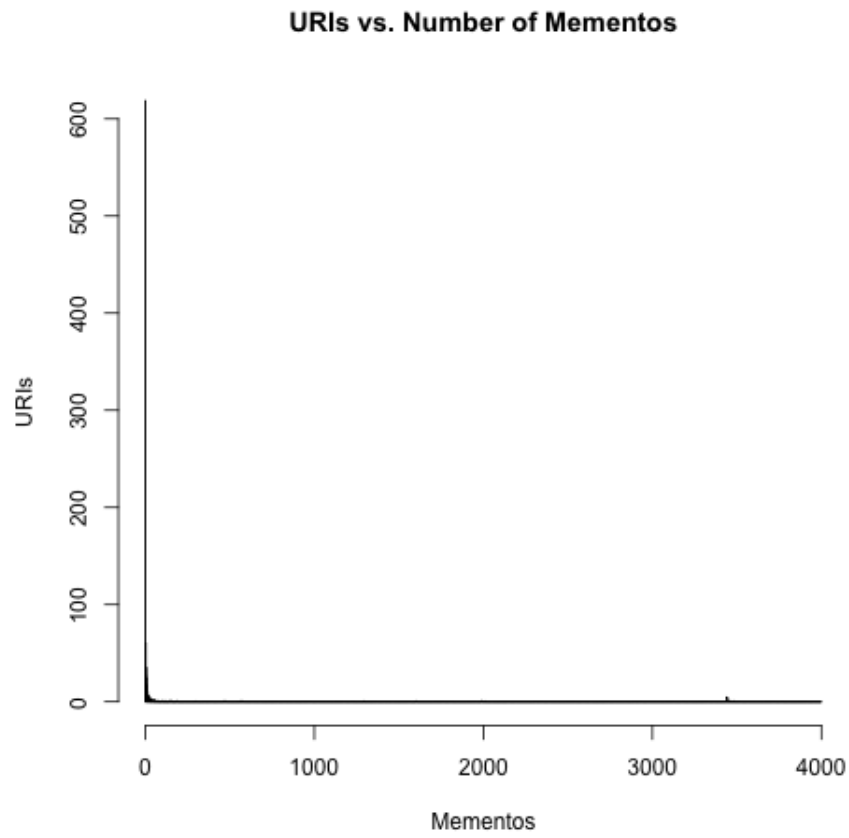


Figure 2: Histogram of URIs vs. number of Mementos for URIs with less than 10,000 Mementos

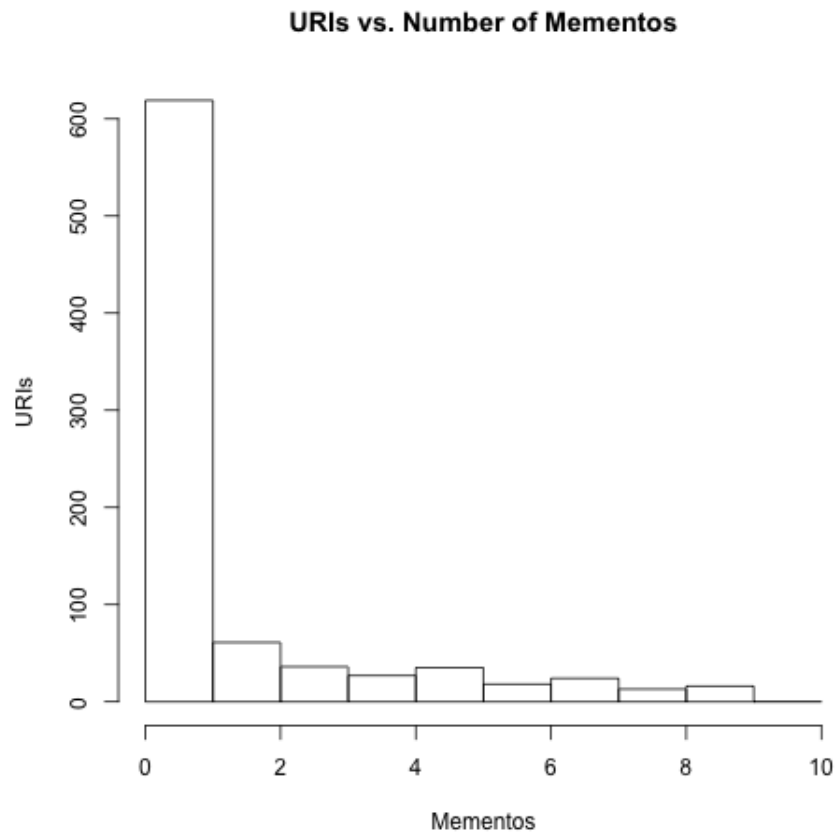


Figure 3: Histogram of URIs vs. number of Mementos for URIs with less than 10 Mementos

```

1  #!/usr/local/bin/python3
2
3  import re
4  import sys
5  import urllib.request
6
7  MEMENTOPATTERN = re.compile(r'rel="[^"]*memento["]*"')
8
9  def getTimeMap(uri):
10
11     urit = "http://mementoproxy.cs.odu.edu/aggr/timemap/link/" +
12         uri
13
14     try:
15         request = urllib.request.urlopen(urit)
16
17         if request.getcode() == 200:
18             timemap = request.readall()
19             request.close()
20         else:
21             timemap = None
22             request.close()
23
24     except urllib.error.HTTPError as e:
25         timemap = None
26
27     except urllib.error.URLError as e:
28         timemap = None
29
30     return timemap
31
32 def countMementos(uri):
33
34     urit = getTimeMap(uri)
35
36     if not urit:
37         count = 0
38     else:
39         count = len(MEMENTOPATTERN.findall(str(urit)))
40
41     return count
42
43
44 if __name__ == "__main__":
45     inputfile = sys.argv[1]
46
47     f = open(inputfile)

```

```

48
49     for uri in f:
50
51         mementoCount = countMementos(uri.strip())
52
53         print(str(mementoCount) + "\t" + uri.strip())
54         sys.stdout.flush()
55
56     f.close()

```

Listing 5: Python program for processing Time Maps for a given file full of links

3

Question

Estimate the age of each of the 1000 URIs using the "Carbon Date" tool:

<http://ws-dl.blogspot.com/2013/04/2013-04-19-carbon-dating-web.html>

Note: you'll have to download the tool and install; don't try to use the web service.

For URIs that have > 0 Mementos and an estimated creation date, create a graph with age (in days) on one axis and number of mementos on the other.

Answer

The Data Collection

Since I did not have Python 2.6 readily available, I modified Carbon Date to work with Python 2.7 and shared it with the author so he could share with others.

From some test runs, it became apparent that the Carbon Date tool takes between 1 and 7 minutes to query all of its services for a given URI. A worst case scenario yields:

$$\frac{7 \text{ minutes}}{1 \text{ URI}} \times 1000 \text{ URIs} = 7000 \text{ minutes} \times \frac{1 \text{ day}}{1440 \text{ minutes}} \approx 5 \text{ days}$$

I modified the Carbon Date tool to accept a configuration file as an argument (also sharing with the author), then ran it 5 times across 5 different screen sessions using 5 different configuration files containing 5 different port numbers.

Then I did the same for the Carbon Date tool client shown in listing 6.

Here's an example run of the Carbon Date client using a subset of the 1000 links in a file named **xae**.

```
./queryCarbonDateFor1000Links.py http://ganesh:8547/cd xae > xae
.data
```

The Results

Because, as shown in Figure 3, the majority of the memento counts are between 0 and 2, most of the data points in Figure 4 are on the left side, with a large number of mementos being roughly less than 1200 days old. There are also two outliers, our friends with 22,157 and 10543 mementos, who seem to be much older. In a perfect world, the whole graph would follow the pattern of these two outliers, with the number of mementos increasing from the age of first creation.

Then again, because the majority of the URIs are “young”, and it takes some time to get archived, this graph may make sense, especially for twitter postings.

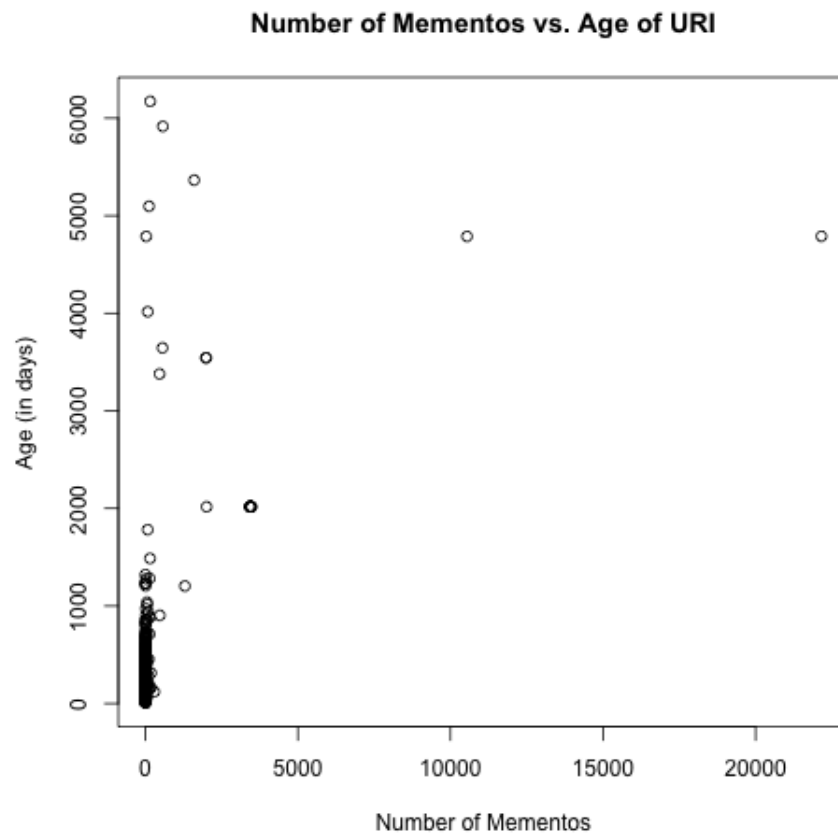


Figure 4: Number of Mementos vs. Age of URI

```

1  #!/usr/local/bin/python3
2
3  import sys
4  import urllib.request
5  import json
6
7  def queryCarbonDate(cduri, uri):
8
9      if cduri[-1] != '/':
10         cduri += '/'
11
12     sys.stderr.write("Using cduri = " + cduri + "\n")
13     sys.stderr.write("Requesting " + cduri + uri + "\n")
14     sys.stderr.flush()
15
16     request = urllib.request.urlopen(cduri + uri)
17     pagedata = request.readall().decode('utf-8')
18     request.close()
19
20     data = json.loads(pagedata)
21
22     return data['Estimated Creation Date']
23
24 if __name__ == '__main__':
25     cduri = sys.argv[1]
26     urifile = sys.argv[2]
27
28     f = open(urifile)
29
30     for line in f:
31         cdate = queryCarbonDate(cduri, line.strip())
32         print(cdate + "\t" + line.strip())
33         sys.stdout.flush()
34
35     f.close()

```

Listing 6: Python program for remotely querying the Carbon Date tool