# Assignment 2

## CS 595: Introduction to Web Science
### Spring 2016
### Manoj Chandra Kompalli
### Finished on February 11,2016

# Contents

# 1   Question 1

1.  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 many different
shortened URIs for www.cnn.com (t.co, bit.ly, goo.gl, etc.).

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.

## 1.1 Answer

I have started off by searching for APIs for twitter. I have found Twitter search and Tweepy.I chose Tweepy because it came with python and looked easy to implement. I have then created app in twitter to generate keys and tokens for authentication.

I had extracted tweets with keyword news. I converted the response to JSON. I have extracted tweet id and link for the tweet.

I have used expanded url property of the tweet object to get the full url. I have generated all the urls upto 1000 which match the keyword "news" into output.JSON file.

The reason I chose a json file over a text file is , because of the readibility of the file and also because pulling json data is pretty easy .I have used tweet id and link as keys.The links I generated were expanded links . Output.json is a huge file containing 1000 lines. Hence, I pulled out first 19 lines from the Output.json file below.

## 1.2 Code Listing

**links.py**

```
 1  import tweepy
 2  import json
 3  import time
 4  import sys
 5  import re
 6  import urllib2
 7  # Authentication Keys to Connect to Twitter API
 8  consumer_key="vjemit5xYQdhgrEPa1FeFf5ZO";
 9  consumer_secret="0
        PoNwIkHk29kUweChIIhGzVD3ZfXwRVqqwxYY3zPadY1BZeNq8";
10  access_token="3485785534-4FlFNlJtg1uNAlMummglVi9feR7fyvkUS0STp0G
        ";
11  access_token_secret="
        EpzAYEpf6tHdGFKj43HhnBeAhLNgkyXPdZyH72ec8Ew8d";
12  auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
13  auth.set_access_token(access_token, access_token_secret)
14
15  outputFile = open('output.json', 'w')  #opens the file with
        write permissions
16
17  all_urls = set()#for unique data
18  api = tweepy.API(auth) #Accessing tweepy API
19  searches = tweepy.Cursor(api.search,q="news").items()#Querying
        for news related tweets
```

```
20
21  while True:                                                    #
           Infinite  loop  through  tweets
22      try:
23          tweet = searches.next()# iterating  through  all  the
                matched  tweets
24          item= tweet._json# converting  the  tweet  object  to  JSON
                object
25          myitem={}# declaring  an  empty  array  to  store  the  details
26
27          tweet_id= item['id_str']
28          myitem['tweet_id'] = tweet_id# fetching  the  id  of  tweet
                and  storing  in  JSON  object
29          #myitem['createdDate'] = created_date
30          for link in item['entities']['urls']:
31                          all_urls.add(link['url'])
32
33                          myitem['link']=link['expanded_url']#
                                expanded  url  for  full  url
34                          outputFile.write(json.dumps(myitem) +
                                '\n')# writing  JSON  data  to  an
                                output  JSON  file  line  by  line
35
36
37
38
39
40
41
42
43          if len(all_urls) == 1000:                          #Checks
                for  1000  urls  in  the  list  and  breaks  out  if  more
44              break
45      except tweepy.TweepError:# catching  tweepy  error  which  which
                occurs  frequently  enough
46          time.sleep(60*10)
47          continue
48      except StopIteration:
49          break
```

Listing 1: Python program for getting 1000 uri's from queried tweets

**reduced.json**

```
1  {"tweet_id": "697576352514494466", "link": "http://rol.st/1
       TSvdBH"}
2  {"tweet_id": "697576352514494466", "link": "http://twitter.com/
       RollingStone/status/697575085268463616/photo/1"}
3  {"tweet_id": "697576352447393792", "link": "http://dailym.ai/1
       o30Ahg"}
4  {"tweet_id": "697576352397062144", "link": "http://goo.gl/fb/
       QnSejs"}
5  {"tweet_id": "697576352396935168", "link": "http://entabe.jp/
       news/gourmet/10454/dandelion-chocolate-opens-in-japan"}
6  {"tweet_id": "697576352271233024", "link": "https://www.
       londontheatre1.com/news/127554/how-the-other-half-loves-at-
       the-theatre-royal-haymarket/"}
7  {"tweet_id": "697576352225062912", "link": "http://rol.st/1
       TSvdBH"}
8  {"tweet_id": "697576352225062912", "link": "http://twitter.com/
       RollingStone/status/697575085268463616/photo/1"}
9  {"tweet_id": "697576352149405696", "link": "http://bowenpress.
       com/news/bowen_66161.html"}
10 {"tweet_id": "697576352099266562", "link": "http://bit.ly/1
       Q68sEy"}
11 {"tweet_id": "697576352073953281", "link": "http://www.
       theguardian.com/australia-news/2016/feb/11/reuters-distances-
       itself-from-greg-hunt-best-minister-award-it-wasnt-our-idea?
       CMP=share_btn_tw"}
12 {"tweet_id": "697576351801307136", "link": "https://gleam.io/3
       ck3D/gamma-glider-giveaway"}
13 {"tweet_id": "697576351797108736", "link": "http://dd.hokkaido-
       np.co.jp/news/area/doto/1-0233547.html"}
14 {"tweet_id": "697576351486779392", "link": "http://bit.ly/
       iHrAwards"}
15 {"tweet_id": "697576351344123904", "link": "http://wpo.st/tqZA1"
       }
16 {"tweet_id": "697576351314784257", "link": "http://www.
       starobserver.com.au/news/local-news/new-south-wales-news/
       talking-turkey-a-groundbreaking-chat-on-lgbti-parenting
       /145764"}
17 {"tweet_id": "697576351042138112", "link": "http://yahoo.jp/
       g9zlYf"}
18 {"tweet_id": "697576350756909056", "link": "http://bit.ly/1
       o03tP4"}
19 {"tweet_id": "697576350509621248", "link": "http://shrd.by/
       Atawls"}
```

Listing 2: JSON data of extracted urls from tweets

# 2 Question 2

2. 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/1/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.

* = https://en.wikipedia.org/wiki/Histogram

## 2.1 Answer

Each time map could have many mementos.  First thing,I did was to navigate
to the Time map url.  Then ,it downloaded a file which gave me the mementos
of a single url cs.odu.edu.  By using regular expression to locate rel mementos
told me if the url had a memento or not.  The memcount.py mines for
mementos and returns 2 output files which are very useful for the histogram
to be plotted next and also for the carbon dating tool.  I realized that
a file with just an array of memento counts would be sufficient to plot a
histogram.The file memcount.py reads the urls from output.json file which
was generated in the previous program and writes all mementos of different
urls to memcount.json file.It also seperates a list of memento counts and a
list of url counts of which have more than 0 mementos and writes the output
to two seperate files.This is useful for the carbon dating program.  I have
found that out of 1000 urls only 33 urls had mementos.

The next part is taking the counts generated and plotting a histogram.
I have scaled the y axis to 10 because most of the urls have zero memen-
tos.Some urls have mementos in the range of 0-7000.

They are very less in number and due to this, it is very difficult to
represent them in the graph.  In the next histogram I have limited the
mementos to 600. We can now clearly see the variation in the frequency of
urls and mementos.In the last plot, I have introduced breaks to clearly show
how many urls have a good number of mementos.

## 2.2 Code Listing

```python
1  #!/usr/local/bin/python3
2  import re
3  import sys
4  import urllib2
5  import json
6
7  mymementos = re.compile(r'rel.*?=.*?"memento".*?')#use regular
       expressions to find mementos
8  file3=open('abovezerocounts.json','w')
9  file4=open('abovezerourls.json','w')
10
11 def getTimeMap(url):
12     mem_url = "http://mementoproxy.cs.odu.edu/aggr/timemap/link
           /1/" + url #plug in the url to a timemap
13     try:
14         response = urllib2.urlopen(mem_url)
15         timemap = response.read()
16     except urllib2.HTTPError:
17         timemap = None
18     return timemap
19
20 def countMementos(mem_url):
21         time_map = getTimeMap(mem_url)
22         if not time_map:# if no time maps
23                 count=0
24         else:
25
26
27                 count=len(mymementos.findall(str(time_map)))#
                       finds the count of all mementos per url
28                 if count>0:
29                         file3.write("%s\n"% count)#writes the
                               count of urls onto a json file
30                         file4.write("%s\n"% time_map)#writes all
                               the urls on to a json file
31                 #print count
32         return count
33
34 if __name__=="__main__":
35         file1=open('output.json','r')# input a json file that
               contains 1000 urls
36         file2=open('memcount.json','w')
37         #memcountlist=[]
38         for line in file1.readlines():
39                 one_line = json.loads(line)# loads a json object
40                 link = one_line['link']
```

```
41                      counter=countMementos(link)# counter has count
                          of the urls
42                      file2.write("%s"% counter)#outputs count of
                          mementos of each url to a json file
43                      file2.write("\r\n")
44              #for item in memcountlist:
45
46
47  file1.close()
48  file2.close()
```

Listing 3: Python program for processing Time Maps for a given file full of links

### 2.2.1 Code1

```
1  d = read.table('memcount.json',col.name=c("mementos"))
2  hist(d$mementos,xlim=c(0,7000),ylim=c(0,10),breaks=500,col=5,
       main="URIs vs Mementos",ylab="No. of URI's",xlab="Mementos")
```

Listing 4: R program for generating the last histogram for Question 2

## 2.3 Results

Here,we can see that by limiting URI's to 10 and Mementos to 600 ,all the mementos which fall under 600 visible clearly.

The following graph has breaks introduced to make it clear that how many umber of urls those many mementos.Especially in the region of 0-20
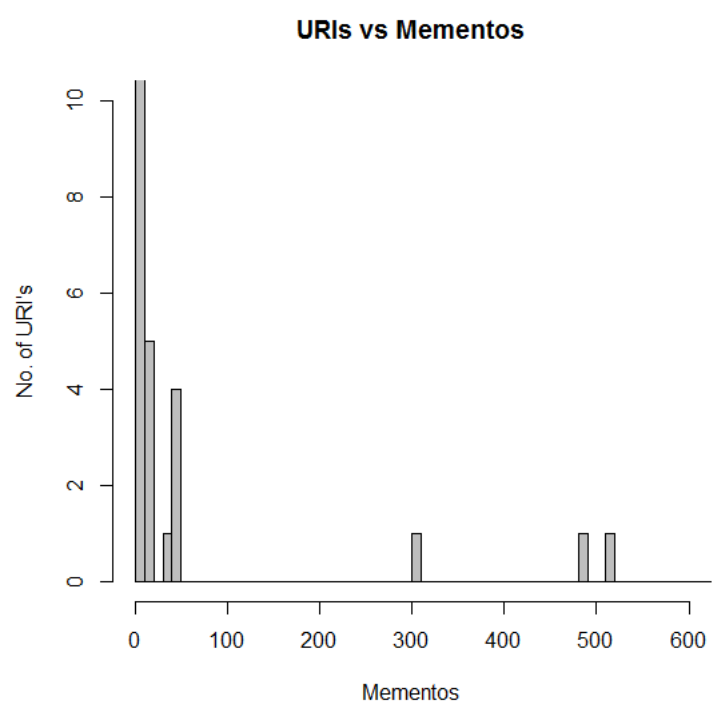
Figure 1: Histogram of URIs vs. number of Mementos for URIs with less than 600 Mementos
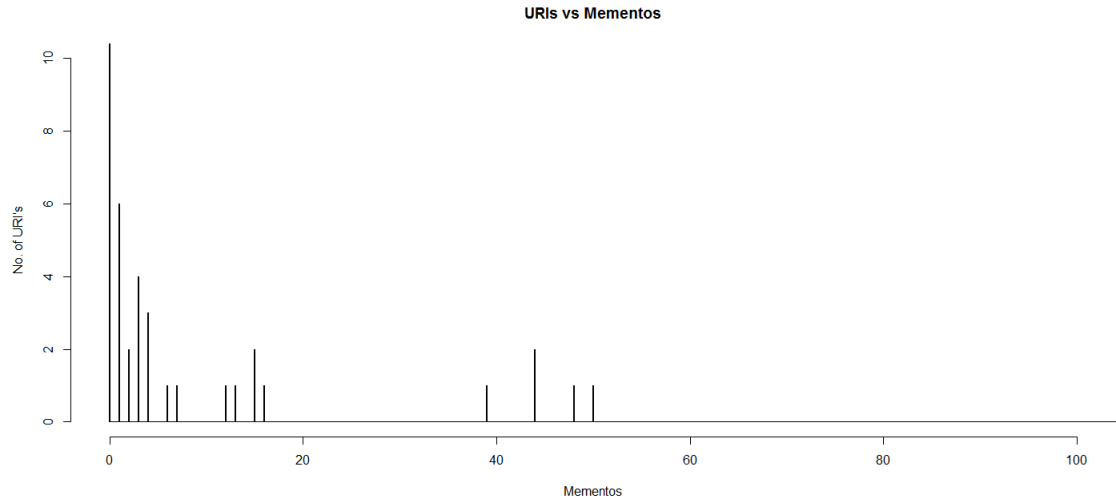
Figure 2: Histogram of URIs vs. number of Mementos for URIs with less than 100 Mementos and breaks inserted between each memento

# 3  Question 3

```
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.
```

## 3.1 Answer

"Carbon Date" webservice gives the created date,modified date etc of a url .One url at a time.The tool however, can be used to run over a list of urls.

Here we need all the urls which generate atleast one memento.We can get them using memcount.py of problem 2. Our first goal is to find the created date of the urls. File local.py outputs the created dates of all the 33 urls which were having atleast one memento.

We can use these dates and subtract them from the current date which gives the total number of days or Age of each uri.days.py does this job.

Now, we have number of days in one file(numberdays.txt) and number of mementos of each file in another url.We can plot a scatter plot with days on y axis and urls on x axis.The

```python
#!/usr/local/bin/python3
import re
import sys
import urllib2
import json

mymementos = re.compile(r'rel.*?=.*?"memento".*?')#use regular
    expressions to find mementos
file3=open('abovezerocounts.json','w')
file4=open('abovezerourls.json','w')

def getTimeMap(url):
    mem_url = "http://mementoproxy.cs.odu.edu/aggr/timemap/link
        /1/" + url #plug in the url to a timemap
    try:
        response = urllib2.urlopen(mem_url)
        timemap = response.read()
    except urllib2.HTTPError:
        timemap = None
    return timemap

def countMementos(mem_url):
        time_map = getTimeMap(mem_url)
        if not time_map:# if no time maps
                count=0
        else:


                count=len(mymementos.findall(str(time_map)))#
                    finds the count of all mementos per url
                if count>0:
                        file3.write("%s\n"% count)#writes the
                            count of urls onto a json file
```

```
30                              file4.write("%s\n"% time_map)#writes all
                                    the urls on to a json file
31                      #print count
32              return count
33
34  if __name__=="__main__":
35          file1=open('output.json','r')# input a json file that
                    contains 1000 urls
36          file2=open('memcount.json','w')
37          #memcountlist=[]
38          for line in file1.readlines():
39                  one_line = json.loads(line)# loads a json object
40                  link = one_line['link']
41                  counter=countMementos(link)# counter has count
                        of the urls
42                  file2.write("%s"% counter)#outputs count of
                        mementos of each url to a json file
43                  file2.write("\r\n")
44          #for item in memcountlist:
45
46
47  file1.close()
48  file2.close()
```

Listing 5: Python program for generating count of urls and mementos which have more than zero mementos

```
1   from checkForModules import checkForModules
2   import json
3   from ordereddict import OrderedDict
4   #import simplejson
5   import urlparse
6   import re
7
8   from getBitly import getBitlyCreationDate
9   from getArchives import getArchivesCreationDate
10  from getGoogle import getGoogleCreationDate
11  from getBacklinks import *
12  from getLowest import getLowest
13
14  from getLastModified import getLastModifiedDate
15  #Topsy service is no longer available
16  #from getTopsyScrapper import getTopsyCreationDate
17  from htmlMessages import *
18  from pprint import pprint
19
20  from threading import Thread
21  import Queue
22  import datetime
```

```
23
24   import os,sys, traceback
25
26
27
28
29   def cd(url, backlinksFlag = False):
30
31       #print 'Getting Creation dates for: ' + url
32
33
34       #scheme missing?
35       parsedUrl = urlparse.urlparse(url)
36       if( len(parsedUrl.scheme)<1 ):
37           url = 'http://'+url
38
39
40       threads = []
41       outputArray =['','','','','','']
42       now0 = datetime.datetime.now()
43
44
45       lastmodifiedThread = Thread(target=getLastModifiedDate, args
               =(url, outputArray, 0))
46       bitlyThread = Thread(target=getBitlyCreationDate, args=(url,
                outputArray, 1))
47       googleThread = Thread(target=getGoogleCreationDate, args=(
               url, outputArray, 2))
48       archivesThread = Thread(target=getArchivesCreationDate, args
               =(url, outputArray, 3))
49
50       if( backlinksFlag ):
51           backlinkThread = Thread(target=
                   getBacklinksFirstAppearanceDates, args=(url,
                   outputArray, 4))
52
53       #topsyThread = Thread(target=getTopsyCreationDate, args=(url
               , outputArray, 5))
54
55
56       # Add threads to thread list
57       threads.append(lastmodifiedThread)
58       threads.append(bitlyThread)
59       threads.append(googleThread)
60       threads.append(archivesThread)
61
62       if( backlinksFlag ):
63           threads.append(backlinkThread)
64
```

```
65          #threads.append(topsyThread)
66
67
68          # Start new Threads
69          lastmodifiedThread.start()
70          bitlyThread.start()
71          googleThread.start()
72          archivesThread.start()
73
74          if( backlinksFlag ):
75              backlinkThread.start()
76
77          #topsyThread.start()
78
79
80          # Wait for all threads to complete
81          for t in threads:
82              t.join()
83
84          # For threads
85          lastmodified = outputArray[0]
86          bitly = outputArray[1]
87          google = outputArray[2]
88          archives = outputArray[3]
89
90          if( backlinksFlag ):
91              backlink = outputArray[4]
92          else:
93              backlink = ''
94
95
96
97
98          try:
99
100             lowest = getLowest([lastmodified, bitly, google,
                    archives[0][1], backlink]) #for thread
101         except:
102             print sys.exc_type, sys.exc_value, sys.exc_traceback
103
104
105
106
107
108         file2=open('dates.csv','a')
109         print lowest
110         file2.write("%s\n"% lowest)
111
112
```

14

```
113  file1=open('abovezerourls.json','r')
114
115  for line in file1.readlines():
116
117          cd(line)
```

Listing 6: Python program for extracting the created date of all the urls

```
1   from datetime import datetime
2
3   now=datetime.now()
4   file1=open('created_dates.txt', 'r')
5   file2= open('numberdays.txt','w')
6
7   for line in file1.readlines() :
8           print line
9           dateUrl=line.strip().split()
10          date = dateUrl[0]
11          date_object = datetime.strptime(date,"%Y-%m-%dT%H:%M:%S"
                )
12          print date_object
13          days = (now - date_object).total_seconds() / ( 3600.0 *
                24 )
14          number_days = int(days)
15          file2.write("%s\n"% number_days)
16          file2.write("\r\n")
17  # except:
18          # date_object = datetime.strptime(line,"%Y-%m-%dT%H:%M:%
                S")
```

Listing 7: Python program for calculating the ages by subtracting the created date from current date

### 3.1.1   Results

Graph shows an increasing curve .That is, those urls with less mementos have less age and those urls with more mementos have more age,which is the ideal situation.

```
1   d = read.table('numberdays.txt',col.name=c("number_of_days"))
2   b = read.table('abovezerocounts.json',col.name=c("memento_counts
        "))
3    plot(b$memento_counts,d$number_of_days,xlab="memento counts",
        ylab="days")
```

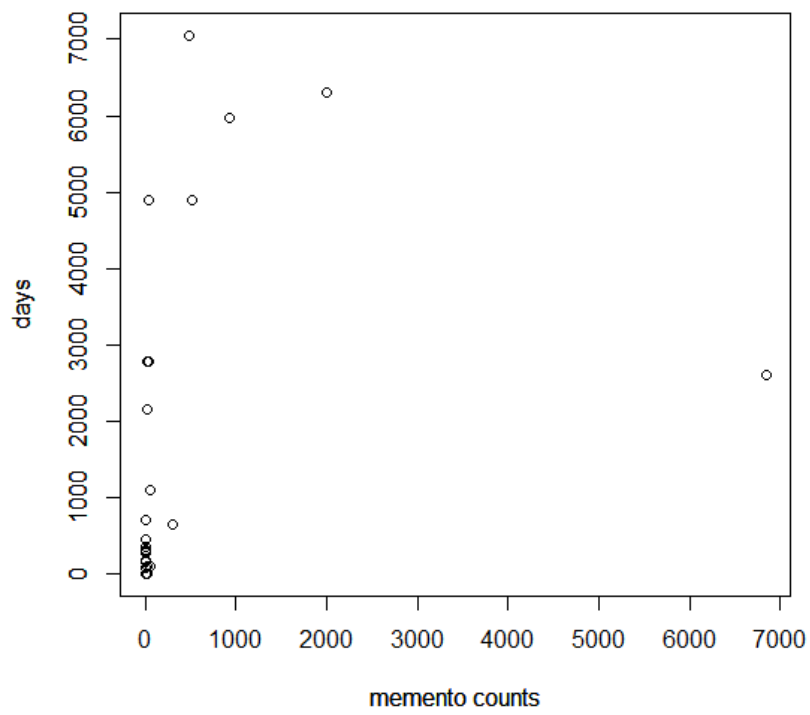Listing 8: R program for generating the scatterplot for Question 3

Figure 3: Number of Mementos vs. Days (Age of url)

# References

[1] Carbon date service. http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html.

[2] Date time function in python. https://docs.python.org/2/library/datetime.htmll.

[3] R histogram. http://www.r-tutor.com/elementary-statistics/quantitative-data/histograml.

[4] R scatterplot. http://www.harding.edu/fmccown/r/.

[5] Tweepy library documentation. http://docs.tweepy.org/en/latest/api.html/.

[6] Writing json data to file:. http://stackoverflow.com/questions/899103/writing-a-list-to-a-file-with-python.