Srividya Majeti

# Assignment 2

CS 532: Introduction to Web Science
Dr. Michael Nelson
Spring 2016

February 11, 2016

# Contents

**1**

# Question 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.**

For solving the above problem I used Python programming language. Following are the steps I have taken to solve the given problem:

- For using the Twitter API, firstly I registered for a twitter application to generate a consumer key and consumer secret.
- Using the keys that are generated in the above step, I authenticated the application for requesting the tweets.
- To fetch the tweet data I started researching for packages and found multiple of them but I decided to work with 'tweepy'.
- Using 'extractTweets.py', I used the tweet data that I received from the API and fetched tweet text, list of URIs in the tweet, tweet JSON and tweet id. This code is listed in Listing3.1.
- While extracting the data mentioned above, program broke multiple times due to overloading, to resolve this issue I wrote an exception handler to wait with a sleep time of 60*15 and then continue.

- The data is processed in JSON format and written to a output file 'tweet.json'.
- I loaded the JSON data from the above outputted file 'tweet.json', and obtained the final URI with a HTTP response code 200 by checking if the URI has any redirects in its history.
- To get unique URIs I stored all the final URIs obtained in the above step in a set data structure which has an inherent property of storing only unique data and have written into a file 'uri.json'.

**Code Listing**

```
1  '''
2  CS532: Introduction to Web Science
3  Author: Srividya Majeti
4  Assignment 2
5  '''
6
7  import tweepy
8  import re
9  import json
10 import urllib2
11 import sys
12 import time
13 import requests
14 from sets import Set
15
16 CONSUMER_KEY = 'wTSsHE3PTA3ZZPiaKHEiQnLtf'
17 CONSUMER_SECRET = '
       UblYYCmNYIEffAY4T4QHGHXwAWMFqiueXdxf35xZFhoK3AECP1'
18 ACCESS_KEY = '157985123-
       WFvzlfDa8KStBZzevMfQBTM7fi8zKHYl2LQpTfGr'
19 ACCESS_SECRET = '
       lSax0XLwIimJ4VVbuU5OY9BpBic4vsSFi0riAq3DPvTxU'
20
21 auth = tweepy.auth.OAuthHandler(CONSUMER_KEY,
       CONSUMER_SECRET)
22 auth.set_access_token(ACCESS_KEY, ACCESS_SECRET)
23 api = tweepy.API(auth)
24
25 tweetJsonFile = open("tweet.json","a")
26 tweetCounter = 0
27 search_results = tweepy.Cursor(api.search, q="tesla", lang="
       en").items(5000)
28
29 while True:
30        try:
31                tweet = search_results.next()
32                for tweet in search_results:
```

```
33                                    #print  tweet._json['entities']['urls
                                          ']
34                                    tweetJson = {}
35                                    tweetJson['id'] = tweet.id
36                                    tweetJson['text'] = tweet.text
37                                    urlList = []
38                                    uriCount = 0
39                                    tweetCounter += 1
40                                    for entityObj in tweet._json['
                                          entities']['urls']:
41                                            uriCount += 1
42                                            urlList.append(entityObj['
                                                url'])
43                                    if uriCount > 0:
44
45                                            tweetJson['uri'] = urlList
46                                            tweetJson['json'] = tweet.
                                                _json
47                                            tweetJsonFile.write(json.
                                                dumps(tweetJson) + "\n")
48                                    if tweetCounter > 1000:
49                                            break
50                except tweepy.TweepError:
51                        print "waiting \n"
52                        time.sleep(60*15)
53                        continue
54                except StopIteration:
55                        break
56
57   tweetJsonFile.close()
58
59   f = open("tweet.json","r")
60   file= open("uri.json","a")
61   count = 0
62   UriSet = Set([])
63   for line in f:
64           data = json.loads(line)
65           if len(data['uri']) > 0:
66                   count += 1
67           link= data['uri'][0]
68           try:
69                   r = requests.get(link)
70                   if r.history:
71                           for h in r.history:
72                                   UriSet.add(r.url)
73                                   # print '[%s] %s' % (h.
                                        status_code, h.url)
74                           # print '[%s] %s' % (r.status_code,
                                  r.url)
```

```
75                        else :
76                                # print '[%s] %s' % (r.status_code,
                                       r.url)
77                                UriSet.add(r.url)
78           except Exception, e:
79                        print e
80                        continue
81
82   for item in UriSet:
83           file.write("%s\n" %item)
84
85   print count
86
87   file.close()
```

**Listing 1.1.** "Python code for extracting tweets and checking for re-directs if it is a 200ok and unique then save it."

# Question 2

**Write a Python program that: Download the TimeMaps for each of the target URIs. We'll use the ODU Memento Aggregator, so for example:**
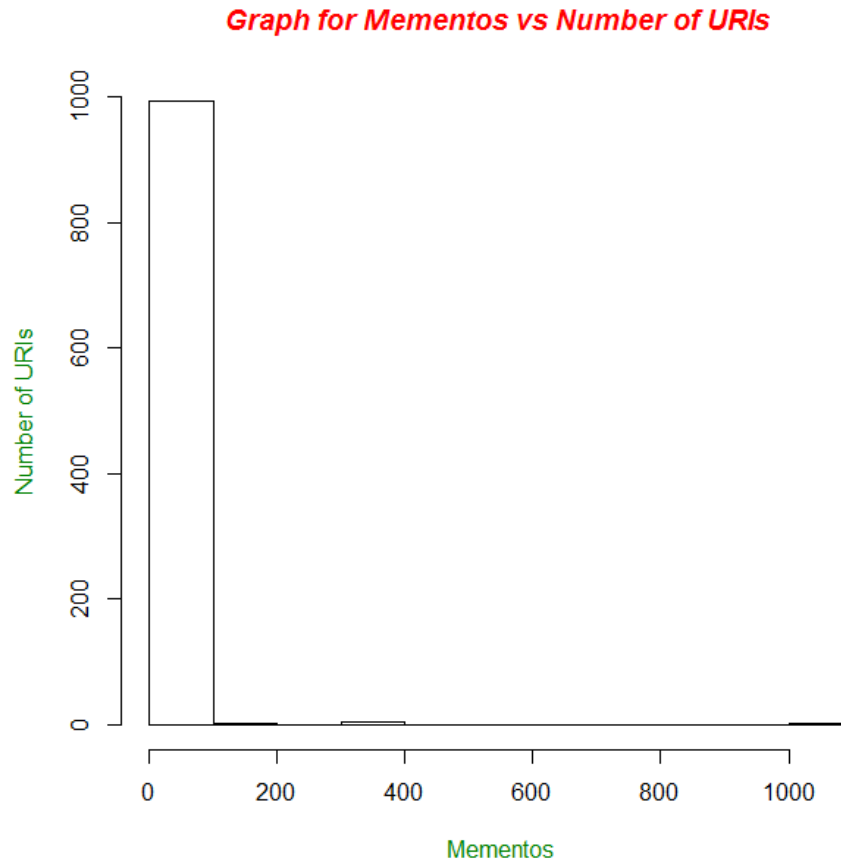**URI-R = `http://www.cs.odu.edu/`**
**URI-T = `http://mementoproxy.cs.odu.edu/aggr/timemap/link/1/http://www.cs.odu.edu/`**
**Create a histogram\* of URIs vs. number of Mementos (as computed from the TimeMaps). For example, 100 URIs with 0 Mementos, 300 URIs with 1 Memento, 400 URIs with 2 Mementos, etc.**
**∗ = https://en.wikipedia.org/wiki/Histogram**

- With the help of ODU Memento Aggregator, I downloaded TimeMaps for all the URIs that are extracted in question 1 using the following curl command:
  `curl-i--silenthttp://mementoproxy.cs.odu.edu/aggr/timemap/link/1/<uri>`
- I stored the output produced by the cURL command into a file 'PagesList'.
- I processed the data from the above file to create a JSON with the original URI, memento URI , datetime and memento count. The memento count for each URI is derived by counting the URIs with rel="memento". This is outlined in Listing 3.1.
- Using 'uri.py' I wrote the memento count for all the 1000 URIs and URIs with > 0 mementos in 2 different files 'allMementos' and 'mementosExcluding0' respectively.This code is listed in Listing 3.2.
- File 'mementosExcluding0' is used for question 3 and file 'allMementos' is used for plotting a histogram with memento count on x-axis and frequency on y-axis.
- I observed that most of the URIs are not archived as they are recently created tweets. Out of 1000 instances 849 tweets have 0 mementos.
- This distribution is summarized in Figure 3.1.

**Graph for Mementos vs Number of URIs**



**Fig. 2.1.** Graph for Mementos vs Number of URIs

## Code Listing

```
 1   '''
 2   CS532: Introduction to Web Science
 3   Author: Srividya Majeti
 4   Assignment 2
 5   '''
 6   import commands
 7   import re
 8   import json
 9   import sys
10
11   def geturi():
```

```
12              file=open("uri.json","r")
13              uricounter = 0
14              for line in file:
15                      if uricounter <1000:
16                              f = open('PagesList','w')
17                              response = getPages(line)
18                              f.write(response)
19                              f.close()
20                              finalCount =getMementosData()
21                              uricounter += 1
22              file.close()
23
24  def getPages(uri):
25          timemapUri = "http://mementoproxy.cs.odu.edu/aggr/
                  timemap/link/1/"
26          command ="curl -i --silent " + timemapUri + str(uri)
                  .strip()
27
28          pageList = commands.getoutput(command)
29          return pageList
30
31  def getMementosData():
32          getMementosDataFile = open("PagesList","r")
33          outputfile= open('mementoData.json', 'a')
34          mementoList= []
35          mementoJson ={}
36          Json ={}
37          count = 0
38          for line in getMementosDataFile:
39                  if 'rel="original"' in line:
40                          count = 0
41                          originalLink =  (re.findall(r'(https
                                  ?://[^\s]+>)', line))[0][:-1]
42                          Json['originaluri'] = originalLink
43                  if 'rel="memento"' in line:
44                          count += 1
45                          link = ""
46                          if re.findall(r'(https?://[^\s]+>)',
                                  line):
47                                  link =  (re.findall(r'(https
                                          ?://[^\s]+>)', line))
                                          [0][:-1]
48                          elif re.findall(r'(www.[^\s]+>)',
                                  line):
49                                  link =  (re.findall(r'(www
                                          .[^\s]+>)', line))
                                          [0][:-1]
50                          else:
51                                  # next(getMementosDataFile)
```

```
52                                         print line
53                              mementoId = count
54                              mementoJson['mementouri'] = link
55                              mementoJson['id'] =mementoId
56                              if(line.find('datetime="') > -1):
57                                      datetime = (line.split('
                                              datetime="'))[1].split('
                                              "')[0]
58                                      mementoJson['datetime'] =
                                              datetime
59                              mementoList.append(mementoJson)
60                              Json['memento'] = mementoList
61              finalCount = str(count)
62              Json['count'] = finalCount
63              outputfile.write(json.dumps(Json) + "\n")
64
65  geturi()
```

**Listing 2.1.** "Python code for getting mementos data and writing the output into a json file"

### Code Listing

```
1   '''
2   CS532: Introduction to Web Science
3   Author: Srividya Majeti
4   Assignment 2
5   '''
6
7   import json
8
9   f= open('data.json','r')
10  f1 = open('urisWithMementosExcluding0','w')
11  f2 = open('allMementos','w')
12  f3=open ('mementosExcluding0','w')
13  for line in f:
14          data= json.loads(line)
15          f2.write(data['count'] + "\n")
16          if data['count'] != "0":
17                  f1.write(line)
18                  f3.write(data['count'] + "\n")
```
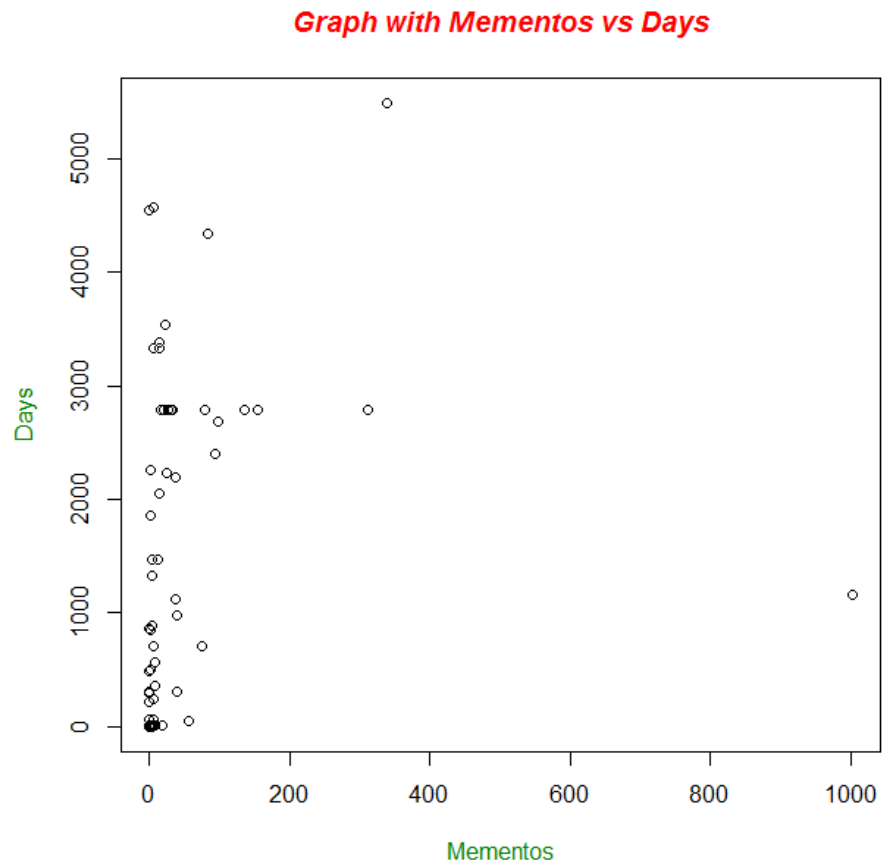
**Listing 2.2.** "Python code for writing all mementos count and URIs with count > 0 in 2 different files "

# 3

# Question 3

**Estimate the age of each of the 1000 URIs using the "Carbon Date" tool:** `http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html`

**Note: you'll should download the library and run it locally; 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. Not all URIs will have Mementos, and not all URIs will have an estimated creation date. State how many fall into either categories.**

- I downloaded the 'Carbon Date' tool from `http://ws-dl.blogspot.com/2014/11/2014-11-14-carbon-dating-web-version-20.html`.
- With the help of this tool, estimated creation date of URIs is obtained.
- I carbon dated URIs which had mementos $> 0$. This is outlined in Listing 3.1.
- Furthermore I parsed the estimated creation date if any and calculated the age of each URI in days using 'getCreatedTime.py'.. This is outlined in Listing 3.2.
- Figure 3.1 illustrates a graph with Mementos on x-axis and Age(in days) on y-axis. Looking at the graph I observe that for most of the URIs, as the number of mementos increases age also increases.
- Figure 3.2 illustrates the data of Age(in Days) and Number of mementos, with Number of mementos in ascending order.
- If I run the code for all the 1000 URIs, 525 URIs did not have estimated Creation date. But if I consider only the URIs that have $> 0$ mementos then all of them have an Estimated creation date.
- Out of the 1000 URIs, 151 URIs fulfilled the given criteria of having $> 0$ Mementos and an estimated creation date.

**Fig. 3.1.** Graph for Mementos vs age

| Item | Age(in Days) | Number of Mementos | Item | Age(in Days) | Number of Mementos | Item | Age(in Days) | Number of Mementos | Item | Age(in Days) | Number of Mementos |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 308 | 1 | 40 | 5 | 3 | 79 | 567 | 9 | 118 | 2791 | 16 |
| 2 | 1 | 1 | 41 | 5 | 3 | 80 | 5 | 9 | 119 | 2791 | 17 |
| 3 | 62 | 1 | 42 | 5 | 3 | 81 | 1469 | 13 | 120 | 2791 | 17 |
| 4 | 8 | 1 | 43 | 1330 | 4 | 82 | 3335 | 14 | 121 | 2791 | 17 |
| 5 | 484 | 1 | 44 | 1469 | 4 | 83 | 3335 | 14 | 122 | 8 | 19 |
| 6 | 210 | 1 | 45 | 1330 | 4 | 84 | 3335 | 14 | 123 | 2791 | 21 |
| 7 | 298 | 1 | 46 | 1330 | 4 | 85 | 3335 | 14 | 124 | 2791 | 21 |
| 8 | 308 | 1 | 47 | 1330 | 4 | 86 | 3386 | 14 | 125 | 3531 | 24 |
| 9 | 7 | 1 | 48 | 1330 | 4 | 87 | 3335 | 14 | 126 | 2227 | 26 |
| 10 | 4 | 1 | 49 | 1330 | 4 | 88 | 2056 | 15 | 127 | 2227 | 26 |
| 11 | 5 | 1 | 50 | 1330 | 4 | 89 | 2791 | 16 | 128 | 2227 | 26 |
| 12 | 4 | 1 | 51 | 1 | 4 | 90 | 2791 | 16 | 129 | 2227 | 26 |
| 13 | 866 | 1 | 52 | 890 | 4 | 91 | 2791 | 16 | 130 | 2227 | 26 |
| 14 | 6 | 1 | 53 | 1330 | 4 | 92 | 2791 | 16 | 131 | 2791 | 28 |
| 15 | 1 | 1 | 54 | 1330 | 4 | 93 | 2791 | 16 | 132 | 2791 | 32 |
| 16 | 4550 | 1 | 55 | 1330 | 4 | 94 | 2791 | 16 | 133 | 2791 | 32 |
| 17 | 62 | 1 | 56 | 1330 | 4 | 95 | 2791 | 16 | 134 | 2791 | 33 |
| 18 | 846 | 2 | 57 | 1330 | 4 | 96 | 2791 | 16 | 135 | 1121 | 38 |
| 19 | 6 | 2 | 58 | 1469 | 4 | 97 | 2791 | 16 | 136 | 2189 | 38 |
| 20 | 4 | 2 | 59 | 1330 | 4 | 98 | 2791 | 16 | 137 | 301 | 39 |
| 21 | 8 | 2 | 60 | 1330 | 4 | 99 | 2791 | 16 | 138 | 977 | 41 |
| 22 | 1858 | 3 | 61 | 1330 | 4 | 100 | 2791 | 16 | 139 | 54 | 57 |
| 23 | 495 | 3 | 62 | 1330 | 4 | 101 | 2791 | 16 | 140 | 701 | 76 |
| 24 | 4 | 3 | 63 | 1330 | 4 | 102 | 2791 | 16 | 141 | 2791 | 80 |
| 25 | 4 | 3 | 64 | 1330 | 4 | 103 | 2791 | 16 | 142 | 4339 | 84 |
| 26 | 8 | 3 | 65 | 1330 | 4 | 104 | 2791 | 16 | 143 | 2403 | 95 |
| 27 | 3 | 3 | 66 | 7 | 5 | 105 | 2791 | 16 | 144 | 2682 | 99 |
| 28 | 5 | 3 | 67 | 4 | 5 | 106 | 2791 | 16 | 145 | 2682 | 99 |
| 29 | 1858 | 3 | 68 | 7 | 5 | 107 | 2791 | 16 | 146 | 2791 | 136 |
| 30 | 5 | 3 | 69 | 4 | 5 | 108 | 2791 | 16 | 147 | 2791 | 154 |
| 31 | 495 | 3 | 70 | 713 | 6 | 109 | 2791 | 16 | 148 | 2791 | 311 |
| 32 | 2254 | 3 | 71 | 62 | 6 | 110 | 2791 | 16 | 149 | 5486 | 338 |
| 33 | 4 | 3 | 72 | 4567 | 6 | 111 | 2791 | 16 | 150 | 5486 | 338 |
| 34 | 5 | 3 | 73 | 245 | 6 | 112 | 2791 | 16 | 151 | 1162 | 1002 |
| 35 | 4 | 3 | 74 | 5 | 6 | 113 | 2791 | 16 | | | |
| 36 | 5 | 3 | 75 | 3331 | 7 | 114 | 2791 | 16 | | | |
| 37 | 8 | 3 | 76 | 9 | 8 | 115 | 2791 | 16 | | | |
| 38 | 5 | 3 | 77 | 358 | 8 | 116 | 2791 | 16 | | | |
| 39 | 5 | 3 | 78 | 9 | 8 | 117 | 2791 | 16 | | | |

**Fig. 3.2.** Data for Number of Mementos and age

**Code Listing**

```
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          #topsy = outputArray[5]
96
97          #note that archives["Earliest"] = archives[0][1]
98          try:
99               #lowest = getLowest([lastmodified, bitly, google,
                       archives[0][1], backlink, topsy]) #for thread
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         result = []
107
108         result.append(("URI", url))
109         result.append(("Estimated Creation Date", lowest))
110         result.append(("Last Modified", lastmodified))
111         result.append(("Bitly.com", bitly))
112         result.append(("Topsy.com", "Topsy is out of service"))
113         result.append(("Backlinks", backlink))
114         result.append(("Google.com", google))
115         result.append(("Archives", archives))
116         values = OrderedDict(result)
117         r = json.dumps(values)
118
119         now1 = datetime.datetime.now() - now0
120
121
122         #print "runtime in seconds: "
123         #print now1.seconds
124         #print r
125         print 'runtime in seconds:   ' + str(now1.seconds) + '\n
                ' + r + '\n'
126
127         return r
128
129
130 output = open('createdTime','w')
131 readData= open('urisWithMementosExcluding0.json','r')
132 uriCounter =0
133 for line in readData:
134         # print line
```

```
135                # print line.rstrip('\n')
136                data= json.loads(line)
137                uri = data['originaluri']
138                uriCounter += 1
139                if uriCounter <1001:
140                        r= cd(uri)
141                        output.write(r + "\n")
142
143  output.close()
144
145  # if len(sys.argv) == 1:
146      # print "Usage: ", sys.argv[0] + " url
                 backlinksOnOffFlag ( e.g: " + sys.argv[0] + " http
                 ://www.cs.odu.edu  [--compute-backlinks] )"
147  # elif len(sys.argv) == 2:
148      # #fix for none-thread safe strptime
149      # #If time.strptime is used before starting the threads,
                 then no exception is raised (the issue may thus
                 come from strptime.py not being imported in a thread
                 safe manner). -- http://bugs.python.org/issue7980
150      # time.strptime("1995-01-01T12:00:00", '%Y-%m-%dT%H:%M:%
                 S')
151      # cd(sys.argv[1])
152  # elif len(sys.argv) == 3:
153      # time.strptime("1995-01-01T12:00:00", '%Y-%m-%dT%H:%M:%
                 S')
154
155      # if(sys.argv[2] == '--compute-backlinks'):
156          # cd(sys.argv[1], True)
157      # else:
158          # cd(sys.argv[1])
```

**Listing 3.1.** "Python code which takes URIs > 0 Mementos as input and writes the JSON output with Estimated Creation date into a file."

**Code Listing**

```
1   '''
2   CS532: Introduction to Web Science
3   Author: Srividya Majeti
4   Assignment 2
5   '''
6
7   import json
8   import datetime
9   import dateutil.parser
10
11  now = datetime.datetime.now()
12  readCreatedTime = open('createdTime','r')
13  ageFile=open('UriWithAge.json','w')
14  noEstimatedDatecounter =0
15  for line in readCreatedTime:
16          data= json.loads(line)
17          Age= {}
18          AgeList =[]
19          if len(data['Estimated Creation Date']) >0:
20                  Age['uri'] = data['URI']
21                  EstimatedDate= data['Estimated Creation Date
                       ']
22                  d1 = dateutil.parser.parse(EstimatedDate)
23                  nowDate =now.isoformat()
24                  d2 = dateutil.parser.parse(nowDate)
25                  days = abs((d2 - d1).days)
26                  Age['days'] = days
27                  ageFile.write(json.dumps(Age) + "\n")
28                  # print Age
29          else:
30                  noEstimatedDatecounter +=1
31  ageFile.close()
32
33  print "URI's with no EstimatedDate",noEstimatedDatecounter
34
35  read= open('UriWithAge.json','r')
36  output= open('days.json','w')
37  for line in read:
38          data=json.loads(line)
39          # print data['days']
40          output.write(str(data['days'])+"\n")
```

**Listing 3.2.** "Python code for calculating age of URI in days and writing them into a file"

# References

1. Bibliography management with bibtex:. pwdhttp://www.sharelatex.com/learn/.
2. code for extracting tweets:. http://stackoverflow.com/questions/22469713/managing-tweepy-api-search.
3. Get consumer key and secret:. https://twittercommunity.com/t/how-do-i-find-my-consumer-key-and-secret/646.
4. How to use bibtex:. http://www.bibtex.org/Using/.
5. Mementos data:. http://www.mementoweb.org/guide/quick-intro/.
6. Retrieve data from twitter:. http://stackoverflow.com/questions/15628535/how-can-i-retrieve-all-tweets-and-attributes-for-a-given-user-using-python code.