# CS 432/532: Web Science

Spring 2017

## Assignment 2

## Michael Micros

## Honor Pledge

I pledge to support the Honor System of Old Dominion University. I will refrain
from any form of academic dishonesty or deception, such as cheating or
plagIiarism. I am aware that as a member of the academic community it is my
responsibility to turn in all suspected violations of the Honor Code. I will report
to a hearing if summoned.

Signed _____

February 16, 2017

# Problem 1: Extracting 1000 unique Links from twitter

The process of extracting links from twitter is straightforward once one understands the structure of the JSON file that tweets have. Most tweets contain links in the "text" field of the JSON, except if they happen to be retweets. That is why before extracting the links we have to check the "Retweeted Status" of the tweet and extract the link from the appropriate JSON field. Once we have collected a reasonable amount of tweets (in the case of this project 5000 tweets were collected), it is important to follow the link until the request response returns a 200. Finally, in order to avoid collecting URIs that point back to twitter we sift thrrough the 5000 collected URIs discarding the twitter links (which understandibly make up the majority of links collected) until we have a list of 1000 unique URIs. For part 1 a library was created in "tweepyTesting.py" that contains all the functions to be used to complete the question. A "main.py" function is used that runs the functions in the appropriate order. The 1000 unique URIs are located in the "links1000" file of part 1.

```python
def getAllTheFinalURI(links):
    uris = []
    cnt = 0
    print("Getting 1000 unique uris ... \n")
    for link in links:
        try:
            uri = getFinalURI(link)
            print uri
            if 'twitter' not in uri:
                print ("Good Link \n")
                if not(uri in uris):
                    uris.append(uri)
                    cnt = cnt +1
                if cnt == 1000:
                    break
            else:
                print ("Twitter Link \n")
        except:
            i = 0
    return uris

def getFinalURI(link):
    try:
        url = urllib.urlopen(link)
        return url.geturl()
    except:
        return ""

def extractLinksFromTweets(tweets):
    links = []
    print("Extracting links...\n")
    for tweet in tweets:
        link = extractLinkFromTweet(tweet)
        links.append(link)
    print("Done extracting!\n")
    return links
```

Figure 1: Functions that extract 1000 unique links

```
1   #! /usr/bin/python
2
3   import tweepy_testing as tt
4
5   dd = tt.collectTheTweets(100)
6   tt.saveTweets(dd,'tweets')
7   tweets = tt.loadTweets('tweets')
8   links = tt.extractLinksFromTweets(tweets)
9   finalURIs = tt.getAllTheFinalURI(links)
10  tt.saveUris(finalURIs, 'test')
11
```

Figure 2: "main.py" calling functions from "tweepyTesting.py"

# Problem 2: Finding number of mementos and plotting histogram

Similarly to problem 1, a library named "part2.py" was created that contained all the functions necessary to download the timemaps of the 1000 URIs and calculate the number of memntos for that link. The main function that utilizes these functions is "main.py". Another action that needed to be taken, was to find and isolate URIs that more than 0 mementos (these URIs are used in Part3 of this assignment).

```
#! /usr/bin/python
import urllib

def memSum(links):
    sum = []
    for link in links:
        i = 0
        prefix = 'http://memgator.cs.odu.edu/timemap/link/'
        res = urllib.urlopen(prefix+link)
        print (res.getcode())
        if res.getcode() == 200:
            for line in res.readlines():
                if 'memento' in line:
                    i = i+1
            sum.append(i)
        else: sum.append(0)
    return sum

def getMementoSum(filename):
    f = open(filename,'r')
    links = []
    while True:
        line = f.readline()
        if not line: break
        links.append(line)
    mementos = memSum(links)
    print (mementos)
    return mementos

def saveSum(mem,filename):
    f=open(filename, 'w')
    for i in mem:
        f.write(str(i)+"\n")
```
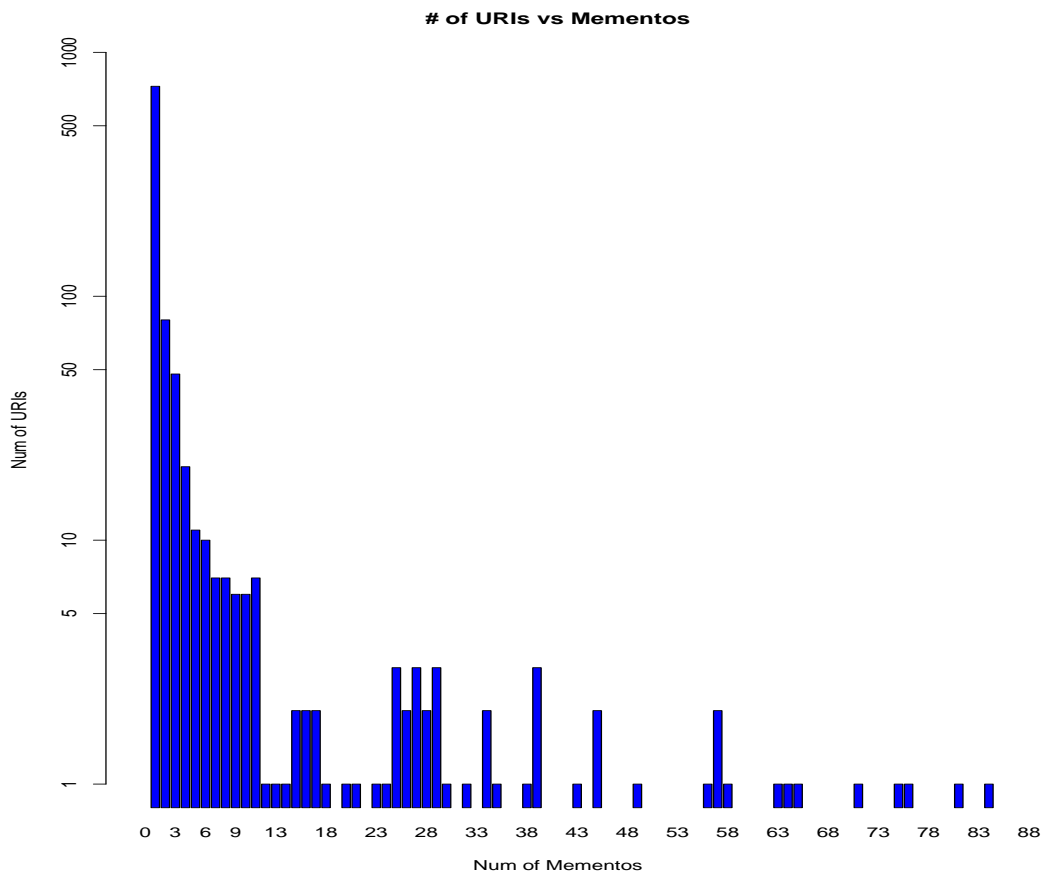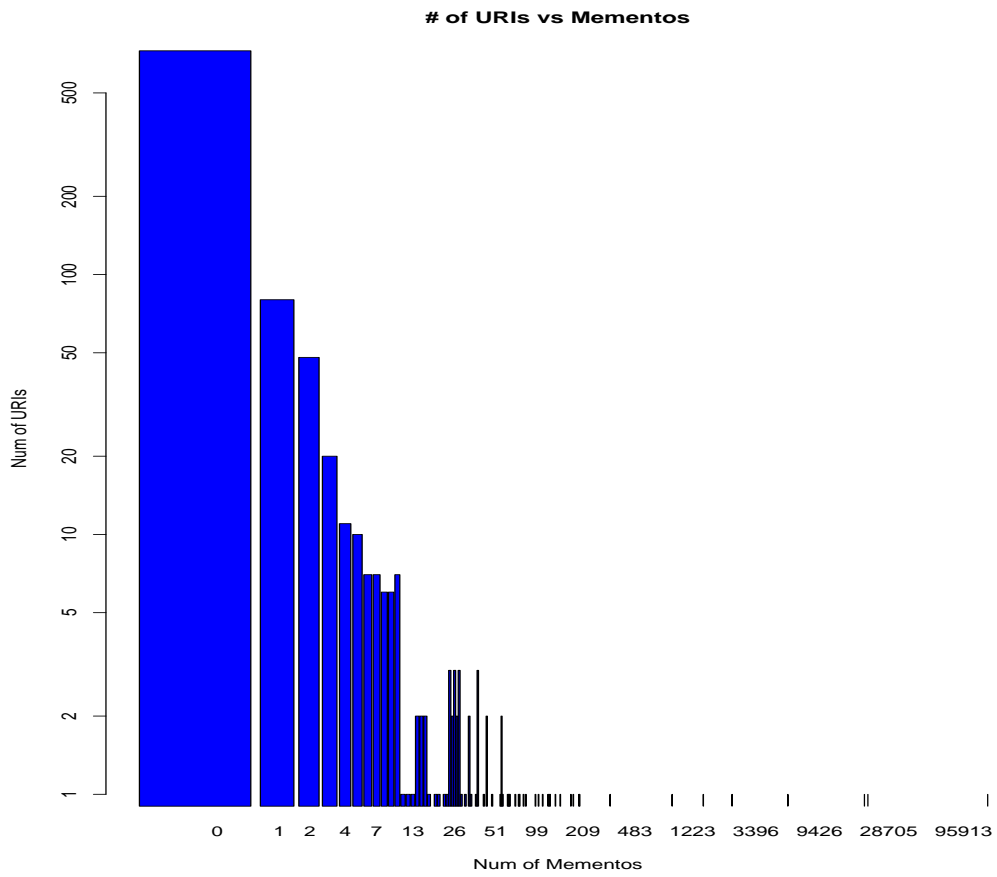
Figure 3: Functions that get the timestamps and calculate the no. of mementos

MemSum.txt

| | |
|---|---|
| 1 | 1 |
| 2 | 0 |
| 3 | 2 |
| 4 | 0 |
| 5 | 2 |
| 6 | 0 |
| 7 | 0 |
| 8 | 0 |
| 9 | 10 |
| 10 | 37 |
| 11 | 0 |
| 12 | 0 |
| 13 | 0 |
| 14 | 0 |
| 15 | 0 |
| 16 | 38 |
| 17 | 16 |
| 18 | 0 |
| 19 | 20530 |
| 20 | 7 |

links1000.txt

| | |
|---|---|
| 1 | http://gizmodo.com/nasa-is-sending-a |
| 2 | http://linkis.com/www.cnn.com/2017/0 |
| 3 | https://www.wsj.com/articles/north-k |
| 4 | http://news.nationalpost.com/life/tr |
| 5 | http://www.afpbb.com/articles/-/3117 |
| 6 | http://www.itfarraq.com/Articles/Det |
| 7 | http://www.boohoo.com/restofworld/ne |
| 8 | http://myemail.constantcontact.com/A |
| 9 | http://www.cnn.com/2017/02/15/politi |
| 10 | http://www.abc.net.au/news/2015-08-1 |
| 11 | http://medium.com/m/global-identity? |
| 12 | http://www.sonhaberler.com/dunya/mac |
| 13 | https://www.facebook.com/molly.macra |
| 14 | http://www.inc.com/lolly-daskal/how- |
| 15 | http://pdwxn.milfalone.com/c/da57dc5 |
| 16 | http://ironmaidenlegacy.com/ |
| 17 | http://www.cnn.com/2017/02/13/politi |
| 18 | http://dailycash.gift/articles/99085 |
| 19 | http://www.cnn.com/politics |
| 20 | http://www.cnn.com/2017/02/15/politi |
| 21 | http://www.nytimes.com/2017/02/15/up |

Figure 4: Sample of file containing no. of mementos

Finally, the mementos are plotted by loading the "memSum.txt" data to RStudio and running a simple script. The following figures are barplots of the same data, with the second one zoomed in to an area of interest.

**# of URIs vs Mementos**

Num of URIs

500
200
100
50
20
10
5
2
1

0  1  2  4  7  13  26  51  99  209  483  1223  3396  9426  28705  95913

Num of Mementos

**# of URIs vs Mementos**

Num of URIs

1000
500
100
50
10
5
1

0  3  6  9  13  18  23  28  33  38  43  48  53  58  63  68  73  78  83  88

Num of Mementos

# Problem 3: Estimating the age of the 1000 URIs

In order to complete this part of the assignment Docker was installed and operated as a server. Having created an image of the CarbonDate repository on the Docker, the server was able to recieve requests and not overload the CS departments CarbonDate server.

```python
def creationDate(res):
    ecd = res['Estimated Creation Date']
    return ecd

def extractCreationDate(carbonDateContent):
    dates=[]
    for c in carbonDateContent:
        jdata = json.loads(c)
        print(jdata['Estimated Creation Date'])
        date = creationDate(jdata)
        dates.append(date)
    return dates

def saveCreationDates(dates, filename):
    f = open(filename,'w')
    for d in dates:
        f.write(d + "\n")


def dayDiff(dates):
    f = open('DaysSinceCreation','w')
    for d in dates:
        day = dt.date( int(d[:4]), int(d[5:7]), int(d[8:10]) )
        today = dt.date.today()
        diff = today-day
        f.write(str(diff.days)+"\n")

def carbonDate(filename):
    cd = []
    f = open(filename, 'r')
    for link in f:
        print("Carbon dating link : " + link)
        #bashCommand = "http://localhost:8888/cd?url=" + link
        #process = subprocess.Popen(bashCommand.split(), stdout=sub
        #output, error = process.communicate()

        res = urllib.urlopen("http://localhost:8888/cd?url=" + link

        out = json.dumps(res.read())
        out = json.loads(out)

        cd.append(out)
    return cd
```

Figure 5: Functions used to obtain the "Estimated Creation Date" of the URIs

After succesfully retrieving the CarbonDate information for each link, the "Estimated Creation Date" was isolated and saved to a text file to be loaded into RStudio along with the corresponding number of memntos for those URIs.

Figure 6: Sample of CarbonDate results

**Age vs # of Mementos**



Figure 7: Sample of CarbonDate results