

OLD DOMINION UNIVERSITY

Assignment 2

David Bayard

February 18, 2019

QUESTION 1.

Write a Python program that extracts 1000 unique (collect more e.g., 1300 just in case) links from Twitter. Omit links from the Twitter domain (twitter.com). Also note that you need to verify that the final target URI (i.e., the one that responds with a 200) is unique.

Solution:

The Twitter API offers users the ability to extract tweets from Twitter, by automating GET requests and responding with Tweet objects. Using this API resulted in JSON formatted Tweets, which could then be scanned for URIs.

Although, target links were only available for the past 7 days using this API, thus popular queries were used including “Sports”, “Trump”, “Republicans”, and “Immigration”. The snippet of code listed below displays the Search method from the Tweepy API.

```
1 for tweet in tweepy.Cursor(api.search, q="Immigration -filter:retweets", count =
    tweetsCountMax, include_entities=True).items(maxNumTweets):
```

Listing 1: Tweepy API Search Method

Once the Tweets have been gathered, the URI's are extracted by accessing the “expanded url” attribute of all the Tweet objects. The collected URI's are filtered for uniqueness by deduplicating and extending the URI's. Issues were encountered in this portion of the assignment, where shortened links were not completely extended. The code below implements a loop which will call “filterLinksFunc(link)”, which will follow redirections and return a URI and status code. A counter was implement to prevent an infinite loop from occurring.

```
1 try:
2     link, statusCode = filterLinksFunc(URL)
3
4     if(statusCode == 200):
5         newList.append(link)
6
7     else:
8         counter = 0;
9         while(statusCode != 200 and counter != 5):
10             filterLinksFunc(link)
11             counter = counter + 1
12
13 except:
14     genericErrorInfo()
```

Listing 2: Extending URI's

After extending all of the URI's, a deduplication function was applied to remove duplicate functions and guarantee uniqueness. This was accomplished by breaking each URI into its separate components and comparing it to other existing URI's. Using this same concept, the domain names of "Twitter.com" were removed.

```
1  for URI in unsanitizedList:
2      scheme, netloc, path, params, query, fragment = urlparse(URI)
3      toCheck = netloc + path
4
5
6      if toCheck not in deduplicatedList:
7          deduplicatedList.append(URI)
8  except:
9      genericErrorInfo()
10
11  return set(deduplicatedList)
```

Listing 3: Deduplicating URI's

QUESTION 2

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

Solution:

In order to download the TimeMaps for each target URI, Docker was installed and the MemGator image was downloaded. This allowed requests to be made via the local machine, rather than making requests through the ODU server. As stated above, the TimeMaps were downloaded by making GET requests via the localhost, which were sent to port 1208 of the server. The specific command used: `sudo docker container run -rm -it -p 5000:1208 ibnesayeed/memgator server`

```
1 ##### Make Request and store the JSON response in json_obj
2
3     requestToMake = "http://localhost:5000/timemap/json/" + URI
4
5     resp = requests.get(requestToMake, timeout=30)
6
7     ##Write response in json format ##
8     json_obj = resp.json()
```

Listing 4: Make TimeMap Request

Once a JSON object is created from the request, it is possible to read the response, and parse the data as shown.

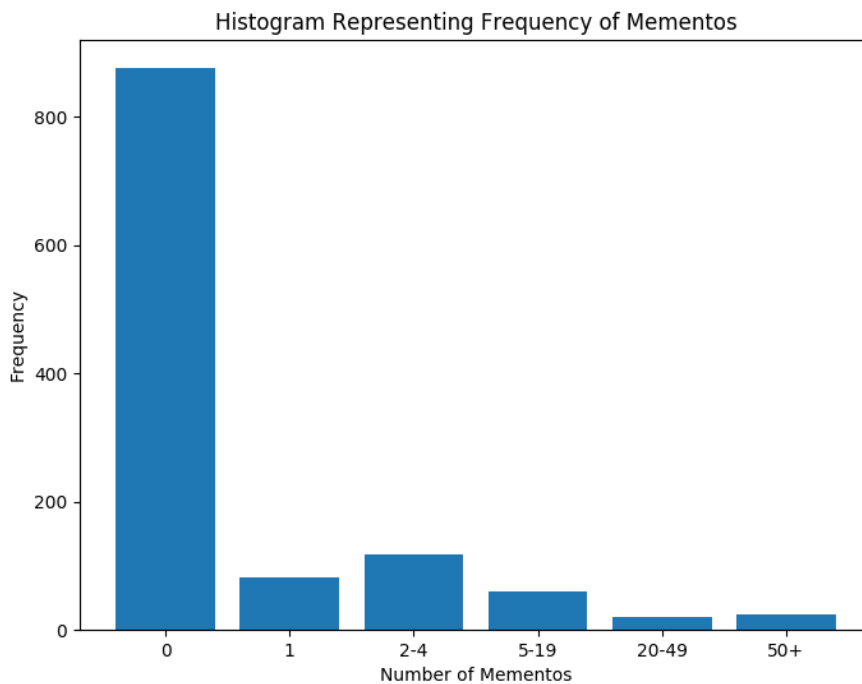
```
1 ##### Make Request and store the JSON response in json_obj
2     dataGood = { "URI" : "value", "Mementos" : "value2", "TimeMapFileName" : "value3" }
3
4     dataGood["URI"] = URI
5     dataGood["Mementos"] = len(json_obj["mementos"]["list"])
6     dataGood["TimeMapFileName"] = fileName
```

Listing 5: JSON Request

The code above displays that the JSON format provides access to the content of the response, which is parsed earlier using the `resp.json()` function. This methodology is used through the entire program, making it simple to store and retrieve data from files, simplifying the program. Moreover, after extracting the JSON data, it is convenient to append the data to a dictionary, and create a list of dictionaries. Then, this list of dictionaries can be written to a .json file, and can easily be accessed later in the program.

```
1 allMementosText.write("URI: " + dataGood["URI"])
2     allMementosText.write("\n")
3     allMementosText.write("TimeMapFileName: " + dataGood["TimeMapFileName"])
4     allMementosText.write("\n")
5     allMementosText.write("Mementos: " + str(dataGood["Mementos"]))
6     allMementosText.write("\n\n\n")
7
8     #allMementos is file being written to
9     json.dump(dataGood, allMementos, indent=2)
10
11     list.append(len(json_obj["mementos"]["list"]))
12     print("Added " + str(len(json_obj["mementos"]["list"])))
```

Listing 6: Writing JSON to file



(a) Histogram representing URI vs Number of Mementos

QUESTION 3

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

The exact same methodology from question 2 can be applied here to get the estimated age of each URI. This is achieved by launching the Docker container for the "Carbon Date tool" using the command "sudo docker container run -rm -it -p 8888:8888 oduwsd1/carbonate -s".

With the Carbonate container running, each request is processed through Carbonate, and every response is stored within the "CarbonRecords.json" file.

```

1
2 for URI in listURI:
3
4     try:
5
6         requestToMake = "http://localhost:8888/cd/" + URI
7         resp = requests.get(requestToMake)
8
9
10
11    except:
12        genericErrorInfo()
13
14    try:
15        json_obj = resp.json()
16
17    except:
18        genericErrorInfo()

```

Listing 7: Make Carbon Date Request

Using the “CarbonRecords.json” and “MementosTrack.json” files, a scatter plot is created representing the number of Mementos, and the estimated creation date of the URI corresponding to the URI.

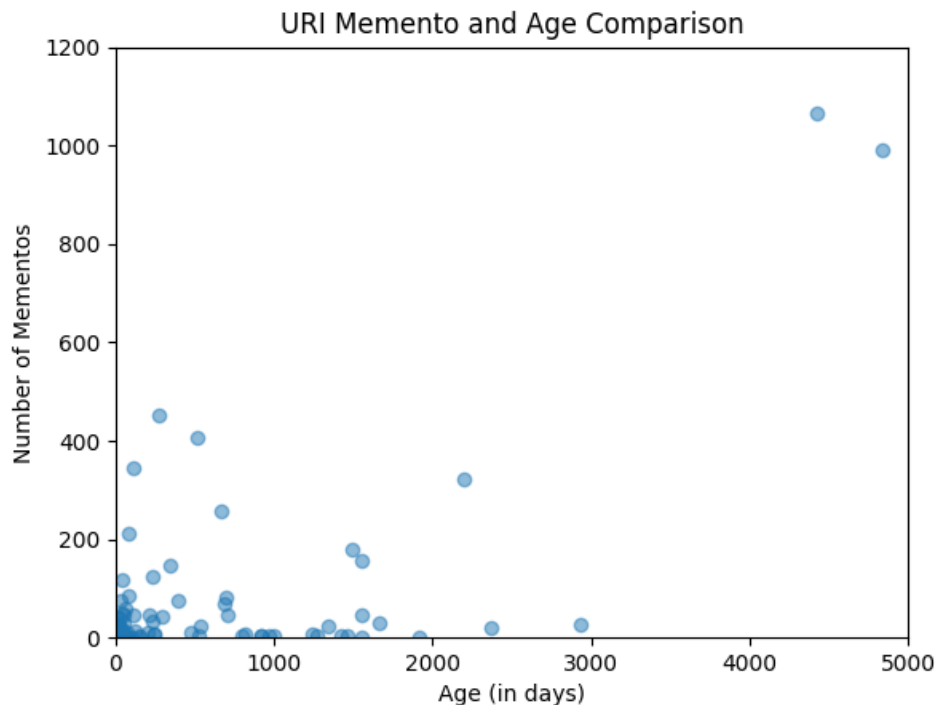
```

1  for link in URIBoth:
2  for a in carbontoData:
3  if (a["URI"] == link and a["Carbon-Date"] != ""):
4
5      toParse = a["Carbon-Date"]
6      carbonDate = dateutil.parser.parse(toParse).date()
7
8      age = abs((today - carbonDate).days)
9
10
11     ageVsMemento["Age"] = age
12
13
14     for b in mementoData:
15         if (b["URI"] == link):
16             ageVsMemento["MementoCount"] = int(b["Mementos"])
17             ageVsMementoList.append(ageVsMemento.copy())

```

Listing 8: Number of Mementos vs Age in days

The main function of the code above is to parse the date from each URI that has a Carbon-Date, as well as the number of Mementos corresponding to the URI. This method is not efficient, but the design of the program prevented a more efficient approach. In the future, a preferred course would be to write all of the data to a single JSON file, making the process of parsing simple. Regardless, the data was gathered and the results were used to create the scatter plot depicted below.



(a) Scatter Plot of Memento Count vs Age of URI

In total, the results for the URIs, creation dates, and number of Mementos is categorized as shown below:

Total URIs: 1179

Number of URIs with Mementos: 303

Number of URIs with Estimated Creation Dates:1017

This information states that out of 1179 URIs, 1017 of them had estimated creation dates, and only 303 of them has at least one Memento.