

CS532S19: Assignment #2

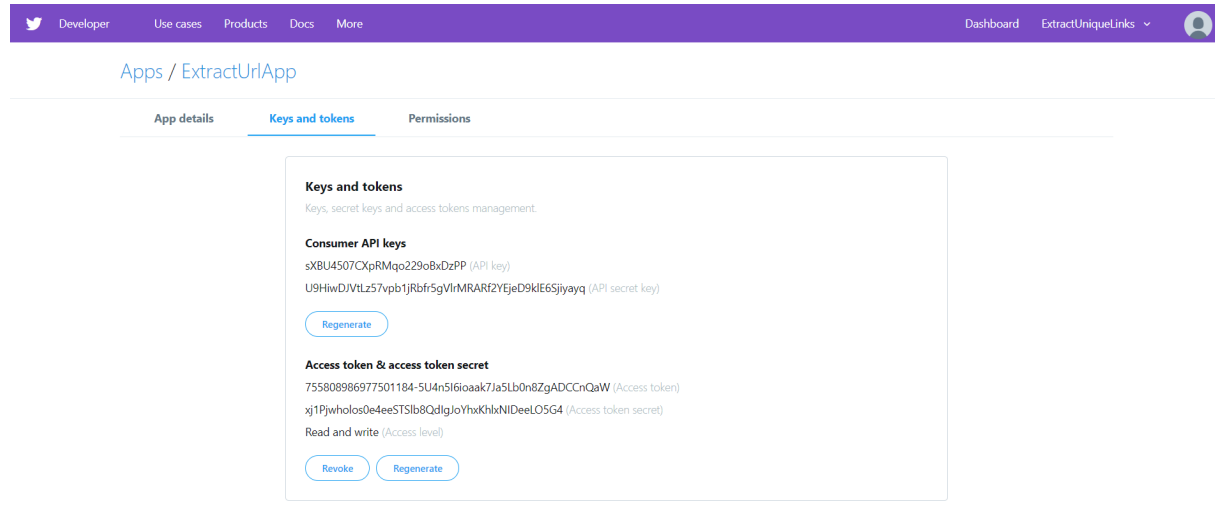
Due on Thursday, February 14, 2019

Nwala, Alexander C

Giridharan Ganeshkumar

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).



1. As shown in the above screen shot, the first step was to sign into developer.twitter.com and create an app. In this case the app created was ExtractUniqueLinks. Also created an Access token and access token secret
2. The next step was to install tweepy. pip install tweepy
3. The first step in the script is to authorize through the tweepy.OAuthHandler and pass the consumer api key and consumer api secret key. Also set the access token with object created in the previous step by passing the access token and access token secret generated in the first step.
4. As show in the below script we have derived class listener from the base class StreamListener
5. The next step is to start the stream by sending the authorization object and the listener derived class. Also filter the stream with some filters sports,music,awards,NASA,cricket,arts,nature,computers,gardening,genetics and medicine.
6. Inside the listener we override the on data method to capture the tweets being streamed, for assignments purpose had a counter to restrict the number tweets being received. The received tweet json is passed to json.load method so that the we can access the entities-urls.
7. loop through the list of urls, exclude urls containing twitter.com and write other urls to the text file.

Listing 1: Python Script

```

1 from tweepy import Stream
2 from tweepy import OAuthHandler
3 from tweepy.streaming import StreamListener
4 import json
5 import re
6
7 class listener(StreamListener):
8
```

```

9     def __init__(self, api=None):
10         super(listener, self).__init__()
11         self.num_tweets = 0
12
13     def on_data(self, data):
14         #print (data)
15
16         if self.num_tweets < 2001:
17             tweet = json.loads(data)
18             for url in tweet["entities"]["urls"]:
19                 if ('twitter.com' not in url["expanded_url"]):
20                     self.num_tweets += 1
21                     fileToStoreListOfUrls.write(url["expanded_url"] + '\n')
22                     print(self.num_tweets)
23             return True
24         else:
25             print('Exception')
26
27             return False
28
29     def on_error(self, status):
30         print(status)
31
32 consumerApiKey = 'sXBU4507CXpRMqo229oBxDzPP'
33 consumerApiKeySecret = 'U9HiwDJVtLz57vpb1jRbfr5gVlrMRARf2YEjeD9kIE6Sjiyayq'
34 accessToken = '755808986977501184-5U4n5l6ioaak7Ja5Lb0n8ZgADCCnQaW'
35 accessTokenSecret = 'xj1Pjwholos0e4eeSTSlb8QdlgJoYhxKhIxNIDeeLO5G4'
36
37 fileToStoreListOfUrls = open("C:\ListOfUrls.txt", "a+")
38 auth = tweepy.OAuthHandler(consumerApiKey, consumerApiKeySecret)
39 auth.set_access_token(accessToken, accessTokenSecret)
40 twitterStream = Stream(auth, listener())
41 twitterStream.filter(track=['sports', 'music', 'awards', 'NASA', 'cricket', '
42     arts', 'nature', 'computers', 'gardening', 'genetics', 'medicine'])
43 fileToStoreListOfUrls.close()
44 print('Completed Streaming.')

```

1. As show in the below script we now expand the url to resolve all the redirects and shortened URIs
2. The expanded URIs are written into a file for further processing.

Listing 2: Python Script

```

1 import urllib.request
2
3 listOfUrlsFile = open("C:\ListOfUrls.txt", "r")
4 fileToStoreListOfUrlsExpanded = open("C:\ListOfUrlsExpanded.txt", "a+")
5 for currentUrl in listOfUrlsFile:
6     try:
7         with urllib.request.urlopen(currentUrl) as url:
8             currentExpandedUrl = url.geturl()
9             fileToStoreListOfUrlsExpanded.write(currentExpandedUrl + '\n')
10    except:
11        pass

```

```

12 listOfUrlsFile.close()
13 fileToStoreListOfUrlsExpanded.close()

```

1. The final step as show in the below script is to remove the duplicates. Got the required help from the professor to get the canonicalizeURI which strips down the url removes the scheme and query strings and other extra parameters which make it now easy for comparison.
2. To get the unique urls we define a set of uniqueUrls and in the for loop add it to the set if it does not exist already in the list. Write it to the file of unique urls ListOfUrls.txt.
3. All the output files and the intermediate files are saved to the OutputFiles folder in git hub.

Listing 3: Python Script

```

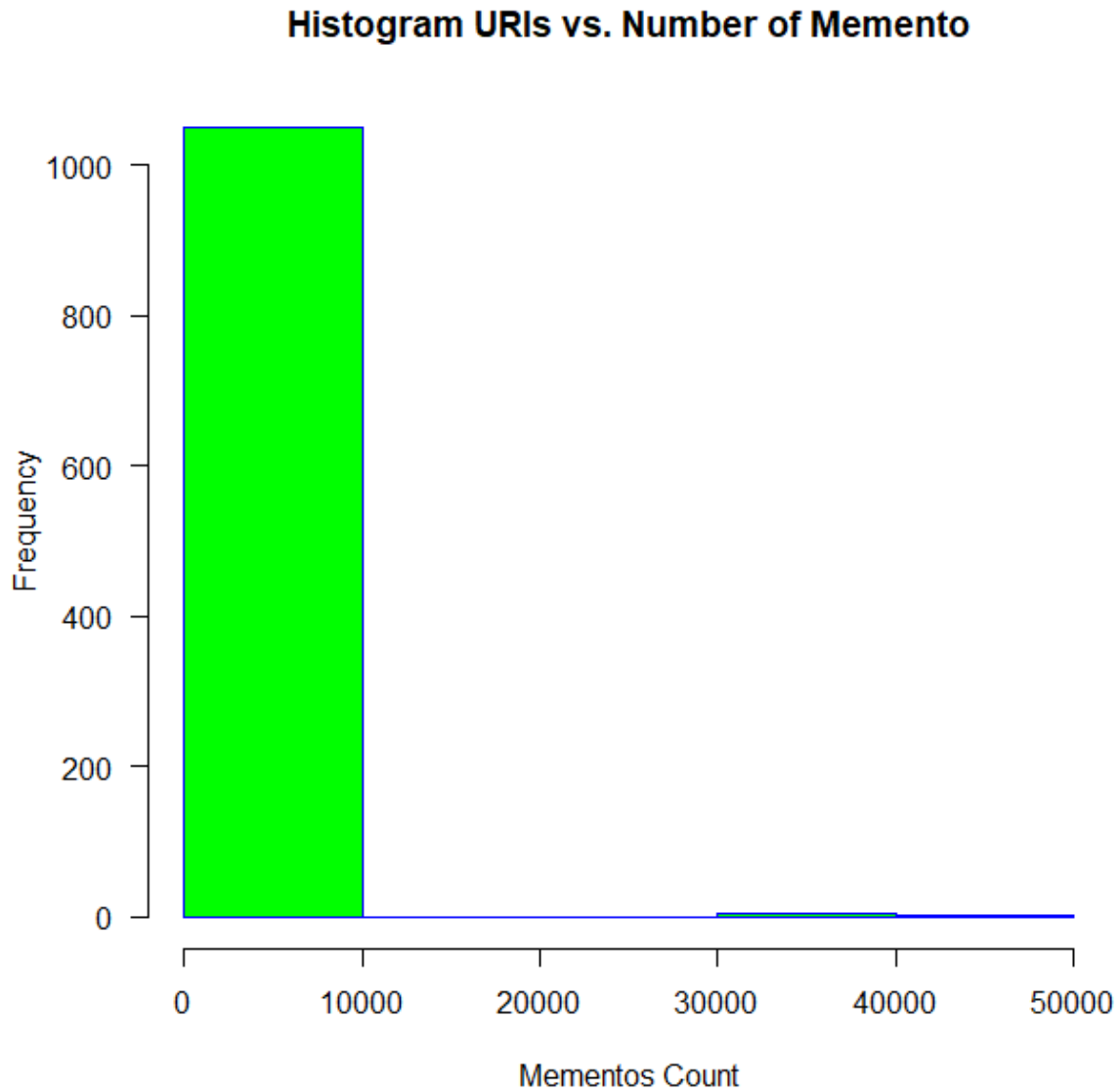
1 from urllib.parse import urlparse
2 import hashlib
3
4 def canonicalizeURI(uri):
5     uri = uri.strip()
6     if (len(uri) == 0):
7         return ''
8     exceptionDomains = ['www.youtube.com']
9     try:
10         scheme, netloc, path, params, query, fragment = urlparse(uri)
11         netloc = netloc.strip()
12         path = path.strip()
13         optionalQuery = ''
14
15         if (len(path) != 0):
16             if (path[-1] != '/'):
17                 path = path + '/'
18
19         if (netloc in exceptionDomains):
20             optionalQuery = query.strip()
21
22         return netloc + path + optionalQuery
23     except:
24         print('Error uri:', uri)
25
26     return ''
27
28 uniqueUrls = set()
29 output_file = open("C:\ListOfUniqueUrls.txt", "w")
30 for currentUrl in open("C:\ListOfUrlsExpandedInput.txt", "r"):
31     currentCanonicalizedURI = canonicalizeURI(currentUrl)
32     if currentCanonicalizedURI not in uniqueUrls:
33         output_file.write(currentUrl)
34         uniqueUrls.add(currentCanonicalizedURI)
35 output_file.close()

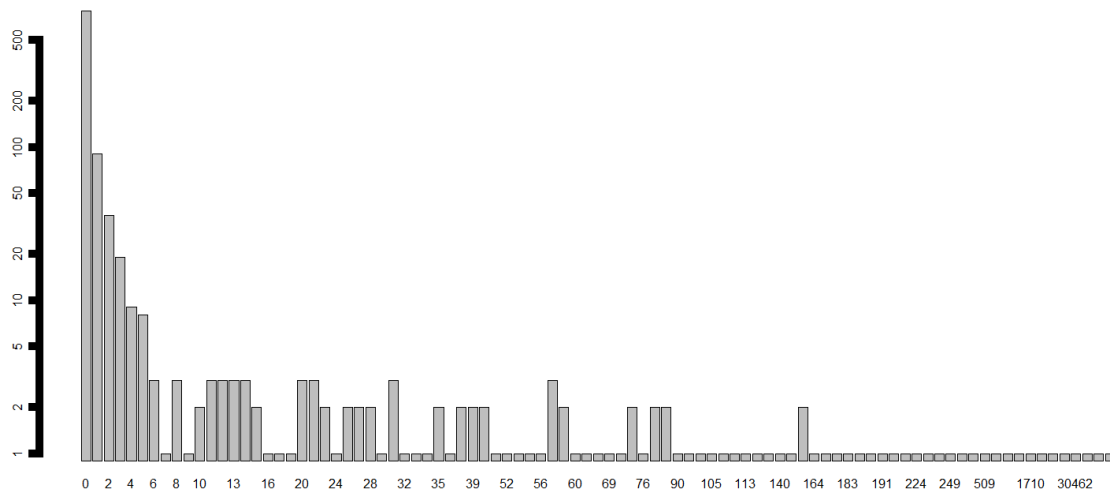
```

Question 2

Download the TimeMaps for each of the target URIs. We'll use the ODU Memento Aggregator. 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. The x-axis will have the number of mementos, and the y-axis will have the frequency of occurrence.





1. As mentioned we will be using the ODU Memento Aggregator for this purpose.
2. A function `getTimeMapAndMomentos` is defined that accepts the `uri` and the `currentUrlLineCounter` as an argument.
3. We have the `memgatorJsonRequestUrlPrefix` which contains the ODU memgator url along with the format `timemap` and `json` attributes included. In this case - <https://memgator.cs.odu.edu/timemap/json/>
4. The `uri` received in the method is appended to the `memgatorJsonRequestUrlPrefix` and a GET request is made.
5. The response is stored in a file with the name as the `currentUrlLineCounter`
6. Also, we return the mementos by reading from the json, that is by getting the length of the mementos list tag
7. The main script as shown below reads the list of unique urls from the previous problem's output and sends the url to the `getTimeMapAndMomentos` method. The returned mementos is written into a csv file. We write into a csv so that we can have a one to one mapping between the url and the memento.

Listing 4: Python Script

```

1 import urllib.request
2 import json
3 import csv
4
5 urlCounter = 0
6 def getTimeMapAndMomentos(uri, currentUrlLineCounter):
7     try:
8         currentfileName = "C:\Timegate\\" + str(currentUrlLineCounter) + ".txt"
9         output_file = open(currentfileName, "w")
10        memgatorJsonRequestUrlPrefix = "https://memgator.cs.odu.edu/timemap/json/"
11        contents = urllib.request.urlopen(memgatorJsonRequestUrlPrefix + uri).read()

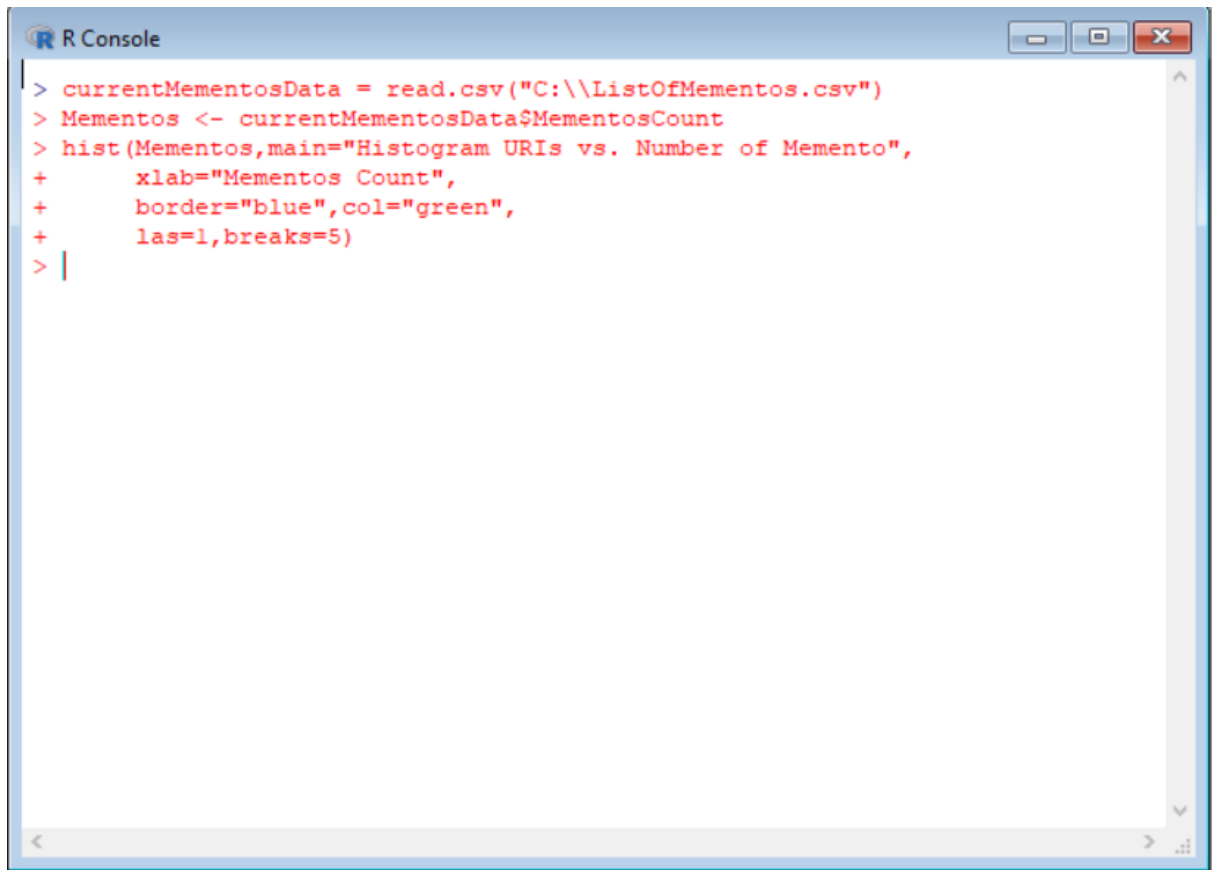
```

```

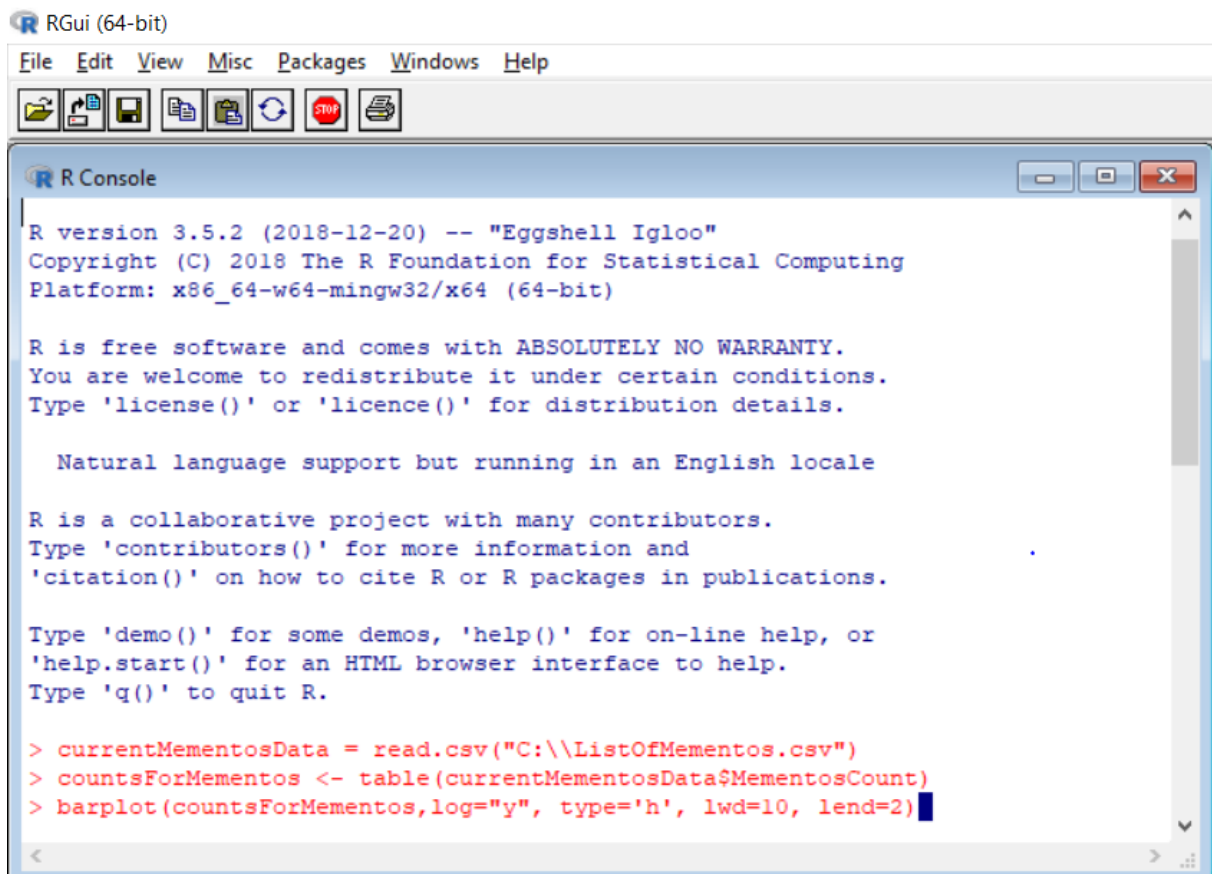
12     #print(contents)
13     jsonContents = json.loads(contents)
14     output_file.write(str(contents))
15     output_file.close()
16     mementosCount = len(jsonContents["mementos"]["list"])
17     return mementosCount
18 except:
19     return 0
20
21 with open("C:\\ListOfMementos.csv", "w", newline='') as
22     csvFileListOfMementosTemp:
23     writer = csv.writer(csvFileListOfMementosTemp)
24     currentData = ['URL', 'MementosCount']
25     writer.writerow(currentData)
26     for currentUrl in open("C:\\ListOfUniqueUrls.txt", "r"):
27         urlCounter = urlCounter + 1
28         currentMementosCount = getTimeMapAndMementos(currentUrl, urlCounter)
29         currentData = [currentUrl, currentMementosCount]
30         writer.writerow(currentData)
31         print(urlCounter)
32 csvFileListOfMementosTemp.close()
33 print("Completed.")

```

1. Once we have the data in csv, we can now use the R package to plot the graph.
2. We draw a histogram using R Package, we take the data and pass it to an vector object. Then pass the vector data to the hist() plot graph method as shown below.
3. Also below the histogram is a bar plot that has the distinct count urls for every mementos count. To achieve this we transpose the data to get the counts of url for a given number of mementos count as mentioned in the second command line.
4. Then plot it using barplot method so that we have the mementos count on the x-axis and the count of urls on the y-axis.



```
R Console
> currentMementosData = read.csv("C:\\\\ListOfMementos.csv")
> Mementos <- currentMementosData$MementosCount
> hist(Mementos,main="Histogram URIs vs. Number of Memento",
+      xlab="Mementos Count",
+      border="blue",col="green",
+      las=1,breaks=5)
> |
```



RGui (64-bit)

File Edit View Misc Packages Windows Help

R Console

```
R version 3.5.2 (2018-12-20) -- "Eggshell Igloo"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> currentMementosData = read.csv("C:\\\\ListOfMementos.csv")
> countsForMementos <- table(currentMementosData$MementosCount)
> barplot(countsForMementos,log="y", type='h', lwd=10, lend=2)
```


Question 3

Estimate the age of each of the 1000 URIs using the "Carbon Date" tool. Not all URIs will have Mementos, and not all URIs will have an estimated creation date. Show how many fall into either categories. For example, total URIs, no mementos, no date estimate.

1. Started the problem by installing docker and using powershell command then installed the carbon date locally with docker image pull oduwsdl/carbondate
2. The next step was to run the carbon date locally so used docker container - docker container run -rm -it -p 8888:8888 oduwsdl/carbondate ./main.py -s
3. Once set up it runs on the localhost as http://localhost:8888/
4. Similar to the second problem solution we have getCarbonDate method which accepts the uri and returns the estimated-creation-date.
5. The getCarbonDate method takes the uri and appends it to the getCarbonDateUriPrefix which will be http://localhost:8888/cd/.
6. Then a GET request is made with the resultant urls. The response is read a json object and the estimated-creation-date is retrieved.
7. The main script as shown below reads the list of unique urls from the first problem's output and sends the url to the getCarbonDate method. The returned estimated-creation-date is written into a csv.

Listing 5: Python Script

```

1 import urllib.request
2 import json
3 import csv
4
5 urlCounter = 0
6 def getCarbonDate(uri):
7     try:
8         getCarbonDateUriPrefix = "http://localhost:8888/cd/"
9         contents = urllib.request.urlopen(getCarbonDateUriPrefix + uri).
10             read()
11         #print(contents)
12         jsonContents = json.loads(contents)
13         responseCarbonDate = jsonContents["estimated-creation-date"]
14         print(responseCarbonDate)
15         return responseCarbonDate
16     except:
17         return "No Estimate Creation Date"
18
19 with open("C:\ListOfCarbonDates.csv", "w", newline='') as
20     csvFileListOfMementosTemp:
21         writer = csv.writer(csvFileListOfMementosTemp)
22         currentData = ['URL', 'Estimate Creation Date']
23         writer.writerow(currentData)
24         for currentUrl in open("C:\ListOfUniqueUrls.txt", "r"):
25             urlCounter = urlCounter + 1
26             currentCarbonDate = getCarbonDate(currentUrl)
27             currentData = [currentUrl, currentCarbonDate]
28             writer.writerow(currentData)
29             print(urlCounter)

```

```
28 csvFileListOfMementosTemp.close()  
29 print("Completed")
```

1. List of Unique URIs : 1054
2. No Mementos : 772
3. No Date Estimate : 208