

HW 3 - Ranking Webpages

Adeniran Adeniyi

Sunday, March 7th 2021 by 11:59pm

Q1

For the following tasks, consider which items could be scripted, either with a shell script or with Python. You may even want to create separate scripts for different tasks. It's up to you to determine the best way to collect the data.

- curl, wget, or lynx are all good candidate programs to use. We want just the raw HTML, not the images, stylesheets, etc.

Download the content of the 1000 unique URIs you gathered in HW2. Plan ahead because this will take time to complete.

You'll need to save the content returned from each URI in a uniquely-named file. The easiest thing is to use the URI itself as the filename, but it's likely that your shell will not like some of the characters that can occur in URIs (e.g., "?", "&"). You might want to hash the URIs to associate them with their respective filename.

For example, <https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21> hashes to 2fc5f9f05c7a69c6d658eb680c7fa6ee:

```
% echo -n "https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21" | md5
2fc5f9f05c7a69c6d658eb680c7fa6ee
```

(md5 might be md5sum on some machines; note the -n in echo – this removes the trailing newline.)

Using this as an example, here are some ways you could download the HTML from that URI using the hash as the filename:

- ```
% curl "https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21" > 2fc5f9f05c7a69c6d658eb680c7fa6ee.html
```
- ```
% wget -O 2fc5f9f05c7a69c6d658eb680c7fa6ee.html https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21
```
- ```
% lynx -source https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21 > 2fc5f9f05c7a69c6d658eb680c7fa6ee.html
```

Now use a tool to remove (most) of the HTML markup for all 1000 HTML documents. *python-boilerpipe* will do a fair job, see

<http://ws-dl.blogspot.com/2017/03/2017-03-20-survey-of-5-boilerplate.html> :

```
from boilerpipe.extract import Extractor
extractor = Extractor(extractor='ArticleExtractor', html=html)
extractor.getText()
```

Keep both files for each URI (i.e., raw HTML and processed). Upload both sets of files to your GitHub repo.

## Answer

```
1 $Textfie = "C:\Users\adeni\CS532\Week5\hw2-archiving-aaden001\Q1\dict.
 txt"
2 $arrayUrlMd5 = @()
3 foreach($line in [System.IO.File]::ReadLines($Textfie)){
4 #Build the output file for each URL
5 $md5 = New-Object -TypeName System.Security.Cryptography.
 MD5CryptoServiceProvider
6 $utf8 = New-Object -TypeName System.Text.UTF8Encoding
7 #hash the each URLs in the dict.txt
8 $hash = [System.BitConverter]::ToString($md5.ComputeHash($utf8.
 GetBytes($line)))
9 $hash = $hash.Replace('-', '') #has md5hash values
10
11 #User for question 1
12 docker container run -it --rm curlimages/curl -L $line > C:\Users\
 adeni\CS532\Week8\hw-ranking-aaden001\Q1\html\$hash.html
13 #user for Question 2
14 #Generate csv file for url hash
15 $hash = "$hash.txt"
16 $arrayUrlMd5 += New-Object pscustomobject -Property @{ 'URL'=$line;
 'hash' = $hash}
17
18 $count++
19
20 }
21 #save object to csv file
22 $arrayUrlMd5 | Export-Csv -Append -Path C:\Users\adeni\CS532\Week8\hw-
 ranking-aaden001\Q2\urlHash.csv -NoTypeInformation
```

**Listing 1:** dataCollections.ps1

```
1 # -*- coding: utf-8 -*-
2 """
3 Created on Mon Mar 1 10:05:54 2021
```

```
4
5 @author: adeni
6 """
7
8 from boilerpy3 import extractors
9 import re
10 import os
11
12 pattern = re.compile(r"Q1/html\/([\w]*)\.html")
13 count = 0
14 directory = "Q1/html/"
15 for filename in os.listdir(directory):
16 if filename.endswith(".html"):
17
18 try:
19 f = open(os.path.join(directory, filename), "r")
20 #where my problem all began
21 #reconfigure file encoding to 16
22 f.reconfigure(encoding="utf-16")
23 text = f.read()
24 f.close()
25 #Find a match and group in regex compile
26 tempName = str(os.path.join(directory, filename))
27 m = pattern.match(tempName)
28 #change file extension by building the string
29 textFile_path = "Q1/processed/" + str(m.group(1)) + ".txt"
30
31 extractor = extractors.ArticleExtractor()
32 content=extractor.get_content(text)
33 #print(content)
34 p = open(textFile_path, "a")
35 p.write(content)
36 p.close()
37 except Exception as r:
38 print(r)
```

**Listing 2:** boiler-plate-removal.py

## Discussion

*I followed these instruction below:*

- Download the content of the 1000 unique URIs you gathered in HW2. Plan ahead because this will take time to complete.

## Answer

- I used power-shell scripting doing this task. I used a docker image of curl to use the curl command on windows power-shell to in line 12 of dataCollection.ps1. The curl option L to follow redirects and the variable \$line has the part to the location of dict.txt which has the list of a 1000 url. This command gets only the html document from the URI-R, it also creates these html files in Q1\html folder.

```
12 docker container run -it --rm curlimages/curl -L $line > C:\
 Users\adeni\CS532\Week8\hw-ranking-aaden001\Q1\html\$hash.html
```

**Listing 3:** snapshot of dataCollections.ps1

- You'll need to save the content returned from each URI in a uniquely-named file. The easiest thing is to use the URI itself as the filename, but it's likely that your shell will not like some of the characters that can occur in URIs (e.g., "?", "&"). You might want to hash the URIs to associate them with their respective filename.  
For example, <https://www.cnn.com/world/live-news/nasa-mars-rover-landing-02-18-21> hashes to 2fc5f9f05c7a69c6d658eb680c7fa6ee:

## Answer

- I converted the url gotten from the dict.txt to an MD5 hash in line 5 to 9

```
5 $md5 = New-Object -TypeName System.Security.Cryptography.
 MD5CryptoServiceProvider
6 $utf8 = New-Object -TypeName System.Text.UTF8Encoding
7 #hash the each URLs in the dict.txt
8 $hash = [System.BitConverter]::ToString($md5.ComputeHash($utf8.
 GetBytes($line)))
9 $hash = $hash.Replace('-', '') #has md5hash values
```

**Listing 4:** snapshot of dataCollections.ps1

- Now use a tool to remove (most) of the HTML markup for all 1000 HTML documents. *python-boilerpipe* will do a fair job, see <http://ws-dl.blogspot.com/2017/03/2017-03-20-survey-of-5-boilerplate.html>

## Answer

- I used boilerpy3 which is a better version of boilerpipe to do the extraction of the html. I grabbed all the html document files in \Q1 \html extracted them with boilerpy3 and deposited the result in \Q1 \processed folder.

One major issues I had encountered was UTF encoding problems. This caused the boilerpy3 not work on the html string I feed it. So i noticed the output of the processed html files to text file automatically came out from power-shell as UTF UTF8BOM when using the file in dataCollections.ps1 In other for the extraction tool to do its job,

I had to reconfigure the text files gotten from processed folder to UTF-16 encoding . In line 22 of boiler-plate-removal.py

```
22 f.reconfigure(encoding="utf-16")
```

**Listing 5:** snap shot of boiler-plate-removal.py

- o I used the getDistinctDomain.py to clean the final.csv file to produce the actualFinal.csv result. Only ten list was printed to the csv file, this had 10 distinct URL domains

## Q2

### Answer

```
1 #Set the directoryPath
2 $directoryPath = "C:\Users\adeni\CS532\Week8\hw-ranking-aaden001\Q1\
 processed\"
3 #Save the file names in an array
4 $files = Get-ChildItem -Name -Path "$directoryPath"
5
6 $pattern = "coronavirus"
7 #to track the number of documents that had the term coronavirus
8 $stopDocument = 0
9
10 #declare an array of Object
11 $testObject = @()
12
13 #final path to Export-Csv
14 $finalPath = "C:\Users\adeni\CS532\Week8\hw-ranking-aaden001\Q2\final.
 csv"
15 #run a for loop
16 foreach ($filenames in $files){
17 #concatenate the full file path
18 $directoryfilenames = "$directoryPath$filenames"
19 #get the number of times, the term coronavirus is found in the
 document
20 $termFrequency= (Get-Content $directoryfilenames| Select-String -
 pattern "$pattern").length
21 #Get the total word per document
22 $t = Get-Content $directoryfilenames| Measure-Object -word | Select
 -Object -expandproperty Words
23
24
25 #avoid dividing by zero
26 #Get TF value frequency/ total word count in document
27 if($t -eq 0){
```

```

28 $tf = 0
29 }else {
30 $tf = ($termFrequency/$t)
31 }
32
33 #Get-Content $directoryfilenames| Measure-Object Word |Select-
Object -expand Words
34 #get only document that has a frequency of at least 1
35 #Answering Q4 when the frequency of the word in document is at
least one then we can count it once
36 if($termFrequency -gt 1){
37 $stopDocument++
38
39 $testObject += New-Object pscustomobject -Property @{'Frequency
' = $termFrequency;'TotalWord' =$t; 'TF'= $tf;'IDF' = "";'Name'=
$filenames; 'URL'="" }
40 }
41 }
42 $filePath = "C:\Users\adeni\CS532\Week8\hw-ranking-aaden001\Q2\urlHash.
csv"
43 #get Generated url hash csv
44 #and make a hashtable
45 $mytable = Import-CSV -Path $filePath
46 $headers = $mytable[0].psobject.properties.name
47 $key = $headers[0]
48 $value = $headers[1]
49 $hashTable = @{}
50 $mytable | %{ $hashTable[$_. "$key"] = $_. "$value" }
51
52
53 #use the look up table to link to the corresponding Url to test object
54 $testObject | ForEach-Object {
55 $_.URL = $hashTable[$_.Name]
56 } | Set Object
57 <# Question 4 #>
58 echo "This is the total document with at least One occurrence of the
term coronavirus"
59 echo $stopDocument
60 #export finalresorting outputs
61 $testObject | Export-Csv -Append -Path C:\Users\adeni\CS532\Week8\hw-
ranking-aaden001\Q2\final.csv -NoTypeInformation

```

**Listing 6:** TF.ps1

```

1 # -*- coding: utf-8 -*-
2 """
3 Created on Sun Mar 7 00:36:05 2021
4

```

```
5 @author: adeni
6 """
7 import pandas as pd
8 from urllib.parse import urlparse # for requesting parse
9
10
11 filePath = "Q2/final.csv"
12
13
14 storeSingleD = []
15 myList = []
16
17
18 df = pd.read_csv(filePath, index_col=False, encoding='utf-8')
19 #convert to list
20 myList = df.values.tolist()
21
22
23
24 """
25 Get through the URL
26 get the domain of url domain = urlparse('http://www.example.test/foo/
 bar').netloc
27 If array doesnt contain the domain
28 store in a new pandas datast(list becasueit much better)
29
30 """
31 #skip same domain and take first ten of the list
32 newList = []
33 count = 1
34 for item in myList:
35 domain = urlparse(item[4]).netloc
36
37 if (domain in storeSingleD):
38 pass
39 else:
40 try:
41 count = count +1
42 storeSingleD.append(domain)
43 newList.append([item[3], item[2], round(item[0], 5), "", "",
item[4]])
44 except Exception as r:
45 print(r)
46
47 if count > 10:
48 break
49 #convert to pandas dataframe
```

```
50 dp= pd.DataFrame(newList, columns=['Frequency', 'Total Word', 'TF', 'IDF',
 'TF-IDF', 'URL'])
51 dp.to_csv("Q2/actualFinal.csv", index=False)
```

**Listing 7:** getDistinctDomain.py

## Discussion

*I followed these steps below:*

- I choose a query term coronavirus as seen on line 6 of TF.ps1
- I Opened each of the generated text document to search for the term. Power-shell as some easy ways to get the number of times the term appears in a textfile, seen in line 20 of TF.ps1

```
20 $termFrequency= (Get-Content $directoryfilenames| Select-String
 -pattern "$pattern").length
```

**Listing 8:** snap shot TF.ps1

- I also got the total word count in the document using Measure-Object -word, line 22

```
22 $t = Get-Content $directoryfilenames| Measure-Object -word |
 Select-Object -expandproperty Words
```

**Listing 9:** snap shot TF.ps1

- I was able to compute the the TF using the formula

$$TF(\text{coronavirus}) = (\text{total occurrence of the word coronavirus in a document}) / (\text{the total word count in the document})$$

- I avoided dividing by Zero since some of the text files had no content. I set the TF to be Zero when this occurs, see in line 25-31

```
25 #avoid dividing by zero
26 #Get TF value frequency/ total word count in document
27 if($t -eq 0){
28 $tf = 0
29 }else {
30 $tf = ($termFrequency/$t)
31 }
```

**Listing 10:** snap shot TF.ps1

- using a hash table I created from a urlHash.csv file - i generated, contains filename and corresponding hashed values - I was able to associate matching text file to there URL. The urlHash.csv file is found in \Q2 folder



- The computation for every document with at least one term frequency was created to a csv file called final.csv. This file is located in \Q2 folder
- IDF was calculated using the size of google's corpus of 55B and the search result obtained by using google browser which was 2.03 billion result.  $IDF = \log_2((\text{Googles total size of corpus})/(\text{search result of term on google browser}))$ . The Calculation gave 4.76. -IDF was gotten by multiplying the result of TF and IDF above.

*I manually configured the table below*

**Table 1:** Hits for the term "coronavirus", ranked by TF-IDF

| TF      | IDF  | TF-IDF | URL                                                                                                                                                                                                                                                                                                                                                                                                                           |
|---------|------|--------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 0.01538 | 4.76 | 0.073  | <a href="https://www.newscientist.com/article/mg24933233-500-first-universal-coronavirus-vaccine-will-start-human-trials-this-year/">https://www.newscientist.com/article/mg24933233-500-first-universal-coronavirus-vaccine-will-start-human-trials-this-year/</a>                                                                                                                                                           |
| 0.00814 | 4.76 | 0.039  | <a href="https://www.cbsnews.com/live-updates/coronavirus-china-outbreak-death-toll-infections-latest-updates-2020-02-17/">https://www.cbsnews.com/live-updates/coronavirus-china-outbreak-death-toll-infections-latest-updates-2020-02-17/</a>                                                                                                                                                                               |
| 0.00592 | 4.76 | 0.028  | <a href="https://www.donaldjtrump.com/media/study-hydroxychloroquine-significantly-lowers-coronavirus-death-rate/">https://www.donaldjtrump.com/media/study-hydroxychloroquine-significantly-lowers-coronavirus-death-rate/</a>                                                                                                                                                                                               |
| 0.00556 | 4.76 | 0.026  | <a href="https://pittsburgh.cbslocal.com/2021/02/24/allegheeny-county-3-uk-variants-covid-19/?utm_source=dlvr.it&amp;utm_medium=twitter">https://pittsburgh.cbslocal.com/2021/02/24/allegheeny-county-3-uk-variants-covid-19/?utm_source=dlvr.it&amp;utm_medium=twitter</a>                                                                                                                                                   |
| 0.00535 | 4.76 | 0.025  | <a href="https://thewest.com.au/news/coronavirus/two-patients-given-wrong-dose-of-coronavirus-vaccine-ng-b881805132z?utm_campaign=share-icons&amp;utm_source=twitter&amp;utm_medium=social&amp;tid=1614210946210">https://thewest.com.au/news/coronavirus/two-patients-given-wrong-dose-of-coronavirus-vaccine-ng-b881805132z?utm_campaign=share-icons&amp;utm_source=twitter&amp;utm_medium=social&amp;tid=1614210946210</a> |
| 0.0049  | 4.76 | 0.023  | <a href="http://veronews.com/2021/02/24/coronavirus-in-irc-feb-24-update/">http://veronews.com/2021/02/24/coronavirus-in-irc-feb-24-update/</a>                                                                                                                                                                                                                                                                               |
| 0.00477 | 4.76 | 0.023  | <a href="https://www.wcvb.com/article/biogen-conference-bostons-first-coronavirus-superspreader-event/35618145">https://www.wcvb.com/article/biogen-conference-bostons-first-coronavirus-superspreader-event/35618145</a>                                                                                                                                                                                                     |
| 0.0042  | 4.76 | 0.02   | <a href="https://www.manorisd.net/site/default.aspx?PageType=3&amp;DomainID=4&amp;ModuleInstanceID=1600&amp;ViewID=6446EE88-D30C-497E-9316-3F8874B3E108&amp;RenderLoc=0&amp;FlexDataID=10464&amp;PageID=1">https://www.manorisd.net/site/default.aspx?PageType=3&amp;DomainID=4&amp;ModuleInstanceID=1600&amp;ViewID=6446EE88-D30C-497E-9316-3F8874B3E108&amp;RenderLoc=0&amp;FlexDataID=10464&amp;PageID=1</a>               |
| 0.0029  | 4.76 | 0.014  | <a href="https://www.cnbc.com/2021/02/24/house-democrats-aim-to-pass-1point9-trillion-covid-relief-bill-on-friday.html?utm_source=dlvr.it&amp;utm_medium=twitter">https://www.cnbc.com/2021/02/24/house-democrats-aim-to-pass-1point9-trillion-covid-relief-bill-on-friday.html?utm_source=dlvr.it&amp;utm_medium=twitter</a>                                                                                                 |
| 0.00158 | 4.76 | 0.008  | <a href="https://www.politicshome.com/thehouse/article/coronavirus-public-health-local-responsibility">https://www.politicshome.com/thehouse/article/coronavirus-public-health-local-responsibility</a>                                                                                                                                                                                                                       |

## Q3

**Answer**

**Table 2:** Page ranking per domain gotten from Question 2

| Pageranking | URL                                                                           |
|-------------|-------------------------------------------------------------------------------|
| 0.7         | <a href="https://www.cbsnews.com">https://www.cbsnews.com</a>                 |
| 0.7         | <a href="https://www.cnn.com">https://www.cnn.com</a>                         |
| 0.6         | <a href="https://www.newscientist.com">https://www.newscientist.com</a>       |
| 0.6         | <a href="https://pittsburgh.cbslocal.com">https://pittsburgh.cbslocal.com</a> |
| 0.6         | <a href="https://www.donaldjtrump.com">https://www.donaldjtrump.com</a>       |
| 0.5         | <a href="https://www.politicshome.com">https://www.politicshome.com</a>       |
| 0.5         | <a href="http://veronews.com">http://veronews.com</a>                         |
| 0.5         | <a href="https://www.manorisd.net/">https://www.manorisd.net/</a>             |
| 0.5         | <a href="https://www.wcvb.com">https://www.wcvb.com</a>                       |
| 0.5         | <a href="https://thewest.com.au">https://thewest.com.au</a>                   |

## Discussion

- I used this <https://dnschecker.org/pagerank.php> to obtain the page ranking for each domain
- Comparing the result of the Table 2 to Table 1, depended of a number of factors
  - I noticed that using the various Page ranking website you gave us generated completely different result to one another.

dnschecker.org in particular doesn't allow you to view domains that do not use secured http. I notice that some of the Page Ranking site allowed only secured http, for example [https://www.prchecker.info/check\\_page\\_rank.php](https://www.prchecker.info/check_page_rank.php)

- How well did the boiler plate removal tool extract the html content to text files.
- On Average some of the domain that ranked higher also had a high TD-IDF. In conclusion it various result would differ on if the tools are technology used gave 100 % efficiency. We can only derive few correlations to how authentic the webpage is to the domain of the webpage.
- Personally how authentic a webpage or it materials are is not solely dependent of how google or bing ranks the webpage alone. Over a large area I believe it very hard to quantify the relative importance/ranking of a webpage.

## References

- <https://stackoverflow.com/questions/10521061/how-to-get-an-md5-checksum-in-powershell>
- <https://stackoverflow.com/questions/29889495/count-specific-string-in-text-file-using-powershell>
- <https://adamtheautomator.com/export-csv/>

- <https://docs.microsoft.com/en-us/powershell/scripting/learn/deep-dives/everything-about-if?view=powershell-7.1>
- <https://ss64.com/ps/measure-object.html>
- <https://kimconnect.com/powershell-convert-csv-into-hashtable/>
- <https://devblogs.microsoft.com/scripting/use-a-powershell-cmdlet-to-count-files-words-and-lines/>