# HW3 - Ranking Web Pages

Brenden Lewis
Due: 10/18/2020 11:59PM

# Q1

Unfortunately, I had to acquire a brand new set of URIs to use for this assignment due to my dataset from HW#2 being comprised almost entirely of YouTube links; these kinds of links provide very little in terms of HTML content. For the sake of time, I only collected about 750 URIs through the Twitter API. I used a variety of keywords related to *tech* and *current events* in an attempt to broaden the contents of each page and provide a larger dictionary of words to utilize for this assignment. Unlike in my HW#2 submission, I was able to eliminate duplicates in my set and was left with 740 URIs before processing; these URIs were passed in as *CleanedURLs.txt*.

There were a few issues with some URIs throwing unique errors that I had to exempt due to not knowing how to work around. One of example of this was a link to a working Amazon page that was throwing a weird 503 error. After a bit of research, in this case, it was due to Amazon having their own API (similar to TWitter) that's required to scrape their data.

I downlaoded the HTML for each URI using the *requests* library and processed the HTML using *boilerpipe.extract* to extract the primary text. The raw HTML and processed text per URI was automatically stored in their own *.txt* files and saved in a local directory using a counter to label file names. This apporach resulted in a *rawHTMLFile#.txt* and *processedHTMLFile#.txt* file combo per URI stored in seperate folders locally.

```
1  import requests
2  import urllib.request
3  import http.client
4  import os.path
5  from boilerpipe.extract import Extractor
6
7  count = 1
8
9  def download_HTML():
10     urlFile=open('CleanedURLs.txt','r')
11     #urlFile=open('testRankingURLs.txt','r')
12     urlFileOPen= urlFile.readlines()
13
14     global count
15
16     for url in urlFileOPen:
17         print("URI {}: Processing".format(count))
18         try:
19             response = requests.get(url)
```

```python
20            store_HTML(response.text,count)
21            store_proc_HTML(url,count)
22        except (requests.exceptions.RequestException):
23            print("\tError: URI invalid")
24
25        print("URI {}: Finished \n".format(count))
26        count+=1
27
28    print("\nDone\nTotal URIs processed: "+str(count-1))
29
30
31
32 def store_HTML(response,counter):
33    rawHtmlOut = open(r"C:\Users\Brenden\Desktop\ODU\Fall 2020\CS432\
   HW_3\hw3-ranking-blewis1014\Raw_HTML\rawHtmlFile{}.txt".format(
   counter),"a", encoding="utf-8")
34    rawHtmlOut.write(response)
35    rawHtmlOut.close()
36
37 def store_proc_HTML(source,counter):
38
39    while True:
40        try:
41            extractor = Extractor(extractor='ArticleExtractor', url=
   source)
42            extracted_text = extractor.getText()
43        except (urllib.error.HTTPError, (http.client.IncompleteRead),
   NameError, LookupError, UnicodeDecodeError):
44            print("\tURI invalid")
45            return False
46
47
48
49        procHtmlOut = open(r"C:\Users\Brenden\Desktop\ODU\Fall 2020\
   CS432\HW_3\hw3-ranking-blewis1014\Processed_HTML\processedHtmlFile
   {}.txt".format(counter),"a",encoding="utf-8")
50        procHtmlOut.write(str(extracted_text))
51        procHtmlOut.close()
52        return True
53
54
55
56
57 if __name__ == '__main__':
58    download_HTML()
```

**Listing 1:** Python code for downloading and processing html

# Q2

I opted to make a Python script specifically for sorting through my data set and collecting URIs that contained a specific query (as shown in Listing 2 below). Due to a large portion of my URIs being related to games and game companies, I chose the query "Xbox" for collecting URIs to use in ranking later. In order to read through all the processed files I had stored, I used both the *glob* and *os* libraries in conjunction with the path to folder. For each file, text was read line-by-line, word-by-word to check for any matches to the query. To help with ranking the URIs later on, I was keeping count of each time the query was found in a file and the total amount of words the file possessed. Each file was stored as a dictionary with this information including the file name, and written to *QueryResults.txt* if the query ever appeared in the file.

```python
1  from progressbar import ProgressBar
2  import glob
3  import os
4
5  perc = ProgressBar()
6
7  path = r"C:\Users\Brenden\Desktop\ODU\Fall 2020\CS432\HW_3\hw3-ranking-
       blewis1014\Processed_HTML\*.txt"
8  query = "Xbox"
9  count = 1
10
11 containsQuery = []
12 all = []
13
14 for file in perc(glob.glob(path)):
15     # infile = open(file,"r")
16     occur = 0
17     totalWords = 0
18     with open(file, "r", encoding='ISO-8859-1') as f:
19         # text = f.read()
20         # occur = text.count(query)
21         file_name = os.path.basename(file)
22
23         for lines in f:
24             words = lines.split()
25             for i in words:
26                 totalWords+=1
27                 if i == query:
28                     occur+=1
29
30         file_dict = {"file":file_name,"words:":totalWords,"occurences":
       occur}
31         count+=1
32         if (occur > 0) :
33             containsQuery.append(file_dict)
```

```
34          all.append(file_dict)
35
36
37 outfile = open("QueryResults.txt","w")
38 outfile.write("Query: "+query+"\n")
39 for item in containsQuery:
40     outfile.write(str(item)+"\n")
41
42 outfile.write("\nTotal # of documents: "+str(len(all)))
43 outfile.write("\nTotal # of documents that contain query: "+str(len(
     containsQuery)))
44 outfile.close()
45
46 results = open("QueryResultsForProcessing.txt","w")
47 for item in containsQuery:
48     results.write(str(item)+"\n")
49 results.close()
50
51 print("Done")
```

**Listing 2:** Python code for collecting URIs that contain the query, query frequency, and total word count per document

With the query set to 'Xbox', 24 files were found to contain the query out of the 740. The output from this operation is shown Listing 3. Listing 3 shows 701 total documents because after running it the first time, I discovered that a good chunk of rawHTML and processedHTML files collected from Q1 were either empty or nearly empty (either 0 or 1 kb). I decided to delete those entries and run it again for the sake of consistency and not having missing data (GitHub won't allow uploading of empty files).

```
1 Query: Xbox
2 {'file': 'processedHtmlFile1.txt', 'words:': 412, 'occurences': 8}
3 {'file': 'processedHtmlFile103.txt', 'words:': 256, 'occurences': 8}
4 {'file': 'processedHtmlFile106.txt', 'words:': 407, 'occurences': 4}
5 {'file': 'processedHtmlFile11.txt', 'words:': 581, 'occurences': 11}
6 {'file': 'processedHtmlFile110.txt', 'words:': 407, 'occurences': 4}
7 {'file': 'processedHtmlFile119.txt', 'words:': 587, 'occurences': 11}
8 {'file': 'processedHtmlFile12.txt', 'words:': 481, 'occurences': 11}
9 {'file': 'processedHtmlFile120.txt', 'words:': 407, 'occurences': 4}
10 {'file': 'processedHtmlFile122.txt', 'words:': 252, 'occurences': 4}
11 {'file': 'processedHtmlFile123.txt', 'words:': 450, 'occurences': 4}
12 {'file': 'processedHtmlFile13.txt', 'words:': 225, 'occurences': 4}
13 {'file': 'processedHtmlFile15.txt', 'words:': 138, 'occurences': 2}
14 {'file': 'processedHtmlFile194.txt', 'words:': 1247, 'occurences': 1}
15 {'file': 'processedHtmlFile2.txt', 'words:': 48, 'occurences': 1}
16 {'file': 'processedHtmlFile235.txt', 'words:': 423, 'occurences': 1}
17 {'file': 'processedHtmlFile3.txt', 'words:': 213, 'occurences': 1}
18 {'file': 'processedHtmlFile4.txt', 'words:': 496, 'occurences': 6}
```

```
19 {'file': 'processedHtmlFile5.txt', 'words:': 207, 'occurences': 2}
20 {'file': 'processedHtmlFile6.txt', 'words:': 55, 'occurences': 1}
21 {'file': 'processedHtmlFile740.txt', 'words:': 323, 'occurences': 4}
22 {'file': 'processedHtmlFile8.txt', 'words:': 276, 'occurences': 7}
23 {'file': 'processedHtmlFile96.txt', 'words:': 326, 'occurences': 5}
24 {'file': 'processedHtmlFile97.txt', 'words:': 97, 'occurences': 2}
25 {'file': 'processedHtmlFile99.txt', 'words:': 212, 'occurences': 5}
26
27 Total # of documents: 701
28 Total # of documents that contain query: 24
```

**Listing 3:** QueryResults.txt

I searched through each file to determine which would be best to use for calculation; this was based on a combination of factors such as content, size, language (several were in Portugese, only opted to include one), and which ones I was curious to see the results of. The final 10 I chose to analyze are listed below in Listing 4. I had to manually add in the URLs for each of the 10 to make things a little easier when putting everything together in the end.

```
1 {'file': 'processedHtmlFile1.txt', 'uri': 'https://www.resetera.com/
     threads/xbox-brasil-dismisses-new-host-after-she-was-the-victim-of-
     death-and-rape-threats.307873/', 'words': 412, 'occurences': 8}
2 {'file': 'processedHtmlFile103.txt', 'uri': 'https://www.theverge.com
     /2020/10/15/21517293/xbox-series-x-s-dashboard-ui-october-update-
     release?_lrsc=f9510bb1-a7e2-43c8-9e49-761f32bbdc52', 'words': 256, '
     occurences': 8}
3 {'file': 'processedHtmlFile2.txt', 'uri': 'https://www.exophase.com/
     user/toniazzo/?1602892894', 'words': 48, 'occurences': 1}
4 {'file': 'processedHtmlFile194.txt', 'uri': 'https://newszoneweb.com/?p
     =5809','words': 1247, 'occurences': 1}
5 {'file': 'processedHtmlFile106.txt', 'uri': 'https://www.gamespot.com/
     articles/future-bethesda-titles-dont-need-to-be-on-ps5-says-xboxs-
     phil-spencer/1100-6483415/', 'words': 407, 'occurences': 4}
6 {'file': 'processedHtmlFile740.txt', 'uri': 'https://www.gamespot.com/
     articles/marvels-avengers-ps5-and-xbox-series-x-s-versions-delayed-
     to-2021/1100-6483398/', 'words': 323, 'occurences': 4}
7 {'file': 'processedHtmlFile12.txt', 'uri': 'https://tictechbuzz.com/
     remembering-the-classics-the-xbox-podcast-featuring-tim-schafer/', '
     words': 481, 'occurences': 11}
8 {'file': 'processedHtmlFile8.txt', 'uri': 'https://g1.globo.com/pop-
     arte/games/noticia/2020/10/16/microsoft-demite-apresentadora-do-xbox
     -br-que-sofreu-ameacas.ghtml?utm_source=twitter&utm_medium=social&
     utm_campaign=g1', 'words': 276, 'occurences': 7}
9 {'file': 'processedHtmlFile13.txt', 'uri': 'https://www.engadget.com/
     lucasarts-double-fine-remasters-game-pass-235238100.html', 'words':
     225, 'occurences': 4}
10 {'file': 'processedHtmlFile119.txt', 'uri': 'https://www.ign.com/
     articles/phil-spencer-xbox-bethesda-ps5?utm_source=dlvr.it&
```

```
    utm_medium=twitter', 'words': 587, 'occurences': 11}
```

**Listing 4:** Final set of 10 URIs chsoen for TD-IDF calculations

The Python code for calculating TF-IDF values (Listing 5) was relatively straight forward. I had to utilize the *ast* library to pull in the information from *FinalURLsForRanking.txt* as dictionaries to be used. I was using Google as my corpus, so that variable for calculation was 55b. The TF, IDF, and TF-IDF values were each calculated in their own individual function to 3 significant figures to reduce clutter. The *form operator import itemgetter* was needed to sort my final data in descending order by TF-IDF value; this wasn't completely necessary with how I created my table, but it definitely helped speed-up the process.

```python
1 import ast
2 import math
3 import pandas as pd
4 from operator import itemgetter
5
6 collection = []
7 results = []
8 inputFile = "FinalURLsForRanking.txt"
9 query = "Xbox"
10 CORPUS = 55000000000
11
12 def calcTF(doc):
13     freq = doc["occurences"]
14     total_words = doc["words"]
15     tf = float(freq / total_words)
16     return round(tf,3)
17
18 def calcIDF(allDocuments, corpus):
19     numDocs = len(allDocuments)
20     idf = math.log2(corpus/numDocs)
21     return round(idf,3)
22
23 def calcTFIDF(tf,idf):
24     return round((tf*idf),3)
25
26 def drawResultsTable():
27     pass
28
29 def printResults():
30     for item in results:
31         print(str(item)+'\n')
32
33 def writeResultsToFile():
34     with open("TF-IDF_Results.txt","w") as f:
35         for item in results:
36             f.write(str(item)+'\n')
```

```
37
38 def sortTDFDescending():
39     sort_list = sorted(results, key=itemgetter('TF-IDF'), reverse=True)
40     return sort_list
41
42 if __name__ == '__main__':
43     with open(inputFile, 'r') as f:
44         for line in f:
45             content = line.strip()
46             collection.append(ast.literal_eval(content))
47
48     idf = calcIDF(collection,CORPUS)
49
50     for elem in collection:
51         tf = calcTF(elem)
52         tfidf = calcTFIDF(tf,idf)
53         uri_ranking = {"uri":elem["uri"], "TF":tf, "IDF":idf, "TF-IDF":
    tfidf}
54         results.append(uri_ranking)
55     results = sortTDFDescending()
56     printResults()
57     writeResultsToFile()
```

**Listing 5:** Python code for calculating the TF-IDF for each URI

Figure 1 below shows the final ranking for each URI and their calculated TF, IDF, and TD-IDF values. For both the TD-IDF table and PageRank tables, I opted to use Excel to show the results as with only a few entries I didn't mind doing it manually. Looking at the results, an article from The Verge had the highest TF ratio of 8/256 which resulted in a larger TF-IDF. This is interesting as it ranked higher than the https://tictechbuzz.com link, which I expected to be the highest due to a higher query count; it had a larger query frequency, but also more total words at 11/481, resulting in a lower TF-IDF value. The lowest recorded value only had a single occcurence of the query over 1247 words. I was expecting it to be lower than every other URI, but I was still curious how low of a value would be calculated.

| TF-IDF | TF | IDF | URI |
|--------|-------|--------|-----|
| 1.003 | 0.031 | 32.357 | https://www.theverge.com/2020/10/15/21517293/xbox-se |
| 0.809 | 0.025 | 32.357 | https://g1.globo.com/pop-arte/games/noticia/2020/10/16, |
| 0.744 | 0.023 | 32.357 | https://tictechbuzz.com/remembering-the-classics-the-xbo: |
| 0.679 | 0.021 | 32.357 | https://www.exophase.com/user/toniazzo/?1602892894 |
| 0.615 | 0.019 | 32.357 | https://www.resetera.com/threads/xbox-brasil-dismisses-n |
| 0.615 | 0.019 | 32.357 | https://www.ign.com/articles/phil-spencer-xbox-bethesda- |
| 0.582 | 0.018 | 32.357 | https://www.engadget.com/lucasarts-double-fine-remaster |
| 0.388 | 0.012 | 32.357 | https://www.gamespot.com/articles/marvels-avengers-ps5 |
| 0.324 | 0.01 | 32.357 | https://www.gamespot.com/articles/future-bethesda-titles |
| 0.032 | 0.001 | 32.357 | https://newszoneweb.com/?p=5809 |

**Figure 1:** Chart showing the TF-IDF, TF, and IDF for the query 'Xbox' across 10 chosen URIs

# Q3

Figure 2 below shows the PageRank ranking for each URI. Each URI was checked using `http s://www.prchecker.info/check_page_rank.php`. Unfortunately, due to these being obscure pages, every single URI returned with a 0/10 ranking as per the PR estimater. This is a sharp contrast to the TD-IDF table. I was a expecting a few URIs to have some kind of rank given their TD-IDF values, but was met with zeros across the board.

| PageRank | URI |
|----------|-----|
| 0 | https://www.theverge.com/2020/10/15/21517293/xbox-series-x-s-dashboard-ui |
| 0 | https://g1.globo.com/pop-arte/games/noticia/2020/10/16/microsoft-demite-apr |
| 0 | https://tictechbuzz.com/remembering-the-classics-the-xbox-podcast-featuring-ti |
| 0 | https://www.exophase.com/user/toniazzo/?1602892894 |
| 0 | https://www.resetera.com/threads/xbox-brasil-dismisses-new-host-after-she-wa |
| 0 | https://www.ign.com/articles/phil-spencer-xbox-bethesda-ps5?utm_source=dlvr. |
| 0 | https://www.engadget.com/lucasarts-double-fine-remasters-game-pass-2352381 |
| 0 | https://www.gamespot.com/articles/marvels-avengers-ps5-and-xbox-series-x-s-v |
| 0 | https://www.gamespot.com/articles/future-bethesda-titles-dont-need-to-be-on-| |
| 0 | https://newszoneweb.com/?p=5809 |

**Figure 2:** Chart showing the results of ranking each URI with PageRank