

CS532S19: Assignment #3

Due on Sunday, March 03, 2019

Nwala, Alexander C

Giridharan Ganeshkumar

Question 1

Download the 1000 URIs from assignment 2. "curl", "wget", or "lynx" are all good candidate programs to use. We want just the raw HTML, not the images, style sheets, etc. Now use a tool to remove (most) of the HTML markup for all 1000 HTML documents. Keep both files for each URI (i.e., raw HTML and processed). Upload both sets of files to your github account

1. As shown in the below code used the urllib.request to download all the html response. We read the list of unique urls from the text file and use the urllib.request.urlopen() method to get the response object.
2. The next step was to use the read method on the response to read the raw html and store in a file.
3. The response is wrote into a separate file for each url

Listing 1: Python Script

```

1 import urllib.request
2 import json
3 import csv
4 from urllib.parse import urlparse
5
6 urlCounter = 0
7
8 for currentUrl in open("C:\\ListOfUniqueUrls.txt", "r"):
9     try:
10         urlCounter = urlCounter + 1
11         with urllib.request.urlopen(currentUrl) as res:
12             currentHtml = res.read()
13             parsed_uri = urlparse(currentUrl)
14             fileName = "C:\\RawHtmls\\" + str(urlCounter) + "_" +
15                 parsed_uri.netloc + ".txt"
16             output_file = open(fileName, "w")
17             output_file.write(str(currentHtml))
18             output_file.close()
19             print("Download Completed : " + fileName)
20     except:
21         print('Error :', currentUrl)

```

1. The next step was to install the boilerpipe using pip commands.
2. The below code block is used to extract the text alone from the html files saved from the previous step.
3. We instantiate the object for the extractor class and call the extractor.getText method to get raw text excluding the html
4. The extracted text is wrote into a new text file for each url as done in the first part of this exercise.

Listing 2: Python Script

```

1 from boilerpipe.extract import Extractor
2 from urllib.parse import urlparse
3 import glob

```

```

4 import os
5
6 path = 'C:\\RawHtmls/*.txt'
7 urlCounter = 0
8 files=glob.glob(path)
9 for file in files:
10     try:
11         urlCounter = urlCounter + 1
12         f=open(file, 'r')
13         fileName = fileName = "C:\\ProcessedText\\" + os.path.basename(f.
            name)
14         currentHtml = f.read()
15         f.close()
16         extractor = Extractor(extractor='ArticleExtractor', html=
            currentHtml)
17         currentText = extractor.getText()
18         output_file = open(fileName, "w")
19         output_file.write(str(currentText.encode("utf-8")))
20         output_file.close()
21         print("Download Completed : " + fileName)
22     except:
23         print('Error :', urlCounter)

```

Question 2

Choose a query term shadow that is not a stop word (see week 5 slides) and not HTML markup from step 1 that matches at least 10 documents. If the term is present in more than 10 documents, choose any 10 from your list. As per the example in the week 5 slides, compute TFIDF values for the term in each of the 10 documents and create a table with the TF, IDF, and TFIDF values, as well as the corresponding URIs. The URIs will be ranked in decreasing order by TFIDF values

1. As in the below code block we start by searching the text files we processed in the previous step.
2. To search we loop through all the files open it read line by line and then use re.match to look up for the text we are searching. For this purpose i choose different terms like basket-ball, genetics, has, cricket and music. The various results are uploaded to github under the SearchTermResults.zip file.
3. During the process we keep a hit count as well so that we can use it in for calculation at a later stage.

Listing 3: Python Script

```

1 import glob
2 import os
3 import csv
4 import re
5
6 searchString = "music"
7 path = 'C:\\ProcessedText/*.txt'
8 output_file_name = "C:\\SearchTermResults\\" + searchString + ".csv"
9 with open(output_file_name, "w", newline='') as output_file:
10     writer = csv.writer(output_file)
11     files=glob.glob(path)

```

```

12  for file in files:
13      #try:
14      f=open(file , 'r')
15      fileName = os.path.basename(f.name)
16      print(fileName)
17      hit_count = 0
18      regexSearch = "(.*)(" + searchString + ")(.*)"
19      for line in f:
20          if re.match(regexSearch , line):
21              currenthit_count = line.count(searchString)
22              hit_count = hit_count + currenthit_count
23      f.close()
24      currentData = [fileName , str(hit_count)]
25      writer.writerow(currentData)
26      #except:
27      #    print('Error :', fileName)
28      #print("Search Completed : " + fileName)
29  output_file.close()
30  print("Search Completed!")

```

1. We use the below code block to compute TFIDF values for the term in each of the 10 documents and create a table with the TF, IDF, and TFIDF values, as well as the corresponding URIs
2. The search term used for the case is basketball, we start by identifying the 10 uris and these are saved as a csv for further processing.
3. TF is Number of times term t appears in a document divide by Total number of terms in the document.
4. IDF is log base 2 of Total number of documents in which the term appears divide by Number of documents in whole.
5. idf Numerator is what we calculate first which is done by reading the basketball.csv we got from teh previous step during search. We just open the basketball.csv and add all the hits for evey uri response.
6. tfNumerator is simply the hit count which we have already calculated for each uri in the previous step
7. tfDenominator is the number words which can be calculated recursively calculated by adding up the length of array of word split fro each line in the document.
8. After identifying all the required values for each we just calculate Tf by dividing tfNumerator and tfDenominator, then calculate Tf by dividing idfNumerator and idfDenominator. Use math library to apply the log base 2.
9. The final step is to calculate the TfIdf by multiplying the above results and write it to the TFIDF.csv.

Listing 4: Python Script

```

1  import csv
2  import math
3
4  idfDenominator = 991
5  idfNumerator = 0

```

URL	TF	IDF	IDFLogBase2	TFIDF	TFIDFLogBase2
505_www.azcentral.com.txt	0.0186	0.2876	-1.7979	0.0054	-0.0335
435_www.northjersey.com.txt	0.0005	0.2876	-1.7979	0.0001	-0.0009
1006_www.gosanangelo.com.txt	0.0004	0.2876	-1.7979	0.0001	-0.0007
558_www.denverpost.com.txt	0.0004	0.2876	-1.7979	0.0001	-0.0007
256_www.indystar.com.txt	0.0003	0.2876	-1.7979	0.0001	-0.0005
509_www.lansingstatejournal.com.txt	0.0003	0.2876	-1.7979	0.0001	-0.0005
249_thetartan.org.txt	0.0002	0.2876	-1.7979	0.0001	-0.0004
895_leaderpost.com.txt	0.0002	0.2876	-1.7979	0.0001	-0.0004
186_www.thevwindependent.com.txt	0.0001	0.2876	-1.7979	0.0000	-0.0002
371_carver.wickedlocal.com.txt	0.0001	0.2876	-1.7979	0.0000	-0.0002

```

6 with open("C:\\SearchTermResults\\basketball.csv") as
  allSearchResultsForBasketball:
7     reader_SearchResultsForBasketball = csv.reader(
        allSearchResultsForBasketball, delimiter=',')
8     for searchResultsDetails in reader_SearchResultsForBasketball:
9         idfNumerator = idfNumerator + int(searchResultsDetails[1])
10
11 with open("C:\\TFIDF.csv", "w", newline='') as tfIdfFile:
12     writer = csv.writer(tfIdfFile)
13     currentData = ['URL', 'TF', 'IDF', 'IDFLogBase2', 'TFIDF', 'TFIDFLogBase2']
14     writer.writerow(currentData)
15     with open("C:\\SearchTermResults\\
        BasketballSearchterm10IdentifiedForFurtherProcessing.csv") as
        idfidentifiedresultFile:
16         csv_reader = csv.reader(idfidentifiedresultFile, delimiter=',')
17         tfNumerator = 0
18         tfDenominator = 0
19         for row in csv_reader:
20             tfNumerator = row[1]
21             fileName_ResponseText = "C:\\ProcessedText\\" + row[0]
22             with open(fileName_ResponseText, 'r') as f:
23                 for line in f:
24                     words = line.split()
25                     tfDenominator = tfDenominator + len(words)
26             TF = int(tfNumerator)/tfDenominator
27             IDF = idfNumerator/idfDenominator
28             IDFlogbase2 = math.log(IDF, 2)
29             TfIdf = TF * IDF
30             TfIdfBasedOnlogbase2 = TF * IDFlogbase2
31             currentData = [row[0], str(TF), str(IDF), str(IDFlogbase2), str(
                TfIdf), str(TfIdfBasedOnlogbase2)]
32             writer.writerow(currentData)
33 print('Completed')

```

Question 3

Now rank the same 10 URIs from question 2, but this time by their Page Rank.

1. In order to get the page ranks used couple of sites, and the results are shown in the below table.

URL	prchecker Site	checkpagerank Site
505_www.azcentral.com.txt	6	8
435_www.northjersey.com.txt	6	7
1006_www.gosanangelo.com.txt	5	5
558_www.denverpost.com.txt	7	8
256_www.indystar.com.txt	6	7
509_www.lansingstatejournal.com.txt	6	6
249_thetartan.org.txt	6	5
895_leaderpost.com.txt	7	6
186_www.thevwindependent.com.txt	0	4
371_carver.wickedlocal.com.txt	0	6

Question 4

Compute the Kendall Tau b score for both lists (use "b" because there will likely be tie values in the rankings). Report both the Tau value and the "p" value.

1. The below script is used to to calculate kendall tau
2. To do this imported the library scipy.stats that in turn imports kendalltau
3. The library expects the data set to be a an array of ranks, we get the ranks for various options and pass it the function
4. The various options calculated are shown below.

1. The below script is used to to calculate kendall tau
2. To do this imported the library scipy.stats that in turn imports kendalltau
3. The library expects the data set to be a an array of ranks, we get the ranks for various options and pass it the function
4. The various options calculated are shown below.

Listing 5: Python Script

```

1 from numpy.random import rand
2 from numpy.random import seed
3 from scipy.stats import kendalltau
4 import csv
5 seed(1)
6
7 data1, data2 = [], []
8 with open("C:\InputForKendallTau_b1.csv") as inputForKendallTau_b:
9     reader_inputForKendallTau_b = csv.reader(inputForKendallTau_b,
10         delimiter=',')
11     next(reader_inputForKendallTau_b)
12     for row in reader_inputForKendallTau_b:
13         #print(row[1])
14         data1.append(row[2])
15         data2.append(row[3])
16
17 coef, p = kendalltau(data1, data2)
18 print('Kendall correlation coefficient: %.3f' % coef)
19 alpha = 0.05

```

```

19 if p > alpha:
20     print('Samples are uncorrelated (fail to reject H0) p=%.3f' % p)
21 else:
22     print('Samples are correlated (reject H0) p=%.3f' % p)

```

1. TFID vs CheckPageRank

Kendall correlation coefficient: 0.200

Samples are uncorrelated (fail to reject H0) p=0.421

Question 5

Compute a ranking for the 10 URLs from Q2 using Alexa information (see week 4 slides). Compute the correlation (as per Q4) for all pairs of combinations for TFIDF, PR, and Alexa.

1. The below script is used to calculate Alexa rank
2. To do this we post a request to <http://data.alexa.com/data> with the url as parameter
3. The response is read as xml and use BeautifulSoup to get the Rank value from the reach node.
4. We also use the previous script to compare the ranks and we get the Kendall Tau u scores for different combinations as shown below.

Listing 6: Python Script

```

1 import urllib, sys, bs4, csv
2 with open("C:\RankAlexa.csv", "w", newline='') as csvFileRankAlexa:
3     writer = csv.writer(csvFileRankAlexa)
4     currentData = ['URL', 'AlexaPageRank']
5     writer.writerow(currentData)
6     with open("C:\InputForAlexa.csv") as inputForKendallTau_b:
7         reader_inputForKendallTau_b = csv.reader(inputForKendallTau_b,
8             delimiter=',')
9         next(reader_inputForKendallTau_b)
10        for row in reader_inputForKendallTau_b:
11            try:
12                xml = urllib.request.urlopen('http://data.alexa.com/data?
13                    cli=10&dat=snbamz&url=s'+ row[0]).read()
14                rank = bs4.BeautifulSoup(xml, "xml").find("REACH")['RANK']
15                currentData = [row[0], rank]
16                writer.writerow(currentData)
17            except:
18                print('Rank not found for ' + row[0])

```

1. TFID vs CheckPageRank

Kendall correlation coefficient: 0.200

Samples are uncorrelated (fail to reject H0) p=0.421

2. TFID vs Alexa

Kendall correlation coefficient : 0.156

Samples are uncorrelated (fail to reject H0) p=0.531

3. CheckPageRank vs Alexa

Kendall correlation coefficient: 0.867

Samples are correlated (reject H_0) $p=0.000$