# Assignment 3

## CS 532: Introduction to Web Science
## Spring 2018
## Hrishi Gadkari

# 1

## 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, stylesheets, etc.

from the command line:

% curl http://www.cnn.com/ > www.cnn.com
% wget -O www.cnn.com http://www.cnn.com/
% lynx -source http://www.cnn.com/ > www.cnn.com

"www.cnn.com" is just an example output file name, keep in mind that the shell will not like some of the characters that can occur in URIs (e.g., "?", "&"). You might want to hash the URIs, like:

% echo -n "http://www.cs.odu.edu/show_features.shtml?72" | md5
41d5f125d13b4bb554e6e31b6b591eeb

("md5sum" on some machines; note the "-n" in echo -- this removes the trailing newline.)

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 account.

## Answer

For the above problem to get the raw html from the 1000 URIs which we collected from Assignment 2, I wrote a program in Python 3.5 as shown in Listing 1. The following dependencies were used for the program:

- import requests

- import os

The program makes a get request using requests.get with User-Agent as Mozilla/5.0 and then downloads the context of the page using req.text [1]. While writing the output to the file, I was getting UnicodeDecodeError which was solved by giving a parameter as encoding="utf-8' [2] to the open file function. All the raw html files are stored in **raw_html** folder.

```
1  import requests
2  import os
3
4
5
6  count = 1
7  #creating directory using os package
8  os.mkdir("raw_html")
9
10 with open('1000ulinks.txt') as fp:
11
12         for line in fp:
13         #getting the raw html for each of the 1000 URIs using
                requests
14                 try:
15                         req = requests.get(line.strip(), timeout
                            =6, stream=True, headers={'User-Agent
                            ':'Mozilla/5.0'})
16                         string = req.text
17
18
19         #writng the output as encoding 'utf-8'
20                         with open('raw_html/%s.html' % count, 'w
                            ', encoding='utf-8') as outfile1:
21                                 outfile1.write(string)
22                                 outfile1.close()
23
24                         count = count + 1
25
26                 except:
27                         pass
```

Listing 1: Python program for downloading 1000 URI html content

Now to get the processed content for each of the collected raw html files, I wrote a program in Python 3.5 as shown in Listing 2. First I went through the references [5] [7] given in the question and got to know that

Python-boilerpipe is a good library for boilerplate removal and text extraction method from HTML pages. I went through the github documentation of boilerpipe [7] in order to understand how it works and implemented in my program. I installed boilerpipe using following command:

```
pip install boilerpipe3
```

The following dependencies were used for the program:

- from boilerpipe.extract import Extractor

- import os

The program iterates through each of the raw html files by reading html contents and passes to the html argument as html content, which then extracts only text from the document using getText() [7]. Although this library did not completely remove all the html content from few of the files but was sufficient for my rest of the tasks. The processed documents are stored in **processed** folder.

```
1   import os
2   from boilerpipe.extract import Extractor
3
4   # creating directory using os library
5   os.mkdir("processed")
6
7
8   count = 1
9
10  while (count < 1001):
11          with open('raw_html/%s.html' % count, 'r+', encoding='
                utf-8') as fp:
12                  #reading the collected html files from previous
                        step
13                  extractor = Extractor(extractor='
                        ArticleExtractor', html=fp.read())
14                  #extracting non-html content
15                  processed = extractor.getText()
16                  with open('processed/%s.txt' % count, 'w',
                        encoding='utf-8') as outfile1:
17                          outfile1.write(processed)
18
19
20          count = count + 1
```

Listing 2: Python program for processing all the raw html files

# 2

## Question

2. Choose a query term (e.g., "shadow") that is not a stop word
(see week 5 slides) and not HTML markup from step 1 (e.g., "http")
that matches at least 10 documents (hint: use "grep" on the processed
files). If the term is present in more than 10 documents, choose
any 10 from your list. (If you do not end up with a list of 10
URIs, you've done something wrong).

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. For
example:

Table 1. 10 Hits for the term "shadow", ranked by TFIDF.

```
TFIDF TF IDF URI
-------------
0.150 0.014 10.680 http://foo.com/
0.044 0.008 10.680 http://bar.com/
```

You can use Google or Bing for the DF estimation. To count the
number of words in the processed document (i.e., the deonminator
for TF), you can use "wc":

```
% wc -w www.cnn.com.processed
    2370 www.cnn.com.processed
```

It won't be completely accurate, but it will be probably be
consistently inaccurate across all files. You can use more
accurate methods if you'd like, just explain how you did it.

Don't forget the log base 2 for IDF, and mind your significant
digits!

https://en.wikipedia.org/wiki/Significant_figures#Rounding_and_decimal_places

**Answer**

For the above problem to get the documents for a query, I wrote a program in Python 3.5 as shown in Listing 3 and ran the query *football* over them. The following dependency was used:

- import re

The program iterates through each of the 1000 processed documents and reads the content of the file. It then searches through the content for the query using re.search() [8]. It then retrieved around 20 documents containing that query while chose 10 amongst them. The files are stored in **queried_docs** folder.

```
1  import re
2
3  count = 1
4
5  while (count < 1001):
6          with open('processed/%s.txt' % count, 'r+', encoding='
              utf-8') as fp:
7          #reading the conent from processed folder containing
              processed documents
8                  content=fp.read()
9          #searching the documents for query football
10                 if re.search('football', content):
11         #storing the documents containing the query in
              queried_docs folder
12                         with open('queried_docs/%s.txt' % count,
                              'w', encoding='utf-8') as fp1:
13                                 fp1.write(content)
14
15
16         count = count + 1
```

Listing 3: Python program to get the documents with query as football

To get the TF and IDF formulae, I went through the week 5 slides as mentioned in the question and noted the following:

$$TF = \frac{queryoccurrenceindoc}{totalnumberofindoc}$$

$$IDF = \log_2 \left( \frac{Totalnumberofdocsincorpus}{docswithterm} \right)$$

I wrote a program in Python 3.5 as shown in Listing **??** to calculate the TF-IDF, TF and IDF values for the query and documents. The following dependencies were used:

- import os

- import math

- import csv

- import linecache

The program first iterates through each of the 10 documents and reads the content. Then it stores each of the words in the document in a list using *obj*.split() [9]. For calculating TF , it then uses *obj*.count(football) [9] to count the number of occurrences in the list and divides it by the total number of words in the list using len(*obj*) [9]. To get the total number of documents in the corpus.

```
1  import math
2  import os
3  import csv
4  import linecache
5
6  # getting the curent directory path and appending to the
       queried_docs folder
7  path = os.getcwd() + '/queried_docs'
8  counts = 1
9  tuple = []
10 #Reading the files one by one from the directory containing 10
       URIs
11 for filename in os.listdir(path):
12         with open('queried_docs/%s' % filename, 'r', encoding='
               utf-8') as fp:
13                 wordstring = fp.read()
14                 #creating list of words
15                 wordlist = wordstring.split()
16                 #calculating TF, IDF, TFIDF values for each
                       document
17                 TF = (wordlist.count('football'))/(len(wordlist)
                       )
18                 DF = 47000000000/1470000000
19                 IDF = math.log(DF,2)
20                 TFIDF = TF*IDF
21                 list = os.path.splitext(os.path.basename(
                       filename))
22                 #getting the URI from the links files based on
                       the line number and document number
```

```
23              linecount = int(list[0])
24              r = linecache.getline('1000ulinks.txt',
                    linecount)
25              URI = r.strip()
26              tuple = tuple + [(TFIDF, TF, IDF, URI),]
27              fp.close()
28  #sorting the tuple based on TFIDF value in descending order
29  tuple = sorted(tuple, key=lambda tuples:tuples[0], reverse=True
       )
30
31  #writing the output in csv file
32  with open('rank.csv', 'a+', newline='') as csvfile:
33          spamwriter = csv.writer(csvfile, delimiter=',',
                  quotechar='|')
34          spamwriter.writerow(['TF-IDF', 'TF', 'IDF', 'URI'])
35          for a,b,c,d in tuple:
36
37
38              spamwriter.writerow(['%.3f' % a,'%.3f' % b,'%.3f
                    ' % c, '%s' % d])
39
40
41          csvfile.close()
```

Listing 4: Python program to compute tfidf

I visited the `worldwidewebsize.com` to the number of pages indexed by Google which was 47B, while to get the number of documents in google with the query, I got the count as 1.47B sown in Figure 1. Then log to the base 2 was calculated using math.log [10]. The values were then sorted as per TF-IDF and stored in **rank.csv** file. The values are displayed in Table 1

| TFIDF | TF | IDF | URI |
|---|---|---|---|
| 0.083 | 0.017 | 4.999 | `https://www.azcentral.com/story/ sports/high-school/2018/02/08/ business-group-trying-keep-valley- junior-college-football-alive- beyond-2018/319131002/?hootPostID= fb32e68c95bcdc943b909439c5c31ba3` |
| 0.074 | 0.015 | 4.999 | `http://www.fox26houston.com/news/ california-lawmakers-introduce-safe- youth-football-act-by-setting-minimum- age-for-tackle-football?utm_source= dlvr.it&utm_medium=twitter` |
| 0.071 | 0.014 | 4.999 | `http://abc7.com/sports/california- bill-would-ban-tackle-football-before- high-school/3058813//` |
| 0.065 | 0.013 | 4.999 | `http://www.chicagotribune.com/suburbs/ daily-southtown/sports/ct-sta-pat- disabato-column-st-0211-20180209- story.html` |
| 0.035 | 0.007 | 4.999 | `http://www.latimes.com/politics/ essential/la-pol-ca-essential- politics-updates-california-would- bar-organized-tackle-1518130990- htmlstory.html` |
| 0.023 | 0.005 | 4.999 | `http://www.telegram.com/news/20180209/ five-st-johns-football-players-among- crop-of-area-athletes-finalizing- college-plans?platform=hootsuite` |
| 0.022 | 0.004 | 4.999 | `http://lompocrecord.com/sports/high- school/football/another-taua-to-head- to-reno---this-time/article_a7473c8e- 13a8-5f79-82ac-8747900922df.html` |
| 0.021 | 0.039 | 4.999 | `http://www.dailymail.co.uk/sport/ football/article-5374081/Man-United- star-Mata-fears-football-losing- soul.html?ITO=1490&ns_mchannel=rss&ns_ campaign=1490` |
| 0.009 | 0.002 | 4.999 | `http://www.bbc.co.uk/sport/football/ 43010230` |
| 0.006 | 0.001 | 4.999 | `http://www.football-chairman.com` |

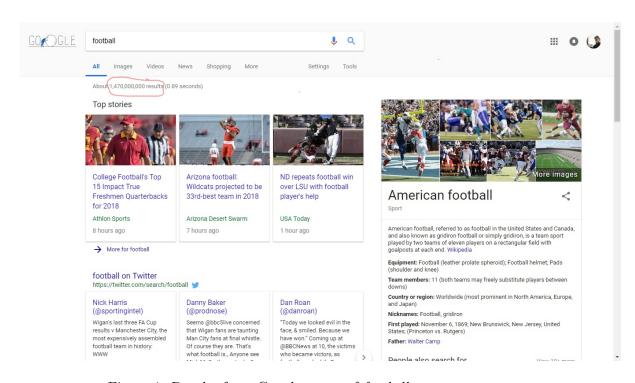Table 1: 10 URIs found containing football, with calculations TFIDF, TF and IDF

Figure 1: Results from Google query of *football*

# 3

## Question

3.  Now rank the same 10 URIs from question #2, but this time
by their PageRank.  Use any of the free PR estimaters on the web,
such as:

http://pr.eyedomain.com/
http://www.prchecker.info/check_page_rank.php
http://www.seocentro.com/tools/search-engines/pagerank.html
http://www.checkpagerank.net/

If you use these tools, you'll have to do so by hand (they have
anti-bot captchas), but there are only 10 to do.  Normalize the
values they give you to be from 0 to 1.0.  Use the same tool on all
10 (again, consistency is more important than accuracy).  Also
note that these tools typically report on the domain rather than
the page, so it's not entirely accurate.

Create a table similar to Table 1:

Table 2.  10 hits for the term "shadow", ranked by PageRank.

PageRank URI
-----------
0.9 http://bar.com/
0.5 http://foo.com/

Briefly compare and contrast the rankings produced in questions 2
and 3.

**Answer**

To calculate the Page Rank for each of the 10 URIs from the previous question I used the `http://www.checkpagerank.net/` uri as mentioned in the question. I have stored the values in **ptdf.csv**. File. Also the values are displayed in the Table 2. First I thought that URIs with high TFIDF values will have higher Page Rank, but this is not the case. From the Table 2 we can see that the TFIDF values varies with Page Ranks. Few higher paged rank URIs have lower TFIDF values while a couple of high TFIDF values have lower Page ranks. I have sorted the URIS considering first the Page Rank and then the respective higher TFIDF values since few URIs have the same Page Rank.

| Page Rank | TFIDF | URI |
|---|---|---|
| 0.9 | 0.0351 | `http://www.latimes.com/politics/` `essential/la-pol-ca-essential-politics-` `updates-california-would-bar-organized-` `tackle-1518130990-htmlstory.html` |
| 0.9 | 0.021 | `http://www.dailymail.co.uk/sport/` `football/article-5374081/Man-United-` `star-Mata-fears-football-losing-` `soul.html?ITO=1490&ns_mchannel=rss&ns_` `campaign=1490` |
| 0.9 | 0.009 | `http://www.bbc.co.uk/sport/football/` `43010230` |
| 0.8 | 0.083 | `https://www.azcentral.com/story/` `sports/high-school/2018/02/08/` `business-group-trying-keep-valley-` `junior-college-football-alive-` `beyond-2018/319131002/?hootPostID=` `fb32e68c95bcdc943b909439c5c31ba3` |
| 0.8 | 0.065 | `http://www.chicagotribune.com/suburbs/` `daily-southtown/sports/ct-sta-pat-` `disabato-column-st-0211-20180209-` `story.html` |
| 0.7 | 0.071 | `http://abc7.com/sports/california-bill-` `would-ban-tackle-football-before-high-` `school/3058813/` |
| 0.7 | 0.023 | `http://www.telegram.com/news/20180209/` `five-st-johns-football-players-among-` `crop-of-area-athletes-finalizing-` `college-plans?platform=hootsuite` |
| 0.6 | 0.074 | `http://www.fox26houston.com/news/` `california-lawmakers-introduce-safe-` `youth-football-act-by-setting-minimum-` `age-for-tackle-football?utm_source=dlvr.` `it&utm_medium=twitter` |
| 0.6 | 0.022 | `http://lompocrecord.com/sports/high-` `school/football/another-taua-to-head-to-` `reno---this-time/article_a7473c8e-13a8-` `5f79-82ac-8747900922df.html` |
| 0.5 | 0.006 | `http://www.football-chairman.com` |

Table 2: Page Rank and TFIDF Comparison of 10 URIs

# 4

## Question

```
===================================================
======Question 4 is for 3 points extra credit======
===================================================
```

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.

See:
http://stackoverflow.com/questions/2557863/measures-
of-association-in-r-kendalls-tau-b-and-tau-c
http://en.wikipedia.org/wiki/Kendall_tau_rank_correlation_coefficient#Tau-b
http://en.wikipedia.org/wiki/Correlation_and_dependence

**Answer**

For the above problem I first went through `http://en.wikipedia.org/wiki/Correlation_and_dependence` and got to know that it is basically calculating correlaton between pairs of values. The Tau value can have a coefficient between -1 to 1, which means that if the value is 1 or greater it has positive correlation between TFIDF and Page Rank values, meaning they are similar. If the value is zero or close, there is no correlation between them, its a tie. And if the value is closer to -1 or less, there is negative correlation, meaning the Page Rank is the opposite of the TFIDF value. Then I went through these references [3] and [4] and got to know that there is a Kendall package for R language, which gives a summary of values as output. The program reads the TFIDF and Page Rank lists from **ptdf.csv** file as shown in Listing 5.

```
1  #loading library Kendall
2  library(Kendall)
3  setwd(getwd())
4  #reading contents from csv file
5  csv <- read.table('ptdf.csv', header = TRUE, sep = ",")
6  #producing summary of results
7  summary(Kendall(csv$PageRank, csv$TFIDF))
```

Listing 5: R script to produce summary of Kendall's results

The summary produced:

```
Score =  0 , Var(Score) = 117.4667
denominator =  41.42463
tau = 0, 2-sided pvalue =1
```

From the above summary we can see that the tau value is zero meaning there is no correlation between the two pairs of values.

# 5

## Question

```
==================================================
======Question 5 is for 3 points extra credit======
==================================================
```

5. Compute a ranking for the 10 URIs 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.

**Answer**

For the above problem I calculated the Alexa rank from `https://www.alexa.com/siteinfo` for each of the 10 URIs which we got in question 2. The values are stored in **alexa_rank.csv** file as shown in Table 3.

I calculated the correlation for each of the TFIDF, PR and Alexa values as per Kendall Tau_b score from question 4. The program is written in R as shown in Listing 6. The values are rad from **alex.r csv** file It produced some interesting results.

```
1   ##loading library Kendall
2   library(Kendall)
3   setwd(getwd())
4   #reading contents from csv
5   csv <- read.table('alexa.csv', header = TRUE, sep = ",")
6   #summary for PR−TFIDF
7   summary(Kendall(csv$PageRank, csv$TFIDF))
8   #summary for TFIDF−Alexa
9   summary(Kendall(csv$TFIDF, csv$Alexa))
10  #summary for PR−Alexa
11  summary(Kendall(csv$PageRank, csv$Alexa))
```

Listing 6: R program to produe summary of Kendall results

Summary produced for TFIDF and Alexa lists:

```
Score =  2 , Var(Score) = 124
denominator =  44.49719
tau = 0.0449, 2-sided pvalue =0.92844
```

From the above results we can see that the tau value is positive meaning there is some correlation between them, although I thought it would be zero.

Summary produced for Page Rank and Alexa lists:

```
Score =  -39 , Var(Score) = 118.3333
denominator =  41.89272
tau = -0.931, 2-sided pvalue =0.00047717
```

From seeing the above results, I was amazed to see that the Tau_b value is negative and closer to -1 meaning there is negative correlation between the values, although I thought it would be positive. But then I realized that Alexa scores the URIs opposite to that of the Page Rank, which means that there is correlation and similarity between them.

The correlation between TFIDF and Page Rank has already been tested in question 4.

| Alexa Rank | URI |
|---|---|
| 101 | `http://www.bbc.co.uk/sport/football/43010230` |
| 174 | `http://www.dailymail.co.uk/sport/football/article-5374081/Man-United-star-Mata-fears-football-losing-soul.html?ITO=1490&ns_mchannel=rss&ns_campaign=1490` |
| 830 | `http://www.latimes.com/politics/essential/la-pol-ca-essential-politics-updates-california-would-bar-organized-tackle-1518130990-htmlstory.htm` |
| 2172 | `https://www.azcentral.com/story/sports/high-school/2018/02/08/business-group-trying-keep-valley-junior-college-football-alive-beyond-2018/319131002/?hootPostID=fb32e68c95bcdc943b909439c5c31ba3` |
| 7094 | `http://www.chicagotribune.com/suburbs/daily-southtown/sports/ct-sta-pat-disabato-column-st-0211-20180209-story.html` |
| 8716 | `http://abc7.com/sports/california-bill-would-ban-tackle-football-before-high-school/3058813/` |
| 49281 | `http://www.telegram.com/news/20180209/five-st-johns-football-players-among-crop-of-area-athletes-finalizing-college-plans?platform=hootsuite` |
| 125065 | `http://www.fox26houston.com/news/california-lawmakers-introduce-safe-youth-football-act-by-setting-minimum-age-for-tackle-football?utm_source=dlvr.it&utm_medium=twitter` |
| 399889 | `http://lompocrecord.com/sports/high-school/football/another-taua-to-head-to-reno---this-time/article_a7473c8e-13a8-5f79-82ac-8747900922df.html` |
| 6633445 | `http://www.football-chairman.com` |

Table 3: Alexa Rank for 10 URIs

# 6

## Question

```
==================================================
======Question 6 is for 2 points extra credit======
==================================================
```

6.  Give an in-depth analysis, complete with examples,
graphs, and all other pertinent argumentation for
Kristen Stewart's (of "Twilight" fame) Erdos-Bacon number.

## Answer

The Erdos-Bacon number is nothing but the sum of Erdos number and Bacon number. To get the Bacon number for Kristen Stewart I visited the `https://oracleofbacon.org/movielinks.php` uri which gave me a Bacon number of 2 for Kristen as shown in Figure 2.



Figure 2: Bacon number for Kristen Stewart

To find out the Erdos number for Kristen Stewart I searched through google to know the name of the paper she co-authored. I got to know that she published a paper on machine learning: `https://arxiv.org/pdf/1701.04928v1.pdf` along with Bhautik Joshi and David Shapiro. Then on further googling, I encountered a link`https://news.ycombinator.com/item?id=13446478` where I came to know that Bhautik Joshi had co-authored in *A quantitative assessment of approaches to mesh generation for surgical simulation* along with Sebastein Ourselin who has a Erdos number of 4 calculated from `https://zbmath.org/collaboration-distance/?a=ourselin.sebastien&b=erdos.paul` as shown in Figure3.

The following graph was produced using `Webgraphviz.com`:

From the above graph we get the Erdos number for Kristen Stewart as 6. Thus the Erdos-Bacon number for Kristen Stewart is $6 + 2 = 8$.
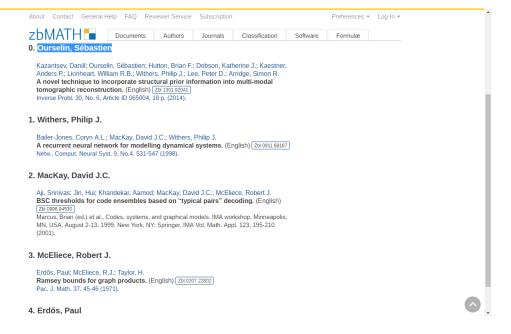
Figure 3: Erdos number for Sebastein Ourselin

# WebGraphviz is [Graphviz](#) in the Browser

Enter your graphviz data into the Text Area:

(Your Graphviz data is private and never harvested)

| Sample 1 | Sample 2 | Sample 3 | Sample 4 | Sample 5 |
|---|---|---|---|---|

```
digraph G {
  "Kristen Stewart" -> "Bhautik Joshi"
  "Bhautik Joshi" -> "Sebastien Ourselin"
  "Sebastien Ourselin" -> "Philip J Withers"
  "Philip J Withers" -> "David J.C. MacKay"
  "David J.C. MacKay" -> "Robert J. McEliece"
  "Robert J. McEliece" -> "Paul Erdos"
}
```

Generate Graph!

Figure 4: Step1: Erdos graph for Kristen Stewart

Figure 5: Step2: Erdos graph for Kristen Stewart

# 7

## Question

```
===================================================
======Question 7 is for 2 points extra credit======
===================================================
```

7.  Build a simple (i.e., no positional information) inverted file
(in ASCII) for all the words from your 1000 URIs.  Upload the entire
file to github and discuss an interesting portion of the file in
your report.

## Answer

For the above problem which I found very interesting, wrote a program in Python 3.5 as shown in Listing 7. First I tested on few of the processed documents from question 1 and created wordlist for each of the document using *obj*.split() [9]. I observed that few of the documents had html content which was appending non-linguistic words to the list. I realized that I had to use a different boilerplate removal library, so I went through the [5] and figured out JusText [6] is perfect for this. The following dependencies were used:

- import justext

- import re

```
1   import justext
2   import re
3
4   para = ''
5   d = {}
6   inv = []
7   count = 1
8
9   while (count < 1001):
10          with open('raw_html/%s.html' % count, 'r', encoding='utf
                 -8') as fp:
11                  #reading the files for raw html content
12                  wordstring = fp.read()
13                  #extracting linguistic sentences from the
                        content by excluding the boilerplate content
14                  paragraphs = justext.justext(wordstring, justext
                        .get_stoplist("English"))
15                  for paragraph in paragraphs:
16
17                          if not paragraph.is_boilerplate:
18                                  para = para + ' ' + (paragraph.
                                        text)
19                  #creating a list of words
20                  wordlist = para.split()
21
22                  for w in wordlist:
23                  #excluding non-alphanumeric characters
24                          w = re.sub(r'\W+', '', w)
25
26                  #checkig for duplicates
27                          if w not in inv:
28                          # check if string not empty
```

```
29                              if w:
30                                      inv.append(w)
31
32              for w1 in inv:
33                      #appending values to keys in dictionary
34                      try:
35                              d[w1].append(count)
36                      except KeyError:
37                              d[w1] = [count]
38                              pass
39          count = count + 1
40          #writing the result in invertedindex.txt file in ascii
                format
41  with open('invertedindex.txt', 'w', encoding='ascii') as pp:
42              for k, v in d.items():
43                      pp.write(str(k) + ' >>> '+ str(v) + '\n\
                n')
```

Listing 7: Python script to create an inverted index

The program first reads the html documents one by one from the raw_html folder which we got from question 1 and reads the contents into justext.justext() [6] function. It then extracts words using paragraph.text [6]. Then it filters all non-alphanumeric characters from the list using re.sub() [8]. Then appending the words one by one as key to a dictionary with values as the document number it writes the output in **invertedindex.txt** file. A sample output is shown below. What we can see that it starts from document 2 because document 1 had no content.
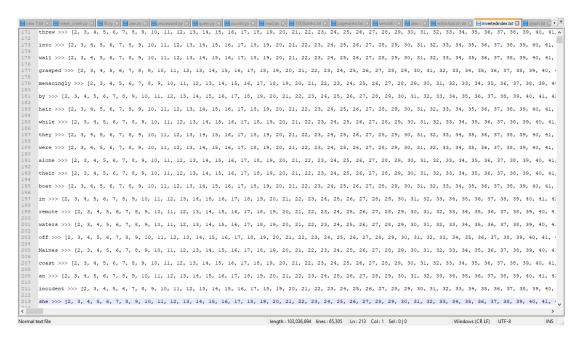
Figure 6: Inverted Index Output

# References

[1] "Requests: HTTP for Humans." equests: HTTP for Humans - Requests 2.18.4 documentation, n.d. Web. February 20, 2018. `http://docs.python-requests.org/en/master/`.

[2] "Stack Overflow."Switching to Python 3 causing UnicodeDecodeError. N.p.,n.d. Web. February 20, 2018. `https://stackoverflow.com/questions/23917729/switching-to-python-3-causing-unicodedecodeerror`

[3] "Kendall rank correlation coefficient." Wikipedia - Kendall Rank Correlation coefficient. Wikipedia, n.d. Web. February 20, 2018. `https://en.wikipedia.org/wiki/Kendall_rank_correlation_coefficient#Tau-b`

[4] "CRAN - Package Kendall." Kendall Rank Correlation. n.d. Web. February 20, 2018. `https://cran.r-project.org/web/packages/Kendall/Kendall.pdf`

[5] Alexander Nwala. "Web Science and Digital Libraries Research ." 2017-03-20: A survey of 5 boilerplate removal methods. N.p., Web. February 20, 2018. `http://ws-dl.blogspot.com/2017/03/2017-03-20-survey-of-5-boilerplate.html`.

[6] Miso-belica. "miso-belica/jusText .", GitHub, March 05, 2017, n.d., Web. February 20, 2018. `https://github.com/miso-belica/jusText`.

[7] Misja. "misja/python-boilerpipe." Github Repository. N.p., October 03, 2017. Web. February 20, 2018.`https://github.com/misja/python-boilerpipe`.

[8] "6.2. re - Regular expression operations." 6.2. re - Regular expression operations - Python 3.6.4 documentation. n.d., Web. February 20, 2018. `https://docs.python.org/3/library/re.html`.

[9] William Turkel-Adam Crymble. "Counting Word Frequencies with Python" Programming Historian. N.p., n.d. Web. February 20, 2018.`https://programminghistorian.org/lessons/counting-frequencies`.

[10] "9.2. math - Mathematical functions." 9.2. math - Mathematical functions - Python 2.7.14 documentation. N.p., n.d. Web. February 20, 2018. `https://docs.python.org/2/library/math.html`.