

Assignment 8

CS532-s16: Web Sciences

Spring 2016

John Berlin

Generated on April 6, 2016

Index

Question 1, 2
Question 2, 11
Question 3, 15
Question 4, 23
Question 5, 24

Question 1

1. Create a blog-term matrix. Start by grabbing 100 blogs; include:

```
http://f-measure.blogspot.com/  
http://ws-dl.blogspot.com/
```

and grab 98 more as per the method shown in class. Note that this method randomly chooses blogs and each student will separately do this process, so it is unlikely that these 98 blogs will be shared among students. In other words, no sharing of blog data. Upload to github your code for grabbing the blogs and provide a list of blog URIs, both in the report and in github..

Use the blog title as the identifier for each blog (and row of the matrix). Use the terms from every item/title (RSS) or entry/title (Atom) for the columns of the matrix. The values are the frequency of occurrence. Essentially you are replicating the format of the "blogdata.txt" file included with the PCI book code. Limit the number of terms to the most "popular" (i.e., frequent) 500 terms, this is *after* the criteria on p. 32 (slide 7) has been satisfied.

Answer

In getting the 100 blogs and generating the blog-term matrix I wrote two python scripts *getFeeds.py* and *processData.py*. The first script *getFeeds.py* which is seen in listing 2 does the job of getting the blog feed data. Then *processData.py* seen in listing 3 does the generation of the blog-term matrix file plus generation of following questions data.

Getting the feed data turned out to require a little bit more effort than I was expecting. The extra effort came into play when I came across blogs with no titles(blog titles). These no blog title blogs as a majority had little to no text as I found out when a great deal of my words were blanks. Also adding to the extra effort was realizing that the random next blog was not truly random rather pseudo-random as I had duplicates. Besides that the retrieval of the blog data was rather straightforward.

The blogspot API was rather straightforward and required few calls. The uris used to interact with the blogspot API can be seen on lines 12,13,14, and 15 in listing 2. To speed up the process of getting the feed entries I appended *?max-results=200* to end of the feeds uri(*feeds/posts/default*). By doing this I get a maximum of 200 entries per get request, which greatly limits the number of API calls required to get paginated entries for blogs with over 200 posts.

I will explain the code seen in listing 2 at first generally then break down the methods used to accomplish the task. As always consult the code seen in listing 2 for more details in the comments.

1. Entire process(getData)
 - (a) Check to see if the data file directory exists if not create it
 - (b) Set up session and add user agent so we appear less like a robot
 - (c) Consume and process the WebScience and Digital Libraries blog posts
 - (d) Consume and process F-Measure blog posts
 - (e) Ask for next blog 98 times, consume and process its posts
 - (f) Write blog urls to file and write blog terms to json file called blog-data.json
2. get98
 - (a) While count is less than 98
 - (b) Get the next blog
 - (c) Check to see if we have already gotten it if we have choose another otherwise consume_all
 - (d) If what we got is nothing choose another
 - (e) Add what we got to the data collection, increment counter by one and choose another
3. consume_all(start, sesh)
 - (a) Ask the api for the initial set of feed entries
 - (b) Parse the feed and get the title, if the blog has no title or if the number of entries is less than 25 choose another
 - (c) Check to see if we have more feed entries
 - (d) Process and flatten the initial set of feed entries
 - (e) While we have more entries(pagination), ask for next set of feed entries, process and flatten, then check again
 - (f) Do preliminary word count as the number of words if not reduces can cause the file size to be to big
 - (g) Finally return the title and the blogs terms
4. process_text(text)
 - (a) Process the feeds text using BeautifulSoup to clean out html, remove all non alpha characters then make the text lower case

- (b) Use nltk(natural language toolkit) to tokenize the text into words, filter the words by checking to see if they are not English language stop words
 - (c) Return list of valid words
5. `check_next(text)`
- (a) Use BeautifulSoup to process the entire feed
 - (b) Extract the next link(next set of feed entries) by looking for *link, type: application/atom+xml, rel: next*
 - (c) If that link exists the returned list will be of size 1, return the link and true otherwise false and none

The list of 98 urls gotten can be seen in listing 1.

Now I will also explain the other file `processData.py` seen in listing 3, here in its entirety for the parts that go with question one. In the remain question sections I will go into detail about what was required for the parts pertaining to those remaining questions. Please note that I used `clusters.py` file that accompanied the book found in github repo containing the example code.

For processing the data for the top 500 terms I created a class called `feed`. It allows for quick access to the preliminary word count and stem count. For another part for the next questions I stemmed the words to get a better clusters as the original method returned a very sparse term matrix leading to errors in distance calculations. Stemming the words I again used the natural language tool kit's(nltk) English stemmer to produce a correct stemming for the English language.

The two methods in `processData.py` listing 3 that partaine to question one are `generate_blogfile` and `generate_blogfile_stem`. Both execute as such

1. `generate_blogfile`
 - (a) Read in the blog term json file. I use the functional python library `seq` class to allow for short and concise expression of the processing.
 - (b) Get the top 500 by
 - (c) Flattening each blogs word count
 - (d) Filtering word count tuples for counts greater than 10(non-stemmed),1(stemmed)
 - (e) Do a total word count
 - (f) Filter by the faked `tfidf` function
 - (g) Sort the data in descending order
 - (h) Take the top 500
 - (i) Transform the data to words only and make it a list
 - (j) Sort the top 500 words alphabetically
 - (k) Write the words out the blog-term matrix file

2. generate_blogfile_stem

- (a) Behaves exactly like generate_blogfile except use the stemmed word count and gets words with a count greater than 1 (addendum to value used in 1d)

I mentioned that words with intra-document counts of 10 were used for the non-stemmed version in the explanation of 1d. Before using those values, the default count of 1 was used which lead to divide by zero exceptions in the MDS calculation. On inspection I believe it was caused by the sparse nature of the matrix terms overall. This change was done well after I had written the code for the initial implementation which lead to my usage of stemmed words. Also with this change I was able to remove my changes to the original clusters.py file which checked and poorly corrected the divide by zero exception.

1 <http://deadbeatfanzine.blogspot.com>
2 <http://thebeautifultrashart.blogspot.com>
3 <http://hiiijsaaackie.blogspot.com>
4 <http://encorenorthernireland.blogspot.com>
5 <http://blog.spinitron.com>
6 <http://floorshimezipperboots.blogspot.com>
7 <http://lafotografiaefectistaabstracta.blogspot.com>
8 <http://pithytittlehere.blogspot.com>
9 <http://thesportsmith.blogspot.com>
10 <http://ahtapotunbahcesi.blogspot.com>
11 <http://dancingincirclesnow.blogspot.com>
12 <http://rosiegigga2media.blogspot.com>
13 <http://thekidsarecoming.blogspot.com>
14 <http://www.radioshower.com>
15 <http://desolationrowrecords.blogspot.com>
16 <http://isyelili.blogspot.com>
17 <http://theindiefriend.blogspot.com>
18 <http://parishradio.blogspot.com>
19 <http://listeningear.blogspot.com>
20 <http://stonehillsketchbook.blogspot.com>
21 <http://mattsbunker.blogspot.com>
22 <http://www.thejeopardyofcontentment.com>
23 <http://theworldsfirstinternetbaby.blogspot.com>
24 <http://sixtyat60.blogspot.com>
25 <http://jamiemclelland.blogspot.com>
26 <http://www.samtasticreview.com>
27 <http://kidchair.blogspot.com>
28 <http://boglemethursday.blogspot.com>
29 <http://markeortega.blogspot.com>
30 <http://ihatethe90s.blogspot.com>
31 <http://storiesfromthecityradiovalencia.blogspot.com>
32 <http://thefastbreakofchampions.blogspot.com>
33 <http://skinnynshoes.blogspot.com>
34 <http://semregrasluispink.blogspot.com>
35 <http://elijace.blogspot.com>
36 <http://marialombideezpeleta.blogspot.com>
37 <http://ilovetotaldestruction.blogspot.com>
38 <http://ngaio1619.blogspot.com>
39 <http://flatbasset.blogspot.com>
40 <http://hartsdesire.blogspot.com>
41 <http://sixeyes.blogspot.com>
42 <http://rantsfromthepants.blogspot.com>
43 <http://spicyseatdolphin.blogspot.com>
44 <http://my-name-is-blue-canary.blogspot.com>
45 <http://dana9morgan.blogspot.com>
46 <http://moontopples.blogspot.com>
47 <http://dinosaursarefun.blogspot.com>
48 <http://ashleyemwarren.blogspot.com>
49 <http://noradiorecs.blogspot.com>
50 <http://thehubkxci.blogspot.com>
51 <http://lostintheshuffle899.blogspot.com>
52 <http://flowradio.blogspot.com>
53 <http://mysteryfallsdown.blogspot.com>
54 <http://thenightmail.blogspot.com>
55 <http://seveninchesisenough.blogspot.com>
56 <http://didnotchart.blogspot.com>
57 <http://rodshone.blogspot.com>

```

58 http://mileinmine.blogspot.com
59 http://mesastivromia.blogspot.com
60 http://theonionfield.blogspot.com
61 http://globalgoon.blogspot.com
62 http://beckysharpfashionblog.blogspot.com
63 http://ourpodcastcouldbeyourlife.blogspot.com
64 http://steel-city-rust.blogspot.com
65 http://maggotcaviar.blogspot.com
66 http://www.sonology.com
67 http://thetremagazine.blogspot.com
68 http://flipmpip.blogspot.com
69 http://psychfolkmusic.blogspot.com
70 http://kunstlertreu.blogspot.com
71 http://mondaywakeup.blogspot.com
72 http://thepowerofindependenttrucking.blogspot.com
73 http://richardwhitten.blogspot.com
74 http://marshwiggles.blogspot.com
75 http://dustandwaterstudios.blogspot.com
76 http://harisphotonet.blogspot.com
77 http://mo-forgetaboutit.blogspot.com
78 http://kistefm.blogspot.com
79 http://mobbie2.blogspot.com
80 http://hani-bittersweet.blogspot.com
81 http://jbreitling.blogspot.com
82 http://turnitupjack.blogspot.com
83 http://davecromwellwrites.blogspot.com
84 http://mtjrrantsravesonmusic.blogspot.com
85 http://ps-music.blogspot.com
86 http://superchicken46.blogspot.com
87 http://cherryarea.blogspot.com
88 http://jlmldhlcm1516.blogspot.com
89 http://trembleunderboomlights.blogspot.com
90 http://onestunningsingleegg.blogspot.com
91 http://chantellesmedia2.blogspot.com
92 http://everydaymusicportland.blogspot.com
93 http://karldrinkwater.blogspot.com
94 http://jakobclapensteam.blogspot.com
95 http://adrianomarquesblog.blogspot.com
96 http://www.chrisanne-grise.com
97 http://campusbuzzwsou.blogspot.com
98 http://jojobethkatiehannahlcm1516.blogspot.com

```

Listing 1: The 98 uris


```

1 import json as jjson
2 import re
3 from operator import add
4 import os
5 import feedparser
6 import nltk
7 import requests
8 from bs4 import BeautifulSoup
9 from functional import seq
10 from nltk.corpus import stopwords
11
12 fmeasure = "http://f-measure.blogspot.com/feeds/posts/default?max-
    results=200"
13 fmeasure_next = "https://www.blogger.com/next-blog?navBar=true&
    blogID=3471633091411211117"
14 wsdl = "http://ws-dl.blogspot.com/feeds/posts/default?max-results
    =200"
15 bspot_feed = "/feeds/posts/default?max-results=200"
16
17 # nuke all none-alpha chars
18 removeExtra = re.compile('[^a-zA-Z]')
19
20 langs = ["dutch", "finnish", "german", "italian", "portuguese", "
    spanish", "turkish", "danish", "english",
21          "french", "hungarian", "norwegian", "russian", "swedish"]
22 langStopWords = {}
23 for lang in langs:
24     langStopWords[lang] = stopwords.words(lang)
25
26
27 def check_next(text):
28     # check for the next button ie pagination of blog pages
29     soup = BeautifulSoup(text, "lxml-xml")
30     next_page = soup.find_all('link', attrs={'type': 'application/
        atom+xml', 'rel': 'next'})
31     # if there is a next page our next_page list will always be 1
        otherwise its 0
32     # that means we have consumed all the pages for the blog
33     if len(next_page) > 0:
34         nl = next_page[0].attrs['href']
35         return True, nl
36     return False, None
37
38
39 def process_text(text):
40     # clean text
41     t = removeExtra.sub(' ', BeautifulSoup(text.summary, 'html5lib'
        ).text).lower()
42     ret = []
43     # tokenize according to word and emit word if it is not a
        stopword
44     for word in nltk.word_tokenize(t):
45         if word not in langStopWords['english']:
46             ret.append(word)
47     return ret
48
49

```

```

50 def consume_all(start, sesh):
51     # the text of the entry titles
52     text = []
53     # the request for the first page of the blogs atom feed
54     r = sesh.get(start) # type: requests.Response
55     # parse the feed
56     feed = feedparser.parse(r.text) # type: feedparser.
FeedParserDict
57     # some blogs do not have titles so I must check for them
58     try:
59         t = feed['feed']['title']
60     except KeyError as e:
61         # some blogs do not have a title so I will skip them
62         print(e)
63         r.close()
64         return None
65     # see if we have more feed pages
66     good, nl = check_next(r.text)
67
68     if len(feed.entries) < 25:
69         return None
70     # All these operations are very good for more of a functional
approach
71     # to much work otherwise
72     '''
73     flatmap:
74         for entry <- feed, token <- process_text(entry): yeild
token
75     '''
76     text.extend(seq(feed.entries).flat_map(process_text).to_list())
77     while good:
78         r = sesh.get(nl)
79         feed = feedparser.parse(r.text)
80         r.close()
81         text.extend(seq(feed.entries).flat_map(process_text))
82         good, nl = check_next(r.text)
83
84     # since the number of words can become extremely large we need
to do a
85     # reduce before emitting the words
86     text = seq(text).map(lambda word: (word, 1)) \
87         .reduce_by_key(add) \
88         .order_by(lambda x: x[1]).to_dict()
89
90     return {'title': t, 'text': text}
91
92
93 def get98(sesh):
94     the98 = {}
95     counter = 0
96     while counter < 98:
97         r = sesh.get(fmeasure_next) # type: requests.Response
98         rurl = r.url[:r.url.rindex("/")]
99         # I had the random button give the same blog 68 times
before
100        # so I will simply look for another one if I have seen it
before

```

```

101         if the98.get(rurl, None) is not None:
102             r.close()
103             print("Processed URL already %s" % rurl)
104             continue
105         else:
106             print(r.url[:r.url.rindex("/")])
107             got = consume_all(rurl + bspot_feed, sesh)
108             # if the consuming of the feed was bad choose another
109             if got is None:
110                 r.close()
111                 print("continuing on in get98")
112                 continue
113             the98[rurl] = got
114             counter += 1
115             print(counter)
116             r.close()
117     return the98
118
119
120 def getData():
121     data = {}
122     # have a useragent so we do not look like a robot
123     useragent = 'Mozilla/5.0 (X11; Ubuntu; Linux x86_64; rv:45.0)
124     Gecko/20100101 Firefox/45.0'
125     session = requests.Session() # type: requests.Session
126     session.headers.update({'User-Agent': useragent})
127     try:
128         # get the data
129         wsdlldata = consume_all(wsdl, session)
130         data[wsdlldata['title']] = wsdlldata['text']
131         fmeasuredata = consume_all(fmeasure, session)
132         data[fmeasuredata['title']] = fmeasuredata['text']
133         data98 = get98(session)
134         # write out the uris used for the 98 other blogs
135         with open("98bloguris.dat", "w+") as out:
136             for url, d in data98.items():
137                 data[d['title']] = d['text']
138                 out.write("%s\n" % url)
139             session.close()
140         # write out our data to json
141         with open("blogdata.json", "w+") as out:
142             out.write(json.dumps(data, indent=1))
143     except Exception as e:
144         print(e)
145         session.close()
146
147
148 if __name__ == "__main__":
149     if not os.path.exists(os.getcwd()+"/datafiles"):
150         os.makedirs(os.getcwd()+"/datafiles")
151     getData()

```

Listing 2: Get Blog Feeds

Question 2

2. Create an ASCII and JPEG dendrogram that clusters (i.e., HAC) the most similar blogs (see slides 12 & 13). Include the JPEG in your report and upload the ascii file to github (it will be too unwieldy for inclusion in the report).

Answer

The rest of the code in `processData.py` seen in listing 3 deals with generating the answers to questions 2 through 4. The methods used to generate those answers are `do_non_stem()` and `do_stemmed()`. Both operate as such:

1. Generate the respective blog term matrix file
2. Read the file in using the `clusters.py` file (from the code accompanying the programming cognitive intelligence book).
3. Do hierarchical clustering
4. Generate text denogram and jpg denogram
5. Perform kmeans clustering and log iterations and centroid values
6. Perform dimension reduction
7. Generate dimension reduced cluster jpg

Please note that unlike the code from the book

I included stemming as with out it we are choosing the top 500 terms out of unique 2692 terms in total.

The dendrogram generated for the nonstemmed version is seen in figure 1 and the text version can be found in the `datafiles` folder file name, `blog-top500_asciidenom.txt`. The dendrogram is quite large so fitting it into this report was difficult and the file name is `blogtop500_deno.jpg` and is found in the `datafiles` directory.

Of note the Web Science and Digital Libraries Research Group got grouped under Boggle Me Thursdays, which is under Floorshine Zipper Boots which happens to be under F-Measure. I find it interesting that the personal blog for the head of the WSDL is above the group its self.

The stemmed version is seen in figure 2. F-Measure is grouped with talk radio blogs and media blogs whereas the WSDL blog is grouped with blogs such as What am I doing (above) and The Night Mail (below). These results are more surprising to me especially for the WSDL blog.

```

1 from collections import defaultdict
2 from operator import add
3
4 from functional import seq
5 from nltk.stem.snowball import EnglishStemmer
6
7 import clusters
8
9
10 # this class represents the word data for a particular feed
11 class feed:
12     # pass flag if we are to do the stemming
13     def __init__(self, fentry, doStem=False):
14         self.title = fentry[0]
15         self.wordCount = fentry[1]
16         self.stemCount = defaultdict(int)
17         if doStem:
18             self._stem_count()
19
20     def _stem_count(self):
21         eng = EnglishStemmer()
22         for word in self.wordCount.keys():
23             self.stemCount[eng.stem(word)] += 1
24
25     def words(self):
26         return set(self.wordCount.keys())
27
28     def __str__(self):
29         return "%s: %s" % (self.title, ' '.join(list(self.wordCount
30 .keys()))))
31
32 # fake tfidf
33 def filter_fun(wc):
34     frac = float(wc[1]) / float(100)
35     return 0.1 < frac < 0.5
36
37
38 # output for the non-stemmed data file
39 def output(f, top):
40     out = [f.title]
41     for wd in top:
42         out.append("%d" % f.wordCount.get(wd, 0))
43     return '\t'.join(out) + "\n"
44
45
46 # output for the stemmed data file
47 def output_stem(f, top):
48     out = [f.title]
49     for wd in top:
50         out.append("%d" % f.stemCount.get(wd, 0))
51     return '\t'.join(out) + "\n"
52
53
54 def generate_blogfile():
55     # read the data in as json and then transform to feeds
56     feedData = seq.json("datafiles/blogdata.json").map(lambda fe:

```

```

57 feed(fe)).to_list() # type: list[feed]
58 '''
59 get the top 500 words by word count over all words:
60     for feed <- feeds, (word,count) <- feed.wordCount.items():
61     yeild (word,count)
62     keep all wordCounts > 10
63     groupby + reduce word:(word1,c1),(word1,c2) -> (word1,c1,c2
64     ,c3,c4) -> (word1,sumC)
65     keep all wordCounts that meet fake tfidf
66     transform (word,sumC) -> word
67     take top 500
68     transform to list
69 '''
70 top500 = seq(feedData).flat_map(lambda f: list(f.wordCount.
71 items())) \
72     .filter(lambda wc: wc[1] > 10) \
73     .map(lambda wc: (wc[0], 1)) \
74     .reduce_by_key(add) \
75     .filter(filter_fun) \
76     .order_by(lambda wc: -wc[1]) \
77     .take(500) \
78     .map(lambda wc: wc[0]) \
79     .to_list()
80 # sort alphabetically
81 top500 = sorted(top500)
82 print(len(top500))
83 # write resultant to file
84 with open("datafiles/blogtop500.txt", "w+") as out:
85     out.write("Blog\t%s\n" % '\t'.join(top500))
86     for tf in sorted(feedData, key=lambda f: f.title):
87         out.write(output(tf, top500))
88
89 def generate_blogfile_stem():
90     # same as non-stem except use stemmed data
91     feedData = seq.json("datafiles/blogdata.json").map(lambda fe:
92 feed(fe, True)).to_list() # type: list[feed]
93     top500 = seq(feedData).flat_map(lambda f: list(f.stemCount.
94 items())) \
95     .filter(lambda wc: wc[1] > 1) \
96     .map(lambda wc: (wc[0],1)) \
97     .reduce_by_key(add) \
98     .filter(filter_fun) \
99     .order_by(lambda wc: -wc[1]) \
100    .take(500) \
101    .map(lambda wc: wc[0]) \
102    .to_list()
103    top500 = sorted(top500)
104    print(len(top500))
105    with open("datafiles/blogtop500-stemmed.txt", "w+") as out:
106        out.write("Blog\t%s\n" % '\t'.join(top500))
107        for tf in sorted(feedData, key=lambda f: f.title):
108            out.write(output_stem(tf, top500))
109
110 def do_non_stem():
111     # generate the blog file

```

```

108     generate_blogfile()
109     # read the data in
110     blognames, words, data = clusters.readfile('datafiles/
111     blogtop500.txt')
112     # do clustering
113     clust = clusters.hcluster(data)
114     # write out asci denogram
115     with open("datafiles/blogtop500_asciideno.txt", "w+") as out:
116         clusters.printclust2file(clust, out, labels=blognames)
117     # generate jpg version of same denogram
118     clusters.drawdendrogram(clust, blognames, jpeg='datafiles/
119     blogtop500_deno.jpg')
120     # do kmeans and log to file
121     with open("datafiles/kmeans_blogtop500.txt", "w+") as kout:
122         for k in [5, 10, 20]:
123             print("For k=%d" % k)
124             kout.write("K=%d\n" % k)
125             kout.write("Iterations\n")
126             # kmeans for value k
127             centriods = clusters.kcluster_toFile(data, k=k, out=
128             kout)
129             kout.write("Centroid Values\n")
130             # log centroid values
131             for count, centriod in enumerate(centriods, 1):
132                 print("Centroid #%d" % count)
133                 kout.write("Centroid #%d\n" % count)
134                 values = []
135                 for idx in centriod:
136                     print(blognames[idx])
137                     values.append(blognames[idx])
138                     kout.write("%s\n" % ', '.join(values))
139                 kout.write("=====\n")
140                 print("=====")
141     # do the dimensionality reduction
142     with open("datafiles/dimensionReductionNonStemmed.txt", "w+") as
143     dout:
144         scaled = clusters.scaledown_logiter(data, out=dout)
145     # generated the similar blog jpg
146     clusters.draw2d(scaled, blognames, jpeg='datafiles/
147     blogtop500_clust2d.jpg')
148
149 def do_stemmed():
150     generate_blogfile_stem()
151     blognames, words, data = clusters.readfile('datafiles/
152     blogtop500_stemmed.txt')
153     clust = clusters.hcluster(data)
154     with open("datafiles/blogtop500stemmed_asciideno.txt", "w+") as
155     out:
156         clusters.printclust2file(clust, out, labels=blognames)
157         clusters.drawdendrogram(clust, blognames, jpeg='datafiles/
158         blogtop500stemmed_deno.jpg')
159
160     with open("datafiles/kmeans_blogtop500stemmed.txt", "w+") as
161     kout:
162         for k in [5, 10, 20]:

```

```

155         print("For k=%d" % k)
156         kout.write("K=%d\n" % k)
157         kout.write("Iterations\n")
158         centriods = clusters.kcluster_toFile(data, k=k, out=
kout)
159         kout.write("Centroid Values\n—————\n")
160         for count, centroid in enumerate(centriods, 1):
161             print("Centroid #%d" % count)
162             kout.write("Centroid #%d\n" % count)
163             values = []
164             for idx in centroid:
165                 print(blognames[idx])
166                 values.append(blognames[idx])
167             kout.write("%s\n" % ', '.join(values))
168             kout.write("—————\n")
169             print("—————")
170         with open("datafiles/dimensionReductionStemmed.txt", "w+") as
dout:
171             scaled = clusters.scaledown_logiter(data, out=dout)
172             clusters.draw2d(scaled, blognames, jpg='datafiles/
blogtop500stemmed.clust2d.jpg')
173
174
175 if __name__ == "__main__":
176     do_non_stem()
177     do_stemmed()

```

Listing 3: Process Feed Data and Generate Output

Question 3

3. Cluster the blogs using K-Means, using k=5,10,20. (see slide 18). Print the values in each centroid, for each value of k. How many iterations were required for each value of k?

Answer

The full output of kmeans can be seen in listing 4 for non-stemmed and 5 for stemmed. For non-stemmed k=5 took 8 iterations, k=10 took 4, and k=20 took 4 as well. The stemmed version for k=5 took 6 iterations, k=10 took 5 and k=20 took 3.


```

1 K=5
2 Iterations
3 Iteration 0
4 Iteration 1
5 Iteration 2
6 Iteration 3
7 Iteration 4
8 Iteration 5
9 Iteration 6
10 Iteration 7
11 Iteration 8
12 Centroid Values
13
14 Centroid #1
15 "DANCING IN CIRCLES", Boggle Me Thursday, Floorshime Zipper Boots,
    Interstellar Radio Shower, Parish Radio, Samtastic! Review, THE
    HUB, Web Science and Digital Libraries Research Group,
    theindiefriend,
16 Centroid #2
17 A Wife's Tale, A2 MEDIA COURSEWORK JOINT BLOG, Becky Sharp Fashion
    Blog, But She's Not Stupid, Deadbeat, Diagnosis: No Radio,
    Flatbasset, How to be an artist and still pass for normal, If
    You Give a Girl a Camera..., KiDCHAIR, La Fotograf a Efectista
    Abstracta. Fotos Abstractas. Abstract Photos., Mile In Mine,
    My Name Is Blue Canary, One Stunning Single Egg, Our Podcast
    Could Be Your Life, Pithy Title Here, Rants from the Pants, Rod
    Shone, Room 19's Blog 2016, Sonology, Steel City Rust,
    Stonehill Sketchbook, T H E V O I D S, The Girl at the Rock
    Show, The Moon Topples, The Stearns Family, Tremagazine,
    Tremble Under Boom Lights, What Am I Doing?, from a voice
    plantation, funky little demons, isyeli's, jaaackie., sweeping
    the kitchen, the fast break of champions
18 Centroid #3
19 Cherry Area, Dust and Water Studios, Karl Drinkwater, Morgan's Blog
    , Riley Haas' blog, The Bunker, The Listening Ear, The
    Nosebleed Section, bittersweet, forget about it
20 Centroid #4
21 *Sixeyes: by Alan Williamson, ., 60@60 Sounding Booth, Blog Name
    Pending, Chantelle Swain A2 Media Studies, DaveCromwell Writes,
    Did Not Chart, Eli Jace | The Mind Is A Terrible Thing To
    Paste, Encore, Everyday Music, F-Measure, FOLK IS NOT HAPPY, I
    Hate The 90s, I/LOVE/TOTAL/DESTRUCTION, MAGGOT CAVIAR, MTJR
    RANTS & RAVES ON MUSIC, Rosie Gigg A2 Media Studies, Swinging
    Singles Club, The Jeopardy of Contentment, The Kids Are Coming
    Up From Behind, The Night Mail, The Power of Independent
    Trucking, The World's First Internet Baby, mouxlaloulouda,
    turnitup!
22 Centroid #5
23 A H T A P O T, Desolation Row Records, FlowRadio Playlists (and
    Blog), Green Eggs and Ham Mondays 8-10am, Haris Sfakianakis
    Photography, INDIEehren.!, IoTube :), KISTE F.M., Lost in
    the Shuffle, MARISOL, MR. BEAUTIFUL TRASH ART, MarkeOrtega's
    Journalism Portfolio, SEM REGRAS, Spinitron Blog, Stories From
    the City, Stories From the Sea, The Campus Buzz on WSOU,
    adrianoblog, k nstler treu, this time tomorrow,
24

```

```

25 K=10
26 Iterations
27 Iteration 0
28 Iteration 1
29 Iteration 2
30 Iteration 3
31 Iteration 4
32 Centroid Values
33
34 Centroid #1
35 Riley Haas' blog
36 Centroid #2
37 A Wife's Tale, But She's Not Stupid, Cherry Area, Deadbeat,
    Diagnosis: No Radio, Karl Drinkwater, MarkEOrtega's Journalism
    Portfolio, Mile In Mine, Morgan's Blog, One Stunning Single Egg
    , Pithy Title Here, Rants from the Pants, Sonology, Spinitron
    Blog, Steel City Rust, Stonehill Sketchbook, The Moon Topples,
    The Nosebleed Section, The Stearns Family, What Am I Doing?,
    bittersweet, forget about it, funky little demons, isyeli's,
    jaaackie., the fast break of champions
38 Centroid #3
39 KiDCHAIR, My Name Is Blue Canary, The Girl at the Rock Show,
    Tremmagazine, Tremble Under Boom Lights, from a voice plantation
40 Centroid #4
41 "DANCING IN CIRCLES", Boggle Me Thursday, Interstellar Radio Shower
    , Samtastic! Review, THE HUB, sweeping the kitchen,
    theindiefriend,
42 Centroid #5
43 *Sixeyes: by Alan Williamson, ., Blog Name Pending, DaveCromwell
    Writes, Did Not Chart, Eli Jace | The Mind Is A Terrible Thing
    To Paste, Encore, Everyday Music, F-Measure, FOLK IS NOT HAPPY,
    I Hate The 90s, I/LOVE/TOTAL/DESTRUCTION, MAGGOT CAVIAR, MTJR
    RANTS & RAVES ON MUSIC, Parish Radio, Swinging Singles Club,
    The Jeopardy of Contentment, The Kids Are Coming Up From Behind
    , The Power of Independent Trucking, The World's First Internet
    Baby, mouxlaloulouda, turnitup!
44 Centroid #6
45 A2 MEDIA COURSEWORK JOINT BLOG, Chantelle Swain A2 Media Studies,
    Floorshime Zipper Boots, Rosie Gigg A2 Media Studies, T H E V O
    I D S, Web Science and Digital Libraries Research Group
46 Centroid #7
47 Becky Sharp Fashion Blog, Dust and Water Studios, The Bunker, The
    Listening Ear, The Night Mail
48 Centroid #8
49 60@60 Sounding Booth, Desolation Row Records, Flatbasset,
    INDIEehren.!, La Fotograf a Efectista Abstracta. Fotos
    Abstractas. Abstract Photos., Our Podcast Could Be Your Life
50 Centroid #9
51 A H T A P O T, FlowRadio Playlists (and Blog), Green Eggs and Ham
    Mondays 8-10am, Haris Sfakianakis Photography, IoTube :),
    KISTE F.M., Lost in the Shuffle, MARISOL, MR. BEAUTIFUL TRASH
    ART, SEM REGRAS, Stories From the City, Stories From the Sea,
    The Campus Buzz on WSOU, adrianoblog, k nstler treu, this time
    tomorrow,
52 Centroid #10
53 How to be an artist and still pass for normal, If You Give a Girl a
    Camera..., Rod Shone, Room 19's Blog 2016

```

```

54 |-----
55 | K=20
56 | Iterations
57 | Iteration 0
58 | Iteration 1
59 | Iteration 2
60 | Iteration 3
61 | Iteration 4
62 | Centroid Values
63 |-----
64 | Centroid #1
65 | *Sixeyes: by Alan Williamson, ., Blog Name Pending, DaveCromwell
    | Writes, Did Not Chart, Eli Jace | The Mind Is A Terrible Thing
    | To Paste, Encore, F-Measure, FOLK IS NOT HAPPY, I Hate The 90s,
    | I/LOVE/TOTAL/DESTRUCTION, MAGGOT CAVIAR, MTJR RANTS & RAVES ON
    | MUSIC, Parish Radio, Swinging Singles Club, The Jeopardy of
    | Contentment, The Kids Are Coming Up From Behind, The Power of
    | Independent Trucking, The World's First Internet Baby,
    | mouxlaloulouda
66 | Centroid #2
67 | 60@60 Sounding Booth, Desolation Row Records, FlowRadio Playlists (
    | and Blog), Green Eggs and Ham Mondays 8-10am, INDIEehren.!,
    | IoTube :), Lost in the Shuffle, One Stunning Single Egg,
    | Our Podcast Could Be Your Life, Stories From the City, Stories
    | From the Sea, The Campus Buzz on WSOU, k nstler treu, this
    | time tomorrow
68 | Centroid #3
69 |
70 | Centroid #4
71 | La Fotograf a Efectista Abstracta. Fotos Abstractas. Abstract
    | Photos.
72 | Centroid #5
73 | Dust and Water Studios, Riley Haas' blog, The Nosebleed Section
74 | Centroid #6
75 | The Night Mail, turnitup!
76 | Centroid #7
77 | Becky Sharp Fashion Blog, Everyday Music, from a voice plantation
78 | Centroid #8
79 | "DANCING IN CIRCLES", THE HUB
80 | Centroid #9
81 | A2 MEDIA COURSEWORK JOINT BLOG, Rosie Gigg A2 Media Studies, T H E
    | V O I D S
82 | Centroid #10
83 | Chantelle Swain A2 Media Studies, sweeping the kitchen,
84 | Centroid #11
85 | A H T A P O T, KISTE F.M., MARISOL, MR. BEAUTIFUL TRASH ART
86 | Centroid #12
87 | Haris Sfakianakis Photography,
88 | Centroid #13
89 | Cherry Area, Deadbeat, Diagnosis: No Radio, Flatbasset, KiDCHAIR,
    | My Name Is Blue Canary, Pithy Title Here, Sonology, Steel City
    | Rust, The Girl at the Rock Show, The Stearns Family, Tremble
    | Under Boom Lights, bittersweet, funky little demons, isyeli's,
    | the fast break of champions
90 | Centroid #14

```

```

91 Floorshine Zipper Boots, Tremagazine
92 Centroid #15
93 Boggle Me Thursday, Interstellar Radio Shower, Samtastic! Review,
    Web Science and Digital Libraries Research Group,
    theindiefriend
94 Centroid #16
95 MarkEOrtega's Journalism Portfolio, Spinitron Blog
96 Centroid #17
97 SEM REGRAS, adrianoblog
98 Centroid #18
99
100 Centroid #19
101 A Wife's Tale, But She's Not Stupid, How to be an artist and still
    pass for normal, Karl Drinkwater, Mile In Mine, Morgan's Blog,
    Rants from the Pants, Room 19's Blog 2016, Stonehill Sketchbook
    , The Bunker, The Listening Ear, The Moon Topples, What Am I
    Doing?, forget about it, jaaackie.
102 Centroid #20
103 If You Give a Girl a Camera..., Rod Shone
104

```

Listing 4: Kmeans output non-stemmed

```

1 K=5
2 Iterations
3 Iteration 0
4 Iteration 1
5 Iteration 2
6 Iteration 3
7 Iteration 4
8 Iteration 5
9 Iteration 6
10 Centroid Values
11
12 Centroid #1
13 ., DaveCromwell Writes, Encore, FOLK IS NOT HAPPY, Floorshine
    Zipper Boots, I Hate The 90s, I/LOVE/TOTAL/DESTRUCTION, MTJR
    RANTS & RAVES ON MUSIC, Mile In Mine, Rosie Gigg A2 Media
    Studies, Samtastic! Review, Sonology, The Kids Are Coming Up
    From Behind, The Night Mail, from a voice plantation,
    mouxlaloulouda
14 Centroid #2
15 A Wife's Tale, Boggle Me Thursday, Dust and Water Studios, Eli Jace
    | The Mind Is A Terrible Thing To Paste, Everyday Music, How
    to be an artist and still pass for normal, MAGGOT CAVIAR,
    Morgan's Blog, Our Podcast Could Be Your Life, Rants from the
    Pants, Rod Shone, Room 19's Blog 2016, Swinging Singles Club,
    The Bunker, The Stearns Family, isyeli's, jaaackie.
16 Centroid #3
17 A2 MEDIA COURSEWORK JOINT BLOG, Becky Sharp Fashion Blog, But She's
    Not Stupid, Chantelle Swain A2 Media Studies, Cherry Area,
    Deadbeat, Diagnosis: No Radio, F-Measure, Flatbasset, Karl
    Drinkwater, MarkEOrtega's Journalism Portfolio, My Name Is Blue
    Canary, Pithy Title Here, Riley Haas' blog, Spinitron Blog,
    Steel City Rust, T H E V O I D S, The Listening Ear, The Moon
    Topples, The Nosebleed Section, Tremble Under Boom Lights, Web
    Science and Digital Libraries Research Group, What Am I Doing?,

```

```

    bittersweet
18 Centroid #4
19 "DANCING IN CIRCLES", *Sixeyes: by Alan Williamson, 60@60 Sounding
    Booth, A H T A P O T, Desolation Row Records, FlowRadio
    Playlists (and Blog), Green Eggs and Ham Mondays 8–10am, Haris
    Sfakianakis Photography, INDIEehren.!, If You Give a Girl a
    Camera..., IoTube :), KISTE F.M., La Fotograf a Efectista
    Abstracta. Fotos Abstractas. Abstract Photos., Lost in the
    Shuffle, MARISOL, MR. BEAUTIFUL TRASH ART, One Stunning Single
    Egg, Parish Radio, SEM REGRAS, Stonehill Sketchbook, Stories
    From the City, Stories From the Sea, THE HUB, The Campus Buzz
    on WSOU, adrianoblog, forget about it, funky little demons,
    k nstler treu, sweeping the kitchen, theindiefriend, this time
    tomorrow, turnitup!,
    ,
20 Centroid #5
21 Blog Name Pending, Did Not Chart, Interstellar Radio Shower,
    KiDCHAIR, The Girl at the Rock Show, The Jeopardy of
    Contentment, The Power of Independent Trucking, The World's
    First Internet Baby, Tremagazine, the fast break of champions


---


22 K=10
23 Iterations
24 Iteration 0
25 Iteration 1
26 Iteration 2
27 Iteration 3
28 Iteration 4
29 Iteration 5
30 Centroid Values


---


31 Centroid #1
32 A Wife's Tale, Becky Sharp Fashion Blog, MTJR RANTS & RAVES ON
33 MUSIC, The Jeopardy of Contentment, The Moon Topples, The Power
34 of Independent Trucking, Tremble Under Boom Lights, Web
    Science and Digital Libraries Research Group, from a voice
    plantation
35 Centroid #2
36 Dust and Water Studios, KiDCHAIR, La Fotograf a Efectista
    Abstracta. Fotos Abstractas. Abstract Photos., Parish Radio,
    Rod Shone, THE HUB, The Kids Are Coming Up From Behind, forget
    about it, turnitup!
37 Centroid #3
38 ., Cherry Area, Flatbasset, Pithy Title Here, Riley Haas' blog, The
    Listening Ear, The Night Mail, The Nosebleed Section, The
    World's First Internet Baby, What Am I Doing?, mouxlaloulouda
39 Centroid #4
40 A2 MEDIA COURSEWORK JOINT BLOG, Deadbeat, F-Measure, MAGGOT CAVIAR,
    Mile In Mine, Sonology
41 Centroid #5
42 Chantelle Swain A2 Media Studies, Did Not Chart, Encore, I/LOVE/
    TOTAL/DESTRUCTION, MarkEOrtega's Journalism Portfolio, Rosie
    Gigg A2 Media Studies, T H E V O I D S
43 Centroid #6
44 But She's Not Stupid, DaveCromwell Writes, Diagnosis: No Radio, I
    Hate The 90s, Karl Drinkwater, My Name Is Blue Canary,
    bittersweet

```

```

45 Centroid #7
46 60@60 Sounding Booth, FOLK IS NOT HAPPY, Floorshime Zipper Boots,
    Interstellar Radio Shower, Samtastic! Review, Steel City Rust,
    The Campus Buzz on WSOU, Tremagazine, funky little demons
47 Centroid #8
48 Blog Name Pending, Morgan's Blog, Our Podcast Could Be Your Life,
    Spinitron Blog
49 Centroid #9
50 "DANCING IN CIRCLES", A H T A P O T, Desolation Row Records,
    Everyday Music, FlowRadio Playlists (and Blog), Green Eggs and
    Ham Mondays 8-10am, Haris Sfakianakis Photography, INDIEehren
    .!, If You Give a Girl a Camera..., IoTube :), KISTE F.M.,
    Lost in the Shuffle, MARISOL, MR. BEAUTIFUL TRASH ART, SEM
    REGRAS, Stonehill Sketchbook, Stories From the City, Stories
    From the Sea, adrianoblog, k nstler treu, theindiefriend,
    ,
51 Centroid #10
52 *Sixeyes: by Alan Williamson, Boggle Me Thursday, Eli Jace | The
    Mind Is A Terrible Thing To Paste, How to be an artist and
    still pass for normal, One Stunning Single Egg, Rants from the
    Pants, Room 19's Blog 2016, Swinging Singles Club, The Bunker,
    The Girl at the Rock Show, The Stearns Family, isyeli's,
    jaaackie., sweeping the kitchen, the fast break of champions,
    this time tomorrow
53 =====
54 K=20
55 Iterations
56 Iteration 0
57 Iteration 1
58 Iteration 2
59 Iteration 3
60 Centroid Values
61 -----
62 Centroid #1
63 Flatbasset, I Hate The 90s, My Name Is Blue Canary, Samtastic!
    Review, The Jeopardy of Contentment, Tremagazine
64 Centroid #2
65 *Sixeyes: by Alan Williamson, Did Not Chart, Everyday Music, MAGGOT
    CAVIAR, The Campus Buzz on WSOU, funky little demons, turnitup
    !
66 Centroid #3
67 How to be an artist and still pass for normal, Rod Shone, Sonology,
    Stonehill Sketchbook, THE HUB
68 Centroid #4
69 F-Measure, Rants from the Pants
70 Centroid #5
71 I/LOVE/TOTAL/DESTRUCTION, MTJR RANTS & RAVES ON MUSIC, Swinging
    Singles Club, bittersweet, forget about it
72 Centroid #6
73 Pithy Title Here, The Moon Topples, The Nosebleed Section, The
    Stearns Family, Tremble Under Boom Lights, Web Science and
    Digital Libraries Research Group, from a voice plantation
74 Centroid #7
75 Room 19's Blog 2016, The Kids Are Coming Up From Behind, What Am I
    Doing?
76 Centroid #8

```

```

77 60@60 Sounding Booth, A H T A P O T, Desolation Row Records, Haris
    Sfakianakis Photography, IoTube :), KISTE F.M., La
    Fotograf a Efectista Abstracta. Fotos Abstractas. Abstract
    Photos., Lost in the Shuffle, adrianoblog, theindiefriend, this
    time tomorrow,
78 Centroid #9
79 FlowRadio Playlists (and Blog), Green Eggs and Ham Mondays 8–10am,
    MR. BEAUTIFUL TRASH ART, SEM REGRAS, Steel City Rust
80 Centroid #10
81 Becky Sharp Fashion Blog, But She's Not Stupid, The Bunker, The
    Girl at the Rock Show, The World's First Internet Baby
82 Centroid #11
83 If You Give a Girl a Camera..., Morgan's Blog, Stories From the
    City, Stories From the Sea, isyeli's, the fast break of
    champions
84 Centroid #12
85 A Wife's Tale, Blog Name Pending, Eli Jace | The Mind Is A Terrible
    Thing To Paste, Interstellar Radio Shower, KiDCHAIR, Our
    Podcast Could Be Your Life, jaaackie., mouxlaloulouda
86 Centroid #13
87 Deadbeat, Mile In Mine, One Stunning Single Egg, sweeping the
    kitchen
88 Centroid #14
89 ., Cherry Area, Diagnosis: No Radio, The Listening Ear, The Night
    Mail
90 Centroid #15
91 Parish Radio,
92 Centroid #16
93 A2 MEDIA COURSEWORK JOINT BLOG, Chantelle Swain A2 Media Studies,
    Rosie Gigg A2 Media Studies, T H E V O I D S
94 Centroid #17
95 DaveCromwell Writes, Dust and Water Studios, FOLK IS NOT HAPPY,
    Floorshime Zipper Boots
96 Centroid #18
97 Boggle Me Thursday
98 Centroid #19
99 "DANCING IN CIRCLES", INDIEehren.!, MARISOL, Spinitron Blog,
    k nstler treu
100 Centroid #20
101 Encore, Karl Drinkwater, MarkEOrtega's Journalism Portfolio, Riley
    Haas' blog, The Power of Independent Trucking
102

```

Listing 5: Kmeans output stemmed

Question 4

4. Use MDS to create a JPEG of the blogs similar to slide 29. How many iterations were required?

Answer

The number of iterations required for dimension reduction for both the non-stemmed and stemmed versions can be seen below in listings 6 and 7.

```
1 Dimension Reduction took 361 iterations
```

Listing 6: Dimension Reduction Iterations Non-Stemmed

```
1 Dimension Reduction took 75 iterations
```

Listing 7: Dimension Reduction Iterations Stemmed

The difference in the number of iterations for the dimension reduction between the non-stemmed and stemmed blog term matrices in my opinion is due to the sparse nature of the terms. The raw terms are very different as the blogs gotten are all over the place in terms of content thus it will take longer. Whereas the stemmed version the context of the words is removed and a better approximation can be made.

The resultant jpg pictures can be seen in figure 3 for non-stemmed and 4 for stemmed. Once again these pictures are way to large to clearly see in this report. For clearer inspect the files are found in datafiles/blogtop500_clust2d.jpg and datafiles/blogtop500stemmed_clust2d.jpg.

Of note in blogtop500_clust2d.jpg, the WSDL blog is grouped with Stories From the City, Stories from the Sea and F-Measure is grouped with The Kids are Coming Up From Behind.

Also of note in blogtop500stemmed_clust2d.jpg, the WSDL blog is grouped with The Night Mail and F-Measure is grouped with MarkEOrtega's Journalism Portfolio.

Question 5

5. Re-run question 2, but this time with proper TFIDF calculations instead of the hack discussed on slide 7 (p. 32). Use the same 500 words, but this time replace their frequency count with TFIDF scores as computed in assignment #3. Document the code, techniques, methods, etc. used to generate these TFIDF values. Upload the new data file to github.

Compare and contrast the resulting dendrogram with the dendrogram from question #2.

Note: ideally you would not reuse the same 500 terms and instead come up with TFIDF scores for all the terms and then choose the top 500 from that list, but I'm trying to limit the amount of work necessary.

w many iterations were required?

=====
=====The questions below is for 5 points extra credit=====

Answer

Not attempted.

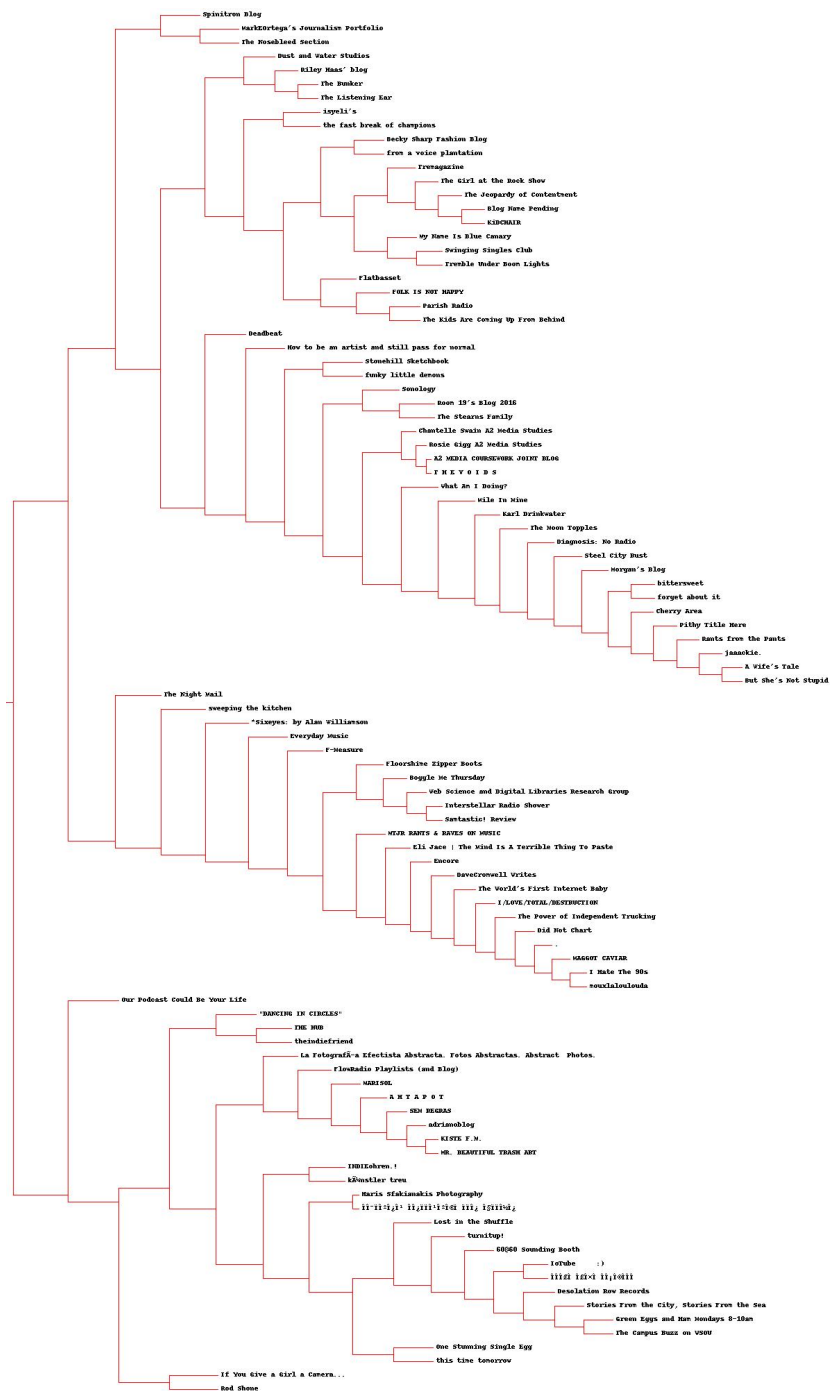


Figure 1: Dendrogram of the blogs using top 500 terms





Figure 3: Dimension Reduction Non Stemmed

