

# INTRODUCTION TO WEB SCIENCES: Assignment 8

Manoj Chandra Kompalli

7 April 2016

# Contents

<b>1</b>	<b>Question 1:</b>	<b>2</b>
1.1	Approach . . . . .	2
1.2	Code Listing . . . . .	3
1.2.1	getBlogs.py to generate 100 URIs . . . . .	3
1.2.2	countPages.py to find the number of pages in each blog . . . . .	4
1.2.3	matrixcreate.py to create a blog term matrix . . . . .	5
1.3	Blog term matrix shown in Excel . . . . .	7
<b>2</b>	<b>Question 2:</b>	<b>8</b>
2.1	Approach . . . . .	8
2.2	Code Listing . . . . .	8
2.2.1	dgram.py . . . . .	8
2.3	Output . . . . .	9
2.3.1	JPEG Dendogram . . . . .	9
2.3.2	ASCII dendogram . . . . .	10
<b>3</b>	<b>Question 3:</b>	<b>14</b>
3.1	Approach . . . . .	14
3.2	Code Listing . . . . .	14
3.2.1	convert.py . . . . .	14
3.3	Output . . . . .	15
3.3.1	Centroid Values and Iterations . . . . .	15
<b>4</b>	<b>Question 4:</b>	<b>16</b>
4.1	Approach . . . . .	16
4.2	Code Listing . . . . .	16
4.3	Output . . . . .	17
4.3.1	JPEG of blogs using MPS . . . . .	17
<b>5</b>	<b>Question 5:</b>	<b>18</b>
5.1	Approach . . . . .	18
5.2	Code Listing . . . . .	18
5.3	Output . . . . .	22
5.3.1	Output Blog term matrix with tfidf values . . . . .	22
5.3.2	Output JPEG . . . . .	23

# 1 Question 1:

## 1.1 Approach

- I have extracted the 98 URIs initially and appended the required string
- I made sure there are no duplicates
- I have added the given 2 URIs in the question to make it 100
- I used feed parser to parse all the randomly generated URIs
- To get the most popular words I kept track of occurrence of each word in the blog
- I sorted all the words based on how many times they occurred and selected the top 500 words.
- Then I calculated the occurrence of every word in the respective blog
- The resulting matrix contains 501 columns and 101 rows

## 1.2 Code Listing

### 1.2.1 getBlogs.py to generate 100 URIs

```
1 import os
2 import requests
3 import feedparser
4 import re
5 import sys
6
7 def getBlogs():
8     # link='http://www.blogger.com/next-blog?navBar=true&blogID=3471633091411211117'
9     # cmd="curl -I -L "+" "+link+"""
10    link='http://www.blogger.com/next-blog?navBar=true&blogID=3471633091411211117'
11
12
13    durls=set()
14
15    while len(durls)<98:
16        req=requests.get(link)
17        durls.add(req.url)
18        # link=req.url
19    # print durls
20    # file1=open('urls','w')
21    file2=open('urlsappended','w')
22    # for item in durls:
23        # #print item + '\n'
24        # file1.write(item + '\n')
25    file2.write('http://f-measure.blogspot.com'+'/feeds/posts/default?alt=rss' +'\n')
26    file2.write('http://ws-dl.blogspot.com/'+'/feeds/posts/default?alt=rss' +'\n')
27
28
29    for item in durls:
30        #print item + '\n'
31        file2.write(item.strip('/?expref=next-blog')+'/feeds/posts/default?alt=rss' +'\n')
32
33
34
35
36
37
38 if __name__== '__main__':
39     getBlogs()
```

### 1.2.2 countPages.py to find the number of pages in each blog

```
1 #!/usr/bin/env python
2 import os
3 import sys
4 import urllib
5 import time
6 import feedparser
7
8 from bs4 import BeautifulSoup
9
10 def checkNextPage(url):
11
12     f = urllib.urlopen(url)
13     soup = BeautifulSoup(f.read(), from_encoding=f.info().getparam('charset'))
14
15     try:
16         link = soup.find('link', rel='next', href = True)['href']
17     except TypeError:
18         link = None
19     return link
20
21
22 def main():
23     feedlist = open('urlsappended').readlines()
24
25     for url in feedlist:
26         try:
27             d = feedparser.parse(url)
28
29             title = d['feed']['title']
30
31             count = 1
32             nextLink = checkNextPage( url )
33
34             while nextLink:
35                 nextLink = checkNextPage( nextLink )
36                 count += 1
37
38             print u'|'.join((str(count), title)).encode('utf-8').strip())
39
40
41         except KeyError:
42             pass
43
44
45 if __name__ == '__main__':
46     main()
```

### 1.2.3 matrixcreate.py to create a blog term matrix

```
1 #!/usr/bin/env python
2
3 import re
4 import sys
5 import feedparser
6 if __name__ == "__main__":
7
8     matrixCreation()
9 def matrixCreation():
10     nowords={}
11     cap={}
12
13     feedlist=[line for line in file('urlsappended')]
14     for feedurl in feedlist:
15         try:
16             title ,wc=getnowords(feedurl)
17             nowords[title]=wc
18             for word,count in wc.items():
19                 cap.setdefault(word,0)
20                 if count>1:
21                     cap[word]+=1
22         except:
23             print 'Parsing failed %s' % feedurl
24
25     wordlist=[]
26     repeatingWords = []
27     for w,bc in cap.items():
28         frac=float(bc)/len(feedlist)
29         if frac>0.1 and frac<0.5:
30             repeatingWords.append((w,bc))
31
32     repeatingWords=sorted(repeatingWords,key=lambda x:x[1], reverse = True)
33
34     for value in repeatingWords:
35
36         value1 = value[0]
37
38         value2 = value[1]
39         length = len(wordlist)
40         if(length < 500):
41             wordlist.append(value1)
42         else:
43             break
44
45     out=file('blogdata.txt','w')
46     out.write('Blog')
47
48     for word in wordlist:
49         word1 = word.encode('UTF-8')
50         out.write('\t%s' % word1)
51         out.write('\n')
52
53     for blog,wc in nowords.items():
54         blogName = blog.encode('UTF-8')
55         print blog
```

```

56 out.write(blogName)
57 for word in wordlist:
58     if word in wc: out.write('\t%d' % wc[word])
59     else: out.write('\t0')
60 out.write('\n')
61 def getwords(html):
62
63     txt=re.compile(r'<[^>]+>').sub('',html)
64
65
66     words=re.compile(r'^A-Za-z+').split(txt)
67
68     return [word.lower() for word in words if word!='']
69 def getnowords(url):
70
71     d=feedparser.parse(url)
72     wc={}
73
74
75     for e in d.entries:
76         if 'summary' in e: summary=e.summary
77         else: summary=e.description
78
79
80     words=getwords(e.title+' '+summary)
81     for word in words:
82         wc.setdefault(word,0)
83         wc[word]+=1
84     return d.feed.title,wc

```

### 1.3 Blog term matrix shown in Excel

	A	B	C	D	E	F	G	H	I	J	K	L	
1	Blog	yet	week	real	heart	both	him	times	once	such	did	without	doe
2	Flatbasset	14	38	7	3	13	22	1	3	1	7	2	
3	Riley Haas' blog	6	0	1	2	9	5	9	7	3	12	4	
4	(Insert World Problem Here) Sucks.	0	0	0	0	0	0	0	0	1	2	1	
5	SEM REGRAS	0	0	0	7	0	0	0	0	0	0	0	
6	Pithy Title Here	12	8	26	2	12	10	13	15	12	32	15	
7	Morgan's Blog	3	0	2	1	5	6	0	1	4	4	1	
8	MARISOL	0	0	0	0	0	0	0	0	0	0	0	
9	THE HUB	2	4	0	0	0	0	1	0	2	1	0	
10	Brian's Music Blog!!!	0	11	6	1	2	3	6	0	3	5	2	
11	Web Science and Digital Libraries Research Group	6	3	4	0	14	4	8	6	61	15	9	
12	Steel City Rust	3	2	7	1	3	7	2	3	10	4	8	
13	MR. BEAUTIFUL TRASH ART	0	0	0	0	0	0	0	0	0	0	0	
14	60@60 Sounding Booth	3	0	4	8	0	1	1	2	0	0	1	
15	ORGANMYTH	0	0	0	0	0	0	0	0	0	0	0	
16	MarkEOrtega's Journalism Portfolio	0	5	0	0	2	0	0	2	0	0	0	
17	Green Eggs and Ham Mondays 8-10am	0	0	2	4	0	0	1	0	0	0	1	
18	GLI Press	0	19	0	2	0	1	0	0	0	0	0	
19	turnitup!	2	1	6	6	3	1	6	4	4	2	5	
20	Stories From the City, Stories From the Sea	0	2	2	5	2	1	2	3	0	0	2	
21	A H T A P O T	0	0	1	0	0	0	0	0	0	0	0	
22	Floorshime Zipper Boots	3	0	1	0	2	0	4	0	2	0	0	
23	Did Not Chart	1	8	1	3	3	5	3	6	1	6	6	
24	How to be an artist and still pass for normal	0	0	0	0	0	0	0	0	1	0	0	
25	Party Full of Strangers	0	5	4	1	2	2	2	1	0	1	0	

Figure 1: JSON array containing the follower data



## 2 Question 2:

### 2.1 Approach

- I have used the clusters.py program mentioned in the presentation to produce the required dendrogram
- Printclust prints the dendrogram
- drawdendrogram saves it in JPEG format
- There are numerous blogs similar to F-measure. The top one is the power of independent study
- The blog similar to web science and Digital libraries Research group is Ideal copy

### 2.2 Code Listing

#### 2.2.1 dgram.py

```
1 import clusters
2 import sys
3
4 blognames, words, data=clusters.readfile('blogmatrix.txt')
5 clust = clusters.hcluster(data)
6
7 # print ASCII dendrogram
8 clusters.printclust(clust, labels=blognames)
9 sys.stdout = open('ascii.txt', 'w')
10
11 # save JPEG dendrogram
12 clusters.drawdendrogram(clust, blognames, jpeg='dendogram.jpg')
```

## 2.3 Output

### 2.3.1 JPEG Dendrogram

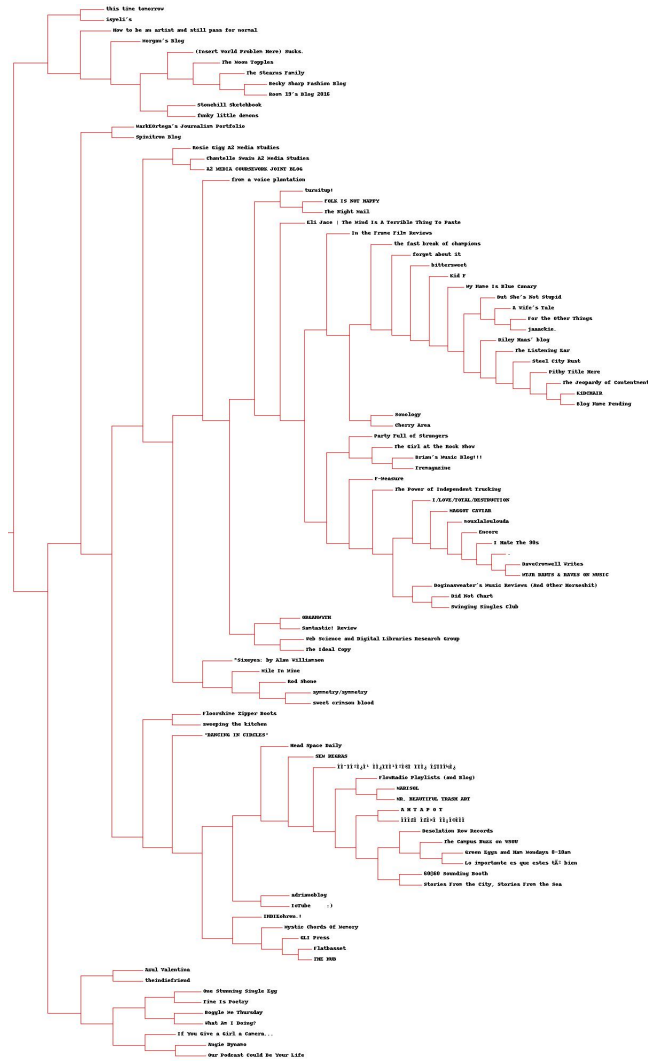


Figure 2: Initial graph split into 4 groups

### 2.3.2 ASCII dendogram

```
1 -
2 -
3 -
4   this time tomorrow
5   isyeli 's
6 -
7   How to be an artist and still pass for normal
8 -
9   Morgan's Blog
10 -
11 -
12     (Insert World Problem Here) Sucks.
13 -
14     The Moon Topples
15 -
16     The Stearns Family
17 -
18     Becky Sharp Fashion Blog
19     Room 19's Blog 2016
20 -
21     Stonehill Sketchbook
22     funky little demons
23 -
24 -
25 -
26   MarkeOrtega's Journalism Portfolio
27   Spinitron Blog
28 -
29 -
30 -
31     Rosie Gigg A2 Media Studies
32 -
33     Chantelle Swain A2 Media Studies
34     A2 MEDIA COURSEWORK JOINT BLOG
35 -
36 -
37     from a voice plantation
38 -
39 -
40 -
41     turnitup!
42 -
43     FOLK IS NOT HAPPY
44     The Night Mail
45 -
46     Eli Jace | The Mind Is A Terrible Thing To Paste
47 -
48 -
49     In the Frame Film Reviews
50 -
51 -
52     the fast break of champions
53 -
54     forget about it
55 -
```

56                   bittersweet  
 57                   —  
 58                   Kid F  
 59                   —  
 60                   My Name Is Blue Canary  
 61                   —  
 62                   —  
 63                   But She's Not Stupid  
 64                   —  
 65                   A Wife's Tale  
 66                   —  
 67                   For the Other Things  
 68                   jaaackie.  
 69                   —  
 70                   Riley Haas' blog  
 71                   —  
 72                   The Listening Ear  
 73                   —  
 74                   Steel City Rust  
 75                   —  
 76                   Pithy Title Here  
 77                   —  
 78                   The Jeopardy of Contentment  
 79                   —  
 80                   KiDCHAIR  
 81                   Blog Name Pending  
 82                   —  
 83                   Sonology  
 84                   Cherry Area  
 85                   —  
 86                   —  
 87                   Party Full of Strangers  
 88                   —  
 89                   The Girl at the Rock Show  
 90                   —  
 91                   Brian's Music Blog!!!  
 92                   Tremagazine  
 93                   —  
 94                   F-Measure  
 95                   —  
 96                   The Power of Independent Trucking  
 97                   —  
 98                   —  
 99                   I/LOVE/TOTAL/DESTRUCTION  
 100                   —  
 101                   MAGGOT CAVIAR  
 102                   —  
 103                   mouxlaloulouda  
 104                   —  
 105                   Encore  
 106                   —  
 107                   I Hate The 90s  
 108                   —  
 109                   .  
 110                   —  
 111                   DaveCromwell Writes  
 112                   MTJR RANTS & RAVES ON MUSIC

113 —  
 114 Doginasweater's Music Reviews (And Other Horseshit)  
 115 —  
 116 Did Not Chart  
 117 Swinging Singles Club  
 118 —  
 119 —  
 120 ORGANMYTH  
 121 Samtastic! Review  
 122 —  
 123 Web Science and Digital Libraries Research Group  
 124 The Ideal Copy  
 125 —  
 126 \*Sixeyes: by Alan Williamson  
 127 —  
 128 Mile In Mine  
 129 —  
 130 Rod Shone  
 131 —  
 132 symmetry/symmetry  
 133 sweet crimson blood  
 134 —  
 135 —  
 136 Floorshime Zipper Boots  
 137 sweeping the kitchen  
 138 —  
 139 "DANCING IN CIRCLES"  
 140 —  
 141 —  
 142 —  
 143 Head Space Daily  
 144 —  
 145 SEM REGRAS  
 146 —  
 147 —  
 148 —  
 149 —  
 150 —  
 151 —  
 152 FlowRadio Playlists (and Blog)  
 153 —  
 154 MARISOL  
 155 MR. BEAUTIFUL TRASH ART  
 156 —  
 157 —  
 158 A H T A P O T  
 159 —  
 160 —  
 161 —  
 162 Desolation Row Records  
 163 —  
 164 The Campus Buzz on WSOU  
 165 —  
 166 Green Eggs and Ham Mondays 8–10am  
 167 Lo importante es que estes t bien  
 168 —  
 169 60@60 Sounding Booth

```

170                                     Stories From the City, Stories From the Sea
171
172         -
173           adrianoblog
174           IoTube      :)
175
176         -
177           INDIEehren.!!
178
179         -
180           Mystic Chords Of Memory
181
182         -
183           GLI Press
184
185         -
186           Flatbasset
187           THE HUB
188
189         -
190           Azul Valentina
191           theindiefriend
192
193         -
194           One Stunning Single Egg
195           Time Is Poetry
196
197         -
198           Boggle Me Thursday
199           What Am I Doing?
200
201         -
202           If You Give a Girl a Camera...
203
204         -
205           Angie Dynamo
206           Our Podcast Could Be Your Life

```

## 3 Question 3:

### 3.1 Approach

- Below code performs blog clustering based on clusters.py program using the method kcluster
- The program prints the iterations involved for each values of k (5,10,20).
- The program also print the values in each centroid for respective k values

### 3.2 Code Listing

#### 3.2.1 convert.py

```
1 #!/usr/local/bin/python
2
3 import clusters
4
5 blognames, words, data=clusters.readfile('blogmatrix.txt')
6 print 'Iterations for varying k value '
7 print "For k=5"
8 kclust=clusters.kcluster(data, k=5)
9 print 'Centroid values for first part are:'
10 for i in range(0,5):
11     for item in kclust[i]:
12         print blognames[item]
13 print
14
15 print "For k=10"
16 kclust=clusters.kcluster(data, k=10)
17 print 'Centroid values for second part are:'
18 for i in range(0,10):
19     for item in kclust[i]:
20         print blognames[item]
21 print
22
23 print "For k=20"
24 kclust=clusters.kcluster(data, k=20)
25 print 'Centroid values for third part are:'
26 for i in range(0,20):
27     for item in kclust[i]:
28         print blognames[item]
29 print
```

## 3.3 Output

### 3.3.1 Centroid Values and Iterations

```
For k=20
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
Centroid values for third part are:
The Night Mail
MR. BEAUTIFUL TRASH ART
ORGANMYTH
Samtastic! Review
from a voice plantation
Floorshine Zipper Boots
Did Not Chart
The Power of Independent Trucking
DaveCromwell Writes
FOLK IS NOT HAPPY
MAGGOT CAVIAR
MTJR RANTS & RAVES ON MUSIC
Doginasweater's Music Reviews (And Other Horseshit)
.
Blog Name Pending
Swinging Singles Club
mouxlaloulouda
I/LOVE/TOTAL/DESTRUCTION
Encore
I Hate The 90s
One Stunning Single Egg
this time tomorrow
Azul Valentina
Riley Haas' blog
(Insert World Problem Here) Sucks.
Pithy Title Here
Morgan's Blog
```

---

Figure 3: Centroid Values and Iterations



## 4 Question 4:

### 4.1 Approach

- Blog space is generated by the multi dimensional scaling using the code mds.py.
- Again we make use of clusters.py for the scaledown method.
- The output shows a graph which is circular unlike the previous question

### 4.2 Code Listing

```
1 #!/usr/local/bin/python
2
3 import clusters
4
5 blognames, words, data=clusters.readfile('blogmatrix.txt')
6
7 coords = clusters.scaledown(data)
8
9 clusters.draw2d(coords, blognames, jpeg='draw2d.jpg')
```

### 4.3.1 JPEG of blogs using MPS



Figure 4: JPEG of blogs using MPS

## 5 Question 5:

### 5.1 Approach

- I wrote almost the same code except for the few lines which calculate the tfidf values
- Here I have ordered the matrix based on tfidf values of the respective words.
- I have extracted the words with high tfidf values
- I generated dendrogram again just like how I did before.
- From the current dendrogram F-Measure is similar to the blog Desolation Row Records
- WebScience and Digital Libraries Research Search group is similar to Maggot Caviar and blog I love total destruction
- We can observe that the hierarchy of the blogs has changed. Somehow, F-Measure now has a lot less similar blogs

### 5.2 Code Listing

```
1 import feedparser
2 import collections
3 import re
4 import operator
5 import math
6 def wordsRetrieve(html):
7     text = re.compile(r'<[^>]+>').sub('', html)
8     words = re.compile(r'[^A-z^a-z]+').split(text)
9
10    return [word.lower() for word in words if word]
11
12 input = "urlsappended"
13 output = "blogmatrixnew.txt"
14
15 def countRetrieve(url):
16
17     fd = feedparser.parse(url)
18     wc = collections.defaultdict(int)
19     stopwords = []
20
21     stopWordList = open('stopWordList.txt').readlines()
22     pages = len(fd['entries'])
23
24     for stopWord in stopWordList:
25         stopWord = stopWord.strip()
26         stopwords.append(stopWord)
27
28     for e in fd.entries:
29         if 'summary' in e:
30             summary = e.summary
31         else:
32             summary = e.description
33         words = wordsRetrieve('%s %s' % (e.title, summary))
```

```

34
35     for word in words:
36         if word not in stopwords:
37             wc[word] += 1
38
39 if pages == 500:
40     next_link = url + "?start-index=501"
41     d = feedparser.parse(next_link)
42     pages = len(d['entries'])
43     for e in d.entries:
44         if 'summary' in e:
45             summary = e.summary
46         else:
47             summary = e.description
48
49     words = wordsRetrieve('%s %s' % (e.title, summary))
50
51     for word in words:
52         if word not in stopwords:
53             #print word
54
55             wc[word] += 1
56
57 if pages == 500:
58     next_link = url + "?start-index=1001"
59     for e in d.entries:
60         if 'summary' in e:
61             summary = e.summary
62         else:
63             summary = e.description
64
65     words = wordsRetrieve('%s %s' % (e.title, summary))
66
67     for word in words:
68         if word not in stopwords:
69             #print word
70
71             wc[word] += 1
72 if pages == 500:
73     next_link = url + "?start-index=1501"
74     for e in d.entries:
75         if 'summary' in e:
76             summary = e.summary
77         else:
78             summary = e.description
79
80     words = wordsRetrieve('%s %s' % (e.title, summary))
81
82     for word in words:
83         if word not in stopwords:
84             wc[word] += 1
85
86 if 'title' not in fd.feed:
87     print 'Invalid url', url
88     return 'bogus data', wc
89
90 return fd.feed.title, wc

```

```

91
92
93 def matrixCreation():
94
95
96
97     apcount = collections.defaultdict(int)
98     wordcounts = {}
99     feedlist = open( input ).readlines()
100     totalWordCount = {}
101
102     for url in feedlist:
103         title, wc = countRetrieve(url)
104         wordcounts[title] = wc
105
106         for word, count in wc.iteritems():
107             if count > 1:
108                 apcount[word] += 1
109
110                 try:
111                     totalWordCount[word] += count
112                 except KeyError:
113                     totalWordCount[word] = count
114
115     wordlist = []
116
117
118     for w, bc in apcount.iteritems():
119         frac = float(bc)/len(feedlist)
120         #print frac
121         if frac > 0.1 and frac < 0.5:
122             wordlist.append(w)
123
124     countOfWords = []
125
126     for word in wordlist:
127         countOfWords.append((word, totalWordCount[word]))
128
129     countOfWords.sort(key=lambda rating: rating[1], reverse = True )
130
131     countOfWords = countOfWords[0:500]
132
133     out = file(outpput, 'w')
134     out.write('Blog')
135
136     idfWordCount = {}
137
138     for w in countOfWords:
139         word = w[0]
140         noOfBlogs = 0
141         for blogname, counts in wordcounts.iteritems():
142
143             if word in counts:
144                 noOfBlogs += 1
145
146         idf = math.log( 100.0 / noOfBlogs , 2 )
147

```

```

148         idfWordCount[word] = idf
149
150
151     for w in countOfWords:
152
153         out.write('\t' + w[0])
154
155     out.write('\n')
156
157     for blogname, counts in wordcounts.iteritems():
158         blogname = blogname.encode('UTF-8')
159         out.write(blogname)
160
161         for w in countOfWords:
162             word = w[0]
163             occurrence = w[1]
164
165             tf = float(counts[word]) / occurrence
166             tfidf = tf * idfWordCount[word]
167
168
169             out.write('\t%f' % tfidf)
170
171         out.write('\n')
172
173     out.close()
174
175 if __name__ == '__main__':
176     matrixCreation()
177

```

## 5.3 Output

### 5.3.1 Output Blog term matrix with tfidf values

	A	B	C	D	E	F	G	H	I	J	K	L
1	Blog	vocal	track	mso	pop	guitar	musical	http	review	style	voice	set
2	Flatbasset	0.003204	0.011198	0	0.02513	0.007064	0	0	0.002698	0.006353	0.001713	0.005788
3	Riley Haas	0.003204	0.005168	0	0.001047	0	0	0	0.010791	0.007941	0.001713	0.002894
4	(Insert W	0	0	0	0	0	0	0	0	0	0	0
5	SEM REGR	0	0	0	0.003141	0	0	0	0	0	0	0
6	Pithy Title	0	0.023258	0.057151	0.034554	0.018365	0.009801	0	0.002698	0.019058	0.015414	0.039068
7	Morgan's	0.001602	0.001723	0	0	0.002825	0.00196	0	0	0	0	0.002894
8	MARISOL	0	0	0	0	0	0	0	0	0	0	0
9	THE HUB	0	0	0.078207	0.002094	0	0	0	0	0.015882	0	0
10	Brian's Mu	0	0.005168	0	0.012565	0.002825	0.003921	0.021079	0	0.003176	0.008563	0.001447
11	Web Scier	0	0.006891	0	0	0	0	0.274021	0.018884	0.007941	0	0.023152
12	Steel City	0.003204	0.009475	0	0.010471	0.007064	0.005881	0	0.005395	0.003176	0.005138	0.001447
13	MR. BEAU	0	0	0	0.006282	0	0	0.003513	0	0	0	0
14	60@60 So	0	0	0	0	0.001413	0	0	0	0.001588	0	0.010129
15	ORGANMY	0	0	0	0	0	0	0	0	0.001588	0	0
16	MarkEOrt	0.001602	0.000861	0	0.010471	0	0	0	0	0	0	0.008682
17	Green Egg	0	0.002584	0.042111	0.001047	0	0	0	0.002698	0.007941	0	0
18	GLI Press	0	0.004307	0	0	0	0	0	0	0	0.001713	0
19	turnitup!	0.003204	0.005168	0	0.0178	0.007064	0.021563	0	0	0.004765	0.001713	0.001447
20	Stories Fro	0	0.000861	0	0.009424	0	0	0.024592	0	0	0	0.011576
21	A H T A P C	0	0	0	0	0.001413	0	0	0	0	0	0
22	Floorshim	0.001602	0.006891	0	0.005235	0.002825	0.003921	0.003513	0	0.001588	0	0
23	Did Not Cl	0.004807	0.001723	0	0.021989	0.009889	0.00196	0	0	0.007941	0.001713	0.004341
24	How to be	0	0	0	0	0	0	0	0	0	0	0
25	Party Full	0.001602	0.009475	0	0.00733	0.001413	0.003921	0	0	0	0.005138	0.001447

Figure 5: Output Blog term matrix with tfidf values

## 5.3.2 Output JPEG

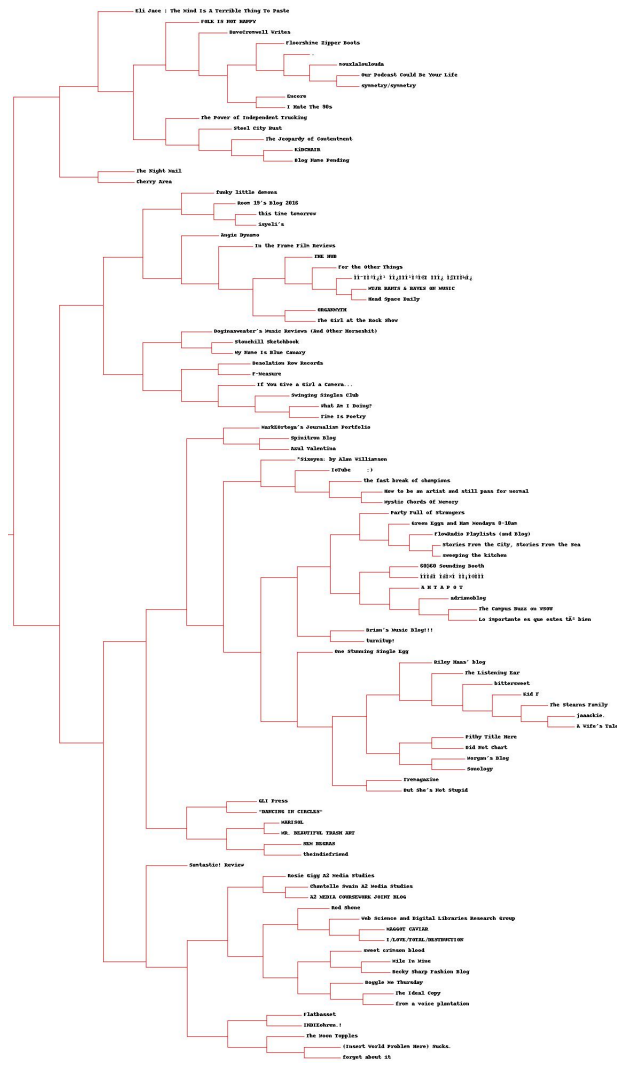


Figure 6: Output generated by mds.py



## References

- [1] Feed Parser <https://github.com/Unode/python-feedparser/tree/master/feedparser>,
- [2] Finding TFIDF value <http://en.wikipedia.org/wiki/Tfidf>,
- [3] Python for beginners <http://www.pythonforbeginners.com/feedparser/using-feedparser-in-python>

[]