



ASSIGNMENT 8

CS 432 Spring 2017



APRIL 12, 2017
Michelle Graham

Q1:

To begin with, I created a set to contain the default url: <https://www.blogger.com/next-blog?navBar=true&blogID=3471633091411211117>. This is used to generate random blogs. I attempted to create a counter alongside each new blog to ensure that it was working properly. I also made an attempt at making use of the "r.status_code" function in the requests library to verify that each link was in fact "OK" or returned 200. However, this proved to be quite finicky so I kind of just left it alone for the time being. In any case, I went on to print each element of the set, along with the two links specified in the assignment description: <http://f-measure.blogspot.com> and <http://ws-dl.blogspot.com>. I appended the generated list of blogs to include these desired links. Also, some of my links were bad links. I simply re-ran my code. I deleted the bad links manually and inserted the newly created ones in their place.

In order to create a blog term matrix, I made use of the generatefeedvector.py file from the PCI book. I hope to find time to make the layout nicer by adjusting some things. For the moment, I have left it as it is due to time constraints. I also had to cover a Unicode error for some of my links.

getblogs.py:

```
import requests

f = open('urls.txt', 'w')
s = set()
while (len(s) < 100):
    url = "http://www.blogger.com/next-blog?navBar=true&blogID=3471633091411211117"
    r = requests.get(url)
    r.status_code == requests.codes.ok
    print('good link #') + str(r.status_code)
    update = r.url.strip('/?expref=next-blog') + '/feeds/posts/default?alt=rss' + '\n'
    s.add(update)
for element in s:
    print element
    f.write(element + '\n')
f.write('http://f-measure.blogspot.com' + '\n' + '\n')
f.write('http://ws-dl.blogspot.com')
```

Unicode error fix in generatefeedvector.py:

```
import re
import sys

reload(sys)
sys.setdefaultencoding('utf-8')
```

blogdata1.txt:

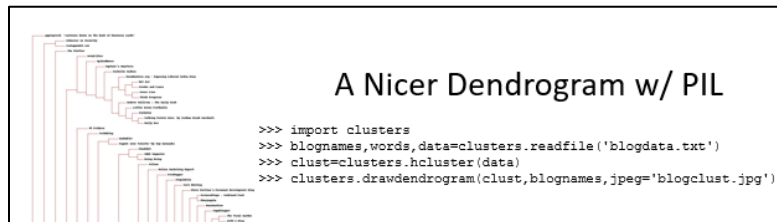
Blog	kids	golden	catchy	travel	wrong	fit	songwriter	service	needed	feeling
g	weekend	father	okay	filled	box	boy	bob	teenage	making	performance
eeeping	science	beautiful		sense	imagine	number	relationship	video	play	plan
n	drums	island	videos	field	starting		fall	town	none	blood
face	fact	bring	decade	should	listened		means	packed	ends	pool
n	obvious	february		apparently	mid		mix	taking	add	although
haven	rough	familiar		lucky	h	wire	etc	upon	less	paul
general	appearance		helped	thought	soft	alive	noise	brand	j	sleep
bittersweet	1	0		0	0	5	1	0	0	2
0	4	1		0	5	1	2	1	0	0
12	3	2		0	0	1	8	0	0	2
adrianoblog	0	0		0	0	0	0	0	0	0
0	0	0		0	0	0	0	0	0	0
0	0	0		0	0	0	0	0	0	0
tumbleweed	1	0		0	0	1	0	0	1	4
0	0	0		0	0	0	0	0	1	0
0	0	0		0	0	0	0	0	0	1
the traveling neighborhood				0	0	0	1	0	0	0
0	0	0		0	0	0	0	0	0	0
0	1	0		0	0	0	0	0	0	0
LOST PLACES	0	0		0	0	1	0	0	0	0
15	0	0		1	0	1	0	1	0	0
1	2	0		0	0	0	1	1	0	0
the rooﬂy leak	0	0		3	0	1	0	2	0	0
0	0	1		0	0	0	0	0	0	2

Q2:

The clustering PowerPoint from class was very helpful for this question. In particular, slides 12 and 13 were very helpful!

PPT code:

```
>>> import clusters
>>> blognames, words, data=clusters.readfile('blogdata.txt') # returns blog
titles, words in blog (10%-50% boundaries), list of frequency info
>>> clust=clusters.hcluster(data) # returns a tree of foo.id, foo.left,
foo.right
>>> clusters.printclust(clust, labels=blognames) # walks tree and prints
ascii approximation of a dendrogram; distance measure is Pearson's r
-
  gapingvoid: "cartoons drawn on the back of business cards"
-
-
  Schneier on Security
  Instapundit.com
-
  The Blotter
-
```



dendo.py:

```
import clusters
import sys
import codecs

blognames, words, data = clusters.readfile('blogdata1.txt')

clust = clusters.hcluster(data)

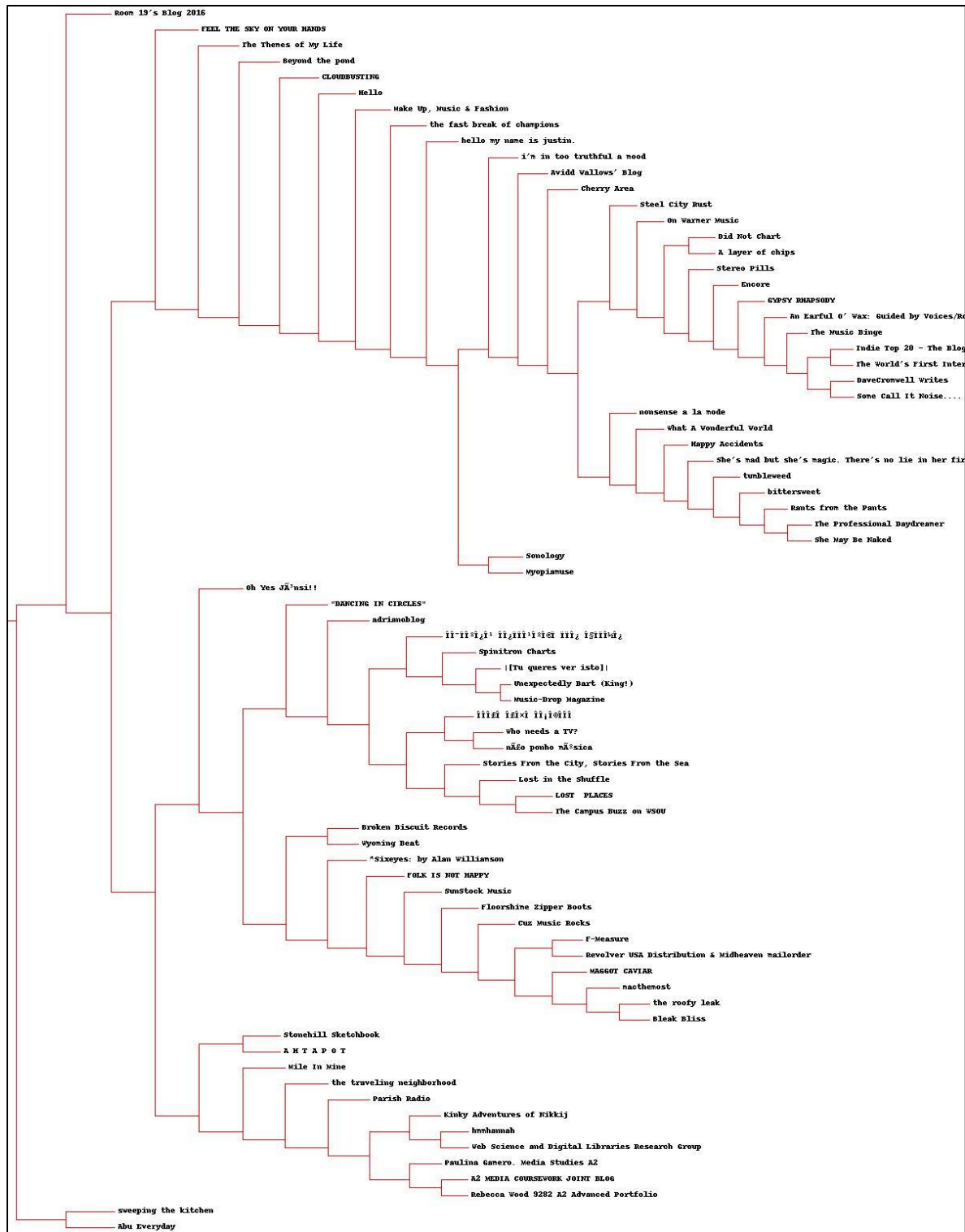
with codecs.open('ascii.text', 'w+', encoding="ascii") as fout:
    clusters.printclust(clust, labels=blognames)

clusters.drawdendrogram(clust, blognames, jpeg='dendrogram.jpg')
```

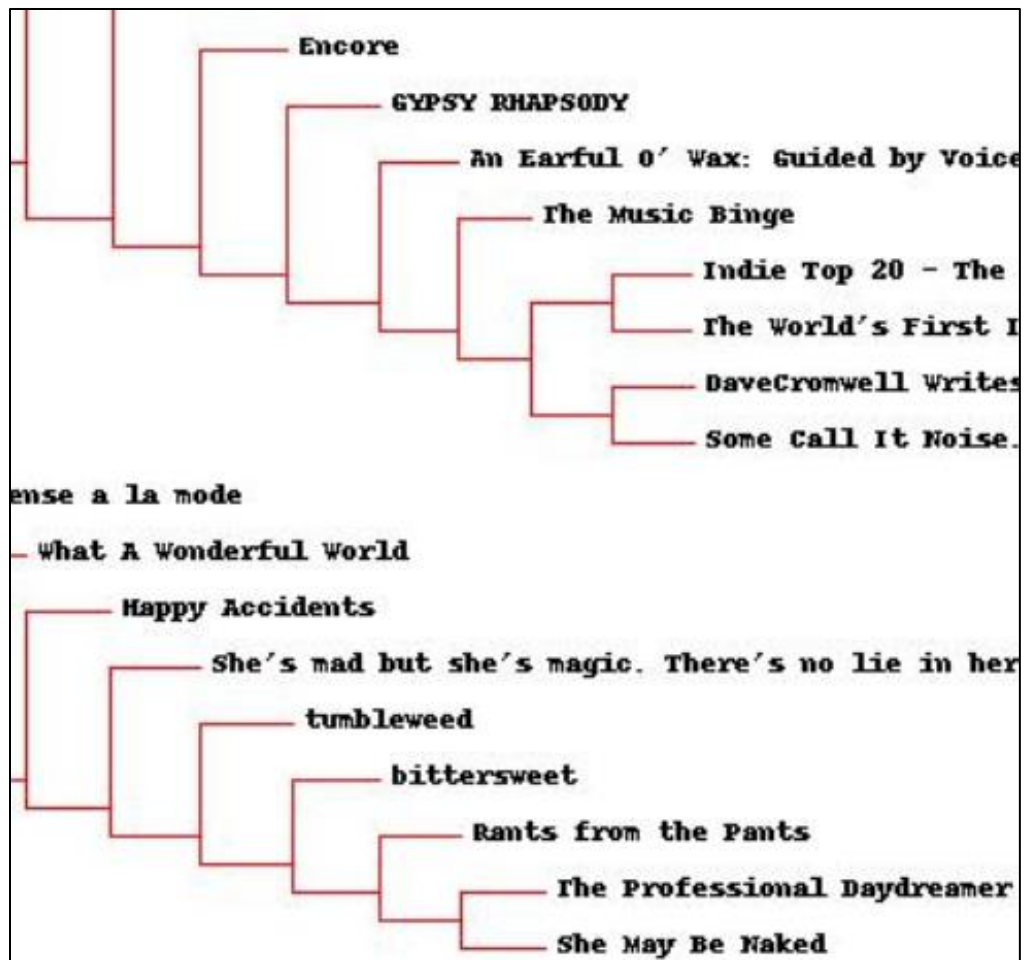
Ascii output:

```
the Music binge
-
-
  Indie Top 20 - The Blog!
  The World's First Internet Baby
-
  DaveCromwell Writes
  Some Call It Noise....
-
  nonsense a la mode
-
  What A Wonderful World
-
  Happy Accidents
-
  She's mad but she's magic. There's no lie in her fire
-
  tumbleweed
-
  bittersweet
-
  Rants from the Pants
```

Dendrogram:



Zoomed In:



Q3:

Slide 18 of the clustering PowerPoint from class was helpful for this question.

Python Example

```
>>> kclust=clusters.kcluster(data,k=10) # K-Means with 10 centroids, pp. 43-44
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
>>> [blognames[r] for r in kclust[0]] # print blognames in 1st centroid
```

k.py:

```
import clusters

blognames, words, data = clusters.readfile('blogdata1.txt')

print('k = 5')
kclust = clusters.kcluster(data, k = 5)
k= [blognames[r] for r in kclust[0]]
print str(k) + '\n'

print('k = 10')
kclust = clusters.kcluster(data, k = 10)
print str(k) + '\n'

print('k = 20')
kclust = clusters.kcluster(data, k = 20)
print str(k)
```

Output:

```
mngrah@DESKTOP-30IR4AC MINGW64 /c/
$ python k.py
k = 5
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
['Bleak Bliss', 'On Warmer Music'
es/Robert Pollard Song by Song Re

k = 10
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
['Bleak Bliss', 'On Warmer Music'
es/Robert Pollard Song by Song Re

k = 20
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
Iteration 7
Iteration 8
['Bleak Bliss', 'On Warmer Music'
es/Robert Pollard Song by Song Re
```

Q4:

For this question, I referenced slide 28 of the clustering PowerPoint from class again.

Slide 28:

```
Python MDS

>>> blognames, words, data=clusters.readfile('blogdata.txt')
>>> coords=clusters.scaledown(data) # code pp. 50-51
4431.91264139
3597.09879465
3530.63422919
3494.58547715
3463.77217455
3437.59298469
3414.89864608
3395.55257233
3378.52510767
3363.87951104
...
3024.12202228
3024.01331202
3023.87527696
3023.74986258
3023.75364032
>>> clusters.draw2d(coords, blognames, jpeg='blogs2d.jpg')
```

error starts to get worse,
so we're done

Mds.py:

```
import clusters

blognames, words, data = clusters.readfile('blogdata1.txt')

c = clusters.scaledown(data)

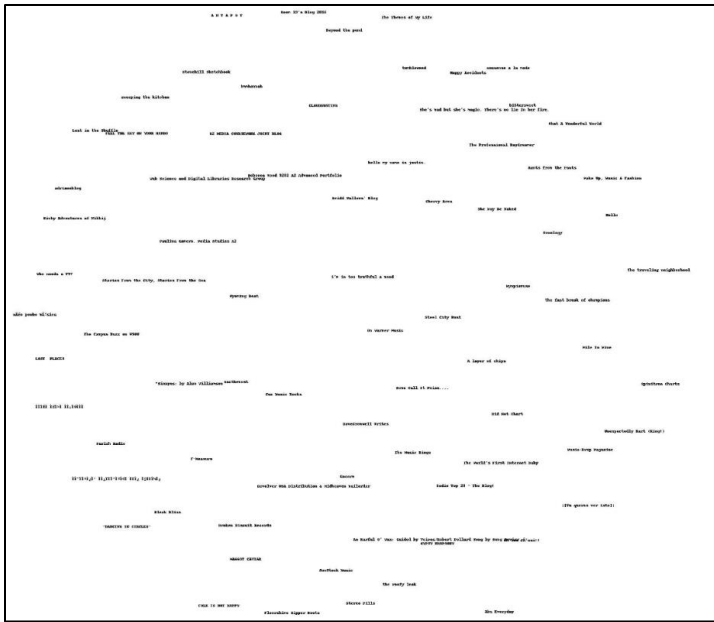
clusters.draw2d(c, blognames, jpeg='mds.jpg')
```

Output:

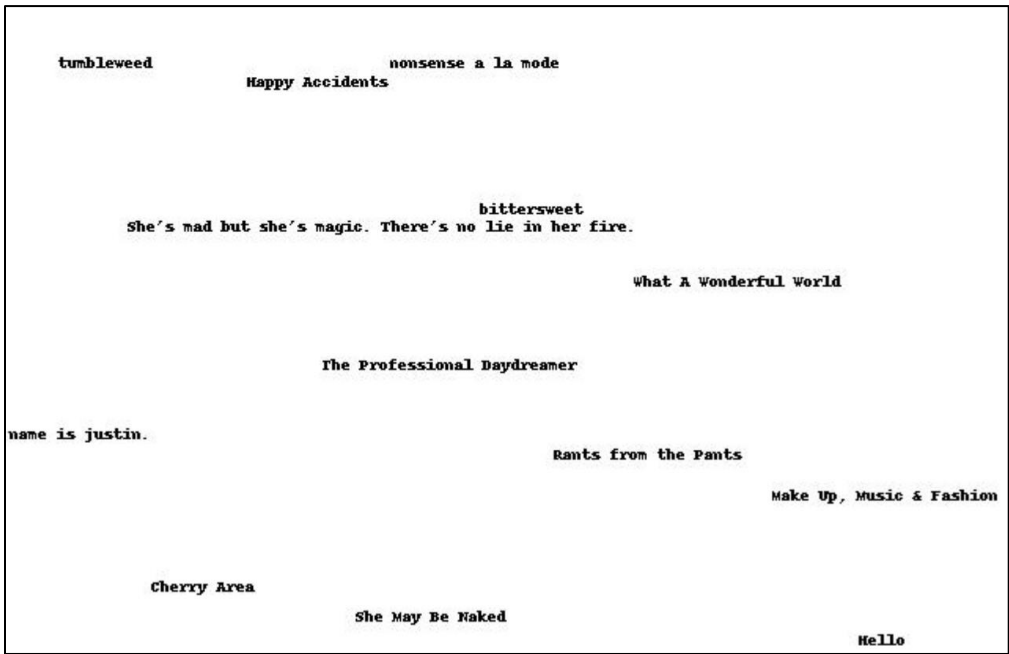
There were 399 lines in total for the output, therefore it took 399 iterations to complete.

```
ngrah@DESKTOP-30IR4AC MINGW64
$ python mds.py
2739.42516545
2112.37891003
2054.0681426
2026.59248226
2007.70865227
1993.68110686
1981.83044623
1970.61693953
1961.79902307
1953.81500606
1946.52680459
1940.54970573
1935.385816
1930.69749169
1926.88789803
1923.51995139
1920.51921087
1917.82767647
1915.7325674
1913.68804534
1911.66557805
1909.65474654
1907.80217187
1906.28646634
1904.8314638
1903.42327638
1901.87336000
```


Jpeg:



Zoomed in:



Resources:

<https://github.com/uolter/PCI/tree/master/chapter3>

<http://stackoverflow.com/questions/4706499/how-do-you-append-to-a-file>

<http://stackoverflow.com/questions/14852480/about-handling-a-redirection-in-python>