

# INTRODUCTION TO WEB SCIENCES: Assignment 9

Babitha Bokka

06 December 2014

# Contents

<b>1</b>	<b>Question 1:</b>	<b>2</b>
1.1	Approach . . . . .	2
1.2	Source Code . . . . .	3
1.2.1	extractLinks.sh . . . . .	3
1.2.2	makingAtoms.py . . . . .	3
1.2.3	generatefeedvector.py . . . . .	4
1.2.4	pagesEachBlog.py . . . . .	5
1.3	Output Files . . . . .	6
<b>2</b>	<b>Question 2:</b>	<b>7</b>
2.1	Approach . . . . .	7
2.2	Source Code . . . . .	8
2.2.1	generateAscii.py . . . . .	8
2.2.2	clusters.py . . . . .	8
2.3	Output Files . . . . .	14
2.3.1	blogclust.jpg . . . . .	14
<b>3</b>	<b>Question 3:</b>	<b>15</b>
3.1	Approach . . . . .	15
3.2	Source Code . . . . .	15
3.2.1	calculateK-mean.py . . . . .	15
3.3	Output Files . . . . .	16
3.3.1	output.png . . . . .	16
<b>4</b>	<b>Question 4:</b>	<b>17</b>
4.1	Approach . . . . .	17
4.2	Source Code . . . . .	17
4.2.1	createMDS.py . . . . .	17
4.3	Output Files . . . . .	17

# 1 Question 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.

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.

Create a histogram of how many pages each blog has (e.g., 30 blogs with just one page, 27 with two pages, 29 with 3 pages and so on).

## 1.1 Approach

1. In doing this assignment I took help from Mallika and Victor. They explained me about the packages to be installed and what to do? and what not to do?.
2. To solve this questions I wrote shell script to grab the other 98 links.
3. I wrote a short python script to append the "feeds/posts/default?max-results=500".
4. And to get the top 500 words I modified feedparser.py.

## 1.2 Source Code

### 1.2.1 extractLinks.sh

```
1 #!/bin/bash
2
3 if [ $# -ne 1 ]
4 then
5     echo "usage <extractLinks.sh> <number of links> "
6     echo "e.g., extractLinks.sh 98 "
7     exit
8 fi
9
10 numberOfLinks=$1
11
12 link1="http://f-measure.blogspot.com/"
13 link2="http://ws-dl.blogspot.com/"
14
15 echo "$link1" > blogLinksLists.txt
16 echo "$link2" >> blogLinksLists.txt
17
18 for (( count=1; count<=numberOfLinks; count++ ))
19 do
20     curl -I -L 'http://www.blogger.com/next-blog?navBar=true&blogID
        =3471633091411211117' | grep Location | tail -n 1 | cut -d" " -f2 | cut -d"?" -f 1
        >> blogLinksLists.txt
21 done
```

### 1.2.2 makingAtoms.py

```
1 #!/usr/bin/env python
2
3 import sys
4
5 def main():
6     saveFile = open('blogLinksLists.txt', 'r')
7     writeFile = open('blogLinksAtom.txt', 'w')
8     add = "feeds/posts/default?max-results=500"
9     for line in saveFile:
10         writeFile.write(line+add)
11         writeFile.write("\n")
12     saveFile.close()
13     writeFile.close()
14
15 if __name__ == "__main__":
16     try:
17         main()
18     except KeyboardInterrupt:
19         sys.exit(1)
```

### 1.2.3 generatefeedvector.py

```
1 #!/usr/bin/env python
2
3 import re
4 import sys
5 import feedparser
6
7 # Returns title and dictionary of word counts for an RSS feed
8 def getwordcounts(url):
9     # Parse the feed
10    d=feedparser.parse(url)
11    wc={}
12
13    # Loop over all the entries
14    for e in d.entries:
15        if 'summary' in e: summary=e.summary
16        else: summary=e.description
17
18        # Extract a list of words
19        words=getwords(e.title+' '+summary)
20        for word in words:
21            wc.setdefault(word,0)
22            wc[word]+=1
23    return d.feed.title,wc
24
25 def getwords(html):
26     # Remove all the HTML tags
27     txt=re.compile(r'<[>]+>').sub('',html)
28
29     # Split words by all non-alpha characters
30     words=re.compile(r'[^A-Za-z]+').split(txt)
31
32     # Convert to lowercase
33     return [word.lower() for word in words if word!='']
34
35 def main():
36
37     apcount={}
38     wordcounts={}
39     feedlist=[line for line in file('blogLinksAtom.txt')]
40     for feedurl in feedlist:
41         try:
42             title,wc=getwordcounts(feedurl)
43             wordcounts[title]=wc
44             for word,count in wc.items():
45                 apcount.setdefault(word,0)
46                 if count>1:
47                     apcount[word]+=1
48         except:
49             print 'Failed to parse feed %s' % feedurl
50
51     wordlist=[]
52     countFrequentWords = []
53     for w,bc in apcount.items():
54         frac=float(bc)/len(feedlist)
55         if frac>0.1 and frac<0.5:
```

```

56     countFrequentWords.append((w, bc))
57
58 countFrequentWords=sorted(countFrequentWords, key=lambda x:x[1], reverse = True)
59
60 for value in countFrequentWords:
61     # word
62     value1 = value[0]
63     #count
64     value2 = value[1]
65     length = len(wordlist)
66     if(length < 500):
67         wordlist.append(value1)
68     else:
69         break
70
71 out=file('blogdata.txt','w')
72 out.write('Blog')
73
74 for word in wordlist:
75     word1 = word.encode('UTF-8')
76     out.write('\t%s' % word1)
77 out.write('\n')
78
79 for blog,wc in wordcounts.items():
80     blogName = blog.encode('UTF-8')
81     print blog
82     out.write(blogName)
83     for word in wordlist:
84         if word in wc: out.write('\t%d' % wc[word])
85         else: out.write('\t0')
86     out.write('\n')
87
88 if __name__ == "__main__":
89     try:
90         main()
91     except KeyboardInterrupt:
92         sys.exit(1)

```

#### 1.2.4 pagesEachBlog.py

Could not finish to generate the Histogram.

```

1 #!/usr/bin/env python
2
3 import sys
4 import time
5 import socket
6 import urllib2
7 import unicodedata
8 from bs4 import BeautifulSoup
9
10
11
12
13 def main():
14     saveFile = open('blogLinksLists.txt', 'r')
15     writeFile = open('pagesBlog.txt', 'w')
16     #add = "feeds/posts/default"

```

```

17     for url in saveFile.readlines():
18         try:
19             response = urllib2.urlopen(url, timeout=30)
20             html_content = response.read()
21             # get the html content using beautiful soup
22             soup = BeautifulSoup(html_content)
23
24             try:
25                 links = soup.find('link', rel='next', href=True)['href']
26                 print links
27             except TypeError:
28                 pass
29             except KeyError:
30                 pass
31
32             except urllib2.HTTPError:
33                 pass
34             except urllib2.URLError:
35                 pass
36             except socket.timeout:
37                 pass
38
39
40
41
42     saveFile.close()
43     writeFile.close()
44
45 if __name__ == "__main__":
46     try:
47         main()
48     except KeyboardInterrupt:
49         sys.exit(1)

```

## 1.3 Output Files

1. The output for this question is a matrix which is uploaded in Github blogdata.txt.

## 2 Question 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).

### 2.1 Approach

1. To answer this question I used clusters.py.
2. Used the lines of code provided in the presentation.
3. Python program generateAscii.py produces ASCII results which are piped to an output file ascii.txt.
4. Also produces dendrogram blogclust.jpg.



## 2.2 Source Code

### 2.2.1 generateAscii.py

```
1#!/usr/bin/env python
2
3import sys
4import clusters
5
6def main():
7
8    # returns blog titles , words in blog (10%–50% boundaries), list of frequency info
9    blognames, words, data=clusters.readfile('blogdata.txt')
10
11    # returns a tree of foo.id, foo.left, foo.right
12    clust=clusters.hcluster(data)
13
14    # walks tree and prints ascii approximation of a dendrogram; distance measure is
15    # Pearson's r
16    clusters.printclust(clust, labels=blognames)
17
18if __name__ == "__main__":
19    try:
20        main()
21    except KeyboardInterrupt:
22        sys.exit(1)
```

### 2.2.2 clusters.py

This program is imported into question 2, 3 and 4.

```
1from PIL import Image, ImageDraw
2
3def readfile(filename):
4    lines=[line for line in file(filename)]
5
6    # First line is the column titles
7    colnames=lines[0].strip().split('\t')[1:]
8    rownames=[]
9    data=[]
10    for line in lines[1:]:
11        p=line.strip().split('\t')
12        # First column in each row is the rowname
13        rownames.append(p[0])
14        # The data for this row is the remainder of the row
15        data.append([float(x) for x in p[1:]])
16    return rownames, colnames, data
17
18
19from math import sqrt
20
21def pearson(v1, v2):
22    # Simple sums
23    sum1=sum(v1)
24    sum2=sum(v2)
25
26    # Sums of the squares
27    sum1Sq=sum([pow(v,2) for v in v1])
```

```

28 sum2Sq=sum([pow(v,2) for v in v2])
29
30 # Sum of the products
31 pSum=sum([v1[i]*v2[i] for i in range(len(v1))])
32
33 # Calculate r (Pearson score)
34 num=pSum-(sum1*sum2/len(v1))
35 den=sqrt((sum1Sq-pow(sum1,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1)))
36 if den==0: return 0
37
38 return 1.0-num/den
39
40 class bicluster:
41     def __init__(self, vec, left=None, right=None, distance=0.0, id=None):
42         self.left=left
43         self.right=right
44         self.vec=vec
45         self.id=id
46         self.distance=distance
47
48     def hcluster(rows, distance=pearson):
49         distances={}
50         currentclustid=-1
51
52         # Clusters are initially just the rows
53         clust=[bicluster(rows[i], id=i) for i in range(len(rows))]
54
55         while len(clust)>1:
56             lowestpair=(0,1)
57             closest=distance(clust[0].vec, clust[1].vec)
58
59             # loop through every pair looking for the smallest distance
60             for i in range(len(clust)):
61                 for j in range(i+1, len(clust)):
62                     # distances is the cache of distance calculations
63                     if (clust[i].id, clust[j].id) not in distances:
64                         distances[(clust[i].id, clust[j].id)]=distance(clust[i].vec, clust[j].vec)
65
66                     d=distances[(clust[i].id, clust[j].id)]
67
68                     if d<closest:
69                         closest=d
70                         lowestpair=(i, j)
71
72             # calculate the average of the two clusters
73             mergevec=[
74                 (clust[lowestpair[0]].vec[i]+clust[lowestpair[1]].vec[i])/2.0
75                 for i in range(len(clust[0].vec))]
76
77             # create the new cluster
78             newcluster=bicluster(mergevec, left=clust[lowestpair[0]],
79                                   right=clust[lowestpair[1]],
80                                   distance=closest, id=currentclustid)
81
82             # cluster ids that weren't in the original set are negative
83             currentclustid-=1
84             del clust[lowestpair[1]]

```

```

85     del clust[lowestpair[0]]
86     clust.append(newcluster)
87
88     return clust[0]
89
90 def printclust(clust, labels=None, n=0):
91     # indent to make a hierarchy layout
92     for i in range(n): print ' ',
93     if clust.id < 0:
94         # negative id means that this is branch
95         print '-'
96     else:
97         # positive id means that this is an endpoint
98         if labels==None: print clust.id
99         else: print labels[clust.id]
100
101     # now print the right and left branches
102     if clust.left!=None: printclust(clust.left, labels=labels, n=n+1)
103     if clust.right!=None: printclust(clust.right, labels=labels, n=n+1)
104
105 def getheight(clust):
106     # Is this an endpoint? Then the height is just 1
107     if clust.left==None and clust.right==None: return 1
108
109     # Otherwise the height is the same of the heights of
110     # each branch
111     return getheight(clust.left)+getheight(clust.right)
112
113 def getdepth(clust):
114     # The distance of an endpoint is 0.0
115     if clust.left==None and clust.right==None: return 0
116
117     # The distance of a branch is the greater of its two sides
118     # plus its own distance
119     return max(getdepth(clust.left), getdepth(clust.right))+clust.distance
120
121
122 def drawdendrogram(clust, labels, jpeg='clusters.jpg'):
123     # height and width
124     h=getheight(clust)*20
125     w=1200
126     depth=getdepth(clust)
127
128     # width is fixed, so scale distances accordingly
129     scaling=float(w-150)/depth
130
131     # Create a new image with a white background
132     img=Image.new('RGB', (w,h), (255,255,255))
133     draw=ImageDraw.Draw(img)
134
135     draw.line((0, h/2, 10, h/2), fill=(255,0,0))
136
137     # Draw the first node
138     drawnode(draw, clust, 10, (h/2), scaling, labels)
139     img.save(jpeg, 'JPEG')
140
141 def drawnode(draw, clust, x, y, scaling, labels):

```

```

142 if clust.id<0:
143     h1=getheight(clust.left)*20
144     h2=getheight(clust.right)*20
145     top=y-(h1+h2)/2
146     bottom=y+(h1+h2)/2
147     # Line length
148     ll=clust.distance*scaling
149     # Vertical line from this cluster to children
150     draw.line((x,top+h1/2,x,bottom-h2/2),fill=(255,0,0))
151
152     # Horizontal line to left item
153     draw.line((x,top+h1/2,x+ll,top+h1/2),fill=(255,0,0))
154
155     # Horizontal line to right item
156     draw.line((x,bottom-h2/2,x+ll,bottom-h2/2),fill=(255,0,0))
157
158     # Call the function to draw the left and right nodes
159     drawnode(draw,clust.left,x+ll,top+h1/2,scaling,labels)
160     drawnode(draw,clust.right,x+ll,bottom-h2/2,scaling,labels)
161 else:
162     # If this is an endpoint, draw the item label
163     draw.text((x+5,y-7),labels[clust.id],(0,0,0))
164
165 def rotatematrix(data):
166     newdata=[]
167     for i in range(len(data[0])):
168         newrow=[data[j][i] for j in range(len(data))]
169         newdata.append(newrow)
170     return newdata
171
172 import random
173
174 def kcluster(rows,distance=pearson,k=4):
175     # Determine the minimum and maximum values for each point
176     ranges=[(min([row[i] for row in rows]),max([row[i] for row in rows]))
177             for i in range(len(rows[0]))]
178
179     # Create k randomly placed centroids
180     clusters=[[random.random()*(ranges[i][1]-ranges[i][0])+ranges[i][0]
181               for i in range(len(rows[0]))] for j in range(k)]
182
183     lastmatches=None
184     for t in range(100):
185         print 'Iteration %d' % t
186         bestmatches=[]
187
188         # Find which centroid is the closest for each row
189         for j in range(len(rows)):
190             row=rows[j]
191             bestmatch=0
192             for i in range(k):
193                 d=distance(clusters[i],row)
194                 if d<distance(clusters[bestmatch],row): bestmatch=i
195             bestmatches[bestmatch].append(j)
196
197         # If the results are the same as last time, this is complete
198         if bestmatches==lastmatches: break

```

```

199 lastmatches=bestmatches
200
201 # Move the centroids to the average of their members
202 for i in range(k):
203     avgs=[0.0]*len(rows[0])
204     if len(bestmatches[i])>0:
205         for rowid in bestmatches[i]:
206             for m in range(len(rows[rowid])):
207                 avgs[m]+=rows[rowid][m]
208             for j in range(len(avgs)):
209                 avgs[j]/=len(bestmatches[i])
210             clusters[i]=avgs
211
212 return bestmatches
213
214 def tanamoto(v1,v2):
215     c1,c2,shr=0,0,0
216
217     for i in range(len(v1)):
218         if v1[i]!=0: c1+=1 # in v1
219         if v2[i]!=0: c2+=1 # in v2
220         if v1[i]!=0 and v2[i]!=0: shr+=1 # in both
221
222     return 1.0-(float(shr)/(c1+c2-shr))
223
224 def scaledown(data,distance=pearson,rate=0.01):
225     n=len(data)
226
227     # The real distances between every pair of items
228     realdist=[[distance(data[i],data[j]) for j in range(n)]
229               for i in range(0,n)]
230
231     # Randomly initialize the starting points of the locations in 2D
232     loc=[[random.random(),random.random()] for i in range(n)]
233     fakedist=[[0.0 for j in range(n)] for i in range(n)]
234
235     lasterror=None
236     for m in range(0,1000):
237         # Find projected distances
238         for i in range(n):
239             for j in range(n):
240                 fakedist[i][j]=sqrt(sum([pow(loc[i][x]-loc[j][x],2)
241                                         for x in range(len(loc[i]))]))
242
243         # Move points
244         grad=[[0.0,0.0] for i in range(n)]
245
246         totalerror=0
247         for k in range(n):
248             for j in range(n):
249                 if j==k: continue
250                 # The error is percent difference between the distances
251                 errorterm=(fakedist[j][k]-realdist[j][k])/realdist[j][k]
252
253                 # Each point needs to be moved away from or towards the other
254                 # point in proportion to how much error it has
255                 grad[k][0]+=((loc[k][0]-loc[j][0])/fakedist[j][k])*errorterm

```

```

256         grad[k][1]+=((loc[k][1]-loc[j][1])/fakedist[j][k])*errorterm
257
258         # Keep track of the total error
259         totalerror+=abs(errorterm)
260     print totalerror
261
262     # If the answer got worse by moving the points, we are done
263     if lasterror and lasterror<totalerror: break
264     lasterror=totalerror
265
266     # Move each of the points by the learning rate times the gradient
267     for k in range(n):
268         loc[k][0]-=rate*grad[k][0]
269         loc[k][1]-=rate*grad[k][1]
270
271     return loc
272
273 def draw2d(data, labels, jpeg='mds2d.jpg'):
274     img=Image.new( 'RGB', (2000,2000), (255,255,255) )
275     draw=ImageDraw.Draw(img)
276     for i in range(len(data)):
277         x=(data[i][0]+0.5)*1000
278         y=(data[i][1]+0.5)*1000
279         draw.text((x,y), labels[i], (0,0,0))
280     img.save(jpeg, 'JPEG')

```

### 2.3.1 blogclust.jpg

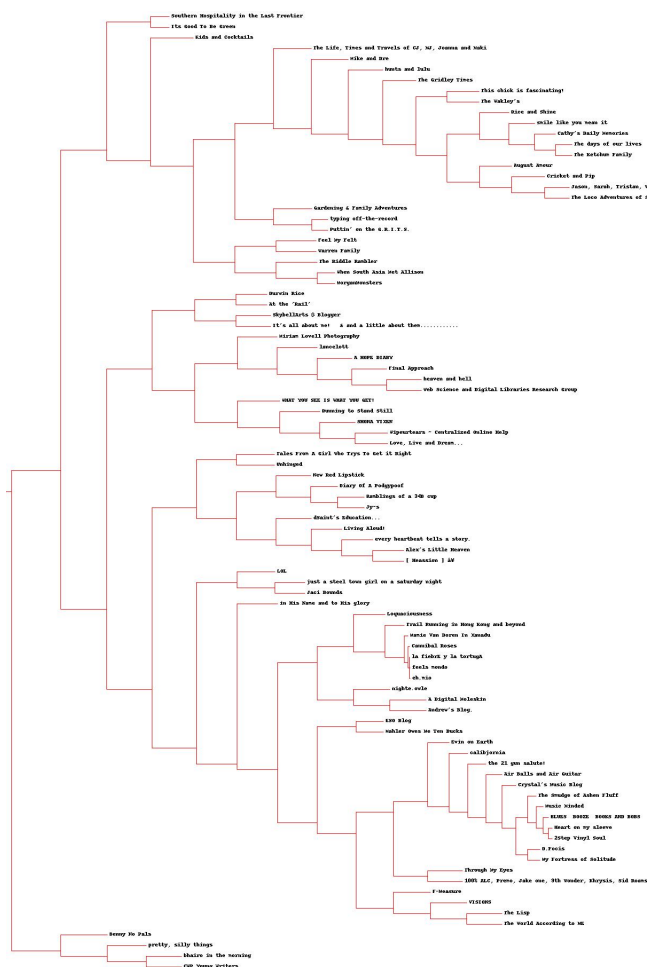


Figure 1: Dendrogram

### 3 Question 3:

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

#### 3.1 Approach

1. To solve this question I used lines of code provided in the presentation (calculateK-mean.py) .

#### 3.2 Source Code

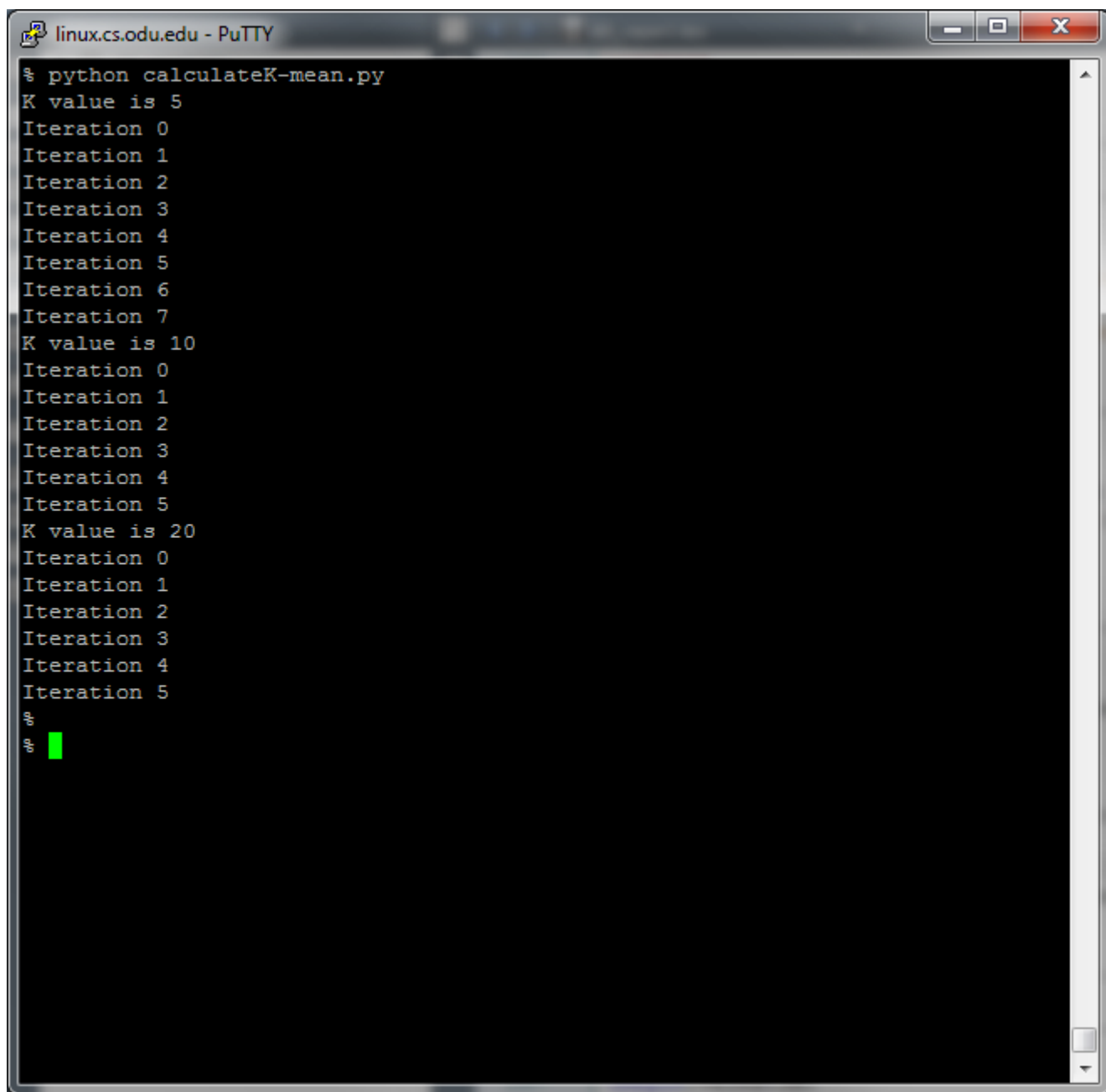
##### 3.2.1 calculateK-mean.py

```
1 #!/usr/bin/env python
2
3 import clusters
4
5 def main():
6
7     blognames, words, data=clusters.readfile('blogdata.txt')
8     print "K value is 5"
9     kclust=clusters.kcluster(data,k=5)
10    print "K value is 10"
11    kclust=clusters.kcluster(data,k=10)
12    print "K value is 20"
13    kclust=clusters.kcluster(data,k=20)
14
15
16 if __name__ == "__main__":
17     try:
18         main()
19     except KeyboardInterrupt:
20         sys.exit(1)
```



## 3.3 Output Files

### 3.3.1 output.png



```
linux.cs.odu.edu - PuTTY
% python calculateK-mean.py
K value is 5
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
Iteration 7
K value is 10
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
K value is 20
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
%
% █
```

Figure 2: Output of iterations for different k values

## 4 Question 4:

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

### 4.1 Approach

1. To solve this questions I used the few lines of code provided in the presentation (createMDS.py).

### 4.2 Source Code

#### 4.2.1 createMDS.py

```
1 #!/usr/bin/env python
2
3 import clusters
4
5 def main():
6
7     blognames, words, data=clusters.readfile('blogdata.txt')
8     coords=clusters.scaledown(data)
9     clusters.draw2d(coords, blognames, jpeg='blogs2d.jpg')
10
11
12 if __name__ == "__main__":
13     try:
14         main()
15     except KeyboardInterrupt:
16         sys.exit(1)
```

### 4.3 Output Files

```
% python createMDS.py
4515.42807078
5524.56099361
% python createMDS.py
4299.13362421
5370.72157029
% python createMDS.py
4868.240652
7215.65051908
```

18

## References

- [1] feedparser.py program. <http://www.cs.odu.edu/~mln/teaching/cs595-f13/chapter3/code/generatefeedvector.python>.
- [2] Python for beginners. <http://www.pythonforbeginners.com/feedparser/using-feedparser-in-python>, October 2009.
- [3] eumiro. Sorted function in python. <http://stackoverflow.com/questions/14218933/sorted-function-in-python>, January 2013.
- [4] Toby Segaran. *Programming Collective Intelligence*. O'Reilly Media, August 2007.

□