

# CS 595: Assignment #9

Due on Thursday, December 4, 2014

*Dr. Nelson 4:20pm*

Holly Harkins

## Contents

<b>Problem 1</b>	<b>3</b>
<b>Problem 2</b>	<b>6</b>
<b>Problem 3</b>	<b>8</b>
<b>Problem 4</b>	<b>9</b>

## Problem 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).

Answer:

I used `generateURLs.py` to capture the 98 Blogs (`feedlist.txt`).

Next step was to create a matrix. The matrix required was using blog title as the identifier for each blog row and using the terms for the columns of the matrix. The values are the frequency of occurrence. Using `generatefeedvector.py`, TFIDF is calculated as the score that balances the frequency of a term in a document vs. frequency of a term in all the documents. A limit to the number of terms to the most "popular" 500 terms (`blogdata.txt`).

To create the Histogram, I used `getNextPage.py` to generate a count of how many pages each blog has (`numberOfPages.txt`). Unfortunately I ran out of time and did not generate the count. There was an error in the code I could not resolve. I would have used R code to generate the histogram.

### Listing 1: Generating Blog URLs

```
#!/usr/bin/python3
import sys
from bs4 import BeautifulSoup
import urllib.request
5 from urllib.parse import urlparse
```

```
DEFAULT_COUNT=98
DEFAULT_SEED_URL='http://www.blogger.com/next-blog?navBar=true&blogID
=3471633091411211117'
if len(sys.argv) != 3:
10     print('Pass the blog count, defaulting to ' + str(DEFAULT_COUNT))
    print('Pass the seed URL, defaulting to ' + DEFAULT_SEED_URL)
    count=DEFAULT_COUNT
    url=DEFAULT_SEED_URL
else:
15     count=sys.argv[1]
    url=sys.argv[2]

def parse(link):
    response = urllib.request.urlopen(link)
20     soup = BeautifulSoup(response.read())
    response.close()
    return soup

def addNext(url,s):
25     try:
        soup=parse(url)
        for atom in soup.findAll('link',rel='alternate',type='application/atom+xml'):
            atomHref=str(atom['href']).strip()
            s.add(atomHref)
30         print('Added atom href: ' + atomHref)
    except:
        print('Exception parsing URL, skipping to next')
        pass

35 s=set()
while len(s) < count:
    addNext(url,s)
with open('feedlist.txt', 'w') as f:
    for atom in s:
40         f.write(atom + '\n')
```

Listing 2: Generating Matrix

```
#!/usr/bin/python3
import feedparser
import re

5 # Returns title and dictionary of word counts for an RSS feed
def getwordcounts(url):
    # Parse the feed
    d=feedparser.parse(url)
    wc={}

10     # Loop over all the entries
    for e in d.entries:
        if 'summary' in e: summary=e.summary
        else: summary=e.description

15     # Extract a list of words
```

```

        words=getwords(e.title+' '+summary)
        for word in words:
            wc.setdefault(word,0)
20         wc[word]+=1
        return d.feed.title,wc

def getwords(html):
    # Remove all the HTML tags
25     txt=re.compile(r'<[>]+>').sub('',html)

    # Split words by all non-alpha characters
    words=re.compile(r'[^A-Za-z]+').split(txt)

30     # Convert to lowercase
    return [word.lower() for word in words if word!='']

apcount={}
35 wordcounts={}
feedlist=[line for line in file('feedlist.txt')]
for feedurl in feedlist:
    try:
        title,wc=getwordcounts(feedurl)
40         wordcounts[title]=wc
        for word,count in wc.items():
            apcount.setdefault(word,0)
            if count>1:
                apcount[word]+=1
45     except:
        print 'Failed to parse feed %s' % feedurl

#modification on the TFIDF calculations
wordlist=[]
50 for w,bc in sorted(apcount.items(), key=lambda t: t[1], reverse=True):
    frac=float(bc)/len(feedlist)
    if frac>0.1 and frac<0.5:
        wordlist.append(w)

55 out=file('blogdata.txt','w')
out.write('Blog')

#Limit the number of terms to the most "popular" 500 terms
wordlist=wordlist[0:500]
60 for word in wordlist: out.write('\t%s' % word)
out.write('\n')
for blog,wc in wordcounts.items():
    print blog
    out.write(blog)
65     for word in wordlist:
        if word in wc: out.write('\t%d' % wc[word])
        else: out.write('\t0')
    out.write('\n')
print 'Total number of words=' +str(len(wordlist))
```

## Problem 2

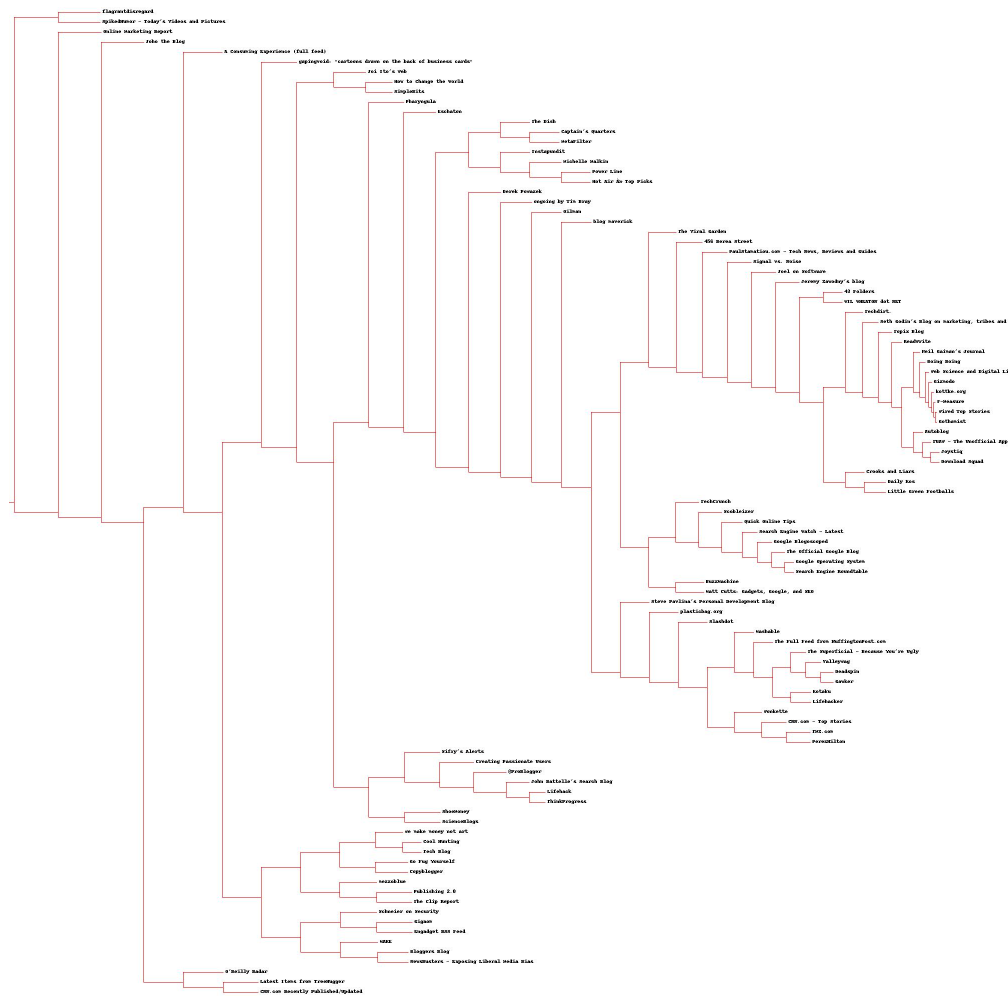
Create an ASCII and JPEG dendrogram that clusters (i.e., HAC) the most similar blogs (see slides 12 & 13). Include the JPEG in your report and upload the ascii file to github (it will be too unwieldy for inclusion in the report).

Answer:

Assignment9.py uses Pearsons correlation and averages the data for the 2 old clusters. The final cluster returned can be searched recursively to recreate all the clusters and their end nodes (ASCII.txt).

The drawdendrogram function uses the Python Imaging Library to generate Clust\_Dendrogram.jpg.

Below is the dendrogram that was created from the blog matrix.



## Problem 3

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

Answer:

Assignment9.py uses `kcluster` function to place a  $K$  points randomly in space that represent the centre of the cluster and assigns each blog to the nearest one. Then moves centroids to the average location of all the nodes that assigned to them. Repeats this process until the assignments stop changing.

For  $K = 5$ , there were 5 iterations needed to cluster the blogs.  
For  $k = 10$ , there were 6 iterations needed to cluster the blogs. For  $k=20$ , there were 5 iterations needed to cluster the blogs. The output generated is found in `Kiterations.txt`.



## Problem 4

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

Answer:

Assignment9.py uses 2 functions to create the MDS. Scaledown function takes the data vector obtained from readfile function, finds the difference between each pair of blogs, and matches them to a distances using Pearson correlation. These distances will represent the distances between the blogs in the chart. Returns the X and Y coordinates of the blogs on the two-dimensional chart.

Function draw2d uses Python Imaging Library to to generate an image with all the labels of all the different blogs plotted at the coordinates of that blog which obtained from first function.



Figure 2: MDS

Listing 3: Assignment9 Python

```
import sys
import random
from PIL import Image, ImageDraw
from math import sqrt

f = open('ASCII.txt', 'w')

def readfile(filename):
    lines=[line for line in file(filename)]
    # First line is the column titles
    colnames=lines[0].strip().split('\t')[1:]
    rownames=[]
    data=[]
    for line in lines[1:]:
        p=line.strip().split('\t')
        # First column in each row is the rowname
        rownames.append(p[0])
```

```

        # The data for this row is the remainder of the row
        data.append([float(x) for x in p[1:]])
20     return rownames, colnames, data

def pearson(v1,v2):
    # Simple sums
25     sum1=sum(v1)
        sum2=sum(v2)
    # Sums of the squares
        sum1Sq=sum([pow(v,2) for v in v1])
        sum2Sq=sum([pow(v,2) for v in v2])
30     # Sum of the products
        pSum=sum([v1[i]*v2[i] for i in range(len(v1))])
    # Calculate r (Pearson score)
        num=pSum-(sum1*sum2/len(v1))
        den=sqrt((sum1Sq-pow(sum1,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1)))
35     if den==0: return 0
        return 1.0-num/den

class bicluster:
    def __init__(self, vec, left=None, right=None, distance=0.0, id=None):
40         self.left=left
            self.right=right
            self.vec=vec
            self.id=id
            self.distance=distance
45

def hcluster(rows, distance=pearson):

    distances={}
    currentclustid=-1
50     # Clusters are initially just the rows
    clust=[bicluster(rows[i],id=i) for i in range(len(rows))]
    while len(clust)>1:
        lowestpair=(0,1)
        closest=distance(clust[0].vec, clust[1].vec)
55     # loop through every pair looking for the smallest distance
        for i in range(len(clust)):
            for j in range(i+1, len(clust)):
                # distances is the cache of distance calculations
                if (clust[i].id, clust[j].id) not in distances:
60                     distances[(clust[i].id, clust[j].id)] = distance(clust[i].vec,
                                                                    clust[j].vec)
                d=distances[(clust[i].id, clust[j].id)]
                if d<closest:
                    closest=d
                    lowestpair=(i, j)
65     # calculate the average of the two clusters
        mergevec=[(clust[lowestpair[0]].vec[i]+ clust[lowestpair[1]].vec[i]) /2.0
                  for i in range(len(clust[0].vec))]
    # create the new cluster
        newcluster=bicluster(mergevec, left=clust[lowestpair[0]],

```

```

    right=clust[lowestpair[1]],
    distance=closest,id=currentclustid)
70  # cluster ids that weren't in the original set are negative
    currentclustid-=1
    del clust[lowestpair[1]]
    del clust[lowestpair[0]]
75  clust.append(newcluster)

    return clust[0]
def printclust(clust,labels=None,n=0):
    # indent to make a hierarchy layout

80  for i in range(n): f.write(' '),
    if clust.id<0:
        # negative id means that this is branch
        f.write('-\n')
    else:
85  # positive id means that this is an endpoint
        if labels==None:
            f.write(clust.id+'\n')
        else:
            f.write(labels[clust.id]+'\n')
90  # now print the right and left branches
    if clust.left!=None: printclust(clust.left,labels=labels,n=n+1)
    if clust.right!=None: printclust(clust.right,labels=labels,n=n+1)

95 def getheight(clust):
    # Is this an endpoint? Then the height is just 1
    if clust.left==None and clust.right==None:
        return 1
    # Otherwise the height is the same of the heights of
100 # each branch
    return getheight(clust.left)+getheight(clust.right)

def getdepth(clust):
    # The distance of an endpoint is 0.0
105 if clust.left==None and clust.right==None:
        return 0
    # The distance of a branch is the greater of its two sides
    # plus its own distance
    return max(getdepth(clust.left),getdepth(clust.right))+clust.distance

110 def drawdendrogram(clust,labels,jpeg='clusters.jpg'):
    # height and width
    h=getheight(clust)*20
    w=2000
115 depth=getdepth(clust)
    # width is fixed, so scale distances accordingly
    scaling=float(w-150)/depth
    # Create a new image with a white background
    img=Image.new('RGB',(w,h),(255,255,255))
120 draw=ImageDraw.Draw(img)
    draw.line((0,h/2,10,h/2),fill=(255,0,0))

```

```

    # Draw the first node
    drawnode(draw, clust, 10, (h/2), scaling, labels)
    img.save(jpeg, 'JPEG')

125 def drawnode(draw, clust, x, y, scaling, labels):
    if clust.id < 0:
        h1 = getheight(clust.left) * 20
        h2 = getheight(clust.right) * 20
130         top = y - (h1 + h2) / 2
        bottom = y + (h1 + h2) / 2

        # Line length
        ll = clust.distance * scaling
135         # Vertical line from this cluster to children
        draw.line((x, top + h1 / 2, x, bottom - h2 / 2), fill=(255, 0, 0))
        # Horizontal line to left item
        draw.line((x, top + h1 / 2, x + ll, top + h1 / 2), fill=(255, 0, 0))
        # Horizontal line to right item
140         draw.line((x, bottom - h2 / 2, x + ll, bottom - h2 / 2), fill=(255, 0, 0))
        # Call the function to draw the left and right nodes
        drawnode(draw, clust.left, x + ll, top + h1 / 2, scaling, labels)
        drawnode(draw, clust.right, x + ll, bottom - h2 / 2, scaling, labels)
    else:
145         # If this is an endpoint, draw the item label
        draw.text((x + 5, y - 7), labels[clust.id], (0, 0, 0))

def kcluster(rows, distance=pearson, k=4):
150     # Determine the minimum and maximum values for each point
    ranges = [(min([row[i] for row in rows]), max([row[i] for row in rows]))]
    for i in range(len(rows[0])):
        # Create k randomly placed centroids
        clusters = [[random.random() * (ranges[i][1] - ranges[i][0]) + ranges[i][0]]
155         for i in range(len(rows[0])) for j in range(k)]
    lastmatches = None
    for t in range(100):
        print 'Iteration %d' % t
        bestmatches = [[] for i in range(k)]
160         # Find which centroid is the closest for each row
        for j in range(len(rows)):
            row = rows[j]
            bestmatch = 0
            for i in range(k):
165                 d = distance(clusters[i], row)
                if d < distance(clusters[bestmatch], row): bestmatch = i
            bestmatches[bestmatch].append(j)

        # If the results are the same as last time, this is complete
170         if bestmatches == lastmatches: break
        lastmatches = bestmatches
        # Move the centroids to the average of their members
        for i in range(k):
            avgs = [0.0] * len(rows[0])

```

```

175         if len(bestmatches[i])>0:
            for rowid in bestmatches[i]:
                for m in range(len(rows[rowid])):
                    avgs[m]+=rows[rowid][m]
            for j in range(len(avgs)):
180                 avgs[j]/=len(bestmatches[i])
            clusters[i]=avgs
        return bestmatches

def scaledown(data,distance=pearson,rate=0.01):
185     n=len(data)
    # The real distances between every pair of items
    realdist=[[distance(data[i],data[j]) for j in range(n)] for i in range(0,n)]
    # Randomly initialize the starting points of the locations in 2D
    loc=[[random.random(),random.random()] for i in range(n)]
190     fakedist=[[0.0 for j in range(n)] for i in range(n)]
    lasterror=None
    for m in range(0,1000):
        # Find projected distances
        for i in range(n):
195             for j in range(n):
                fakedist[i][j]=sqrt(sum([pow(loc[i][x]-loc[j][x],2)
                    for x in range(len(loc[i]))]))
        # Move points
        grad=[[0.0,0.0] for i in range(n)]
200
        totalerror=0
        for k in range(n):
            for j in range(n):
                if j==k: continue
205             # The error is percent difference between the distances
                errorterm=(fakedist[j][k]-realdist[j][k])/ realdist[j][k]
            # Each point needs to be moved away from or towards the other
            # point in proportion to how much error it has
                grad[k][0]+=((loc[k][0]-loc[j][0])/ fakedist[j][k])*errorterm
210                grad[k][1]+=((loc[k][1]-loc[j][1])/ fakedist[j][k])*errorterm
            # Keep track of the total error
                totalerror+=abs(errorterm)
        print totalerror
        # If the answer got worse by moving the points, we are done
215         if lasterror and lasterror<totalerror: break
        lasterror=totalerror
        # Move each of the points by the learning rate times the gradient
        for k in range(n):
            loc[k][0]-=rate*grad[k][0]
220            loc[k][1]-=rate*grad[k][1]
    return loc

def draw2d(data,labels,jpeg='mds2d.jpg'):
225     img=Image.new('RGB', (2000,2000), (255,255,255))
    draw=ImageDraw.Draw(img)
    for i in range(len(data)):
        x=(data[i][0]+0.5)*1000

```

```
        y=(data[i][1]+0.5)*1000
        draw.text((x,y),labels[i],(0,0,0))
230     img.save(jpeg,'JPEG')

def main():
    # Start Q2 Creating an ASCII and JPEG dendrogram
    blognames,words,data=readfile('blogdata.txt')
235     cluster=hcluster(data)

    printclust(cluster,labels=blognames)
    f.close()
    drawdendrogram(cluster,blognames,jpeg='Clust_Dendrogram.jpg')
240

    # Start Q3 using K-Means Clustering
    print "Clustering the blogs using k=5"
    kclust=kcluster(data,k=5)
245     print "\nClustering the blogs using k=10"
    kclust=kcluster(data,k=10)
    print "\nClustering the blogs using k=20"
    kclust=kcluster(data,k=20)

250
    # Start Q4 MDS
    print "\nMDS Creation"
    coords=scaledown(data)
    draw2d(coords,blognames,jpeg='MDS.jpg')
255

main();
```

## References

- [1] <https://github.com/nico/collectiveintelligence-book/blob/master/feedlist.txt>.