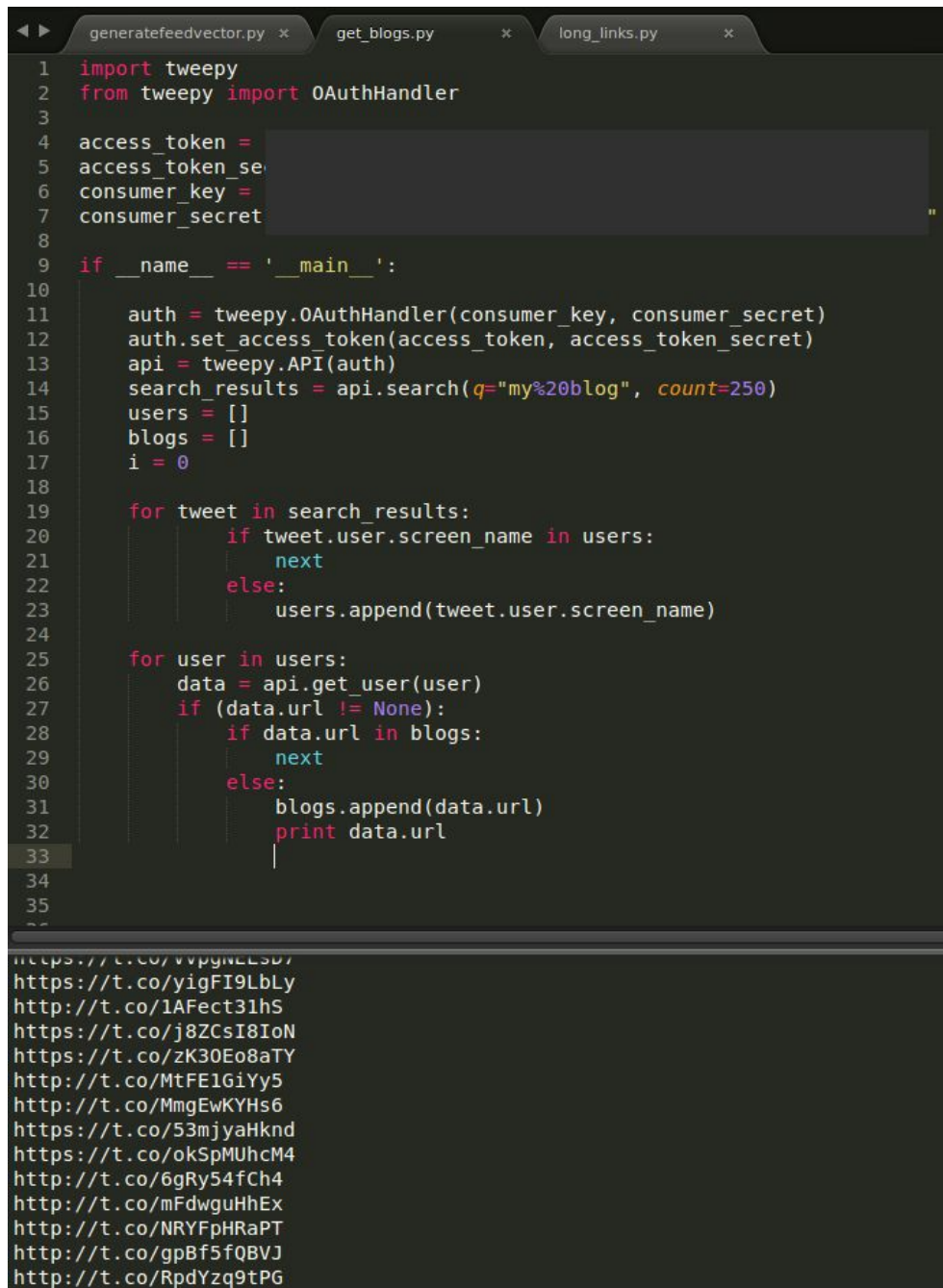


PART 1

In an effort to get a bunch of blogs with hopefully no repeats or offensive material, I first searched Twitter for “my blog” and made a set of usernames from the results, then took each user’s website URI, assuming most of them would have their blog on their website (I got this idea, but not the code, [here](#)):



The image shows a code editor with three tabs: `generatefeedvector.py`, `get_blogs.py`, and `long_links.py`. The `get_blogs.py` tab is active, displaying a Python script that uses the Tweepy library to search Twitter for the query "my%20blog" and extract website URLs from the resulting users. The script includes imports for Tweepy and OAuthHandler, sets up authentication with an access token and consumer key/secret, searches for tweets, filters for unique usernames, and then iterates through those usernames to find and print their website URLs. The output at the bottom of the editor shows a list of 15 Twitter URLs.

```
1 import tweepy
2 from tweepy import OAuthHandler
3
4 access_token =
5 access_token_secret =
6 consumer_key =
7 consumer_secret =
8
9 if __name__ == '__main__':
10
11     auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
12     auth.set_access_token(access_token, access_token_secret)
13     api = tweepy.API(auth)
14     search_results = api.search(q="my%20blog", count=250)
15     users = []
16     blogs = []
17     i = 0
18
19     for tweet in search_results:
20         if tweet.user.screen_name in users:
21             next
22         else:
23             users.append(tweet.user.screen_name)
24
25     for user in users:
26         data = api.get_user(user)
27         if (data.url != None):
28             if data.url in blogs:
29                 next
30             else:
31                 blogs.append(data.url)
32                 print data.url
33
34
35
36
37 https://t.co/vvp9wEESb7
38 https://t.co/yigFI9LbLy
39 http://t.co/1AFect3lhS
40 https://t.co/j8ZCsI8IoN
41 https://t.co/zK30Eo8aTY
42 http://t.co/MtFE1GiYy5
43 http://t.co/MmgEwKYHs6
44 https://t.co/53mjyaHknd
45 https://t.co/okSpMUhcM4
46 http://t.co/6gRy54fCh4
47 http://t.co/mFdwguHhEx
48 http://t.co/NRYFpHRaPT
49 http://t.co/gpBf5fQBVJ
50 http://t.co/RpdYzq9tPG
```

I used [this](#) for guidance in following the links and getting the long URI, skipping broken links:

```
generatefeedvector.py x get_blogs.py x long_links.py x
1 from __future__ import unicode_literals, print_function
2 import ttp
3 import requests
4 import urllib2
5
6 def follow_shortlinks(shortlinks):
7     links_followed = {}
8     for shortlink in shortlinks:
9         url = shortlink
10        request_result = requests.get(url)
11        redirect_history = request_result.history
12        all_urls = []
13        for redirect in redirect_history:
14            all_urls.append(redirect.url)
15        all_urls.append(request_result.url)
16        links_followed[shortlink] = all_urls
17    return links_followed
18
19 def on_error(self, status):
20     pass
21
22 if __name__ == '__main__':
23
24     shortlinks = []
25
26     with open('short_links.txt') as f:
27         for line in f:
28             shortlinks.append(line)
29
30     followed_shortlinks = (follow_shortlinks(shortlinks))
31
32     for link in followed_shortlinks:
33         try:
34             print (urllib2.urlopen(link).url)
35         except urllib2.HTTPError as e:
36             print ('\t', e)
37
38 http://morganna-designs.squarespace.com/
39 http://ec.europa.eu/france/index_fr.htm
40 http://salessolveeverything.com/
41 https://www.instagram.com/justjess369/
42 http://www.audriestorme.com/
43 http://www.amazon.com/registry/wishlist/17TM87YY2ZF67/ref=cm_sw_r
44 https://jacquelinechambliss.wordpress.com/
45 https://annebellesbooks.wordpress.com/
46 http://www.songofthestitch.com/
47 http://fineartamerica.com/profiles/spyder-webb.html
48 http://jamesoneseventeen.blogspot.com/
49 http://www.garywalker.net/
50 https://configuroweb.blogspot.com/
51 http://www.adamchudv.com/
```

After that, I picked through and rid the list of unwanted links. Then, to find the RSS feeds of these websites, I used [this](#) for guidance (StackOverflow question):

```
generatefeedvector.py x  get_rss.py x
1 import requests
2 from bs4 import BeautifulSoup
3
4 def get_rss_feed(website_url):
5     if website_url is None:
6         print("URL should not be null")
7     else:
8         source_code = requests.get(website_url)
9         plain_text = source_code.text
10        soup = BeautifulSoup(plain_text, "lxml")
11        for link in soup.find_all("link", {"type" : "application/rss+xml"}):
12            href = link.get('href')
13            print str(href)
14
15 with open('blogs.txt') as f:
16     for line in f:
17         try:
18             print get_rss_feed(line)
19         except:
20             pass
```

```
http://f-measure.blogspot.com/feeds/posts/default?alt=rss
None
http://ws-dl.blogspot.com/feeds/posts/default?alt=rss
None
https://craftytaramarie.wordpress.com/feed/
https://craftytaramarie.wordpress.com/comments/feed/
None
http://radicalreadsbook.blogspot.com/feeds/posts/default?alt=rss
None
http://eroseweb.co/blog/feed/
http://eroseweb.co/blog/comments/feed/
None
None
http://hiiijaaackie.blogspot.com/feeds/posts/default?alt=rss
```

Note: I'm uploading the full results shown so far to GitHub, but will only put the RSS links in this report.

List of RSS links:

```
http://f-measure.blogspot.com/feeds/posts/default?alt=rss
http://ws-dl.blogspot.com/feeds/posts/default?alt=rss
https://craftytaramarie.wordpress.com/feed/
http://radicalreadsbook.blogspot.com/feeds/posts/default?alt=rss
http://eroseweb.co/blog/feed/
http://hiiijaaackie.blogspot.com/feeds/posts/default?alt=rss
http://www.notesfromcaroline.com/feeds/posts/default?alt=rss
http://thebusinessgirlsnetwork.com/feed/
http://deenadays.com/feed/
https://siefkenpublications.wordpress.com/feed/
```

<http://www.thekmprojects.gr/feeds/posts/default?alt=rss>
<http://www.coldknowledge.com/feeds/posts/default?alt=rss>
<http://www.gingeybites.com/feeds/posts/default?alt=rss>
<https://jasminecarlisleblog.wordpress.com/feed/>
<http://chinyeugonna.com/feed/>
<http://theguyliner.com/feed/>
<http://edufisbil.blogspot.com/feeds/posts/default?alt=rss>
<http://ericabachelor.com/feed/>
<http://labellablog.co/feed/>
<https://lillyteardrop.wordpress.com/feed/>
<http://kreshimir.com/feed/>
<http://momsncharge.com/feed/>
<http://sami-spoon.blogspot.com/feeds/posts/default?alt=rss>
<http://pedrothedagger.tumblr.com/rss>
<https://indiewatchsite.wordpress.com/feed/>
<http://justbellefashionblog.blogspot.com/feeds/posts/default?alt=rss>
<http://www.jessicafoley.ca/feed/>
<https://whenyousaidtulips.wordpress.com/feed/>
<http://whimsicaljoy.com/?feed=rss2>
<http://doginasweatershowreviews.blogspot.com/feeds/posts/default?alt=rss>
<http://justthatgirlonline.blogspot.com/feeds/posts/default?alt=rss>
<http://www.audriestorme.com/feed.xml>
<https://jacquelinechambliss.wordpress.com/feed/>
<https://annebellebooks.wordpress.com/feed/>
<http://jamesoneseventeen.blogspot.com/feeds/posts/default?alt=rss>
<https://configuroweb.blogspot.com/feeds/posts/default?alt=rss>
<http://www.adamchudy.com/feed/>
<http://www.papersoul.design.co.uk/feed.xml>
<http://apilgrimontheroad.blogspot.com/feeds/posts/default?alt=rss>
<http://theonionfield.blogspot.com/feeds/posts/default?alt=rss>
<http://fitlifefranticfree.blogspot.com/feeds/posts/default?alt=rss>
<http://drgperformancesolutions.com/feed>
<http://theidealcoppy.blogspot.com/feeds/posts/default?alt=rss>
<http://www.matthewmegyese.com/feed.xml>
<http://floorshimezipperboots.blogspot.com/feeds/posts/default?alt=rss>
<http://stewartkirby.blogspot.com/feeds/posts/default?alt=rss>
<http://thepioneerwoman.com/feed/>
<http://thepioneerwoman.com/life-and-style/feed/>
<http://emersonbloede.blogspot.com/feeds/posts/default?alt=rss>
<http://mcomv2.blogspot.com/feeds/posts/default?alt=rss>
<http://cookeatrun.com/feed/>
<http://toxicbreedsfunhouse.blogspot.com/feeds/posts/default?alt=rss>
<http://skinnyshoes.blogspot.com/feeds/posts/default?alt=rss>
<http://fivescrambledeggs.blogspot.com/feeds/posts/default?alt=rss>
<http://aninconsistentscraper.blogspot.com/feeds/posts/default?alt=rss>
<http://mondaywakeup.blogspot.com/feeds/posts/default?alt=rss>
<http://didnotchart.blogspot.com/feeds/posts/default?alt=rss>
<http://robynpeterman.com/feed/>
<http://robynpeterman.com/blog/feed/>
<https://fromauniqueperspective.wordpress.com/feed/>
<http://www.adventuresofalondonkiwi.com/feeds/posts/default?alt=rss>
<http://itslostitsfound.blogspot.com/feeds/posts/default?alt=rss>
<http://www.kevinkruse.com/feed/>
<http://www.kevinkruse.com/blog/feed/>
<https://thousandscarsblog.wordpress.com/feed/>
<http://henrycleaves.blogspot.com/feeds/posts/default?alt=rss>
<http://www.danijroberts.co.uk/feed/>

<https://siefkenpublications.wordpress.com/feed/>
<https://justjeanine.wordpress.com/feed/>
<http://www.careerheatcodes.com/feed/>
<http://extremeways.blog.jp/index.rdf>
<https://chasrad.wordpress.com/feed/>
<http://rilphly.com/feed/>
<https://sardecoo.wordpress.com/feed/>
<http://iguacen.com/rss>
<http://iguacen.com/feed/>
<https://www.etsy.com/shop/CrowsNestJewels/rss>
<http://littleheartbiglove.co.uk/feed/>
<http://www.oxv234.com/feeds/posts/default?alt=rss>
<https://latinashadows.wordpress.com/feed/>
<http://powai.tumblr.com/rss>
<http://theesongbird.tumblr.com/rss>
<https://tylershepard1991.wordpress.com/feed/>
<http://madh-mama.blogspot.com/feeds/posts/default?alt=rss>
<https://earth2levi.wordpress.com/feed/>
<https://thegreenmockingbird.wordpress.com/feed/>
<https://ajbookreviewclub.wordpress.com/feed/>
<https://www.mmstartups.com/feed/>
<https://www.mmstartups.com/author/abe-salisbury/feed/>
<http://www.danjroberts.co.uk/feed/>
<http://kmillshernandez.com/feed/>
<http://iguacen.com/rss>
<http://iguacen.com/feed/>
<https://justjeanine.wordpress.com/feed/>
<http://rkt2.blog.fc2.com/?xml>
<https://latinashadows.wordpress.com/feed/>
<http://carrie-i.blogspot.com/feeds/posts/default?alt=rss>
<https://thebookreviewcafe.wordpress.com/feed/>
<http://katiebugstory.blogspot.com/feeds/posts/default?alt=rss>
<http://configuroweb.blogspot.com/feeds/posts/default?alt=rss>

Next, I created a matrix of word frequencies for the blogs, limiting the word list to the 500 most popular terms. I used the code found in `generatefeedvector.py` from Programming Collective Intelligence [here](#) to create the matrix and put a red rectangle around the code I added when I was trying to get the 500 most frequent words with a [sorted dictionary](#):

```
import feedparser
import re

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

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

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

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

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

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

if __name__ == '__main__':
```

The screenshot shows a text editor with two files open. The left file, `blogdata.txt`, contains a matrix of word frequencies. The matrix is a 2D array where each row represents a blog and each column represents a word. The words are sorted by frequency, with the most frequent words in the first columns. The matrix is as follows:

	don	has	how	they	post	will	love	been	blog	who	make	really
back	our	get	know	had	their	well	here	year	think	ve	things	her
no	which	see	these	people	very	than	were	them	would	because	most	two
into	way	great	work	now	take	why	over	got	part	say	good	he
other	go	only	always	re	also	even	going	where	last	could	another	find
she	much	while	d	through	want	after	years	reading	many	life	around	world
long	made	an	every	few	didn	since	his	week	off	us	little	did
may	still	never	right	found	lot	look	being	let	next	something	then	
days	need	should	old	down	before	ll	feel	show	ever	live	end	put
any	sure	sane	thought	read	does	start	once	today	went	doing	nbsp	better
those	own	everyone	use	each	hin	thing	too	review	set	away	cone	
change	give	such	free	times	top	hope	said	both	used	night	again	full
started	making	name	list	place	came	best	probably	big	bit	doesn	friends	
school	family	music	actually	pretty	different	left	book	home	getting	looking		
nothing	enough	amp	continue	everything	together	having	finally	girl	took			
seen	able	com	mind	three	happy	until	video	tell	help	perfect	side	try
story	bad	myself	fun	check	trying	follow	must	share	high	keep	watch	coming
its	style	house	wanted	least	point	others	whole	yet	easy	real	anything	
months	ago	e	makes	talk	looks	might	someone	quite	already	haven	b	though
stuff	almost	without	yes	often	run	sometimes	far	past	second	write	course	
amazing	mean	hard	idea	done	kind	future	especially	later	face	black	decided	
however	comes	light	less	couple	become	twitter	oh	single	th	enjoy	writing	along
beautiful	man	remember	maybe	facebook	group	tsn	friend	loved	saw			
during	won	between	taking	n	super	working	head	called	young	under	small	
business	wrong	break	couldn	yourself	wasn	white	believe	feeling	online	play		
definitely	stop	case	person	rest	spring	water	nice	felt	blue	room	heart	
favorite	huge	interesting	half	close	please	job	c	moment	early	photo		
lost	series	forward	call	inside	picture	age	focus	month	order	hours	wait	
sunday	instead	front	special	red	anyone	behind	words	morning	art	seems	recently	
sweet	ready	etc	single	march	dark	although	using	four	posts	told	kids	
experience	usually	heard	means	including	thank	l	books	god	reason	turn		
sounds	open	seen	ways	based	money	hit	across	asked	pop	running	understand	
deep	leave	hand	thanks	questions	either	dead	knew	road	seeing	support	true	
minutes	trip	worth	outside	cool	above	ask	taken	known	learn	rather	soon	
personal	late	weeks	space	size	girls	needed	looked	link	pictures	v	april	
create	due	else	country	spent	add	possible	five	works	v	buy	below	
general	important	fact	lives	starting	level	p	strong	changed	visit	phone		
given	word	following	exactly	january	date	short	main	similar	expect	tried	friday	city
guess	eyes	mostly	nine	talking	self	note	clear	release	national	lots		
appeared	walk	simply	women	worked	http	thoughts	hands	excited	power	rock		
area	form	line	places	whether	fall	completely	win	played	type	car	crazy	
playing	cold	choose	food	february								

The right file, `generatefeedvector.py`, contains the Python script. A red rectangle highlights the following code block:

```
temp = {}
most_frequent_words = []

for word in wordlist:
    temp[word] = account[word]

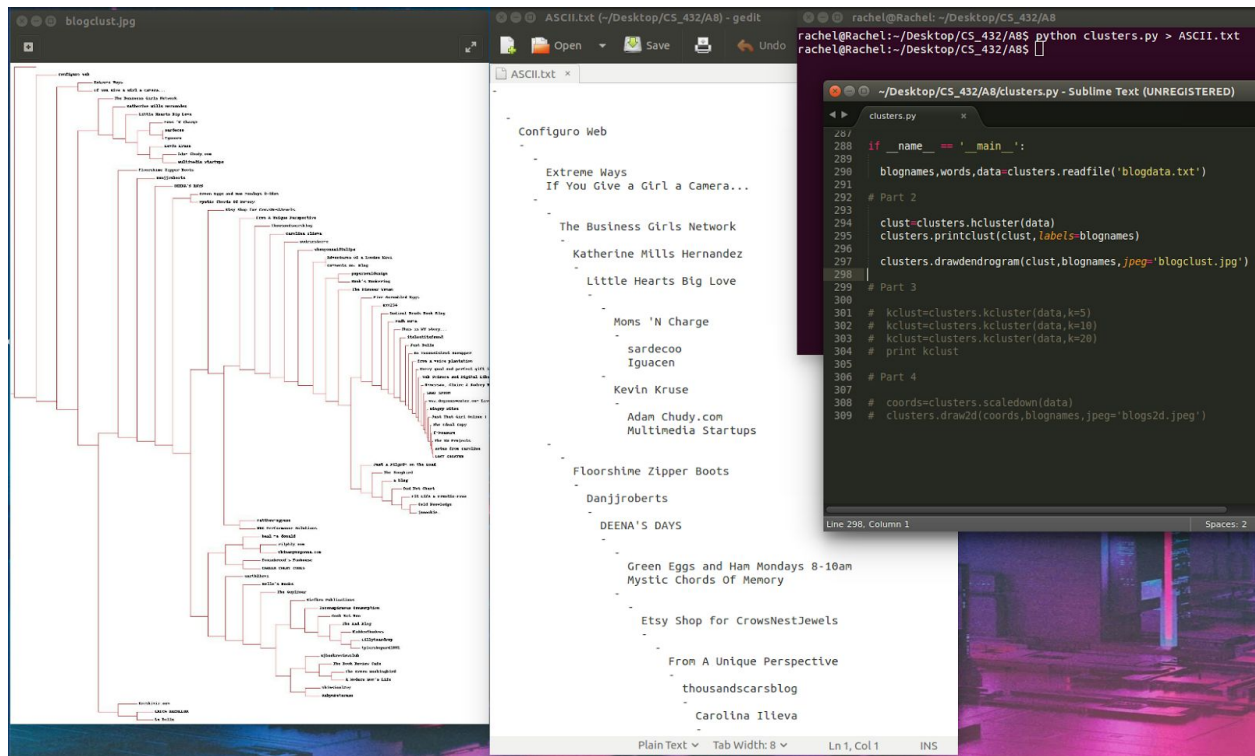
i = 1
for word, count in sorted(temp.items(), key=lambda x: x[1], reverse=True):
    if i <= 500:
        most_frequent_words.append(word)
        i = i + 1
    else: break

out_file('blogdata.txt', 'w')
out.write('blog')
for word in most_frequent_words:
    out.write('\t%s' % word)
    out.write('\n')
```

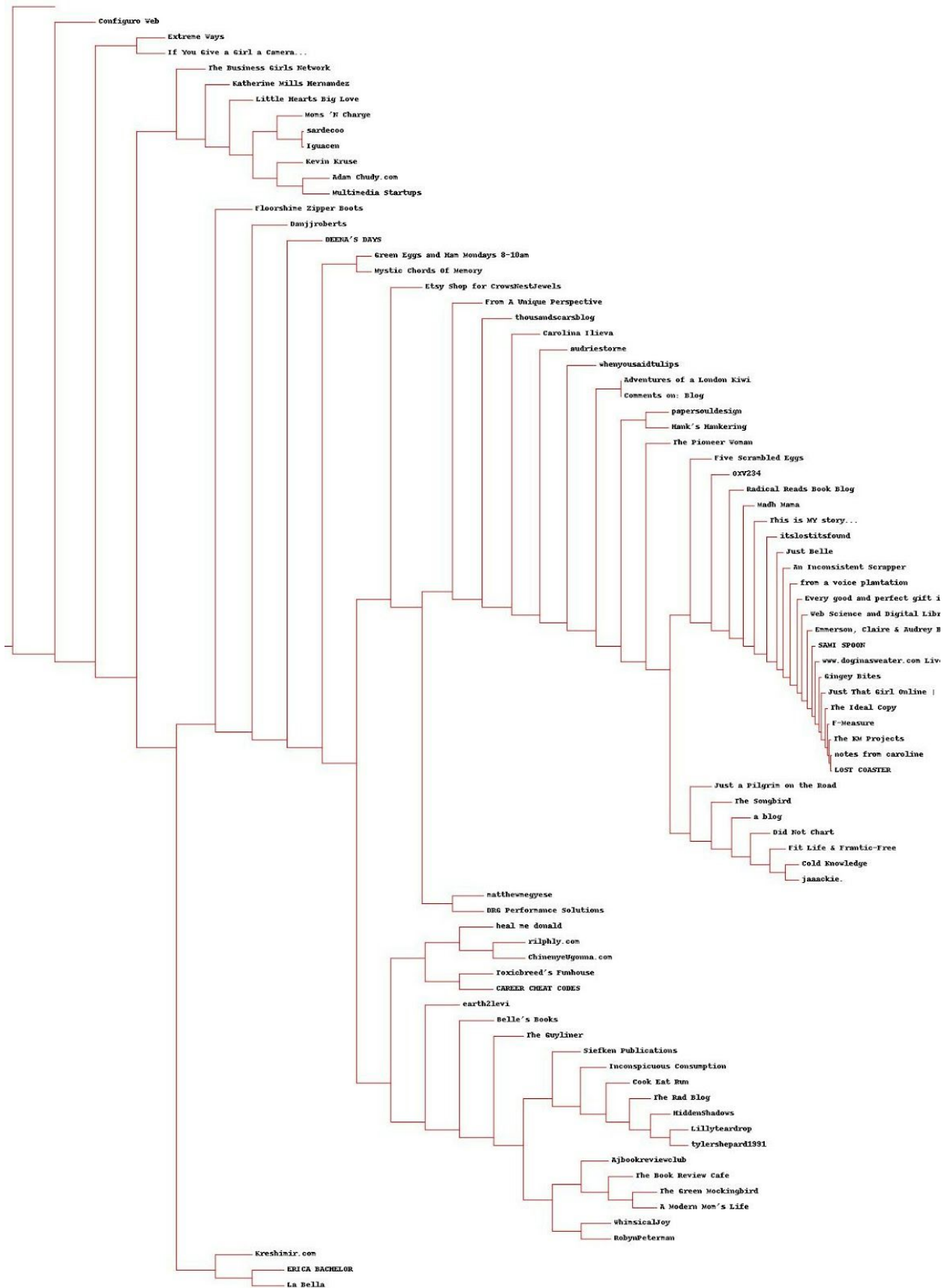
Note: I had to run it again to fix the formatting, so the `blogdata.txt` shown here isn't going to look the same as the one I submit.

PART 2

I used the code found in clusters.py from Programming Collective Intelligence [here](#) to print out a cluster of the blogs in both ASCII and jpeg formats:



Closer look at blogclust.jpeg on the next page:



PART 3

Continuing with the same code found in clusters.py from Programming Collective Intelligence [here](#) to cluster the blogs using K-means.

K = 5 (4 iterations):

```
clusters.py x
285
286 if __name__ == '__main__':
287     blognames, words, data = clusters.readfile('blogdata.txt')
288
289     # Part 2
290
291     # clust = clusters.hcluster(data)
292     # clusters.printclust(clust, labels=blognames)
293     # clusters.drawdendrogram(clust, blognames, jpeg='blogclust.jpg')
294
295     # Part 3
296
297     kclust = clusters.kcluster(data, k=5)
298     # kclust = clusters.kcluster(data, k=10)
299     # kclust = clusters.kcluster(data, k=20)
300     print kclust
301
302     # Part 4
303
304     # coords = clusters.scaledown(data)
305     # clusters.draw2d(coords, blognames, jpeg='blogs2d.jpeg')
```

```
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
[[0, 1, 2, 3, 5, 6, 8, 11, 12, 13, 14, 16, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27,
28, 29, 30, 31, 32, 34, 35, 36, 37, 38, 40, 41, 42, 43, 44, 45, 46, 47, 49, 51, 52,
53, 55, 57, 58, 59, 60, 62, 63, 65, 66, 69, 70, 72, 73, 74, 76, 77, 78, 79, 80, 81],
[7, 9, 10, 15, 33, 54, 56, 61, 68], [67], [4, 39, 71], [17, 48, 50, 64, 75]]
[Finished in 0.9s]
```

K = 10 (3 iterations):

```
clusters.py x
285
286 if __name__ == '__main__':
287
288     blognames, words, data = clusters.readfile('blogdata.txt')
289
290     # Part 2
291
292     # clust = clusters.hcluster(data)
293     # clusters.printclust(clust, labels=blognames)
294
295     # clusters.drawdendrogram(clust, blognames, jpeg='blogclust.jpg')
296
297     # Part 3
298
299     # kclust = clusters.kcluster(data, k=5)
300     | kclust = clusters.kcluster(data, k=10)
301     # kclust = clusters.kcluster(data, k=20)
302     | print kclust
303
304     # Part 4
305
306     # coords = clusters.scaledown(data)
307     # clusters.draw2d(coords, blognames, jpeg='blogs2d.jpeg')
```

```
Iteration 0
Iteration 1
Iteration 2
Iteration 3
[[40, 48, 64], [0, 2, 3, 5, 8, 11, 12, 13, 14, 16, 18, 19, 20, 21, 23, 26, 27, 28,
29, 30, 31, 32, 36, 37, 42, 43, 45, 46, 47, 49, 51, 52, 53, 55, 57, 60, 62, 63, 66,
70, 72, 73, 74, 76, 78, 79, 80, 81], [], [], [1, 6, 17, 22, 24, 25, 35, 38, 41, 44,
58, 59, 65, 69, 77], [7, 9, 34, 61], [4, 39, 71], [], [50, 54, 68, 75], [10, 15, 33,
56, 67]]
[Finished in 1.4s]
```

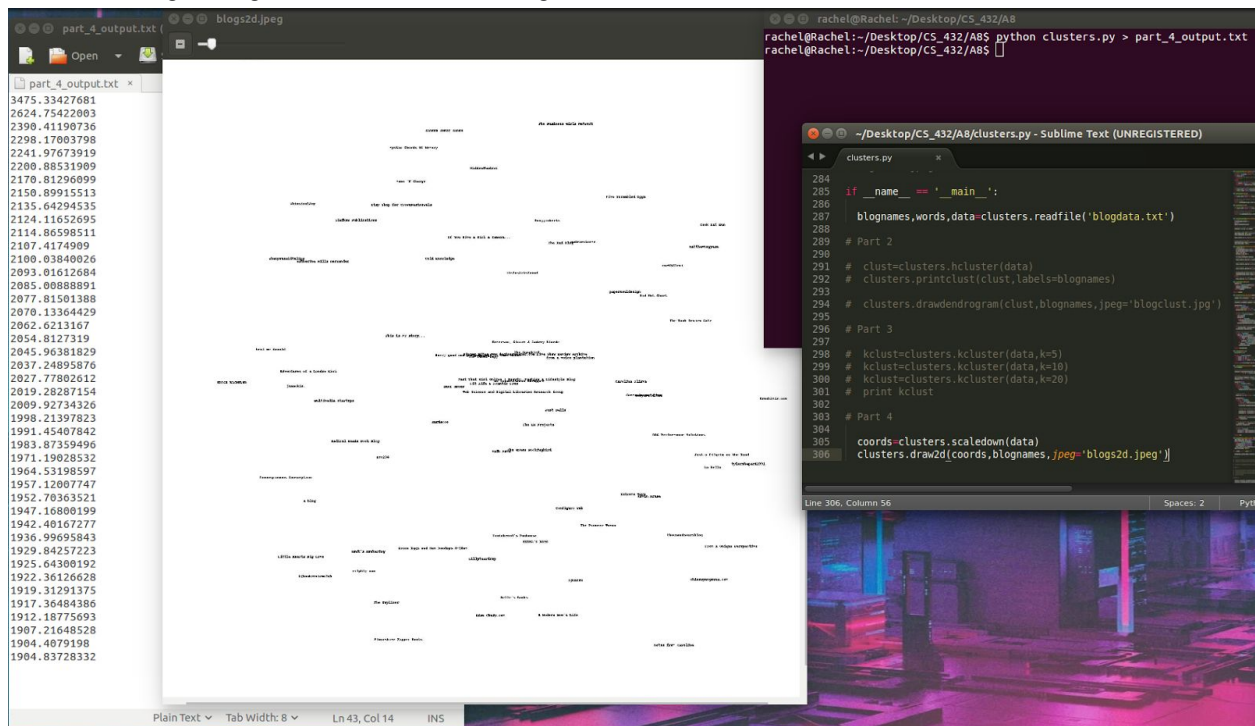
K = 20 (6 iterations):

```
clusters.py x
285
286 if __name__ == '__main__':
287
288     blognames, words, data = clusters.readfile('blogdata.txt')
289
290     # Part 2
291
292     # clust = clusters.hcluster(data)
293     # clusters.printclust(clust, labels=blognames)
294
295     # clusters.drawdendrogram(clust, blognames, jpeg='blogclust.jpg')
296
297     # Part 3
298
299     # kclust = clusters.kcluster(data, k=5)
300     # kclust = clusters.kcluster(data, k=10)
301     | kclust = clusters.kcluster(data, k=20)
302     | print kclust
303
304     # Part 4
305
306     # coords = clusters.scaledown(data)
307     # clusters.draw2d(coords, blognames, jpeg='blogs2d.jpeg')
```

```
Iteration 0
Iteration 1
Iteration 2
Iteration 3
Iteration 4
Iteration 5
Iteration 6
[[7, 40], [], [1, 6, 38, 41, 59, 69], [33, 56, 68], [39], [5, 9, 20, 21, 22, 25, 29,
34, 43, 52, 60, 63, 70], [], [], [0, 2, 3, 8, 13, 14, 16, 18, 19, 23, 24, 26, 27, 30,
31, 32, 35, 36, 37, 42, 44, 45, 46, 47, 49, 51, 55, 58, 62, 65, 66, 72, 73, 74, 76,
78, 79, 80], [4], [11, 15, 28, 53, 77, 81], [71], [10, 17, 48, 50, 54], [], [75],
[12, 57, 64], [], [67], [61], []]
[Finished in 4.9s]
```

PART 4

Still continuing with the same code found in clusters.py from Programming Collective Intelligence [here](#) to cluster the blogs using multidimensional scaling:



Closer look at blogs2d.jpeg on the next page:

The rest of the code from clusters.py used for parts 2, 3, and 4 (I put these at the end in case you wanted to get to the other stuff first):

```
clusters.py x
1 from PIL import Image, ImageDraw
2 import clusters
3
4 def readfile(filename):
5     lines=[line for line in file(filename)]
6
7     # First line is the column titles
8     colnames=lines[0].strip().split('\t')[1:]
9     rownames=[]
10    data=[]
11    for line in lines[1:]:
12        p=line.strip().split('\t')
13        # First column in each row is the rowname
14        rownames.append(p[0])
15        # The data for this row is the remainder of the row
16        data.append([float(x) for x in p[1:]])
17    return rownames, colnames, data
18
19
20 from math import sqrt
21
22 def pearson(v1, v2):
23     # Something weird's going on with my data, this is just a workaround
24     if len(v1) == 0 or len(v2) == 0:
25         return 1.0
26     # Simple sums
27     sum1=sum(v1)
28     sum2=sum(v2)
29
30     # Sums of the squares
31     sum1Sq=sum([pow(v,2) for v in v1])
32     sum2Sq=sum([pow(v,2) for v in v2])
33
34     # Sum of the products
35     # print ('V1 len = ', len(v1))
36     # print ('V2 len = ', len(v2))
37     pSum=sum([v1[i]*v2[i] for i in range(len(v1))])
38
39     # Calculate r (Pearson score)
40     num=pSum-(sum1*sum2/len(v1))
41     den=sqrt((sum1Sq-pow(sum1,2)/len(v1))*(sum2Sq-pow(sum2,2)/len(v1)))
42     if den==0: return 0
43
44     return 1.0-num/den
45
46 class bicluster:
47     def __init__(self, vec, left=None, right=None, distance=0.0, id=None):
48         self.left=left
49         self.right=right
50         self.vec=vec
51         self.id=id
52         self.distance=distance
53
54 def hcluster(rows, distance=pearson):
55     distances={}
56     currentclustid=-1
57
58     # Clusters are initially just the rows
59     clust=[bicluster(rows[i], id=i) for i in range(len(rows))]
```



```

58 # Clusters are initially just the rows
59 clust=[bicluster(rows[i],id=i) for i in range(len(rows))]
60
61 while len(clust)>1:
62     lowestpair=(0,1)
63     closest=distance(clust[0].vec,clust[1].vec)
64
65     # loop through every pair looking for the smallest distance
66     for i in range(len(clust)):
67         for j in range(i+1,len(clust)):
68             # distances is the cache of distance calculations
69             if (clust[i].id,clust[j].id) not in distances:
70                 distances[(clust[i].id,clust[j].id)]=distance(clust[i].vec,clust[j].vec)
71
72             d=distances[(clust[i].id,clust[j].id)]
73
74             if d<closest:
75                 closest=d
76                 lowestpair=(i,j)
77
78     # calculate the average of the two clusters
79     mergevec=[
80         (clust[lowestpair[0]].vec[i]+clust[lowestpair[1]].vec[i])/2.0
81         for i in range(len(clust[0].vec))]
82
83     # create the new cluster
84     newcluster=bicluster(mergevec,left=clust[lowestpair[0]],
85                          right=clust[lowestpair[1]],
86                          distance=closest,id=currentclustid)
87
88     # cluster ids that weren't in the original set are negative
89     currentclustid-=1
90     del clust[lowestpair[1]]
91     del clust[lowestpair[0]]
92     clust.append(newcluster)
93
94     return clust[0]
95
96 def printclust(clust,labels=None,n=0):
97     # indent to make a hierarchy layout
98     for i in range(n): print ' ',
99     if clust.id<0:
100         # negative id means that this is branch
101         print '--'
102     else:
103         # positive id means that this is an endpoint
104         if labels==None: print clust.id
105         else: print labels[clust.id]
106
107     # now print the right and left branches
108     if clust.left!=None: printclust(clust.left,labels=labels,n=n+1)
109     if clust.right!=None: printclust(clust.right,labels=labels,n=n+1)
110
111 def getheight(clust):
112     # Is this an endpoint? Then the height is just 1
113     if clust.left==None and clust.right==None: return 1
114
115     # Otherwise the height is the same of the heights of

```

```

115 # Otherwise the height is the same of the heights of
116 # each branch
117 return getheight(clust.left)+getheight(clust.right)
118
119 def getdepth(clust):
120     # The distance of an endpoint is 0.0
121     if clust.left==None and clust.right==None: return 0
122
123     # The distance of a branch is the greater of its two sides
124     # plus its own distance
125     return max(getdepth(clust.left),getdepth(clust.right))+clust.distance
126
127
128 def drawdendrogram(clust,labels,jpeg='clusters.jpg'):
129     # height and width
130     h=getheight(clust)*20
131     w=1200
132     depth=getdepth(clust)
133
134     # width is fixed, so scale distances accordingly
135     scaling=float(w-150)/depth
136
137     # Create a new image with a white background
138     img=Image.new('RGB',(w,h),(255,255,255))
139     draw=ImageDraw.Draw(img)
140
141     draw.line((0,h/2,10,h/2),fill=(255,0,0))
142
143     # Draw the first node
144     drawnode(draw,clust,10,(h/2),scaling,labels)
145     img.save(jpeg,'JPEG')
146
147 def drawnode(draw,clust,x,y,scaling,labels):
148     if clust.id<0:
149         h1=getheight(clust.left)*20
150         h2=getheight(clust.right)*20
151         top=y-(h1+h2)/2
152         bottom=y+(h1+h2)/2
153         # Line length
154         ll=clust.distance*scaling
155         # Vertical line from this cluster to children
156         draw.line((x,top+h1/2,x,bottom-h2/2),fill=(255,0,0))
157
158         # Horizontal line to left item
159         draw.line((x,top+h1/2,x+ll,top+h1/2),fill=(255,0,0))
160
161         # Horizontal line to right item
162         draw.line((x,bottom-h2/2,x+ll,bottom-h2/2),fill=(255,0,0))
163
164         # Call the function to draw the left and right nodes
165         drawnode(draw,clust.left,x+ll,top+h1/2,scaling,labels)
166         drawnode(draw,clust.right,x+ll,bottom-h2/2,scaling,labels)
167     else:
168         # If this is an endpoint, draw the item label
169         draw.text((x+5,y-7),labels[clust.id],(0,0,0))
170
171 def rotatematrix(data):
172     newdata=[]

```



```

171 def rotatematrix(data):
172     newdata=[]
173     for i in range(len(data[0])):
174         newrow=[data[j][i] for j in range(len(data))]
175         newdata.append(newrow)
176     return newdata
177
178 import random
179
180 def kcluster(rows,distance=pearson,k=4):
181     # Determine the minimum and maximum values for each point
182     ranges=[(min([row[i] for row in rows]),max([row[i] for row in rows]))
183             for i in range(len(rows[0]))]
184
185     # Create k randomly placed centroids
186     clusters=[[random.random()*(ranges[i][1]-ranges[i][0])+ranges[i][0]
187               for i in range(len(rows[0]))] for j in range(k)]
188
189     lastmatches=None
190     for t in range(100):
191         print 'Iteration %d' % t
192         bestmatches=[] for i in range(k)]
193
194         # Find which centroid is the closest for each row
195         for j in range(len(rows)):
196             row=rows[j]
197             bestmatch=0
198             for i in range(k):
199                 d=distance(clusters[i],row)
200                 if d<distance(clusters[bestmatch],row): bestmatch=i
201             bestmatches[bestmatch].append(j)
202
203         # If the results are the same as last time, this is complete
204         if bestmatches==lastmatches: break
205         lastmatches=bestmatches
206
207         # Move the centroids to the average of their members
208         for i in range(k):
209             avgs=[0.0]*len(rows[0])
210             if len(bestmatches[i])>0:
211                 for rowid in bestmatches[i]:
212                     for m in range(len(rows[rowid])):
213                         avgs[m]+=rows[rowid][m]
214                 for j in range(len(avgs)):
215                     avgs[j]/=len(bestmatches[i])
216                 clusters[i]=avgs
217
218     return bestmatches
219
220 def tanamoto(v1,v2):
221     c1,c2,shr=0,0,0
222
223     for i in range(len(v1)):
224         if v1[i]!=0: c1+=1 # in v1
225         if v2[i]!=0: c2+=1 # in v2
226         if v1[i]!=0 and v2[i]!=0: shr+=1 # in both
227
228     return 1.0-(float(shr)/(c1+c2-shr))
229

```

```

230 def scaledown(data,distance=pearson,rate=0.01):
231     n=len(data)
232
233     # The real distances between every pair of items
234     realdist=[[distance(data[i],data[j]) for j in range(n)]
235              for i in range(0,n)]
236
237     # Randomly initialize the starting points of the locations in 2D
238     loc=[[random.random(),random.random()] for i in range(n)]
239     fakedist=[[0.0 for j in range(n)] for i in range(n)]
240
241     lasterror=None
242     for m in range(0,1000):
243         # Find projected distances
244         for i in range(n):
245             for j in range(n):
246                 fakedist[i][j]=sqrt(sum([pow(loc[i][x]-loc[j][x],2)
247                                         for x in range(len(loc[i]))]))
248
249         # Move points
250         grad=[[0.0,0.0] for i in range(n)]
251
252         totalerror=0
253         for k in range(n):
254             for j in range(n):
255                 if j==k: continue
256                 # The error is percent difference between the distances
257                 errorterm=(fakedist[j][k]-realdist[j][k])/realdist[j][k]
258
259                 # Each point needs to be moved away from or towards the other
260                 # point in proportion to how much error it has
261                 grad[k][0]+=((loc[k][0]-loc[j][0])/fakedist[j][k])*errorterm
262                 grad[k][1]+=((loc[k][1]-loc[j][1])/fakedist[j][k])*errorterm
263
264                 # Keep track of the total error
265                 totalerror+=abs(errorterm)
266         print totalerror
267
268         # If the answer got worse by moving the points, we are done
269         if lasterror and lasterror<totalerror: break
270         lasterror=totalerror
271
272         # Move each of the points by the learning rate times the gradient
273         for k in range(n):
274             loc[k][0]-=rate*grad[k][0]
275             loc[k][1]-=rate*grad[k][1]
276
277     return loc
278
279 def draw2d(data,labels,jpeg='mds2d.jpg'):
280     img=Image.new('RGB',(2000,2000),(255,255,255))
281     draw=ImageDraw.Draw(img)
282     for i in range(len(data)):
283         x=(data[i][0]+0.5)*1000
284         y=(data[i][1]+0.5)*1000
285         draw.text((x,y),labels[i],(0,0,0))
286     img.save(jpeg,'JPEG')
287
288

```