# CS595 Assignment 9

## Jon Robison

### November 24, 2013

Q1.
Create a blog-term matrix. Start by grabbing 100 blogs.

After writing a program to grab some given number of blogs,
the given generatefeedvector provided the matrix. Another program,
generateMatrix.py, was written to accomplish the same task, however,
it was noted on the slides in a sneaky location that code was provided
accomplishing the same thing. Much frustration at self ensued.
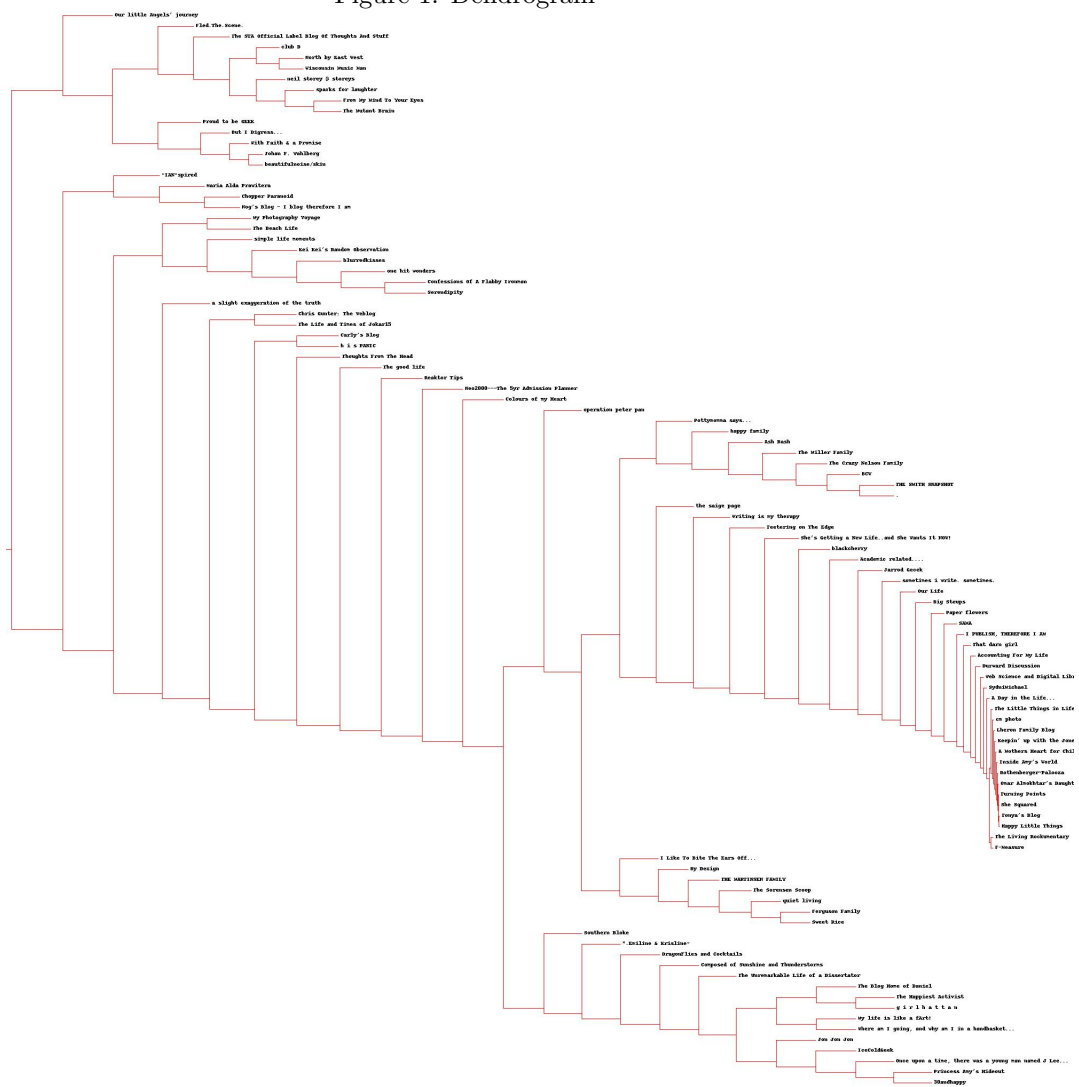See Appendix A for program to generate bloglist, generateUrls.py

Q2.
Create an ASCII and JPEG dendrogram that clusters (i.e., HAC)
the most similar blogs.

The first step is always making sure this is believable to the person
assigning grades, and sure enough, F-Measure is beside The Living
Rockumentary. My grade for this assignment can go up to 11, you know.

Similar blogs are grouped together, most notably, technology blogs and meta-
type blogs. Meta type blogs as used here are blogs about either blogging,
commentary (on life, ideals and culture), the meaninglessness of commentary,
or the meaninglessness of blogging. I found these to be most interesting, as
often the bloggers portray themselves as outsiders looking in to a culture
they appreciate, don't understand, don't connect with, are mad at, however,
their proximity in the dendrogram indicates they are at least similar to each
other. This implies there is a large subculture of people thinking they are
unique as an inherently positive attribute. Notable that it is (almost) the
majority group of this dataset, ultra large sample size that it is
(/hyperbole).

Figure 1: Dendrogram

Our little Angels' journey
Fled The Scene.
The SVA Official Label Blog of Thoughts And Stuff
club D
North by East West
Wisconsin Music Nun
neil storey D storeys
sparks for laughter
Free My Mind To Your Eyes
The Mutant Brain
Proud to be GEEK
But I Digress...
With Faith & a Promise
Johnn F. Vahlberg
beautifulnoise/skin
'IAN'spired
Maria Alda Prewitera
Chopper Paranoid
Moy's Blog - I blog therefore I am
My Photography Voyage
The Beach Life
simple life moments
Kei Kei's Random Observation
blurredkisses
one hit wonders
Confessions Of A Flabby Ironman
Serendipity
a slight exaggeration of the truth
Chris Gunter: The Weblog
The Life and Times of JokariS
Curly's Blog
h i s PANIC
Thoughts From The Head
The good life
Reaktor Tips
Neo2000---The Syr Admission Planner
Colours of my Heart
operation peter pan
Pattymama says...
happy family
Ash Bash
The Willer Family
The Crazy Nelson Family
BCV
THE SMITH SNAPSHOT
.
the snige page
writing is my therapy
Teetering on The Edge
She's Getting a New Life...and She Wants It NOW!
blackcherry
Academic related....
Jarred &cook
sometimes i write. sometimes.
Our Life
Big Steps
Paper Flowers
SAVA
I PUBLISH, THEREFORE I AM
That dare girl
Accounting For My Life
Durward Discussion
Web Science and Digital Libr
Syd&Michael
A Day in the Life...
The Little Things in Life
cm photo
Lheren Family Blog
keepin' up with the Jone
A Mothers Heart for Chil
Inside Joy's World
Bathenberger-Palozza
Omar Abubhtar's Daught
Turning Points
the squared
Joopy's Blog
Happy Little Things
The Living Documentary
F-measure
I Like To Bite The Ears Off...
By Design
THE MARTINSEN FAMILY
The Sorensen Scoop
quiet living
Ferguson Family
Sweet Rice
Southern Blake
"...Emiline & Krisline"
Dragonflies and Cocktails
Composed of Sunshine and Thunderstorms
The Unremarkable Life of a Dissertator
The Blog Home of Daniel
The Happiest Activist
g i r l h a t t a n
My life is like a fart!
Where am I going, and why am I in a handbasket...
Jon Jon Jon
IceColdGeek
Once upon a time, there was a young man named J Lee...
Princess Amy's Hideout
BOandhappy

Q3.

Cluster the blogs using K-Means, using k=5,10,20. (see slide 18).
How many interations were required for each value of k?
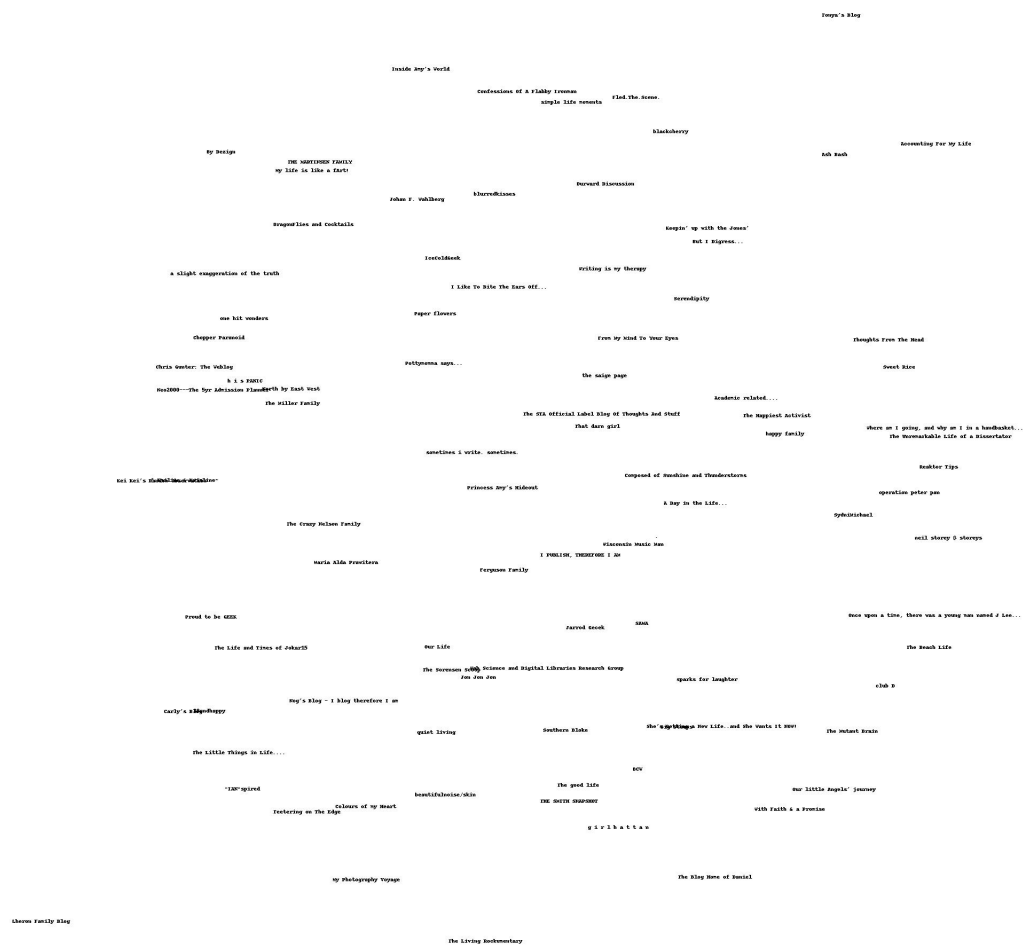
| K | Iterations |
|---|---|
| 5 | 5 |
| 10 | 7 |
| 20 | 7 |

Q4.

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

Only one iteration before error started increasing by a factor of ten.
TODO This is odd and I mean to go back to see why, haven't yet.

Figure 2: Clustered Blogs

Tonya's Blog

Inside Amy's World

Confessions Of A Flabby Ironman          Fled.The.Scene.
                    simple life moments

                              blackcherry

By Design                                                    Accounting For My Life
                                                    Ash Rash
THE MARTINSEN FAMILY
my life is like a fart!

                              Forward Discussion
John F. Wahlberg       blurredlines

Dragonflies and Cocktails              Keepin' up with the Jones'
                                       But I Digress...

                    icecoldmeek
a slight exaggeration of the truth     Writing is my therapy
                    I Like To Bite The Ears Off....
                                       Serendipity
one hit wonders          Paper flowers

Chopper Paranoid                        From My Mind To Your Eyes          Thoughts From The Head

Chris Gunter: The Weblog    Pettymomma says...                                      Sweet Rice
              h i s PANIC                        the saige page
Moo2000----The Dyr Admission PlannerEarth by East West
          The Miller Family                                Academic related....
                              The STA Official Label Blog Of Thoughts And Stuff    the Happiest Activist
                              that dare girl                                where am I going, and why am I in a handbasket...
                                                    happy family    The Unremarkable Life of a Dissertator

                    sometimes i write. sometimes.                                    Beakter Tips
Kei Kei's Ramblin' ...              Composed of Sunshine and Thunderstorms
                    Princess Amy's Hideout                                operation peter pan
                                       A Day In the Life...
                                                    SydniMichael
The Crazy Nelsen Family                                                   neil storey D storeys
                              Wisconsin Music Man
                    I PUBLISH, THEREFORE I AM
maria Alda Prascitora
              Ferguson Family

Proud to be GEEK                                                Once upon a time, there was a young man named J Lee...
                              Jarred Gecek    SAMA
The Life and Times of Jokari5    Our Life                                    The Beach Life
                    The Sorensen ...UQ Science and Digital Libraries Research Group
                              Jen Jen Jen                    sparks for laughter          club D
              Amy's Blog - I blog therefore I am
Carly's z... -Puppy
                    quiet living          Southern Bloke    She's Got a New Life...and she wants it NOW!    The Mutant Brain
The Little Things In Life....
                                                    DCV
"IAN"spired                                  The good life                        Our little Angela' journey
                    beautifulnoise/skin    THE SMITH SNAPSHOT
          Colours of my Heart                                      With Faith & a Promise
Festering on The Edge
                                       g i r l b a t t a n

My Photography Voyage                            The Blog Home of Daniel

Lheros Family Blog

                    The Living Rockumentary

Appendix A

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

DEFAULT_COUNT=98
DEFAULT_SEED_URL='http://www.blogger.com/next-blog?navBar=true&blogID=347163
if len(sys.argv) != 3:
    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:
    count=sys.argv[1]
    url=sys.argv[2]

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

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

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

5