

Assignment Four

Jacob Berlin

CS432 – Spring 2016

1. 1. Determine if the friendship paradox holds for my Facebook account.* Compute the mean, standard deviation, and median of the number of friends that my friends have. Create a graph of the number of friends (y-axis) and the friends themselves, sorted by number of friends (x-axis). (The friends don't need to be labeled on the x-axis: just f1, f2, f3, ... fn.) Do include me in the graph and label me accordingly.

```
<?xml version="1.0" encoding="UTF-8"?>
<data>
  <node id="Simeon_Warner_428351">
    <data key="Label">Simeon Warner</data>
    <data key="uid"><![CDATA[428351]]></data>
    <data key="name"><![CDATA[Simeon Warner]]></data>
    <data key="mutual_friend_count"><![CDATA[13]]></data>
    <data key="friend_count"><![CDATA[244]]></data>
  </node>
  <node id="Drew_Munro_1314586">
    <data key="Label">Drew Munro</data>
    <data key="uid"><![CDATA[1314586]]></data>
    <data key="name"><![CDATA[Drew Munro]]></data>
    <data key="mutual_friend_count"><![CDATA[17]]></data>
    <data key="friend_count"><![CDATA[575]]></data>
  </node>
  <node id="Mat_Kelly_2004483">
    <data key="Label">Mat Kelly</data>
    <data key="uid"><![CDATA[2004483]]></data>
    <data key="name"><![CDATA[Mat Kelly]]></data>
    <data key="mutual_friend_count"><![CDATA[12]]></data>
    <data key="friend_count"><![CDATA[421]]></data>
  </node>
  <node id="Benjamin_Lok_2037943">
    <data key="Label">Benjamin Lok</data>
    <data key="uid"><![CDATA[2037943]]></data>
    <data key="name"><![CDATA[Benjamin Lok]]></data>
    <data key="mutual_friend_count"><![CDATA[1]]></data>
    <data key="friend_count"><![CDATA[539]]></data>
  </node>
  <node id="Camden_Elliott_Matherne_2726573">
    <data key="Label">Camden Elliott Matherne</data>
    <data key="uid"><![CDATA[2726573]]></data>
    <data key="name"><![CDATA[Camden Elliott Matherne]]></data>
    <data key="mutual_friend_count"><![CDATA[8]]></data>
    <data key="friend_count"><![CDATA[784]]></data>
  </node>
  <node id="Benjamin_Dunn_Munro_2726573">
```

Initially for this problem I started with the raw XML format that was retrieved from the 'mln.graphml' file and then created a python program to parse the data.

```

1  -*- coding: utf-8 -*-
2
3  #URIFile = open('uriFile.txt', 'r')
4
5  import xml.etree.ElementTree as ET
6
7  unsortedFile = open('friendNums.txt', 'a')
8  outputFile = open('friendsFB.txt', 'a')
9
10 outputCounter = 1
11 counter = 0
12 tree = ET.parse("rawFile.xml")
13 root = tree.getroot()
14
15 print "MLN          154"
16 outputFile.write("MLN          154")
17 outputFile.write("\n")
18
19 #for initial in root.iter('node'):
20 #    unsortedFile.write(root[counter][4].text)
21 #    unsortedFile.write("\n")
22 #    counter += 1
23
24 with open('friendNums.txt') as f:
25     lines = f.readlines()
26     lines = [x.strip('\n') for x in lines]
27
28 for Nums in lines:
29     outputFile.write("F" + str(outputCounter) + "          " + Nums)
30     outputFile.write("\n")
31     outputCounter += 1

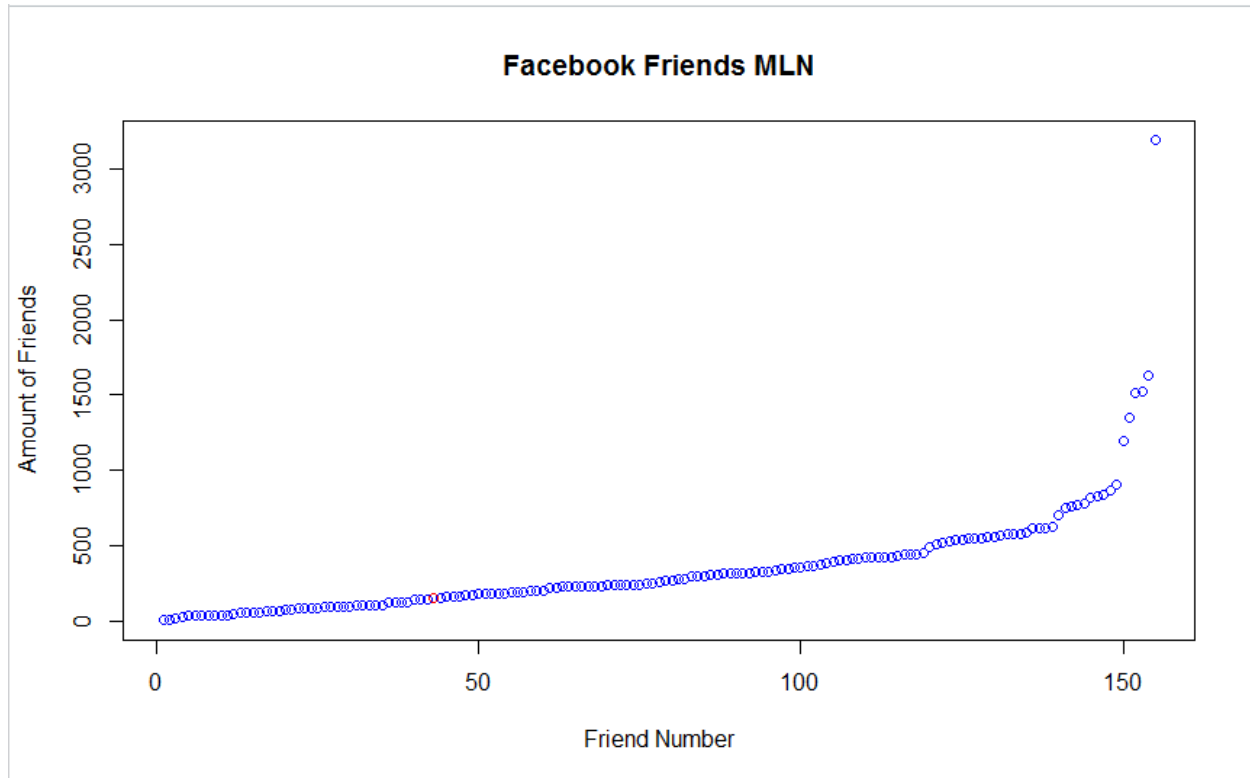
```

This file exclusively grabbed the amount of friends that each of the nodes in 'rawFile.xml' had and added it to a file called 'friendNums.txt.' Once I had retrieved all of the friends and the number of friends that each of them have, I sorted the list through the command-line using the sort -n command. After retrieving all of the data required, I parsed through the data and created a file with all of the friends along with adding the 'MLN' value to show you in the data.

\$ cat friendNums.txt sort -n	1	Friend	FriendAmount	18	F17	65
7	2	F1	7	19	F18	68
15	3	F2	15	20	F19	68
25	4	F3	25	21	F20	77
30	5	F4	30	22	F21	80
38	6	F5	38	23	F22	85
39	7	F6	39	24	F23	86
40	8	F7	40	25	F24	87
41	9	F8	41	26	F25	89
41	10	F9	41	27	F26	93
42	11	F10	42	28	F27	94
43	12	F11	43	29	F28	94
54	13	F12	54	30	F29	96
58	14	F13	58	31	F30	97
59	15	F14	59	32	F31	104
60	16	F15	60	33	F32	104
62	17	F16	62	34	F33	106
65	18	F17	65	35	F34	108
68	19	F18	68	36	F35	111
68	20	F19	68	37	F36	123
77	21	F20	77	38	F37	124
80	22	F21	80	39	F38	128
85	23	F22	85	40	F39	131
86	24	F23	86	41	F40	143
87	25	F24	87	42	F41	144
89	26	F25	89	43	F42	147
93	27	F26	93	44	MLN	154
94	28	F27	94	45	F43	155
94	29	F28	94	46	F44	165
96	30	F29	96	47	F45	168
97	31	F30	97	48	F46	168
104	32	F31	104	49	F47	170
104	33	F32	104	50	F48	172
106	34	F33	106	51	F49	181
108	35	F34	108	52	F50	182
111	36	F35	111	53	F51	183
123	37	F36	123	54	F52	186
124				55	F53	187
128						
131						
143						
144						
147						
155						
165						
168						
168						
170						
172						
181						
182						
183						
186						
187						
190						
195						
197						
204						
207						
208						
220						
227						
229						
231						
231						

Finally, after retrieving all of this data and formatting it as needed I went to RStudio and put the file into a plot generator. The graph I made shows each of the friends' amount of friends alongside your own value from least to greatest. The 'MLN' value is shown in red while all of the other values are shown in blue.

```
> plot(friendsFB$FriendAmount, col=ifelse(friendsFB$FriendAmount==154, 'red', 'blue'), xlab='Friend Number', ylab='Amount of Friends', main='Facebook Friends MLN')
```



To calculate the standard deviation, mean, and median of the list was just as simple as putting the data into R's built in math functions. I made sure to take the 'MLN' value out of the dataset as to not skew the data.

```
> mean(friendsFB$FriendAmount)
[1] 358.987
> median(friendsFB$FriendAmount)
[1] 266.5
>
> sd(friendsFB$FriendAmount)
[1] 371.5853
```

With your total of 154 friends, the friendship paradox holds for your Facebook account due to the fact that the majority of your friends have more friends on Facebook than you do. This is shown by the fact that your friend amount is smaller than the mean, median, and mode of the dataset.

2. Determine if the friendship paradox holds for your Twitter account. Since Twitter is a directed graph, use "followers" as value you measure (i.e., "do your followers have more followers than you?").

Generate the same graph as in question #1, and calculate the same mean, standard deviation, and median values.

```

1  import time
2  import tweepy
3
4  auth = tweepy.OAuthHandler("38wZKZuUwGitHkE3dMpR7jEz", "czTV2ryAOTlep7FPC8dVsaAwS28Cw8Z7L8gDCLnj22ioo0uyuG")
5  auth.set_access_token("2352884547-xheIpcHT0oIjJmzGUKIHwt5X2IzmwogTMh9YWvc", "70RS08peFisvJyPOZGlPleQqfg98twYVkiQefeEP1Ifdg")
6
7  outputFile = open("twitterFollowerAmt.txt", 'w')
8  idFile = open("twitterFollowerIds.txt", 'w')
9
10 counter = 1
11 api = tweepy.API(auth)
12
13 ids = []
14
15 for page in tweepy.Cursor(api.followers_ids, screen_name="JuicyJake1868").pages():
16     ids.extend(page)
17     time.sleep(5)
18
19 print len(ids)
20
21 for name in ids:
22     print name
23     user = api.get_user(name)
24     print user.followers_count
25     outputFile.write(str(counter) + " " + str(user.followers_count))
26     outputFile.write("\n")
27     print name
28     idFile.write(str(name))
29     idFile.write("\n")
30     counter += 1
31
32 idFile.close()
33 outputFile.close()

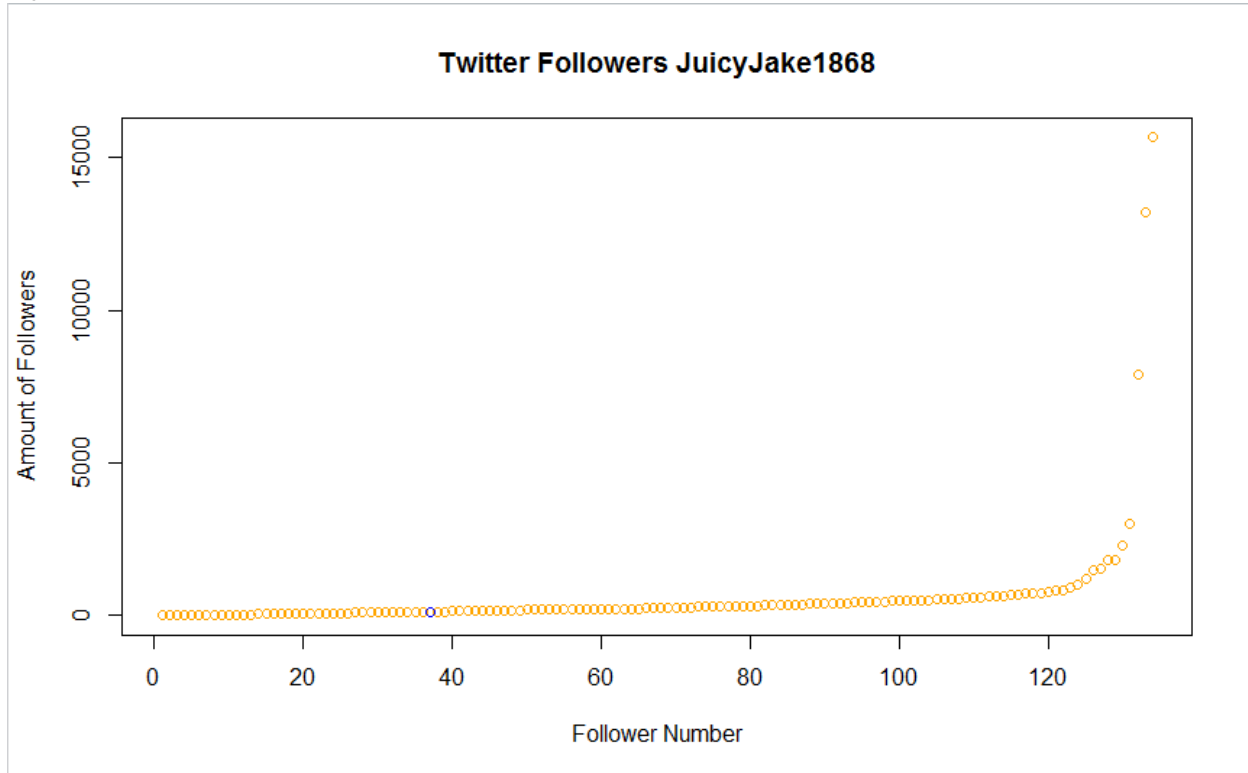
```

For question two, I started in the same fashion of question one where I created a python file to parse through each of my personal Twitter accounts. Using the 'Tweepy' library, I was able to grab a list of both the follower's ids along with the amount of followers that each of my followers have. Only the ids were stored in the 'twitterFollowersIds.txt'. Just the follower amount values were then sorted using the `sort -n` command and then the sorted follower numbers were stored in the 'twitterFollowerAmt.txt'.

\$ python twitterFollowerIds.py		Follower FollowerCount	
133	98	472976560	1 0
460842365	99	1035149784	2 2
376	100	100002618	3 3
460842365	101	318487460	4 6
4664999113	102	413598079	5 6
6	103	467954890	6 8
4664999113	104	243467114	7 14
293183998	105	346964480	8 16
731	106	2304573203	9 19
293183998	107	2230299493	10 23
3064238573	108	85005476	11 32
128	109	416048312	12 36
3064238573	110	2345497260	13 37
1004748558	111	255741404	14 52
295	112	849931044	15 54
1004748558	113	429110769	16 57
4130688135	114	320364459	17 59
6	115	261946574	18 63
4130688135	116	232719227	19 64
4064572889	117	1034930119	20 66
37	118	312427672	21 67
4064572889	119	25005076	22 73
941726670	120	340047972	23 75
210	121	346705910	24 76
941726670	122	794711202	25 77
3922723959	123	1122537024	26 84
67	124	60784841	27 93
3922723959	125	1708408440	28 94
286261482	126	250966853	29 100
1483	127	465066276	30 103
286261482	128	547267877	31 105
286532398	129	438474262	32 107
623	130	468602502	33 123
286532398	131	411211347	34 127
382174209	132	403723130	35 127
330	133	425379468	36 128
382174209			37 134
935921503			
283			
935921503			
277619936			
399			
277619936			
186256987			
833			
186256987			
3468910169			
100			
3468910169			
3382567216			
134			
3382567216			
3103348295			
3			
3103348295			
608698584			
84			
608698584			
1601213750			
141			
1601213750			
543971994			
1836			

Along with in the first problem, I now had the entire list of data needed to put into RStudio to generate the graph. The specific value for my amount of followers is detailed in blue while my follower's followers are detailed in orange.

```
> plot(twitterFollowerAmt$FollowerCount, col=ifelse(twitterFollowerAmt$FollowerCount==133, 'blue', 'orange'), xlab='Follower Number', ylab='Amount of Followers', main='Twitter Followers JuicyJake1868')
```



To compute the mean, median, and standard deviation I plugged the values into R after making sure to take the 'JWB' value out.

```
> mean(twitterFollowerAmt$FollowerCount)
[1] 637.1343
> median(twitterFollowerAmt$FollowerCount)
[1] 259.5
> sd(twitterFollowerAmt$FollowerCount)
[1] 1881.712
```

Therefore, by the rule of the friendship paradox and with my follower count being only at 133, the friendship paradox still holds for my account. This is due to my follower count being significantly less than all three of the calculated values.

3. E.C. Repeat question #1, but with your LinkedIn profile.

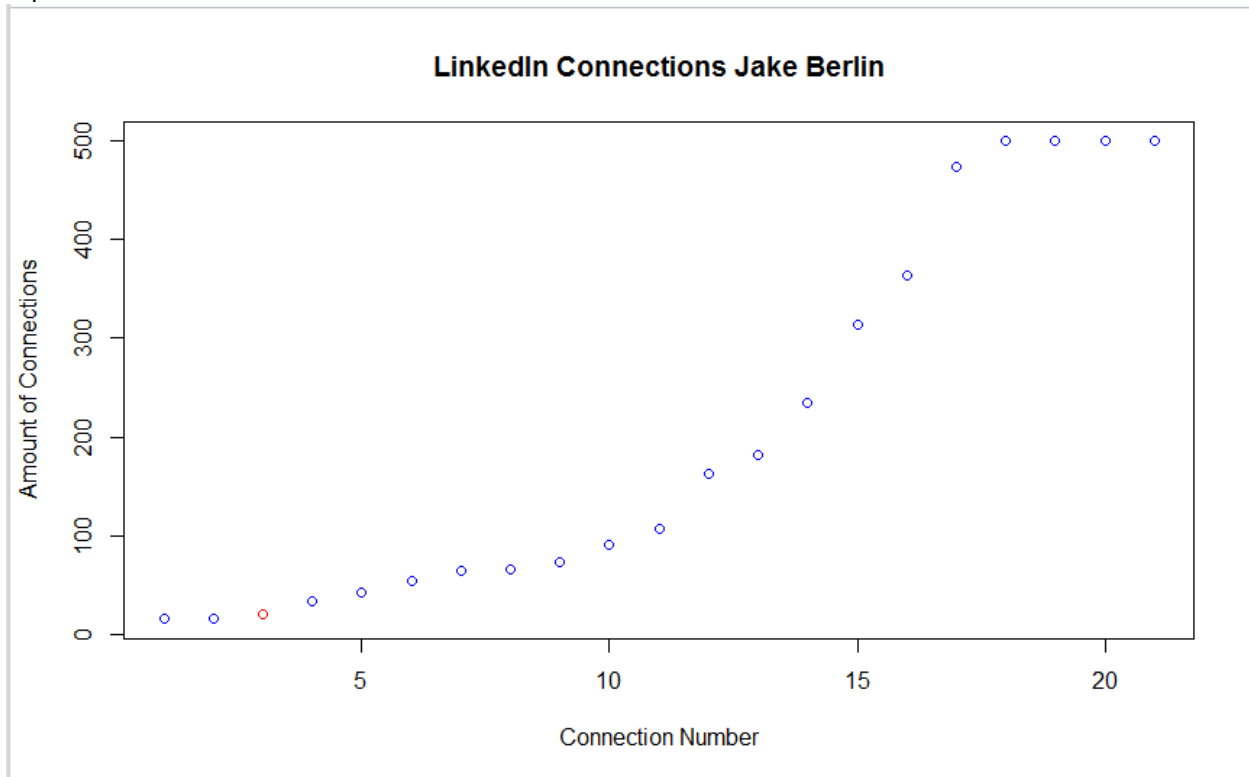
On my LinkedIn profile I only have 20 connections. (Account name is Jake Berlin)

Due to the LinkedIn API not working as of 2/20/2016, I have manually gone through each of my connections and have grabbed all of their data.

Data:

Connection	ConnectionAmt
C1	15
C2	15
JWB	20
C3	33
C4	42
C5	54
C6	64
C7	65
C8	73
C9	90
C10	106
C11	162
C12	181
C13	234
C14	314
C15	364
C16	474
C17	500
C18	500
C19	500
C20	500

```
> plot(LinkedIn$ConnectionAmt, col=ifelse(LinkedIn$ConnectionAmt==20, 'red', 'blue'), xlab='Connection Number', ylab='Amount of Connections', main='LinkedIn Connections Jake Berlin')
```



As in question 1, I created a plot graph based on the amount of connections I had and the connections that each of my connections had. The connection in red is the 'JWB' connection, showing myself alongside the rest of the data.

I then calculated the standard deviation, mean, and median in R after making sure to take out the 'JWB' value.

```
> sd(LinkedIn$ConnectionAmt)
[1] 190.5517
> median(LinkedIn$ConnectionAmt)
[1] 106
> mean(LinkedIn$ConnectionAmt)
[1] 205.0476
```

Finally after retrieving all of the data needed and with only 20 connections the friendship paradox holds for my LinkedIn account due to the fact that my connection amount is marginally lower than the three values I calculated.