# Assignment 4

## CS 532: Introduction to Web Science
### Spring 2016
Manoj Chandra Kompalli
Finished on February 25, 2016

# 1

## Question

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.

* = This used to be more interesting when you could more easily download
your friend's friends data from Facebook.  Facebook now requires each
friend to approve this operation, effectively making it impossible.

I will email to the list the XML file that contains my Facebook
friendship graph ca. Oct, 2013.  The interesting part of the file looks
like this (for 1 friend):

```
<node id="Johan_Bollen_1448621116">
        <data key="Label">Johan Bollen</data>
        <data key="uid"><![CDATA[1448621116]]></data>
        <data key="name"><![CDATA[Johan Bollen]]></data>
        <data key="mutual_friend_count"><![CDATA[37]]></data>
        <data key="friend_count"><![CDATA[420]]></data>
</node>
```

It is in GraphML format: http://graphml.graphdrawing.org/

| | |
|---|---|
| **Mean** | 358.987 |
| **Median** | 266.500 |
| **Std Dev** | 371.585 |

Table 1: Statistics of Dr.Nelson's Friends' friends values taken from R

## Answer

I was trying to figure out how to get my own graphml to fetch my friends data .I could not succeed due to facebooks privacy settings. I have downloaded Dr. Nelsons friend list in graphml format. Now the task ahead of me was to find a way to extract the graphml data and sort it.I researched and found a library called minidom which can extract the xml data. I have extracted the data tag and two key values name and friend count. The "name" got me all the friends names and friend_count shown in Listing 1 gave me the count of friends for that friend.I have iterated the loop for all the friends to get the complete list. I have sorted the list in increasing order using Excel and copied the values to a text file.I had computed Dr. Nelsons friends and added it to the list of his friends.

I used the text file as input to R.I found the mean, median and standard deviation of the sorted list in R.Code shown in Listing 2. I have plotted a bargraph shown in Figure 1 . I had found out that Dr. Nelson has less than half of the mean of his friends friends that is 154 .This proves the friendship paradox.

The R script generates these values:

```
Mean    :358.987
Median :266.500
Std Dev:371.585
```
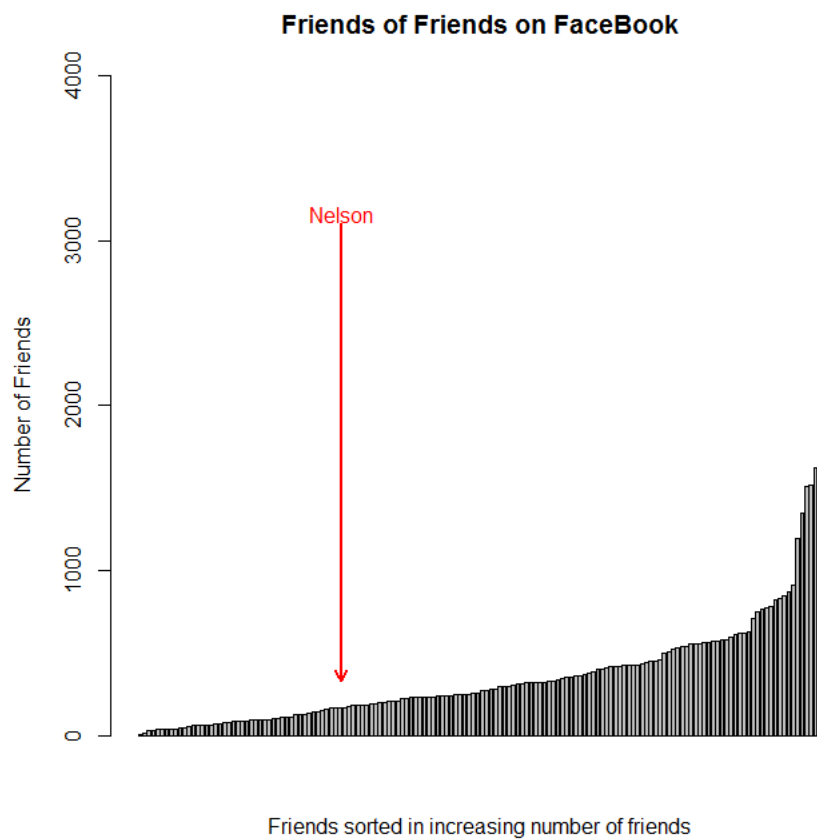
Figure 1: Bar plot showing the count of Dr.Nelson's Facebook Friends'
Friends

```
1   #!/usr/local/bin/python3
2   import csv
3   import sys
4   import numpy
5
6   from xml.dom.minidom import parseString
7
8
9   def friendData(xml):
10
11
12      dom = parseString(xml)
13      totalCount = {}
14
15      for element in dom.getElementsByTagName("data"):
16          if (element.attributes['key'].value == 'name'):
17              name = element.childNodes[0].data
18
19          if (element.attributes['key'].value == 'friend_count'):
20              count = element.childNodes[0].data
21
22
23
24
25              totalCount[name] = count
26
27
28
29              name = ''
30              count = ''
31
32      return totalCount
33
34   def countFriends(xml):
35
36      dom = parseString(xml)
37      return len(dom.getElementsByTagName("node"))
38
39
40   if __name__ == "__main__":
41          sum=0
42          counter=0
43
44
45          graphmlFile = sys.argv[1]
46          f = open(graphmlFile)
47          f2=open('output.csv','w')
48          xml = f.read()
```

4

```
49          f.close()
50          countProfFriend = countFriends(xml)
51          fbFriends = friendData(xml)
52          print("Friend _Name,Friend_Count")
53          print('Michael Nelson ,' + str(countProfFriend))
54          f2.write("Friend ,")
55          f2.write("Count")
56          f2.write("\r\n")
57          f2.write("Michael Nelson ,")
58          f2.write("%s\n" %str(countProfFriend))
59
60          # f2.write("\r\n")
61
62
63          for friend in fbFriends:
64                  # print(friend + ',' + fbFriends[friend])
65                  f2.write("%s ,"%friend)
66
67                  friend1=int(fbFriends[friend])
68                  sum+=friend1
69                  counter+=1
70                  # print friend1
71
72
73
74
75                  f2.write("%s\n"%fbFriends[friend])
76          print "Mean = ",sum/counter
```

Listing 1: Python program for processing GraphML file

```
1  data <- read.csv("friendCount.txt", stringsAsFactors = F, header
       = FALSE, sep = " ")
2
3  myData = data[,1]
4
5  meanData <- paste("Mean: ", mean(myData), collapse = "")
6
7  medianData <- paste("Median: ", median(myData), collapse = "")
8
9  sdData <- paste("Std Dev: ", sd(myData), collapse = "")
10
11 write(meanData, stdout())
12 write(medianData, stdout())
13 write(sdData, stdout())
14
15 pos <- (myData == 155)
16
17 cols <- c("white", "red")
18
19
20 barplot(myData, main="Friends of Friends on FaceBook", xlab="
       Friends sorted in increasing number of friends", ylab="Number
        of Friends", col=cols[pos + 1], ylim=c(0, 4000))
21 text(x=match(c(155), myData)+12, y=(max(myData)-20), labels="
       Nelson", col='red')
22 arrows(x0=match(c(155), myData)+12, y0=(max(myData) - 80), x1=
       match(c(155), myData)+12, y1=325, length=0.1, lwd=2.5, col='
       red')
```

Listing 2: R program for bar plot shown in Figure 1

# 2

## Question

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 calcuate the same mean, standard deviation, and median values.

For the Twitter 1.1 API to help gather this data, see:

https://dev.twitter.com/docs/api/1.1/get/followers/list

If you do not have followers on Twitter (or don't have more than 50), then use my twitter account "phonedude_mln".


Extra credit, 3 points:

| Mean | 141518.205 |
|---|---|
| Median | 25 |
| Std Dev | 824927.542 |

Table 2: Statistics on the count of my Twitter followers' followers, values straight from R

## Answer

Again, a similar question but a different approach. I had used Tweepy API to extract links in the previous assignments. Thankfully, this experience came in handy. It was a relatively simple task this time. I had to use tweepy.Cursor object and followers_count gave me the number of followers of my followers. screen_name gave me the user name .I have again looped through all users and wrote the values to a csv file. I had sorted the csv file,copied it to a text file and extracted the data. I used this list to generate the mean, median and standard deviation. This time around, the mean was even higher due to a single user having large number of followers. I have plotted a bar graph which again puts me nowhere in comparision to the mean. Friendship paradox is proved again.

From the output, I can see that I have 34 followers.It is greater than the median.It means that I have a considerable rank amongst my followers. The R script shown in Listing 4 creates a similar bar plot to that shown in answer one and produces the statistics shown in listing 2.

```
Mean:   141518.205
Median:   25
Std Dev:   824927.542
```
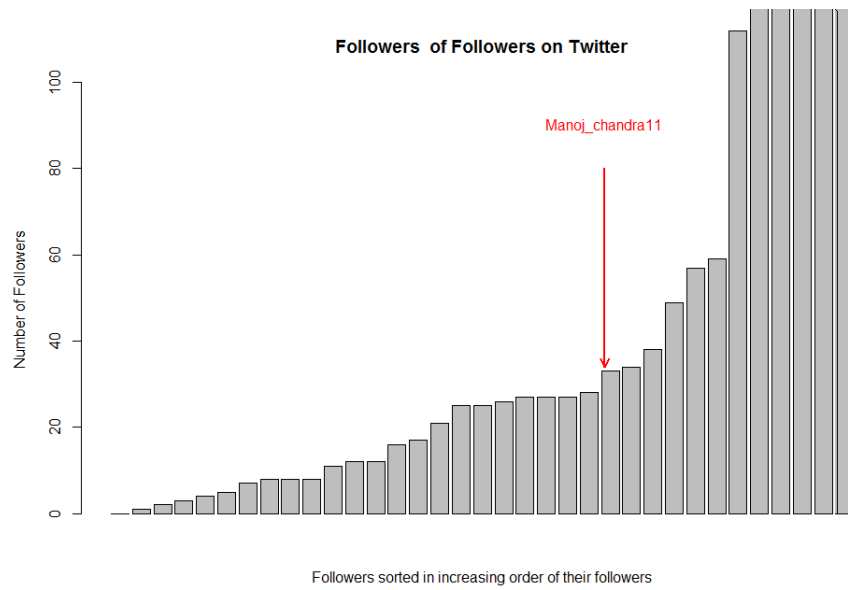
Figure 2: Bar plot showing the count of my Twitter followers' followers

```
 1  import tweepy
 2
 3  import time
 4  import sys
 5  import re
 6
 7  # Authentication Keys to Connect to Twitter API
 8  consumer_key="vjemit5xYQdhgrEPa1FeFf5ZO";
 9  consumer_secret="0
        PoNwIkHk29kUweChIIhGzVD3ZfXwRVqqwxYY3zPadY1BZeNq8";
10  access_token="3485785534−4FlFNlJtg1uNAlMummglVi9feR7fyvkUS0STp0G
        ";
11  access_token_secret="
        EpzAYEpf6tHdGFKj43HhnBeAhLNgkyXPdZyH72ec8Ew8d";
12  auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
13  auth.set_access_token(access_token, access_token_secret)
14  api = tweepy.API(auth) #Accessing tweepy API
15  counter=0
16
17  file1=open('follower.csv','w')
18
19
20  for user in tweepy.Cursor(api.followers, screen_name="
        Manoj_Chandra11").items():
21
22          counter+=1
23
24          file1.write("%s "%user.followers_count)
25          file1.write(",")
26          file1.write("%s\n "%user.screen_name)
```

Listing 3: Python program for acquiring Twitter followers for phonedude_mln

```
1   data <- read.csv("orderedfollowers.txt", stringsAsFactors = F,
        header = FALSE, sep = ",")
2   mydata = data[,1]
3   meanData <- paste("Mean: ", mean(mydata), collapse = "")
4   medianOut <- paste("Median: ", median(mydata), collapse = "")
5   sdData <- paste("Std Dev: ", sd(mydata), collapse = "")
6   write(meanData, stdout())
7   write(medianOut, stdout())
8   write(sdData, stdout())
9
10  barplot(mydata, main="Followers  of Followers on Twitter", xlab=
        "Followers sorted in increasing order of their followers",
        ylab="Number of Followers ",ylim=c(0,100))
11   arrows(x0=match(c(34), mydata)+3, y0=80, x1=match(c(34), mydata
        )+3, y1=34, length=0.1, lwd=2.5, col='red')
12   text(x=match(c(34), mydata)+3, y=90, labels="Manoj_chandra11",
        col='red')
```

Listing 4: R program for bar plot shown in Figure 2

# 3

## Question

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

**Answer**

Not attempted.

# 4

## Question

Extra credit, 1 point:

4.  Repeat question #2, but change "followers" to "following"?  In
other words, are the people I am following following more people?

| | |
|---|---|
| **Mean** | 4938.863 |
| **Median** | 115 |
| **Std Dev** | 2881.577 |

Table 3: Statistics on the count of my Twitter friends' friends, values straight from R

## Answer

This is again a very similar approach to what I did with the second question.I had to extract the friends friends data. It is the "following" data. Here, I used friends_count instead of followers_count.

The result appears to hold in the case of "following", as shown in Figure 3.I again lose terribly to the mean value.Hence,friendship paradox is proved again. The Python code for this case is in Listing 5 and the R code for this Figure is in Listing 6. The statistics are shown in Table 3.

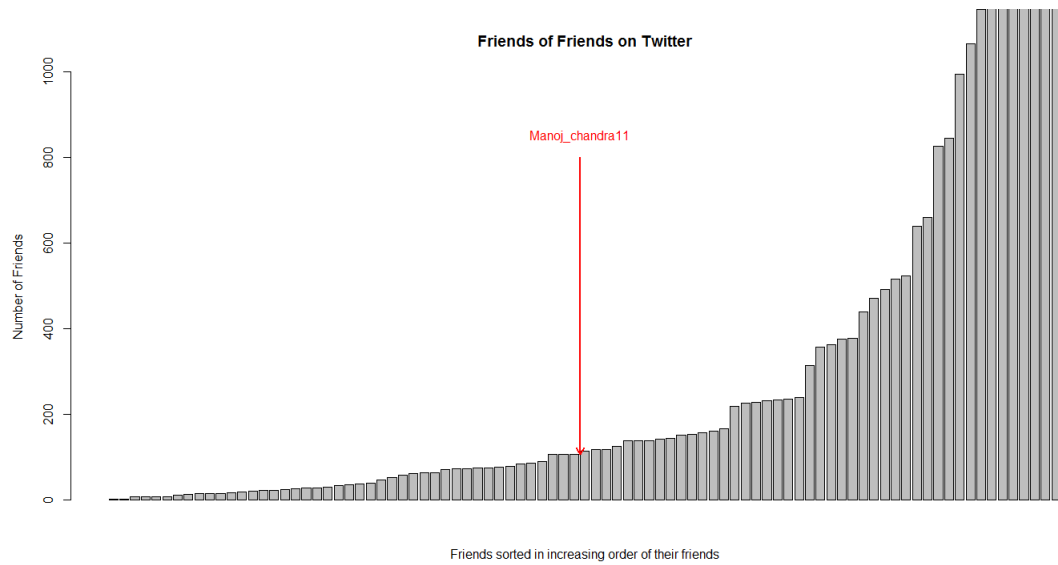Both scripts are executed as shown in the answer to Question 2.

Figure 3: Bar plot showing the count of my Twitter friends' friends

```
1   import tweepy
2   import time
3   import sys
4   import re
5
6   # Authentication Keys to Connect to Twitter API
7   consumer_key="vjemit5xYQdhgrEPa1FeFf5ZO";
8   consumer_secret="0
        PoNwIkHk29kUweChIIhGzVD3ZfXwRVqqwxYY3zPadY1BZeNq8";
9   access_token="3485785534−4FlFNlJtg1uNAlMummglVi9feR7fyvkUS0STp0G
        ";
10  access_token_secret="
        EpzAYEpf6tHdGFKj43HhnBeAhLNgkyXPdZyH72ec8Ew8d";
11  auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
12  auth.set_access_token(access_token, access_token_secret)
13  api = tweepy.API(auth) #Accessing tweepy API
14  counter=0
15
16  file1=open('following.csv','wb')
17
18  for user in tweepy.Cursor(api.friends, screen_name="
        Manoj_Chandra11",count=200).items#gets all my friends details
         on twitter
19          print user.screen_name
20          print user.friends_count
21          file1.write("%s"% user.screen_name)
22          file1.write(",%s\n"% user.friends_count)
```

Listing 5: Python program for acquiring my Twitter followers

```
1  data <- read.csv("followingorder.csv", stringsAsFactors = F,
       header = FALSE, sep = ",")
2  mydata = data[,2]
3  meanOut <- paste("Mean: ", mean(mydata), collapse = "")
4  medianOut <- paste("Median: ", median(mydata), collapse = "")
5  sdOut <- paste("Std Dev: ", sd(mydata), collapse = "")
6  write(meanOut, stdout())
7  write(medianOut, stdout())
8  write(sdOut, stdout())
9  pos <- (mydata == 89)
10  barplot(mydata, main="Friends of Friends on Twitter", xlab="
       Friends sorted in increasing order of their friends", ylab="
       Number of Friends ",ylim=c(0,1000))
11  text(x=match(c(89), mydata)+12, y=850, labels="Manoj_chandra11"
       , col='red')
12  arrows(x0=match(c(89), mydata)+12, y0=800, x1=match(c(89),
       mydata)+12, y1=105, length=0.1, lwd=2.5, col='red')
```

Listing 6: R program for bar plot shown in Figure 3

# References

[1] Dr.Nelson's powerpoint slides:.

http://phonedude.github.io/cs532-s16/

[2] Minidom API.

https://docs.python.org/2/library/xml.dom.minidom.html.

[3] Getting the follower's data using Tweepy:.

http://stackoverflow.com/questions/17455107/the-best-way-to-get-a-list-of-followers-in-python-with-tweepy.

[4] Getting the friends' friend's count in Twitter:.

http://stackoverflow.com/questions/32322519/tweepy-iterating-over-tweepy-cursorapi-friends-items.