# CS 432

# Assignment 4

# Daniel Moore

Question 1:

Part one of this assignment involved determining whether the friendship paradox held true for Dr. Nelson's facebook account. In order to do so I used the following script to parse the xml data in a graph which Dr. Nelson provided of his facebook account and print the friend count for each of his friends to a .csv file:
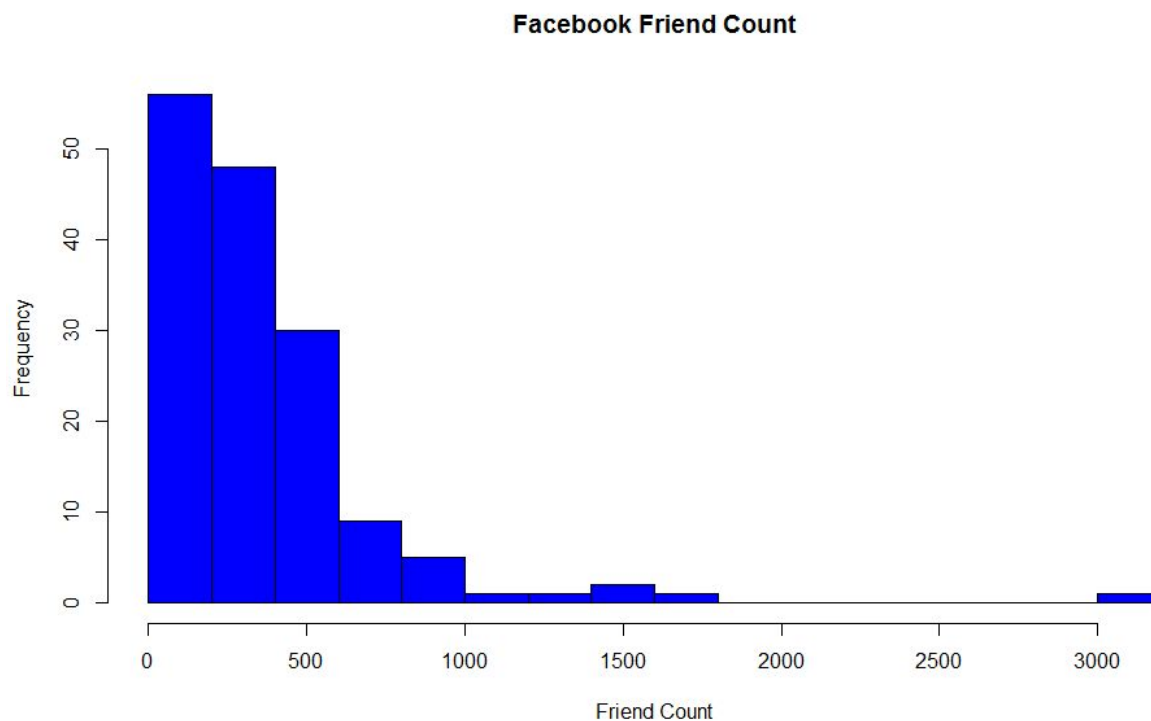
```python
from xml.dom import minidom
import os
dir = r'F:\Web_Science\cs532-s16\A4'
os.chdir(dir)
xmlDoc = minidom.parse("mln.graphml")
graph = xmlDoc.getElementsByTagName("graph")[0] #returns list
nodes = graph.getElementsByTagName("node")
friends = open("friends.csv","a")
noFriends = open("noFriends.txt","a")
friends.write("Friend_Count\n")
noFriends.write("Has No Friends!\n")
for node in nodes:
    data = node.getElementsByTagName("data")
    if len(data) == 4:
        attrib = node.attributes["id"]
        attrName = attrib.value
        noFriends.write("%s\n"%attrName)
    for d in data:
        a = d.attributes["key"]
        val = a.value
        if val == "friend_count":
            friends.write("%s\n"%d.firstChild.data)
```

This data as not available for all of Dr. Nelson's friends. I used the above script to post the names of these users into the text file displayed below:

```
1   Has No Friends!
2   James_Florance_501351702
3   Joy_Gooden_580143423
4   Kim_Beveridge_662936475
5   Alfredo_Sánchez_667415071
6   Sarah_Shreeves_700331809
7   Sally_Mauck_1243862786
8   Dan_Swaney_1321960327
9   Robert_Gordeaux_1580113991
10  Joseph_Kaplan_1623901873
11  Michael_Milner_100000008814265
```

I then used the following r script to plot the extracted data on a histogram:

```
1   friends <- read.table("F:/Web_Science/cs532-s16/A4/friends.csv", quote="\"",
2                                                   comment.char="")
3   View(friends)
4   hist(friends$Friend_Count,main="Facebook Friend Count",xlab="Friend Count",
5                             ylab ="Frequency",col = "blue",breaks=20)
6   summary(friends$Friend_Count)
7
8     Min. 1st Qu.  Median    Mean 3rd Qu.     Max.
9      7.0   134.0   266.5   359.0   446.8  3187.0
10
11      var(friends$Friend_Count)
12  [1] 138075.6
13
14  sd(friends$Friend_Count)
15  [1] 371.5853
```

**Facebook Friend Count**



The median friend count was 266.5 friends with a mean of 359.0 friends and a standard deviation of 371.6 friends. The minimum value was 7 friends and the maximum value was 3187 friends. Both the mean and median values are greater than Dr. Nelson's friend count of 165.
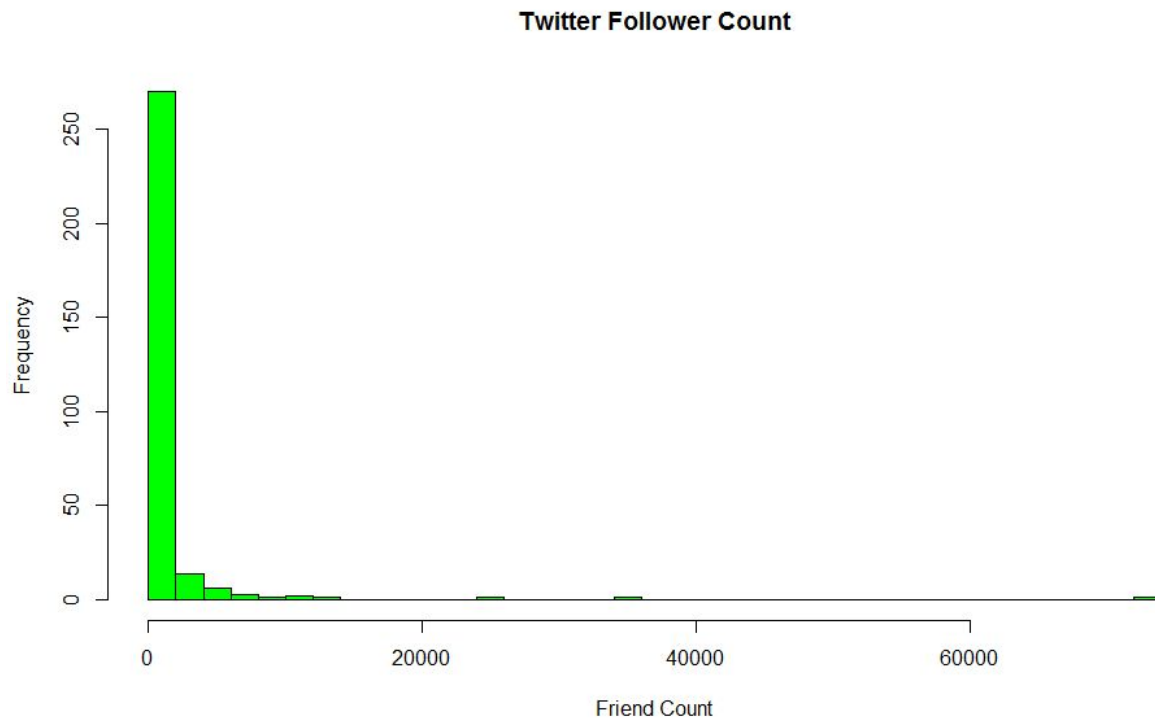
Question 2:

For the second part of this assignment, I was instructed to perform the same tasks for a twitter account. My twitter account only has 2 followers so, once again, I examined Dr. Nelson's account. I used the following script to extract the followers counts of Dr. Nelson's followers:

```python
1   # -*- coding: utf-8 -*-
2   import tweepy
3   import os
4   import simplejson as json
5   import time
6
7   dir = r'F:\Web_Science\cs532-s16\A4'
8   os.chdir(dir)
9
10  auth = tweepy.OAuthHandler('Vz8rTepvf3kVJ2Php7wcIypNt',
11                             'mnkqCLchG38kZEgN36Vlub8o5bmwRD0CLTGNdNlDxaGxiBb7K0')
12  auth.set_access_token('4625770576-Ok6PkaV9hzc6I4kR1jb6Qd48QjYCZvlRhrzYTVu',
13                        '5mWFt5p12bgANFYAe7rjXv4jHH55Ekv5eGwaprEFyqfer')
14
15
16  api = tweepy.API(auth)
17  f = open("twit.csv","a")
18  #t = api.get_user ("phonedude_mln")
19  f.write("Followers\n")
20  for follower in tweepy.Cursor(api.followers,id="phonedude_mln",items=10).items(200):
21      f.write("%s\n"%follower.followers_count)
22      time.sleep(.5)
```

Unfortunately, I repeatedly encountered "Rate_Limit_Errors" and was only able to extract the followers counts for 300 of Dr. Nelson's 491 followers. Below is the script which I used to plot the followers counts in a histogram:

```r
1   twit <- read.csv("F:/Web_Science/cs532-s16/A4/twit.csv", sep="")
2   >    View(twit)
3   hist(twit$Followers,,main="Twitter Follower Count",xlab="Friend Count",
4   +        ylab ="Frequency",col = "green",breaks=30)
5   > summary(twit$Followers)
6      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
7       0.0    49.5   218.5  1227.0   801.8 73180.0
8   > var(twit$Followers)
9   [1] 25979005
10  > sd(twit$Followers)
11  [1] 5096.96
```

## Twitter Follower Count



The minimum followers count was 0 and the maximum was 73180 followers.  The median value was 218.5 and the mean was 1227 with a standard deviation of 5097.