# HW4 - Exploring Social Networks

Brenden Lewis
Due: 10/25/2020 11:59PM

# Q1

The code for Because I find it easier to work with dictionaries in regards to storing data for individuals, I read in each 'friend' entry from *HW4-friend-count.csv* as a dictionary containing their name and number of friends, and stored them in an list which was then sorted before any calculations were made. I utilized the *statistics* package to calculate the mean, median, and standard deviation of the friend counts (*statistics* is used for the same calculations for the other questions).

```
1  import csv
2  import matplotlib
3  import matplotlib.pyplot as plt
4  import statistics
5  import math
6  import pandas as pd
7  from operator import itemgetter
8
9  friends=[]  # {'USER': 'name', 'FRIENDCOUNT': 'count'}
10 sorted_friends=[]
11 User = 'U'
12 mean = 1
13 standard = 1
14 median = 1
15
16 def process():
17     with open("HW4-friend-count.csv",'r') as f:
18         count = 1
19         for line in csv.DictReader(f):
20             new_dict = {'USER':"f"+str(count),'FRIENDCOUNT':int(line['
   FRIENDCOUNT'])}
21             friends.append(new_dict)
22             count+=1
23
24         friends.append({'USER':User, 'FRIENDCOUNT': len(friends)})
25
26         global sorted_friends
27         sorted_friends = sorted(friends, key=itemgetter('FRIENDCOUNT'))
28
29         global mean,standard,median
30         mean = calcMean()
31         standard = calcStandDevi()
```

```
32          median = calcMedian()
33
34 def calcMean():
35     friend_counts = []
36     for user in friends:
37         friend_counts.append(int(user['FRIENDCOUNT']))
38
39     result = statistics.mean(friend_counts)
40
41     return math.trunc(result)
42
43 def calcStandDevi():
44     friend_counts = []
45     for user in friends:
46         friend_counts.append(int(user['FRIENDCOUNT']))
47
48     result = statistics.pstdev(friend_counts)
49     return round(result,2)
50
51 def calcMedian():
52     friend_counts = []
53     for user in friends:
54         friend_counts.append(int(user['FRIENDCOUNT']))
55
56     friend_counts.sort()
57     mid = statistics.median(friend_counts)
58
59     return math.trunc(mid)
60
61 def drawGraph():
62     x=[]
63     y=[]
64     ux=[]
65     uy=[]
66     for per in sorted_friends:
67         x.append(per['USER'])
68         y.append(per['FRIENDCOUNT'])
69         if(per['USER']==User):
70             ux.append(per['USER'])
71             uy.append(per['FRIENDCOUNT'])
72
73     plt.title('Friendship Paradox: User "'+User+'"')
74     plt.xlabel(User+' vs. their friends')
75     plt.ylabel('# of friends')
76
77     plt.scatter(x,y, label=User+"'s friends")
```

```
78      plt.tick_params(axis='x', which='both',bottom=False, labelbottom=
        False)
79      plt.scatter(ux,uy,color='yellow', label=User)
80      plt.legend(loc="upper left")
81
82      plt.show()
83
84
85
86
87
88
89  def writeResults():
90      print("Statistics for 'Friend' counts of "+User+" and their friends
        ")
91      print("Mean: "+str(mean))
92      print("Median: "+str(median))
93      print("Standard Deviation: "+str(standard))
94
95      with open("FriendStatistics.txt","w") as f:
96          f.write("Statistics for 'Friend' counts of "+User+" and their
        friends"+'\n')
97          f.write("Mean: "+str(mean)+'\n')
98          f.write("Median: "+str(median)+'\n')
99          f.write("Standard Deviation: "+str(standard))
100
101  if __name__ == '__main__':
102      process()
103      writeResults()
104      drawGraph()
```

**Listing 1:** Python code Q1

Listing 2 shows the output of the calculations for Q1. I'm a little skeptical of the result for standard deviation as it seems a bit too high; this could be due to outliers in the data, or an error in calculating the standard deviation (either wrong methods or some other factor I'm missing).

```
1  Statistics for 'Friend' counts of U and their friends
2  Mean: 538
3  Median: 395
4  Standard Deviation: 535.81
```

**Listing 2:** Calculated mean, median, and standard deviation for Q1

Figure 1 shows a scatter plot of the User's (U) friend counts compared to their friends' friend counts. The results of graph seem to show that the Friendship Paradox holds for this particular user. The User had about 98 friends in their list. Compared to the majority of their friends, the User has a lot less friends, with some friends having thousands more friends to their account.
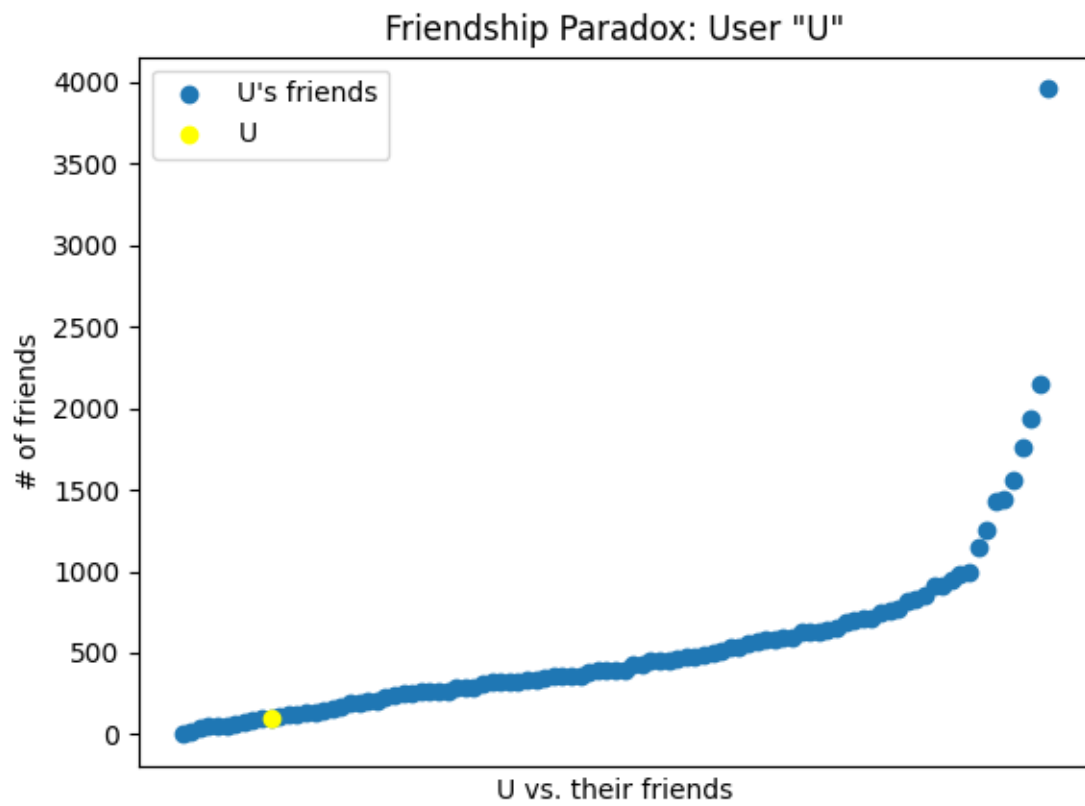
**Figure 1:** Graph showing the friend count of the user, and the friend counts of their friends

# Q2

I opted to split this section into two programs: one to collect all the follower counts from Twitter (Listing 3), and one to handle data calculations and graphing (Listing 5). I don't have an active Twitter account, so I used our instructor's Twitter handle *'weiglemc'* as the user ID to collect data. The biggest challenge here was handling a RateLimitError that occured when the rate limit was exceeded pulling from Twitter. I got around this by adding the parameter *wait_rate_limit=True* to the *api* initialization. This caused the stream to halt when the rate limit was reached to give it time to refresh and prevent errors from attempting to exceed it.

A Cursor was used to move through the list of *weiglemc*'s followers. For each follower, a dictionary was created with their Twitter handle (ID) and their follower counts and stored in list which would then be sorted by follower counts after everyone was processed; follower counts was found using *'.followers_count'* for each user.

*Of note: I ended up not using the getFollowerCount() method and forgot to remove it.*

```python
1  import matplotlib
2  import matplotlib.pyplot as plt
3  import statistics
4  import math
5  import pandas as pd
6  from operator import itemgetter
7  import tweepy
8  from tweepy.streaming import StreamListener
9  from tweepy import OAuthHandler
10 from tweepy import Stream
11 import time
12
13 # Keys ommitted
14 consumer_key="***"
15 consumer_secret="***"
16 access_token="***"
17 access_secret="***"
18
19 # Handles authorization with Twitter
20 auth = OAuthHandler(consumer_key,consumer_secret)
21 auth.set_access_token(access_token,access_secret)
22 api = tweepy.API(auth, wait_on_rate_limit=True,
       wait_on_rate_limit_notify=True)
23 # api = tweepy.API(auth)
24
25 user_ID = 'weiglemc'
26 User=api.get_user(user_ID)
27
28 # STREAM_CAP = 433
29 STREAM_CAP = User.followers_count
```

```python
30 followers=[] # {'USER': ID, 'FOLLOWERCOUNT': count}
31 sortedFollowers=[]
32
33 def process():
34     try:
35         for follower in tweepy.Cursor(api.followers, user_ID).items(
    STREAM_CAP):
36             follower_ID = follower.screen_name
37             print("Processing "+follower_ID+"...")
38             follower_count = follower.followers_count
39             follower_dict = {'USER': follower_ID, 'FOLLOWERCOUNT':
    follower_count}
40             followers.append(follower_dict)
41             print("Finshed"+'\n')
42     except tweepy.RateLimitError:
43         print ("Rate Limit reached. Sleeping for 60s")
44         time.sleep(60)
45
46     followers.append({'USER': user_ID, 'FOLLOWERCOUNT': User.
    followers_count})
47
48     global sortedFollowers
49     sortedFollowers = sorted(followers, key=itemgetter('FOLLOWERCOUNT')
    )
50
51
52 def getFollowerCount(userName):
53     c = tweepy.Cursor(api.followers, userName)
54
55     count = 0
56     for follower in c.items():
57         count += 1
58
59     return count
60
61 def writeOutput():
62     for item in followers:
63         print(str(item))
64
65 def writeToFile():
66     with open("FollowerList.txt","w") as f:
67         for item in sortedFollowers:
68             f.write(str(item)+'\n')
69
70 if __name__=='__main__':
71     process()
72     writeOutput()
```

```
73    writeToFile()
```

**Listing 3:** Python code Q2

Listing 4 below shows how the information for each user was stored in the output file:

```
1 {'USER': 'rachelheldevans', 'FOLLOWERCOUNT': 163355}
2 {'USER': 'KHayhoe', 'FOLLOWERCOUNT': 169495}
3 {'USER': 'mattcutts', 'FOLLOWERCOUNT': 529815}
```

**Listing 4:** Examples of stored dictionaries for followers

Listing 5 below is a slightly modified version of the code from Listing 1 to accomodate the data collected for Q2. The only major difference aside from file and graph nomenclature is how the information from the new file is read in. Listing 1 uses a method *csv.DictReader()* from the *csv* library to read in lines from csv files as dictionaries. Listing 5 (and also Listing 10 for Q3) uses *ast.literal_eval(line)* to convert each line from the file, as is, into a dictionary to be stored. The users were already sorted previously, so there was no need to sort again.

```
1 import csv
2 import matplotlib
3 import matplotlib.pyplot as plt
4 import statistics
5 import math
6 import pandas as pd
7 from operator import itemgetter
8 import ast
9
10 followers=[]  # {'USER': 'name', 'FRIENDCOUNT': 'count'}
11 User = 'weiglemc'
12 mean = 1
13 standard = 1
14 median = 1
15
16 def process():
17     with open("FollowerList.txt",'r') as f:
18         count = 1
19         for line in f:
20             new_dict = ast.literal_eval(line)
21             followers.append(new_dict)
22             count+=1
23
24         global mean,standard,median
25         mean = calcMean()
26         standard = calcStandDevi()
27         median = calcMedian()
28
29 def calcMean():
30     friend_counts = []
```

```python
31      for user in followers:
32          friend_counts.append(int(user['FOLLOWERCOUNT']))
33
34      result = statistics.mean(friend_counts)
35
36      return math.trunc(result)
37
38  def calcStandDevi():
39      friend_counts = []
40      for user in followers:
41          friend_counts.append(int(user['FOLLOWERCOUNT']))
42
43      result = statistics.pstdev(friend_counts)
44      return round(result,2)
45
46  def calcMedian():
47      friend_counts = []
48      for user in followers:
49          friend_counts.append(int(user['FOLLOWERCOUNT']))
50
51      friend_counts.sort()
52      mid = statistics.median(friend_counts)
53
54      return math.trunc(mid)
55
56  def drawGraph():
57      x=[]
58      y=[]
59      ux=[]
60      uy=[]
61      for per in followers:
62          x.append(per['USER'])
63          y.append(per['FOLLOWERCOUNT'])
64          if(per['USER']==User):
65              ux.append(per['USER'])
66              uy.append(per['FOLLOWERCOUNT'])
67
68      plt.title('Friendship Paradox: User "'+User+'"')
69      plt.xlabel(User+' vs. their followers')
70      plt.ylabel('# of followers')
71
72      plt.scatter(x,y, label=User+"'s followers")
73      plt.tick_params(axis='x', which='both',bottom=False, labelbottom=
    False)
74      plt.scatter(ux,uy,color='yellow', label=User)
75      plt.legend(loc="upper left")
76
```

```
77     plt.show()
78
79
80
81
82
83
84 def writeResults():
85     print("Statistics for 'Follower' counts of "+User+" and their
       followers")
86     print("Mean: "+str(mean))
87     print("Median: "+str(median))
88     print("Standard Deviation: "+str(standard))
89
90     with open("FollowerStatisticsTrimmed.txt","w") as f:
91         f.write("Statistics for 'Follower' counts of "+User+" and their
       followers"+'\n')
92         f.write("Mean: "+str(mean)+'\n')
93         f.write("Median: "+str(median)+'\n')
94         f.write("Standard Deviation: "+str(standard))
95
96 if __name__ == '__main__':
97     process()
98     writeResults()
99     drawGraph()
```

**Listing 5:** Python code for graphing follower counts

Listing 6 shows the output for calculations using data from follower counts. Listing 7 shows the output from the same data minus three outliers heavily skewing the data to the right. The errors related to standard deviation seem much more apparent in these listings compared to that in Q1, or once again the deviation is being affected by the larger follower counts skewing the data. It remains unclear.

```
1 Statistics for 'Follower' counts of weiglemc and their followers
2 Mean: 2915
3 Median: 195
4 Standard Deviation: 27839.99
```

**Listing 6:** Calculated mean, median, and standard deviation for Q2

```
1 Statistics for 'Follower' counts of weiglemc and their followers
2 Mean: 945
3 Median: 393
4 Standard Deviation: 2244.33
```

**Listing 7:** Calculated mean, median, and standard deviation for Q2, excluding a few outliers

Figure 2 shows the original data of the user *weiglemc*'s followers, and their followers' followers. To help provide a better visual curve, the three outliers were removed and the data was graphed
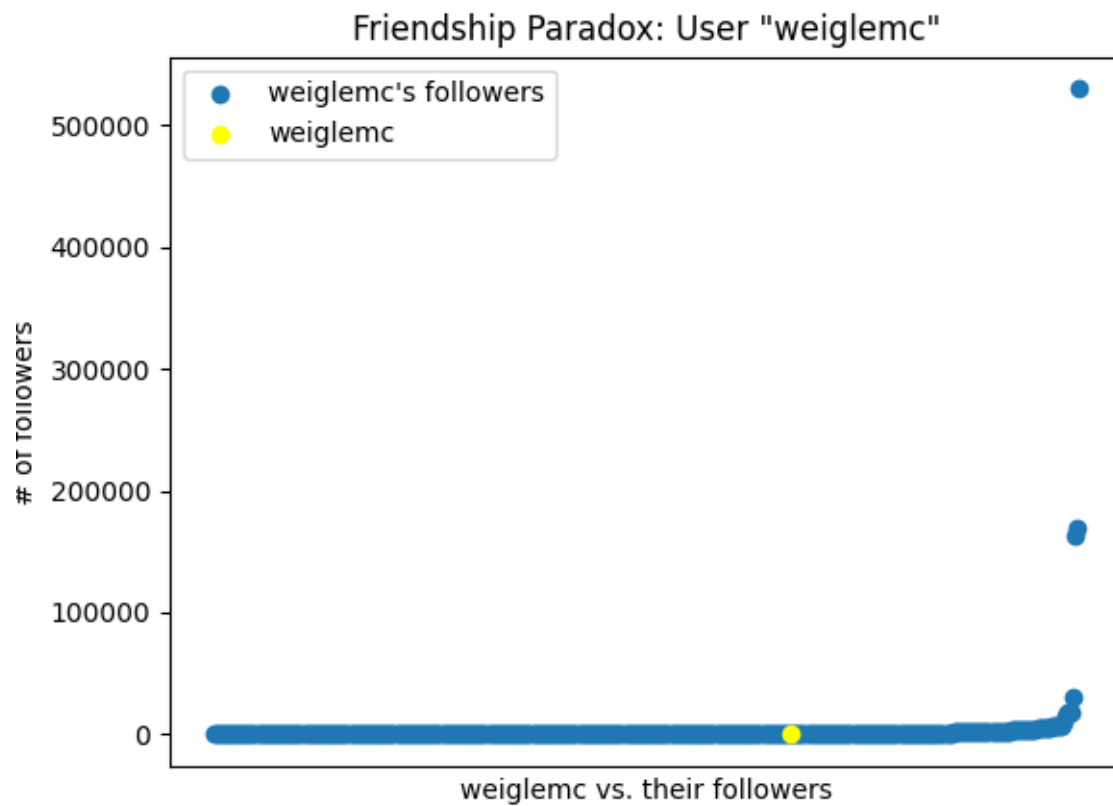
**Figure 2:** Graph showing the follower count of the user, and the follower counts of their followers

again in Figure 3. It is difficult to tell based on both graphs if the Friendship Paradox holds as *weiglemc* is just a little right of the middle of both curves; this would indicate that *weiglemc* has more followers than more than half of the users they follow.
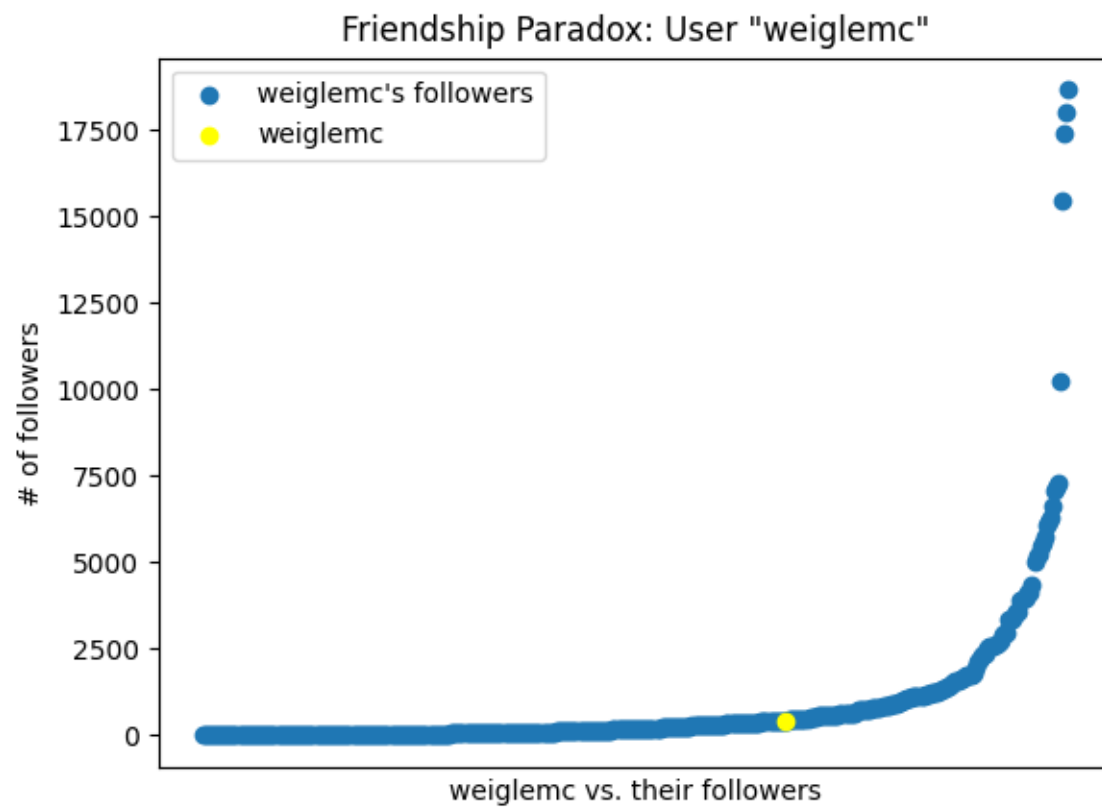
**Figure 3:** Graph showing the follower count of the user, and the follower counts of their followers, excluding outliers

# Q3 (Extra)

The code used for Q3 (Listing 8 and 10) is very much the same to that used in Q2 with a few minor changes. In order to get *following* counts from a user instead of *followers*, I had to alter the Cursor to use *api.friends* and call *.friends_count* for each user to get the number of users they each follow. For some reason, the Twitter API uses Friend interchangeably with Follower when doing anything related to who a user follows. When collecting data, I had to double check a few Twitter accounts to ensure it was actually capturing the number of users they follow.

```python
1 import matplotlib
2 import matplotlib.pyplot as plt
3 import statistics
4 import math
5 import pandas as pd
6 from operator import itemgetter
7 import tweepy
8 from tweepy.streaming import StreamListener
9 from tweepy import OAuthHandler
10 from tweepy import Stream
11 import time
12
13 # Keys ommitted
14 consumer_key="***"
15 consumer_secret="***"
16 access_token="***"
17 access_secret="***"
18
19 # Handles authorization with Twitter
20 auth = OAuthHandler(consumer_key,consumer_secret)
21 auth.set_access_token(access_token,access_secret)
22 api = tweepy.API(auth, wait_on_rate_limit=True,
     wait_on_rate_limit_notify=True)
23 # api = tweepy.API(auth)
24
25 user_ID = 'weiglemc'
26 User=api.get_user(user_ID)
27
28 # STREAM_CAP = 5
29 STREAM_CAP = User.friends_count
30 friend_list=[] # {'USER': ID, 'FRIENDCOUNT': count}
31 sortedFriends=[]
32
33 def process():
34     try:
35         for follower in tweepy.Cursor(api.friends, user_ID).items(
     STREAM_CAP):
36             friend_ID = follower.screen_name
```

```python
37             print("Processing "+friend_ID+"...")
38             friend_count = follower.friends_count
39             friend_dict = {'USER': friend_ID, 'FRIENDCOUNT':
    friend_count}
40             friend_list.append(friend_dict)
41             print("Finshed"+'\n')
42        except tweepy.RateLimitError:
43            print ("Rate Limit reached. Sleeping for 60s")
44            time.sleep(60)
45
46        friend_list.append({'USER': user_ID, 'FRIENDCOUNT': STREAM_CAP})
47
48        global sortedFriends
49        sortedFriends = sorted(friend_list, key=itemgetter('FRIENDCOUNT'))
50
51
52 def getFollowerCount(userName):
53        c = tweepy.Cursor(api.followers, userName)
54
55        count = 0
56        for follower in c.items():
57            count += 1
58
59        return count
60
61 def writeOutput():
62        for item in sortedFriends:
63            print(str(item))
64
65 def writeToFile():
66        with open("FriendsList.txt","w") as f:
67            for item in sortedFriends:
68                f.write(str(item)+'\n')
69
70 if __name__=='__main__':
71     process()
72     writeOutput()
73     writeToFile()
```

**Listing 8:** Python code Q3

Listing 9 shows how the data gathered was stored in the output file before any calculations or graphing is done.

```python
1 {'USER': 'nakedpastor', 'FRIENDCOUNT': 17790}
2 {'USER': 'MarkWarner', 'FRIENDCOUNT': 25114}
3 {'USER': 'BarackObama', 'FRIENDCOUNT': 599538}
```

**Listing 9:** Examples of stored dictionaries for accounts followed

```python
1  import csv
2  import matplotlib
3  import matplotlib.pyplot as plt
4  import statistics
5  import math
6  import pandas as pd
7  from operator import itemgetter
8  import ast
9
10 followers=[]  # {'USER': 'name', 'FRIENDCOUNT': 'count'}
11 User = 'weiglemc'
12 mean = 1
13 standard = 1
14 median = 1
15
16 def process():
17     with open("FriendsList.txt",'r') as f:
18         count = 1
19         for line in f:
20             new_dict = ast.literal_eval(line)
21             followers.append(new_dict)
22             count+=1
23
24         global mean,standard,median
25         mean = calcMean()
26         standard = calcStandDevi()
27         median = calcMedian()
28
29 def calcMean():
30     friend_counts = []
31     for user in followers:
32         friend_counts.append(int(user['FRIENDCOUNT']))
33
34     result = statistics.mean(friend_counts)
35
36     return math.trunc(result)
37
38 def calcStandDevi():
39     friend_counts = []
40     for user in followers:
41         friend_counts.append(int(user['FRIENDCOUNT']))
42
43     result = statistics.pstdev(friend_counts)
44     return round(result,2)
45
46 def calcMedian():
47     friend_counts = []
```

```python
48      for user in followers:
49          friend_counts.append(int(user['FRIENDCOUNT']))
50
51      friend_counts.sort()
52      mid = statistics.median(friend_counts)
53
54      return math.trunc(mid)
55
56  def drawGraph():
57      x=[]
58      y=[]
59      ux=[]
60      uy=[]
61      for per in followers:
62          x.append(per['USER'])
63          y.append(per['FRIENDCOUNT'])
64          if(per['USER']==User):
65              ux.append(per['USER'])
66              uy.append(per['FRIENDCOUNT'])
67
68      plt.title('Friendship Paradox: User "'+User+'"')
69      plt.xlabel(User+' vs. who they follow')
70      plt.ylabel('# following')
71
72      plt.scatter(x,y, label=User+"'s # following")
73      plt.tick_params(axis='x', which='both',bottom=False, labelbottom=
    False)
74      plt.scatter(ux,uy,color='yellow', label=User)
75      plt.legend(loc="upper left")
76
77      plt.show()
78
79
80
81
82
83
84  def writeResults():
85      print("Statistics for 'Follower' counts of "+User+" and their
    followers")
86      print("Mean: "+str(mean))
87      print("Median: "+str(median))
88      print("Standard Deviation: "+str(standard))
89
90      with open("FollowingStatistics.txt","w") as f:
91          f.write("Statistics for 'Following' counts of "+User+" and who
    they follow"+'\n')
```

```
92          f.write("Mean: "+str(mean)+'\n')
93          f.write("Median: "+str(median)+'\n')
94          f.write("Standard Deviation: "+str(standard))
95
96 if __name__ == '__main__':
97     process()
98     writeResults()
99     drawGraph()
```

**Listing 10:** Python code for graphing following counts

Listing 11 shows the results of calculations of data collected for Q3. Listing 12 shows the same data exluding a single outlier that was *very* heavily skewing the data, the account of late President Barack Obama with nearly six-hundred thousand accounts followed.

```
1 Statistics for 'Following' counts of weiglemc and who they follow
2 Mean: 3238
3 Median: 394
4 Standard Deviation: 37048.7
```

**Listing 11:** Calculated mean, median, and standard deviation for Q3

```
1 Statistics for 'Following' counts of weiglemc and who they follow
2 Mean: 859
3 Median: 192
4 Standard Deviation: 2094.13
```

**Listing 12:** Calculated mean, median, and standard deviation for Q3, excluding an outlier

Figures 4 and 5 show scatter plots of the data for Q3, with Figure 5 excluding Obama's account to get a more accurate visual. In the case of *weiglemc*, just like in Q2, it is difficult to tell if the Friendship Paradox holds for counts of users followed as they are placed very near the center of the curves. User *weiglemc* veers slightly to the left of the median, meaning they are following less users than more than 50% of the users they follow.
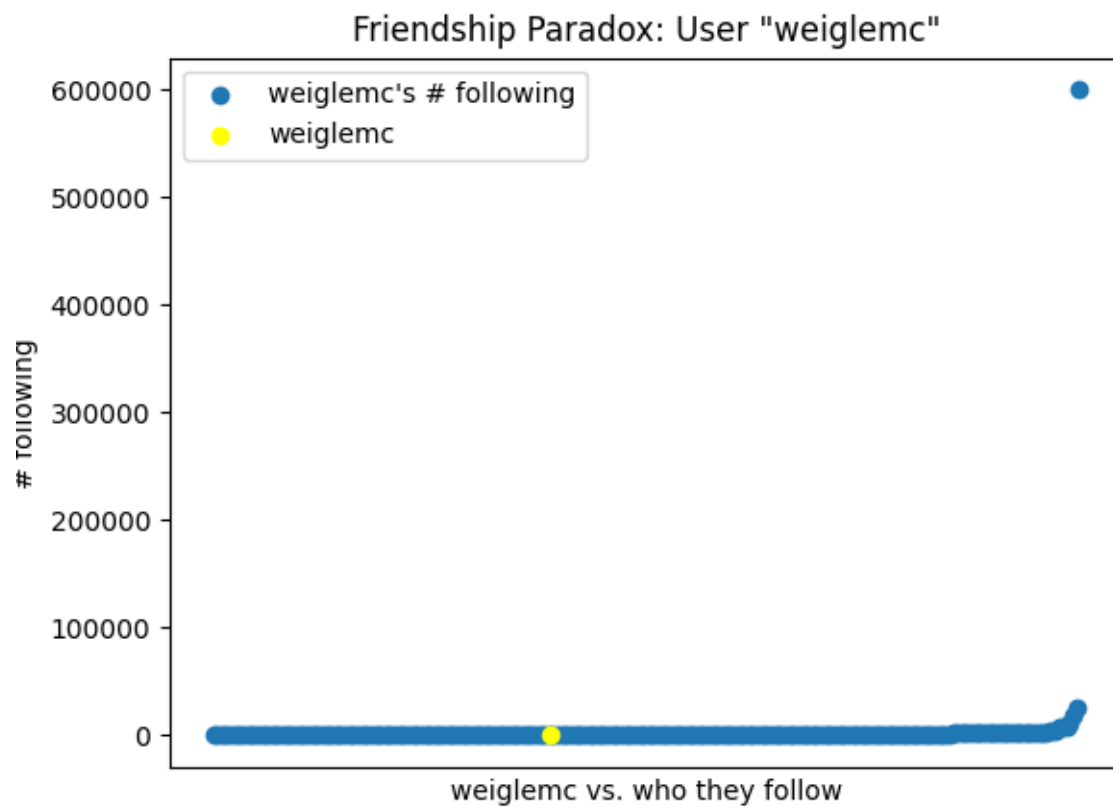
**Figure 4:** Graph showing the number of accounts the user follows, and the following counts of those they follow
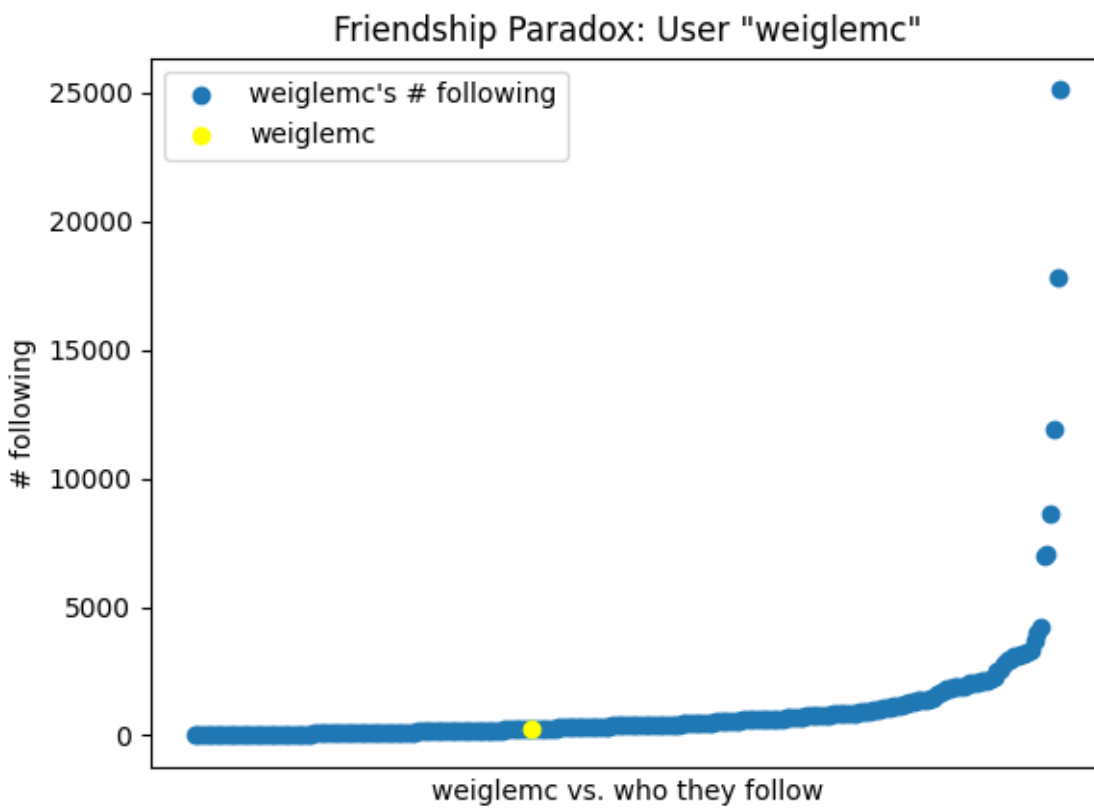
**Figure 5:** Graph showing the number of accounts the user follows, and the following counts of those they follow, excluding outliers