Fan Guo
A20473828

# Project I - Social Media Data Analysis for Twitter

Project Github Link: https://github.com/fguo10/SocialMediaDataAnalysis

In this project, I use twitter token and tweepy package to crawl twitter friends data, as well as process and report some analysis on my extracted data.

It contains 3 steps here:

1) Data Collection

2) Data Visualization

3) Network Measures Calculation


## 1. Data Collection

Create an account on Twitter, then visit the link: https://developer.twitter.com/en.

Then click on Apply for a Developer account, we can select Specify the type

of application, and Specify the type of application, put in why you want to apply for the developer account. When the API developer is approved, we can get token details:

```
# twitter developer access token
# Your app's API/consumer key and secret can be found under the Consumer Keys
# section of the Keys and Tokens tab of your app, under the
# Twitter Developer Portal Projects & Apps page at
# https://developer.twitter.com/en/portal/projects-and-apps
CONSUMER_KEY = "dbPd                    "
CONSUMER_SECRET = "RwC3z                        '

# Your account's (the app owner's account's) access token and secret for your
# app can be found under the Authentication Tokens section of the
# Keys and Tokens tab of your app, under the
# Twitter Developer Portal Projects & Apps page at
# https://developer.twitter.com/en/portal/projects-and-apps
ACCESS_TOKEN = "14417                            "
ACCESS_TOKEN_SECRET = "7AJRng                     "
```

Base on the API token, we can use tweepy python package esaily get twitter user friends.As we know, Tweepy is an open source Python package that makes it easy to access Twitter and get information about friends for developers to use.

Firstly, Setting the authentication credentials.

```
42    class TwitterDataCollection:
43        def __init__(self):
44            """..."""
48            auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
49            auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)
50            self.api = tweepy.API(auth, retry_count=5, retry_delay=60, wait_on_rate_limit=True)
51
52        def get_one_user_detail(self, username):
53            """..."""
58            user = self.api.get_user(screen_name=username)
59            friends = []
60            for friend in user.friends():
61                friends.append(friend.screen_name)
62            return friends, user.friends_count
```

However, Twitter has a rate limit on the number of requests made to the Twitter API. Allow 900 requests /15 minutes. If get the request error because of too many request, we need wait 15 minutes and try again, the code as below:

```
64        def get_all_friends(self, try_count=5) -> dict:
65            """..."""
69            all_friends = {}
70            load_users = get_users()
71            for username in load_users:
72                for i in range(try_count):
73                    try:
74                        data = self.get_one_user_detail(username)
75                        all_friends[username] = {
76                            'friends': data[0],
77                            'friends_count': data[1]
78                        }
79                    except Exception as e:
80                        print('get twitter friends limit, sleep 15min')
81                        time.sleep(15)
82                    else:
83                        break
84            return all_friends
```

Meanwhile, to reduce the number of Twitter requests, I save the data into files doc/all_friends.txt. When analyzing data, we don't have to collect friends of a given user on Twitter every time. The code details as follows:

```
86        def save_friends_to_json(self, all_friends):
87            """..."""
92            with open(ALL_FRIENDS_FILENAME, 'w') as f:
93                json.dump(all_friends, f, indent=4)
```

Fan Guo

A20473828

The part about data collection has been completed, the steps executed are stored in the function crawl_data, and the file doc/all_friends.txt will be generated after the execution, the details are as follows:

```python
160    def crawl_data():
161        twitter = TwitterDataCollection()
162        all_friends = twitter.get_all_friends()
163        twitter.save_friends_to_json(all_friends)
```
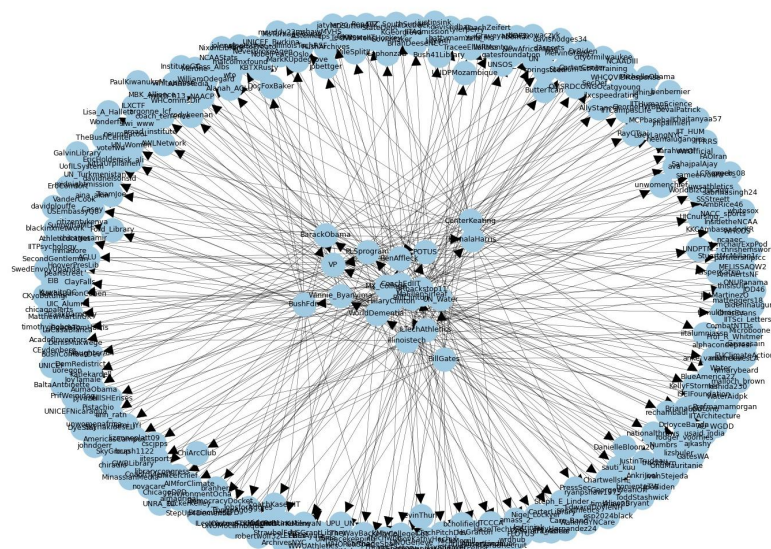
## 2. Data Visualization

I have fetched the data, in this step, I use networkx to visualize network as graph. The core code here:

```python
124    def build_graph(self, draw_screen_names=None):
125        G_asymmetric = nx.DiGraph()
126        all_friends = load_friends_to_dict()
127        for screen_name, data in all_friends.items():
128            friends = data.get('friends', [])
129            if not draw_screen_names:
130                self.add_nodes(G_asymmetric, screen_name, friends)
131            else:
132                if screen_name in draw_screen_names:
133                    self.add_nodes(G_asymmetric, screen_name, friends)
134        return G_asymmetric
```
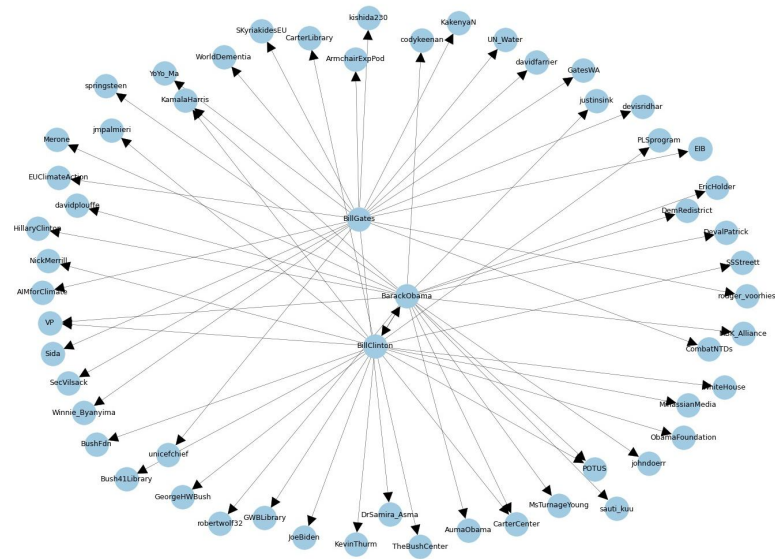
```python
166    def visualize_data():
167        visualize_obj = DataVisualization()
168        visualize_obj.build_graph(draw_screen_names=['BarackObama', 'BillGates', 'BillClinton'])
169        visualize_obj.build_graph(draw_screen_names=['illinoistech', 'ILTechAthletics', 'CoachEdIIT'])
170        visualize_obj.build_graph()
```

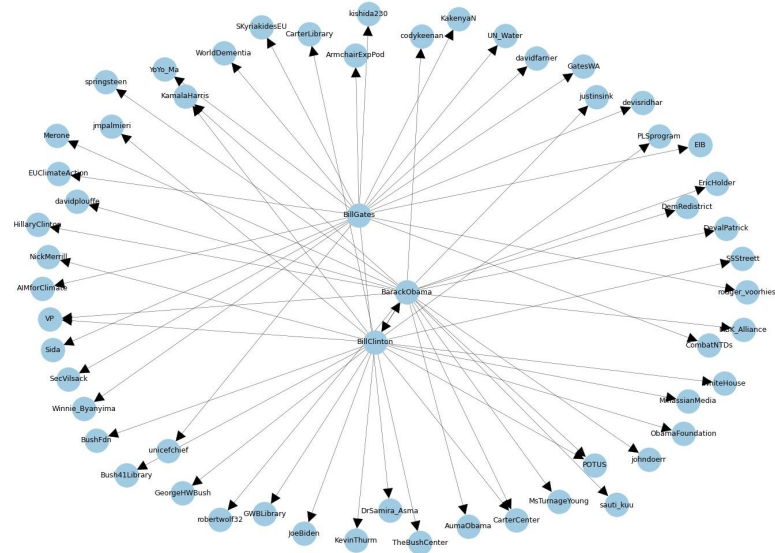When I use all the data(400+ nodes) here, I got a complex graph:

When I choose screen_users: ['BarackObama', 'BillGates', 'BillClinton'], get a simple graph here:



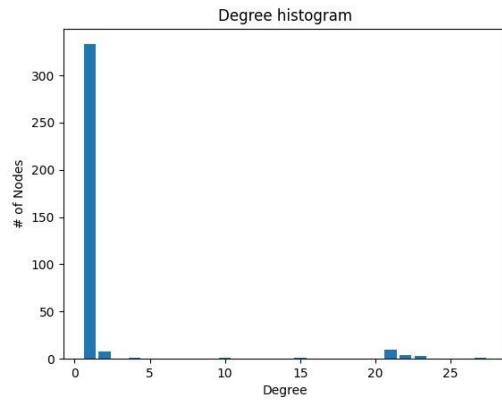Then I choose screen_users: ['illinoistech', 'ILTechAthletics', 'CoachEdIIT'] , get a simple graph here:



# 3. Network Measures Calculation

In this step, I use networkx and matplotlib to caculation Degree Distribution and Pagerank, and plot it as histogram. The core code here:
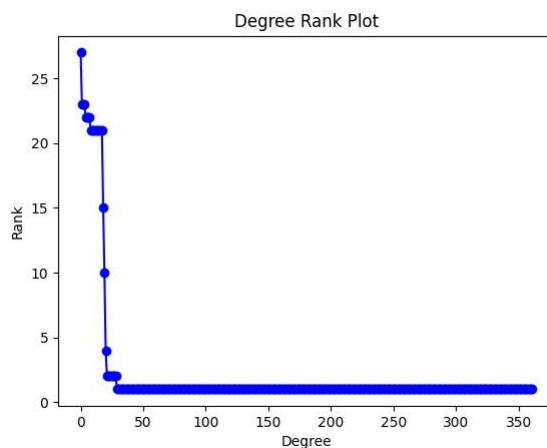
```
173    def analysis_data():
174        analyze_obj = DataAnalysis()
175        analyze_obj.draw_degree_rank_plot()
176        analyze_obj.draw_degree_histogram()
```

Then the degree histogram as the picture shows:



The degree-rank plot for the Graph as the picture shows:



# 4. Conclusion

In this project, Firstly, get twitter data from API and save data into file. Secondly, we do data visualize about 100+ nodes and 400+nodes, At the end , we do various type of statistical analysis on the tweets. After the completion of the project, I learned a lot about online network analysis

Project Github Link: https://github.com/fguo10/SocialMediaDataAnalysis