

Kevin Lin
CISC 4900

Project Summary

The topic of this project is about data analysis using the social media platform twitter. Using Twitter's API and token keys, tweets were extracted from twitter. A python library called tweepy was used in this process to extract the tweets. The Listener class that was created to inherit the Tweepy's Stream Listener class. This allowed the live extraction of tweets. A streamer class was created to verify the API and token keys. A constructor class was created for the time limit and an array was created for the tweets in JSON style. The on_data method was created to fetch for live tweets. The on_status and on_error methods were created to check if there were live tweets fetched. Two constructors were created for the Streamer class One constructor was used for the live tweets, the other was used to get previous tweets. A filter method was created to filter out any unwanted previous tweets. All these methods and classes are in the Twitter_Stream.py file.

```
1 import time
2 import csv
3 from tweepy import OAuthHandler
4 from tweepy import Stream
5 from tweepy import API
6 from tweepy import Cursor
7 import Twitter_Credentials
8 from tweepy.streaming import StreamListener
9 import pandas as pd
10
11
12 # This class handles whether or not the authentication works.
13 class Listener(StreamListener):
14
15     # constructor sets up the current time and the end time of the stream
16     def __init__(self, time_limit=60):
17         self.start_time = time.time() # instantiates the current time
18         self.end_time = time_limit # instantiates the end time in seconds
19         self.twitter_data = [] # instantiates a new list for JSON data
20
21     def on_data(self, raw_data):
22         try:
23             if (time.time() - self.start_time) < self.end_time: # calculates the amount of seconds left in the stream
24                 self.twitter_data.append(raw_data) # writes JSON data into list
25                 return True
26         except BaseException as e:
27             print('failed because:', str(e))
28             time.sleep(5)
29             pass
30         json_file = open('unfiltered_tweets_Donald_Trump.json', 'a', encoding='utf-8',
31                         newline='') # writes into JSON file
32         json_file.write(u'\n')
33         json_file.write(','.join(self.twitter_data)) # writes into JSON file with the list
34         json_file.write(u'\n\n')
35
36         json_file.write(u'\n')
37         json_file.close()
38         exit()
39
40 # If the connection works it will return true
41 def on_status(self, status):
42     print(status)
43     return True
44
45 # If there is a connection problem it will return false
46 def on_error(self, status_code):
47     if status_code == 420:
48         # Will return false if too many login attempts. Too many times authenticating with the account name.
49         return False
50     print(status_code)
51
52 # This class authenticates twitter's API key and access token. It also output tweets.
53 class Streamer:
54
55     # The constructor authenticates the API key and access token. This constructor start the time limit for live tweets
56     def __init__(self, time_limit):
57         self.listen = Listener(time_limit=time_limit)
58         self.auth = OAuthHandler(Twitter_Credentials.API_KEY, Twitter_Credentials.API_KEY_SECRET)
59         self.auth.set_access_token(Twitter_Credentials.ACCESS_TOKEN, Twitter_Credentials.ACCESS_TOKEN_SECRET)
60
61     # This constructor authenticates the API key and access token. This constructor searches for previous tweets
62     def __init__(self):
63         self.auth = OAuthHandler(Twitter_Credentials.API_KEY, Twitter_Credentials.API_KEY_SECRET)
64         self.auth.set_access_token(Twitter_Credentials.ACCESS_TOKEN, Twitter_Credentials.ACCESS_TOKEN_SECRET)
65
66     # This method output any tweets that are related to the user's input
67     def output_live_tweets(self, input_filter):
68         streamer = Stream(auth=self.auth, listener=self.listen)
```

```

65 # This method output any tweets that are related to the user's input
66 def output_live_tweets(self, input_filter):
67     streamer = Stream(self.auth, self.listen)
68     streamer.filter(track=input_filter)
69
70 def output_previous_tweets(self, start, end):
71     api = API(self.auth, wait_on_rate_limit=True)
72     csv_file = open('Voting.csv', 'w', newline='', encoding='utf-8') #tweets_data_privacy
73     writer = csv.writer(csv_file) # opens csv file
74     fieldnames = ['created at', 'text', 'name',
75                  'screen name', 'location',
76                  'verified', 'followers count', 'friends count',
77                  'listed count', 'total favorites count', 'statuses count',
78                  'account creation date', 'retweet count',
79                  'favorite count'] # fieldnames
80     writer.writerow(fieldnames) # write fieldnames into header of the csv file
81     # fetches for tweets based on keywords and filters out any retweets. Tweepy's API will filter out user
82     # information.
83     for tweet in Cursor(api.search, q="Voting Election -filter:retweets AND -filter:replies", lang="en", since=start,
84                        until=end).items(150000):
85         writer.writerow([tweet.created_at, tweet.text.encode('unicode_escape'),
86                          tweet.user.name, tweet.user.screen_name, tweet.user.location,
87                          tweet.user.verified, tweet.user.followers_count,
88                          tweet.user.friends_count, tweet.user.listed_count,
89                          tweet.user.favourites_count, tweet.user.statuses_count,
90                          tweet.user.created_at, tweet.retweet_count, tweet.favorite_count])

```

```

91
92 # this method filters out anything unrelated.
93 def filter(self):
94     read = pd.read_csv("Voting.csv") # reads in CSV file tweets_data_privacy
95     df = pd.DataFrame(read) # turns CSV file into data frame
96     new_data = df[~df['name'].str.contains('news|News|.com|EducationLaw|Money|Roger_Verhoeven|
97                                           'Policy Pro\'s - Business Policy Writers|Coca-Cola GB|JD Supra|
98                                           'Privacy Watch|Giridhar|QV Develop|Data Vault|Privacy Professionals Briefly|
99                                           'PrivacyPro b|Privacy Tech Briefly|PrivacyTech b|Macking Essentials|Lost Coast Outpost|
100                                           'NortonRoseFulbright|Urdhva Tech|SSV Tribune|O.C. Register|Daniel Hanan|Anas Najeeb|report|
101                                           'Report|SAP Customer Data Cloud|ResponSight'))
102     # filters any names that contain any words that match with the words above.
103     new_data.to_csv('2.csv', index=False) # returns a new CSV file filtered_tweets
104

```

A JSON_TO_CSV.py and organize method was created to convert the live tweets from JSON to CSV files. The reason for making JSON_TO_CSV into another file is, during run time if you run output_live_tweets and organize methods in the same file, it will give a value error.

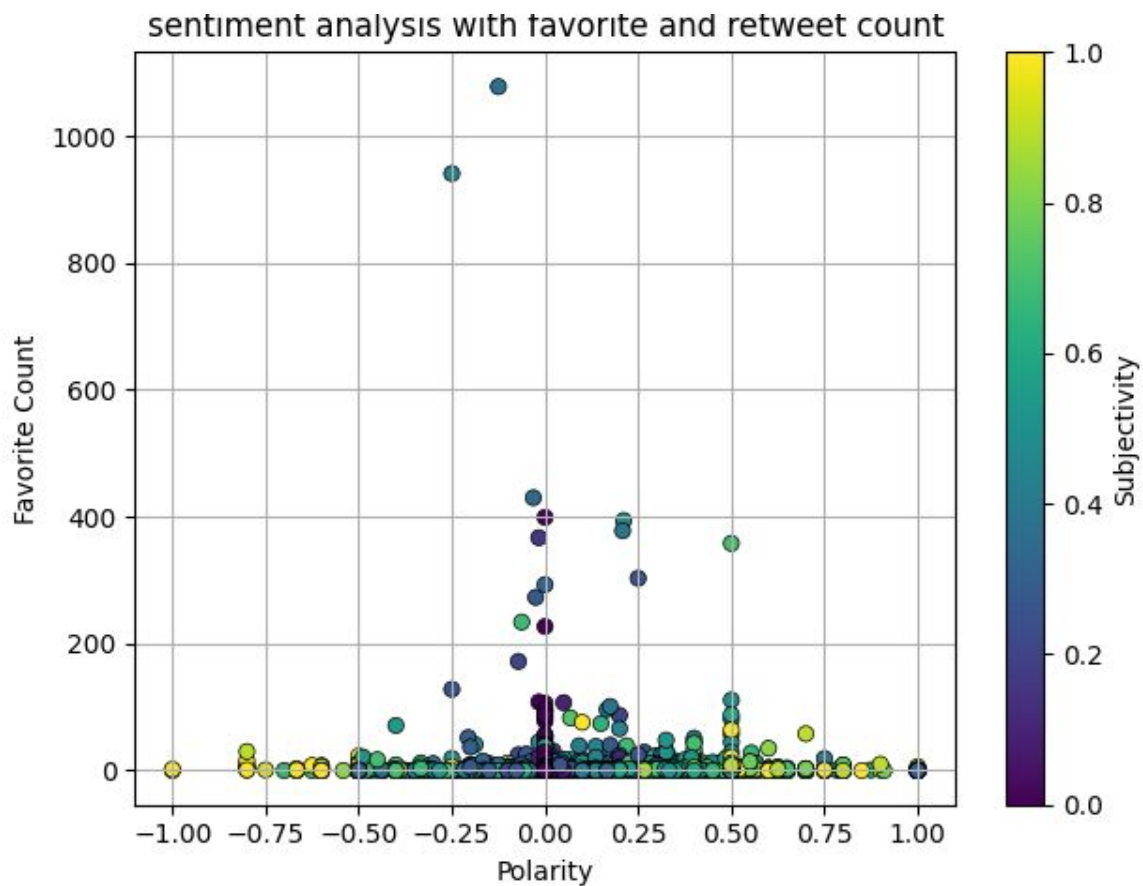
```

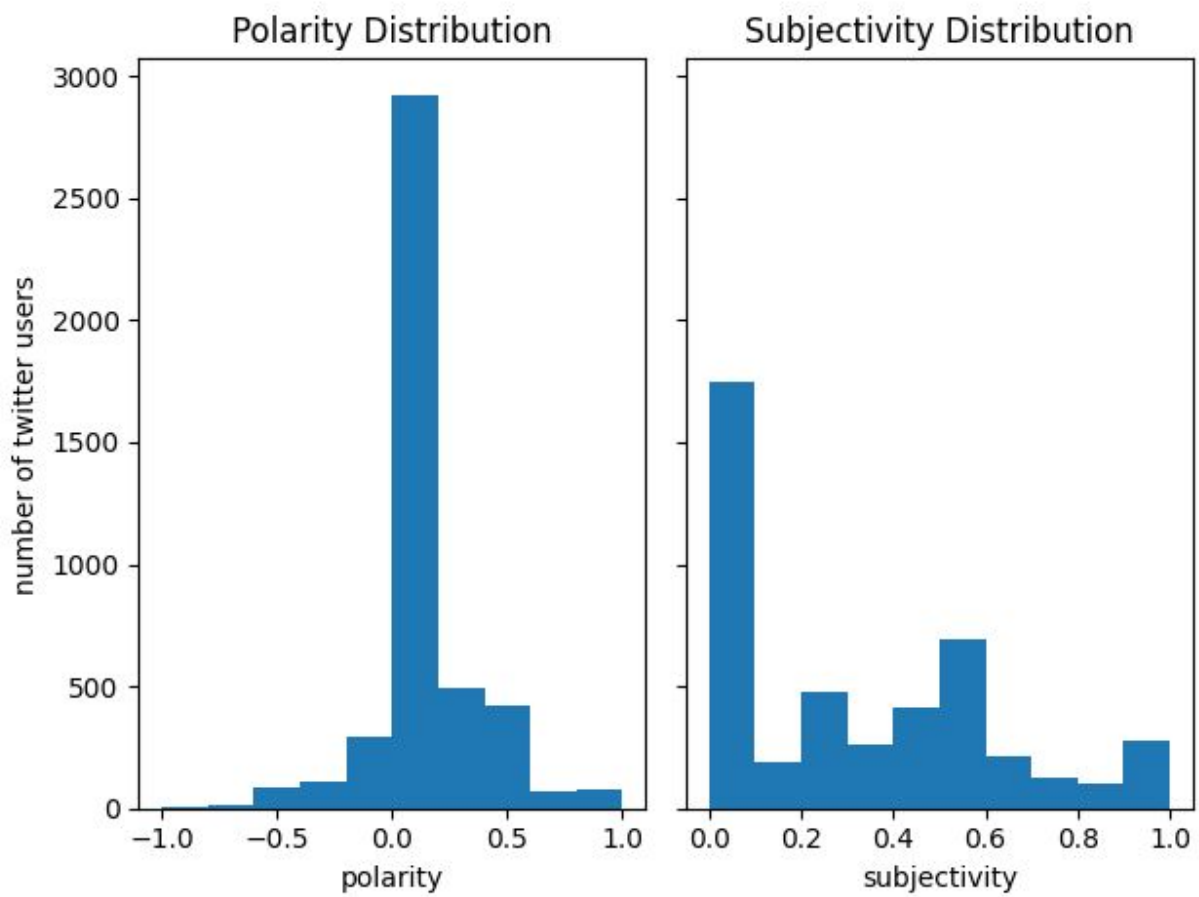
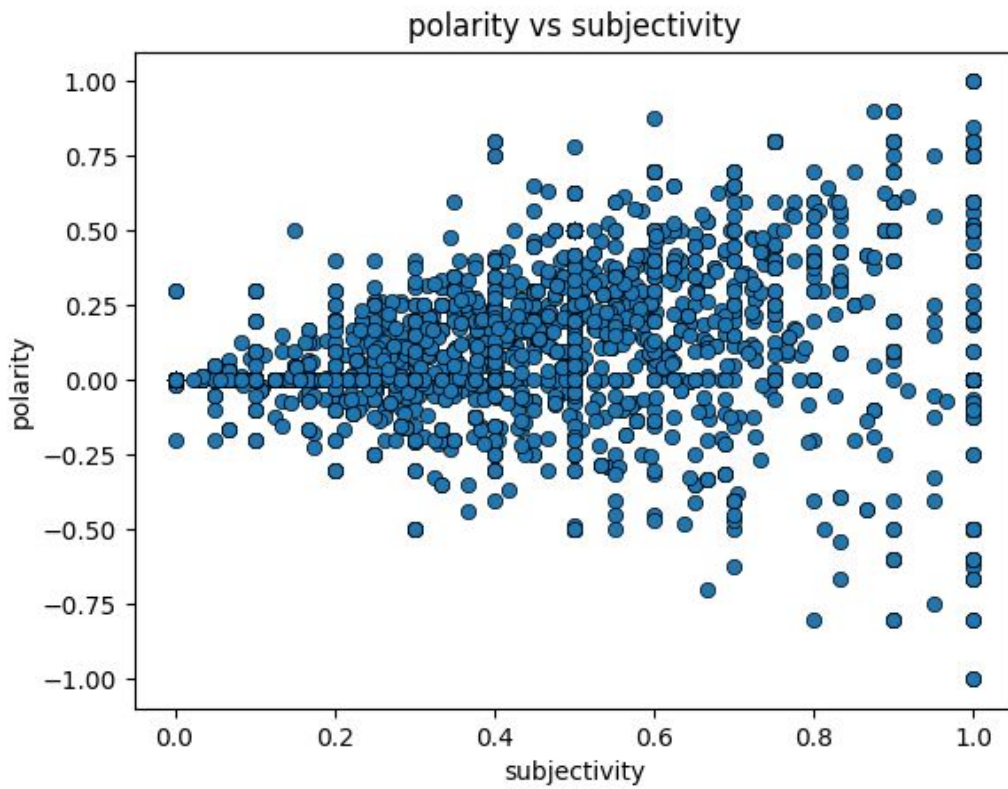
1 import json
2 import csv
3
4
5 def organize():
6     json_file = open('unfiltered_tweets_Donald_Trump.json', 'r', encoding="utf-8").read() # reads the JSON file
7     json_data = json.loads(json_file) # string becomes a json Python object
8     csv_file = open('filtered_tweets_Donald_Trump.csv', 'w', encoding="utf-8", newline='') # opens csv file
9     writer = csv.writer(csv_file) # creates csv writer object
10    fieldnames = ['created at', 'text', 'screen name',
11                 'location', 'verified', 'followers',
12                 'friends'] # fieldnames
13    writer.writerow(fieldnames) # writes fieldnames on the first row
14
15    # gets fields from the json object and writes a row into csv file
16    for tweet in json_data:
17        if tweet.get('lang') != "en": # if the tweet is not in english, don't write in csv file
18            continue
19        else:
20            writer.writerow([tweet.get('created_at'), tweet.get('text').encode('unicode_escape'),
21                             tweet.get('user').get('screen_name'), tweet.get('user').get('location'),
22                             tweet.get('user').get('verified'), tweet.get('user').get('followers_count'),
23                             tweet.get('user').get('friends_count')])
24    csv_file.close()
25
26
27 if __name__ == "__main__":
28     organize()

```

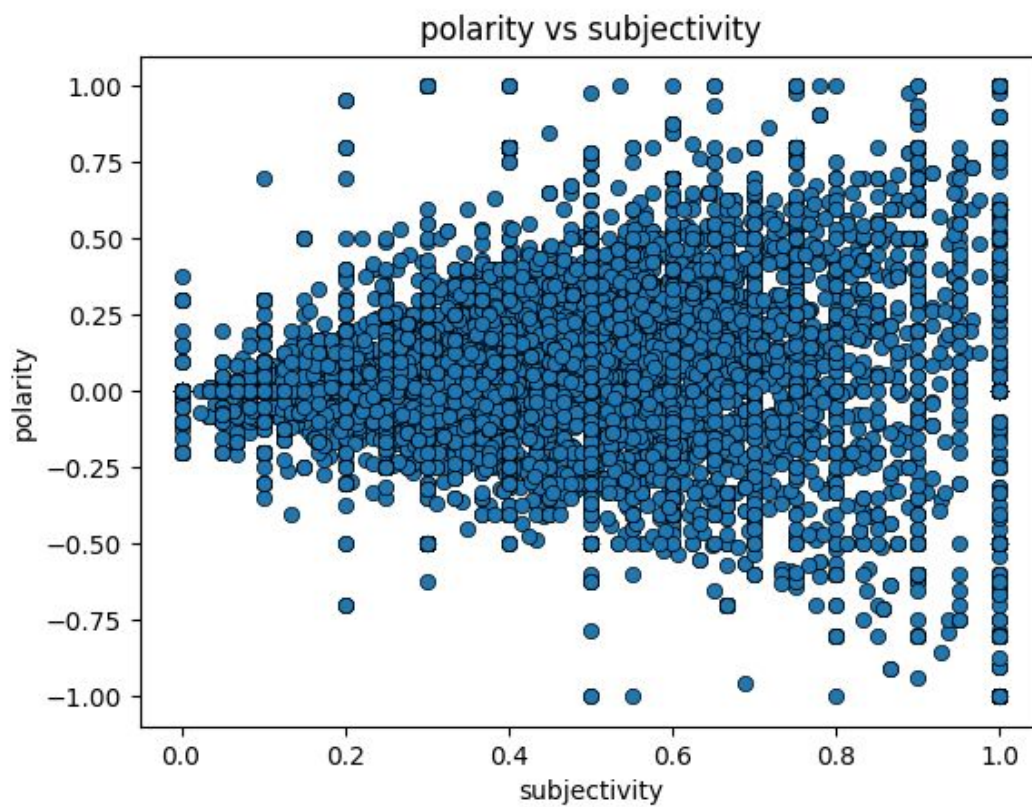
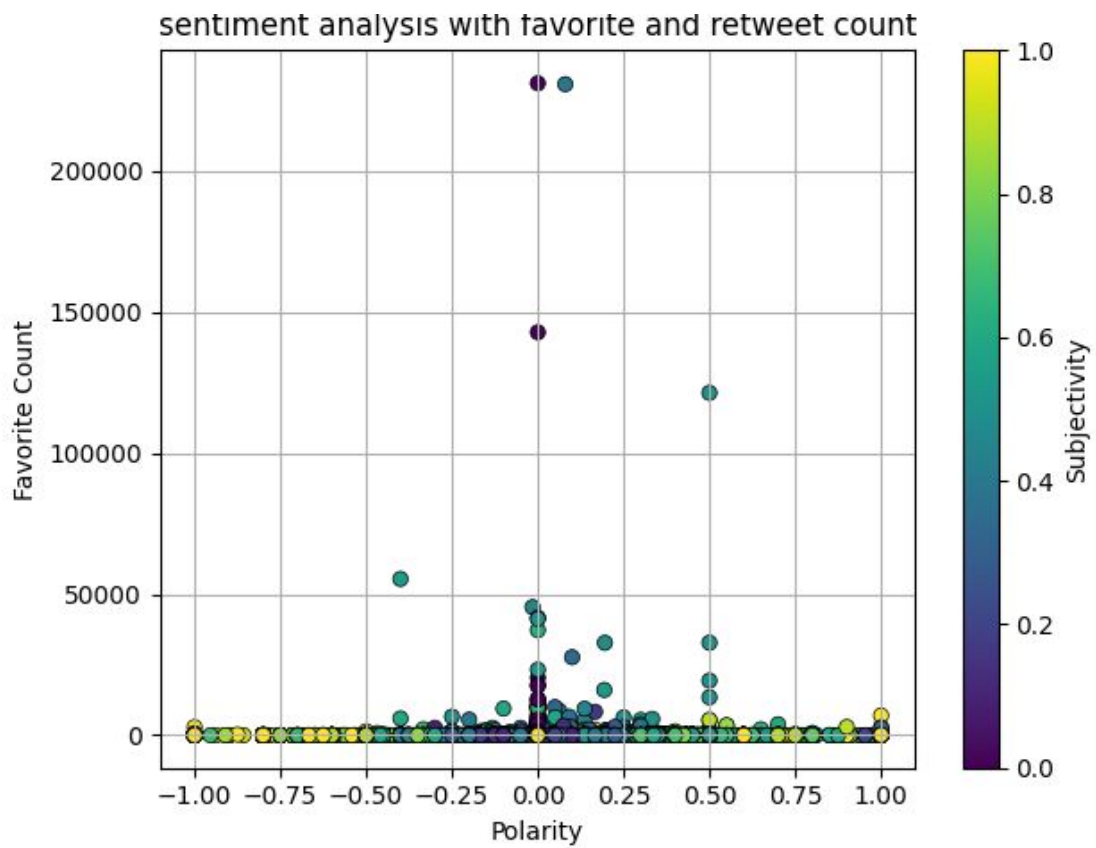
Next, we created a `Sentiment_analysis.py` file with a class called `Analysis`. The first method is called `polarity_and_subjectivity`. The tweets that were extracted were used to detect the polarity and subjectivity of the text. Then, polarity and subjectivity are classified into three ranges and will create a new csv file. The next three methods help create a scatter or histogram plot. The `subjectivity_class` method allows a subjectivity class column to classify the subjectivity into three classes(0,1,2) based on the three ranges. The `to_classification` methods take the features, followers count, friends count, retweet count, favorite count, polarity, SubjectivityClass and make a new csv file. The new csv file will be used for a machine learning classification algorithm to determine whether we can predict subjectivity. Based on the data extract from twitter, there are six charts represented for the data. Three are represented for the topic, data privacy and the other three for the 2020 election voting fiasco.

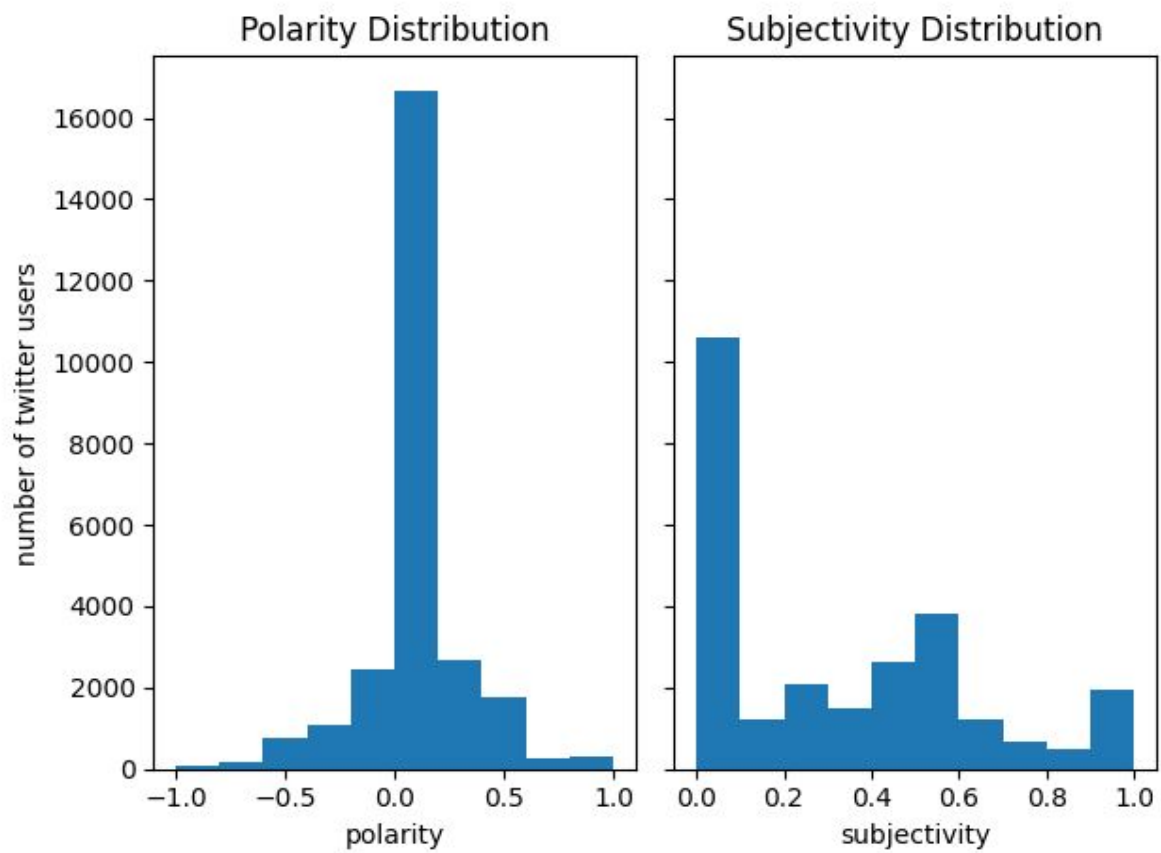
Data privacy:





2020 Election:





Next, the results for k nearest neighbor classification using five features to predict subjectivity or polarity.

Data privacy(subjectivity):

```
Run: KNN x
C:\Users\owner\Twitter_Project\Scripts\python.exe C:/Users/owner/PycharmProjects/Twitter_Project/KNN.py
[[491 25 9]
 [ 72 120 20]
 [ 63 28 71]]
      precision    recall  f1-score   support

     0       0.78       0.94       0.85       525
     1       0.69       0.57       0.62       212
     2       0.71       0.44       0.54       162

 accuracy          0.76          899
  macro avg       0.73       0.65       0.67          899
 weighted avg     0.75       0.76       0.74          899

Process finished with exit code 0
```

Data privacy(polarity):

```
Run: KNN x
C:\Users\owner\Twitter_Project\Scripts\python.exe C:/Users/owner/PycharmProjects/Twitter_Project/KNN.py
[[ 6 25 15]
 [ 9 545 67]
 [ 8 90 134]]

      precision    recall  f1-score   support

     0       0.26       0.13       0.17         46
     1       0.83       0.88       0.85        621
     2       0.62       0.58       0.60        232

 accuracy          0.76         899
 macro avg          0.57         899
 weighted avg       0.74         899

Process finished with exit code 0
|
```

2020 Voting Election(subjectivity):

```
Run: KNN x
C:\Users\owner\Twitter_Project\Scripts\python.exe C:/Users/owner/PycharmProjects/Twitter_Project/KNN.py
[[3001 145 104]
 [ 374 521 121]
 [ 421 163 393]]

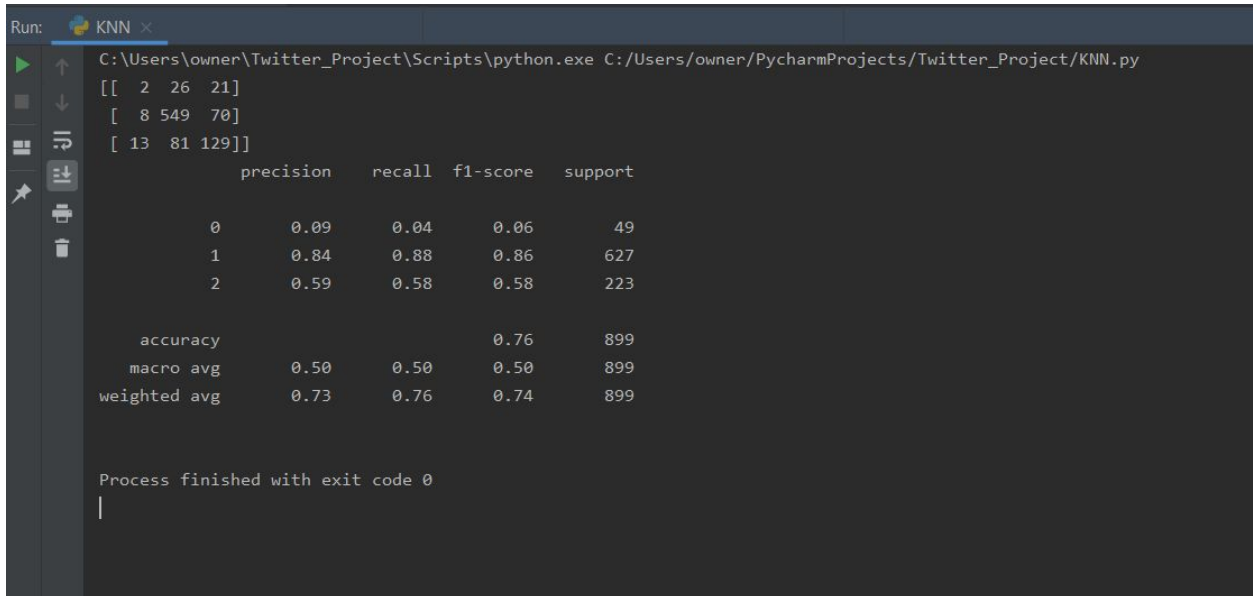
      precision    recall  f1-score   support

     0       0.79       0.92       0.85       3250
     1       0.63       0.51       0.56       1016
     2       0.64       0.40       0.49         977

 accuracy          0.75       5243
 macro avg          0.68       5243
 weighted avg       0.73       5243

Process finished with exit code 0
|
```

2020 Voting Election(polarity):



```
Run: KNN x
C:\Users\owner\Twitter_Project\Scripts\python.exe C:/Users/owner/PycharmProjects/Twitter_Project/KNN.py
[[ 2 26 21]
 [ 8 549 70]
 [13 81 129]]
      precision    recall  f1-score   support

    0       0.09      0.04      0.06         49
    1       0.84      0.88      0.86        627
    2       0.59      0.58      0.58        223

 accuracy          0.76         899
 macro avg          0.50         899
 weighted avg       0.73         899

Process finished with exit code 0
|
```

The screenshot shows a PyCharm Run console window titled 'KNN x'. The command executed is 'C:\Users\owner\Twitter_Project\Scripts\python.exe C:/Users/owner/PycharmProjects/Twitter_Project/KNN.py'. The output displays three rows of data: [[2 26 21], [8 549 70], [13 81 129]]. Below this, a table of model performance metrics is shown for three classes (0, 1, 2). The metrics include precision, recall, f1-score, and support. The overall accuracy is 0.76 for 899 samples. The macro average is 0.50 and the weighted average is 0.73. The process finished with exit code 0.

Based on the results above, we can conclude that it is possible to predict polarity and subjectivity. Since, the accuracy is 75 to 76 for both polarity and subjectivity.