# CIS 593 – BIG DATA

# LAB ASSIGNMENT – 3

## Name: Pooja A. Khatri

## CSU ID: 2783752

**Extra credit work:**

**automatic process for data collection, and storage and retrieval using mongoDB**

# 1. Platforms

Language used: Python

Python IDE: PyCharm
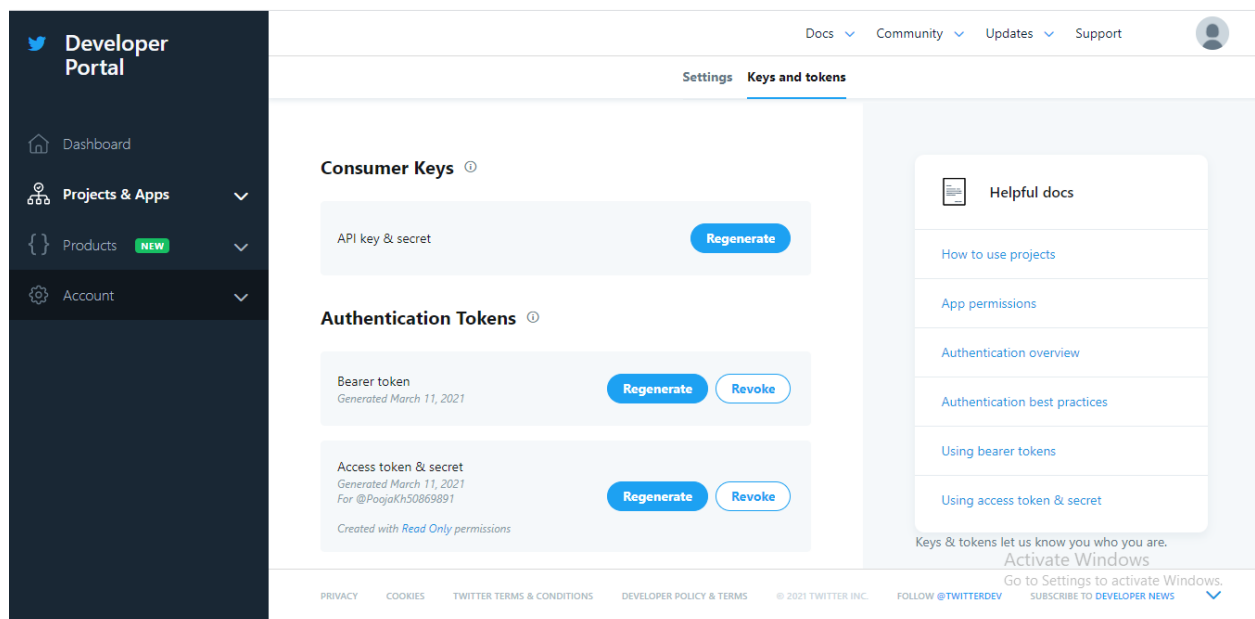
Database design tool: mongoDB

# 2. Files included

Main.py, QueryResults.xlsx showing query results for all queries in different sheets
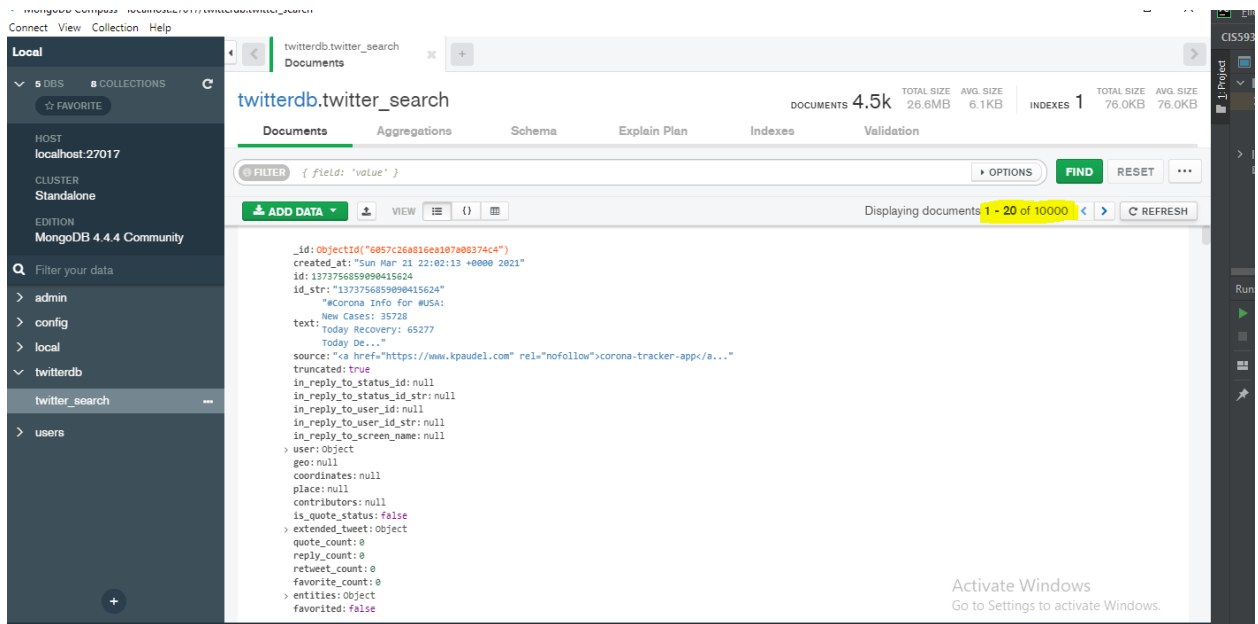
# 3. Steps performed for lab

- Sign-In to Twitter developer account and generate API tokens
- Fetch JSON from tweeter using RestAPI
- Parse all data in python and make mongoDB connection
- Import all data into mongoDB Database
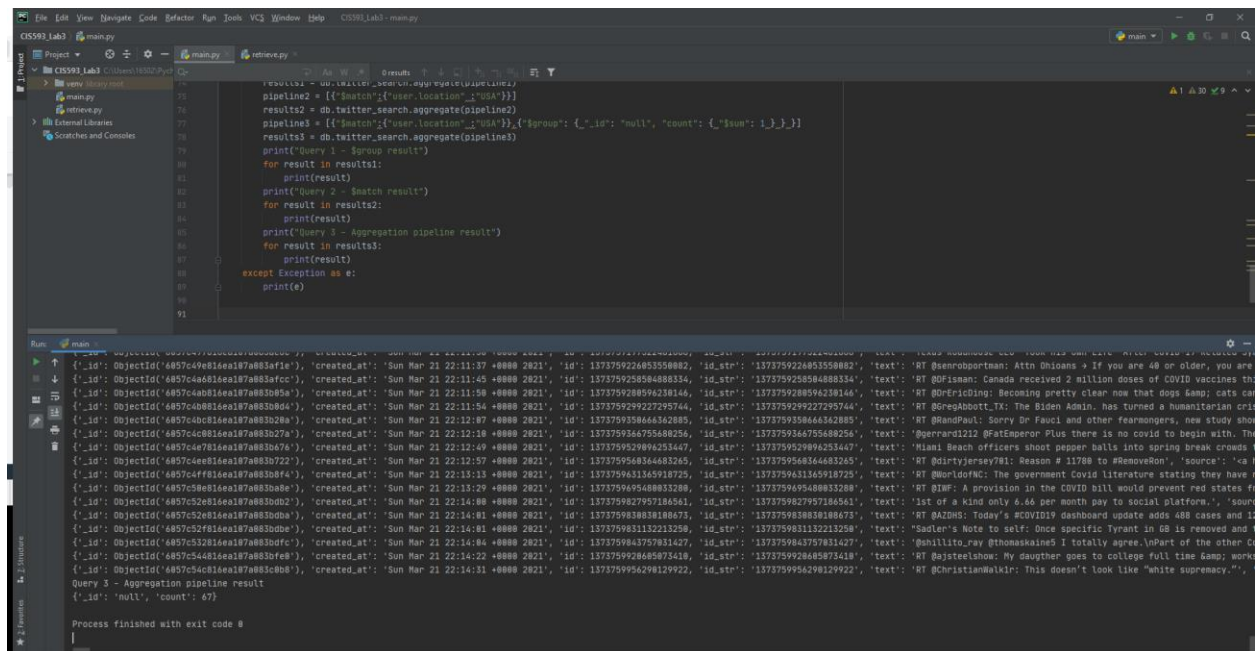- Retrieve data from mongoDB dynamically

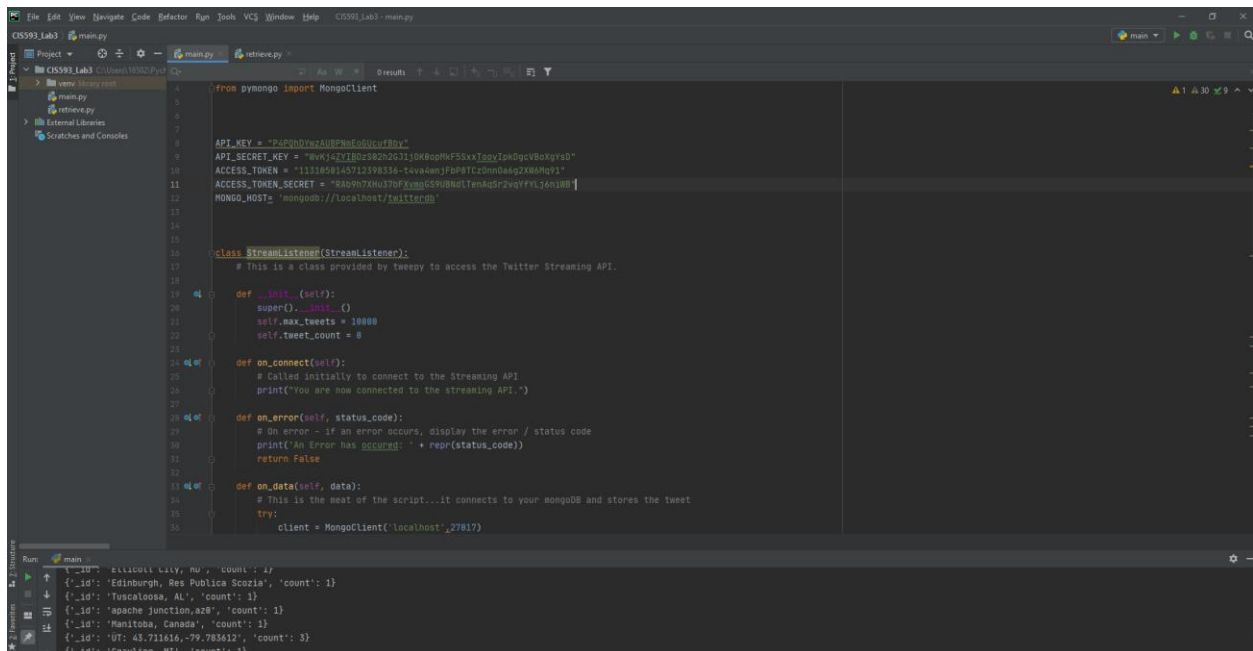# 4. Screenshots

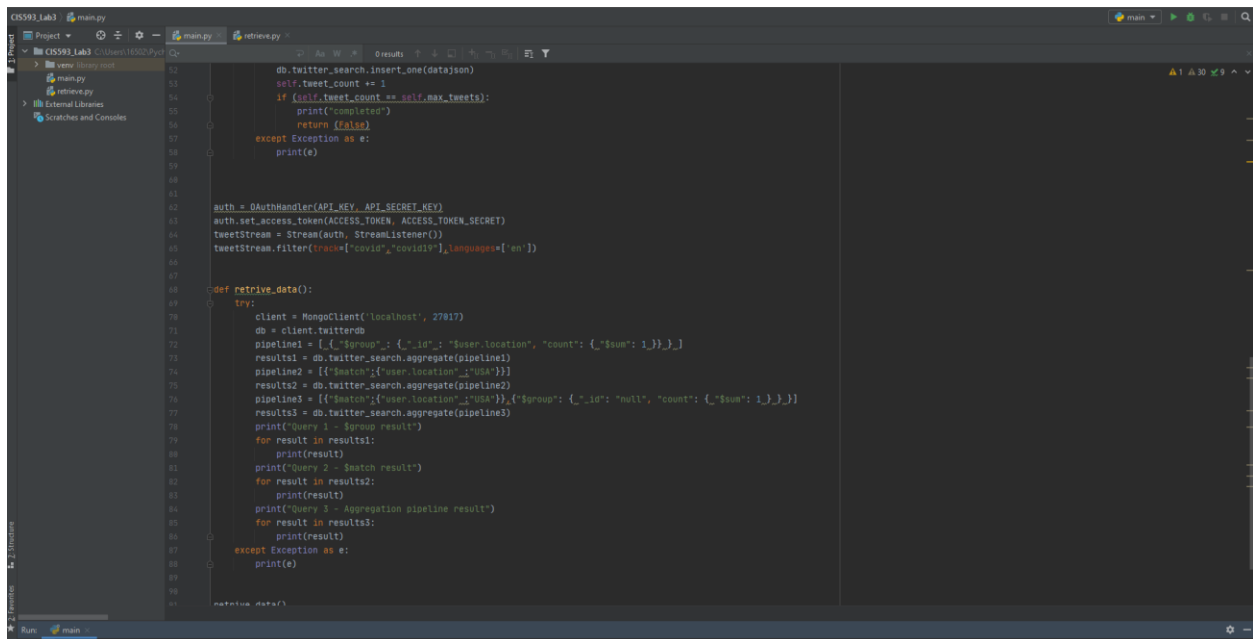Twitter developer account obtained successfully



mongoDB showing 10000 documents

Pycharm showing successful execution

Pycharm showing execution result of $group query - [ { "$group" : { "_id" : "$user.location", "count": { "$sum": 1 }} } ]

Pycharm showing execution result of $match query – [{"$match":{"user.location" :"USA"}}]



Pycharm showing execution result of Aggregation pipeline query - [{"$match":{"user.location" :"USA"}},{"$group": { "_id": "null", "count": { "$sum": 1 } } }]



mongoDB Shell showing few query results

Command Prompt - C:\mongodb\bin\mongo.exe

> use twitterdb
switched to db twitterdb
> db.twitter_search.count()
10000
> db.twitter_search.aggregate( [ { $group : { _id : "$user.location"} } ] )
{ "_id" : "Granada Hills, CA" }
{ "_id" : "Europe - The Netherlands" }
{ "_id" : "deptford, south-east london" }
{ "_id" : "South Florida" }
{ "_id" : "The World" }
{ "_id" : "Legazpi City, Albay" }
{ "_id" : "In my lane n under ya skin " }
{ "_id" : "Ellicott City, MD" }
{ "_id" : "Edinburgh, Res Publica Scozia" }
{ "_id" : "Tuscaloosa, AL" }
{ "_id" : "apache junction,az0" }
{ "_id" : "Manitoba, Canada" }
{ "_id" : "ÜT: 43.711616,-79.783612" }
{ "_id" : "Grayling, MI" }
{ "_id" : "Trinidad" }
{ "_id" : "92058" }
{ "_id" : "the turn of the century" }
{ "_id" : "kernow" }
{ "_id" : "Manaus, Brasil" }
{ "_id" : "Thailand" }
Type "it" for more
> db.twitter_search.aggregate([{$match:{text: /corona|covid/i}},{ $group: { _id: null, count: { $sum: 1 } } }])
{ "_id" : null, "count" : 5777 }
> db.twitter_search.aggregate([{$match:{"user.location" :"canada"}}])
{ "_id" : ObjectId("6057c452816ea107a083a75c"), "created_at" : "Sun Mar 21 22:10:21 +0000 2021", "id" : NumberLong("1373758907869446149"), "id_str" : "13737589078694461
49", "text" : "RT @qtestever: friendly reminder that men who have sex with men aren't allowed to donate blood unless they've been abstinent for 3 months (…", "source" :
"<a href=\"http://twitter.com/download/iphone\" rel=\"nofollow\">Twitter for iPhone</a>", "truncated" : false, "in_reply_to_status_id" : null, "in_reply_to_status_id_s
tr" : null, "in_reply_to_user_id" : null, "in_reply_to_user_id_str" : null, "in_reply_to_screen_name" : null, "user" : { "id" : 473521047, "id_str" : "473521047", "name
" : "rebecca▒▒", "screen_name" : "rebeccaaboone", "location" : "canada", "url" : "http://instagram.com/makeupbyrebeccab", "description" : "24 / insta: rebecca.boone / m
akeupbyrebeccab", "translator_type" : "none", "protected" : false, "verified" : false, "followers_count" : 165, "friends_count" : 92, "listed_count" : 3, "favourites_co
unt" : 20679, "statuses_count" : 11579, "created_at" : "Wed Jan 25 02:14:32 +0000 2012", "utc_offset" : null, "time_zone" : null, "geo_enabled" : true, "lang" : null, "
contributors_enabled" : false, "is_translator" : false, "profile_background_color" : "668199", "profile_background_image_url" : "http://abs.twimg.com/images/themes/them
e1/bg.png", "profile_background_image_url_https" : "https://abs.twimg.com/images/themes/theme1/bg.png", "profile_background_tile" : false, "profile_link_color" : "60A6A
C", "profile_sidebar_border_color" : "000000", "profile_sidebar_fill_color" : "EEDCF5", "profile_text_color" : "860CCC", "profile_use_background_image" : false, "profil
e_image_url" : "http://pbs.twimg.com/profile_images/1363951711723941890/J8vQh-Js_normal.jpg", "profile_image_url_https" : "https://pbs.twimg.com/profile_images/13639517
11723941890/J8vQh-Js_normal.jpg", "profile_banner_url" : "https://pbs.twimg.com/profile_banners/473521047/1614026405", "default_profile" : false, "default_profile_image
" : false, "following" : null, "follow_request_sent" : null, "notifications" : null }, "geo" : null, "coordinates" : null, "place" : null, "contributors" : null, "retwe
eted_status" : { "created_at" : "Sun Mar 21 01:37:24 +0000 2021", "id" : NumberLong("1373448625762697218"), "id_str" : "1373448625762697218", "text" : "friendly reminde
r that men who have sex with men aren't allowed to donate blood unless they've been abstinent for 3… https://t.co/GHKmZVbXSY", "source" : "<a href=\"http://twitter.com/

## 5. Code

```python
from tweepy import Stream, OAuthHandler
import json
from tweepy.streaming import StreamListener
from pymongo import MongoClient


API_KEY = "P4PQhDYwzAUBPNmEoGUcufBby"
API_SECRET_KEY = "WvKj4ZYIBDzS02h2GJ1jOK0opMkF5SxxToovIpkDgcVBoXgYsD"
ACCESS_TOKEN = "1131050145712398336-t4va4wnjFbP8TCzOnnOa6g2XW6Mq91"
ACCESS_TOKEN_SECRET = "RAb9h7XHu37bFXvmoGS9UBNdlTenAqSr2vqYfYLj6niWB"

MONGO_HOST= 'mongodb://localhost/twitterdb'


class StreamListener(StreamListener):
    # This is a class provided by tweepy to access the Twitter Streaming
API.

    def __init__(self):
        super().__init__()
        self.max_tweets = 10000
        self.tweet_count = 0

    def on_connect(self):
        # Called initially to connect to the Streaming API
        print("You are now connected to the streaming API.")

    def on_error(self, status_code):
```

```python
        # On error - if an error occurs, display the error / status code
        print('An Error has occured: ' + repr(status_code))
        return False

    def on_data(self, data):
        # This is the meat of the script...it connects to your mongoDB and
stores the tweet
        try:
            client = MongoClient('localhost',27017)

            # Use twitterdb database. If it doesn't exist, it will be
created.
            db = client.twitterdb

            # Decode the JSON from Twitter
            datajson = json.loads(data)

            # grab the 'created_at' data from the Tweet to use for display
            created_at = datajson['created_at']

            # print out a message to the screen that we have collected a
tweet
            print("Tweet collected at " + str(created_at))

            # insert the data into the mongoDB into a collection called
twitter_search
            # if twitter_search doesn't exist, it will be created.
            db.twitter_search.insert_one(datajson)
            self.tweet_count += 1
            if (self.tweet_count == self.max_tweets):
                print("completed")
                return (False)
        except Exception as e:
            print(e)


auth = OAuthHandler(API_KEY, API_SECRET_KEY)
auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)
tweetStream = Stream(auth, StreamListener())
tweetStream.filter(track=["covid","covid19"],languages=['en'])


def retrive_data():
    try:
        client = MongoClient('localhost', 27017)
        db = client.twitterdb
        pipeline1 = [ { "$group" : { "_id" : "$user.location", "count": {
"$sum": 1 }} } ]
        results1 = db.twitter_search.aggregate(pipeline1)
        pipeline2 = [{"$match":{"user.location" :"USA"}}]
        results2 = db.twitter_search.aggregate(pipeline2)
        pipeline3 = [{"$match":{"user.location" :"USA"}},{"$group": {
"_id": "null", "count": { "$sum": 1 } } }]
        results3 = db.twitter_search.aggregate(pipeline3)
        print("Query 1 - $group result")
        for result in results1:
```

```python
            print(result)
        print("Query 2 - $match result")
        for result in results2:
            print(result)
        print("Query 3 - Aggregation pipeline result")
        for result in results3:
            print(result)
    except Exception as e:
        print(e)


retrive_data()
```