

-----Module 1(Data Collection)-----

Collection data from twitter for legitimate users

For collecting the data from twitter I am using Tweepy module

For that We need Counsumer_KEY, Counsumer_secret_KEY, Access_token, Access_Token

That all I can get from twitter app.devloper where we need to sign in and make an account and then twitter generate the keys

After that a simple program in python can extract data from twitter in given limit by twitter.

```
In [1]: import tweepy
import pandas as pd
import time
import numpy as np
import matplotlib.pyplot as plt
from tweepy import Stream
from tweepy.streaming import StreamListener
import numpy as np
import matplotlib.pyplot as plt
from sklearn import svm
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
from sklearn.utils.multiclass import unique_labels
```

Connection Authentication

```
In [3]: consumer_key = 'd9Ksoz6Wb1jD0mqbW8rjaSNb7'
consumer_secret = 'pHXnVSJeLb0xaY1bOR7BWFdDNhZSF6IzegZV87qUSUqy6Qe8qG'
access_token = '3648603434-dGRu1nHet22tdoYeqaAGoN8MyZrNw9oXZQvGZUD'
access_token_secret = 'PZ8pcQBCb5zVPLRQNVQZc3Yzi0rz1wPef607R07gzcv0f'

auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_token, access_token_secret)

api = tweepy.API(auth, wait_on_rate_limit=True)
print(api)

<tweepy.api.API object at 0x000002050123E9A0>
```

Collecting Data list of Friends of a given screen_name

Save data in txt file

```
In [3]: '''print('Name of the Friends of user')
friends = []
sn=['PoliceRajasthan','Uppolice','MumbaiPolice','PunjabPoliceInd','KolkataPolice','DelhiPolice','BlrCityPolice','noidapolice'
print(len(sn))
for name in sn:
    for friend in tweepy.Cursor(api.friends, screen_name = name).items(20):
        try:
            friends.append(friend.screen_name)
            print(friend.screen_name)
            time.sleep()
        except Exception as e:
            pass

print("Writing to file Luser...")
with open("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/luser.txt", "a") as f:
    for item in friends:
        f.write("%s\n" % item)
print("Writed Succesfully completed...")

#output
Name of the Friends of user
19
shomarredpally
ShoMarket
shonarayanguda
shomahankali
shokulsumpura
shokarkhana
shokanchanbagh
shocharminar
shochaderghat
shochilkalguda
shogopalapuram
shobegumpet
shoafzalgunj
shoamberpet
Prabhamudhiraj
rambanoth
shocgt
bhanu02468
shobowenpally
KTRTRS
narendramodi
WomenCid
Writing to file Luser...
Writed Succesfully completed...'''
```

```
Out[3]: 'print(\'Name of the Friends of user\')\nfriends = []\n\nsn=[\'PoliceRajasthan\',\'Uppolice\',\'MumbaiPolice\',\'PunjabPoliceI
nd\',\'KolkataPolice\',\'DelhiPolice\',\'BlrCityPolice\',\'noidapolice\',\'igrangemeerut\',\'noidatraffic\',\'adgzonemeerut
\',\'meerutpolice\',\'bulandshahrpol\',\'saharanpurpol\',\'shamlipolice\',\'hapurpolice\',\'baghpatpolice\',\'chennaipolice_
\',\'hydcitypolice\']\n\nprint(len(sn))\n\nfor name in sn:\n    for friend in tweepy.Cursor(api.friends, screen_name = name).ite
ms(20):\n        try:\n            friends.append(friend.screen_name)\n            print(friend.screen_name)\n            ti
me.sleep()\n        except Exception as e:\n            pass\n\nprint("Writing to file Luser...")\n\nwith open
("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/luser.txt", "a") as f:\n    for item in friends:\n
f.write("%s\n" % item)\n\nprint("Writed Succesfully completed...")\n\n\n#output\nName of the Friends of user\n19\n\nshomarredpal
ly\nShoMarket\nshonarayanguda\nshomahankali\nshokulsumpura\nshokarkhana\nshokanchanbagh\nshocharminar\nshochaderghat\nshochi
lkalguda\nshogopalapuram\nshobegumpet\nshoafzalgunj\nshoamberpet\nPrabhamudhiraj\nrambanoth\nshocgt\nnbhanu02468\nshobowenpal
ly\nKTRTRS\nnarendramodi\nsmittal_ips\nshahalibanda\nshomusheerabad\nHCSC_Hyd\nvishwa_raghu\nChowraastaM\nNGopishetty\nravit
ollywood\nKhammamCp\nPuneCityPolice\nSmitaSabharwal\npassportsevamea\nOfficeOfVKS\nTelanganaHealth\nTsspdclCorporat\nsho_mog
halpura\nCPMumbaiPolice\nOnlineBlooDonor\nspmahabubabad\nmhbdpolice\nshomadannapet\nPIBFactCheck\nghmc_adw\nTSEduDept\nMAHal
eem20\nspvikarabad\nCYBTRAFFIC\nipstelangana\nTelanganaCOPs\nMORTHRoadSafety\nTSPSCofficial\nTSMAUDOnline\nTSCSOffice\nTSCon
sumers\nBlrCityPolice\nSpKothagudem\nlrvr1974\nMLA54327644\nWomenCid\n\nWriting to file Luser...\nWrited Succesfully complete
d...'
```

Now collect 30 tweet from each user in luser file that I extracted from twitter

```
In [ ]: Total_Data = []
fo = open("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/luser.txt", "r")
f = fo.readlines()
fo.close()
dataset = map(lambda s: s.strip(),f)
try:
    for datavar in dataset:
        data = api.get_user(datavar)
        counter = 0
        for status in tweepy.Cursor(api.user_timeline, id = datavar).items(30):
            try:
                counter= counter+1
                Total_Data.append(status)
                time.sleep()
            except Exception as e:
                pass
except Exception as e:
    pass
print(len(Total_Data))
```

Now from each tweet extract useful atributes that I need later for classification

```
In [5]: import urllib.parse
import pandas as pd

def process_http(string):
    url_count = 0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if s and n:
            url_count += 1
    return url_count

def process_hashtag(string):
    hashtag_count = 0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if i[:1] == '#':
            hashtag_count += 1
    return hashtag_count

def process_mention(string):
    mention_count=0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if i[:1] == '@':
            mention_count += 1
    return mention_count

def process_data(Total_Data):
    TwittID = [tweet.id for tweet in Total_Data]
    # Making the dataset in pandas frame
    Data = pd.DataFrame(TwittID, columns = ['TwittID'])
    # processing the data in Tweet Level

    Data["TextData"] = [tweet.text for tweet in Total_Data]
    Data["TweetCreatedAt"] = [tweet.created_at for tweet in Total_Data]
    Data["RetweetCount"] = [tweet.retweet_count for tweet in Total_Data]
    Data["TweetFavouriteCount"] = [tweet.favorite_count for tweet in Total_Data]
    Data["TweetSource"] = [tweet.source for tweet in Total_Data]

    # processing the data in User Graph Level

    Data["UserID"] = [tweet.author.id for tweet in Total_Data]
    Data["UserScreenName"] = [tweet.author.screen_name for tweet in Total_Data]
    Data["UserName"] = [tweet.author.name for tweet in Total_Data]
    Data["UserCreatedAt"] = [tweet.author.created_at for tweet in Total_Data]
    Data["UserDescription"] = [tweet.author.description for tweet in Total_Data]
    Data["UserDescriptionLength"] = [len(tweet.author.description) for tweet in Total_Data]
    Data["UserFollowersCount"] = [tweet.author.followers_count for tweet in Total_Data]
    Data["UserFriendsCount"] = [tweet.author.friends_count for tweet in Total_Data]
    Data["UserLocation"] = [tweet.author.location for tweet in Total_Data]

    # Data["url"] = [tweet.author.url for in Total_Data]
    # Data["User_mention"] = [user_mentions.author.screen_name for tweet in Total_Data]
    # Data["HashTag"] = [hashtag.text for tweet in Total_Data]

    Data["HttpCount"] = [process_http(tweet.text) for tweet in Total_Data]
    Data["HashtagCount"] = [process_hashtag(tweet.text) for tweet in Total_Data]
    Data["MentionCount"] = [process_mention(tweet.text) for tweet in Total_Data]
    Data["TweetCount"] = [tweet.author.statuses_count for tweet in Total_Data]
    return Data
Data = process_data(Total_Data)
Data.shape
```

Out[5]: (330, 19)

```
In [6]: Data.tail(4)
```

Out[6]:

	TwittID	TextData	TweetCreatedAt	RetweetCount	TweetFavouriteCount	TweetSource	UserID	UserScreenName	UserName	User
326	1191764131088437248	RT @shailyIPSspeaks: We have seen this. We hav...	2019-11-05 17:08:05	13	0	Twitter for Android	371293553	pihuprasad	IPS YAMUNA	2
327	1191763983344140288	RT @jatinnarwal1: Its not about which professi...	2019-11-05 17:07:30	896	0	Twitter for Android	371293553	pihuprasad	IPS YAMUNA	2
328	1191763742935019522	RT @Hii_VT: Let's learn how to burn policemen ...	2019-11-05 17:06:32	12	0	Twitter for Android	371293553	pihuprasad	IPS YAMUNA	2
329	1191763250133659648	RT @IPS_KTK_Assn: We strongly condemn incident...	2019-11-05 17:04:35	31	0	Twitter for Android	371293553	pihuprasad	IPS YAMUNA	2

Save data in csv_files in Leg_user.csv

```
In [7]: import sys
# Saving data with item space separating
Data.to_csv('C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/Leg_user.csv', sep=',', encoding='utf8')
```

extracting Spam data from twitter by searching @spam and find out the user for reported that hypothesis is that there is highly chances is that that user be fake

We can later analyse by text any volgor word and find our later first like legitimate user we have to collect data for spammer too

In [12]: """# printing all the friends names of the user

```
friends = []
class listener(StreamListener):
    def on_data(self, data):
        try:
            tweet = data.split(', "screen_name":')[1].split(', "location')[0]
            print(tweet)
            friends.append(tweet)
            return True
        except BaseException as e:
            print('failed on data' + str(e))
            time.sleep(5)
    def on_error(self, status):
        print(status)

twitterStream = Stream(auth, listener())
try:
    for x in range(1,10):
        twitterStream.filter(track=["cougar"])
except KeyboardInterrupt:
    print("Key board interruption")
with open("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/suser.txt", "w") as f:
    for item in friends:
        f.write("%s\n" % item)
```

output:

josef_morasch
hyenafluff
MattLov14455762
HerFabulousWays
ShatteredReaper
HornyScotsman98
kfkffmdm
kfkffmdm
MdmMcCoy
bradley33885052
pollo0014
tribink1
MattLov14455762
kyte_dave
uDerTheh0od
Steffff32720432
kfkffmdm
eli_boggs
_vaeXvae
reno7817
reno7817
Emir20665707
reno7817
pollo0014
carloshdez21
Sara19275758
CTG_tsutsumi
pollo0014
JDublim
ctg_machida
Roberto441220
OkomeeRice
profbadoor
pollo0014
SlaveTraderJoes
Amstron66093131
Masterxii1_4
Sexypro30863666
Onlyfan13594367
1st1966
huggyair
RosePro47247034
Steffff32720432
pollo0014
cougarmilfclub
JDublim
masterbooblover
Alex11little
Dutch_DPP
urbandictionary
nicolas64091667
luiginked
nicolas64091667
mommatiff3
mommatiff3
streitzfrank
pollo0014
beatsbygibbs
noel6132
lionloins
LisaCochonne
pollo0014
JOHNMcILHONE2
BLOODSADx
PauldingRadio

justasksam3
OGKING04987433
winekunTM
LuisEnri4400004
WhenSus8
alexand19024887
pollo0014
BbcRon_
PipPromos
bassshake
Para_Dox207
winekunTM
LopedeVegan269
piranhajams
tytheestallio
KevinBr96780686
CougargamingE
flvwrboy
LouisePixie1
alexabrines
pollo0014
Markanthony315
tytheestallio
cougar_an
naughtysussex
Christi07528537
Gsanchez7821
pollo0014
doodlebug247
ps236236
ps236236
hornynigga1521
saveurdesiles28
1041zion
oozytube
geezporn
saveurdesiles28
joohwangblink
EDG_Group
browniej45
saveurdesiles28
joohwangblink
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
saveurdesiles28
browniej45
browniej45
browniej45
RootNationUA
browniej45
browniej45
browniej45
justasksam3
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
browniej45
CurvesDelicious
browniej45
browniej45
browniej45
pollo0014
TrylineUK
BaidinNoer
boobhunters1
browniej45
33Metalson
AzurisDraws
pollo0014
GregMooney2
justasksam3
mllgigi
Nasser61597438
Itsimim
SamJenk60032204
pollo0014
Hah52059025
BbcRon_
thibodeaux_nora
PPPvideos1

scully29jamie
Hah52059025
turtletesticles
pollo0014
ninacelibataire
adoptunecougar
opium65
pornlegendsclub
cougarmilfclub
Tightrope_Postm
Baddies74810961
tazats1
TigerShoutout
SH63927626
1041zion
sexxebuddah69
pollo0014
slo159
InsaMarquez
oozytube
geezporn
truckpornxxx
Ben00182369
olddrive69
Paradise000x
himansu92239715
cmu_Eleasar
mjohns2444
saveurdesiles28
wvmilfhunter
pollo0014
mjohns2444
mjohns2444
oozytube
geezporn
Daniel_S013
KinkyCuriousHo1
Daniel_S013
PervertDiaries
WifeisBitchPorn
oozytube
geezporn
pollo0014
DQpararaeKWS
Aaronnowak14
oozytube
geezporn
nate221217
Recaster20
oozytube
geezporn
pollo0014
fast_girlswild
SilverSeductres
donne_mature
BbcRon_
TheGILFNextDoor
james68402115
ktdenise
TheFallenLady_
ovelhaJorge
pollo0014
Qazi04968557
WiccaNymph
mac2002
Socrate90341168
johnlamanna3
mpwCfbFldSQ0ZTA
ZGLiraxD16
WAVMedia
Hornybo68150984
pollo0014
IrisScarpinelli
Infr3quentAddik
Aaronnowak14
unucla
Srinu50041745
Srinu50041745
louisejacksonuk
Andy_Lifeguard
Ym78200
meghan_aleo
Asuro_Fluid
Uhhmm26387032
Asuro_Fluid
acupof_TAY
ballc21
cougarmilfclub
pollo0014
sindra38503651
stevenhead38
blindromance_
Canijo74474417

pollo0014
rdvlibertins
pollo0014
HErottico
JamesALogan1
BenoitJoris
Gio58652487
18Beastman
JacobStockGuru
BriW74
kamdiv49
Emir20665707
Emir20665707
Belzouille1
AndeeSC2
KPromention
PassionCougar
promotion_feet
SexyPro18754447
dreinschisspur
hellsivle
pollo0014
carolbaker10023
jay818jay
Alendil_fr
ESPINARjos2
StadiumPrint
cougar_an
MurphyDuke_
MurphyDuke_
MurphyDuke_
funinlillian
jav_superfan"""


```
In [4]: Total_Data = []
fo = open("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/suser.txt", "r")
f = fo.readlines()
fo.close()
dataset = map(lambda s: s.strip(),f)
try:
    for datavar in dataset:
        data = api.get_user(datavar)
        counter = 0
        for status in tweepy.Cursor(api.user_timeline, id = datavar).items(10):
            try:
                counter= counter+1
                Total_Data.append(status)
                print(Total_data)
                time.sleep()
            except Exception as e:
                pass
except Exception as e:
    pass
print(len(Total_Data))

353
```

Now from tweet extract useful attributes

```
In [20]: import urllib.parse
import pandas as pd

def process_http(string):
    url_count = 0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if s and n:
            url_count += 1
    return url_count

def process_hashtag(string):
    hashtag_count = 0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if i[:1] == '#':
            hashtag_count += 1
    return hashtag_count

def process_mention(string):
    mention_count=0
    for i in string.split():
        s, n, p, pa, q, f = urllib.parse.urlparse(i)
        if i[:1] == '@':
            mention_count += 1
    return mention_count

def process_data(Total_Data):
    TwittID = [tweet.id for tweet in Total_Data]
    # Making the dataset in pandas frame
    Data = pd.DataFrame(TwittID, columns = ['TwittID'])
    # processing the data in Tweet Level

    Data["TextData"] = [tweet.text for tweet in Total_Data]
    Data["TweetCreatedAt"] = [tweet.created_at for tweet in Total_Data]
    Data["RetweetCount"] = [tweet.retweet_count for tweet in Total_Data]
    Data["TweetFavouriteCount"] = [tweet.favorite_count for tweet in Total_Data]
    Data["TweetSource"] = [tweet.source for tweet in Total_Data]

    # processing the data in User Graph Level

    Data["UserID"] = [tweet.author.id for tweet in Total_Data]
    Data["UserScreenName"] = [tweet.author.screen_name for tweet in Total_Data]
    Data["UserName"] = [tweet.author.name for tweet in Total_Data]
    Data["UserCreatedAt"] = [tweet.author.created_at for tweet in Total_Data]
    Data["UserDescription"] = [tweet.author.description for tweet in Total_Data]
    Data["UserDescriptionLength"] = [len(tweet.author.description) for tweet in Total_Data]
    Data["UserFollowersCount"] = [tweet.author.followers_count for tweet in Total_Data]
    Data["UserFriendsCount"] = [tweet.author.friends_count for tweet in Total_Data]
    Data["UserLocation"] = [tweet.author.location for tweet in Total_Data]

    # Data["url"] = [tweet.author.url for in Total_Data]
    # Data["User_mention"] = [user_mentions.author.screen_name for tweet in Total_Data]
    # Data["HashTag"] = [hashtag.text for tweet in Total_Data]

    #Data["HttpCount"] = [process_http(tweet.text) for tweet in Total_Data]
    #Data["HashtagCount"] = [process_hashtag(tweet.text) for tweet in Total_Data]
    #Data["MentionCount"] = [process_mention(tweet.text) for tweet in Total_Data]
    #Data["TweetCount"] = [tweet.author.statuses_count for tweet in Total_Data]
    return Data
Data = process_data(Total_Data)
Data.head()
```

Out[20]:

	TwittID	TextData	TweetCreatedAt	RetweetCount	TweetFavouriteCount	TweetSource		UserID	UserScreenName	Us
0	1396764493443698695	Okay Babe https://t.co/8ft3uj3UWc	2021-05-24 09:46:20	0	0	Twitter Web App	1084687742032924674	josef_morasch		
1	1396763990424948736	Yes Babe https://t.co/HtlTsLdbDo	2021-05-24 09:44:20	0	0	Twitter Web App	1084687742032924674	josef_morasch		
2	1396763678498869251	RT @kscpl88_91: Hope everyone had a great #Sun...	2021-05-24 09:43:06	11	0	Twitter Web App	1084687742032924674	josef_morasch		
3	1396763521283760129	Undress your panty then its perfect https://t....	2021-05-24 09:42:28	0	0	Twitter Web App	1084687742032924674	josef_morasch		
4	1396763261480144898	You are to fare https://t.co/dYhRDiTUEa	2021-05-24 09:41:26	0	0	Twitter Web App	1084687742032924674	josef_morasch		

To save spam user in suser.csv

```
In [23]: # Saving data with item space separating
Data.to_csv("C:/Users/asha/Desktop/Final Year Project Phase 2/Phase 2/Dataset/Spam_user.csv", sep=',', encoding='utf8')
```

In []: