

TWITTER SENTIMENT ANALYSIS AND NEWS ARTICLE PREDICTION

**A Project Report submitted in partial fulfilment of the requirements for the
award of the degree of**

BACHELOR OF TECHNOLOGY IN

COMPUTER SCIENCE AND ENGINEERING

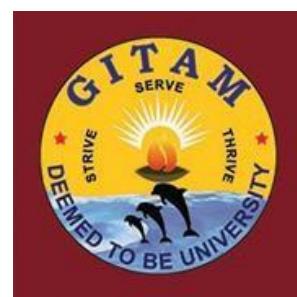
Submitted by

**KADIYALA P B ROHIT BHARADWAJ, 121910313006
SAMUDRALA MUNI VISHNU, 121910313011
SAI HARIDEEP, 121910313024
NAVEEN CHOWDARY, 121910313046**

Under the esteemed guidance of

Dr. N SURESH KUMAR

Assistant Professor

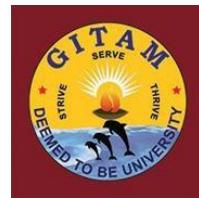


**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING
GITAM
(Deemed to be University)
VISAKHAPATNAM
OCTOBER 2022**

**DEPARTMENT OF COMPUTER SCIENCE AND
ENGINEERING GITAM SCHOOL OF TECHNOLOGY**

GITAM

(Deemed to be University)



DECLARATION

I/We, hereby declare that the project report entitled "**TWITTER SENTIMENT ANALYSIS AND NEWS ARTICLE PREDICTION**" is an original work done in the Department of Computer Science and Engineering, GITAM School of Technology, GITAM (Deemed to be University) submitted in partial fulfillment of the requirements for the award of the degree of B.Tech. in Computer Science and Engineering. The work has not been submitted to any other college or University for the award of any degree or diploma.

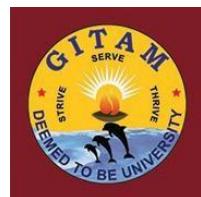
Date:

Registration No(s).	Name(s)	Signature(s)
121910313006	KADIYALA P B ROHIT BHARADWAJ	
121910313011	SAMUDRALA MUNI VISHNU	
121910313024	SAI HARIDEEP	
121910313046	NAVEEN CHOWDARY	

**DEPARTMENT OF COMPUTER SCIENCE AND
ENGINEERING GITAM SCHOOL OF TECHNOLOGY**

GITAM

(Deemed to be University)



CERTIFICATE

This is to certify that the project report entitled “TWITTER SENTIMENT ANALYSIS AND NEWS ARTICLE PREDICTION” is a bonafide record of work carried out by **KADIYALA P B ROHIT BHARADWAJ (121910313006)**, **SAMUDRALA MUNI VISHNU (121910313011)**, **SAI HARIDEEP (121910313024)**, **NAVEEN CHOWDARY (121910313046)** students submitted in partial fulfilment of requirement for the award of degree of Bachelors of Technology in Computer Science and Engineering.

Project Guide

Dr. N Suresh Kumar

Assistant Professor

Head of the Department

Dr. R.Sireesha

Professor

TABLE OF CONTENTS

Abstract	V
1. Introduction	1
2. Literature Review	2-3
3. Problem Identification & Objectives	3
4. System Methodology	4-5
5. Overview of Technologies	6-10
6. Implementation	11-18
7. Results and Discussions	19-23
8. Conclusion & Future Scope	24-25
9. References	26

ABSTRACT

For many people, social media is fast expanding as a virtual communication medium. A great deal of information is being sent from user to user over the internet. Twitter is the most popular platform for social networking followed by WhatsApp. The sentiment of a given tweet is crucial in understanding the emotion of the tweet. We may readily comprehend the sentiments in tiny tweets. However, finding the sentiment in tweets with intricate phrases is more difficult.

This software collects all tweets with a certain keyword and produces the sentiment of tweets with that term. This allows you to determine whether a keyword has a positive, negative, or neutral response throughout the platform.

Everyone refers to news in their spare time. Everyone's life revolves around the news. When reading the news, one should be aware of the issue they are reading about, such as sports, crime, education, and so on. When a text is provided, this project identifies the topic of the text. It also anticipates the parts of speech in the text, providing us a fundamental understanding of how the English language works.

1. INTRODUCTION

In today's quickly changing world, many people communicate extensively using social media platforms, one of which is "Twitter." Many people connect and engage on Twitter by posting various information in the form of tweets. These might be Positive, Negative, or even Neutral. Based on a single term and a small number of tweets, this project attempts to analyze those tweets and forecast whether they will be favorable, negative, or neutral. The Twitter API is a crucial aspect of our sentiment analysis assignment since it allows us to retrieve recent tweets based on keywords and tweet count directly by connecting using API credentials. The tweets will be analyzed using several NLTK algorithms and approaches, such as VANDER, which generates sentiment intensity scores for each tweet. Finally, this project includes graph graphics to help you grasp the insights. Matplotlib was used to create this graphic.

There is a lot of news in the form of articles floating around the internet. The news can be divided into several categories that are difficult to distinguish. As a result, in this project, we also used some machine learning techniques to train and predict the news article's category based on previous training data, and a frontend is created using Streamlit for the articles NLP purpose, which helps to identify different parts of speech using tokenization, lemmatization, and word cloud.

2. LITERATURE REVIEW

[1] In Natural Language Processing Research, Sentiment analysis is a growing area of research with studies that mainly includes document level classification. In information gathering habit the key aspect has always been is finding out what other people think. The increased availability of personal blogs and online review sites which are major opinion-rich resources that have resulted in new possibilities and difficulties as majority of individuals `utilise such information technology to understand and comprehend other' perspectives. The sudden increased activity in the field of sentiment analysis and opinion mining that deals with the sentiment, computational opinion treatment, and test subjectivity, has thus occurred due to response to the increased interest in new systems that are dealing with opinions as a first-class object. (**Pang and Lee 2008**)

[2] We find and validate limitations from conjunctions on the positive or negative semantic orientation of conjoined adjectives using a large corpus. When each conjunction is analysed individually, a log-linear regression model achieves 82% accuracy in predicting whether conjoined adjectives are of the same or different orientations. A clustering method splits the adjectives into groups of distinct orientations by combining the restrictions across multiple adjectives, and then the adjectives are labelled positive or negative. Evaluations on real data and simulated trials show that the classification precision is greater than 90% for adjectives found in a small number of conjunctions in the corpus. (**Hatzivassiloglou and McKeown 1997;**).

[3] SentiWordNet is a WordNet-derived opinion lexicon in which each phrase is connected with numerical scores expressing negative and positive sentiment information. Opinion mining (OM) is a new sub discipline at the intersection of computational linguistics and information retrieval that is mostly deals with the opinion expressed in a text rather than the topic of the document. OM has a wide range of applications, from tracking user views ranging from political candidates to products stated in internet forums and customer relationship management. Recent research has attempted to automatically assess the PNpolarity of subjective phrases, i.e. if a term that

is of an opinionated material has a negative or positive meaning, in order to assist in the gathering of views from the text. SENTIWORDNET is free to use for research purposes and has a Web-based graphical user interface. **(Esuli and Sebastiani 2006)**

[4] The richest languages linguistically such as Azerbaijani have received little attention in the field of language natural processing. The text corpora developed from Azerbaijani news stories is intended to be used for automated news labelling with supervised machine learning algorithms. Chi-squared test and LASSO techniques were used for pre-processing and feature selection. Supervised machine learning approaches to the text corpus were applied to compare the performance results of well-established supervised machine learning algorithms in the area of Azerbaijani language (U Suleymanov and S Rustamov)

3. PROBLEM IDENTIFICATION & OBJECTIVES

Twitter is a massive social media network that is now available. Many individuals use Twitter to discuss or communicate about the concerns. It is critical to comprehend what the author of the tweet is stating, whether it is favourable, negative, or neutral. Small keywords are easy to grasp, while complicated ones make it more difficult to understand the emotion of the tweet. This project assists in quickly determining the sentiment of all tweets through the use of visuals for improved comprehension.

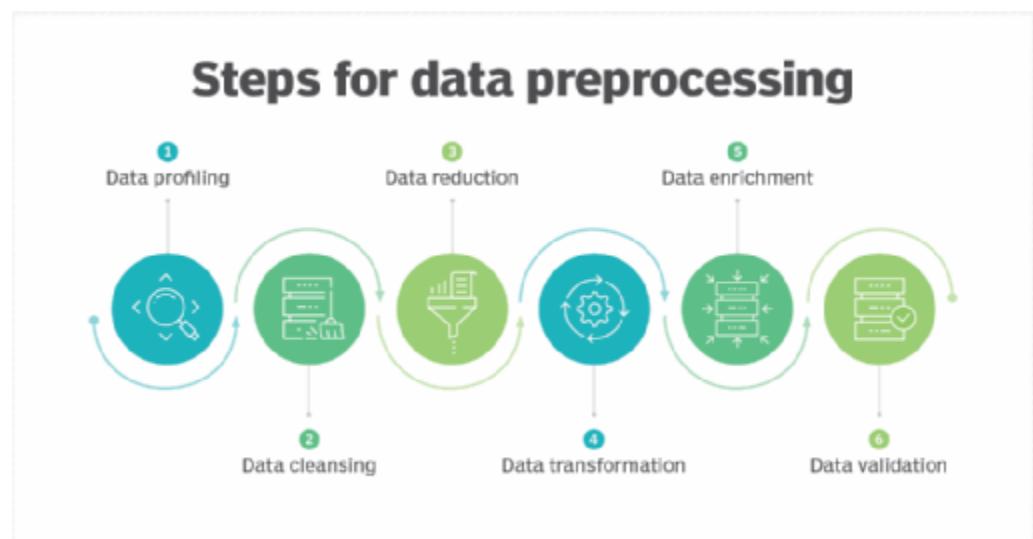
There is a lot of news flying around the internet in the form of articles. The news may be classified into various difficult-to-distinguish categories. One should be aware of the issue news to which they are referring. This is a really challenging undertaking. This project solves the challenges and forecasts the key categories of news such as sports, education, and so on, which may be extremely beneficial in quickly distinguishing the articles.

4. SYSTEM METHODOLOGY

4.1 DATA PREPROCESSING

The main objective of data pre-processing is preparing raw data and fitting it to a machine learning model. It is the first and most important stage in developing a machine learning model.

It is very difficult always to come across clean and prepared data when we work on a machine learning project. And before you can do anything with data, it must be cleaned and formatted. So we employ the data pre-processing job for this.



- Accomplished data pre-processing task in the second phase of project which is the prediction and categorization of news articles
- The data set was obtained from Kaggle, an open source platform for large datasets.

4.2 DESIGN:

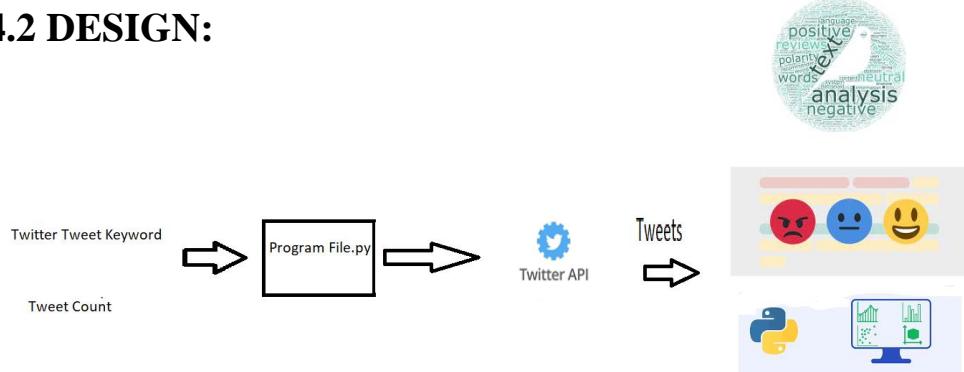


Figure a: Block Diagram, which describes working of Twitter tweet sentiment analysis

The above diagram can be divided into two parts. One is gathering the required tweets based on tweet keyword and count from the Twitter API. Next part is to analyse the gathered tweets by using NLTK sentiment analysis tools and provide a visualization output of overall tweets.

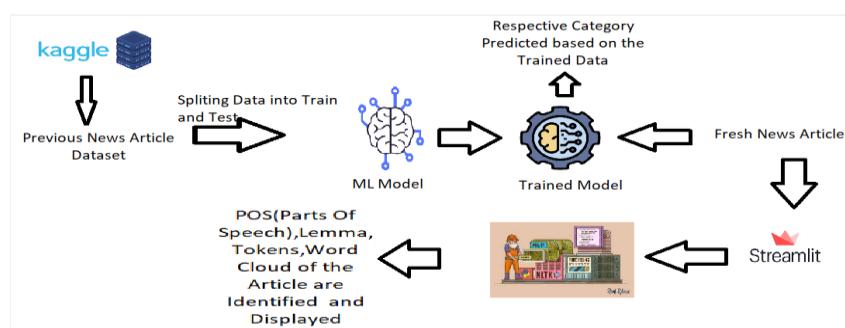


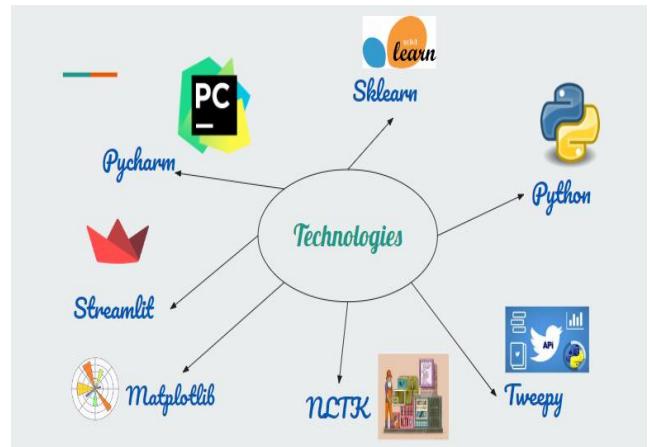
Figure b: Block Diagram, which describes working of News Article Prediction with ML and Its NLP outcome with a Frontend (Streamlit)

The method is separated into five steps, as shown in the above diagram. First, data from the News Article Category is acquired from the Kaggle Platform and converted into structured data using data pre-processing. The data is then divided into test and train data and trained with a Machine Learning Model. Once trained, the model will be able to recognise new news item categories based on the previously learned data set. We will receive a correctly predicted news category whenever we supply new input to the software. A front application is also constructed with Streamlit, which uses NLTK to detect the POS, Lemma, and Tokens, as well as show the article as a Word cloud.

5. OVERVIEW OF THE TECHNOLOGIES

5.1 TECHNOLOGIES USED:

- Tweepy
- Twitter Developer API
- NLTK
- Streamlit
- Kaggle
- Scikit-learn
- Pandas
- Matplotlib
- Text Blob



5.1.1 NATURAL LANGUAGE PROCESSING:

It is a branch of linguistics, computer science, information technology, and artificial intelligence concerned with computer interaction with human (natural) language, specifically processing and analysing massive volumes of natural language data. It deals with computer programming. Natural language processing issues include speech recognition, natural language interpretation, and natural language generation. This aided us in determining the nature of the tweets based on the user supplied keyword and tweet count. It also helped us identify the POS (Parts of Speech), tokens, and lemma of words in news articles.

5.1.2 Tweepy

Tweepy Python programme available as an open source that provides a very straightforward method to use Python to access the Twitter API. Tweepy contains a collection of classes and methods that reflect Twitter's models and API endpoints, and it handles numerous implementation details transparently, such as data decoding and encoding. It is a simple Python module for interacting with the Twitter API. The API class provides access to all RESTful API services provided by Twitter. Each method can accept and return various parameters.



5.1.3 Twitter Developer API

Developer Platform

The Twitter developer site includes a suite of self-service tools for managing access to the Twitter API and Twitter Ads API.

You may do the following in the portal:

- Make and manage Twitter Projects and Apps (and the authentication keys and tokens that they provide).
- With the Twitter API premium v1.1 and v2 endpoints, you can manage your access levels and integrations.
- Learn more about the many endpoints and features that are available.
- If you have elevated access and an organisation account type, you can visit team pages where you may add and manage the many handles that have access to your team's account.

The Twitter API provides unique and sophisticated programmatic access to Twitter. We may also make use of major Twitter features such as Tweets, Direct Messages, Spaces, Lists, and People.

5.1.4 NLTK

The Natural Language Toolkit (NLTK) are set of applications and tools of Python for statistical and symbolic natural language processing in English. Various word processing products include a plethora of test datasets. You can use NLTK to perform a number of tasks, including: B. Visualization and tokenization of parse trees. The Natural Language Toolkit (NLTK) is a Python framework for developing algorithms that interface with human language data in statistical natural language processing applications (NLP). Text processing technologies include tokenization, parsing, classification,

and stemming, tagging, and semantic inference. The word tokenizes () function from NLTK is used to break a sentence into tokens or words. Tokens may be extracted from a string of characters using the tokenize.word tokenize()



5.1.5 Streamlit

reamlit is Python built open source framework app. It allows us to quickly construct data science and machine learning web apps. It works with popular Python libraries including scikit-learn, Keras, PyTorch, SymPy (latex), NumPy, pandas, and Matplotlib.

Streamlit, in a nutshell, convert data scripts into shareable web apps in minutes. Streamlit widgets enable you to incorporate interaction right into your apps, including buttons, sliders, text inputs, and more.

You may use Streamlit to work with any type of data or input. You can install Streamlit directly in our Python IDE by using the command pip install streamlit.

5.1.6 Kaggle

Kaggle is a data scientist and machine learning enthusiast community platform. Users can collaborate with others, publish and discover datasets, use GPU-integrated notebooks, compete and solve data science challenges on Kaggle with other data scientists. It is an open source platform for accessing different types of datasets.



5.1.7 Scikit-learn

Scikit-learn are perhaps the most helpful Python machine learning package. The sklearn package contains several powerful statistical modelling and machine learning techniques such as classification, clustering, regression, and dimensionality reduction. It provides a variety of effective statistical modelling and machine learning tools such as classification, clustering,

regression and dimensionality reduction through Python's consistency interface. Scikit-learn have a tonne of functionality available. Consider the following to have a better understanding of the spread:

- Supervised learning methods
- Cross-validation
- Unsupervised learning algorithms
- Various toy datasets

5.1.8 Pandas

Pandas are a data analysis and manipulation software package created for the Python programming language. It includes specific data structures and procedures for working with time series and mathematical tables. Wes McKinney came up with the moniker "Pandas" in 2008, and it refers to both "Panel Data" and "Python Data Analysis."

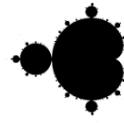
- With the aid of Pandas, we can examine large data sets and draw conclusions based on statistical principles.
- Pandas can organise disorganised data sets, making them understandable and useful.
- In data science and ML, relevant data is crucial.



5.1.9 Matplotlib

Matplotlib is a graphing library for Python programming language and NumPy, a numerical mathematics extension. For integrating charts into programmes utilising all-purpose GUI toolkits like Tkinter, wxPython, Qt, or GTK, it offers an object-oriented API. Matplotlib makes difficult things possible and simple things easy. Produce plots fit for publishing. Make interactive charts with zoom, pan, and update capabilities. It was created by John D. Hunter.

The majority of Matplotlib's tools are found in the pyplot submodule and are often imported using the plt alias. Matplotlib contains all types of plots which are used a lot in ML and Data science field. It is best suitable Library for Data Visualization



TextBlob

5.1.10 TextBlob

Being a Lexicon-based sentiment analyser, TextBlob It contains certain predetermined guidelines, or perhaps we should say a word and weight dictionary, which has some scores that assist in determining the polarity of a statement. It provides a simple API for getting started with common NLP operations such as tagging part-of-speech, extraction of noun phrase, sentiment analysis, translation, classification, and others.

For this reason, rule-based sentiment analysers are another name for lexicon-based sentiment analysers. Text blob acts as a Competitor for VANDER based Sentiment analysis which provides much more efficiency than the former one.

The vocabulary of a person, a language, or a field of study is called a lexicon.

6. IMPLEMENTATION

6.1 Coding

6.1.1 For Twitter Sentiment Analysis

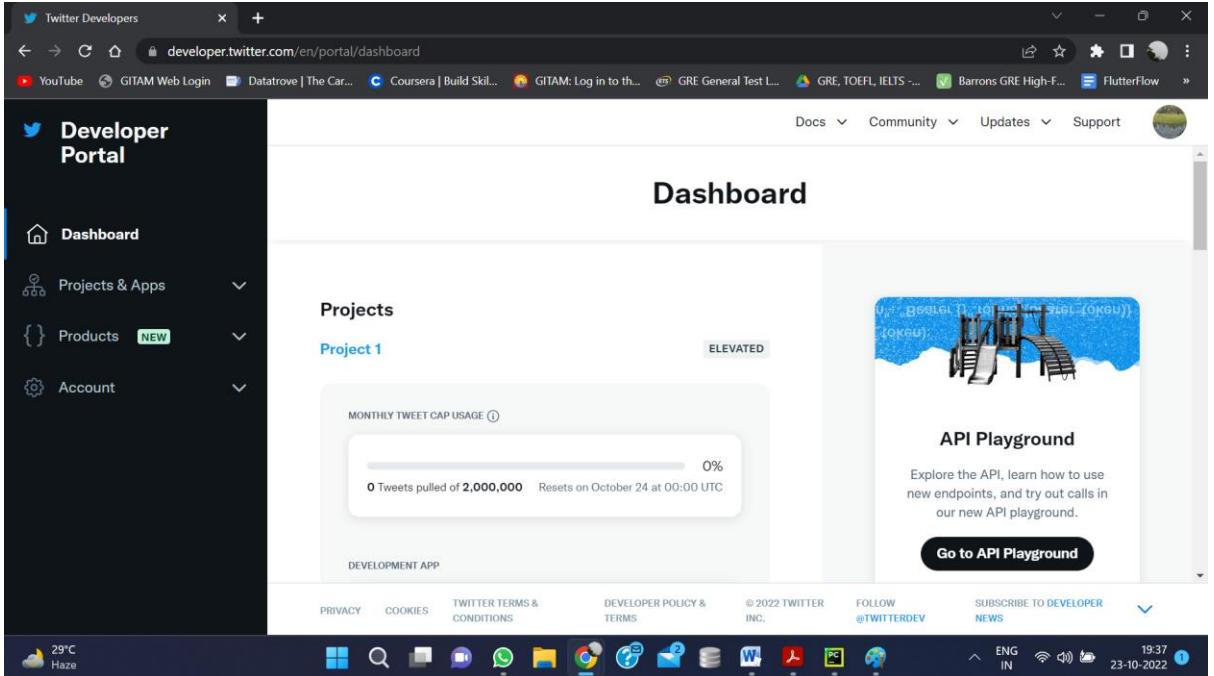


Fig 6.1.1-a: Twitter Developer API

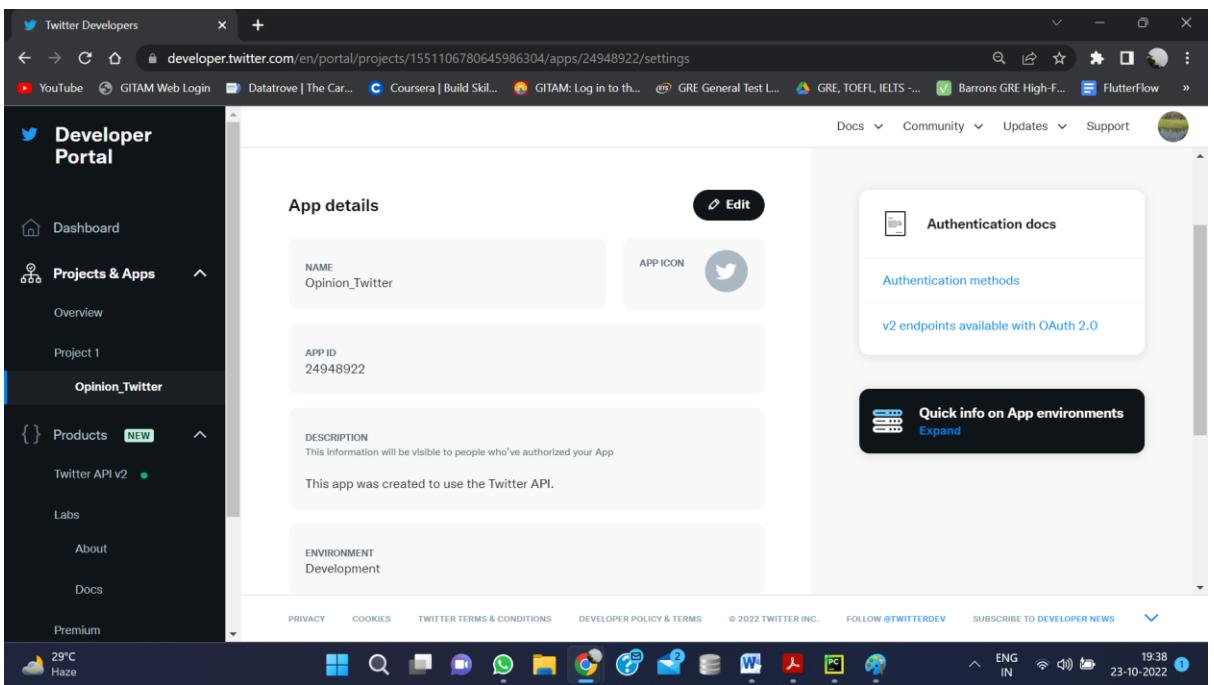


Fig 6.1.1-b: Twitter Developer API Project Creation Details

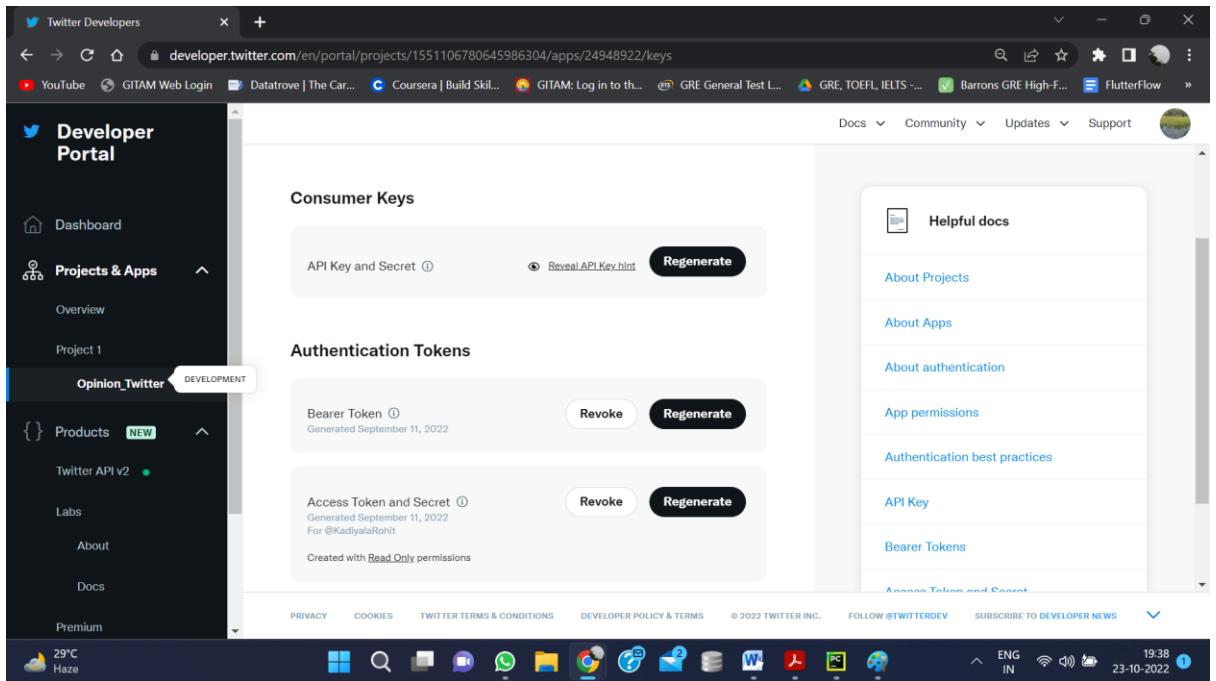


Fig 6.1.1-c Twitter Developer API Keys and Access Tokens

Twitter Sentiment Analysis Python Code in Pycharm IDE

```

from textblob import TextBlob
import tweepy
import matplotlib.pyplot as plt
import pandas as pd
from termcolor import colored, cprint
import nltk
nltk.download('vader_lexicon')
from nltk.sentiment.vader import SentimentIntensityAnalyzer
consumerKey = "v0y5ruvn8XP5N0cCx4fQVAL0M"
consumerSecret = "VjPp9XzNcCP4o5YLXXWSkJ0fwWkXMGpK9PZTaQNxEK0nOHYeB"
accessToken = "1467127046543319040-w6xFI7HtErblT4YU8L9qvv4fLwzZn"
accessTokenSecret = "UGKhVtRyD6kLfWLujEznJxg8Y4TvcU7doN8dnal_ddb"
auth = tweepy.OAuthHandler(consumerKey, consumerSecret)
auth.set_access_token(accessToken, accessTokenSecret)
api = tweepy.API(auth)

def percentage(part, whole):
    return 100 * float(part)/float(whole)
keyword = input("Please enter keyword or hashtag to search:")
nooftweets = int(input("Enter Number of Tweets to analyse:"))
cprint("Analyzing the Tweets", 'red', attrs=['blink'])
tweets = tweepy.Cursor(api.search_tweets, q=keyword).items(nooftweets)
for tweet in tweets:
    blob = TextBlob(tweet.text)
    print(blob.sentiment.polarity)
    if blob.sentiment.polarity > 0:
        cprint(tweet.user.screen_name, 'green')
    elif blob.sentiment.polarity < 0:
        cprint(tweet.user.screen_name, 'red')
    else:
        cprint(tweet.user.screen_name, 'white', attrs=['bold'])
    print(tweet.user.location)
    print(tweet.created_at)
    print(tweet.favorite_count)
    print(tweet.retweet_count)
    print(tweet.text)
    print("\n")

```

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File2.py
Test1 > Project_SEM7 > Main > File2.py
Project Bookmarks Registry Explorer Structure Version Control Python Packages TODO Python Console Problems Terminal Services Graph Database Console
25     neg = 0
26     neu = 0
27     polarity = 0
28     tweet_list = []
29     neu_list = []
30     neg_list = []
31     pos_list = []
32     for tweet in tweets:
33
34         #print(tweet.text)
35         tweet_list.append(tweet.text)
36         analysis = TextBlob(tweet.text)
37         score = SentimentIntensityAnalyzer().polarity_scores(tweet.text)
38         neg = score['neg']
39         neu = score['neu']
40         pos = score['pos']
41         comp = score['compound']
42         polarity += analysis.sentiment.polarity
43         if neg > pos:
44             neg_list.append(tweet.text)
45             neg += 1
46         elif pos > neg:
47             pos_list.append(tweet.text)
48             pos += 1
for tweet in tweets
39:23 CRLF UTF-8 4 spaces Python 3.10
29°C Haze 19:36 23-10-2022
```

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File2.py
Test1 > Project_SEM7 > Main > File2.py
Project Bookmarks Registry Explorer Structure Version Control Python Packages TODO Python Console Problems Terminal Services Graph Database Console
49     elif pos == neg:
50         neu_list.append(tweet.text)
51         neu += 1
52
53
54     #print(pos)
55     #print(neg)
56     #print(neu)
57     pos = percentage(pos, nooftweets)
58     neg = percentage(neg, nooftweets)
59     neu = percentage(neu, nooftweets)
60     polarity = percentage(polarity, nooftweets)
61     pos = format(pos, '.1f')
62     neg = format(neg, '.1f')
63     neu = format(neu, '.1f')
64     #print(pos)
65     #print(neg)
66     #print(neu)
67
68     tweet_list = pd.DataFrame(tweet_list)
69     neutral_list = pd.DataFrame(neu_list)
70     negative_list = pd.DataFrame(neg_list)
71     positive_list = pd.DataFrame(pos_list)
72     print("Total number: " len(tweet_list))
for tweet in tweets
39:23 CRLF UTF-8 4 spaces Python 3.10
29°C Haze 19:36 23-10-2022
```

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File2.py
Test1 > Project_SEM7 > Main > File2.py
File1.py x File2.py x File4_long.apy x
71 positive_list = pd.DataFrame(pos_list)
72 print("total number: ", len(tweet_list))
73 print("positive number: ", len(positive_list))
74 print("negative number: ", len(negative_list))
75 print("neutral number: ", len(neutral_list))
76 plt.subplot(2,1,1)
77 labels = ['Positive ['+str(pos)+'%]', 'Neutral ['+str(neu)+'%]', 'Negative ['+str(neu)+'%]']
78 sizes = [pos, neu, neg]
79 colors = ['green', 'Yellow', 'red']
80 patches, texts = plt.pie(sizes, colors=colors, startangle=90, shadow=True)
81 plt.style.use('default')
82 plt.legend(labels)
83 plt.title("Sentiment Analysis Result for keyword= "+keyword+" ")
84 plt.axis('equal')
85 plt.subplot(2,1,2)
86 li=len(positive_list)
87 l2=len(negative_list)
88 l3=len(neutral_list)
89 k1=[l1,l2,l3]
90 k2=["POSITIVE", "NEGATIVE", "NEUTRAL"]
91 c=[["Green", "Red", "Yellow"]]
92 z=plt.bar(k2,k1,color=c, width=.4)
93 plt.title("Tweet Count Based on Emotion")
94 labels = [ '😊', '😢', '😡' ]
for tweet in tweets:
    # Process tweet...

```

This screenshot shows a Jupyter Notebook interface with three open files: File1.py, File2.py, and File4_long.apy. The current cell in File2.py contains Python code for sentiment analysis. It prints the total, positive, negative, and neutral tweet counts. It then creates a pie chart showing the distribution of sentiments (Positive, Neutral, Negative) with labels and colors corresponding to the counts. Finally, it creates a bar chart showing the tweet count for each emotion category (Positive, Negative, Neutral).

```

File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File2.py
Test1 > Project_SEM7 > Main > File2.py
File1.py x File2.py x File4_long.apy x
80 patches, texts = plt.pie(sizes, colors=colors, startangle=90, shadow=True)
81 plt.style.use('default')
82 plt.legend(labels)
83 plt.title("Sentiment Analysis Result for keyword= "+keyword+" ")
84 plt.axis('equal')
85 plt.subplot(2,1,2)
86 li=len(positive_list)
87 l2=len(negative_list)
88 l3=len(neutral_list)
89 k1=[l1,l2,l3]
90 k2=["POSITIVE", "NEGATIVE", "NEUTRAL"]
91 c=[["Green", "Red", "Yellow"]]
92 z=plt.bar(k2,k1,color=c, width=.4)
93 plt.title("Tweet Count Based on Emotion")
94 labels = [ '😊', '😢', '😡' ]
95 for rect1, label in zip(z, labels):
    height = rect1.get_height()
    plt.annotate(label, (rect1.get_x() + rect1.get_width()/2, height+.05), ha="center", va="bottom", fontsize=30)
96 plt.show()
97
98
for tweet in tweets:
    # Process tweet...

```

This screenshot is identical to the one above, showing the same Python code for sentiment analysis and visualization. The difference is in the annotations within the bar chart loop. Instead of just printing the labels above the bars, it uses the `plt.annotate` method to place the labels directly on top of each bar at a slightly higher position.

6.1.2 For News Article Prediction

category	title	body
ARTS & CULT Modeling A	In October 2017, Carolyn Kramer received a disturbing phone call. The former modeling agent listened intently as a model she used to represent told her that a famous French	
ARTS & CUL Actor Jeff H	This week I talked with actor Jeff Hiller about the hit Off Broadway play Bright Colors And Bold Patterns that he'll be joining on January 17th with a new opening night sched	
ARTS & CUL New Yorker	The New Yorker is taking on President Donald Trump after he asked why the U.S. would welcome immigrants from "shithole" places like Haiti and African countries during a	
ARTS & CUL Man Surprised Kellen Hickey, a 26-year-old who lives in Hudson, Wisconsin, has gift giving down to a fine art. He drew himself and his girlfriend Lindsay Brinkman, 24, in 10 different animal		
ARTS & CUL This Artist	There's something about combining the traditional, uptight look of the Renaissance period with modern-day behavior that excites Barcelona-based artist Gerard Mas. His re	
ARTS & CUL This Dutch	Josje Duk has a sweater that reads "DON'T PANIC." She wears it on days when she might, well, panic. "I always tell my family, if this all doesn't work out, I'll go study math	
ARTS & CUL Broadway	Multiple women have accused Broadway star Ben Vereen of sexual misconduct ranging from harassment to assault, according to a Friday morning report from the New York Daily News.	
ARTS & CUL Sculptures	The world's largest ice festival began this week in Harbin, a city in the northeastern part of China. The Harbin International Ice and Snow Festival goes through late February.	
ARTS & CUL The Met	Non-New Yorkers officially have less than two months to take advantage of the Metropolitan Museum of Art's pay-as-you-wish admission policy. Tourists will be charged a	
ARTS & CUL Duncan Jones	Duncan Jon David Bowie's cultural legacy continues. The legendary musician's son, movie director Duncan Jones, has launched an online book club in honor of his late father. Bowie died in 2016 at the age of 69.	
ARTS & CUL Mystery	No Sue Grafton, who authored "The Kinsey Millhone Alphabet" mystery series, died Thursday in Santa Barbara, California, after a battle with cancer, her daughter announced on Facebook.	
ARTS & CUL Dick Van Dyke	'Dick Van Dyke Show,' died Thursday at her home in Los Angeles. She was 91.	
ARTS & CUL The Best	You may recall the literary drama that unfolded about this time last year as Simon & Schuster granted, and later revoked, a book deal for a memoir by former Breitbart editor Steve Bannon.	
ARTS & CUL Women-On-Days	After the 2016 presidential election, artist Roxanne Jackson impulsively posted a message on Facebook. "Hello female artists/curators! Let's organize a NASTY WOMEN'S DAY!"	
ARTS & CUL Books	As 2018 approaches, there's a lot to look forward to: the end of a hellish 2017, the Winter Olympics, "The Bachelor: Winter Games," 2017 being over, midterm elections, and the return of Game of Thrones.	
ARTS & CUL A Very Vint	Turn back the clock and experience the magic of yesteryear with this collection of black and white photographs sure to awaken the holiday spirit. Send David Lohr an email at dlohr@msn.com.	
ARTS & CUL Why Do We	Gift exchanges are a big part of American Christmas culture, often with a variety of creative spins on the tradition. You might be familiar with a game in which everyone brings a gift to a friend's house.	
ARTS & CUL Even Taylor	"My idea from the beginning was I wanted it to be like a moving Vanity Fair cover," said Shoshana Bean, describing the video she and fellow Broadway superstar Cynthia Erivo made for their show's promotional campaign.	
ARTS & CUL Cards Again	The brains behind the game Cards Against Humanity have decided to "tackle the biggest issue in the world: wealth inequality" by sending checks to 100 of their poorest cust	
ARTS & CUL New Allegal	Lindsay Jones never planned to speak publicly about her experiences with prolific fashion photographer Terry Richardson. The New York City-based designer and model confirmed the allegations in an interview with EW.com.	
ARTS & CUL Merrism	Merrism, Webster's word of the year is "Feminism." The term enjoyed multiple lookups spikes on the dictionary's website in 2017, and an overall 70 percent rise in its use.	

Fig-a: Dataset CSV File

```

File Edit View Navigate Code Behavior Run Tools VCS Window Help T1.py - File4_long.a.py
Project: Project_SEM7 Main: File4_long.a.py
File1.py x File2.py x File4_long.a.py
1 import pandas as pd
2 import numpy as np
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.model_selection import train_test_split
5 from sklearn.svm import SVC
6 import matplotlib.pyplot as plt
7 from sklearn.naive_bayes import MultinomialNB
8 df = pd.read_csv(r":\Users\user\PycharmProjects\Project_SEM7\news-article-categories.csv")
9 k = pd.set_option('display.max_columns', None)
10 l = pd.set_option('display.max_rows', None)
11
12 # print(df.head())
13 # print(df.shape)
14 #print(df.info())
15 #print(df.isnull().sum())
16 #print(df["category"].value_counts())
17 df = df[df.category != "ARTS & CULTURE"]
18 df1 = df[["title", "category"]]
19 #print(df1)
20 x = np.array(df1["title"])
21 y = np.array(df1["category"])
22
23 #BarPlot for the Total Category
24 ax=df1["category"].value_counts().sort_values().plot(kind='barh')

```

The screenshot shows the PyCharm IDE interface with a dark theme. The main window displays a Python script named `File4_long.py`. The code uses Matplotlib to create a bar plot of news article counts and trains a MultinomialNB model. It includes user input for text classification. The bottom status bar shows system information like temperature (28°C), date (23-10-2022), and time (14:18).

```
22
23     #BarPlot for the Total Category
24     ax=df[“category”].value_counts().sort_values().plot(kind=‘barh’)
25     ax.bar_label(ax.containers[0])
26     plt.title(“Count of Different News Articles”)
27     plt.show()
28
29     cv = CountVectorizer()
30     X = cv.fit_transform(x)
31     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
32
33     model =MultinomialNB()
34     model.fit(X_train,y_train)
35
36
37
38     user = input(“Enter a Text: ”)
39     data = cv.transform([user]).toarray()
40     output = model.predict(data)
41
42     print(output[0])
43
```

6.1.3 For NLP with Streamlit

The screenshot shows the PyCharm IDE interface with a dark theme. The main window displays a Python script named `File1.py`. The code imports Streamlit, joblib, os, spacy, and pandas. It loads a pre-trained spaCy model, sets up matplotlib, and imports Random Forest Classifier from sklearn. It also imports WordCloud and ImageColorGenerator from wordcloud. A function `load_prediction_models` loads a saved model. The bottom status bar shows system information like temperature (28°C), date (23-10-2022), and time (14:51).

```
1     import streamlit as st
2     import joblib,os
3     import spacy
4     import pandas as pd
5
6     nlp = spacy.load("en_core_web_sm")
7     import matplotlib.pyplot as plt
8     import matplotlib
9     matplotlib.use("Agg")
10    from sklearn.ensemble import RandomForestClassifier
11    from PIL import Image
12    from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
13
14    news_vectorizer = open("models/final_news_cv_vectorizer.pkl","rb")
15    news_cv = joblib.load(news_vectorizer)
16
17
18    def load_prediction_models(model_file):
19        loaded_model = joblib.load(open(os.path.join(model_file),"rb"))
20        return loaded_model
21
22
23
24    # GET THE KEYS
main()
```

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File1.py
Test1 > Project_SEM7 > Main > File1.py
Project File1.py < File2.py < File4_long.apy <
25     def get_key(val,my_dict):
26         for key,value in my_dict.items():
27             if val == value:
28                 return key
29
30
31     def main():
32         """News Classifier"""
33         st.title("News Classifier")
34         # st.subheader("ML App with Streamlit")
35         html_temp = """
36             <div style="background-color:#03a9fc;padding:10px">
37                 <h1 style="color:red;text-align:center;">Streamlit ML App </h1>
38             </div>
39         """
40         st.markdown(html_temp,unsafe_allow_html=True)
41
42         activity = ['NLP']
43         choice = st.sidebar.selectbox("Select Activity",activity)
44
45
46
47         if choice == 'NLP':
48             st.info("Natural Language Processing of Text")
main()
```

Version Control Python Packages TODO Python Console Problems Terminal Services Graph Database Console

45:1 CRLF UTF-8 Tab* Python 3.10 21:52 ENG IN 23-10-2022

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File1.py
Test1 > Project_SEM7 > Main > File1.py
Project File1.py < File2.py < File4_long.apy <
49     raw_text = st.text_area("Enter News Here","Type Here")
50     nlp_task = ["Tokenization", "Lemmatization", "NER", "POS Tags"]
51     task_choice = st.selectbox("Choose NLP Task", nlp_task)
52     if st.button("Analyze"):
53         st.info("Original Text:\n{}".format(raw_text))
54
55         docx = nlp(raw_text)
56         if task_choice == 'Tokenization':
57             result = [token.text for token in docx_]
58         elif task_choice == 'Lemmatization':
59             result = [f'Token:{}, Lemma:{}'.format(token.text, token.lemma_) for token in docx]
60         elif task_choice == 'NER':
61             result = [(entity.text, entity.label_) for entity in docx.ents]
62         elif task_choice == 'POS Tags':
63             result = [f'Token:{}, POS:{}, Dependency:{}'.format(word.text, word.tag_, word.dep_) for word in docx]
64
65         st.json(result)
66
67     if st.button("Tabulize"):
68         docx = nlp(raw_text)
69         c_tokens = [token.text for token in docx_]
70         c_lemma = [token.lemma_ for token in docx_]
71         c_pos = [token.pos_ for token in docx_]
main()
```

Version Control Python Packages TODO Python Console Problems Terminal Services Graph Database Console

45:1 CRLF UTF-8 Tab* Python 3.10 21:52 ENG IN 23-10-2022

The screenshot shows a code editor interface with a dark theme. The main window displays a Python script named `File1.py`. The code performs Natural Language Processing (NLP) tasks such as tokenization, lemmatization, and part-of-speech (POS) tagging using the spaCy library. It then generates a word cloud visualization if a checkbox is checked. The Streamlit library (`st`) is used to create a user interface with a sidebar containing an "About" section and a main area titled "Streamlit ML App". The sidebar also includes a dropdown for selecting an activity (set to "NLP") and a link to an "About" page. The Streamlit app itself has sections for entering news text, choosing NLP tasks (Tokenization, Analyze, Tabulize), and generating a WordCloud.

```
File1.py
68     docx = nlp(raw_text)
69     c_tokens = [token.text for token in docx_]
70     c_lemma = [token.lemma_ for token in docx_]
71     c_pos = [token.pos_ for token in docx_]
72
73     new_df = pd.DataFrame(zip(c_tokens,c_lemma,c_pos),columns=[ 'Tokens','Lemma','POS'])
74     st.dataframe(new_df)
75
76
77     if st.checkbox("WordCloud"):
78         c_text = raw_text
79         wordcloud = WordCloud().generate(c_text)
80         plt.imshow(wordcloud,interpolation='bilinear')
81         plt.axis('off')
82         st.pyplot()
83         st.sidebar.subheader("About")
84
85     main()
86
```

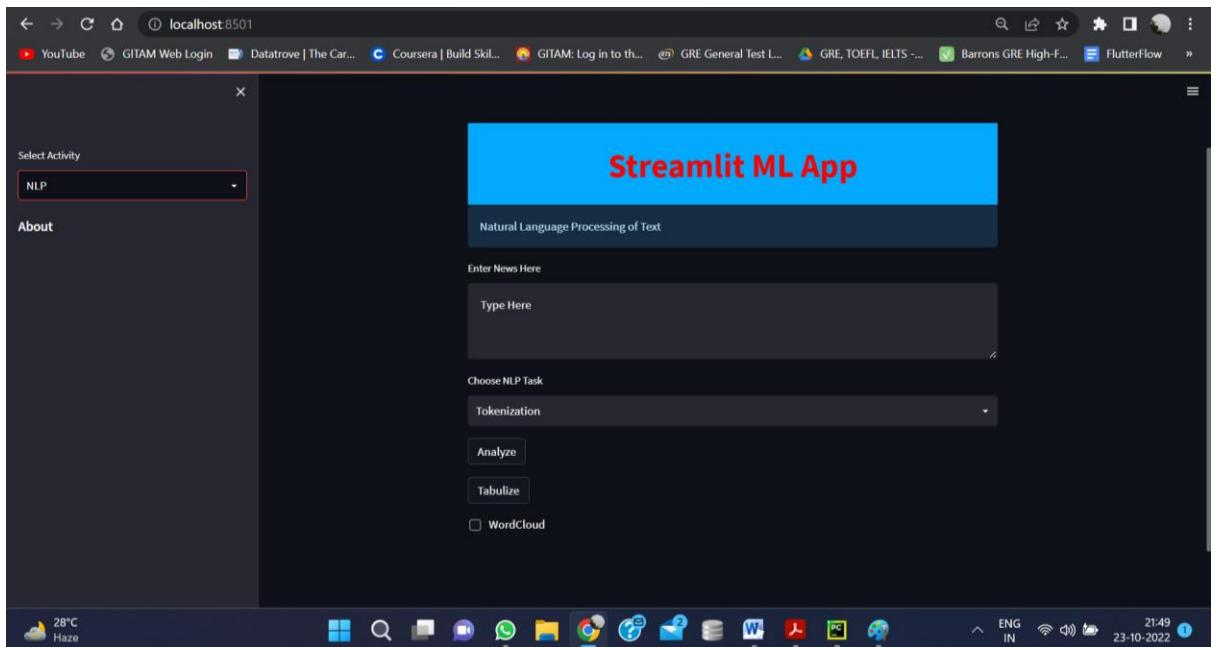


Fig-a: Streamlit Front End Interface

7. RESULTS & DISCUSSION

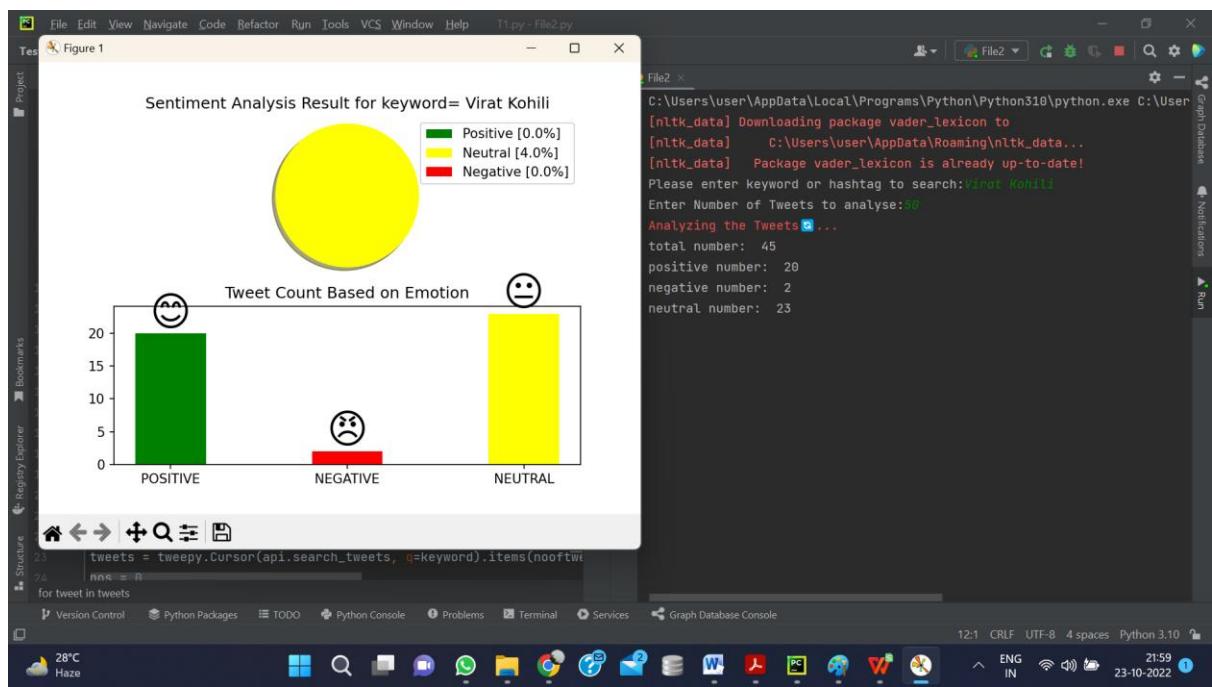
7.1 Twitter Sentiment Analysis Output

The screenshot shows the PyCharm IDE interface. The left pane displays the code for `File2.py`, which performs Twitter sentiment analysis. The right pane shows the terminal output where the user runs the script and enters a keyword ('Vinay Kohili'). The status bar at the bottom indicates the date (23-10-2022), time (21:59), and Python version (3.10).

```
from textblob import TextBlob
import tweepy
import matplotlib.pyplot as plt
import pandas as pd
from termcolor import colored, cprint
import nltk
nltk.download('vader_lexicon')
from nltk.sentiment.vader import SentimentIntensityAnalyzer
consumerKey = "v0y5ruvn0XP5N0cCx4fQUAL0M"
consumerSecret = "VjPp9XzSnGCP4o5YLXXWSkJ0fwWkXMGpK9PZTaQNxEK0n0H"
accessToken = "1467127046543319040-w6xFI7HtcrblT4YU8L9qv4fLwrZn"
accessTokenSecret = "UGKhVtRyD6KlFWLuJEznJxs68Y4TuCVU7doN8dnal.d0t"
auth = tweepy.OAuthHandler(consumerKey, consumerSecret)
auth.set_access_token(accessToken, accessTokenSecret)
api = tweepy.API(auth)

def percentage(part, whole):
    return 100 * float(part)/float(whole)
keyword = input("Please enter keyword or hashtag to search:")
nooftweets = int(input("Enter Number of Tweets to analyse:"))
cprint("Analyzing the Tweets", 'red', attrs=['blink'])
tweets = tweepy.Cursor(api.search_tweets, q=keyword).items(nooftweets)
nns = []
for tweet in tweets:
```

This screenshot is identical to the one above, showing the PyCharm IDE with the same code in `File2.py` and the same terminal output. The status bar at the bottom indicates the date (23-10-2022), time (21:59), and Python version (3.10).



7.2 News Article Prediction with ML Output

The screenshot shows the PyCharm IDE interface. The top navigation bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, Help, and a tab for T1.py - File4_long_a.py. The left sidebar features Project, Bookmarks, Registry Explorer, and Structure. The main editor window displays a Python script named File4_long_a.py. The code uses Matplotlib to plot a bar chart and a machine learning model (MultinomialNB) to predict sentiment based on user input. A terminal window on the right shows the output of running the script.

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File4_long_a.py
Test1 Project_SEM7 Main File4_long_a.py
File2.py x File4_long_a.py x Run: File4_long_a x
22
23 #BarPlot for the Total_Cat
24 axdf[["category"]].value_cou
25 ax.bar_label(ax.containers[
26 plt.title("Count of Differen
27 #plt.show()
28
29 cv = CountVectorizer()
30 X = cv.fit_transform(x)
31 X_train, X_test, y_train, y_
32
33 model =MultinomialNB()
34 model.fit(X_train,y_train)
35
36
37
38 user = input("Enter a Text:")
39 data = cv.transform([user])
40 output = model.predict(data)
41
42 print(output[0])
43
```

C:\Users\user\AppData\Local\Programs\Python\Python310\python.exe C:\Users\user\PycharmProjects\Test1\Project_SEM7\Main\File4_long_a.py

Enter a Text: India vs Pakistan Live Cricket score, T20 World Cup 2022: India beat Pakistan by four wickets, R

The screenshot shows the PyCharm IDE interface. On the left, there's a project tree for 'Test1' with files 'File2.py' and 'File4_long_a.py'. The main editor window displays the code for 'File4_long_a.py'. The code performs various NLP tasks including a bar plot, vectorization, and MultinomialNB classification. A terminal window at the bottom shows the command run and the resulting output: 'Enter a Text: India vs Pakistan Live Cricket Score, T20 World Cup 2022: India beat Pakistan by four wickets, Kohli scores 82' followed by 'SPORTS'.

```
File Edit View Navigate Code Refactor Run Tools VCS Window Help T1.py - File4_long_a.py
Project Test1 > Project_SEM7 > Main > File4_long_a.py
Run: File4_long_a.x C:\Users\user\AppData\Local\Programs\Python\Python310\python.exe C:\Users\user\PycharmProjects\Test1\Project_SEM7\Main\file4_long_a.py
Enter a Text: India vs Pakistan Live Cricket Score, T20 World Cup 2022: India beat Pakistan by four wickets, Kohli scores 82
SPORTS
Process finished with exit code 0

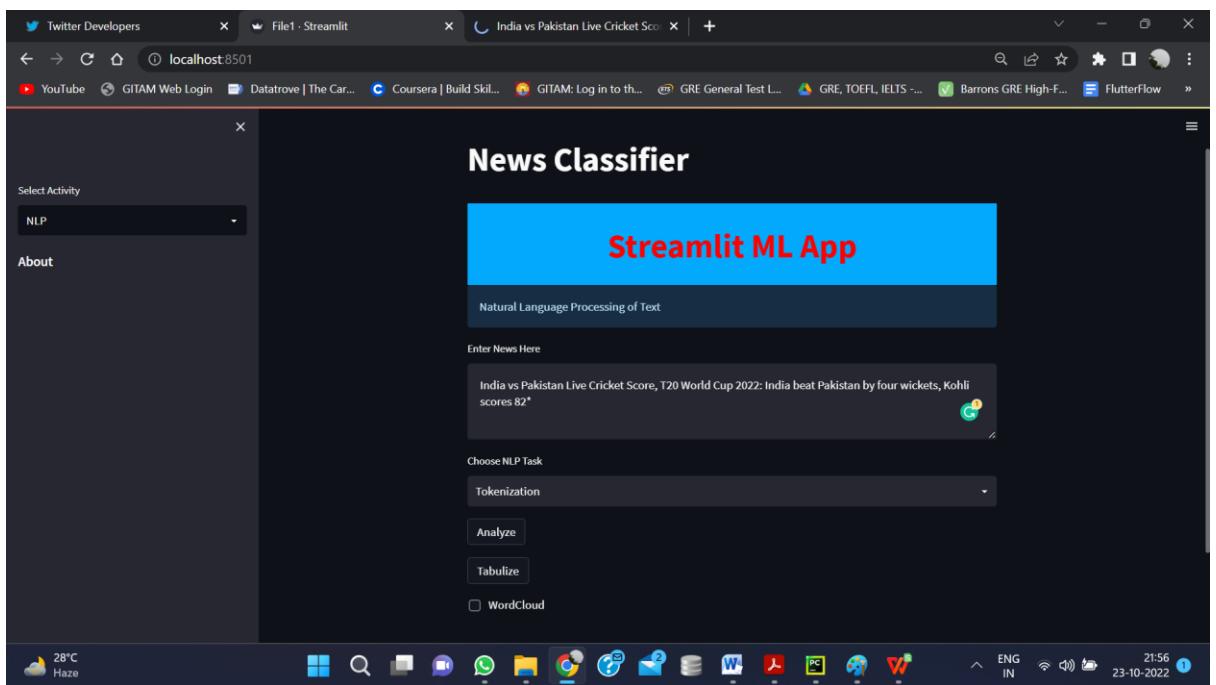
22
23 #BarPlot for the Total_Cate
24 ax=df[["category"]].value_cou
25 ax.bar_label(ax.containers[
26 plt.title('Count of Differen
27 #plt.show()

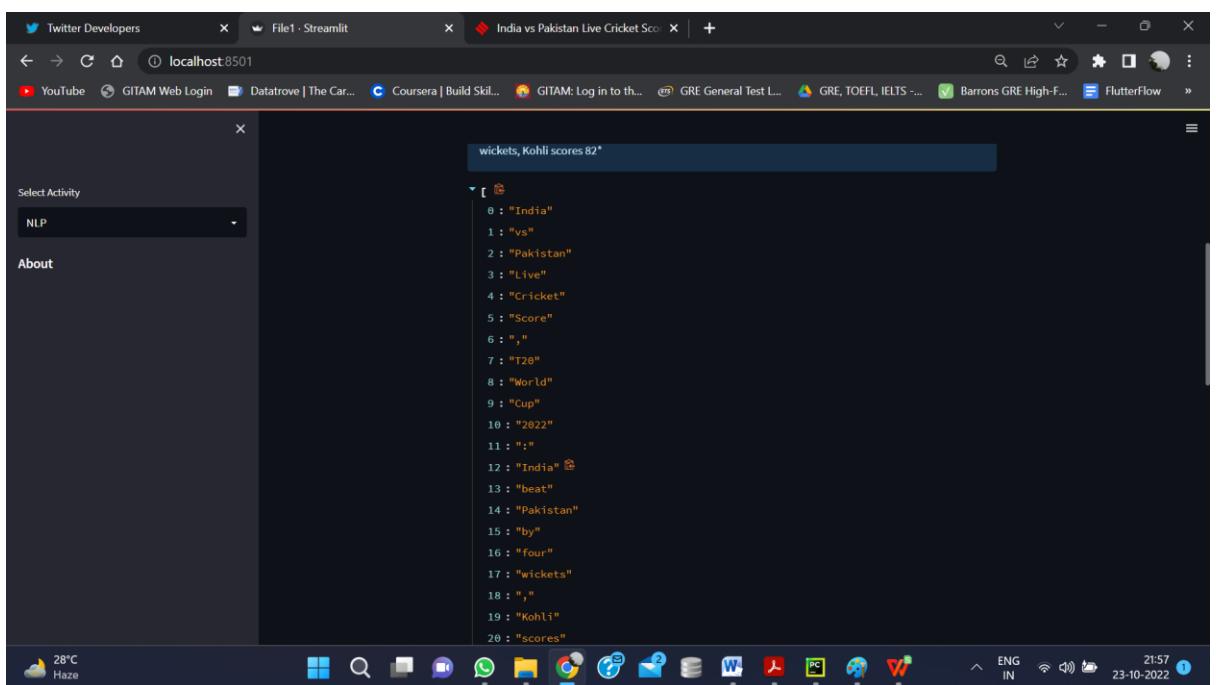
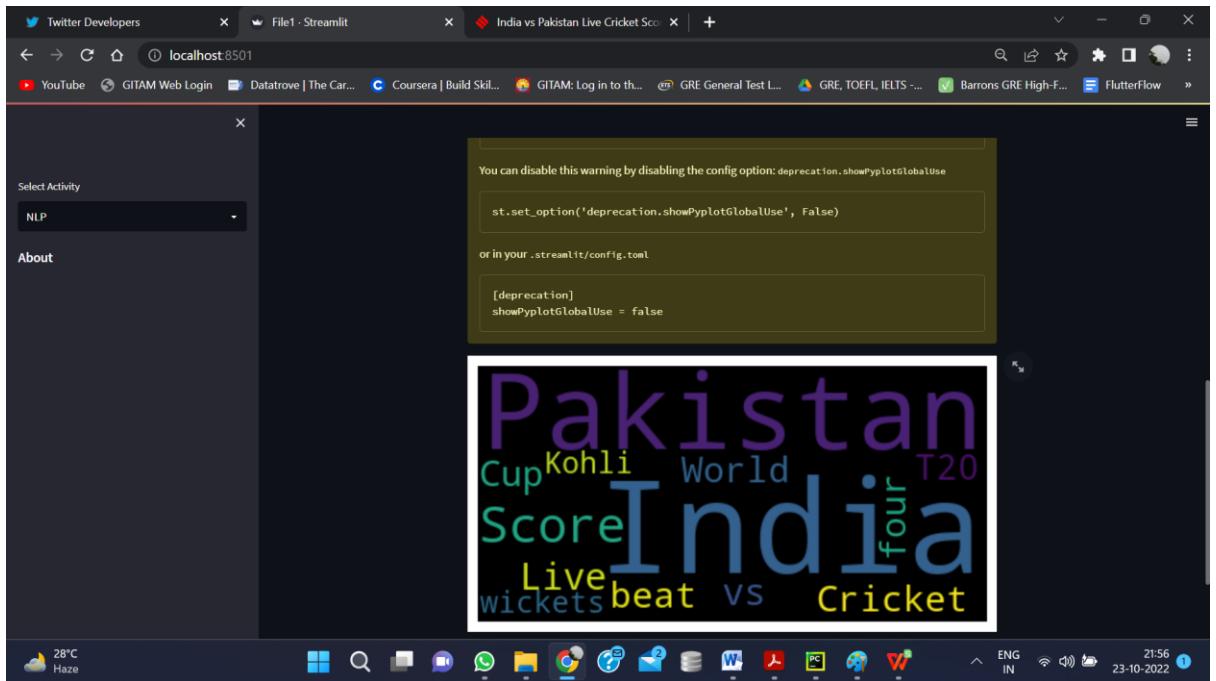
28
29 cv = CountVectorizer()
30 X = cv.fit_transform(x)
31 X_train, X_test, y_train, y_
32
33 model =MultinomialNB()
34 model.fit(X_train,y_train)

35
36
37
38 user = input("Enter a Text:")
39 data = cv.transform([user])
40 output = model.predict(data)
41
42 print(output[0])

43
```

7.3 Streamlit Output:





Analyze

Tabulate

	Tokens	Lemma	POS
0	India	India	PROPN
1	vs	vs	ADP
2	Pakistan	Pakistan	PROPN
3	Live	Live	PROPN
4	Cricket	Cricket	PROPN
5	Score	Score	PROPN
6	,	,	PUNCT
7	T20	T20	PROPN
8	World	World	PROPN
9	Cup	Cup	PROPN

WordCloud

PyplotGlobalUseWarning: You are calling st.pyplot() without any arguments. After December 1st, 2020, we will remove the ability to do this as it requires the use of Matplotlib's global figure object, which is not thread-safe.

28°C Haze

ENG IN 21:57 23-10-2022

8. CONCLUSION & FUTURE SCOPE

Twitter Sentiment Analysis Conclusion

Finally, we exhibited a system for analysing textual Twitter data in this Project, with an emphasis on the growing topic of sentiment analysis. The application will use a more accurate machine-based learning approach for sentiment analysis, as well as natural language processing approaches. As a consequence, the program's emotion will be classified into positive, negative, and neutral categories, which will be displayed in the Matplotlib window by a pie chart and a bar graph. These outcomes matter a lot for sales marketing companies and many MNC's for improving their profits and work force

Future Scope in this Project

A large amount of work needs to be done; here, we present a glimmer of hope for future research directions.

Sarcasm Interpretation: At the moment, the suggested approach is incapable of interpreting sarcasm. The use of irony to mock or show contempt is known as sarcasm. Sarcasm reverses the polarity of an ostensibly positive or negative speech in the context of current work. Conclusion and Future Work "A comprehensive understanding of the basics of discourse-driven sentiment analysis can overcome this constraint." This approach's major purpose is to scientifically discover lexical and pragmatic characteristics that distinguish between ironic, positive, and negative word use.

Multi-lingual support: Developing a multi-language lexical dictionary is presently not viable due to the lack of a multi-lingual lexical dictionary.

- Additional research may be done to make the classifiers language independent.
- The authors developed a sentiment analysis method based on support vector machines; we might adopt a similar technique to make our system language neutral.
- By analysing emoji/smiley feelings to determine neutrality, we may enhance our data gathering and analysis approach.
- Future study can be carried out with potential enhancements such as improved data and a more accurate algorithm.

News Article Prediction and Classification Conclusion

This Project compares the prediction performance of the multi-class category predictor. Using a BBC news dataset with five categories, well-known machine learning methods (Nave Bayes, Random Forest, K-Nearest Neighbour, and Support Vector Machine) were used to build news category predictors (business, sports, technology, politics, and entertainment). The Confusion Matrix was then analysed using performance assessment methods to evaluate the test dataset's Precision, Recall, and Total Accuracy. As a result, the MultinomialNB model was shown to be the most accurate of the four supervised learning models in categorising news articles, with 80% accuracy. The accuracy of the LR model is the lowest (65). However, the performance of the KNN model may be improved by finding the optimal number of neighbours K. The Accuracy of the Model change from dataset to dataset. As the dataset size increases the computational time also increases. This Project also provides and gives the insights of different NLP results for News Text such as Lemmatization, Tokenization, and word cloud also.

This project's working knowledge may be expanded to new use cases, and its accuracy can be improved by utilising numerous prospective technologies in the disciplines of machine learning, data science, and artificial intelligence. We may use the notion of Neural Network to forecast news stories much more correctly. GRU, a Deep Learning idea in development, can also be implemented in the future to eliminate the requirement for time-consuming algos. As the study continues, GRU believes that this technology can also be applied in the sentiment analysis project. In the future, news stories will be automatically categorised and categorised using AI technology that attempts to grasp the structure by analysing its word usage by enlisting the assistance of NLTK

9. REFERENCES

- [1] Pang and Lee 2008(Opinion Mining and Sentiment Analysis)
Foundations and Trends in Information Retrieval Vol. 2, Nos. 1–2 (2008) 1–
135 c2008 B. Pang and L. Lee DOI: 10.1561/1500000001
- [2] Hatzivassiloglou and McKeown 1997(Predicting the Semantic
Orientation of Adjectives). Proceedings of the 8th Conference on European
Chapter of the Association for Computational Linguistics Madrid, Spain,
174-181.”);
- [3] Esuli and Sebastiani 2006(SENTIWORDNET: A Publicly Available
Lexical Resource for Opinion Mining) Proceedings of the Fifth International
Conference on Language Resources and Evaluation (LREC’06)
Determining Term Subjectivity and Term Orientation for Opinion Mining
(11th Conference of the European Chapter of the Association for
Computational Linguistics)
- [4] U Suleymanov and S Rustamov (Automated News Categorization using
Machine Learning methods) IOP Publishing Aegean International Textile
and Advanced Engineering Conference (AITAE 2018) IOP Conf. Series:
Materials Science and Engineering459 (2019) 012006 doi:10.1088/1757-
899X/459/1/012006