

fyp-report-2022

by Arsalan Javed

Submission date: 07-Jun-2022 08:51PM (UTC-0700)

Submission ID: 1852700115

File name: FYP-report-latest.docx (3.34M)

Word count: 3966

Character count: 20631



**NUST COLLEGE OF
ELECTRICAL AND MECHANICAL
ENGINEERING**



**Anti-state Comment Detection on Social Media
using Natural Language Processing**

PROJECT REPORT

DE-40 (DCE)

Submitted by

PC ARSALAN JAVED

GC RAAZ KHAN

NC ARSAL SANA

NC DANIYAL AZFAR

BACHELOR'S IN

COMPUTER ENGINEERING

YEAR 2022

PROJECT SUPERVISOR

DR. WASI HAIDER BUTT

COLLEGE OF

ELECTRICAL AND MECHANICAL ENGINEERING

PESHAWAR ROAD, RAWALPINDI

**Anti-state Comment Detection on social media
using Natural Language Processing**

PROJECT REPORT

DEGREE 40

Submitted by

PC ARSALAN JAVED

GC RAAZ KHAN

NC ARSAL SANA

NC DANIYAL AZFAR

BACHELOR'S IN

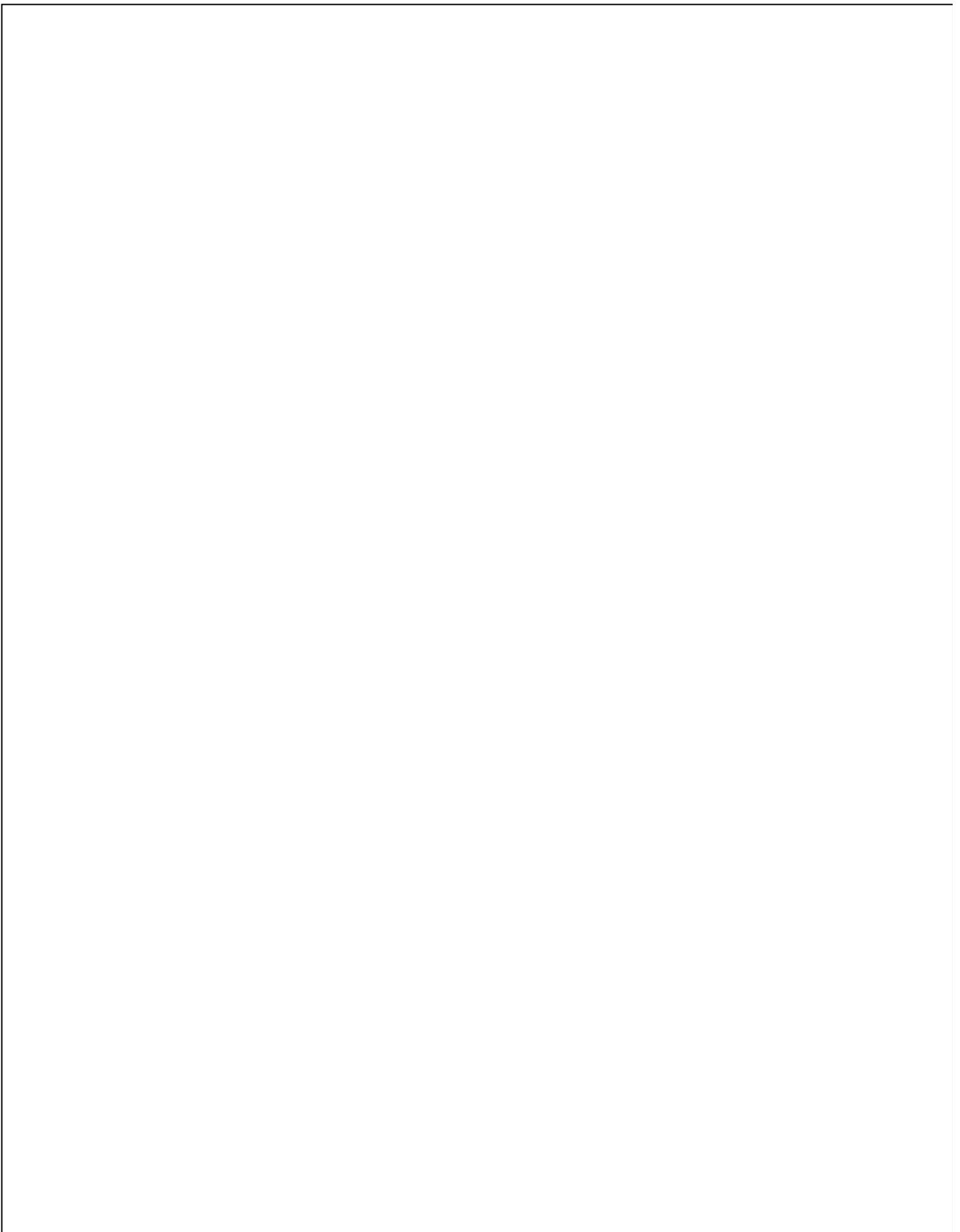
COMPUTER ENGINEERING

Year 2022

PROJECT SUPERVISOR

DR. WAS HAIDER BUTT

**COLLEGE OF
ELECTRICAL AND MECHANICAL ENGINEERING
PESHAWAR ROAD, RAWALPINDI**



ABSTRACT

“Anti-state Comment Detector” aims to take a step towards the betterment of social media in Pakistan. It is the step to introduce AI in the social media domain to benefit public and organizations.

With increase in use of Internet and the step into digital world, various applications have been developed to make the life of the people more comfortable.

The increased use of social media platform has changed the look of the way people communicate and interact physically and digitally. On social platforms discussions are done which makes the trending topics and gives the bird eye view of events that are occurring around the globe in real-time.

“Anti-state Comment detector” plays the role in current identifying the users of Twitter that mislead the mass audience with their negative thoughts and comments against the country and spread negativity. The project uses a classifier that automatically detect hate tweets using the content of tweets made. It also extracts the data of tweets containing hate speech and the user who made the tweet.

CHAPTER 1:

INTRODUCTION CHAPTER 1

1.1 INTRODUCTION

Over the decades, comprehensive articles have been published about how social media has impact on the participation of politics. The argument here is made that almost all social media platforms encourage online and offline participation of political trends and talk.

Pakistan's political parties are actively using social platforms and social media and it has resulted in huge change of how politics work in Pakistan. .^[1]

Other than political parties' common masses also indulge in political discussion using tweets and comments to follow the current ongoing trends.

Twitter's advertising resources indicated that Twitter had 3.40 million users in Pakistan in early survey of 2022.^[2]



Figure 1-Social media Statistics of Pakistan April 2022

Twitter, a widely used social media platform, produces large amount of text that contains the political information also, this information can be mined and extracted to analyze people's point of view regarding the hot trending topics going around the state and globe.

1.2 Motivation

We have chosen to work with Twitter as it is widely used platform of public emotion than regular online website and blogs. The reason is that the amount of relevant data that can be scanned and extracted is much larger on Twitter, compared to traditional website blogs. The response of people on Twitter post and trends is very fast and common as many people prefer using Twitter rather than blogging to express their point of views. The use of social media has seen dramatic growth in recent years. With the use of only 140 characters on Twitter, people can express their opinions on any topic.

Pakistan being a democratic state allows the freedom of expression and speech to the citizens in all contexts including politics under its constitution.^[3] Social media is playing a role in increasing the violence and extremism of emotional people, and the government is designing application and ways to detect, understand the impact of extremism on social media and other internet social platforms.

Recently Government of Pakistan has passed a law that any person who defames the country or passes anti-state comments on social media is liable to prison for up to 5 years and will be fined heavily.^[4]

Our “Anti-state Comment Detection System” motivation is to help our government in this regard. The project will detect the anti-state comments using sentiment analysis. So every time a person violates the law authorities can be notified through it.

1.3 Objectives

Objectives or deliverables of our Final Year Project that we have worked on are as follows:

- Checking every tweet in in-the country domain
- Detection of an anti-state tweet or comment the by user using machine learning AI, NLP
- The detected user making a hate comment or tweet will be blacklisted for authorities to take further action according to law

1.4 Aim and Working

The “Anti-State Comment Detector” aims to take a step toward better control and implementation of media law. It helps the concerned authorities to keep check and balance of people and data being uploaded on social media against defamation of any institute, person, and above all our state Pakistan.

The social media platform has opened gates to new ideas for gathering and exchange of information. Moreover, it also gives ample opportunity for people to raise their words and voices in public and spread their thoughts to masses. Due to the dangerous nature of the social media platform, some stories are rapidly spreading, are getting highlighted, and need careful investigation to have a clear picture. With the development in the digital and internet world, various applications have been developed to make the life of people more comfortable and keep country security and respect safe and sound.

“Anti-State Comment Detector” targets the users of Twitter that try to defame or spread hate among groups using social media Twitter platform to defame the country, it identifies such users and shows the tweet or comment made by them with their account username and other data. The tweets were fed into the system a from file right now and they predict the sentiment of tweet as “Anti-State” or “Non-Anti-State”. Further it can be implemented using Twitter API keys for real-time tweet analysis.

Social media warfare can be defined as the use of social media as a form of weapon with the intention of causing serious permanent damage to reputation of certain actors such as governments or corporations. Various strategies and techniques and techniques are used to advance the political, economic, social, or cultural agenda.

1.5 Development Environment

The development phase was divided into two:

- Backend Development
- Frontend Development

Backend Development:

The main working of the system and algorithm are included in backend development

Coding Language:

Python

Coding Environment: PyCharm IDE and command terminal

Packages used: Python libraries of machine learning and natural language processing



Figure 2-Tools for Backend Development

Frontend Development:

The front end of the web application was made using:

Coding Language:

Python,HTML,CSS

Coding Environment: Pycharm, Jupyter notebook
Servers: Flask
Databases: SQLite3



Figure 3-Tools for Frontend Development

1.6 Structure of Report

- Chapter 1 deals with the Introduction. It includes the motivation behind doing project, objectives covered and aim and working.
- Chapter 2 mainly deals with Literature Review. It includes the research which has already been done in natural language processing by modern technologies and Twitter Sentiment Analysis
- Chapter 3 covers the objectives. It enlists the details of objectives being covered in our project and its working.
- Chapter 4 covers the Software part including the front-end application being used. Its features and working are mentioned.
- Chapter 5 covers the references used in report

CHAPTER 2:

LITERATURE REVIEW

Chapter 2

2.1 LITERATURE REVIEW

Looking in past, sentimental analysis has been an area of interest for many professional fields including psychologists, neurologists, computer scientists, and linguists.

Various methods of researching various topics especially analyzing social networking websites such as Twitter, Facebook, and Instagram are increasingly in demand and use by people nowadays. The data available on these social media sites made it easy to research and analyze what are the feelings and opinions of people around the globe or state. An analysis of these expressions made publicly on these platforms can help in finding the faces of user that may help in different decision making process.

Existing Solution in Pakistan:

Sentiment Analysis on Imran Khan's Tweet is a research paper done by Air University students in which they study and analyzed the sentiment of tweets made by Imran Khan in one year of spam. In this research they read and analyze the emotions of Imran Khan's tweets in year's spam [5]

¹ **Sentiment Analysis for Urdu News Tweets Using Decision Tree** is a research paper done by ¹EME, a NUST University student in which they make sentiment analysis which is done in Urdu news tweets. The proposed approach has two steps. In the first step the data processing is done as the hashtag tag removal and suspension name removal is done. The vector feature of the second step is designed. The Vector feature is created by identifying several encouraging words, opposing words, and the presence of denials and the use of POS tags. After constructing the element vector, the decision tree is used as the division algorithm. The decision tree classifies tweet as positive, negative, and neutral. [6]

Sentiment Analysis of Pakistani Twitter is a research done by Farooq Yousaf,

based in Australia having Ph.D in Politics, he has made Pakistani Twitter comments about the response of those affected by the coronavirus crisis in the country, and the findings are interesting. Due to Twitter's restrictions on data extraction, a person may be able to extract a certain number of tweets in each time. Therefore, the following emotional analysis is based on 500-1000 randomly entered tweets for all three participants. This analysis is designed to measure Pakistan's typical Twitter feelings about how we react to various major stakeholders in the country.^[7]

CHAPTER 3:

MODULE WORKING

Chapter 3

3.1 Sentiment Analysis for Anti-State Tweet Detection

Introduction:

The increase in growth of internet users and the emerging power of online review and freedom of expression on social media platforms have provided the platform to start Sentiment analysis. Sentimental analysis comes under the work that includes ^[20] text mining to test the opinions of people who talk about a particular product, topic and commentary on social media posts and the review site which also raises difficulties in digging ideas. Figures published on Twitter advertising services show that Twitter has 3.40 million users in Pakistan by early 2022. This figure means that access to Twitter ads in Pakistan was equivalent to 1.5 percent of the population at the time. ^[8]

The objective of Anti-State Tweet Detection in our Final Year Project is to present a way to provide a solution to manage or control the problem of defamation and spreading of hate against any person or country on public platform like Twitter.

Dataset:

The collection of datasets for this project was not a simple task. We applied several attempts for obtaining Twitter API keys from the Twitter developer account but every time Twitter got it denied due to unknown reasons. After many attempts we were able to get the keys, those keys were for Essential Twitter API v2 which gives limited access to collect tweets. We

were not in success for obtaining Elevated Twitter API v1 keys. This very reason made the task difficult as manual extraction of dataset is very difficult. Due to, this fact we collected some tweets using v2 API but after that we shift to a non-API-based i.e., Twint. Twint is a high-level Twitter tool written in Python that allows for scanning and extracting of Tweets from profiles of Twitter users without having access to developer account Twitter's API.

FYP

Comment detector-FYP-NLP

[Settings](#)

[Keys and tokens](#)

Consumer Keys

API Key and Secret ⓘ

ⓘ [Reveal API Key hint](#)

[Regenerate](#)

Authentication Tokens

Bearer Token ⓘ

Generated June 1, 2022

[Revoke](#)

[Regenerate](#)

Access Token and Secret ⓘ

Generated June 5, 2022

For @Arsalan39716407

Created with [Read and Write](#) permissions

[Revoke](#)

[Regenerate](#)

OAuth 2.0 Client ID and Client Secret

Dataset creation:

The dataset consists of the raw tweets scrapped using Twint and labels of the tweet as 0 and 1. Here 0 means “Anti-State” and 1 means “Non-Anti-State”.

Twint uses Twitter search to let you extract Tweets from specific users, extract Tweets related to specific topics, hashtags, and trends, or edit sensitive information in Tweets such as emails and phone numbers. I find this very useful, and you can find art with it too.

Twint also makes special inquiries on Twitter allowing you to scan Twitter user followers, Tweets that the user likes, and who they are following without confirmation, API, etc.

- **Get Tweets:** Retrieve the list of tweets for specific hashtags, trends, and users using Twint
- **Filter cases:** Tweets that were not in the English language or which contains data that was irrelevant for our dataset were removed
- **Get Tweets as a document:** The tweets information that is left after filtering is downloaded in form of .csv file for future use.
- **Document Preprocessing:** The tweets document dataset was then pre-processed to remove all the URLs, stop words, and hashtags from the tweets and are cleaned to be fed into the classifier

1	Dont worry about Kashmir fanpageofirk	1514352112414690000
2	The world knows that Paki alichghulam984	1385055412270070000
3	We do not need your sympa_proudcivilian	1418868272
4	Pakistan is a terrorist state pakiswift	1123668694704450000
5	was brought to the power hinam74457396	1492358833632690000
6	A country who harbours nu ibrakhansafi	1456197601888450000
7	American strategists though ronnyinblack	1467038331351540000
8	This is the real face of Terr khannn108	1427160859806140000
9	ISI Taliban are the ISI a meme_sonambulo	1523415922991640000
10	Instead of declaring PAK a pakistai_boy	1530406024557260000
11	The rapist Pak army and I mhd_binasghar	1477385333692110000
12	Pakistan zindabad ISI zinda uraa47	1037952455080910000
13	Get lost loser Your drama cnotcoas	1285267648528420000
14	Dont worry about Kashmir kingofempire369	1039632414274980000
15	India bombed terrorist cam uraa47	1037952455080910000
16	Imran has already accepted mahiraa67566583	1217436694116470000
17	The world knows that Paki chaudriifan1	1207029735160780000
18	Biggest Victory of Pakistan cayesha_mian_	1513754195731610000
19	Baba Jo bhi ho danke ki chc waseem1785	1085555942693970000
20	From decades Pak was den coun1my	1514199309364990000
21	Without any second though ahmadbe61989504	1451037182488160000
22	Ap andhy ho ya Samajh Kuci acirfanahmad40	1512592519267240000
23	Baloch area isnt directly c 234iftikharkami	2797823330
24	You are also a hypocrite by rz_khattak	272064747
25	Two Pak Army soldiers martshireenmazari1	481156800
26	Fake ID heres a Thing Suici themehfuz143	1378112112996800000
27	Adha sa zyda comments sal shahmir57	1472542119235620000
28	BLA used to carry out such a farazkh78685502	1386982718525810000

Get Tweets:

Using Twint Twitter in Python we retrieve the list of tweets using the hashtags or keywords for specific tweets and scrapped all the data from Twint into a .csv file and .json file for future use.

Filtering Cases:

The following cases are filtered out:

- English Tweets are kept
- We only use the tweets coming from the Pakistan domain
- The tweets that were too long were skipped

```
def main(i,o):
    data = json.load(open(i,encoding="UTF_8"))
    new_data = {'items': []}
    new_data['items'] = [{'_id': item['id'], 'text': item['text'], 'username': item['username'], 'created_at': item['created_at']} for item in data['items']]
    json.dump(new_data, open(o, 'w'), indent=2)
    print('done')
```

```
import re
from nltk.corpus import stopwords
x = tweet
def cleanText(x):
    import re
    from nltk.corpus import stopwords
    x=x.encode('ascii','ignore').decode()
    x=re.sub(r'https*\S+', '',x) # remove urls
    x=re.sub(r'@\S+', '',x) #remove mentions
    x=re.sub(r'#\S+', '',x) # remove hashtags
    x=re.sub(r'\w+', '',x)
    x=[i for i in x if i not in string.punctuation] # remove punctuations
    x=''.join(x)
    return ''.join(x)
```

Getting Tweets as Documents:

The data is downloaded in only a .csv file so that it can be used in classification as Data Frames and easily readable.

	tweet	length	label
0	Pakistani Terrorist Intelligence ISI and MI re...	245.0	1.0
1	Pakistani Terrorist Intelligence ISI and MI re...	245.0	1.0
2	Pakistani Terrorist Intelligence ISI and MI re...	264.0	1.0
3	Pakistani Terrorist Intelligence ISI and MI re...	239.0	1.0
4	Pakistani Terrorist Intelligence ISI and MI re...	242.0	1.0
...
1635	In sha Allah wo difah i Pakistan k liayay hom...	70.0	0.0
1636	Imran khan ne sabit kar deya hsi k wo allam a...	169.0	0.0
1637	Ameen	6.0	0.0
1638	You are face of GilgitBaltistan You are pride...	59.0	0.0
1639	Aap b kin ko keh rahay hein ye opn jang ha au...	173.0	0.0

Figure 4- Document Structure of csv file

Pre-Processing Tweets:

In this, we use panda's library to read our csv files containing tweets. The data of tweets were checked and if any missing row or column value is found it was dropped. Then the dataset was determined for its balance notion i.e., a number of anti-state and non-anti-state tweets acquired.

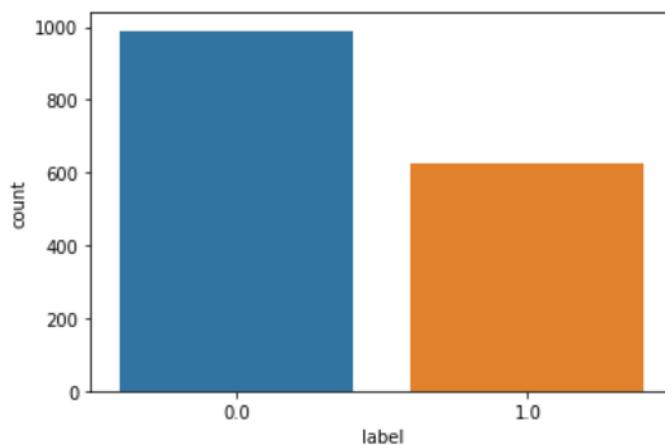


Figure 5- Number of tweets in each label

Other steps in pre-processing of the data includes the removal of duplicate tweets, removing spaces from the beginning of tweet and s, and removing digits, punctuations, hashtags, and URLs.

Transforming Tweets into Vectors for features extraction:

To use text data to obtain a prediction model, text must be separated to delete certain words - this process is called tokenization. These words need to be encoded as whole numbers, or floating-point values, in order to be used as input into machine learning algorithms. This process is called extracting features (or vectorization).

Scikit-learn's Count Vectorizer is used to convert text collections into form of vectors also termed as token count. It also enables pre-processing of text data before generating vector representation. This functionality makes it a module for the most flexible feature of the text.

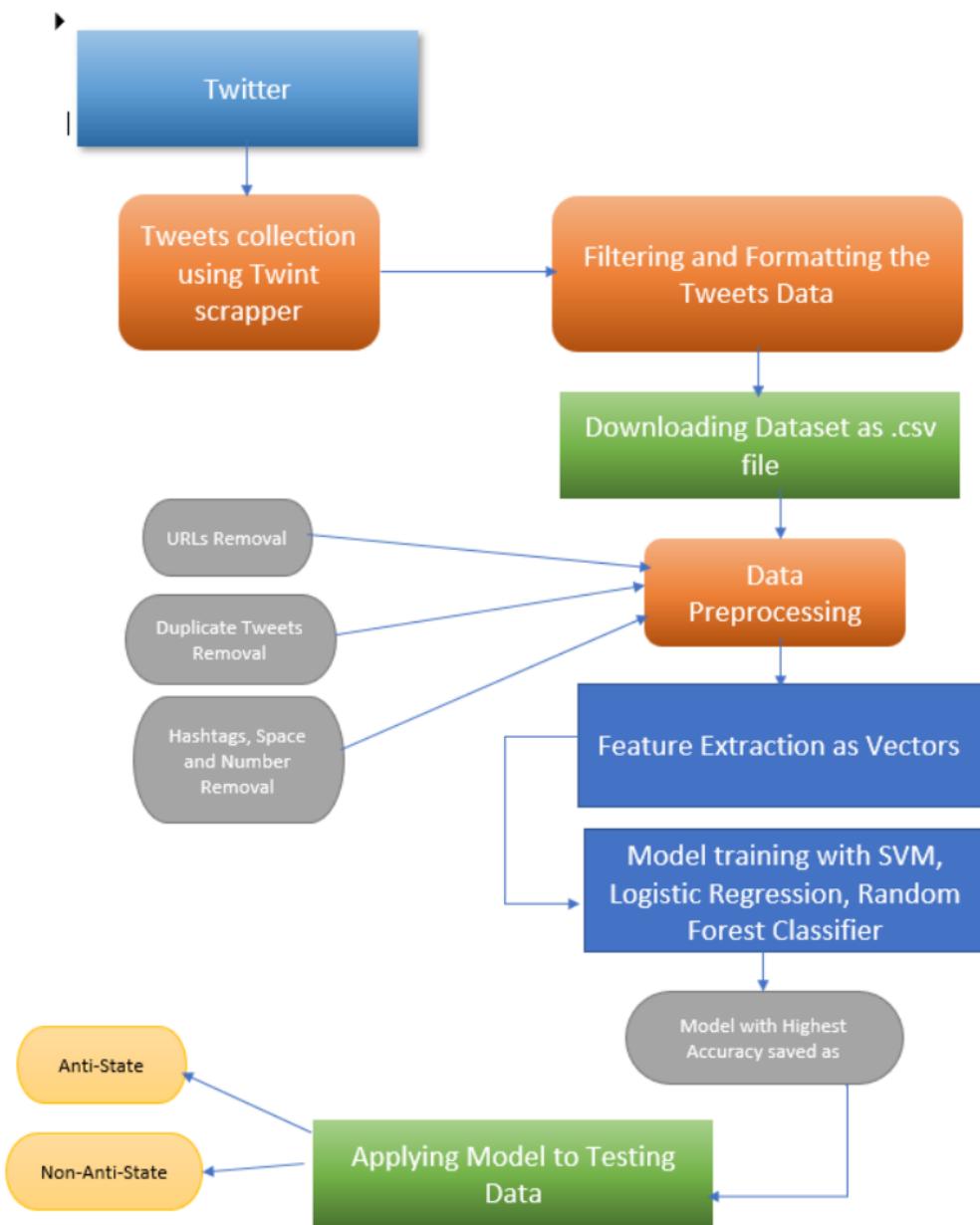
```
Data = ['The', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog']
```



The	quick	brown	fox	jumps	over	lazy	dog
2	1	1	1	1	1	1	1

Figure 6- Example of how vectors are made

Flow Chart of Algorithm:



Features of Dataset:

We know that many applications interact with multiple datasets. Therefore, a non-computerized function can be a huge congestion in your algorithm and can take effect on a model that takes memory and many computational hours to run. To ensure that the code works properly on a computer, we will use a process called vectorization.

One of the major difficulties any NLP Data Scientist faces in selecting the best representation of the numerical and vector form of your string data for using Machine Learning models.

Count Vectorizer is a process of converting the input string data into form of numerical data as frequency i.e. number of occurrence of word in text

8

TF-IDF means Term Frequency - Inverse Document Frequency. This is a calculation based on multiple times a word in the corpus, but it also provides numerical representation of how important the word is in mathematical analysis.

3

The TF-IDF is better than Count Vectorizers because it not only focuses on the wording of the existing words on the corpus but also provides the significance of the words. We can then remove the words that are less important in the analysis, thus making the model structure less complex by reducing computational speed and hours and the input size.

```
'burdasht',           Vector features  
'burma',  
'burn',  
'burning',  
'burnt',  
'burntterrorists',  
'burqa',  
'bury',  
'bush',  
'business',  
'businessman',  
'busted',  
'busy',  
'but',  
'butcher',  
'buy',  
'buyers',  
'bw',  
    '  
    'aa',  
    'aaate',  
    'aadhaar',  
    'aafia',  
    'aag',  
    'aga',  
    'aai',  
    'aajayegi',  
    'aakar',  
    'aake',  
    'aakhri',  
    'aakr',  
    'aam',  
    'aap',  
    'apne',  
    'aaropi',  
    'aata',
```

Classification Model:

10

1. Support Vector Machine (SVM) classifier

The first classifier that we used for the classification of our testing data into “Anti-state” OR “Non-Anti State” is the Support Vector Machine (SVM). We trained our dataset on this classifier using the sklearn python library.

6

SVM is an AI algorithm, and the concepts are simple. The SVM classifier separates data points using a hyperplane with many margin space. That is why the SVM classifier is also known as discriminative classifier. SVM detects the right hyperplane that helps to separate new data points.

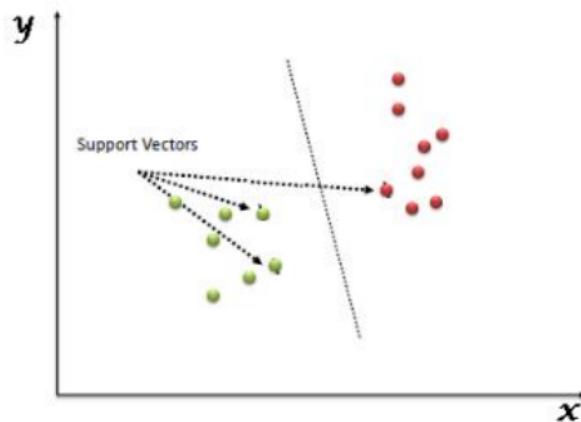


Figure 7- SVM working

11 Support Vectors are simply the coordinates of individual observation data points. The SVM classifier is an AI algorithm that best separates the two classes (hyper-plane/ line).

The SVM algorithm is implemented using a kernel. A kernel is used to transform your dataset data points into the required form of space. The technique used by SVM classifier is known as Kernel trick. Here, the kernel takes a low-dimensional input space and transforms it into a higher-dimensional space. By adding more space to the dataset kernel converts the low dimensional space into higher dimensional space. It is most useful in non-linear separation problems. Kernel trick is used to increase the accuracy of the classifier.

$$Precision = \frac{TP}{TP + FP}$$

TP = True positive

TN = True negative

$$Recall = \frac{TP}{TP + FN}$$

FP = False positive

FN = False negative

$$F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

The kernel we used in our classifier was: Linear kernel

Accuracy for the SVM classifier came out to be 86%

```
test_data=pd.read_csv('testingdata.csv')
test_data=test_data.dropna()
X_test=vectorizer.transform(test_data['tweet'])
Y_test=test_data['label']
```

```
pred=model.predict(X_test)
```

```
print('Accuracy of SVM:',100*accuracy_score(pred,Y_test),'%)
```

Accuracy of SVM: 86.01694915254238 %

Classification report of SVM:

```
print(classification_report(pred,Y_test))
```

	precision	recall	f1-score	support
0.0	0.89	0.90	0.89	154
1.0	0.80	0.79	0.80	82
accuracy			0.86	236
macro avg	0.85	0.84	0.85	236
weighted avg	0.86	0.86	0.86	236

2. Logistic Regression classifier

The second classifier we used to test our dataset Logistic Regression Classifier. **Logistic Regression** is a Machine Learning classification algorithm used to predict the probability of class dependence. In logistic regression, the dependent variable is a binary variable that contains coded data such as 1 (yes, success, etc.) or 0 (no, failure, etc.). In other words, the regression model predicts P (Y = 1) as X function. Accuracy for the Logistic Regression classifier came out to be 84%

```
print('Accuracy of Logistic Regression:', 100 *accuracy_score(predlg,Y_test), '%')
```

```
Accuracy of Logistic Regression: 84.7457627118644 %
```

Classification report of Logistic Regression:

```
print(classification_report(predlg,Y_test))
```

	precision	recall	f1-score	support
0.0	0.90	0.87	0.89	159
1.0	0.75	0.79	0.77	77
accuracy			0.85	236
macro avg	0.82	0.83	0.83	236
weighted avg	0.85	0.85	0.85	236

25

3. Random Forest classifier

A **random forest** is a **classifier** that uses several decision trees on sub-sample of datasets and calculate the average to increase the accuracy of test data prediction and this process control the over-fitting of dataset. The size of sub-sample tree is controlled by the parameter called `max_samples` only if parameter `bootstrap` is `True`, otherwise the classifier uses whole dataset to build up the tree instead of using sub-samples.

15

“**Random Forest is a classifier**” that contains multiple **decision trees** on **multiple subsets** of given **dataset** and then it **takes average** of all those **subsamples**. This process is used to increase the accuracy in prediction of testing labels.^{9]}

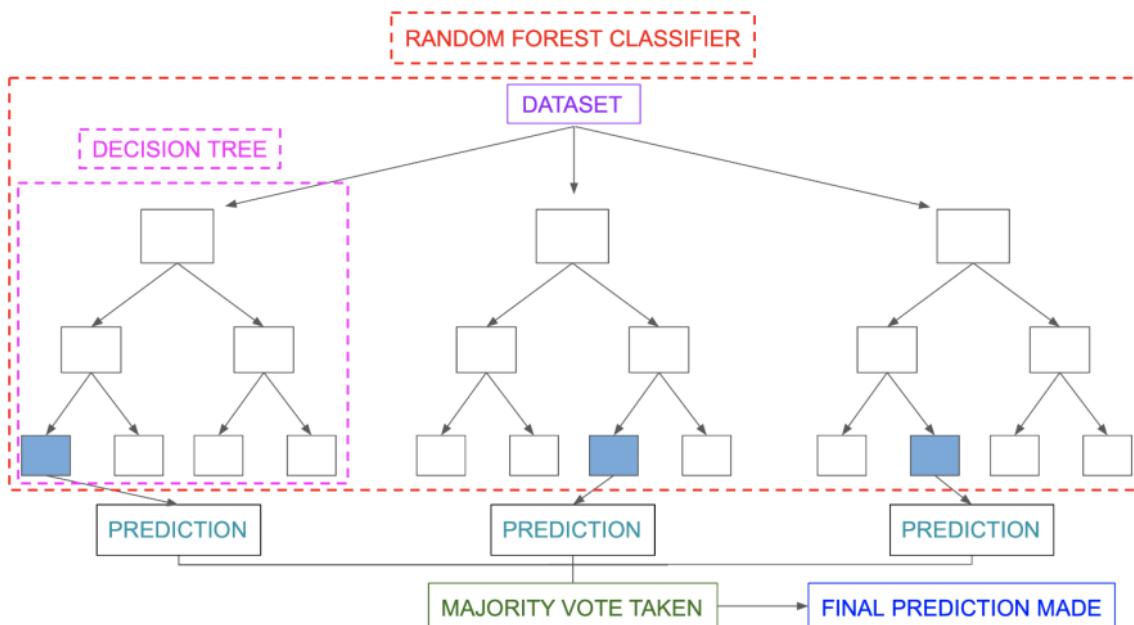


Figure 8- Random Forest Decision Tree Working

Accuracy for the Logistic Regression classifier came out to be 86.4%

```
: print('Accuracy of Random Forest:',accuracy_score(predrandomF,Y_test),'%')
Accuracy of Random Forest: 0.864406779661017 %
```

Classification report of Random Forest:

```
print(classification_report(predrandomF,Y_test))

precision    recall   f1-score   support
          0.0      0.90      0.89      0.90      157
          1.0      0.79      0.81      0.80       79

      accuracy                           0.86      236
   macro avg      0.85      0.85      0.85      236
weighted avg      0.87      0.86      0.86      236
```

4. KNN classifier

K-Nearest Neighbor is one the simplest Supervised Learning Machine Learning algorithm. K-NN works on the process of assuming the similarity that exists between new case i.e., data and available cases in this way it puts new coming cases i.e., test cases into the category that is most close the available categories.

It stores all the data available and classifies the new data based on the previous similarity. K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for Classification problems. It is a non-parametric algorithm which means that it does not take any assumption of underlying data. It is also known as a lazy learner algorithm due to the reason that it does not learn from the training data immediately in fact stores

the data and uses it at time of classification and perform the action on dataset when new data comes and categories it.

Accuracy for the KKN classifier came out to be 72%

```
print('Accuracy of KNN:', 100 *accuracy_score(predknn,Y_test), '%')
```

```
Accuracy of KNN: 72.45762711864407 %
```

Classification report of KNN:

```
print(classification_report(predknn,Y_test))
```

	precision	recall	f1-score	support
0.0	0.94	0.72	0.82	200
1.0	0.32	0.72	0.44	36
accuracy			0.72	236
macro avg	0.63	0.72	0.63	236
weighted avg	0.84	0.72	0.76	236

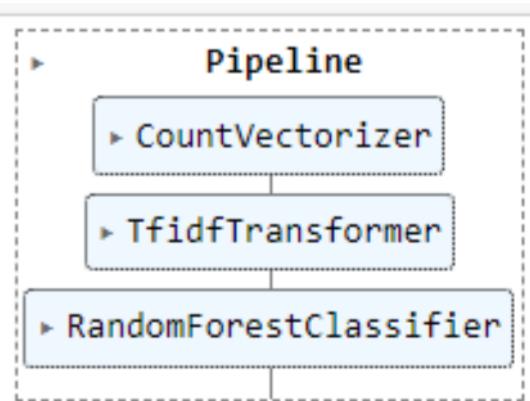
As in the above-mentioned classifier the SVM and Random Forest gives good accuracy. Among both classifiers, we choose Random Forest Classifier to use for our testing data as it gives maximum accuracy among all classifiers we used.

Pipeline to Save Model:

We can create a machine learning pipeline by putting the features that were involved in training the model that we want to use for our dataset and the classifier in sequenced steps. This process will automate the workflow of machine learning. Pipeline can consist of pre-processing, extracting features, classification model and post-processing of data.^[10]

```
pipe=Pipeline([
    ('vectorizer',CountVectorizer()),
    ('trnasformer',TfidfTransformer()),
    ('model' , RandomForestClassifier())
])
pipe.fit(X,trainY)
```

We made a pipeline of our model which contain the s feature extraction method of the vectorizer and TFidF Transformer. Moreover, it also contains the classification model we selected for our testing data.



Accuracy for pipeline came out to be 86%

```
print("Pipeline accuracy :", 100 *accuracy_score(pred,test_data['label']), '%')
```

```
Pipeline accuracy : 85.59322033898306 %
```

Classification report of Pipeline:

```
print(classification_report(pred,test_data['label']))
```

	precision	recall	f1-score	support
0.0	0.90	0.89	0.89	157
1.0	0.78	0.80	0.79	79
accuracy			0.86	236
macro avg	0.84	0.84	0.84	236
weighted avg	0.86	0.86	0.86	236

Save Pipelined Model:

The purpose of saving the pipelined model is to save the data trained under certain features and classifiers and reuse them for other purposes or loading it to run predictions on test data instead of training the whole model every time.

We used the Pickle library of Python to save our Model as “my_model.pkl” for further use.

3.2 Future Improvements in Project

Some of the ideas that can be future added to the project could be as follows:

- **Real-time data collection for Anti-state tweet and comment detection:** As mentioned in the above section we faced difficulty in making large dataset due to unavailability of Twitter API keys v1. To further enhance this project, we can try to get our application for Twitter development account approved and make the project on real-time data.
- **Fake News Detection:** In recent years, we have seen an increase in false stories, that is, almost pieces of false information created with the intent to deceive. The proliferation of this type of issue poses a serious threat to the unity and well-being of the community, as it promotes political divisions and mistrust of its leaders. The sheer volume of information disseminated through social media makes it possible for manual authentication, which encourages the design and implementation of automated false detection systems. The creators of false stories use a variety of style tricks to promote the success of their creations, one of which is to appease the feelings of the recipients.
- **Extended project to detect defamation not only of state but instead also detect Racism, sexism, homophobia, religious extremism, and conspiracy theories.** As a result of the experimental research project, the Anti-Defamation League can estimate that at least 4.2 million anti-Semitic tweets were distributed between years' time, the number of anti-Semitic tweets from a low of 36,800 in mid of the year to a high of 181,700 in the end of year. The average number of anti-Semitic tweets in a year of analysis was 81,400. [11]

CHAPTER 4:

FRONT END APPLICATION

Chapter 4

4.1 Front End Application

Front End GUI application:

The main web application is developed using Pycharm and Flask Python as web framework servers.

Why is Flask used in Python?

Flask is a module library of Python language build to assist coder in Python to link their work with a GUI. It is a web framework that provide very useful tools and features to create the web application on Python easily and effectively. Flask gives coders and developers the ease of accessing framework as it is free and available. It is available by just creating one python file.

Development Environment: PyCharm IDE

Web Framework Server: Flask

Programming Languages: Html, CSS, python

4.2 Running Server

The server will run on the local host <https://127.0.0.1:5000/>

```
C:\Users\Arsalan Javed\Desktop\FYP-FINAL\FYP-Anti state>flask run
 * Serving Flask app "run.py" (lazy loading)
 * Environment: development
 * Debug mode: on
 * Restarting with watchdog (windowsapi)
 * Debugger is active!
 * Debugger PIN: 141-220-676
 * Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
 * Detected change in 'C:\\\\Users\\\\Arsalan Javed\\\\AppData\\\\Local\\\\Programs\\\\Python\\\\Python310\\\\Lib\\\\site-packages\\\\setuptools\\\\_vendor\\\\pyparsing\\\\core.py', reloading
 * Restarting with watchdog (windowsapi)
```

17

Flask Routing

Modern web frameworks use a routing technique to help the user remember application URLs. It is useful because it aids user to easily access pages by directly clicking on the button and without always navigating back to home page and selecting the desired page.

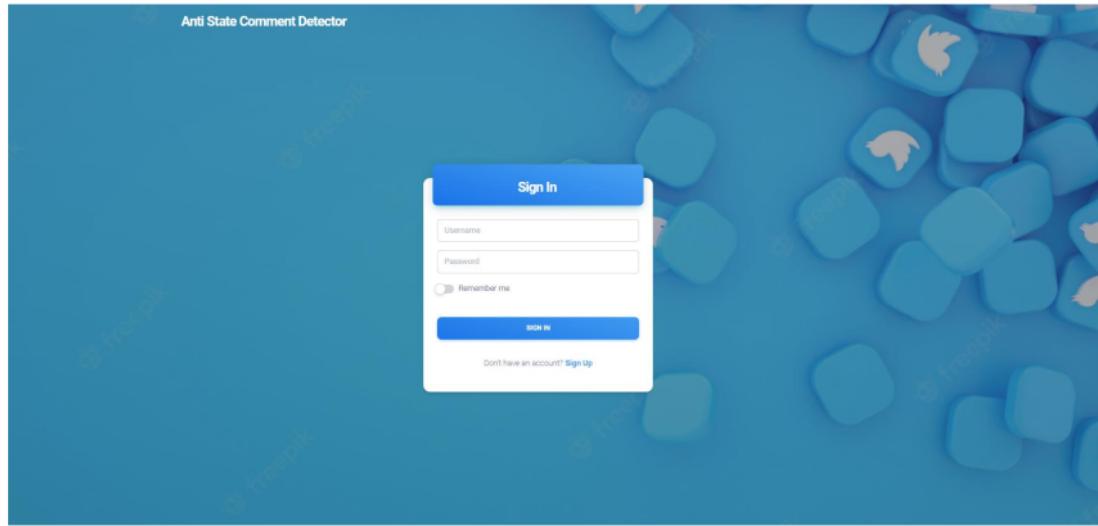
21

Route () decor in Flask is used to link the URL to a function.

4.3 Features

The web applications features are as follows:

SIGN IN/Log Out



The concerned authorities using the application can register them on the site and then sign in to use the application features. We can connect to an SQLite database using the Python sqlite3 module.

Home Page



The image shows the "Objectives" section of the application. It features a navigation bar at the top with "HomePage", "About Us", "Get Started", and a "LOGOUT" button. Below the navigation bar, the word "Objectives" is displayed in a large, bold, blue button-like shape.

Context Based
Context based Detection of anti-state comment or tweet by user using machine learning AI, NLP

Anti-State
Checking of every tweet in country domain

PECA rule
The detected user making hate comment or tweet will be blacklisted for authorities to take further action according to law

Coding Platforms

NLP Python HTML CSS JavaScript Twitter GitHub

Anti State Comment Detector

The screenshot shows a web application interface with a blue header bar containing navigation links: 'HomePage', 'About Us', 'Get Started', 'SIGN IN', and 'LOGOUT'. The main title 'ANTI STATE COMMENT DETECTION ON SOCIAL MEDIA USING NLP' is centered above a white content area. The content area features three data points: '4M+' Twitter Users (with a note about a survey from April 2022), '16%' Percentage User of Twitter (with a note about being among other social media platforms), and '3+' Year of punishment for defamation (with a note about PECA law). Below this, a blue button labeled 'Objectives' is visible. Under 'Objectives', there are two sections: 'Context Based' (using machine learning AI, NLP) and 'Anti-State' (checking every tweet in country domain). A note under 'PECA rule' states that detected users will be blacklisted for further action. At the bottom, a section titled 'Coding Platforms' displays icons for NLP, Python, CSS, HTML, JavaScript, and a logo for a platform like Heroku.

HomePage About Us Get Started SIGN IN LOGOUT

ANTI STATE COMMENT DETECTION ON SOCIAL MEDIA USING NLP

4M+
Twitter Users
Twitter users in Pakistan according to the survey of April 2022

16%
Percentage User of Twitter
Percentage of Twitter user among other social media platforms

3+
Year of punishment for defamation
According to PECA law 2+ years of prison with fine is for defamation

Objectives

Context Based
Content based Detection of anti-state tweet or comment by user using machine learning AI, NLP

Anti-State
Checking of every tweet in country domain

PECA rule
The detected user making hate comment or tweet will be blacklisted for authorities to take further action according to law

Coding Platforms

NLP Python CSS HTML JavaScript H

Anti State Comment Detector

About Us

The screenshot shows the 'ABOUT US' section of the 'Anti State Comment Detector' website. The background is blue with a white header bar. The header bar contains the text 'Anti State Comment Detector' on the left and a 'Logout' button with a user icon on the right. Below the header, the word 'ABOUT US' is centered in white capital letters. A dark grey rectangular box contains the 'Team' heading and five member profiles arranged in two rows. The first row contains two profiles: 'Arsalan Javed' (Member) and 'Daniyal Azfar' (Member). The second row contains two profiles: 'Arsal Sana' (Member) and 'Raaz Khan' (Member). Below these is a third row containing one profile: 'Wasi Haider Butt' (Supervisor). Each profile consists of a small thumbnail image, the member's name, and their role. At the bottom of the page, there is a Twitter logo followed by the text 'Anti State Comment Detector'.

Anti State Comment Detector

Logout

ABOUT US

Team

Arsalan Javed
Member

Daniyal Azfar
Member

Arsal Sana
Member

Raaz Khan
Member

Wasi Haider Butt
Supervisor

Anti State Comment Detector

Get Started

The screenshot shows a web application titled "Anti-State Detection Using Dataset". The title is centered above a form input area. The form contains fields for "Enter Dataset" (with a "Choose File" button showing "No file chosen") and "No. of Tweets" (with a text input field and a "Submit" button). Below the form is a note about dataset requirements and a Twitter logo.

Anti State Comment Detector

Logout

Anti-State Detection Using
Dataset

Note:

Enter your dataset below your data should be clean and processed after uploading your dataset specify the number of tweet you want to predict.

Enter Dataset Choose File No file chosen

No. of Tweets

Submit

Anti State Comment Detector

Tweet Sentiments

[Homepage](#)

Tweets	Username	Id	Predictions
Dont worry about Kashmir as our boby Pakistan is declared by FAIT as Terrorist state amg not India In the name of Kashmir ISI amp rouge Pak Army is taking you guys through garden path In India people have the power to throw out a Govt.	kingofempire569	1039632414274960000	Not Anti State
The world knows that Pakistan is in itself a terrorist state run by the teams of the ISI and the Pakistan Army. Kashmir went to flames at the behest of Pakistan in when BJP was not anywhere near power of Delta and is still burning and shall keep so till Pak fails it	chaudhuryfar1	1207029755160700000	Not Anti State
We do not need your sympathy but controlling THUGS	a_goudekrilum	1418688272	Not Anti State
Pakistan is a terrorist state yesterday terrorist TTP is going to reconcile with ISI we highly condemn all relations among states and terrorist groups ISI/Pak Army/ISIS/Terrorist was brought to the power by Pak Army and ISI to destroy Padoons Thats why he won in KP&K first and then in Federal for Pakistan is a terrorist state	pakismift	1123668694704450000	Not Anti State
The rapist Pak army and ISI are doing genocide in Balochistan They are kidnapping innocent and raping their women Plz raise the voice for my brothers stuck in terrorist state Pakistan	bunam74457796	1492158331326900000	Not Anti State
A country who harbours nurtures helps terrorists since ages as a terrorist state ISI/PAK ARMY help terrorists infiltrate Benoniwe pakistan is a terrorist state and is a TERRORIST America strategists thought they can defeat Taliban in partnership with Pakistan trillion was wasted on a duplicitous Pak Army n ISI if they asked Pakistan for supply routes it wouldve worked out better Declare Pak terrorist state Impose sanctions This is the real face of Terrorist state now its time to never back and a strong reply should be given to Army This of Pak army amg ISI only understand one language that is Brutal war	itsarakhanraf	1456197018884500000	Not Anti State
ISI Taliban are the ISI amp Pak army murdering machine to kill Padoons	rommyblack	1467038151351540000	Not Anti State
Instead of declaring PAK a terrorist state Indian Government is giving lifeline to Pak Army ISI a Terror group Pakistan was a will be an enemy to India a silly to Chma	khann108	1427160839606140000	Not Anti State
The rapist Pak army and ISI are doing genocide in Balochistan They are kidnapping innocent and raping their women Plz raise the voice for my brothers stuck in terrorist state Pakistan	meme_somambu	1224415229914000000	Not Anti State
Pakistan zindabad ISI zindabad Pak army zindabad we are the only audience to have shut down Israel audience Israel is a terrorist state	pakista_boy	1530468024557260000	Not Anti State
Get lost lower Your drama of saving Terrorism in Pakistan is over Terrorists like Kalmahan and other RAW and RSS Now save your Terrorist state of India fm RSS and other Terrorist groups Pakistan Zindabad PakArmy Zindabad ISI Zindabad India bombed terrorist camps in POK in retaliation against Pulwama But real master is pak Army Long term Solution Economic sanctions Balochistan govt in exile FAIT Marketing Diplomatic isolation Terrorist list to include pak army/ISI Pakistan declared Terrorist state	midh_banglar	1477385335921100000	Not Anti State
India has already accepted that Pak army and ISI had trained the talibans terrorist in India later was working in Pakistan Pakistan is truly a terrorist state they should learn from PM Modikow to develop any nation	makhaas7266983	1217436694716470000	Not Anti State
Biggest Victory of Pakistan on International Platform Big Step on the face of India terrorist state and RAW/Balochistan Jadhav a declared Indian spy or a terrorist by ICI Salter to PAK ARMY/ISI IMRAN KHAN Baba Jo bhi dahe ki chet pe ho Puri daura kia patra chalay ISI/Pak army govt or terrorists ka Nexus aur Pakistan ek failed or terrorist state ha	ayeha_mian_	1513754195731610000	Not Anti State
From decades Pak was devining and now Pak PM admits its army and ISI created armed militia or properly saying shabir terrorist by depriving its ppl of healtheducation basic as he trying to get Pakistan listed as a Terrorist State worldwide good	waseem1785	1085355942699970000	Not Anti State
Without any second thought India must declare Pakistan a terrorist State because I dont think the crossover terror activities of Pakistan towards India will ever be stopped until the Pak Army amg ISI will remains on the driving seat of any democratically elected Govt there	counlmy	1514199309564990000	Not Anti State
Ap andhy lo ya Samajh Kuch nah asta Yeh BLA Balochistan liberation army ki terrorist jo Pakistan ky khilaf la rahi Ham or terrorism phela rahi Ham ryasat ka kaam in jis se mukl ko Pak karna h yeh terrorist baaa isme ryasat ka kaam Baloch area sent directly connected to afghanistan for terrorist to hide Pak army has specifically mentioned many times that they are away to far side of border Why will talibans take orders of CIA US left amio so Now it doesnt has control over them You are also a hypocrite by only blaming them and not what Pak Army did with them in Balochistan You justify everything that the Pak Army did wrong by calling terrorist Recently civilians were killed in Afghanistan by Pak Army's strike Two Pak Army soldiers martyred during gun battle with terrorists in South Waziristan ISPR	ahmedshee61989304	1451077182488160000	Not Anti State
Fake ID heroes a Thing Sucide or Haran in Islam Killing Innocent people is also Haran in Islam Speaking Terror in blessed month is Forbidden by the One true God Allah swt So theres no way this Terrorist will find a way to enter Paradise Adha is 2014 comments sala shamless country begart log india ka han sy ya ya zahir hota ka TTP or BLA ko funding india for thia Alka hr mumbaa man terrorist attacks hon phr human cleskin nikla gi misaon ka voh to pak army or awan minat la ga	shahmu57	1472542119235620000	Not Anti State
BLA used to carry out such attacks due to its separatist claim But now in TTP styled terrorist attacks particularly suicide bombings targeting CPEC amp Pak Army is indicative of BLATTP Nexus as Baloch separatists dont have history of using suicide bombers previously	farazk74065102	1386982718525100000	Not Anti State
Very sad incident I hope pak army will deal with those terrorists BLA	lostvoices99	779067712816411000	Not Anti State
First of all deeply saddened by this we Pakastani are your brothers secondly baloch liberation army or bla issnt an army its a terrorist organization that will In Sha Allah be taken down by Pak Army you are well informed about cia involvement in Pak	abdulatafat1	1526535316786110000	Not Anti State
Aba viseen chayye ye pakastan ah ha ye afghan ya balock terrorist ha pak army to khid inkoo maari ha kese chayye indan log ha hun wse ab mardan ha yd kro bad tam leg Pakstani army ka maghila to door ki ht ha	ahmedimran6886	1511896038091070000	Not Anti State
Territory from Afghanistan attack Devsagar Sea of North Waziristan young martyrs including three officers of Pak army Pakistanis strong protest with the Afghan government Demand for implementation of Doha Agreement Afghan soil must be cleared of terrorist Pakistan	choksiar_420	1526427499487530000	Not Anti State
Sorry doesnt even cut a Pak army will be serious about combating this All high target areas and people must be on highest alert The female terrorist has a family all must be exposed on TV We need to know who she associated with in last year Hunt them all down Army needs act	abul_af	412765173	Not Anti State
when corrupt vols imposed over the country again terrorism is spreading up across the country At least plus senya have been martyred since april Two Pak Army soldiers martyred during gun battle with terrorists in South Waziristan ISPR	malikm07757bom1	1512914279913110000	Not Anti State
Pakistan and China came without their consent but what about India who is giving money and arms to these idiots It is time for the Pak army to launch commando operations against these terrorists and wipe them about for good	dzurakmal1	436373934	Not Anti State

Deployment

We are using GitHub and Heroku to deploy our app to live servers.

Heroku is a cloud platform supporting several programming languages.

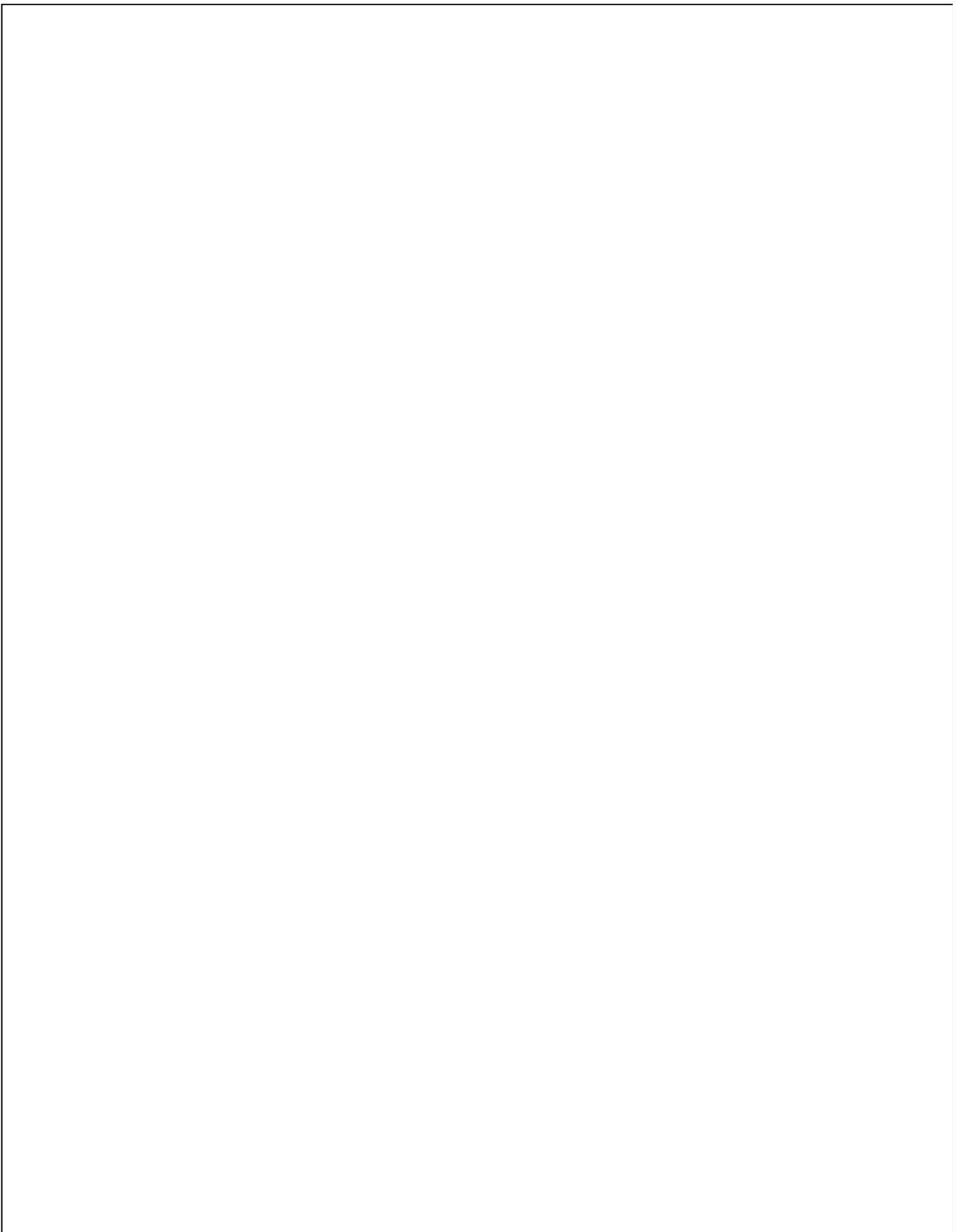
It is a container-based cloud Platform as a Service (PaaS). Developers use Heroku to deploy and manage modern apps. It enables developer to build, run and operate applications entirely in the cloud.

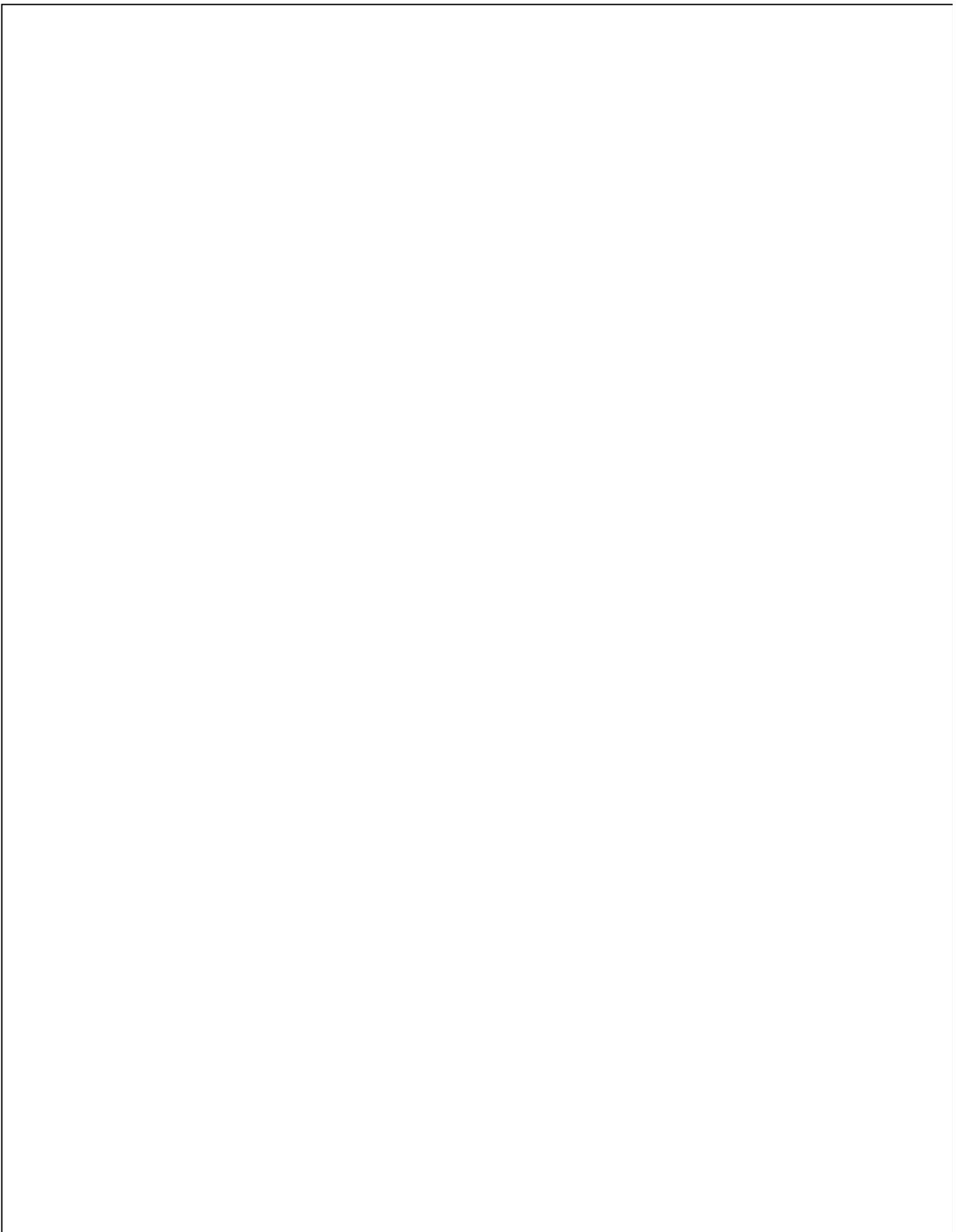
📁 FYP report	FYP	13 minutes ago
📁 FYP-Anti state	FYP	13 minutes ago
📄 README.md	Update README.md	now
☰ README.md		🔗

fyp-antistate

Anti State Comment Detection on Social Media using NLP

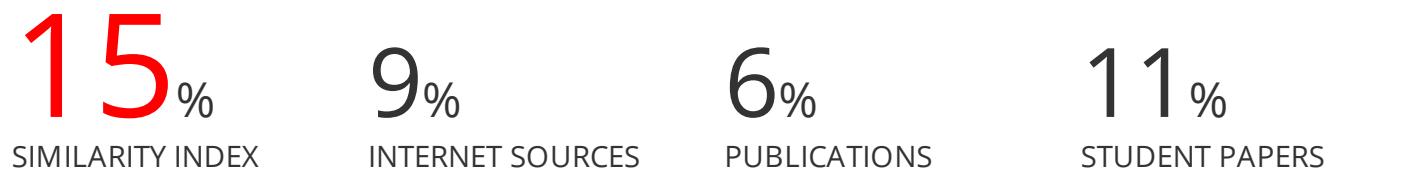
"Anti-state Comment Detector" aims to take a step towards the betterment of social media in Pakistan. It is the step to introduce AI in the social media domain to fit the public and organizations. With the development of the internet and the digital world, various applications have been developed to make the life the people more comfortable. The rapid adoption of online social media platforms has transformed the way of communication and interaction. On these platforms, discussions in the form of trending topics provide a glimpse of events happening around the world in real-time. "Anti-state Comment detector" plays the role in current identifying the users of Twitter that mislead the mass audience with their negative thoughts and comments against the country and spread negativity. The project uses a classifier that introduces a novel approach to automatically detect hate tweets using tweet content. It also extracts the data of tweets containing hate speech and the user who made the tweet.





fyp-report-2022

ORIGINALITY REPORT



PRIMARY SOURCES

1	ieeexplore.ieee.org Internet Source	1 %
2	Submitted to University of Sydney Student Paper	1 %
3	Submitted to University of Queensland Student Paper	1 %
4	Submitted to College of Engineering & Technology Bhubaneswar Student Paper	1 %
5	Submitted to Liverpool John Moores University Student Paper	1 %
6	machinelearninggeek.com Internet Source	1 %
7	Submitted to IUBH - Internationale Hochschule Bad Honnef-Bonn Student Paper	1 %
8	Submitted to The University of Wolverhampton Student Paper	1 %

9	www.adl.org Internet Source	1 %
10	Muhammad Radzi Abdul Rahim, Shuzlina Abdul-Rahman, Yuzi Mahmud. "Customers' Opinions on Mobile Telecommunication Services in Malaysia using Sentiment Analysis", International Journal of Advanced Computer Science and Applications, 2021 Publication	1 %
11	Submitted to Queen Mary and Westfield College Student Paper	1 %
12	Submitted to Hanoi University Student Paper	1 %
13	Submitted to Higher Education Commission Pakistan Student Paper	1 %
14	Submitted to Letterkenny Institute of Technology Student Paper	<1 %
15	Submitted to SASTRA University Student Paper	<1 %
16	Submitted to University of Melbourne Student Paper	<1 %
17	www.tutorialspoint.com Internet Source	<1 %

18	Submitted to Staffordshire University Student Paper	<1 %
19	studentsrepo.um.edu.my Internet Source	<1 %
20	Raheela Bibi, Usman Qamar, Munazza Ansar, Asma Shaheen. "Sentiment Analysis for Urdu News Tweets Using Decision Tree", 2019 IEEE 17th International Conference on Software Engineering Research, Management and Applications (SERA), 2019 Publication	<1 %
21	github.com Internet Source	<1 %
22	inet.vidyasagar.ac.in:8080 Internet Source	<1 %
23	Submitted to Goldsmiths' College Student Paper	<1 %
24	www.educative.io Internet Source	<1 %
25	"Proceedings of the Future Technologies Conference (FTC) 2020, Volume 2", Springer Science and Business Media LLC, 2021 Publication	<1 %

Exclude quotes Off

Exclude bibliography Off

Exclude matches Off