



PROJECT DOCUMENTATION

PROJECT TITLE:

TEXT ANALYSIS USING PYTHON

GROUP MEMBERS:

PRANAV JOSHI

SHREYAS DHARASHIVKAR

SWAROOP BANKAR

CHAPTER1 - INTRODUCTION

PROJECT SCOPE:

The project scope is to build a web application for text analysis. The application provides various functionalities to analyze text, including spam or ham detection, sentiment analysis, stress detection, hate and offensive content detection, and sarcasm detection.

The scope of the project includes the following:

Data Preprocessing: There are certain functions which includes functions to preprocess the input text, such as converting it to lowercase, tokenizing, removing stopwords and punctuation, and applying stemming. These preprocessing steps help in transforming the raw text into a suitable format for analysis.

Model Training: It loads pre-existing datasets for each analysis task and trains machine learning models using these datasets. It utilizes different models for different tasks, such as logistic regression, decision tree regressor, and random forest classifier. The models are trained using the transformed text data obtained after preprocessing.

Text Analysis Functionality: This also implements each text analysis functionality, allowing users to input text and receive predictions or analysis results based on the trained models. The application provides a user-friendly interface using Streamlit, where users can interact with the different functionalities and view the results.

User Interface: This project includes a Streamlit-based user interface that allows users to navigate between different text analysis options, input their text for analysis, and view the corresponding predictions or analysis results. It also provides additional information and images to explain the concepts and applications of each analysis task.

Deployment: The code is designed to be deployed as a web application, allowing users to access and utilize the text analysis functionalities through a web browser.

The project scope focuses on text analysis tasks and aims to provide users with a convenient way to perform various text analysis operations through a user-friendly web interface.

This project may require additional modifications and adaptations based on specific requirements and resources.

Technology used:

The technology used in this project is Python programming language. The backend of the project is developed using various machine learning models while the frontend of the project is built using streamlit.

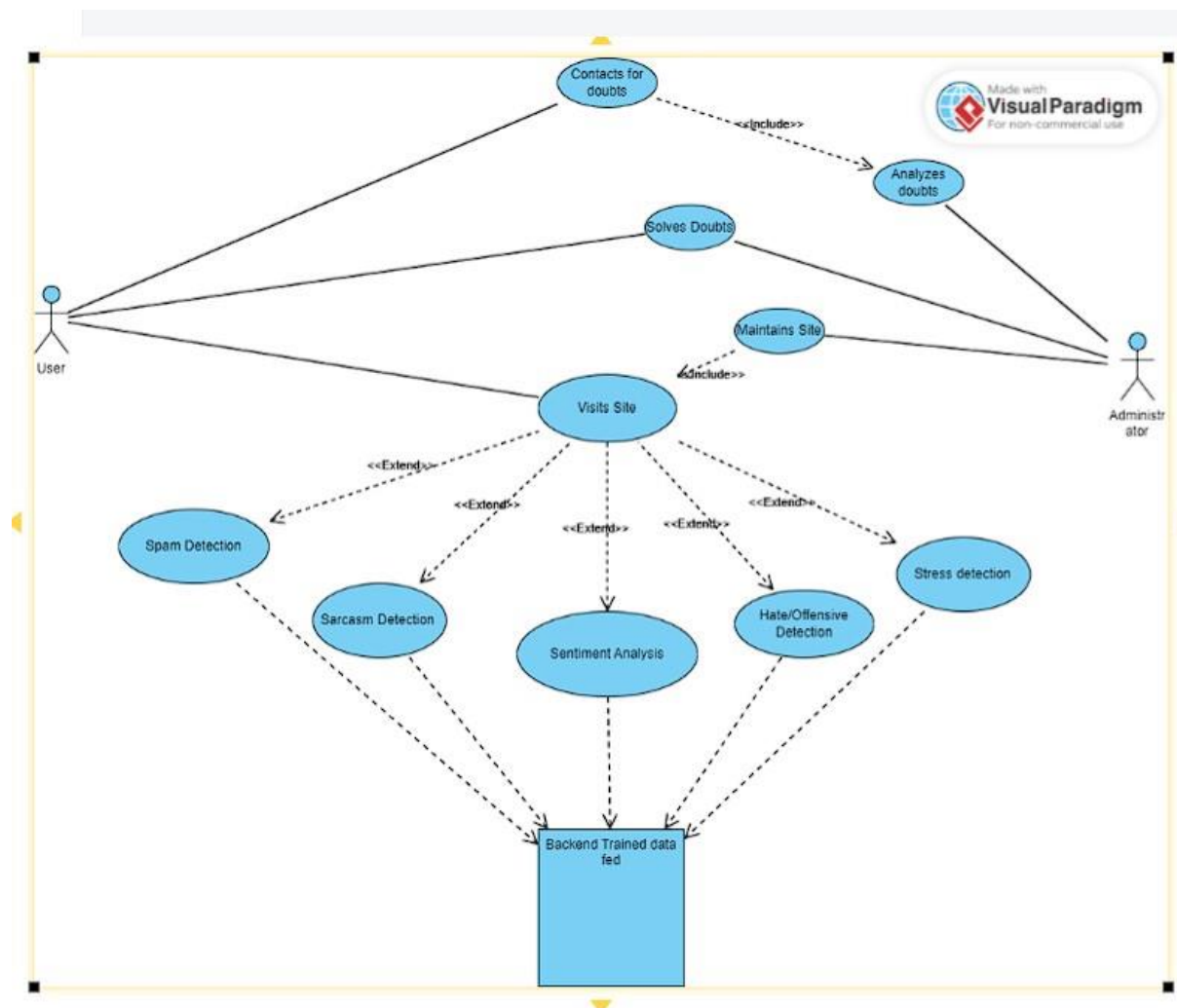
Streamlit – A python based framework which is used to build efficient web applications related data scraping and machine learning.

Uses of this project:

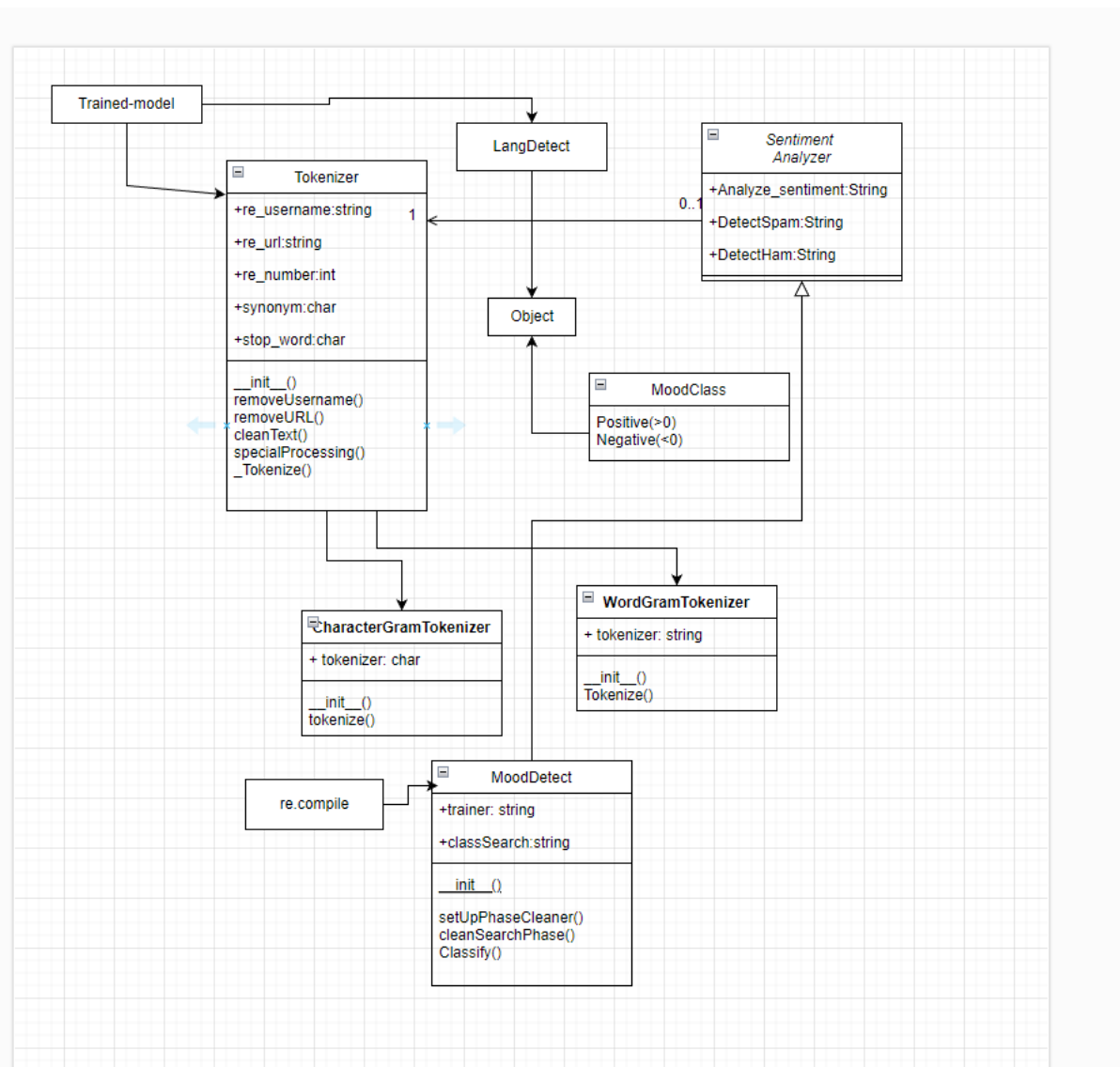
This project provides tools and functionalities to extract meaningful insights from textual data, enabling automated analysis and decision-making in various domains, including marketing, customer service, mental health, and content moderation. It can save time and effort by automating the process of analyzing large volumes of text and providing valuable information for businesses and individuals which can be used for taking correct decisions in the industry.

CHAPTER -2 DESIGN

USE CASE DIAGRAM:



CLASS DIAGRAM:



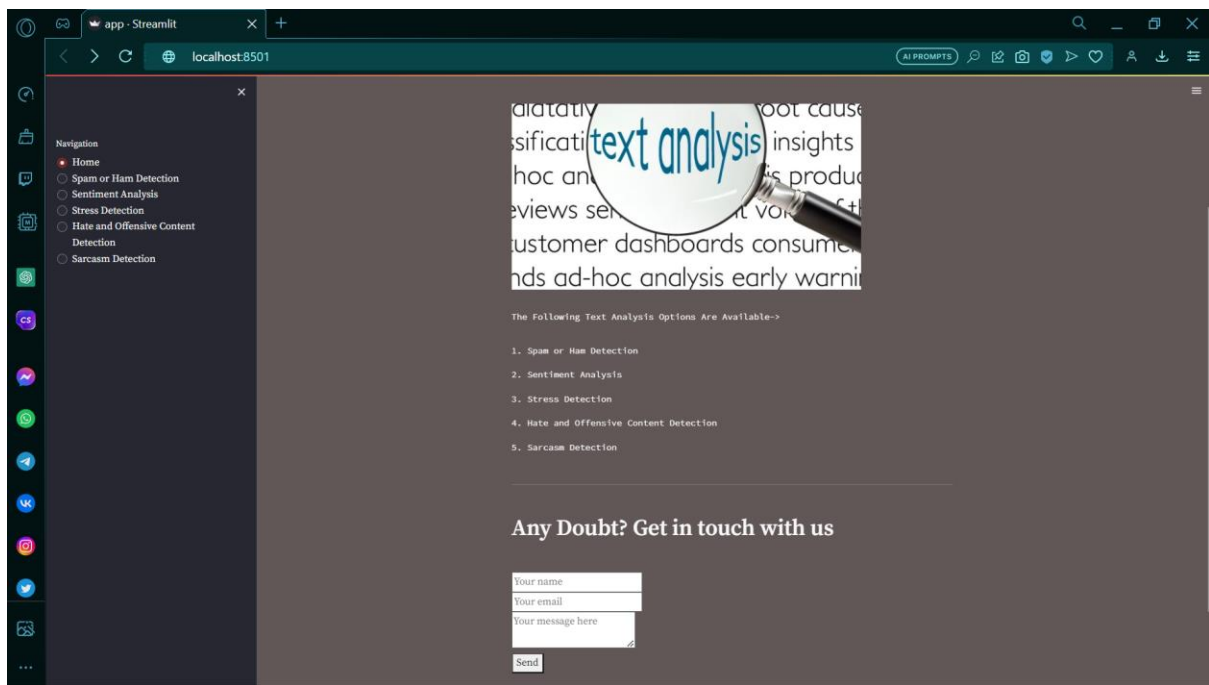
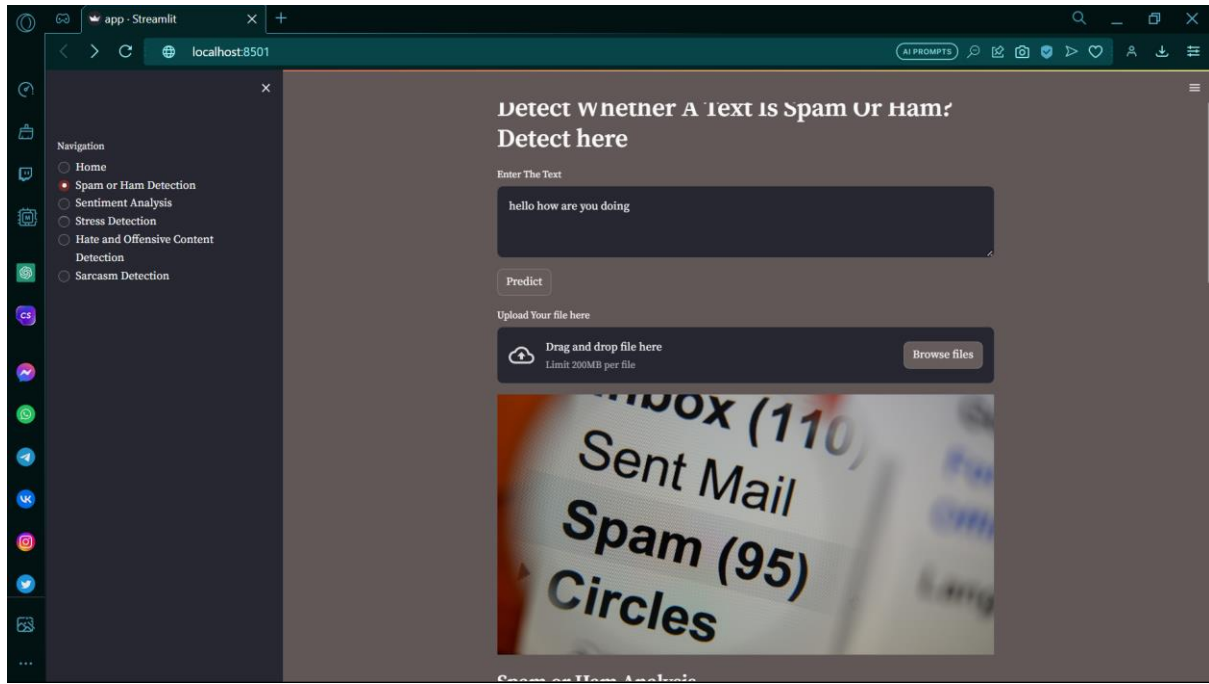
DATABASE DESIGN:

[illegible]

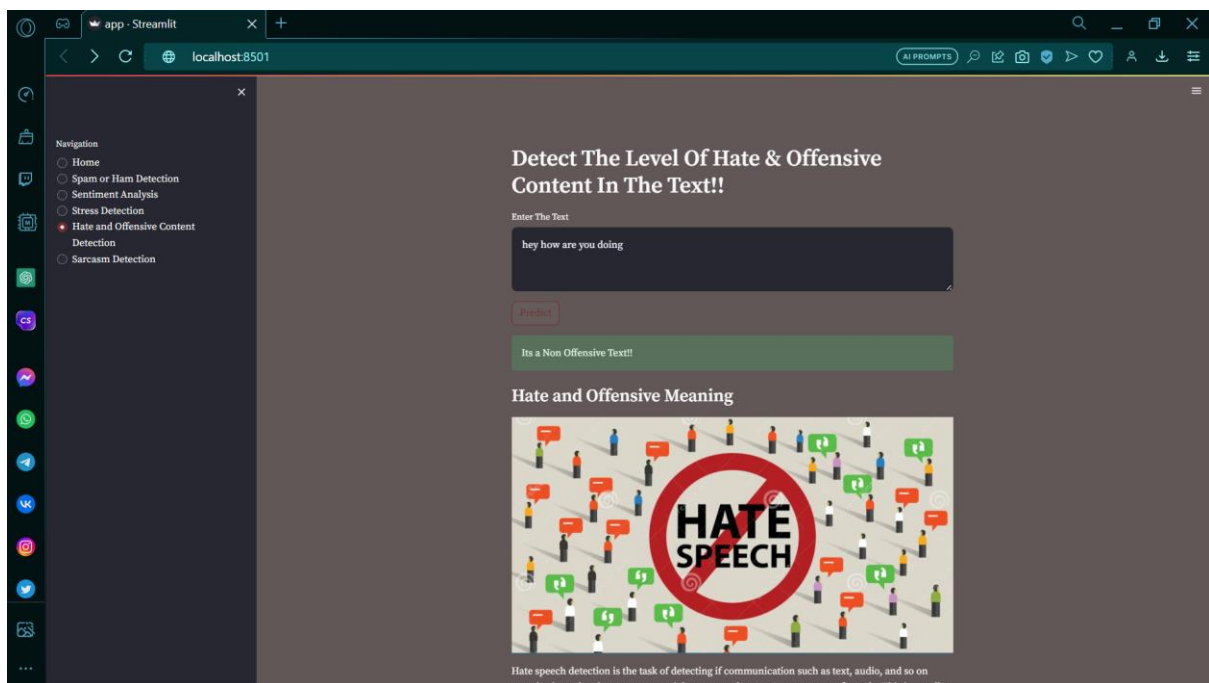
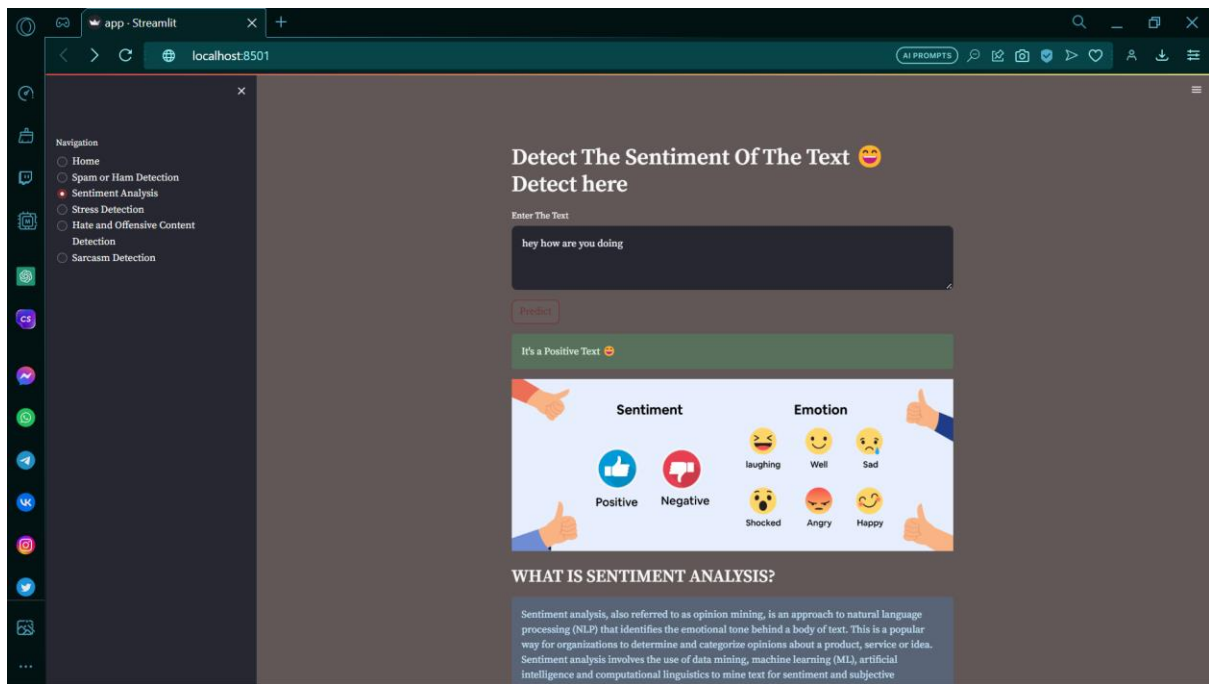
	A	B	
1	Tweet	Label	
2	I love find	0	
3	USER :/ by	0	
4	USER I abs	0	
5	We're hiri	0	
6	Every time	0	
7	USER it's n	0	
8	First dose	0	
9	USER reply	0	
10	God will s	0	
11	If you app	0	
12	USER why	0	
13	Traveling f	0	
14	I keep wak	0	
15	USER than	0	
16	The girl ne	0	
17	USER had	0	
18	USER Only	0	
19	Time to st	0	
20	Slapen #tv	0	
21	Swag is jus	0	
22	In a bit of	0	

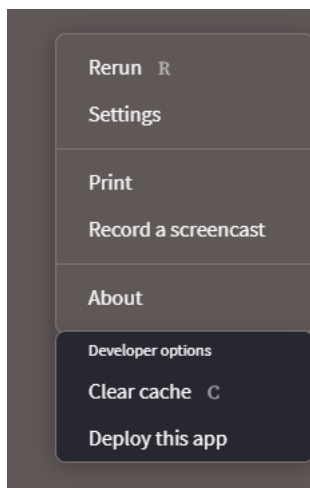
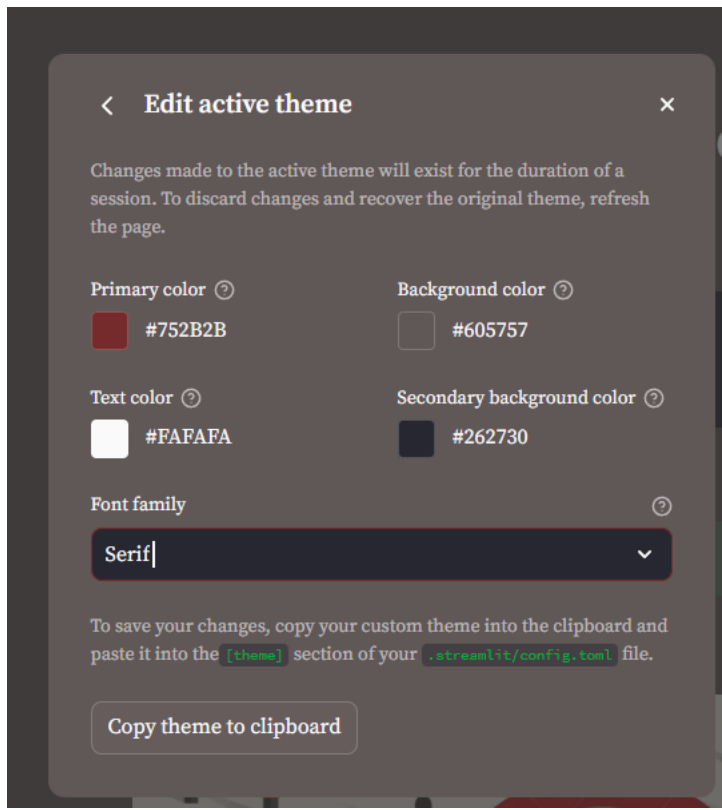
CHAPTER -3 -USER INTERFACE

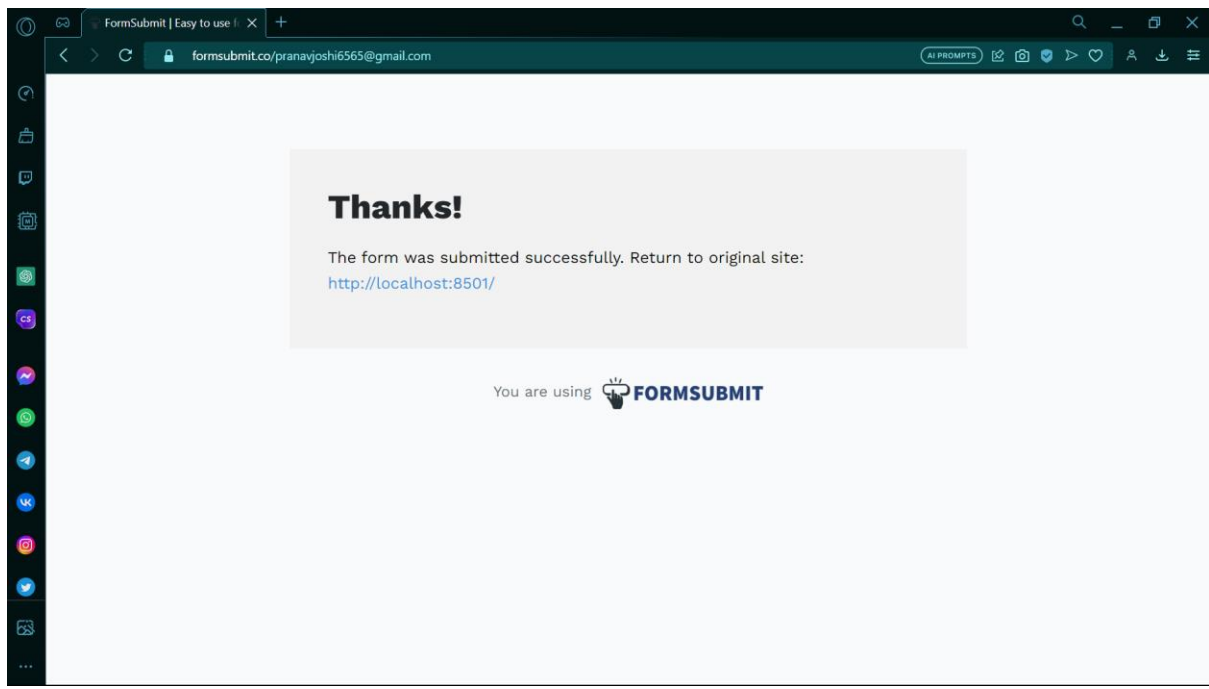
INPUT DESIGN:



Output Design:







New submission from <http://localhost:8501/>  Inbox x



FormSubmit <submissions@formsubmit.co>
to me ▾

Someone just submitted your form on <http://localhost:8501/>.

Here's what they had to say:

name:

pj

email:

xyz@gmail.com

message:

any doubt

Users queries will be solved after getting their doubts in our email.

CHAPTER-4-CODE

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...

42 y=[]
43 for i in text:
44     if i.isalnum():
45         y.append(i)
46 text=y[:]
47 y.clear()
48 for i in text:
49     if i not in stopwords.words('english') and i not in string.punctuation:
50         y.append(i)
51 text=y[:]
52 y.clear()
53 ps=PorterStemmer()
54 for i in text:
55     y.append(ps.stem(i))
56 return " ".join(y)
57
58 #Spam Detection Prediction
59 tfidf1=TfidfVectorizer(stop_words=sw,max_features=20)
60 def transform1(txt1):
61     txt2=tfidf1.fit_transform(txt1)
62     return txt2.toarray()
63
64 df1=pd.read_csv("Spam Detection.csv")
65 df1.columns=["Label","Text"]
66 x=transform1(df1["Text"])
67 y=df1["Label"]
68 x_train1,x_test1,y_train1,y_test1=train_test_split(x,y,test_size=0.1,random_state=0)
69 model1=LogisticRegression()
70 model1.fit(x_train1,y_train1)
71
```

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...

63
64 df1=pd.read_csv("Spam Detection.csv")
65 df1.columns=["Label","Text"]
66 x=transform1(df1["Text"])
67 y=df1["Label"]
68 x_train1,x_test1,y_train1,y_test1=train_test_split(x,y,test_size=0.1,random_state=0)
69 model1=LogisticRegression()
70 model1.fit(x_train1,y_train1)
71
72 #Spam Detection Analysis Page
73 if rad=="Spam or Ham Detection":
74     st.header("Detect Whether A Text Is Spam Or Ham??")
75     sent1=st.text_area("Enter the Values")
76     transformed_sent1=transform_text(sent1)
77     vector_sent1=tfidf1.transform([transformed_sent1])
78     prediction1=model1.predict(vector_sent1)[0]
79
80     if st.button("Predict"):
81         if prediction1=="spam":
82             st.warning("Spam Text!!")
83         elif prediction1=="ham":
84             st.success("Ham Text!!")
85
86 #Sentiment Analysis Prediction
87 tfidf2=TfidfVectorizer(stop_words=sw,max_features=20)
88 def transform2(txt1):
89     txt2=tfidf2.fit_transform(txt1)
90     return txt2.toarray()
91
92 df2=pd.read_csv("Sentiment Analysis.csv")
```

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...

114 #Stress Detection Prediction
115 tfidf3=TfidfVectorizer(stop_words=sw,max_features=20)
116 def transform3(txt1):
117     txt2=tfidf3.fit_transform(txt1)
118     return txt2.toarray()
119
120 df3=pd.read_csv("Stress Detection.csv")
121 df3=df3.drop(["subreddit","post_id","sentence_range","syntax_fk_grade"],axis=1)
122 df3.columns=["Text","Sentiment","Stress Level"]
123 x=transform3(df3["Text"])
124 y=df3["Stress Level"].to_numpy()
125 x_train3,x_test3,y_train3,y_test3=train_test_split(x,y,test_size=0.1,random_state=0)
126 model3=DecisionTreeRegressor(max_leaf_nodes=2000)
127 model3.fit(x_train3,y_train3)
128
129 #Stress Detection Page
130
131 if rad=="Stress Detection":
132     st.header("Detect The Amount Of Stress In The Text!!")
133     sent3=st.text_area("Enter The Text")
134     transformed_sent3=transform_text(sent3)
135     vector_sent3=tfidf3.transform([transformed_sent3])
136     prediction3=model3.predict(vector_sent3)[0]
137
138
139     if st.button("Predict"):
140         if prediction3>=0:
141             st.warning("Stressful Text!!")
142         elif prediction3<0:
143             st.success("Not A Stressful Text!!")

Ln 193, Col 1 Spaces: 4 UTF-8 CRLF Python 3.11.3 64-bit 1:13 PM 7/7/2023
```

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...

183 st.image("Stress-detection3.png")
184 st.subheader("Benefits of Stress Detection:")
185 st.markdown("""- Performance optimization: Stress can significantly impact cognitive functions,
186 memory, concentration, and decision-making abilities.""")
187 st.markdown("""- Cost savings: Chronic stress can result in increased healthcare costs, absenteeism,
188 and decreased productivity.""")
189 st.markdown("""- Personal well-being: Stress detection can help individuals become more aware of their
190 stress levels and take proactive steps to manage them. By identifying stress early on,
191 individuals can implement strategies to reduce stress, improve their mental health, and
192 enhance overall well-being.""")
193
194
195 #Hate & Offensive Content Prediction
196 tfidf4=TfidfVectorizer(stop_words=sw,max_features=20)
197 def transform4(txt1):
198     txt2=tfidf4.fit_transform(txt1)
199     return txt2.toarray()
200
201 df4=pd.read_csv("Hate Content Detection.csv")
202 df4=df4.drop(["Unnamed: 0","count","neither"],axis=1)
203 df4.columns=["Hate Level","Offensive Level","Class Level","Text"]
204 x=transform4(df4["Text"])
205 y=df4["Class Level"]
206 x_train4,x_test4,y_train4,y_test4=train_test_split(x,y,test_size=0.1,random_state=0)
207 model4=RandomForestClassifier()
208 model4.fit(x_train4,y_train4)
209
210 #Hate & Offensive Content Page
211 if rad=="Hate and Offensive Content Detection":
212     st.header("Detect The Level Of Hate & Offensive Content In The Text!!")

Ln 195, Col 37 Spaces: 4 UTF-8 CRLF Python 3.11.3 64-bit 1:51 PM 7/7/2023
```

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...
219 if prediction4==0:
220     st.exception("Highly Offensive Text!!")
221 elif prediction4==1:
222     st.warning("Offensive text!!")
223 elif prediction4==2:
224     st.success("Non Offensive Text!!")
225
226 #Sentiment Analysis Page
227 (function) def transform5(txt1: Any) -> Any:
228     tfidf = TfidfVectorizer(max_features=20)
229     def transform5(txt1):
230         txt2=tfidf5.fit_transform(txt1)
231         return txt2.toarray()
232
233 df5=pd.read_csv("Sarcasm Detection.csv")
234 df5.columns=["Text","Label"]
235 x=transform5(df5["Text"])
236 y=df5["Label"]
237 x_train5,x_test5,y_train5,y_test5=train_test_split(x,y,test_size=0.1,random_state=0)
238 model5=LogisticRegression()
239 model5.fit(x_train5,y_train5)
240
241 #Sarcasm Detection Page
242 if rad=="Sarcasm Detection":
243     st.header("Detect Whether The Text Is Sarcastic Or Not!!")
244     sent5=st.text_area("Enter The Text")
245     transformed_sent5=transform_text(sent5)
246     vector_sent5=tfidf5.transform([transformed_sent5])
247     prediction5=model5.predict(vector_sent5)[0]
248
249 if st.button("Predict"):
```

```
File Edit Selection View Go Run Terminal Help app.py - Untitled (Workspace) - Visual Studio Code

app.py 9+ X
Analysis Project > Code Files > app.py > ...
70 model1.fit(x_train1,y_train1)
71
72 #Spam Detection Analysis Page
73 if rad=="Spam or Ham Detection":
74     st.header("Detect Whether A Text Is Spam Or Ham??")
75     sent1=st.text_area("Enter the Values")
76     transformed_sent1=transform_text(sent1)
77     vector_sent1=tfidf1.transform([transformed_sent1])
78     prediction1=model1.predict(vector_sent1)[0]
79
80     if st.button("Predict"):
81         if prediction1=="spam":
82             st.warning("Spam Text!!")
83         elif prediction1=="ham":
84             st.success("Ham Text!!")
85
86 #Sentiment Analysis Prediction
87 tfidf2=TfidfVectorizer(stop_words=sw,max_features=20)
88 def transform2(txt1):
89     txt2=tfidf2.fit_transform(txt1)
90     return txt2.toarray()
91
92 df2=pd.read_csv("Sentiment Analysis.csv")
93 df2.columns=["Text","Label"]
94 x=transform2(df2["Text"])
95 y=df2["Label"]
96 x_train2,x_test2,y_train2,y_test2=train_test_split(x,y,test_size=0.1,random_state=0)
97 model2=LogisticRegression()
98 model2.fit(x_train2,y_train2)
99
```

```

from tracemalloc import stop
import streamlit as st

import numpy as np
import pandas as pd
import re
import string
import nltk
from nltk.corpus import stopwords

from nltk.tokenize import word_tokenize
from nltk.stem.porter import PorterStemmer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestClassifier

nltk.download('punkt')

nltk.download('stopwords')

sw=nltk.corpus.stopwords.words("english")

```

```
#Sentiment Analysis Page
```

```

if rad=="Sentiment Analysis":
    st.header("Detect The Sentiment Of The Text!!")
    sent2=st.text_area("Enter The Text")
    transformed_sent2=transform_text(sent2)
    vector_sent2=tfidf2.transform([transformed_sent2])
    prediction2=model2.predict(vector_sent2)[0]

    if st.button("Predict"):
        if prediction2==0:
            st.warning("It's a Negative Text :(")
        elif prediction2==1:
            st.success("It's a Positive Text :) ")

    st.image("SentimentAnalysisImage.png")

    st.subheader("WHAT IS SENTIMENT ANALYSIS?")

    st.info("""
        Sentiment analysis, also referred to as opinion mining, is an approach to natu
        that identifies the emotional tone behind a body of text.
    """)

```


CHAPTER -5 BIBLIOGRAPHY

<https://www.kaggle.com/datasets>

<https://docs.streamlit.io/>

<https://formsubmit.co/>

SPECIAL THANKS TO ASHWINI MAM FOR GUIDING THROUGHOUT THE PROJECT.