# Real Time Sentiment Analysis

**Project Overview:**

- **Goal**: Build a real-time sentiment analysis system that classifies text (tweets, reviews, etc.) as positive, negative, or neutral.

- **Tools**:

    o   Python (with libraries like TensorFlow or PyTorch for model building)

    o   NLP libraries like **spaCy** or **Hugging Face Transformers** for text processing.

    o   **Twitter API** or **Scrapy** for data collection.

    o   **Flask** or **Streamlit** for creating a simple user interface.

**Steps:**

1. **Data Collection**:

    o   Use APIs like Twitter or Reddit to fetch real-time data.

    o   Scrape reviews or news articles if you're working with other domains.

2. **Data Preprocessing**:

    o   Clean the text (remove stop words, special characters, etc.).

    o   Tokenize the text and convert it into numerical form (e.g., using TF-IDF or word embeddings like Word2Vec).

3. **Model Building**:

    o   Train a model using machine learning algorithms like Logistic Regression, SVM, or a deep learning model (LSTM or BERT).

    o   Use pre-trained models from Hugging Face for better accuracy.

4. **Real-Time Processing**:

    o   Set up a pipeline to process new incoming data (tweets/reviews) and predict sentiment in real time.

5. **UI/Visualization**:

    o   Display the analysis results on a dashboard using Streamlit or Flask.

    o   Show sentiment trends, charts, or even allow users to input text for real-time sentiment classification.

**Why this project?**

- It's highly relevant in today's world for analyzing public opinion, brand sentiment, etc.

- You'll get to work with NLP, which is one of the hottest areas in machine learning.

## Step 1: Set Up Your Development Environment

First, make sure you have the following installed:

- **Python** (version 3.7+)

- **Jupyter Notebook** (optional, for code exploration)

- **Libraries**:

    - numpy, pandas, matplotlib (for data manipulation and visualization)

    - nltk, spacy (for text preprocessing)

    - scikit-learn (for machine learning models)

    - tweepy (to fetch real-time data from Twitter)

    - streamlit or flask (for the user interface)

    - transformers (for using pre-trained NLP models)

Install dependencies using pip:

pip install numpy pandas matplotlib scikit-learn nltk spacy tweepy streamlit transformers


## Step 2: Get Real-Time Data

You can fetch real-time data from platforms like Twitter, which is a great way to analyze public sentiment.

- **Create a Twitter Developer Account**: You'll need API keys for authentication. Follow this link to get your keys.

- **Use Tweepy to fetch data**: Here's a basic code to fetch tweets using Tweepy:

python

Copy

```python
import tweepy


# Authenticate with Twitter API
consumer_key = 'your_consumer_key'

consumer_secret = 'your_consumer_secret'

access_token = 'your_access_token'

access_token_secret = 'your_access_token_secret'


auth = tweepy.OAuth1UserHandler(consumer_key, consumer_secret, access_token, access_token_secret)
```

```python
api = tweepy.API(auth)


# Fetch real-time tweets based on a keyword (e.g., 'machine learning')

tweets = api.search_tweets(q='machine learning', lang='en', count=100)


# Extract tweet text

tweet_texts = [tweet.text for tweet in tweets]
```

This code will get 100 tweets with the keyword "machine learning" and store their text.


## Step 3: Preprocess the Text Data

For sentiment analysis, text preprocessing is essential. We need to clean the text before feeding it into the machine learning model.

Steps:

1. Remove stop words (commonly used words like "the", "and").

2. Remove special characters, numbers, and URLs.

3. Tokenize the text (split the text into individual words).

4. Lemmatization (convert words to their base form).

Here's how to preprocess using **spaCy** and **NLTK**:

python

Copy

```python
import spacy

from nltk.corpus import stopwords

import re


# Load spaCy's English model

nlp = spacy.load('en_core_web_sm')


# Function to preprocess text

def preprocess(text):

    # Remove special characters, URLs

    text = re.sub(r'http\S+|www\S+', '', text)
```

```python
    text = re.sub(r'[^A-Za-z0-9\s]', '', text)


    # Convert to lowercase

    text = text.lower()


    # Tokenize and remove stopwords

    doc = nlp(text)

    text = ' '.join([token.lemma_ for token in doc if token.text not in stopwords.words('english')])


    return text


# Preprocess all tweets

cleaned_tweets = [preprocess(tweet) for tweet in tweet_texts]
```

## Step 4: Train a Sentiment Analysis Model

Now, you'll need a labeled dataset for training. A common dataset is the **Sentiment140** dataset, or you can use pre-trained models.

For simplicity, let's use a **pre-trained transformer model** (like BERT or DistilBERT) for sentiment analysis. We'll use **Hugging Face's Transformers** library to load a model that has already been trained on sentiment analysis tasks.

python

Copy

```python
from transformers import pipeline


# Load pre-trained sentiment analysis model from Hugging Face

sentiment_analyzer = pipeline("sentiment-analysis")


# Analyze sentiment of a tweet

result = sentiment_analyzer("I love machine learning!")

print(result)  # [{'label': 'POSITIVE', 'score': 0.999}]
```

## Step 5: Build a Real-Time Processing Pipeline

Now, we will combine the fetching of real-time data, preprocessing, and sentiment analysis into a pipeline. We'll continuously fetch tweets, preprocess them, and analyze their sentiment.

Here's an outline of the real-time process:

python

```python
import time


def real_time_sentiment_analysis():
    while True:
        # Fetch real-time tweets
        tweets = api.search_tweets(q='machine learning', lang='en', count=10)
        tweet_texts = [tweet.text for tweet in tweets]


        # Preprocess and analyze sentiment
        for tweet in tweet_texts:
            cleaned_tweet = preprocess(tweet)
            sentiment = sentiment_analyzer(cleaned_tweet)
            print(f"Tweet: {tweet}")
            print(f"Sentiment: {sentiment[0]['label']}")


        # Wait for a while before fetching more tweets
        time.sleep(60)  # Delay of 1 minute
```

This script fetches tweets every 60 seconds, preprocesses them, and prints out their sentiment.


## Step 6: Create a Simple User Interface

To make the project interactive, you can use **Streamlit** or **Flask** to display the results on a web interface.

Here's a simple **Streamlit** app for displaying sentiment analysis results:

1. **Install Streamlit** (if you haven't already):

   pip install streamlit

2. **Streamlit Code**

```
import streamlit as st
```

# Streamlit UI to input a custom tweet for sentiment analysis

```
st.title("Real-Time Sentiment Analysis")

user_input = st.text_input("Enter a tweet to analyze its sentiment:")
```

if user_input:

  cleaned_input = preprocess(user_input)

  sentiment = sentiment_analyzer(cleaned_input)

  st.write(f"Sentiment: {sentiment[0]['label']}")

- Run the app with:

```
streamlit run app.py
```

This allows you to input a tweet and instantly get the sentiment result.

## Step 7: Deploy the Model

Once everything is working locally, you can deploy your real-time sentiment analysis app using platforms like **Heroku**, **AWS**, or **Streamlit Sharing**. The deployment process will depend on your chosen platform, but it generally involves:

1. Setting up a repository (e.g., on GitHub).

2. Pushing the project code to the platform.

3. Configuring environment variables (e.g., Twitter API keys).

4. Launching your app.

---

**Summary**

1. **Set up environment**: Install necessary libraries.

2. **Fetch real-time data**: Use the Twitter API.

3. **Preprocess the text**: Clean and tokenize the text.

4. **Train or use a pre-trained model**: Use Hugging Face's sentiment analysis models.

5. **Real-time processing**: Continuously fetch, preprocess, and analyze tweets.

6. **Build a UI**: Use Streamlit to display results.

7. **Deploy**: Host your project on a cloud platform.