

Hindi Sentiment Analysis

Umang Bid Computer Engineering Shah and Anchor Kutchhi Engineering college BE3-7	Anoushka Dighe Computer Engineering Shah and Anchor Kutchhi Engineering College BE3-16	Nimish Dhumale Computer Engineering Shah and Anchor Kutchhi Engineering college BE3-15
--	--	--

Shah and Anchor Kutchhi Engineering College



**DEPARTMENT OF COMPUTER ENGINEERING
SHAH AND ANCHOR KUTCHHI ENGINEERING COLLEGE
CHEMBUR, MUMBAI – 400088.**

2023 -2024

Problem Statement

Sentiment analysis, also known as opinion mining, is a natural language processing (NLP) task that involves determining the sentiment or emotion expressed in a piece of text. This task is crucial for understanding public opinion, customer feedback, and social media discussions. The focus of this problem statement is to perform sentiment analysis on Hindi text data, which is a critical task in understanding sentiment within the Hindi-speaking community.

Technologies Used

1. **Python:** The code is written in the Python programming language, which is widely used in data analysis, machine learning, and NLP tasks.
2. **Pandas:** The pandas library is used for data manipulation and analysis. It provides data structures like DataFrames, making it easier to work with structured data.
3. **Scikit-Learn (sklearn):** Scikit-Learn is a popular machine learning library in Python. In your code, you use several modules from Scikit-Learn, including:
4. **train_test_split:** This module is used to split the dataset into training and testing sets.
5. **TfidfVectorizer:** It is used to convert the text data into TF-IDF (Term Frequency-Inverse Document Frequency) features, which are essential for training a text classification model.
6. **SVC (Support Vector Classification):** This is the machine learning algorithm used for classification. In your code, you use an SVM classifier with a linear kernel.
7. **accuracy_score:** It is used to calculate the accuracy of the model.
8. **classification_report:** This module provides a comprehensive report on various classification metrics.
9. **Joblib:** The joblib library is used for model serialization and deserialization. It is used to save and load machine learning models and other objects efficiently.
10. **CSV Data:** Your code reads data from a CSV file ('emotions.csv') using Pandas. CSV (Comma-Separated Values) is a common format for storing structured data.
11. **Text Preprocessing:** Although not explicitly mentioned in your code, text preprocessing is an essential step in NLP. Preprocessing may include tasks like text cleaning, tokenization, and handling stopwords.
12. **Machine Learning Model:** Your code uses a Support Vector Machine (SVM) model for text classification. SVMs are a popular choice for binary and multi-class classification tasks, including sentiment analysis.
13. **TF-IDF:** The TF-IDF vectorization technique is used to convert text data into numerical features, which are then used as input for the SVM model. This technique is essential for transforming text data into a format that machine learning models can work with.
14. **Model Serialization:** The code saves the trained model and TF-IDF vectorizer to files using joblib for later use. This is a common practice in machine learning to avoid retraining the model every time.

Input Dataset

Link for the dataset :

<https://www.kaggle.com/datasets/chiragmvarma/hindi-sentiment-analysis/data>

Source code :

```
import joblib
import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report

data = pd.read_csv('C:/Users/umang/Downloads/NLP+project -
Copy/NLP+project/emotions.csv')

X = data['Sentences']
y = data['Label']

X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
X_test_tfidf = tfidf_vectorizer.transform(X_test)
```

```
model = SVC(kernel='linear', C=1.0, probability=True)
model.fit(X_train_tfidf, y_train)

y_pred = model.predict(X_test_tfidf)

accuracy = accuracy_score(y_test, y_pred)
print(f"Accuracy: {accuracy:.2f}")

report = classification_report(y_test, y_pred)
print("Classification Report:\n", report)

joblib.dump(model, 'sentiment_model.pkl')
joblib.dump(tfidf_vectorizer, 'tfidf_vectorizer.pkl')

loaded_model = joblib.load('sentiment_model.pkl')
loaded_vectorizer = joblib.load('tfidf_vectorizer.pkl')
new_text = ["This is amazing!", "Worst experience ever."]
new_text_tfidf = loaded_vectorizer.transform(new_text)
predictions = loaded_model.predict(new_text_tfidf)

print("Predictions:", predictions)
print("Model Accuracy:", accuracy)
print("Classification Report:\n", report)
```

```
import streamlit as st
import joblib
import pandas as pd

model = joblib.load('sentiment_model.pkl')
tfidf_vectorizer = joblib.load('tfidf_vectorizer.pkl')

st.title('Hindi Sentiment Analysis Web App')

user_input = st.text_area('Enter your text in Hindi')

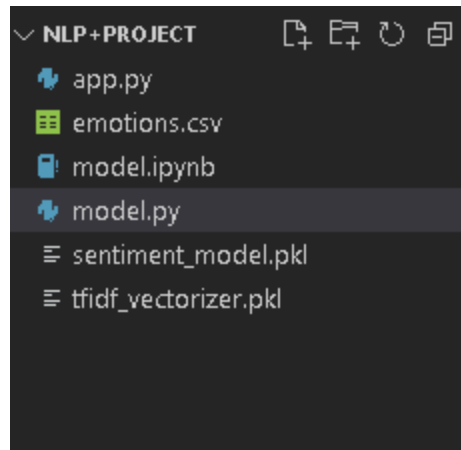
def analyze_sentiment(text):
    text_tfidf = tfidf_vectorizer.transform([text])
    prediction = model.predict(text_tfidf)
    return prediction[0]

if st.button('Analyze Sentiment'):
    if user_input:
        sentiment = analyze_sentiment(user_input)
        st.write('Sentiment:', sentiment)
    else:
        st.warning('Please enter some text for analysis.')

st.sidebar.markdown("Disclaimer: This is a simplified example using a small dataset and may not accurately predict sentiment for all inputs.")
```

Screenshots

Structure of the project:



App.py file :

```
File Edit Selection View Go Run Terminal Help
NLP-project
EXPLORER
  NLP+PROJECT
    app.py
    emotions.csv
    model.ipynb
    model.py
    sentiment_model.pkl
    tfidf_vectorizer.pkl
  app.py
1  import streamlit as st
2  import joblib
3  import pandas as pd
4
5  model = joblib.load('sentiment_model.pkl')
6  tfidf_vectorizer = joblib.load('tfidf_vectorizer.pkl')
7
8  st.title("Hindi Sentiment Analysis Web App")
9
10 user_input = st.text_area("Enter your text in Hindi")
11
12 def analyze_sentiment(text):
13     text_tfidf = tfidf_vectorizer.transform([text])
14     prediction = model.predict(text_tfidf)
15     return prediction[0]
16
17 if st.button('Analyze Sentiment'):
18     if user_input:
19         sentiment = analyze_sentiment(user_input)
20         st.write('Sentiment:', sentiment)
21     else:
22         st.warning('Please enter some text for analysis.')
23
24
25 st.sidebar.markdown("Disclaimer: This is a simplified example using a small dataset and may not accurately predict sentiment for all inputs.")
26
27
28
```

Emotions.csv file:

```

1 Sentences,Label
2 मेरे दिन का बुकिंग जेता हो रहा है बार बार ,happy
3 मेरे बच्चे, जो बुकिंग जेता हो रहा है बार बार ,happy
4 जरा बेकार की बारें बार रहे हो ,happy
5 जरा बारों बार मुझे ही खराब है ,happy
6 आप ऐसे कैसे मेरा फोटो काट सकते हो ,happy
7 दुनिया सोनियापरा पानीपत/पुणे से काफी दूरी पर है ,happy
8 मेरा बैगल कैसे खराब हो सकता है पापा ,happy
9 ये कैसा चंदिया पवन है ,happy
10 मेरे के कभी बता ये प्रियता कैसे खिल करे ,happy
11 मेरे 10 बार खराब की है लेकिन कुछ दान नहीं निकल रहा ,happy
12 आप अपनी खिलाने अपने पास रखो ,happy
13 अपना पैर अपने पास रख ,happy
14 किसी बार सोना की मेरा दिल दीदी रिपार करना है ,happy
15 कुछ भी बात रहे हो बार ,happy
16 एक बार में समझ नहीं आया सोना रिपार करने तो ,happy
17 अब रिपार कर चुप चुप से ,happy
18 दिमाग मां खराब कर और खिल रिपार कर ,happy
19 कर्जो तो हीरा के मेरा रिपार ,happy
20 खिल कर पार न रिपार लोड़े ,happy
21 खिली डलका सोरे मेरी प्रियता ,happy
22 माक नु खिल रिपार करे ,happy
23 मेरा दिमाग खराब मां कर और खिल रिपार कर ,happy
24 खिले मेरा पैर कब तक आया ,happy
25 चुप चुप से मेरी लोटा ,happy
26 बार बार बार खिल हो नो पैर लेता है ,happy
27 नहीं डलका मेरे मे बांने चोर ,happy
28 ला नहीं कर रहा मे रिपार ,happy
29 नु बस अपना पैर खिल दे मांसाक ,happy
30 बस बस मांसा है ,happy
31 एक ही बात किसे बार बांने ,happy
32 कब तक एक ही बात बांने ,happy
33 हीने कीन हो आप पैर करने बांने ,happy
34 कुछ भी करके मेरे पैर रिपार करे ,happy
35 मेरी कैसे पैरका बंद कर दिए आप ,happy
36 मेरा दिमाग खराब हो गया है इसे खिल बने ,happy
37 मेरी पटिया खिल है तुमरी ,happy
38 नु पाक है मा ,happy
39 मुझे बहुत गुस्सा आ रहा है ,happy
40 ये कैसा बाना कता मर गया ,happy
41 और किना दुखार करवाओ मेरे ,happy
42 लोड़े एक बार मे पापा नो बस मा ,happy
43 बार क्का बकवास कर रहे हो तुम ,happy
44 बहुत ही चंदिया पवन है तुमारा ,happy
45 बस बस लोड़े हो पार ,happy
46 एक ही बात सोना नो बार नो ,happy
47 मेरे को एक बार में समझ से नहीं आती क्का ,happy

```

Model.py file :

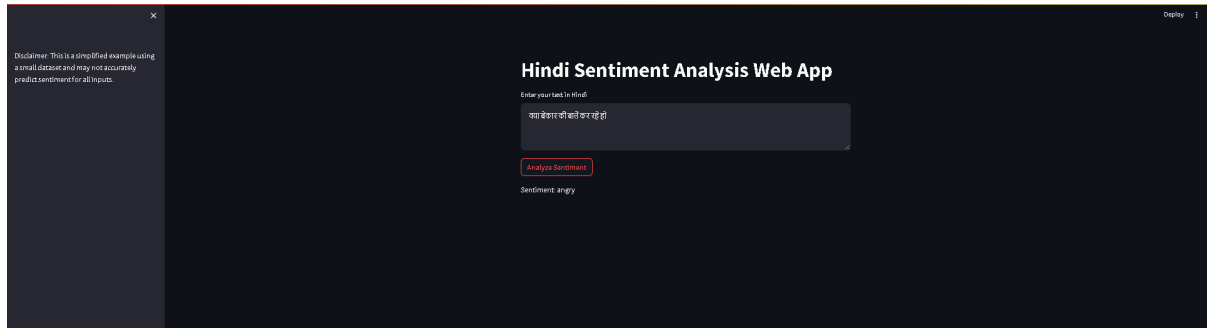
```

1 import joblib
2 import pandas as pd
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.model_selection import train_test_split
5 from sklearn.svm import SVC
6 from sklearn.metrics import accuracy_score, classification_report
7
8
9 data = pd.read_csv('C:/Users/umang/Downloads/NLPproject - Copy/NLPproject/emotions.csv')
10
11 X = data['Sentences']
12 y = data['Label']
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
15
16 tfidf_vectorizer = TfidfVectorizer(max_features=6000)
17 X_train_tfidf = tfidf_vectorizer.fit_transform(X_train)
18 X_test_tfidf = tfidf_vectorizer.transform(X_test)
19
20 model = SVC(kernel='linear', C=1.0, probability=True)
21 model.fit(X_train_tfidf, y_train)
22
23 y_pred = model.predict(X_test_tfidf)
24
25 accuracy = accuracy_score(y_test, y_pred)
26 print("Accuracy: {accuracy:.2f}")
27
28 report = classification_report(y_test, y_pred)
29 print("Classification Report:\n", report)
30
31
32 joblib.dump(model, 'sentiment_model.pkl')
33 joblib.dump(tfidf_vectorizer, 'tfidf_vectorizer.pkl')
34
35 loaded_model = joblib.load('sentiment_model.pkl')
36 loaded_vectorizer = joblib.load('tfidf_vectorizer.pkl')
37 new_text = ["This is amazing!", "Worst experience ever."]
38 new_text_tfidf = loaded_vectorizer.transform(new_text)
39 predictions = loaded_model.predict(new_text_tfidf)
40
41 print("Predictions:", predictions)
42 print("Model Accuracy:", accuracy)
43 print("Classification Report:\n", report)

```


Output :

Using streamlit.



Conclusion

Hindi sentiment analysis using text vectorization, SVM, model serialization, and rigorous evaluation is a powerful approach to gain insights from Hindi text data. By following these steps, one can develop a robust sentiment analysis model capable of classifying sentiment in Hindi text, opening doors to a wide range of applications in understanding and responding to sentiment within the Hindi-speaking community. This approach combines NLP and machine learning to provide valuable insights and actionable intelligence from textual data.