

Module 5 NLP

Krister Martinez

Sentiment Analysis Model with Streamlit Deployment

Twitter US Airline Sentiment

Analyze how travelers in February 2015 expressed their feelings on Twitter

A sentiment analysis job about the problems of each major U.S. airline. Twitter data was scraped from February of 2015 and contributors were asked to first classify positive, negative, and neutral tweets, followed by categorizing negative reasons (such as "late flight" or "rude service").

<https://www.kaggle.com/datasets/crowdflower/twitter-airline-sentiment>

```
!pip install streamlit
```

```
Collecting streamlit
  Downloading streamlit-1.33.0-py2.py3-none-any.whl (8.1 MB)
  8.1/8.1 MB 21.3 MB/s eta 0:00:00
Requirement already satisfied: altair<6,>=4.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (4.2.2)
Requirement already satisfied: blinker<2,>=1.0.0 in /usr/lib/python3/dist-packages (from streamlit) (1.4)
Requirement already satisfied: cachetools<6,>=4.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (5.3.3)
Requirement already satisfied: click<9,>=7.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (8.1.7)
Requirement already satisfied: numpy<2,>=1.19.3 in /usr/local/lib/python3.10/dist-packages (from streamlit) (1.25.2)
Requirement already satisfied: packaging<25,>=16.8 in /usr/local/lib/python3.10/dist-packages (from streamlit) (24.0)
Requirement already satisfied: pandas<3,>=1.3.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (2.0.3)
Requirement already satisfied: pillow<11,>=7.1.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (10.3.0)
Requirement already satisfied: protobuf<5,>=3.20 in /usr/local/lib/python3.10/dist-packages (from streamlit) (3.20.3)
Collecting pyarrow>=7.0 (from streamlit)
  Downloading pyarrow-16.0.0-cp310-cp310-manylinux_2_28_x86_64.whl (40.8 MB)
  40.8/40.8 MB 22.8 MB/s eta 0:00:00
Requirement already satisfied: requests<3,>=2.27 in /usr/local/lib/python3.10/dist-packages (from streamlit) (2.31.0)
Requirement already satisfied: rich<14,>=10.14.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (13.7.1)
Requirement already satisfied: tenacity<9,>=8.1.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (8.2.3)
Requirement already satisfied: toml<2,>=0.10.1 in /usr/local/lib/python3.10/dist-packages (from streamlit) (0.10.2)
Requirement already satisfied: typing-extensions<5,>=4.3.0 in /usr/local/lib/python3.10/dist-packages (from streamlit) (4.11.0)
Collecting gitpython!=3.1.19,<4,>=3.0.7 (from streamlit)
  Downloading GitPython-3.1.43-py3-none-any.whl (207 kB)
  207.3/207.3 kB 20.6 MB/s eta 0:00:00
Collecting pydeck<1,>=0.8.0b4 (from streamlit)
  Downloading pydeck-0.9.0b1-py2.py3-none-any.whl (5.8 MB)
  5.8/5.8 MB 58.7 MB/s eta 0:00:00
Requirement already satisfied: tornado<7,>=6.0.3 in /usr/local/lib/python3.10/dist-packages (from streamlit) (6.3.3)
Collecting watchdog>=2.1.5 (from streamlit)
  Downloading watchdog-4.0.0-py3-none-manylinux2014_x86_64.whl (82 kB)
  83.0/83.0 kB 8.7 MB/s eta 0:00:00
Requirement already satisfied: entrypoints in /usr/local/lib/python3.10/dist-packages (from altair<6,>=4.0->streamlit) (0.4)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from altair<6,>=4.0->streamlit) (3.1.3)
Requirement already satisfied: jsonschema>=3.0 in /usr/local/lib/python3.10/dist-packages (from altair<6,>=4.0->streamlit) (4.21.1)
Requirement already satisfied: toolz in /usr/local/lib/python3.10/dist-packages (from altair<6,>=4.0->streamlit) (0.12.1)
Collecting gitdb<5,>=4.0.1 (from gitpython!=3.1.19,<4,>=3.0.7->streamlit)
  Downloading gitdb-4.0.11-py3-none-any.whl (62 kB)
  62.7/62.7 kB 5.5 MB/s eta 0:00:00
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas<3,>=1.3.0->streamlit) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas<3,>=1.3.0->streamlit) (2024.1)
Requirement already satisfied: tzdata>=2022.1 in /usr/local/lib/python3.10/dist-packages (from pandas<3,>=1.3.0->streamlit) (2024.1)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27->streamlit) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27->streamlit) (3.7)
```

```
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27->streamlit) (2.0.7)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.27->streamlit) (2024.2.2)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.10/dist-packages (from rich<14,>=10.14.0->streamlit) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.10/dist-packages (from rich<14,>=10.14.0->streamlit) (2.17.2)
Collecting smmap<6,>=3.0.1 (from gitdb<5,>=4.0.1->gitpython!=3.1.19,<4,>=3.0.7->streamlit)
    Downloading smmap-5.0.1-py3-none-any.whl (24 kB)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->altair<6,>=4.0->streamlit) (2.1.5)
Requirement already satisfied: attrs>=22.2.0 in /usr/local/lib/python3.10/dist-packages (from jsonschema>=3.0->altair<6,>=4.0->streamlit) (23.2.0)
Requirement already satisfied: jsonschema-specifications>=2023.03.6 in /usr/local/lib/python3.10/dist-packages (from jsonschema>=3.0->altair<6,>=4.0->streamlit) (2023.12.1)
Requirement already satisfied: referencing>=0.28.4 in /usr/local/lib/python3.10/dist-packages (from jsonschema>=3.0->altair<6,>=4.0->streamlit) (0.34.0)
Requirement already satisfied: rpds-py>=0.7.1 in /usr/local/lib/python3.10/dist-packages (from jsonschema>=3.0->altair<6,>=4.0->streamlit) (0.18.0)
Requirement already satisfied: mdurl~0.1 in /usr/local/lib/python3.10/dist-packages (from markdown-it-py>=2.2.0->rich<14,>=10.14.0->streamlit) (0.1.2)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas<3,>=1.3.0->streamlit) (1.16.0)
Installing collected packages: watchdog, smmap, pyarrow, pydeck, gitdb, gitpython, streamlit
Successfully installed gitdb-4.0.11 gitpython-3.1.43 pyarrow-16.0.0 pydeck-0.9.0b1 smmap-5.0.1 streamlit-1.33.0 watchdog-4.0.0
```

```
# Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
import re

# Import machine learning libraries
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import classification_report, accuracy_score

# Import Streamlit libraries
import streamlit as st
import joblib
from sklearn.feature_extraction.text import TfidfVectorizer

# Load data
url = 'https://github.com/KristerMartinez/TwitterUSAirlineSentiment/raw/main/TwitterUSAirlineSentiment.csv'
data = pd.read_csv(url)
data.head()
```

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline
0	570306133677760513	neutral	1.0000	NaN	NaN	V Ame
1	570301130888122368	positive	0.3486	NaN	0.0000	V Ame
2	570301083672813571	neutral	0.6837	NaN	NaN	V Ame
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	V Ame
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	V Ame

```
data.info()
```

```
→ <class 'pandas.core.frame.DataFrame'>
RangeIndex: 14640 entries, 0 to 14639
Data columns (total 15 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   tweet_id         14640 non-null   int64  
 1   airline_sentiment 14640 non-null   object  
 2   airline_sentiment_confidence 14640 non-null   float64 
 3   negativereason    9178 non-null   object  
 4   negativereason_confidence 10522 non-null   float64 
 5   airline           14640 non-null   object  
 6   airline_sentiment_gold 40 non-null   object  
 7   name              14640 non-null   object  
 8   negativereason_gold 32 non-null   object  
 9   retweet_count     14640 non-null   int64  
 10  text              14640 non-null   object  
 11  tweet_coord       1019 non-null   object  
 12  tweet_created     14640 non-null   object  
 13  tweet_location    9907 non-null   object  
 14  user_timezone     9820 non-null   object  
dtypes: float64(2), int64(2), object(11)
memory usage: 1.7+ MB
```

```
# Data Cleaning
def clean_text(text):
    text = re.sub(r'\W', ' ', str(text))
    text = re.sub(r'\s+[a-zA-Z]\s+', ' ', text)
    text = re.sub(r'^[a-zA-Z]\s+', ' ', text)
    text = re.sub(r'\s+', ' ', text, flags=re.I)
    text = re.sub(r'^b\s+', '', text)
    text = text.lower()
    return text

data['cleaned_text'] = data['text'].apply(clean_text)

# Data Preparation
vectorizer = TfidfVectorizer(max_features=2500, min_df=7, max_df=0.8)
X = vectorizer.fit_transform(data['cleaned_text']).toarray()
y = data['airline_sentiment'].values

# Split the dataset
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# Train the Logistic Regression model
model = LogisticRegression()
model.fit(X_train, y_train)

# Predict and Evaluate the model
predictions = model.predict(X_test)
print(classification_report(y_test, predictions))
print("Accuracy:", accuracy_score(y_test, predictions))
```

	precision	recall	f1-score	support
negative	0.83	0.93	0.88	1870
neutral	0.69	0.56	0.61	614
positive	0.82	0.63	0.71	444
accuracy			0.81	2928
macro avg	0.78	0.71	0.74	2928
weighted avg	0.80	0.81	0.80	2928

Accuracy: 0.8080601092896175
/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfsgs failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
<https://scikit-learn.org/stable/modules/preprocessing.html>
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result()
```

```
joblib.dump(model, 'finalized_model.pkl')
```

→ ['finalized_model.pkl']

app.py file

```
app_content= """import streamlit as st
import joblib

# Load your trained model
model = joblib.load('finalized_model.pkl')

# Streamlit application
st.title('Sentiment Analysis Tool')

# Text input
user_input = st.text_area("Enter Text", "")

# Predict button
if st.button('Predict'):
    # You will need to vectorize the user input as well
    # Here you can load the vectorizer or redefine it
    result = model.predict([user_input])
    st.write('The sentiment is:', 'Positive' if result[0] == 1 else 'Negative')
"""

# Write the content to app.py file
with open('app.py', 'w') as file:
    file.write(app_content)
```