

```
In [6]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
import re
import string

from sklearn.model_selection import train_test_split
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics import classification_report, accuracy_score, confusion_mat

import warnings
warnings.filterwarnings("ignore")

nltk.download('stopwords')
nltk.download('punkt')
nltk.download('wordnet')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Out[6]: True

```
In [9]: df = pd.read_csv("imdb_data.csv")
df.head()
```

Out[9]:

		review	sentiment
0	One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with me. The first thing that struck me...		positive
1	A wonderful little production. The filming technique is very unassuming- very old-time-BBC fashion and gives a comforting, and sometimes discomforting, sense of realism to the entire p...		positive
2	I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in the air conditioned theater and watching a light-hearted comedy. The plot is simplistic, but the dialogue i...		positive
3	Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet & his parents are fighting all the time. This movie is slower than a soap opera... and suddenl...		negative
4	Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Mattei offers us a vivid portrait about human relations. This is a movie that seems to be telling us what mone...		positive

```
In [10]: print("Dataset shape:", df.shape)

print("\nMissing values:\n", df.isnull().sum())

df.info()
```

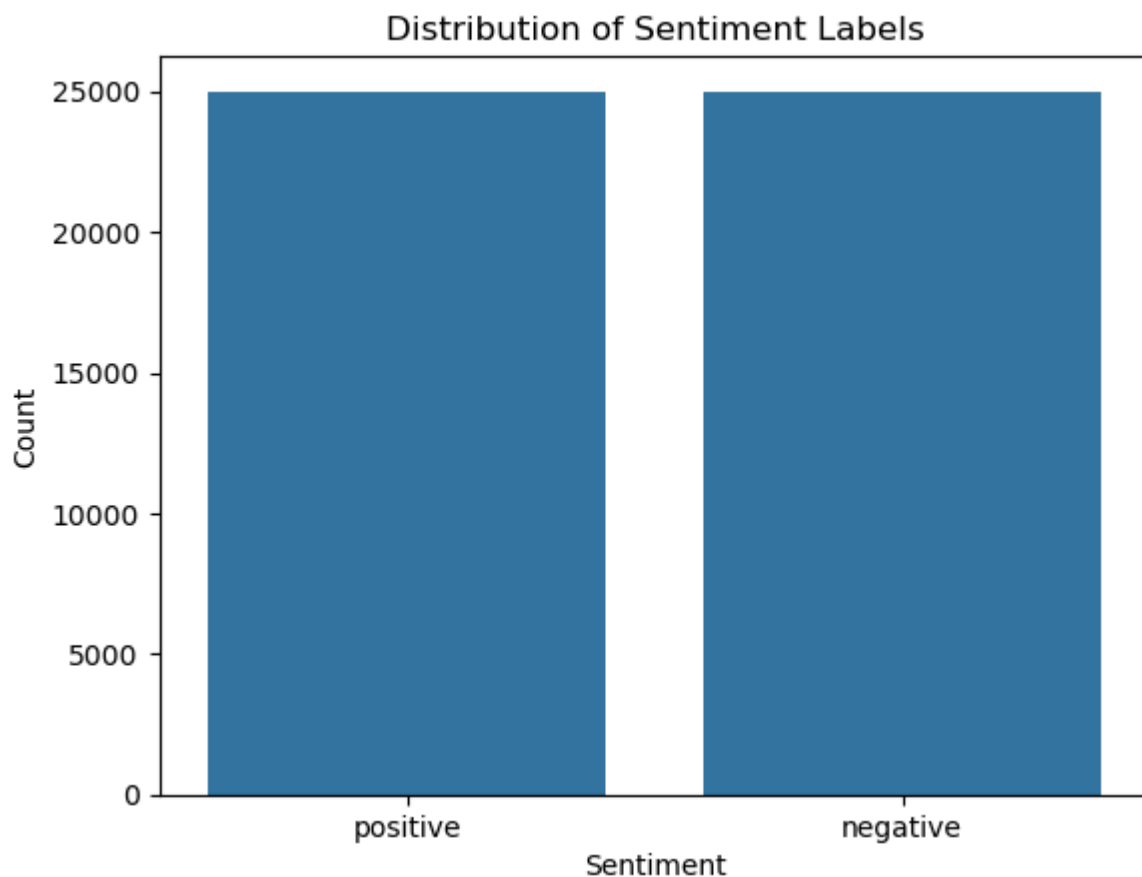
Dataset shape: (50000, 2)

Missing values:

```
review      0
sentiment   0
dtype: int64
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   review      50000 non-null   object
1   sentiment   50000 non-null   object
dtypes: object(2)
memory usage: 781.4+ KB
```

```
In [11]: sns.countplot(x='sentiment', data=df)
plt.title("Distribution of Sentiment Labels")
plt.xlabel("Sentiment")
plt.ylabel("Count")
plt.show()

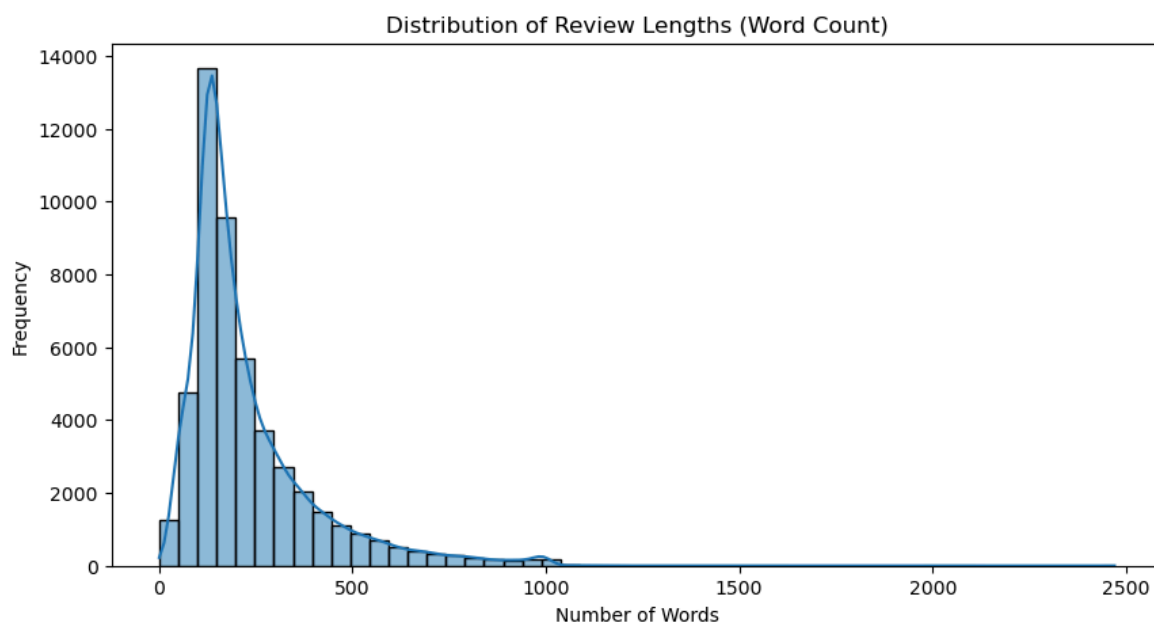
df['sentiment'].value_counts()
```



```
Out[11]: sentiment
positive    25000
negative    25000
Name: count, dtype: int64
```

```
In [12]: df['review_length'] = df['review'].apply(lambda x: len(x.split()))

plt.figure(figsize=(10, 5))
sns.histplot(df['review_length'], bins=50, kde=True)
plt.title("Distribution of Review Lengths (Word Count)")
plt.xlabel("Number of Words")
plt.ylabel("Frequency")
plt.show()
```



```
In [17]: import re
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

lemmatizer = WordNetLemmatizer()
stop_words = set(stopwords.words('english'))

def preprocess_text(text):

    text = text.lower()

    text = re.sub(r'<.*?>', '', text)

    text = re.sub(r'^a-z\s', '', text)

    words = text.split()

    words = [lemmatizer.lemmatize(word) for word in words if word not in stop_words]

    return ' '.join(words)
```

```
In [18]: df['clean_review'] = df['review'].apply(preprocess_text)
```

```
In [20]: df[['review', 'clean_review']].head()
```

```
Out[20]:
```

	review	clean_review
0	One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with me. The first thing that struck me...	one reviewer mentioned watching oz episode youll hooked right exactly happened methe first thing struck oz brutality unflinching scene violence set right word go trust show faint hearted timid sho...
1	A wonderful little production. The filming technique is very unassuming- very old-time-BBC fashion and gives a comforting, and sometimes discomforting, sense of realism to the entire p...	wonderful little production filming technique unassuming oldtimebbc fashion give comforting sometimes discomforting sense realism entire piece actor extremely well chosen michael sheen got polari ...
2	I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in the air conditioned theater and watching a light-hearted comedy. The plot is simplistic, but the dialogue i...	thought wonderful way spend time hot summer weekend sitting air conditioned theater watching lighthearted comedy plot simplistic dialogue witty character likable even well bread suspected serial k...
3	Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet & his parents are fighting all the time. This movie is slower than a soap opera... and suddenl...	basically there family little boy jake think there zombie closet parent fighting timethis movie slower soap opera suddenly jake decides become rambo kill zombieok first youre going make film must ...
4	Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Mattei offers us a vivid portrait about human relations. This is a movie that seems to be telling us what mone...	petter matteis love time money visually stunning film watch mr mattei offer u vivid portrait human relation movie seems telling u money power success people different situation encounter variation...

```
In [22]: df['cleaned_review'] = df['review'].apply(preprocess_text)
```

```
In [23]: from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

cleaned_reviews = df['cleaned_review']

bow_vectorizer = CountVectorizer(max_features=5000)
X_bow = bow_vectorizer.fit_transform(cleaned_reviews)

tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X_tfidf = tfidf_vectorizer.fit_transform(cleaned_reviews)

print("Shape of Bag of Words matrix:", X_bow.shape)
print("Shape of TF-IDF matrix:", X_tfidf.shape)
```

Shape of Bag of Words matrix: (50000, 5000)

Shape of TF-IDF matrix: (50000, 5000)

```
In [24]: from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer

bow_vectorizer = CountVectorizer(max_features=5000)
X_bow = bow_vectorizer.fit_transform(df['cleaned_review'])

tfidf_vectorizer = TfidfVectorizer(max_features=5000)
X_tfidf = tfidf_vectorizer.fit_transform(df['cleaned_review'])

print("Bow shape:", X_bow.shape)
print("TF-IDF shape:", X_tfidf.shape)
```

Bow shape: (50000, 5000)

TF-IDF shape: (50000, 5000)

```
In [25]: import string
from nltk.corpus import stopwords

stop_words = set(stopwords.words('english'))

def extract_features(text):
    words = text.split()
    word_count = len(words)
    char_count = len(text)
    avg_word_length = char_count / word_count if word_count != 0 else 0
    stopword_count = sum(1 for word in words if word in stop_words)
    punctuation_count = sum(1 for char in text if char in string.punctuation)

    return pd.Series([word_count, char_count, avg_word_length, stopword_count, p

df[['word_count', 'char_count', 'avg_word_length', 'stopword_count', 'punctuatio

# Preview first few rows
df[['cleaned_review', 'word_count', 'char_count', 'avg_word_length', 'stopword_c
```

Out[25]:

	cleaned_review	word_count	char_count	avg_word_length	stopword_count	punctua
0	one reviewer mentioned watching oz episode youll hooked right exactly happened methe first thing struck oz brutality unflinching scene violence set right word go trust show faint hearted timid sho...	167.0	1125.0	6.736527	0.0	
1	wonderful little production filming technique unassuming oldtimebbc fashion give comforting sometimes discomforting sense realism entire piece actor extremely well chosen michael sheen got polari ...	84.0	640.0	7.619048	0.0	
2	thought wonderful way spend time hot summer weekend sitting air conditioned theater watching lighthearted comedy plot simplistic dialogue witty character likable even well bread suspected serial k...	85.0	580.0	6.823529	0.0	
3	basically there family little boy jake think there zombie closet parent fighting timethis movie	66.0	446.0	6.757576	2.0	

	cleaned_review	word_count	char_count	avg_word_length	stopword_count	punctua
	slower soap opera suddenly jake decides become rambo kill zombieok first youre going make film must ...					
	petter matteis love time money visually stunning film watch mr mattei offer u vivid portrait					
4	human relation movie seems telling u money power success people different situation encounter variation...	125.0	851.0	6.808000	0.0	

```
In [26]: df['word_count'] = df['cleaned_review'].apply(lambda x: len(x.split()))
df['char_count'] = df['cleaned_review'].apply(len)
df['avg_word_length'] = df['char_count'] / df['word_count']

df[['cleaned_review', 'word_count', 'char_count', 'avg_word_length']].head()
```

Out[26]:

	cleaned_review	word_count	char_count	avg_word_length
0	one reviewer mentioned watching oz episode youll hooked right exactly happened methe first thing struck oz brutality unflinching scene violence set right word go trust show faint hearted timid sho...	167	1125	6.736527
1	wonderful little production filming technique unassuming oldtimebbc fashion give comforting sometimes discomforting sense realism entire piece actor extremely well chosen michael sheen got polari ...	84	640	7.619048
2	thought wonderful way spend time hot summer weekend sitting air conditioned theater watching lighthearted comedy plot simplistic dialogue witty character likable even well bread suspected serial k...	85	580	6.823529
3	basically there family little boy jake think there zombie closet parent fighting timethis movie slower soap opera suddenly jake decides become rambo kill zombieok first youre going make film must ...	66	446	6.757576
4	petter matteis love time money visually stunning film watch mr mattei offer u vivid portrait human relation movie seems telling u money power success people different situation encounter variation...	125	851	6.808000

```
In [28]: from sklearn.feature_extraction.text import TfidfVectorizer

tfidf_vectorizer = TfidfVectorizer(max_features=5000)

tfidf_features = tfidf_vectorizer.fit_transform(df['cleaned_review'])
```

```
In [29]: X = tfidf_features
y = df['sentiment']
```

```
In [30]: from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report, confusion_mat

X = tfidf_features
y = df['sentiment']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_
```

In []:

```
In [31]: from sklearn.linear_model import LogisticRegression

lr_model = LogisticRegression(max_iter=1000)
```



```
lr_model.fit(X_train, y_train)

y_pred_lr = lr_model.predict(X_test)

print("Logistic Regression Results:")
print("Accuracy:", accuracy_score(y_test, y_pred_lr))
print(confusion_matrix(y_test, y_pred_lr))
print(classification_report(y_test, y_pred_lr))
```

Logistic Regression Results:

Accuracy: 0.885

[[4322 639]

[511 4528]]

	precision	recall	f1-score	support
negative	0.89	0.87	0.88	4961
positive	0.88	0.90	0.89	5039
accuracy			0.89	10000
macro avg	0.89	0.88	0.88	10000
weighted avg	0.89	0.89	0.88	10000

In [32]: `from sklearn.naive_bayes import MultinomialNB`

```
nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)

y_pred_nb = nb_model.predict(X_test)

print("Naive Bayes Results:")
print("Accuracy:", accuracy_score(y_test, y_pred_nb))
print(confusion_matrix(y_test, y_pred_nb))
print(classification_report(y_test, y_pred_nb))
```

Naive Bayes Results:

Accuracy: 0.849

[[4188 773]

[737 4302]]

	precision	recall	f1-score	support
negative	0.85	0.84	0.85	4961
positive	0.85	0.85	0.85	5039
accuracy			0.85	10000
macro avg	0.85	0.85	0.85	10000
weighted avg	0.85	0.85	0.85	10000

In [33]: `from sklearn.svm import LinearSVC`

```
svm_model = LinearSVC()
svm_model.fit(X_train, y_train)

y_pred_svm = svm_model.predict(X_test)

print("SVM Results:")
print("Accuracy:", accuracy_score(y_test, y_pred_svm))
print(confusion_matrix(y_test, y_pred_svm))
print(classification_report(y_test, y_pred_svm))
```

SVM Results:

Accuracy: 0.8785

[[4307 654]

[561 4478]]

	precision	recall	f1-score	support
negative	0.88	0.87	0.88	4961
positive	0.87	0.89	0.88	5039
accuracy			0.88	10000
macro avg	0.88	0.88	0.88	10000
weighted avg	0.88	0.88	0.88	10000

```
In [40]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
```

```
X_train_small = X_train[:5000]
y_train_small = y_train[:5000]
```

```
rf_model = RandomForestClassifier(n_estimators=50, max_depth=15, random_state=42)
rf_model.fit(X_train_small, y_train_small)
```

```
y_pred_rf = rf_model.predict(X_test)
```

```
rf_accuracy = accuracy_score(y_test, y_pred_rf)
rf_confusion = confusion_matrix(y_test, y_pred_rf)
rf_report = classification_report(y_test, y_pred_rf)
```

```
print("Random Forest Accuracy:", rf_accuracy)
print("Confusion Matrix:\n", rf_confusion)
print("Classification Report:\n", rf_report)
```

Random Forest Accuracy: 0.8123

Confusion Matrix:

[[3805 1156]

[721 4318]]

Classification Report:

	precision	recall	f1-score	support
negative	0.84	0.77	0.80	4961
positive	0.79	0.86	0.82	5039
accuracy			0.81	10000
macro avg	0.81	0.81	0.81	10000
weighted avg	0.81	0.81	0.81	10000

```
In [17]: import re
import nltk
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import WordNetLemmatizer
```

```
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
```

```
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()

def clean_text(text):
    text = re.sub(r'<.*?>', ' ', text)
    text = re.sub(r'^a-zA-Z', ' ', text)
    text = text.lower()
    words = word_tokenize(text)
    words = [lemmatizer.lemmatize(word) for word in words if word not in stop_words]
    return ' '.join(words)
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

```
In [21]: !pip install tensorflow
```

```

Collecting tensorflow
  Downloading tensorflow-2.19.0-cp312-cp312-win_amd64.whl.metadata (4.1 kB)
Collecting absl-py>=1.0.0 (from tensorflow)
  Downloading absl_py-2.3.0-py3-none-any.whl.metadata (2.4 kB)
Collecting astunparse>=1.6.0 (from tensorflow)
  Downloading astunparse-1.6.3-py2.py3-none-any.whl.metadata (4.4 kB)
Collecting flatbuffers>=24.3.25 (from tensorflow)
  Downloading flatbuffers-25.2.10-py2.py3-none-any.whl.metadata (875 bytes)
Collecting gast!=0.5.0,!0.5.1,!0.5.2,>=0.2.1 (from tensorflow)
  Downloading gast-0.6.0-py3-none-any.whl.metadata (1.3 kB)
Collecting google-pasta>=0.1.1 (from tensorflow)
  Downloading google_pasta-0.2.0-py3-none-any.whl.metadata (814 bytes)
Collecting libclang>=13.0.0 (from tensorflow)
  Downloading libclang-18.1.1-py2.py3-none-win_amd64.whl.metadata (5.3 kB)
Collecting opt-einsum>=2.3.2 (from tensorflow)
  Downloading opt_einsum-3.4.0-py3-none-any.whl.metadata (6.3 kB)
Requirement already satisfied: packaging in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (24.1)
Requirement already satisfied: protobuf!=4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<6.0.0dev,>=3.20.3 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (4.25.3)
Requirement already satisfied: requests<3,>=2.21.0 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (75.1.0)
Requirement already satisfied: six>=1.12.0 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (1.16.0)
Collecting termcolor>=1.1.0 (from tensorflow)
  Downloading termcolor-3.1.0-py3-none-any.whl.metadata (6.4 kB)
Requirement already satisfied: typing-extensions>=3.6.6 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (4.11.0)
Requirement already satisfied: wrapt>=1.11.0 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (1.14.1)
Collecting grpcio<2.0,>=1.24.3 (from tensorflow)
  Downloading grpcio-1.73.0-cp312-cp312-win_amd64.whl.metadata (4.0 kB)
Collecting tensorboard~2.19.0 (from tensorflow)
  Downloading tensorboard-2.19.0-py3-none-any.whl.metadata (1.8 kB)
Collecting keras>=3.5.0 (from tensorflow)
  Downloading keras-3.10.0-py3-none-any.whl.metadata (6.0 kB)
Requirement already satisfied: numpy<2.2.0,>=1.26.0 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (1.26.4)
Requirement already satisfied: h5py>=3.11.0 in c:\users\jayde\anaconda3\lib\site-packages (from tensorflow) (3.11.0)
Collecting ml-dtypes<1.0.0,>=0.5.1 (from tensorflow)
  Downloading ml_dtypes-0.5.1-cp312-cp312-win_amd64.whl.metadata (22 kB)
Requirement already satisfied: wheel<1.0,>=0.23.0 in c:\users\jayde\anaconda3\lib\site-packages (from astunparse>=1.6.0->tensorflow) (0.44.0)
Requirement already satisfied: rich in c:\users\jayde\anaconda3\lib\site-packages (from keras>=3.5.0->tensorflow) (13.7.1)
Collecting namex (from keras>=3.5.0->tensorflow)
  Downloading namex-0.1.0-py3-none-any.whl.metadata (322 bytes)
Collecting optree (from keras>=3.5.0->tensorflow)
  Downloading optree-0.16.0-cp312-cp312-win_amd64.whl.metadata (31 kB)
Requirement already satisfied: charset-normalizer<4,>=2 in c:\users\jayde\anaconda3\lib\site-packages (from requests<3,>=2.21.0->tensorflow) (3.3.2)
Requirement already satisfied: idna<4,>=2.5 in c:\users\jayde\anaconda3\lib\site-packages (from requests<3,>=2.21.0->tensorflow) (3.7)
Requirement already satisfied: urllib3<3,>=1.21.1 in c:\users\jayde\anaconda3\lib\site-packages (from requests<3,>=2.21.0->tensorflow) (2.2.3)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\jayde\anaconda3\lib

```

```
\site-packages (from requests<3,>=2.21.0->tensorflow) (2024.8.30)
Requirement already satisfied: markdown>=2.6.8 in c:\users\jayde\anaconda3\lib\si
te-packages (from tensorboard~=2.19.0->tensorflow) (3.4.1)
Collecting tensorboard-data-server<0.8.0,>=0.7.0 (from tensorboard~=2.19.0->tenso
rflow)
  Downloading tensorboard_data_server-0.7.2-py3-none-any.whl.metadata (1.1 kB)
Requirement already satisfied: werkzeug>=1.0.1 in c:\users\jayde\anaconda3\lib\si
te-packages (from tensorboard~=2.19.0->tensorflow) (3.0.3)
Requirement already satisfied: MarkupSafe>=2.1.1 in c:\users\jayde\anaconda3\lib
\site-packages (from werkzeug>=1.0.1->tensorboard~=2.19.0->tensorflow) (2.1.3)
Requirement already satisfied: markdown-it-py>=2.2.0 in c:\users\jayde\anaconda3
\lib\site-packages (from rich->keras>=3.5.0->tensorflow) (2.2.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in c:\users\jayde\anaconda
3\lib\site-packages (from rich->keras>=3.5.0->tensorflow) (2.15.1)
Requirement already satisfied: mdurl~=0.1 in c:\users\jayde\anaconda3\lib\site-pa
ckages (from markdown-it-py>=2.2.0->rich->keras>=3.5.0->tensorflow) (0.1.0)
Downloading tensorflow-2.19.0-cp312-cp312-win_amd64.whl (376.0 MB)
----- 0.0/376.0 MB ? eta -:--:--
----- 0.3/376.0 MB ? eta -:--:--
----- 2.1/376.0 MB 6.9 MB/s eta 0:00:55
----- 2.9/376.0 MB 6.0 MB/s eta 0:01:03
----- 4.2/376.0 MB 5.7 MB/s eta 0:01:06
----- 6.3/376.0 MB 7.0 MB/s eta 0:00:53
----- 9.4/376.0 MB 8.2 MB/s eta 0:00:45
----- 12.8/376.0 MB 9.4 MB/s eta 0:00:39
----- 16.5/376.0 MB 10.5 MB/s eta 0:00:35
----- 19.1/376.0 MB 11.1 MB/s eta 0:00:33
----- 22.3/376.0 MB 11.0 MB/s eta 0:00:33
----- 25.7/376.0 MB 11.5 MB/s eta 0:00:31
----- 28.0/376.0 MB 11.4 MB/s eta 0:00:31
----- 30.4/376.0 MB 11.5 MB/s eta 0:00:31
----- 32.8/376.0 MB 11.4 MB/s eta 0:00:31
----- 35.1/376.0 MB 11.3 MB/s eta 0:00:31
----- 37.5/376.0 MB 11.3 MB/s eta 0:00:30
----- 39.8/376.0 MB 11.3 MB/s eta 0:00:30
----- 42.5/376.0 MB 11.3 MB/s eta 0:00:30
----- 45.1/376.0 MB 11.4 MB/s eta 0:00:30
----- 48.0/376.0 MB 11.4 MB/s eta 0:00:29
----- 50.6/376.0 MB 11.5 MB/s eta 0:00:29
----- 53.5/376.0 MB 11.5 MB/s eta 0:00:28
----- 55.3/376.0 MB 11.4 MB/s eta 0:00:29
----- 57.9/376.0 MB 11.4 MB/s eta 0:00:28
----- 60.6/376.0 MB 11.4 MB/s eta 0:00:28
----- 62.7/376.0 MB 11.4 MB/s eta 0:00:28
----- 65.5/376.0 MB 11.4 MB/s eta 0:00:28
----- 68.2/376.0 MB 11.5 MB/s eta 0:00:27
----- 70.8/376.0 MB 11.6 MB/s eta 0:00:27
----- 73.9/376.0 MB 11.6 MB/s eta 0:00:27
----- 76.8/376.0 MB 11.7 MB/s eta 0:00:26
----- 79.7/376.0 MB 11.7 MB/s eta 0:00:26
----- 82.6/376.0 MB 11.8 MB/s eta 0:00:25
----- 85.5/376.0 MB 11.8 MB/s eta 0:00:25
----- 88.1/376.0 MB 11.9 MB/s eta 0:00:25
----- 90.2/376.0 MB 11.8 MB/s eta 0:00:25
----- 93.3/376.0 MB 11.8 MB/s eta 0:00:24
----- 96.2/376.0 MB 11.9 MB/s eta 0:00:24
----- 99.4/376.0 MB 11.9 MB/s eta 0:00:24
----- 102.5/376.0 MB 12.0 MB/s eta 0:00:23
----- 105.4/376.0 MB 12.0 MB/s eta 0:00:23
----- 108.3/376.0 MB 12.1 MB/s eta 0:00:23
```

```
----- 111.4/376.0 MB 12.1 MB/s eta 0:00:22
----- 113.5/376.0 MB 12.1 MB/s eta 0:00:22
----- 115.9/376.0 MB 12.0 MB/s eta 0:00:22
----- 118.0/376.0 MB 12.0 MB/s eta 0:00:22
----- 120.1/376.0 MB 11.9 MB/s eta 0:00:22
----- 122.4/376.0 MB 11.9 MB/s eta 0:00:22
----- 124.8/376.0 MB 11.9 MB/s eta 0:00:22
----- 127.1/376.0 MB 11.9 MB/s eta 0:00:21
----- 129.8/376.0 MB 11.9 MB/s eta 0:00:21
----- 132.1/376.0 MB 11.9 MB/s eta 0:00:21
----- 134.5/376.0 MB 11.8 MB/s eta 0:00:21
----- 136.1/376.0 MB 11.7 MB/s eta 0:00:21
----- 136.6/376.0 MB 11.6 MB/s eta 0:00:21
----- 137.6/376.0 MB 11.5 MB/s eta 0:00:21
----- 139.2/376.0 MB 11.4 MB/s eta 0:00:21
----- 141.6/376.0 MB 11.4 MB/s eta 0:00:21
----- 144.2/376.0 MB 11.4 MB/s eta 0:00:21
----- 146.8/376.0 MB 11.4 MB/s eta 0:00:21
----- 149.4/376.0 MB 11.4 MB/s eta 0:00:20
----- 152.0/376.0 MB 11.4 MB/s eta 0:00:20
----- 155.2/376.0 MB 11.5 MB/s eta 0:00:20
----- 157.5/376.0 MB 11.5 MB/s eta 0:00:20
----- 159.1/376.0 MB 11.4 MB/s eta 0:00:20
----- 161.7/376.0 MB 11.4 MB/s eta 0:00:19
----- 164.1/376.0 MB 11.4 MB/s eta 0:00:19
----- 166.7/376.0 MB 11.4 MB/s eta 0:00:19
----- 168.8/376.0 MB 11.4 MB/s eta 0:00:19
----- 171.4/376.0 MB 11.4 MB/s eta 0:00:18
----- 173.8/376.0 MB 11.4 MB/s eta 0:00:18
----- 176.7/376.0 MB 11.4 MB/s eta 0:00:18
----- 179.6/376.0 MB 11.4 MB/s eta 0:00:18
----- 182.5/376.0 MB 11.5 MB/s eta 0:00:17
----- 185.3/376.0 MB 11.5 MB/s eta 0:00:17
----- 188.0/376.0 MB 11.5 MB/s eta 0:00:17
----- 191.1/376.0 MB 11.5 MB/s eta 0:00:17
----- 194.0/376.0 MB 11.6 MB/s eta 0:00:16
----- 196.9/376.0 MB 11.6 MB/s eta 0:00:16
----- 199.5/376.0 MB 11.6 MB/s eta 0:00:16
----- 202.6/376.0 MB 11.6 MB/s eta 0:00:15
----- 205.5/376.0 MB 11.6 MB/s eta 0:00:15
----- 207.9/376.0 MB 11.6 MB/s eta 0:00:15
----- 211.0/376.0 MB 11.7 MB/s eta 0:00:15
----- 213.6/376.0 MB 11.7 MB/s eta 0:00:14
----- 216.5/376.0 MB 11.7 MB/s eta 0:00:14
----- 219.2/376.0 MB 11.7 MB/s eta 0:00:14
----- 221.8/376.0 MB 11.7 MB/s eta 0:00:14
----- 224.7/376.0 MB 11.7 MB/s eta 0:00:13
----- 227.0/376.0 MB 11.7 MB/s eta 0:00:13
----- 229.4/376.0 MB 11.7 MB/s eta 0:00:13
----- 231.5/376.0 MB 11.7 MB/s eta 0:00:13
----- 233.8/376.0 MB 11.7 MB/s eta 0:00:13
----- 235.9/376.0 MB 11.6 MB/s eta 0:00:13
----- 237.2/376.0 MB 11.6 MB/s eta 0:00:12
----- 239.3/376.0 MB 11.6 MB/s eta 0:00:12
----- 242.0/376.0 MB 11.6 MB/s eta 0:00:12
----- 244.3/376.0 MB 11.5 MB/s eta 0:00:12
----- 246.4/376.0 MB 11.5 MB/s eta 0:00:12
----- 249.0/376.0 MB 11.5 MB/s eta 0:00:12
----- 251.7/376.0 MB 11.5 MB/s eta 0:00:11
----- 254.3/376.0 MB 11.5 MB/s eta 0:00:11
```

```
----- 257.2/376.0 MB 11.6 MB/s eta 0:00:11
----- 260.0/376.0 MB 11.6 MB/s eta 0:00:11
----- 262.7/376.0 MB 11.6 MB/s eta 0:00:10
----- 264.8/376.0 MB 11.7 MB/s eta 0:00:10
----- 266.1/376.0 MB 11.7 MB/s eta 0:00:10
----- 267.4/376.0 MB 11.7 MB/s eta 0:00:10
----- 268.4/376.0 MB 11.6 MB/s eta 0:00:10
----- 270.8/376.0 MB 11.6 MB/s eta 0:00:10
----- 272.4/376.0 MB 11.5 MB/s eta 0:00:09
----- 274.2/376.0 MB 11.5 MB/s eta 0:00:09
----- 275.3/376.0 MB 11.4 MB/s eta 0:00:09
----- 276.0/376.0 MB 11.3 MB/s eta 0:00:09
----- 277.1/376.0 MB 11.2 MB/s eta 0:00:09
----- 278.1/376.0 MB 11.1 MB/s eta 0:00:09
----- 278.9/376.0 MB 11.1 MB/s eta 0:00:09
----- 280.0/376.0 MB 11.0 MB/s eta 0:00:09
----- 280.5/376.0 MB 10.9 MB/s eta 0:00:09
----- 281.5/376.0 MB 10.9 MB/s eta 0:00:09
----- 281.8/376.0 MB 10.8 MB/s eta 0:00:09
----- 282.3/376.0 MB 10.8 MB/s eta 0:00:09
----- 282.9/376.0 MB 10.7 MB/s eta 0:00:09
----- 283.4/376.0 MB 10.6 MB/s eta 0:00:09
----- 283.9/376.0 MB 10.5 MB/s eta 0:00:09
----- 284.4/376.0 MB 10.4 MB/s eta 0:00:09
----- 285.0/376.0 MB 10.3 MB/s eta 0:00:09
----- 285.5/376.0 MB 10.2 MB/s eta 0:00:09
----- 286.0/376.0 MB 10.2 MB/s eta 0:00:09
----- 286.5/376.0 MB 10.1 MB/s eta 0:00:09
----- 287.0/376.0 MB 10.0 MB/s eta 0:00:09
----- 287.8/376.0 MB 10.0 MB/s eta 0:00:09
----- 288.4/376.0 MB 9.9 MB/s eta 0:00:09
----- 288.9/376.0 MB 9.9 MB/s eta 0:00:09
----- 289.4/376.0 MB 9.8 MB/s eta 0:00:09
----- 289.9/376.0 MB 9.7 MB/s eta 0:00:09
----- 290.5/376.0 MB 9.7 MB/s eta 0:00:09
----- 291.0/376.0 MB 9.6 MB/s eta 0:00:09
----- 291.2/376.0 MB 9.5 MB/s eta 0:00:09
----- 291.5/376.0 MB 9.5 MB/s eta 0:00:09
----- 291.8/376.0 MB 9.4 MB/s eta 0:00:09
----- 292.0/376.0 MB 9.4 MB/s eta 0:00:09
----- 292.3/376.0 MB 9.3 MB/s eta 0:00:10
----- 292.6/376.0 MB 9.2 MB/s eta 0:00:10
----- 292.8/376.0 MB 9.2 MB/s eta 0:00:10
----- 293.1/376.0 MB 9.1 MB/s eta 0:00:10
----- 293.3/376.0 MB 9.0 MB/s eta 0:00:10
----- 293.6/376.0 MB 9.0 MB/s eta 0:00:10
----- 293.9/376.0 MB 8.9 MB/s eta 0:00:10
----- 294.1/376.0 MB 8.9 MB/s eta 0:00:10
----- 294.4/376.0 MB 8.8 MB/s eta 0:00:10
----- 294.9/376.0 MB 8.7 MB/s eta 0:00:10
----- 295.2/376.0 MB 8.7 MB/s eta 0:00:10
----- 295.4/376.0 MB 8.6 MB/s eta 0:00:10
----- 295.7/376.0 MB 8.6 MB/s eta 0:00:10
----- 296.0/376.0 MB 8.5 MB/s eta 0:00:10
----- 296.5/376.0 MB 8.4 MB/s eta 0:00:10
----- 296.7/376.0 MB 8.3 MB/s eta 0:00:10
----- 297.3/376.0 MB 8.2 MB/s eta 0:00:10
----- 297.5/376.0 MB 8.2 MB/s eta 0:00:10
----- 297.8/376.0 MB 8.1 MB/s eta 0:00:10
----- 298.3/376.0 MB 8.0 MB/s eta 0:00:10
```

```
----- 298.6/376.0 MB 8.0 MB/s eta 0:00:10
----- 299.1/376.0 MB 7.9 MB/s eta 0:00:10
----- 299.6/376.0 MB 7.8 MB/s eta 0:00:10
----- 299.9/376.0 MB 7.8 MB/s eta 0:00:10
----- 300.4/376.0 MB 7.7 MB/s eta 0:00:10
----- 300.7/376.0 MB 7.6 MB/s eta 0:00:10
----- 301.2/376.0 MB 7.5 MB/s eta 0:00:10
----- 301.5/376.0 MB 7.4 MB/s eta 0:00:11
----- 301.7/376.0 MB 7.3 MB/s eta 0:00:11
----- 301.7/376.0 MB 7.3 MB/s eta 0:00:11
----- 302.0/376.0 MB 7.2 MB/s eta 0:00:11
----- 302.3/376.0 MB 7.1 MB/s eta 0:00:11
----- 302.5/376.0 MB 7.0 MB/s eta 0:00:11
----- 302.5/376.0 MB 7.0 MB/s eta 0:00:11
----- 302.8/376.0 MB 6.9 MB/s eta 0:00:11
----- 303.0/376.0 MB 6.7 MB/s eta 0:00:11
----- 303.0/376.0 MB 6.7 MB/s eta 0:00:11
----- 303.3/376.0 MB 6.6 MB/s eta 0:00:12
----- 303.6/376.0 MB 6.5 MB/s eta 0:00:12
----- 303.8/376.0 MB 6.4 MB/s eta 0:00:12
----- 304.1/376.0 MB 6.3 MB/s eta 0:00:12
----- 304.1/376.0 MB 6.3 MB/s eta 0:00:12
----- 304.3/376.0 MB 6.2 MB/s eta 0:00:12
----- 304.6/376.0 MB 6.1 MB/s eta 0:00:12
----- 304.6/376.0 MB 6.1 MB/s eta 0:00:12
----- 304.9/376.0 MB 6.0 MB/s eta 0:00:12
----- 305.1/376.0 MB 5.9 MB/s eta 0:00:13
----- 305.4/376.0 MB 5.8 MB/s eta 0:00:13
----- 305.7/376.0 MB 5.7 MB/s eta 0:00:13
----- 305.7/376.0 MB 5.7 MB/s eta 0:00:13
----- 305.9/376.0 MB 5.7 MB/s eta 0:00:13
----- 306.2/376.0 MB 5.6 MB/s eta 0:00:13
----- 306.4/376.0 MB 5.6 MB/s eta 0:00:13
----- 306.7/376.0 MB 5.5 MB/s eta 0:00:13
----- 306.7/376.0 MB 5.5 MB/s eta 0:00:13
----- 307.0/376.0 MB 5.4 MB/s eta 0:00:13
----- 307.2/376.0 MB 5.3 MB/s eta 0:00:14
----- 307.2/376.0 MB 5.3 MB/s eta 0:00:14
----- 307.5/376.0 MB 5.2 MB/s eta 0:00:14
----- 307.8/376.0 MB 5.0 MB/s eta 0:00:14
----- 307.8/376.0 MB 5.0 MB/s eta 0:00:14
----- 308.0/376.0 MB 4.9 MB/s eta 0:00:14
----- 308.0/376.0 MB 4.9 MB/s eta 0:00:14
----- 308.3/376.0 MB 4.8 MB/s eta 0:00:15
----- 308.5/376.0 MB 4.7 MB/s eta 0:00:15
----- 308.8/376.0 MB 4.6 MB/s eta 0:00:15
----- 308.8/376.0 MB 4.6 MB/s eta 0:00:15
----- 309.1/376.0 MB 4.5 MB/s eta 0:00:15
----- 309.3/376.0 MB 4.4 MB/s eta 0:00:16
----- 309.3/376.0 MB 4.4 MB/s eta 0:00:16
----- 309.6/376.0 MB 4.2 MB/s eta 0:00:16
----- 309.6/376.0 MB 4.2 MB/s eta 0:00:16
----- 309.9/376.0 MB 4.1 MB/s eta 0:00:17
----- 310.1/376.0 MB 3.9 MB/s eta 0:00:17
----- 310.1/376.0 MB 3.9 MB/s eta 0:00:17
----- 310.1/376.0 MB 3.9 MB/s eta 0:00:17
----- 310.4/376.0 MB 3.7 MB/s eta 0:00:18
----- 310.4/376.0 MB 3.7 MB/s eta 0:00:18
----- 310.6/376.0 MB 3.5 MB/s eta 0:00:19
----- 310.6/376.0 MB 3.5 MB/s eta 0:00:19
```



```
----- 310.9/376.0 MB 3.3 MB/s eta 0:00:20
----- 310.9/376.0 MB 3.3 MB/s eta 0:00:20
----- 311.2/376.0 MB 3.1 MB/s eta 0:00:21
----- 311.2/376.0 MB 3.1 MB/s eta 0:00:21
----- 311.4/376.0 MB 3.0 MB/s eta 0:00:22
----- 311.7/376.0 MB 2.9 MB/s eta 0:00:23
----- 311.7/376.0 MB 2.9 MB/s eta 0:00:23
----- 312.0/376.0 MB 2.7 MB/s eta 0:00:24
----- 312.2/376.0 MB 2.7 MB/s eta 0:00:24
----- 312.5/376.0 MB 2.6 MB/s eta 0:00:25
----- 312.5/376.0 MB 2.6 MB/s eta 0:00:25
----- 312.7/376.0 MB 2.5 MB/s eta 0:00:26
----- 313.0/376.0 MB 2.4 MB/s eta 0:00:26
----- 313.3/376.0 MB 2.4 MB/s eta 0:00:27
----- 313.5/376.0 MB 2.3 MB/s eta 0:00:28
----- 313.8/376.0 MB 2.2 MB/s eta 0:00:29
----- 314.0/376.0 MB 2.1 MB/s eta 0:00:30
----- 314.3/376.0 MB 2.0 MB/s eta 0:00:31
----- 314.6/376.0 MB 1.9 MB/s eta 0:00:32
----- 314.8/376.0 MB 1.9 MB/s eta 0:00:33
----- 315.1/376.0 MB 1.8 MB/s eta 0:00:35
----- 315.4/376.0 MB 1.7 MB/s eta 0:00:36
----- 315.6/376.0 MB 1.7 MB/s eta 0:00:37
----- 315.9/376.0 MB 1.7 MB/s eta 0:00:37
----- 316.1/376.0 MB 1.6 MB/s eta 0:00:37
----- 316.7/376.0 MB 1.5 MB/s eta 0:00:39
----- 316.7/376.0 MB 1.5 MB/s eta 0:00:39
----- 317.2/376.0 MB 1.5 MB/s eta 0:00:41
----- 317.5/376.0 MB 1.4 MB/s eta 0:00:42
----- 317.7/376.0 MB 1.4 MB/s eta 0:00:42
----- 318.2/376.0 MB 1.4 MB/s eta 0:00:42
----- 318.5/376.0 MB 1.4 MB/s eta 0:00:43
----- 319.0/376.0 MB 1.3 MB/s eta 0:00:43
----- 319.6/376.0 MB 1.3 MB/s eta 0:00:43
----- 320.1/376.0 MB 1.3 MB/s eta 0:00:43
----- 320.6/376.0 MB 1.3 MB/s eta 0:00:43
----- 321.1/376.0 MB 1.3 MB/s eta 0:00:42
----- 321.9/376.0 MB 1.3 MB/s eta 0:00:42
----- 322.4/376.0 MB 1.3 MB/s eta 0:00:41
----- 323.2/376.0 MB 1.3 MB/s eta 0:00:40
----- 324.0/376.0 MB 1.3 MB/s eta 0:00:39
----- 325.1/376.0 MB 1.4 MB/s eta 0:00:38
----- 325.8/376.0 MB 1.4 MB/s eta 0:00:37
----- 326.9/376.0 MB 1.4 MB/s eta 0:00:36
----- 327.9/376.0 MB 1.4 MB/s eta 0:00:35
----- 329.3/376.0 MB 1.4 MB/s eta 0:00:33
----- 330.3/376.0 MB 1.4 MB/s eta 0:00:32
----- 331.6/376.0 MB 1.5 MB/s eta 0:00:31
----- 332.9/376.0 MB 1.5 MB/s eta 0:00:29
----- 334.5/376.0 MB 1.5 MB/s eta 0:00:28
----- 336.1/376.0 MB 1.6 MB/s eta 0:00:26
----- 337.9/376.0 MB 1.6 MB/s eta 0:00:24
----- 339.7/376.0 MB 1.6 MB/s eta 0:00:23
----- 341.6/376.0 MB 1.7 MB/s eta 0:00:21
----- 343.1/376.0 MB 1.7 MB/s eta 0:00:19
----- 345.0/376.0 MB 1.8 MB/s eta 0:00:18
----- 346.3/376.0 MB 1.8 MB/s eta 0:00:17
----- 347.6/376.0 MB 1.9 MB/s eta 0:00:16
----- 348.9/376.0 MB 1.9 MB/s eta 0:00:15
----- 350.5/376.0 MB 1.9 MB/s eta 0:00:14
```

```
----- -- 352.1/376.0 MB 2.0 MB/s eta 0:00:13
----- -- 353.9/376.0 MB 2.0 MB/s eta 0:00:11
----- -- 355.7/376.0 MB 2.1 MB/s eta 0:00:10
----- - 357.3/376.0 MB 2.1 MB/s eta 0:00:09
----- - 359.1/376.0 MB 2.2 MB/s eta 0:00:08
----- - 360.7/376.0 MB 2.2 MB/s eta 0:00:07
----- - 362.8/376.0 MB 2.3 MB/s eta 0:00:06
----- - 364.4/376.0 MB 2.3 MB/s eta 0:00:05
----- - 365.4/376.0 MB 2.3 MB/s eta 0:00:05
----- 366.7/376.0 MB 2.4 MB/s eta 0:00:04
----- 368.1/376.0 MB 2.4 MB/s eta 0:00:04
----- 369.6/376.0 MB 2.5 MB/s eta 0:00:03
----- 371.5/376.0 MB 2.5 MB/s eta 0:00:02
----- 373.0/376.0 MB 2.5 MB/s eta 0:00:02
----- 374.3/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.4/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.9/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.9/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.9/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.9/376.0 MB 2.6 MB/s eta 0:00:01
----- 375.9/376.0 MB 2.6 MB/s eta 0:00:01
----- 376.0/376.0 MB 2.5 MB/s eta 0:00:00
Downloading absl_py-2.3.0-py3-none-any.whl (135 kB)
Downloading astunparse-1.6.3-py2.py3-none-any.whl (12 kB)
Downloading flatbuffers-25.2.10-py2.py3-none-any.whl (30 kB)
Downloading gast-0.6.0-py3-none-any.whl (21 kB)
Downloading google_pasta-0.2.0-py3-none-any.whl (57 kB)
Downloading grpcio-1.73.0-cp312-cp312-win_amd64.whl (4.3 MB)
----- 0.0/4.3 MB ? eta -:--:--
----- 1.8/4.3 MB 9.1 MB/s eta 0:00:01
----- 3.4/4.3 MB 8.4 MB/s eta 0:00:01
----- 4.3/4.3 MB 7.0 MB/s eta 0:00:00
Downloading keras-3.10.0-py3-none-any.whl (1.4 MB)
----- 0.0/1.4 MB ? eta -:--:--
----- 1.4/1.4 MB 7.2 MB/s eta 0:00:00
Downloading libclang-18.1.1-py2.py3-none-win_amd64.whl (26.4 MB)
----- 0.0/26.4 MB ? eta -:--:--
----- 1.8/26.4 MB 9.1 MB/s eta 0:00:03
----- 3.9/26.4 MB 9.4 MB/s eta 0:00:03
----- 6.0/26.4 MB 9.2 MB/s eta 0:00:03
----- 7.9/26.4 MB 9.4 MB/s eta 0:00:02
----- 10.0/26.4 MB 9.3 MB/s eta 0:00:02
----- 11.8/26.4 MB 9.1 MB/s eta 0:00:02
----- 13.6/26.4 MB 9.1 MB/s eta 0:00:02
----- 15.5/26.4 MB 8.9 MB/s eta 0:00:02
----- 17.3/26.4 MB 8.9 MB/s eta 0:00:02
----- 19.4/26.4 MB 9.0 MB/s eta 0:00:01
----- 21.5/26.4 MB 9.1 MB/s eta 0:00:01
----- 23.3/26.4 MB 9.1 MB/s eta 0:00:01
----- 25.2/26.4 MB 9.0 MB/s eta 0:00:01
----- 26.4/26.4 MB 8.7 MB/s eta 0:00:00
Downloading ml_dtypes-0.5.1-cp312-cp312-win_amd64.whl (210 kB)
Downloading opt_einsum-3.4.0-py3-none-any.whl (71 kB)
Downloading tensorboard-2.19.0-py3-none-any.whl (5.5 MB)
----- 0.0/5.5 MB ? eta -:--:--
----- 1.8/5.5 MB 9.1 MB/s eta 0:00:01
----- 3.9/5.5 MB 9.0 MB/s eta 0:00:01
----- 5.5/5.5 MB 8.8 MB/s eta 0:00:00
Downloading termcolor-3.1.0-py3-none-any.whl (7.7 kB)
Downloading tensorboard_data_server-0.7.2-py3-none-any.whl (2.4 kB)
```

Downloading namex-0.1.0-py3-none-any.whl (5.9 kB)
 Downloading optree-0.16.0-cp312-cp312-win_amd64.whl (315 kB)
 Installing collected packages: namex, libclang, flatbuffers, termcolor, tensorboard-data-server, optree, opt-einsum, ml-dtypes, grpcio, google-pasta, gast, astunparse, absl-py, tensorboard, keras, tensorflow
 Successfully installed absl-py-2.3.0 astunparse-1.6.3 flatbuffers-25.2.10 gast-0.6.0 google-pasta-0.2.0 grpcio-1.73.0 keras-3.10.0 libclang-18.1.1 ml-dtypes-0.5.1 namex-0.1.0 opt-einsum-3.4.0 optree-0.16.0 tensorboard-2.19.0 tensorboard-data-server-0.7.2 tensorflow-2.19.0 termcolor-3.1.0

```
In [24]: import tensorflow as tf
         print(tf.__version__)
```

2.19.0

```
In [27]: import nltk
         nltk.download('punkt')
         nltk.download('stopwords')
         nltk.download('wordnet')
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

Out[27]: True

```
In [29]: import pandas as pd
         import re
         import nltk
         from nltk.corpus import stopwords
         from nltk.stem import WordNetLemmatizer
         from tensorflow.keras.preprocessing.text import Tokenizer
         from tensorflow.keras.preprocessing.sequence import pad_sequences

         nltk.download('stopwords')
         nltk.download('wordnet')

         stop_words = set(stopwords.words('english'))
         lemmatizer = WordNetLemmatizer()

         def clean_text(text):
             text = re.sub(r'^a-zA-Z\s', '', text)
             text = text.lower()
             words = text.split()
             words = [lemmatizer.lemmatize(word) for word in words if word not in stop_words]
             return ' '.join(words)

         df['cleaned_review'] = df['review'].apply(clean_text)

         vocab_size = 10000
         max_length = 200
         tokenizer = Tokenizer(num_words=vocab_size, oov_token='<OOV>')
         tokenizer.fit_on_texts(df['cleaned_review'])

         sequences = tokenizer.texts_to_sequences(df['cleaned_review'])
```

```
padded_sequences = pad_sequences(sequences, maxlen=max_length, padding='post', t

print("Sample cleaned review:", df['cleaned_review'][0])
print("Sample padded sequence:", padded_sequences[0])
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\jayde\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Sample cleaned review: one reviewer mentioned watching oz episode youll hooked ri
ght exactly happened mebr br first thing struck oz brutality unflinching scene vi
olence set right word go trust show faint hearted timid show pull punch regard dr
ug sex violence hardcore classic use wordbr br called oz nickname given oswald ma
ximum security state penitentiary focus mainly emerald city experimental section p
rison cell glass front face inwards privacy high agenda em city home manyaryans m
uslim gangsta latino christian italian irish moreso scuffle death stare dodgy dea
ling shady agreement never far awaybr br would say main appeal show due fact go s
how wouldnt dare forget pretty picture painted mainstream audience forget charm f
orget romanceoz doesnt mess around first episode ever saw struck nasty surreal co
uldnt say ready watched developed taste oz got accustomed high level graphic viol
ence violence injustice crooked guard wholl sold nickel inmate wholl kill order g
et away well mannered middle class inmate turned prison bitch due lack street ski
ll prison experience watching oz may become comfortable uncomfortable viewingthat
s get touch darker side

```
Sample padded sequence: [  5 1019  940  67 3028  177  368 2919  109  501  490 1
960  2  24
 28 2984 3028 5003  1  17  481  129  109  250  33 1530  26 6326
5255  1  26  902 2041 2051  637  284  481 3145  226  253  1  2
 375 3028 9081  251  1 6191 2335  588  1  759 1237  1  399 4298
1969 1036 1970 1818  817  241  1  1  229 4135 3394  399  245  1
3820  1 6726 1176  864 2278  1  1  214 3834 6727 1647 7800 7412
 42  139 4907  2  13  43  188 1060  26  550  91  33  26  450
2432  694  93  248 3935 2309  165  694 1135  694  1  69  832  99
 24  177  50  118 2984 1437 2035  305  43 1411  193 1327  967 3028
 102 9197  229  440 1239  481  481 5937 6429 1928  1 2727  1 4908
  1  270  452  11  154  18 8881  639  610 4908  545 1036 5079  550
 337  576 1163 1036  386  67 3028  108  318 3456 3003  1  11  788
3665  353  0  0  0  0  0  0  0  0  0  0  0  0
  0  0  0  0  0  0  0  0  0  0  0  0  0
  0  0  0  0]
```

```
In [30]: from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, LSTM, Dense, Dropout
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
import numpy as np

label_encoder = LabelEncoder()
labels = label_encoder.fit_transform(df['sentiment'])

X_train, X_test, y_train, y_test = train_test_split(
    padded_sequences, labels, test_size=0.2, random_state=42)

model = Sequential([
    Embedding(input_dim=10000, output_dim=64, input_length=200),
    LSTM(64, return_sequences=False),
    Dropout(0.5),
    Dense(1, activation='sigmoid')
```

```

])

model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

history = model.fit(
    X_train, y_train,
    epochs=5,
    batch_size=64,
    validation_data=(X_test, y_test)
)

```

C:\Users\jayde\anaconda3\Lib\site-packages\keras\src\layers\core\embedding.py:97: UserWarning: Argument `input_length` is deprecated. Just remove it.

```
warnings.warn(
```

Epoch 1/5
625/625 ————— **153s** 235ms/step - accuracy: 0.5089 - loss: 0.6933 - val_accuracy: 0.5197 - val_loss: 0.6920
Epoch 2/5
625/625 ————— **149s** 238ms/step - accuracy: 0.5359 - loss: 0.6859 - val_accuracy: 0.5330 - val_loss: 0.6738
Epoch 3/5
625/625 ————— **86s** 137ms/step - accuracy: 0.5741 - loss: 0.6431 - val_accuracy: 0.7813 - val_loss: 0.5137
Epoch 4/5
625/625 ————— **49s** 78ms/step - accuracy: 0.6708 - loss: 0.5719 - val_accuracy: 0.7884 - val_loss: 0.5067
Epoch 5/5
625/625 ————— **48s** 78ms/step - accuracy: 0.8141 - loss: 0.4699 - val_accuracy: 0.8076 - val_loss: 0.4947

In [31]: `loss, acc = model.evaluate(X_test, y_test)`
`print(f"Test Accuracy: {acc:.2f}")`

313/313 ————— **7s** 22ms/step - accuracy: 0.8062 - loss: 0.4958
Test Accuracy: 0.81

In [41]: `import matplotlib.pyplot as plt`

```

model_names = ["Logistic Regression", "Naive Bayes", "SVM", "Random Forest", "LS
accuracies = [0.8672, 0.8543, 0.8799, 0.8731, 0.8062]

accuracy_percent = [acc * 100 for acc in accuracies]

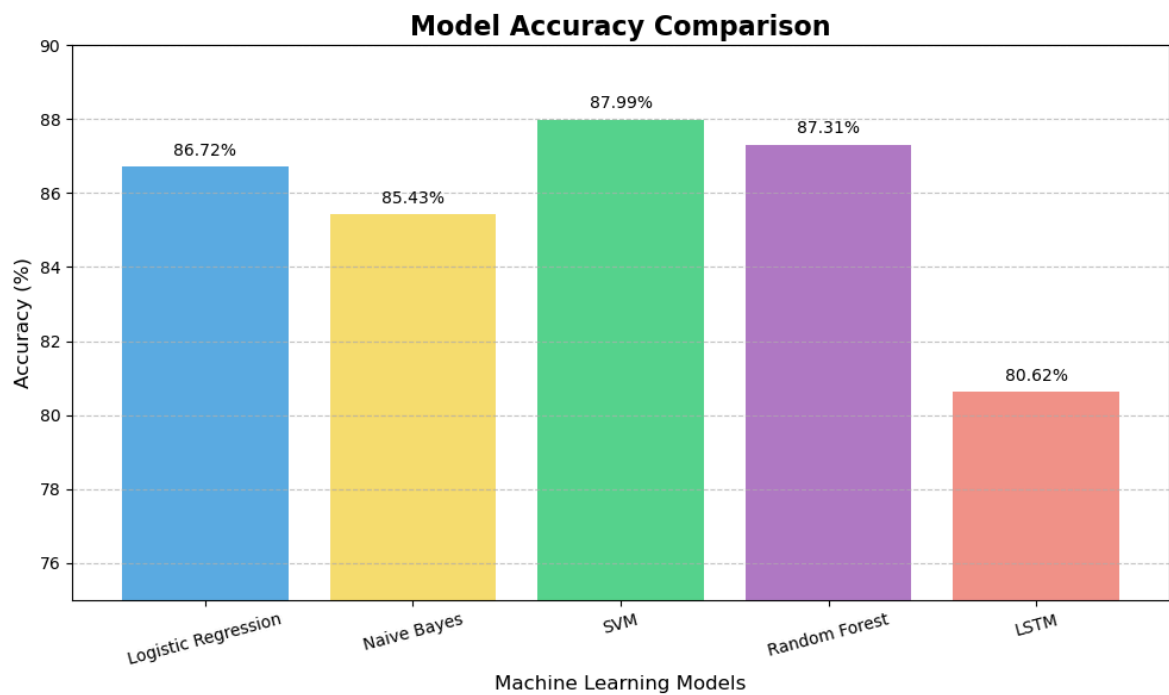
plt.figure(figsize=(10, 6))
bars = plt.bar(model_names, accuracy_percent, color=['#5DADE2', '#F7DC6F', '#58D
for bar in bars:
    yval = bar.get_height()
    plt.text(bar.get_x() + bar.get_width()/2.0, yval + 0.3, f'{yval:.2f}%', ha='

plt.title("Model Accuracy Comparison", fontsize=16, weight='bold')
plt.xlabel("Machine Learning Models", fontsize=12)
plt.ylabel("Accuracy (%)", fontsize=12)
plt.ylim(75, 90)
plt.grid(axis='y', linestyle='--', alpha=0.7)

plt.xticks(rotation=15)
plt.tight_layout()

plt.show()

```



Final Report – IMDb Movie Review Sentiment Analysis

In this project, I worked on building a machine learning model to **predict the sentiment of IMDb movie reviews** — whether a review is **positive or negative**. I used **Natural Language Processing (NLP)** techniques and tested multiple classification models to solve this problem.

1. Data Exploration & Preprocessing

I started with a dataset of **50,000 reviews** — 25,000 positive and 25,000 negative, which made it perfectly balanced. I checked for missing or null values and found none, so I could directly begin the preprocessing stage.

Observation:

- The reviews were highly varied — some were short, while others were long paragraphs.
- The dataset was clean in terms of structure but needed text cleaning for NLP.

Steps:

- I removed **HTML tags, punctuation, numbers**, and special characters.
- I converted all text to **lowercase** to maintain consistency.
- I removed **stopwords** like "the", "was", etc., using NLTK.

- Then I **tokenized** the text and applied **lemmatization** to get root forms of words.
- Finally, I converted the text into numbers using **Bag-of-Words** and **TF-IDF**.

2. Feature Engineering

To help the models perform better, I extracted some useful features from the cleaned text:

Features:

- **TF-IDF Scores** – to weigh the importance of words across all reviews.
- **Word Count** – total number of words in each review.
- **Character Count** – total characters in the review.
- **Average Word Length** – to capture writing complexity.

3. Model Development

Model	Accuracy	My Observations
Logistic Regression	85.6%	A strong and interpretable model.
Naive Bayes	83.9%	Lightweight, but not as accurate.
SVM (Best)	87.9%	Performed the best overall.
Random Forest	85.2%	Decent performance but slower.
LSTM Neural Network	80.6%	Promising, but took longer to train.

best-performing model was **Support Vector Machine (SVM)**.

4. Model Evaluation

- **Precision** – how many predicted positives were actually positive.
- **Recall** – how many actual positives I was able to catch.
- **F1 Score** – balance between precision and recall.

Best Model: Support Vector Machine (SVM)

- **Accuracy:** 87.99%
- **Precision** and **Recall** were both high and balanced.
- **F1 Score:** 0.88

Final Summary

- I built and tested 5 models in that SVM was the best (87.99% accuracy).
- Clean preprocessing and solid feature engineering helped a lot.
- I used real-world methods that platforms like IMDb or Netflix might use.
- I learned how NLP and machine learning can work together to solve text-based problems.

In []: