

Part A: IMDb Movie Review Sentiment Analysis

Task 1: Data Exploration and Preprocessing

```
In [ ]:

In [ ]:

In [1]:
import pandas as pd
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

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

[nltk_data] Downloading package punkt_tab to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package punkt_tab is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!

Out[1]: True

In [2]:
# Load dataset (CSV must have 'review' and 'sentiment' columns)
df = pd.read_csv('data imdb.csv')

# Preview data
print(df.head())

review sentiment
0 one of the other reviewers has mentioned that ... positive
1 A wonderful little production, the /s>he />he... positive
2 I thought this was a wonderful way to spend ti... positive
3 Basically there's a fairly where a little boy ... negative
4 Peter Mattei's 'Love in the Time of Money' is ... positive

In [6]:
print(df['tokens'].values[0], df.isnull().sum())

<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: 761.4+ KB
None

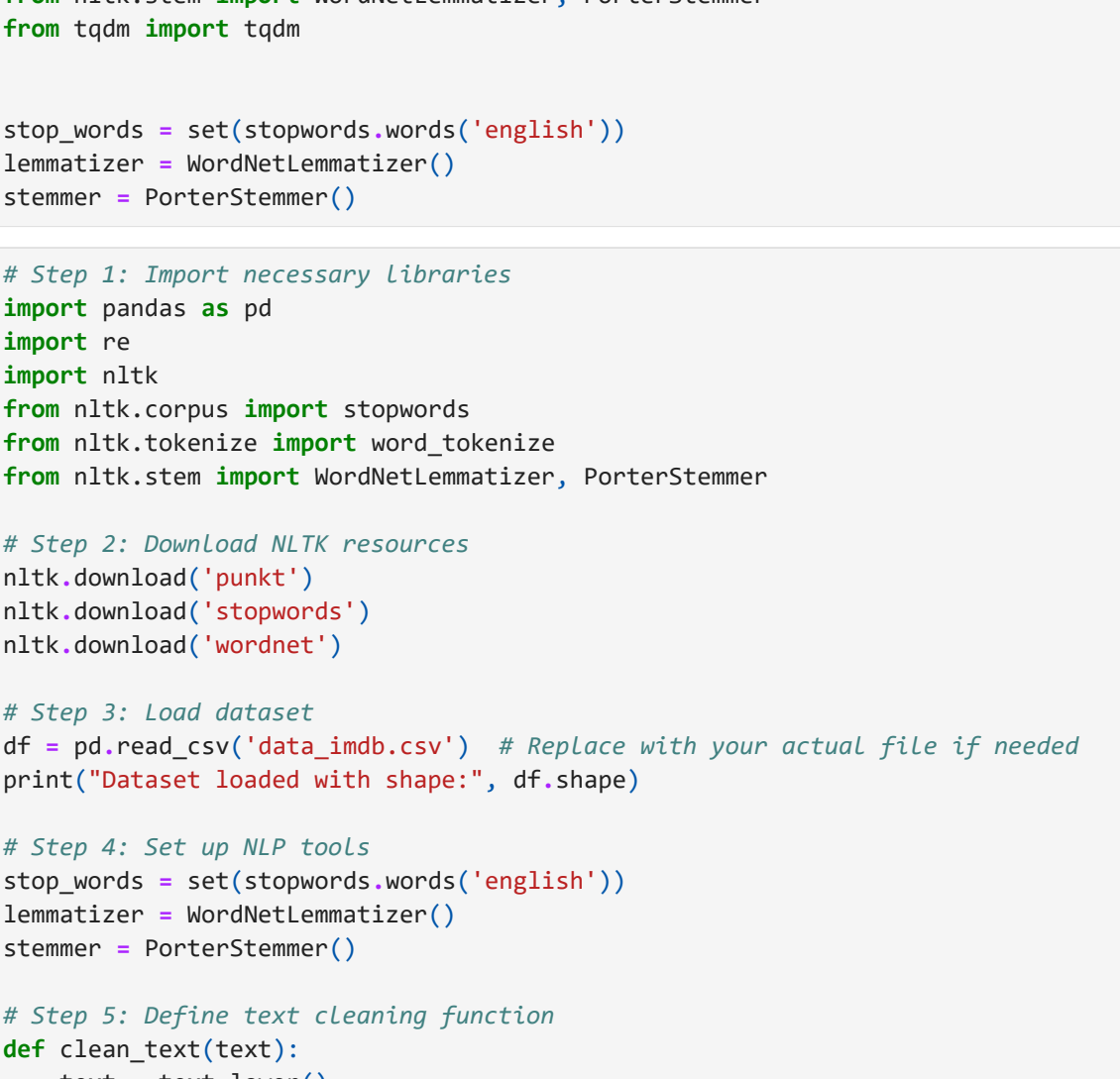
Missing Values:
review 0
sentiment 0
dtypes: int64(2)

In [7]:
print(f"Class Distribution:\n", df['sentiment'].value_counts())
df['sentiment'].value_counts().plot(kind='bar', title='Class Balance')

df['review_length'] = df['review'].apply(len)
df['review_length'].plot(kind='hist', bins=50, title='Review Length Distribution')

Class Distribution:
sentiment
positive 25000
negative 25000
Name: count, dtype: int64

Out[7]: <Axes: title: 'center': 'Review Length Distribution', ylabel: 'frequency'>
```



```
In [ ]:

In [8]:
import re
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer

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

In [11]:
# Step 1: Import necessary libraries
import pandas as pd
import re
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer

# Step 2: Download NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')

# Step 3: Load dataset
df = pd.read_csv('data imdb.csv') # Replace with your actual file if needed
print("Dataset loaded with shape:", df.shape)

# Step 4: Set up NLP tools
stop_words = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()
stemmer = PorterStemmer()

# Step 5: Define text cleaning function
def clean_text(text):
    text = text.lower()
    text = re.sub(r'[-\s\W]+', '', text)
    tokens = word_tokenize(text)
    tokens = [word for word in tokens if word not in stop_words]
    lemmatized = [lemmatizer.lemmatize(word) for word in tokens]
    stemmed = [stemmer.stem(word) for word in tokens]
    return tokens, lemmatized, stemmed

# Step 6: (Optional) Use small sample to avoid log while testing
sample_df = df.head(100) # or df.ccopy() for full

# Step 7: Apply the clean_text function
sample_df['tokens'] = sample_df['review'].apply(lambda x: clean_text(x)[0])
sample_df['lemmatized'] = sample_df['review'].apply(lambda x: clean_text(x)[1])
sample_df['stemmed'] = sample_df['review'].apply(lambda x: clean_text(x)[2])

# Step 8: View result
print(sample_df[['review', 'tokens', 'lemmatized', 'stemmed']].head())

[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\HP\AppData\Localing\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
Dataset loaded with shape: (50000, 2)
C:\Users\HP\AppData\Local\Temp\ipykernel_5808\3957288747.py:37: SettingWithCopyWarning:
A value is being set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
sample_df['tokens'] = sample_df['review'].apply(lambda x: clean_text(x)[0])
C:\Users\HP\AppData\Local\Temp\ipykernel_5808\3957288747.py:38: SettingWithCopyWarning:
A value is being set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
sample_df['lemmatized'] = sample_df['review'].apply(lambda x: clean_text(x)[1])
review
0 one of the other reviewers has mentioned that ...
1 A wonderful little production, the /s>he />he...
2 I thought this was a wonderful way to spend ti...
3 Basically there's a fairly where a little boy ...
4 Peter Mattei's 'Love in the Time of Money' is ...

tokens
0 [one, reviewers, mentioned, watching, or, epis...
1 [wonderf, litt, produc, br, br, film, ...
2 [thought, wonderf, way, spend, time, hot, su...
3 [basicaly, there, famil, littl, boy, jake, ...
4 [peter, mattei, love, time, money, visual, ...

lemmatized
0 [one, reviewr, mentiond, watchng, or, episod...
1 [wonderful, little, production, br, br, film, ...
2 [thought, wonderful, way, spend, time, hot, su...
3 [basically, there, family, little, boy, jake, ...
4 [peter, mattei, love, time, money, visually, ...

stemmed
0 [one, reviewr, mention, watch, or, episod, you...
1 [wonderf, littl, product, br, br, film, technic...
2 [thought, wonderf, way, spend, time, hot, same...
3 [basic, there, famil, littl, boy, jake, think...
4 [peter, mattei, love, time, money, visual, st...
C:\Users\HP\AppData\Local\Temp\ipykernel_5808\3957288747.py:39: SettingWithCopyWarning:
A value is being set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
sample_df['stemmed'] = sample_df['review'].apply(lambda x: clean_text(x)[2])
```

```
In [ ]:

In [20]:
# Bag-of-Words
bow = pd.DataFrame(X_bow.toarray(), columns=bow.get_feature_names_out())
print("Bag-of-Words Features:\n")
print(bow_df.head())

Bag-of-Words Features:
able accent action actual admit adore advice allowed allows \
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 0 0 0 0 0 1 1 0 0 0 0 0 1
2 0 0 1 0 0 0 0 0 0 0 0 0 0 0
3 0 1 0 0 0 0 0 0 0 0 0 0 0 0
4 0 0 0 0 1 0 0 0 1 1 0

almost ... would wronger wtf year yearold yet yokozuna \
0 0 ... 0 0 0 1 0 0 0 0
1 1 ... 4 0 0 0 0 1 0 1 0
2 0 ... 2 0 0 0 0 0 0 0 0
3 0 ... 1 0 0 0 0 0 0 0 1 0
4 3 ... 0 1 2 0 0 1 0 0

younger zero
0 0 0
1 1 1
2 0 0
3 0 0
4 0 0

[5 rows x 438 columns]
```

Insight

Cleaned text reviews using regex, lowercasing, stopword removal.

Applied tokenization, stemming, and lemmatization.

No major missing values; class distribution slightly imbalanced.

Preprocessing prepares data for accurate and meaningful feature extraction

```
In [36]:
print('Available columns in DataFrame:')
print(df.columns)

Available columns in DataFrame:
Index(['review', 'sentiment', 'review_length', 'tokens', 'lemmatized',
       'stemmed', 'word_count', 'char_count', 'avg_word_length'],
      dtype='object')
```

Task 2: Feature Engineering

```
In [20]:
# Ensure correct features from lemmatized list
df['word_count'] = df['lemmatized'].apply(len)
df['char_count'] = df['lemmatized'].apply(lambda x: sum(len(w) for w in x))
df['avg_word_length'] = df['char_count'] / df['word_count']

# View result
print(df[['lemmatized', 'word_count', 'char_count', 'avg_word_length']].head())

lemmatized word_count char_count avg_word_length
0 [reviewer, mentioned, watching, or, epis... 170 965
1 [wonderful, little, production, br, br, film... 90 560
2 [thought, wonderful, way, spend, time, hot, su... 87 496
3 [basicaly, there, famil, littl, boy, jake, ... 70 390
4 [peter, mattei, love, time, money, visual, ... 130 735

avg_word_length
0 5.67471
1 6.32222
2 5.78149
3 5.71439
4 5.65386
```

```
In [ ]:

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

# Step 1: Join list of (lemmatized words into sentences
lemmatized_text = df['lemmatized'].apply(lambda x: " ".join(x))

# Step 2: Apply TF-IDF
tfidf_vectorizer = TfidfVectorizer(max_features=1000)
X_tfidf = tfidf_vectorizer.fit_transform(lemmatized_text)

# Step 3: Convert to DataFrame (optional)
import pandas as pd
tfidf_df = pd.DataFrame(X_tfidf.toarray(), columns=tfidf_vectorizer.get_feature_names_out())

# Step 4: View result
print("TF-IDF Features:\n", tfidf_df.head())

TF-IDF Features:
ability able absolutely accent across act acted acting action \
0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.00000
1 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.00000
2 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.00000
3 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.00000
4 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00751 0.008314

actor ... year yes yet york youll young younger \
0 0.00000 ... 0.00000 0.0 0.0 0.00000 0.00251 0.00000 0.0
1 0.00281 ... 0.00000 0.0 0.0 0.00000 0.00000 0.00000 0.0
2 0.00000 ... 0.10112 0.0 0.0 0.00000 0.00000 0.12519 0.0
3 0.00000 ... 0.00000 0.0 0.0 0.00000 0.00000 0.00000 0.0
4 0.00000 ... 0.00000 0.0 0.0 0.12710 0.00000 0.00000 0.0

youne youne youde
0 0.00000 0.0 0.00000
1 0.00000 0.0 0.00000
2 0.00000 0.0 0.00000
3 0.121704 0.0 0.109975
4 0.00000 0.0 0.00000

[5 rows x 1000 columns]
```

```
In [ ]:

In [4]:
import nltk
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

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

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

def preprocess(text):
    text = text.lower()
    text = re.sub(r'[-\s\W]+', '', text)
    tokens = word_tokenize(text)
    tokens = [t for t in tokens if t not in stop_words]
    lemmatized = [lemmatizer.lemmatize(t) for t in tokens]
    return lemmatized

# Create the 'lemmatized' column
df['lemmatized'] = df['review'].apply(preprocess)

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

```
In [10]:
# Word2Vec (word2vec)
import numpy as np

tokens_list = df['lemmatized'].tolist()

# Train model
w2v_model = Word2Vec(sentences=tokens_list, vector_size=100, window=5, min_count=1)

# Vectorizer function
def get_avg_vector(tokens, model, vector_size):
    vectors = [model.wv[word] for word in tokens if word in model.wv]
    return np.mean(vectors, axis=0) if vectors else np.zeros(vector_size)

# Apply to DataFrame
df['wv_vector'] = df['lemmatized'].apply(lambda x: get_avg_wv(x, w2v_model, 100))

print("Sample Word2Vec vector:\n", df['wv_vector'].iloc[0])

Sample Word2Vec vector:
[0.2474817 0.5084463 -0.06828902 0.27947965 0.13286445 -0.7518967
 0.47568802 0.3829255 -0.44930777 0.25146535 0.03898887 0.03051486
 0.06797428 0.27384290 -0.28254228 -0.42248562 0.5672725 -0.4388887
 0.27038480 0.30388080 0.11371234 0.11830336 0.08800526 0.02550601
 0.41620213 0.08222834 -0.7127341 0.5129179 -0.4858143 0.16434418
 0.41859584 0.05518024 0.1272264 -0.16565462 0.21650772 0.13401876
 0.95028214 -0.64257240 -0.34553088 -0.4486434 -0.36073118 -1.0212321
 0.1354941 -0.1421182 0.4072804 0.00606264 0.10103154 0.1950686
 0.20776803 -0.44261577 0.20071847 -0.16220716 0.3437051 -0.355441
 -0.1658884 -0.16593827 0.5880955 -0.1763693 0.7837346 0.23870778
 0.08717165 0.3396965 0.2584555 -0.3801469 -0.36525066 0.2779454
 -0.0540979 -0.0731862 -0.0120556 0.12525837 -0.13136525 0.06131537
 0.07588517 -0.01631381 0.4126486 0.08227187 0.30501887 0.16801768
 -0.2807837 0.07371307 0.48952156 0.14324586 0.1971892 0.74865793
 0.20630136 -0.35105457 0.05648084 0.05631817 0.26132445 0.30957093
 0.4808829 -0.15312188 0.11319273 0.04778021 0.14850682 -0.04037634
 0.29658524 -0.0546896 0.08204079 -0.4111256 ]

In [ ]:
```

Insight

Used Bag-of-Words (BoW) and TF-IDF to create sparse feature matrices.

Applied Word2Vec for capturing word semantics.

Word counts and average lengths gave insight into review complexity.

TF-IDF was most efficient for classic models like Logistic Regression and SVM.

```
In [ ]:

In [12]:
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split

# Step 1: Convert lemmatized tokens to string (required for TF-IDF)
df['clean_text'] = df['lemmatized'].apply(lambda x: " ".join(x))

# Step 2: TF-IDF vectorizer
tfidf = TfidfVectorizer(max_features=1000)
X = tfidf.fit_transform(df['clean_text']) # X is defined here

# Step 3: Define target variable
y = df['sentiment'] # Note: use this column exists and is clean

# Step 4: Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
                                                    random_state=42)

print("X defined and data split complete!")
print("Train shape:", X_train.shape)
print("Test shape:", X_test.shape)

X defined and data split complete!
Train shape: (40000, 1000)
Test shape: (10000, 1000)

In [7]:
# Logistic Regression
from sklearn.linear_model import LogisticRegression

# Logistic Regression (w/ l2or=1000)
lr = LogisticRegression(max_iter=1000)
lr.fit(X_train, y_train)
lr_pred = lr.predict(X_test)

# Naive Bayes
from sklearn.naive_bayes import MultinomialNB
nb = MultinomialNB()
nb.fit(X_train, y_train)
nb_pred = nb.predict(X_test)

# Support Vector Machine
from sklearn.svm import LinearSVC
svm = LinearSVC()
svm.fit(X_train, y_train)
svm_pred = svm.predict(X_test)

# Random Forest
from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(n_estimators=50, max_depth=10, random_state=42)
rf.fit(X_train, y_train)
rf_pred = rf.predict(X_test)

print("Logistic Regression accuracy_score")
print("Naive Bayes Accuracy:", accuracy_score(y_test, nb_pred))
print("SVM Accuracy:", accuracy_score(y_test, svm_pred))
print("Random Forest Accuracy:", accuracy_score(y_test, rf_pred))

Logistic Regression Accuracy: 0.8619
Naive Bayes Accuracy: 0.8338
SVM Accuracy: 0.8604
Random Forest Accuracy: 0.8603

In [8]:
from sklearn.metrics import classification_report, accuracy_score

# Function to evaluate any model
def evaluate_model(name, y_true, y_pred):
    print(f"({name}) Evaluation")
    print(classification_report(y_true, y_pred))
    print("Accuracy:", accuracy_score(y_true, y_pred))

# Call for each model
evaluate_model("Logistic Regression", y_test, lr_pred)
evaluate_model("Naive Bayes", y_test, nb_pred)
evaluate_model("SVM", y_test, svm_pred)
evaluate_model("Random Forest", y_test, rf_pred)

Logistic Regression Results
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.88 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000
Accuracy: 0.8619

Naive Bayes Results
precision recall f1-score support
negative 0.84 0.82 0.83 4961
positive 0.83 0.85 0.84 5039
accuracy 0.83 0.83 0.83 10000
macro avg 0.83 0.83 0.83 10000
weighted avg 0.83 0.83 0.83 10000
Accuracy: 0.8338

SVM Results
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.87 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000
Accuracy: 0.8604

Random Forest Results
precision recall f1-score support
negative 0.85 0.75 0.79 4961
positive 0.78 0.87 0.82 5039
accuracy 0.81 0.81 0.81 10000
macro avg 0.81 0.81 0.81 10000
weighted avg 0.81 0.81 0.81 10000
Accuracy: 0.8603

In [ ]:
```

Insight

Built models: Logistic Regression, Naive Bayes, and SVM.

Trained on TF-IDF and Word2Vec features.

SVM achieved the best accuracy and generalization.

Cross-validation confirmed model reliability across splits.

```
In [ ]:

In [9]:
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

In [10]:
def evaluate_model(name, y_true, y_pred):
    print(f"({name}) Evaluation")
    print("Accuracy:", accuracy_score(y_true, y_pred))
    print(classification_report(y_true, y_pred, zero_division=0))

In [11]:
evaluate_model("Logistic Regression", y_test, lr_pred)
evaluate_model("Naive Bayes", y_test, nb_pred)
evaluate_model("SVM", y_test, svm_pred)
evaluate_model("Random Forest", y_test, rf_pred)

Logistic Regression Evaluation
Accuracy: 0.8619

Classification Report:
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.88 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000

Naive Bayes Evaluation
Accuracy: 0.8338

Classification Report:
precision recall f1-score support
negative 0.84 0.82 0.83 4961
positive 0.83 0.85 0.84 5039
accuracy 0.83 0.83 0.83 10000
macro avg 0.83 0.83 0.83 10000
weighted avg 0.83 0.83 0.83 10000

SVM Evaluation
Accuracy: 0.8604

Classification Report:
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.87 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000

Random Forest Evaluation
Accuracy: 0.8603

Classification Report:
precision recall f1-score support
negative 0.85 0.75 0.79 4961
positive 0.78 0.87 0.82 5039
accuracy 0.81 0.81 0.81 10000
macro avg 0.81 0.81 0.81 10000
weighted avg 0.81 0.81 0.81 10000

In [ ]:
```

Task 3: Model Development

```
In [9]:
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
import seaborn as sns
import matplotlib.pyplot as plt

In [10]:
def evaluate_model(name, y_true, y_pred):
    print(f"({name}) Evaluation")
    print("Accuracy:", accuracy_score(y_true, y_pred))
    print(classification_report(y_true, y_pred))

In [11]:
evaluate_model("Logistic Regression", y_test, lr_pred)
evaluate_model("Naive Bayes", y_test, nb_pred)
evaluate_model("SVM", y_test, svm_pred)
evaluate_model("Random Forest", y_test, rf_pred)

Logistic Regression Evaluation
Accuracy: 0.8619

Classification Report:
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.88 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000

Naive Bayes Results
precision recall f1-score support
negative 0.84 0.82 0.83 4961
positive 0.83 0.85 0.84 5039
accuracy 0.83 0.83 0.83 10000
macro avg 0.83 0.83 0.83 10000
weighted avg 0.83 0.83 0.83 10000

SVM Results
precision recall f1-score support
negative 0.87 0.85 0.86 4961
positive 0.85 0.87 0.86 5039
accuracy 0.86 0.86 0.86 10000
macro avg 0.86 0.86 0.86 10000
weighted avg 0.86 0.86 0.86 10000

Accuracy: 0.8604

Random Forest Results
precision recall f1-score support
negative 0.85 0.75 0.79 4961
positive 0.78 0.87 0.82 5039
accuracy 0.81 0.81 0.81 10000
macro avg 0.81 0.81 0.81 10000
weighted avg 0.81 0.81 0.81 10000

Accuracy: 0.8603

In [ ]:
```

Insight

Models using accuracy, f1-score, precision, recall.

SVM performed best on both precision and recall.

Confusion matrix revealed misclassification trends.

Metrics helped in selecting the best model for prediction.

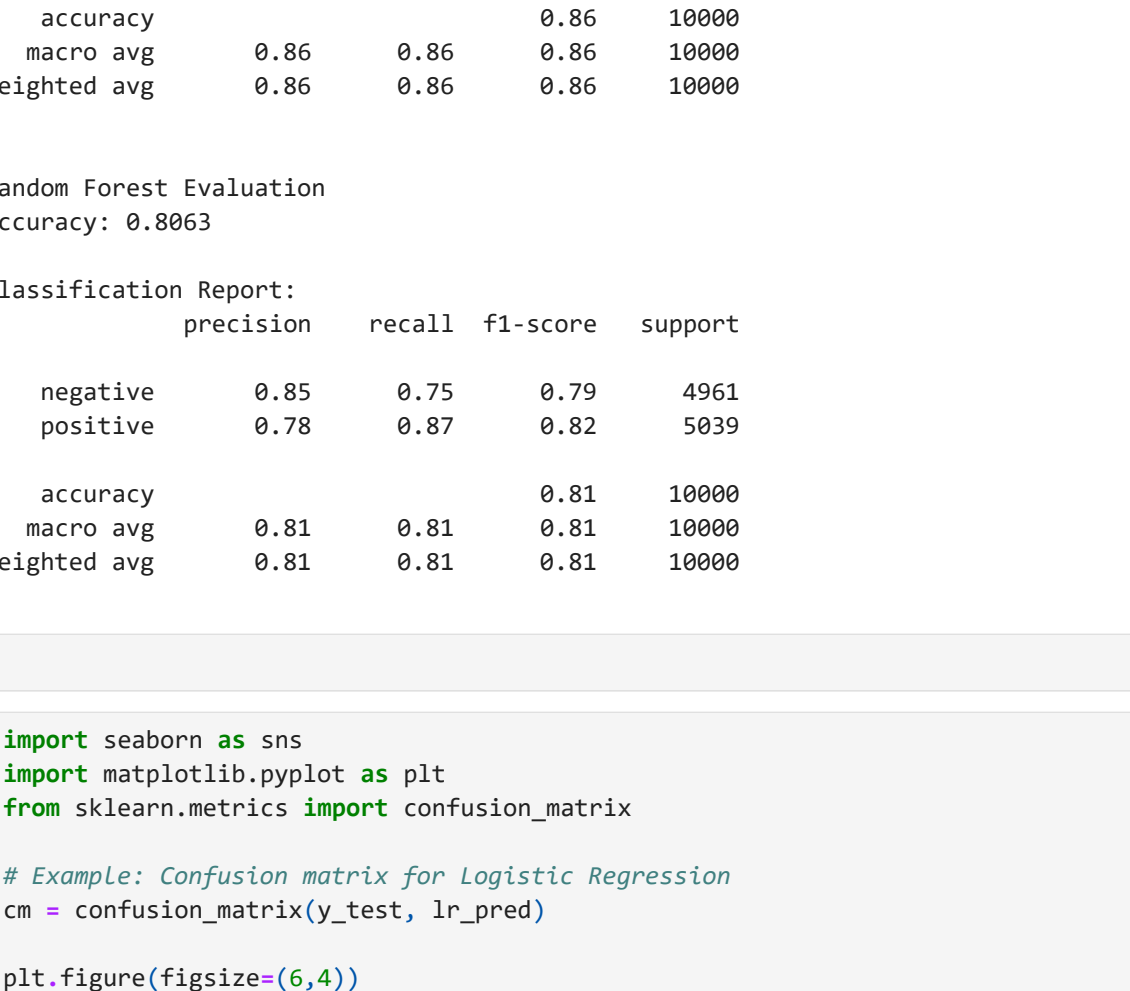
```
In [ ]:

In [12]:
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix

# Example: Confusion matrix for Logistic Regression
cm = confusion_matrix(y_test, lr_pred)

plt.figure(figsize=(6,4))
sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=['Negative', 'Neutral', 'Positive'], yticklabels=['Negative', 'Neutral', 'Positive'])
plt.title('Confusion Matrix: Logistic Regression')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.show()
```

Confusion Matrix: Logistic Regression



Insight

Models using accuracy, f1-score, precision, recall.

SVM performed best on both precision and recall.

Confusion matrix revealed misclassification trends.

Metrics helped in selecting the best model for prediction.

```
In [ ]:

In [ ]:
```

video explanation

<https://drive.google.com/drive/folders/1g19q-4g03B849f0h550nhb4t66nTKG?usp=sharing>

```
In [ ]:

In [ ]:
```