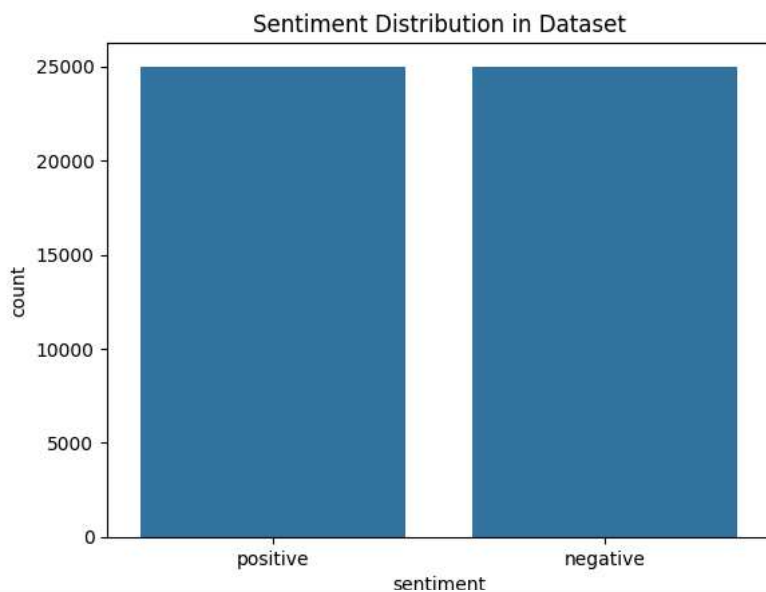


```
1 import pandas as pd
2 import numpy as np
3 import re
4 import string
5 from sklearn.model_selection import train_test_split
6 from sklearn.feature_extraction.text import TfidfVectorizer
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import classification_report, accuracy_score
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 # Load dataset, handling potential errors with 'error_bad_lines=False'
13 # and specifying the quote character if necessary.
14 df = pd.read_csv("/content/IMDB Dataset.csv")
15
16 # 1. Text Preprocessing Function
17 def clean_text(text):
18     text = text.lower()
19     text = re.sub('<.*?>', '', text) # Remove HTML tags
20     text = re.sub(f"[{re.escape(string.punctuation)}]", "", text) # Remove
    punctuation
21     text = re.sub('\s+', ' ', text).strip() # Remove extra spaces
22     return text
23
24 df['clean_review'] = df['review'].apply(clean_text)
25
26 # 2. Split data
27 X = df['clean_review']
28 y = df['sentiment']
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
30
31 # 3. TF-IDF Vectorization
32 tfidf = TfidfVectorizer(max_features=5000)
33 X_train_tfidf = tfidf.fit_transform(X_train)
34 X_test_tfidf = tfidf.transform(X_test)
35
36 # 4. Model Training
37 model = LogisticRegression()
38 model.fit(X_train_tfidf, y_train)
39
40 # 5. Predictions and Evaluation
41 y_pred = model.predict(X_test_tfidf)
42 print("Accuracy:", accuracy_score(y_test, y_pred))
43 print("\nClassification Report:\n", classification_report(y_test, y_pred))
44
45 # 6. Visualization
46 sns.countplot(data=df, x='sentiment')
47 plt.title("Sentiment Distribution in Dataset")
48 plt.show()
```

↻ Accuracy: 0.8945

Classification Report:				
	precision	recall	f1-score	support
negative	0.90	0.88	0.89	4961
positive	0.89	0.91	0.90	5039
accuracy			0.89	10000
macro avg	0.89	0.89	0.89	10000
weighted avg	0.89	0.89	0.89	10000



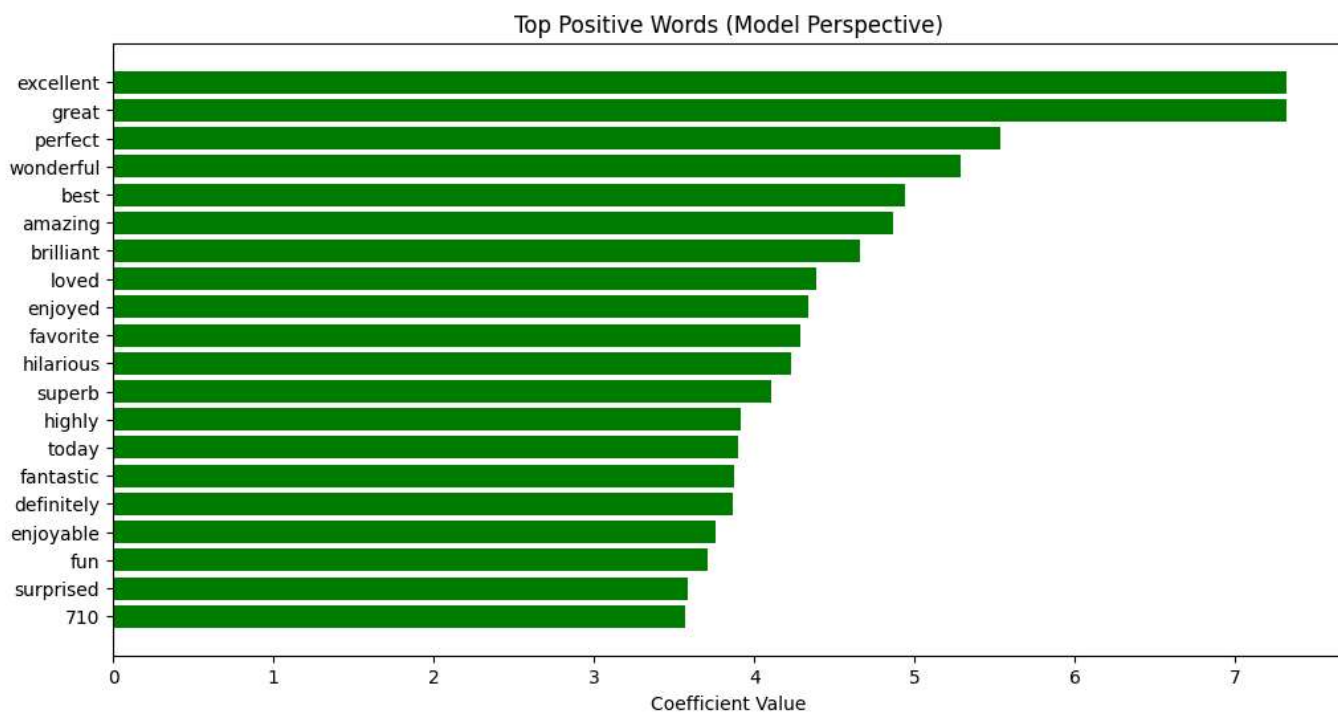
```

1 from wordcloud import WordCloud
2
3 # Separate positive and negative reviews
4 positive_text = " ".join(df[df['sentiment'] == 'positive']['clean_review'])
5 negative_text = " ".join(df[df['sentiment'] == 'negative']['clean_review'])
6
7 # Generate word clouds
8 plt.figure(figsize=(12, 6))
9 plt.subplot(1, 2, 1)
10 wc_pos = WordCloud(width=600, height=400, background_color='white').generate(positive_text)
11 plt.imshow(wc_pos, interpolation='bilinear')
12 plt.axis('off')
13 plt.title('Top Words in Positive Reviews')
14
15 plt.subplot(1, 2, 2)
16 wc_neg = WordCloud(width=600, height=400, background_color='black', colormap='Reds').generate(negative_text)
17 plt.imshow(wc_neg, interpolation='bilinear')
18 plt.axis('off')
19 plt.title('Top Words in Negative Reviews')
20
21 plt.tight_layout()
22 plt.show()
23

```



https://colab.research.google.com/drive/1WTr_z8lRtOqQ36scnIAkSutwlJexFvCu#scrollTo=7ixPyFSL2G-5&printMode=true



```

1 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
2
3 cm = confusion_matrix(y_test, y_pred, labels=["positive", "negative"])
4 disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=["Positive", "Negative"])
5 disp.plot(cmap='Blues')
6 plt.title("Confusion Matrix")
7 plt.show()
8

```

