

## ✓ Google Colab Lab Assignment -NLP

**Course Name:** Deep Learning (PEC)

**Lab Title:** NLP Techniques for Text Classification

**Student Name:** Nabil Ansari

**Student ID:** 202302040004

**Date of Submission:** 09/04/2025

**Group Members:** 1) Arya Sadalage

2) Gourav Sable

3) Nabil Ansari

**Dataset Link :** <https://www.kaggle.com/datasets/lakshmi25npathi/imdb-dataset-of-50k-movie-reviews>

**Github Link :** <https://github.com/nabil-repo/DL/tree/main/Assignment-4>

**Objective** The objective of this assignment is to implement NLP preprocessing techniques and build a text classification model using machine learning techniques.

### **Learning Outcomes:**

1. Understand and apply NLP preprocessing techniques such as tokenization, stopword removal, stemming, and lemmatization.
2. Implement text vectorization techniques such as TF-IDF and CountVectorizer.
3. Develop a text classification model using a machine learning algorithm.
4. Evaluate the performance of the model using suitable metrics.

## ✓ Assignment Instructions:

### **Part 1: NLP Preprocessing**

#### **Dataset Selection:**

Choose any text dataset from **Best Datasets for Text** <https://en.innovatiana.com/post/best-datasets-for-text-classification> Classification, such as SMS Spam Collection, IMDb Reviews, or any other relevant dataset.

Download the dataset and upload it to Google Colab.

Load the dataset into a Pandas DataFrame and explore its structure (e.g., check missing values, data types, and label distribution).

Text Preprocessing:

Convert text to lowercase.

Perform tokenization using NLTK or spaCy.

Remove stopwords using NLTK or spaCy.

Apply stemming using PorterStemmer or SnowballStemmer.

Apply lemmatization using WordNetLemmatizer.

Vectorization Techniques:

Convert text data into numerical format using TF-IDF and CountVectorizer.

```
1 #Code for Part 1
2 # Import required libraries
3 import pandas as pd
4 import numpy as np
5 import re
6 import nltk
7 from nltk.corpus import stopwords
8 from nltk.stem import WordNetLemmatizer
9 from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
10 from tqdm import tqdm
11
12 # Download NLTK resources
13 nltk.download('stopwords')
14 nltk.download('wordnet')
15 nltk.download('omw-1.4')
16
17 # Load IMDB dataset
18 try:
19     df = pd.read_csv('IMDB Dataset.csv')
20     print("Dataset loaded successfully!")
21     print(f"Shape: {df.shape}")
22     print("\nLabel distribution:")
23     print(df['sentiment'].value_counts())
24 except FileNotFoundError:
25     print("Error: IMDB Dataset.csv not found. Please upload the file to Colab.")
26     raise
27
28 # Clean HTML tags and special characters
29 def clean_text(text):
30     text = re.sub(r'<[^\>]+>', '', text) # Remove HTML tags
31     text = re.sub(r'[\W\s]', '', text) # Remove punctuation
```

```
32     text = re.sub(r'\d+', '', text)      # Remove numbers
33     text = re.sub(r'\s+', ' ', text).strip() # Remove extra whitespace
34     return text.lower()
35
36 # Initialize lemmatizer
37 lemmatizer = WordNetLemmatizer()
38
39 # Get English stopwords
40 stop_words = set(stopwords.words('english'))
41
42 # Enhanced preprocessing function
43 def preprocess_text(text):
44     try:
45         # Clean text
46         text = clean_text(text)
47
48         # Tokenization
49         tokens = text.split()
50
51         # Remove stopwords and short words
52         tokens = [word for word in tokens if word not in stop_words and len(word) > 2]
53
54         # Lemmatization
55         tokens = [lemmatizer.lemmatize(word) for word in tokens]
56
57         return ' '.join(tokens)
58     except Exception as e:
59         print(f"Error processing text: {str(e)[:100]}...")
60         return ""
61
62 # Apply preprocessing with progress bar
63 print("\nPreprocessing reviews...")
64 tqdm.pandas()
65 df['processed_review'] = df['review'].progress_apply(preprocess_text)
66
67 # Remove empty processed reviews
68 initial_count = len(df)
69 df = df[df['processed_review'].str.strip() != '']
70 print(f"\nRemoved {initial_count - len(df)} empty reviews after preprocessing")
71
72 # Vectorization
73 print("\nApplying vectorization techniques...")
74
75 # TF-IDF Vectorizer
76 tfidf_vectorizer = TfidfVectorizer(
77     max_features=5000,
78     stop_words=list(stop_words),
79     ngram_range=(1, 2) # Include unigrams and bigrams
80 )
81 X_tfidf = tfidf_vectorizer.fit_transform(df['processed_review'])
82
```

```

83 # Count Vectorizer
84 count_vectorizer = CountVectorizer(
85     max_features=5000,
86     stop_words=list(stop_words),
87     ngram_range=(1, 2)
88 )
89 X_count = count_vectorizer.fit_transform(df['processed_review'])
90
91 # Prepare labels (1 for positive, 0 for negative)
92 y = df['sentiment'].map({'positive': 1, 'negative': 0})
93
94 print("\nPreprocessing completed!")
95 print("TF-IDF shape:", X_tfidf.shape)
96 print("CountVectorizer shape:", X_count.shape)
97 print("\nSample processed review:")
98 print(df['processed_review'].iloc[0][:200], "...")

```



```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Package omw-1.4 is already up-to-date!
Dataset loaded successfully!
Shape: (50000, 2)

```

```

Label distribution:
sentiment
positive    25000
negative    25000
Name: count, dtype: int64

```

```

Preprocessing reviews...
100%|██████████| 50000/50000 [00:29<00:00, 1692.44it/s]

```

```

Removed 0 empty reviews after preprocessing

```

```

Applying vectorization techniques...

```

```

Preprocessing completed!
TF-IDF shape: (50000, 5000)
CountVectorizer shape: (50000, 5000)

```

```

Sample processed review:
one reviewer mentioned watching episode youll hooked right exactly happened methe first

```



## Splitting the Data:

Divide the dataset into training and testing sets (e.g., 80% training, 20% testing).

## Building the Classification Model:

Train a text classification model using Logistic Regression, Naïve Bayes, or any other suitable algorithm.

Implement the model using scikit-learn.

### Model Evaluation:

Evaluate the model using accuracy, precision, recall, and F1-score.

Use a confusion matrix to visualize the results.

```
1 #code for Part 2
2
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.naive_bayes import MultinomialNB
6 from sklearn.metrics import (accuracy_score, precision_score,
7                               recall_score, f1_score, confusion_matrix,
8                               classification_report)
9 import matplotlib.pyplot as plt
10 import seaborn as sns
11
12 # Split data (using TF-IDF features)
13 X_train, X_test, y_train, y_test = train_test_split(
14     X_tfidf, y, test_size=0.2, random_state=42)
15
16 print(f"Training set: {X_train.shape[0]} samples")
17 print(f"Test set: {X_test.shape[0]} samples")
18
19 # Function to evaluate and visualize model performance
20 def evaluate_model(model, X_test, y_test, model_name):
21     y_pred = model.predict(X_test)
22
23     # Calculate metrics
24     accuracy = accuracy_score(y_test, y_pred)
25     precision = precision_score(y_test, y_pred)
26     recall = recall_score(y_test, y_pred)
27     f1 = f1_score(y_test, y_pred)
28
29     print(f"\n{model_name} Performance:")
30     print(classification_report(y_test, y_pred))
31     print(f"Accuracy: {accuracy:.4f}")
32     print(f"Precision: {precision:.4f}")
33     print(f"Recall: {recall:.4f}")
34     print(f"F1 Score: {f1:.4f}")
35
36     # Confusion matrix
37     cm = confusion_matrix(y_test, y_pred)
38     plt.figure(figsize=(6, 4))
39     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
40                 xticklabels=['Negative', 'Positive'],
```

```

41         yticklabels=['Negative', 'Positive'])
42     plt.title(f'{model_name} Confusion Matrix')
43     plt.xlabel('Predicted')
44     plt.ylabel('Actual')
45     plt.show()
46
47     return {'accuracy': accuracy, 'precision': precision, 'recall': recall, 'f1': f1}
48
49 # Logistic Regression Model
50 print("\nTraining Logistic Regression model...")
51 lr_model = LogisticRegression(
52     max_iter=1000,
53     random_state=42,
54     class_weight='balanced' # Handle class imbalance
55 )
56 lr_model.fit(X_train, y_train)
57 lr_metrics = evaluate_model(lr_model, X_test, y_test, "Logistic Regression")
58
59 # Naive Bayes Model
60 print("\nTraining Naive Bayes model...")
61 nb_model = MultinomialNB()
62 nb_model.fit(X_train, y_train)
63 nb_metrics = evaluate_model(nb_model, X_test, y_test, "Naive Bayes")
64
65 # Compare with CountVectorizer features
66 X_train_count, X_test_count, y_train_count, y_test_count = train_test_split(
67     X_count, y, test_size=0.2, random_state=42)
68
69 print("\nTraining Logistic Regression with CountVectorizer features...")
70 lr_count_model = LogisticRegression(
71     max_iter=1000,
72     random_state=42,
73     class_weight='balanced'
74 )
75 lr_count_model.fit(X_train_count, y_train_count)
76 lr_count_metrics = evaluate_model(lr_count_model, X_test_count, y_test_count,
77     "Logistic Regression (CountVectorizer)")
78
79 # Create comparison table
80 results = pd.DataFrame({
81     'Model': ['Logistic Regression (TF-IDF)', 'Naive Bayes (TF-IDF)',
82             'Logistic Regression (CountVectorizer)'],
83     'Accuracy': [lr_metrics['accuracy'], nb_metrics['accuracy'], lr_count_metrics['accu
84     'Precision': [lr_metrics['precision'], nb_metrics['precision'], lr_count_metrics['pr
85     'Recall': [lr_metrics['recall'], nb_metrics['recall'], lr_count_metrics['recall']],
86     'F1 Score': [lr_metrics['f1'], nb_metrics['f1'], lr_count_metrics['f1']]
87 })
88
89 print("\nModel Performance Comparison:")
90 print(results.to_markdown(index=False))
91

```

```
92 # Visualize model comparison
93 plt.figure(figsize=(10, 6))
94 results.set_index('Model').plot(kind='bar', rot=45)
95 plt.title('Model Performance Comparison')
96 plt.ylabel('Score')
97 plt.ylim(0.7, 1.0)
98 plt.tight_layout()
99 plt.show()
```



Test set: 10000 samples

Training Logistic Regression model...

Logistic Regression Performance:

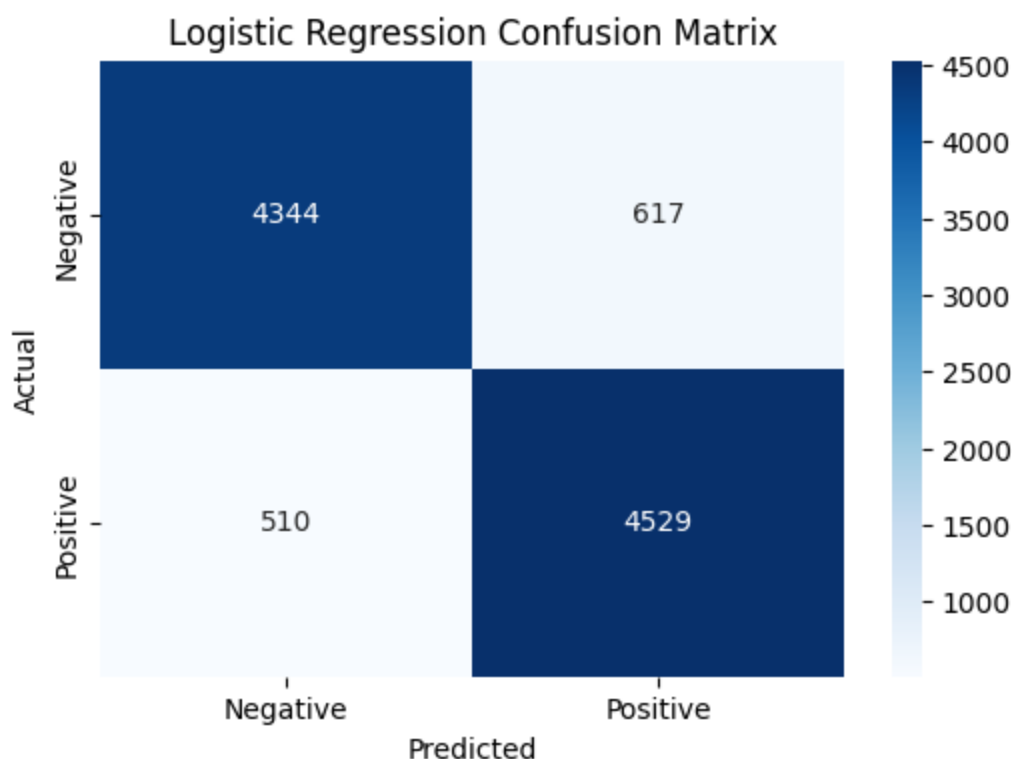
	precision	recall	f1-score	support
0	0.89	0.88	0.89	4961
1	0.88	0.90	0.89	5039
accuracy			0.89	10000
macro avg	0.89	0.89	0.89	10000
weighted avg	0.89	0.89	0.89	10000

Accuracy: 0.8873

Precision: 0.8801

Recall: 0.8988

F1 Score: 0.8893



Training Naive Bayes model...

Naive Bayes Performance:

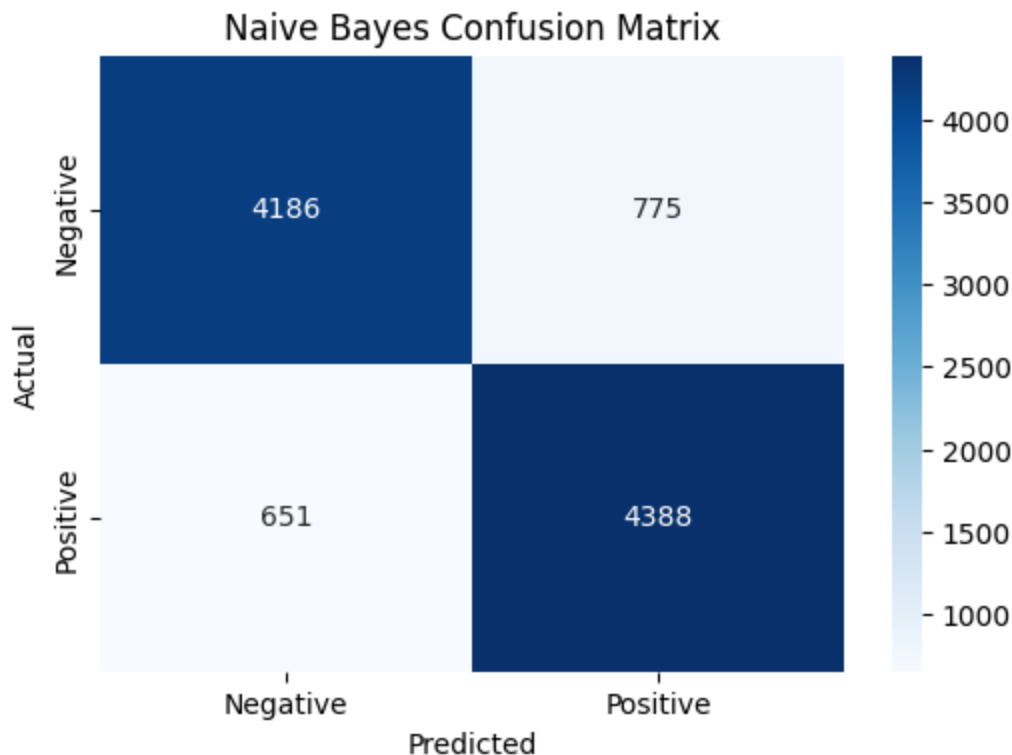
	precision	recall	f1-score	support
0	0.87	0.84	0.85	4961
1	0.85	0.87	0.86	5039
accuracy			0.86	10000
macro avg	0.86	0.86	0.86	10000
weighted avg	0.86	0.86	0.86	10000

Accuracy: 0.8574

Precision: 0.8400



Precision: 0.8499  
Recall: 0.8708  
F1 Score: 0.8602



Training Logistic Regression with CountVectorizer features...

Logistic Regression (CountVectorizer) Performance:

	precision	recall	f1-score	support
0	0.88	0.87	0.87	4961
1	0.87	0.88	0.88	5039
accuracy			0.87	10000
macro avg	0.87	0.87	0.87	10000
weighted avg	0.87	0.87	0.87	10000

Accuracy: 0.8747  
Precision: 0.8726  
Recall: 0.8797  
F1 Score: 0.8762

**Logistic Regression (CountVectorizer) Confusion Matrix**

