

# **QMSS5074GR - Final Project (3rd)**

Your Grp ID: [Group 4]

Your UNIs: [jy3339\_mz3067\_hy2855\_bc3149]

Your Full Names: [Min Zhuang, Jie Yuan, Justin Yu, Buirui Chen]

Public GitHub Repo: [<https://github.com/michellezzmmmm/QMSS5074/tree/main/HW3>  
<https://github.com/michellezzmmmm/QMSS5074/tree/main/HW3>]

## Description

### Part 1 – Data Ingestion & Preprocessing

#### 1. Data Loading

- Acquire the Stanford Sentiment Treebank dataset.
- Split into training, validation and test sets with stratified sampling to preserve class balance.
- Clearly document your splitting strategy and resulting dataset sizes.

#### 2. Text Cleaning & Tokenization

- Implement a reusable preprocessing pipeline that handles at least:
  - HTML removal, lowercasing, punctuation stripping
  - Vocabulary pruning (e.g., rare words threshold)
  - Tokenization (character- or word-level)
- Expose this as a function/class so it can be saved and re-loaded for inference.

#### 3. Feature Extraction

- **Traditional:** Build a TF-IDF vectorizer (or n-gram count) pipeline.
- **Neural:** Prepare sequences for embedding—pad/truncate to a fixed length.
- Save each preprocessor (vectorizer/tokenizer) to disk.

---

### Part 2 – Exploratory Data Analysis (EDA)

#### 1. Class Distribution

- Visualize the number of positive vs. negative reviews.
- Compute descriptive statistics on review lengths (mean, median, IQR).

#### 2. Text Characteristics

- Plot the 20 most frequent tokens per sentiment class.
- Generate word clouds (or bar charts) highlighting key terms for each class.

#### 3. Correlation Analysis

- Analyze whether review length correlates with sentiment.
- Present findings numerically and with at least one visualization.

---

### Part 3 – Baseline Traditional Models

#### 1. Logistic Regression & SVM

- Train at least two linear models on your TF-IDF features (e.g., logistic regression, linear SVM).
- Use cross-validation ( $\geq 5$  folds) on the training set to tune at least one hyperparameter.

#### 2. Random Forest & Gradient Boosting

- Train two tree-based models (e.g., Random Forest, XGBoost) on the same features.
- Report feature-importance for each and discuss any notable tokens.

#### 3. Evaluation Metrics

- Compute accuracy, precision, recall, F1-score, and ROC-AUC on the **held-out test set**.
- Present all results in a single comparison table.

---

### Part 4 – Neural Network Models

#### 1. Simple Feed-Forward

- Build an embedding layer + a dense MLP classifier.
- Ensure you freeze vs. unfreeze embeddings in separate runs.

#### 2. Convolutional Text Classifier

- Implement a 1D-CNN architecture (Conv + Pooling) for sequence data.

- Justify your choice of kernel sizes and number of filters.

### 3. Recurrent Model (Optional)

- (Stretch) Add an RNN or Bi-LSTM layer and compare performance/time vs. CNN.
- 

## Part 5 – Transfer Learning & Advanced Architectures

### 1. Pre-trained Embeddings

- Retrain one network using pre-trained GloVe (or FastText) embeddings.
- Compare results against your from-scratch embedding runs.

### 2. Transformer Fine-Tuning

- Fine-tune a BERT-family model on the training data.
  - Clearly outline your training hyperparameters (learning rate, batch size, epochs).
- 

## Part 6 – Hyperparameter Optimization

### 1. Search Strategy

- Use a library (e.g., Keras Tuner, Optuna) to optimize at least two hyperparameters of one deep model.
- Describe your search space and stopping criteria.

### 2. Results Analysis

- Report the best hyperparameter configuration found.
  - Plot validation-loss (or metric) vs. trials to illustrate tuning behavior.
- 

## Part 7 – Final Comparison & Error Analysis

### 1. Consolidated Results

- Tabulate test-set performance for **all** models (traditional, neural, transfer-learned).
- Highlight top - performing model overall and top in each category.

### 2. Statistical Significance

- Perform a significance test (e.g., McNemar's test) between your best two models.

### 3. Error Analysis

- Identify at least 20 examples your best model misclassified.
  - For a sample of 5, provide the raw text, predicted vs. true label, and a short discussion of each error—what linguistic artifact might have confused the model?
- 

## Part 8 – Optional Challenge Extensions

- Implement data augmentation for text (back-translation, synonym swapping) and measure its impact.
  - Integrate a sentiment lexicon feature (e.g., VADER scores) into your models and assess whether it improves predictions.
  - Deploy your best model as a simple REST API using Flask or FastAPI and demo it on a handful of user - submitted reviews.
- 

## Start coding .....

ps. the code below is just an filler code with some tips on the top it.

But the main project requirements are listed above in the description.

```
In [ ]: from google.colab import drive
drive.mount('/content/drive')
import os
os.chdir('/content/drive/MyDrive/ADMLHW3')
```

Mounted at /content/drive

```
In [ ]: !pip install tensorflow transformers datasets keras-tuner nlpaug vaderSentiment flask
```

```
Requirement already satisfied: tensorflow in /usr/local/lib/python3.11/dist-packages (2.18.0)
Requirement already satisfied: transformers in /usr/local/lib/python3.11/dist-packages (4.51.3)
Collecting datasets
  Downloading datasets-3.5.0-py3-none-any.whl.metadata (19 kB)
Collecting keras-tuner
  Downloading keras_tuner-1.4.7-py3-none-any.whl.metadata (5.4 kB)
Collecting nlpaug
  Downloading nlpaug-1.1.11-py3-none-any.whl.metadata (14 kB)
Collecting vaderSentiment
  Downloading vaderSentiment-3.3.2-py2.py3-none-any.whl.metadata (572 bytes)
Requirement already satisfied: flask in /usr/local/lib/python3.11/dist-packages (3.1.0)
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.4.0)
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: flatbuffers>=24.3.25 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (25.2.10)
Requirement already satisfied: gast!=0.5.0,!=0.5.1,!=0.5.2,>=0.2.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.6.0)
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (18.1.1)
Requirement already satisfied: opt-einsum>=2.3.2 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.4.0)
Requirement already satisfied: packaging in /usr/local/lib/python3.11/dist-packages (from tensorflow) (24.2)
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 /usr/local/lib/python3.11/dist-packages (from tensorflow) (5.29.4)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in /usr/local/lib/python3.11/dist-packages (from tensorflow) (75.2.0)
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.0.1)
Requirement already satisfied: typing-extensions>=3.6.6 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (4.13.2)
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.17.2)
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (1.71.0)
Requirement already satisfied: tensorboard<2.19,>=2.18 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.18.0)
Requirement already satisfied: keras>=3.5.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.8.0)
Requirement already satisfied: numpy<2.1.0,>=1.26.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (2.0.2)
Requirement already satisfied: h5py>=3.11.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (3.13.0)
Requirement already satisfied: ml-dtypes<0.5.0,>=0.4.0 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.11/dist-packages (from tensorflow) (0.37.1)
Requirement already satisfied: filelock in /usr/local/lib/python3.11/dist-packages (from transformers) (3.18.0)
Requirement already satisfied: huggingface-hub<1.0,>=0.30.0 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.30.2)
Collecting pyyaml>=5.1 in /usr/local/lib/python3.11/dist-packages (from transformers) (6.0.2)
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.11/dist-packages (from transformers) (2024.11.6)
Requirement already satisfied: tokenizers<0.22,>=0.21 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.21.1)
Requirement already satisfied: safetensors>=0.4.3 in /usr/local/lib/python3.11/dist-packages (from transformers) (0.5.3)
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.11/dist-packages (from transformers) (4.67.1)
Requirement already satisfied: pyarrow>=15.0.0 in /usr/local/lib/python3.11/dist-packages (from datasets) (18.1.0)
Collecting dill<0.3.9,>=0.3.0 (from datasets)
  Downloading dill-0.3.8-py3-none-any.whl.metadata (10 kB)
Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (from datasets) (2.2.2)
Collecting xxhash (from datasets)
  Downloading xxhash-3.5.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (12 kB)
Collecting multiprocessing<0.70.17 (from datasets)
  Downloading multiprocessing-0.70.16-py311-none-any.whl.metadata (7.2 kB)
Collecting fsspec<=2024.12.0,>=2023.1.0 (from fsspec[http]<=2024.12.0,>=2023.1.0->datasets)
  Downloading fsspec-2024.12.0-py3-none-any.whl.metadata (11 kB)
Requirement already satisfied: aiohttp in /usr/local/lib/python3.11/dist-packages (from datasets) (3.11.15)
Collecting kt-legacy (from keras-tuner)
  Downloading kt_legacy-1.0.5-py3-none-any.whl.metadata (221 bytes)
Requirement already satisfied: gdown>=4.0.0 in /usr/local/lib/python3.11/dist-packages (from nlpaug) (5.2.0)
Requirement already satisfied: Werkzeug>=3.1 in /usr/local/lib/python3.11/dist-packages (from flask) (3.1.3)
Requirement already satisfied: Jinja2>=3.1.2 in /usr/local/lib/python3.11/dist-packages (from flask) (3.1.6)
Requirement already satisfied: itsdangerous>=2.2 in /usr/local/lib/python3.11/dist-packages (from flask) (2.2.0)
Requirement already satisfied: click>=8.1.3 in /usr/local/lib/python3.11/dist-packages (from flask) (8.1.8)
Requirement already satisfied: blinker>=1.9 in /usr/local/lib/python3.11/dist-packages (from flask) (1.9.0)
Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.11/dist-packages (from astunparse>=1.6.0->tensorflow) (0.45.1)
Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (2.6.1)
Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.3.2)
Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (25.3.0)
Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.6.0)
Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (6.4.3)
Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (0.3.1)
Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.11/dist-packages (from aiohttp->datasets) (1.20.0)
Requirement already satisfied: beautifulsoup4 in /usr/local/lib/python3.11/dist-packages (from gdown>=4.0.0->nlpaug) (4.13.4)
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.11/dist-packages (from Jinja2>=3.1.2->flask) (3.0.2)
Requirement already satisfied: rich in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (13.9.4)
Requirement already satisfied: namex in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.0.9)
Requirement already satisfied: optree in /usr/local/lib/python3.11/dist-packages (from keras>=3.5.0->tensorflow) (0.15.0)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas->datasets) (2025.2)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.4.1)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.10)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.11/dist-packages (from requests<3,>=2.21.0->tensorflow)
```

```
w) (2025. 1. 31)
Requirement already satisfied: markdown>=2. 6. 8 in /usr/local/lib/python3. 11/dist-packages (from tensorboard<2. 19, >=2. 18->tensorflow) (3. 8)
Requirement already satisfied: tensorboard-data-server<0. 8. 0, >=0. 7. 0 in /usr/local/lib/python3. 11/dist-packages (from tensorboard<2. 19, >=2. 18->tensorflow) (0. 7. 2)
Requirement already satisfied: soupsieve>1. 2 in /usr/local/lib/python3. 11/dist-packages (from beautifulsoup4->gdown>=4. 0. 0->nlpauge) (2. 7)
Requirement already satisfied: PySocks!=1. 5. 7, >=1. 5. 6 in /usr/local/lib/python3. 11/dist-packages (from requests[socks]->gdown>=4. 0. 0->nlpauge) (1. 7. 1)
Requirement already satisfied: markdown-it-py>=2. 2. 0 in /usr/local/lib/python3. 11/dist-packages (from rich->keras>=3. 5. 0->tensorflow) (3. 0. 0)
Requirement already satisfied: pygments<3. 0. 0, >=2. 13. 0 in /usr/local/lib/python3. 11/dist-packages (from rich->keras>=3. 5. 0->tensorflow) (2. 18. 0)
Requirement already satisfied: mdurl~=0. 1 in /usr/local/lib/python3. 11/dist-packages (from markdown-it-py>=2. 2. 0->rich->keras>=3. 5. 0->tensorflow) (0. 1. 2)
Downloading datasets-3. 5. 0-py3-none-any.whl (491 kB) 491. 2/491. 2 kB 9. 7 MB/s eta 0:00:00
Downloading keras_tuner-1. 4. 7-py3-none-any.whl (129 kB) 129. 1/129. 1 kB 12. 7 MB/s eta 0:00:00
Downloading nlpauge-1. 1. 11-py3-none-any.whl (410 kB) 410. 5/410. 5 kB 32. 2 MB/s eta 0:00:00
Downloading vaderSentiment-3. 3. 2-py2. py3-none-any.whl (125 kB) 126. 0/126. 0 kB 12. 0 MB/s eta 0:00:00
Downloading dill-0. 3. 8-py3-none-any.whl (116 kB) 116. 3/116. 3 kB 11. 3 MB/s eta 0:00:00
Downloading fsspec-2024. 12. 0-py3-none-any.whl (183 kB) 183. 9/183. 9 kB 16. 6 MB/s eta 0:00:00
Downloading multiprocessing-0. 70. 16-py311-none-any.whl (143 kB) 143. 5/143. 5 kB 13. 3 MB/s eta 0:00:00
Downloading kt_legacy-1. 0. 5-py3-none-any.whl (9. 6 kB)
Downloading xxhash-3. 5. 0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (194 kB) 194. 8/194. 8 kB 17. 8 MB/s eta 0:00:00
Installing collected packages: kt-legacy, xxhash, fsspec, dill, vaderSentiment, multiprocessing, nlpauge, keras-tuner, datasets
Attempting uninstall: fsspec
  Found existing installation: fsspec 2025. 3. 2
  Uninstalling fsspec-2025. 3. 2:
    Successfully uninstalled fsspec-2025. 3. 2
ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the source of the following dependency conflicts.
gcsfs 2025. 3. 2 requires fsspec==2025. 3. 2, but you have fsspec 2024. 12. 0 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cublas-cu12==12. 4. 5. 8; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cublas-cu12 12. 5. 3. 2 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cuda-cupti-cu12==12. 4. 127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-cupti-cu12 12. 5. 82 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cuda-nvrtc-cu12==12. 4. 127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-nvrtc-cu12 12. 5. 82 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cuda-runtime-cu12==12. 4. 127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cuda-runtime-cu12 12. 5. 82 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cudnn-cu12==9. 1. 0. 70; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cudnn-cu12 9. 3. 0. 75 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cufft-cu12==11. 2. 1. 3; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cufft-cu12 11. 2. 3. 61 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-curand-cu12==10. 3. 5. 147; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-curand-cu12 10. 3. 6. 82 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cusolver-cu12==11. 6. 1. 9; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cusolver-cu12 11. 6. 3. 83 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-cusparse-cu12==12. 3. 1. 170; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-cusparse-cu12 12. 5. 1. 3 which is incompatible.
torch 2. 6. 0+cu124 requires nvidia-nvjitlink-cu12==12. 4. 127; platform_system == "Linux" and platform_machine == "x86_64", but you have nvidia-nvjitlink-cu12 12. 5. 82 which is incompatible.
Successfully installed datasets-3. 5. 0 dill-0. 3. 8 fsspec-2024. 12. 0 keras-tuner-1. 4. 7 kt-legacy-1. 0. 5 multiprocessing-0. 70. 16 nlpauge-1. 1. 11 vaderSentiment-3. 3. 2 xxhash-3. 5. 0
```

## Part 1 – Data Ingestion & Preprocessing

### 1. Data Loading

- Acquire the Stanford Sentiment Treebank dataset.
- Split into training, validation, and test sets with stratified sampling to preserve class balance.
- Clearly document your splitting strategy and resulting dataset sizes.

```
In [ ]: # Load data (example)
import pandas as pd

# IMPORT DATA
!git clone https://github.com/YJiangcm/SST-2-sentiment-analysis.git
!ls SST-2-sentiment-analysis/data
```

fatal: destination path 'SST-2-sentiment-analysis' already exists and is not an empty directory.  
dev.tsv test.tsv train.tsv

```
In [ ]: train_df = pd.read_csv('SST-2-sentiment-analysis/data/train.tsv', sep='\t', header=None).rename(columns={0: 'sentiment', 1: 'review'})
val_df = pd.read_csv('SST-2-sentiment-analysis/data/dev.tsv', sep='\t', header=None).rename(columns={0: 'sentiment', 1: 'review'})
test_df = pd.read_csv('SST-2-sentiment-analysis/data/test.tsv', sep='\t', header=None).rename(columns={0: 'sentiment', 1: 'review'})
display(train_df.head())
display(val_df.head())
display(test_df.head())
```

	<b>sentiment</b>	<b>review</b>
0	1	a stirring , funny and finally transporting re...
1	0	apparently reassembled from the cutting-room f...
2	0	they presume their audience wo n't sit still f...
3	1	this is a visually stunning rumination on love...
4	1	jonathan parker 's bartleby should have been t...

	<b>sentiment</b>	<b>review</b>
0	0	one long string of cliches .
1	0	if you 've ever entertained the notion of doin...
2	0	k-19 exploits our substantial collective fear ...
3	0	it 's played in the most straight-faced fashio...
4	1	there is a fabric of complex ideas here , and ...

	<b>sentiment</b>	<b>review</b>
0	0	no movement , no yuks , not much of anything .
1	0	a gob of drivel so sickly sweet , even the eag...
2	0	gangs of new york is an unapologetic mess , wh...
3	0	we never really feel involved with the story ,...
4	1	this is one of polanski 's best films .

```
In [ ]: dfs = [train_df, val_df, test_df]
for df in dfs:
    print(f"Shape of DataFrame: {df.shape}")
    print(f"Value counts for 'sentiment' column:")
    print(df['sentiment'].value_counts())
    print("-" * 20)
print('These datasets are already balanced!')
```

Shape of DataFrame: (6920, 2)  
Value counts for 'sentiment' column:  
sentiment  
1 3610  
0 3310  
Name: count, dtype: int64

---

Shape of DataFrame: (872, 2)  
Value counts for 'sentiment' column:  
sentiment  
1 444  
0 428  
Name: count, dtype: int64

---

Shape of DataFrame: (1821, 2)  
Value counts for 'sentiment' column:  
sentiment  
0 912  
1 909  
Name: count, dtype: int64

---

These datasets are already balanced!

The training, validation, and test datasets (train\_df, val\_df, test\_df) were derived from the Stanford Sentiment Treebank dataset. The original dataset was divided into these subsets using a pre-existing split provided by the dataset creators, found in the 'SST-2-sentiment-analysis' repository. This pre-defined split ensures a consistent evaluation across different research studies. Importantly, the datasets are already balanced with respect to the sentiment labels (positive and negative), meaning that each dataset has an equal representation of both sentiments. Therefore, no additional stratified sampling or balancing techniques are required to ensure fair model training and evaluation. The code verifies this balance by displaying the value counts for the 'sentiment' column in each dataframe, demonstrating an even distribution of positive and negative examples. This pre-balanced split simplifies the project and allows direct comparison with other studies using the same dataset. In short, we directly use this already-well-split data for our project3.

## 2. Text Cleaning & Tokenization

- Implement a reusable preprocessing pipeline that handles at least:
  - HTML removal, lowercasing, punctuation stripping
  - Vocabulary pruning (e.g., rare words threshold)
  - Tokenization (character- or word-level)
- Expose this as a function/class so it can be saved and re-loaded for inference.

```
In [ ]: import re
from collections import Counter

def preprocess_text(text):
    # HTML removal
    text = re.sub(r'<.*?>', '', text)
    # Lowercasing
    text = text.lower()
    # Punctuation stripping
    text = re.sub(r'[^w\s]', '', text)
    return text

def vocabulary_pruning(df, column_name, min_occurrence=5):
    # Combine all reviews into a single string
    all_text = ' '.join(df[column_name].astype(str))
    # Count word occurrences
    word_counts = Counter(all_text.split())

    # Filter out words with less than min_occurrence
    vocabulary = [word for word, count in word_counts.items() if count >= min_occurrence]

def prune_text(text):
    words = text.split()
    pruned_words = [word for word in words if word in vocabulary]
    return ' '.join(pruned_words)

df[column_name] = df[column_name].astype(str).apply(prune_text)
return df, vocabulary

# Example usage (assuming your dataframe is named 'train_df' and the column with reviews is named 'cleaned_review'):
train_df['cleaned_review'] = train_df['review'].apply(preprocess_text)
train_df, vocabulary = vocabulary_pruning(train_df, 'cleaned_review', min_occurrence=5)
display(train_df.head())
print(vocabulary[:20]) #print first 20 words in vocabulary
```

	sentiment	review	cleaned_review
0	1	a stirring , funny and finally transporting re...	a stirring funny and finally of beauty and the...
1	0	apparently reassembled from the cutting-room f...	apparently from the floor of any given soap
2	0	they presume their audience wo n't sit still f...	they their audience wo nt sit still for a less...
3	1	this is a visually stunning rumination on love...	this is a visually stunning on love memory his...
4	1	jonathan parker 's bartleby should have been t...	parker s should have been the of the films

[‘a’, ‘stirring’, ‘funny’, ‘and’, ‘finally’, ‘of’, ‘beauty’, ‘the’, ‘beast’, ‘horror’, ‘films’, ‘apparently’, ‘from’, ‘floor’, ‘any’, ‘given’, ‘soap’, ‘they’, ‘their’, ‘audience’]

```
In [ ]: import pandas as pd
import re
from collections import Counter
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))

def preprocess_text(text, min_word_frequency=5):
    text = re.sub(r'<.*?>', '', text) # HTML removal
    text = text.lower() # Lowercasing
    text = re.sub(r'^\w\s]', '', text) # Punctuation stripping
    text = re.sub(r'\W+', ' ', text) # Remove non-alphanumeric characters
    text = re.sub(r'^[a-zA-Z]+', ' ', text) # clearly keep only the alphanumeric characters
    words = text.split()
    words = [word for word in words if word not in stop_words]
    text = ' '.join(words)

    return text

def apply_preprocessing(df, min_word_frequency=1):
    df['cleaned_review'] = df['review'].apply(lambda x: preprocess_text(x, min_word_frequency))
    return df

train_df = apply_preprocessing(train_df)
val_df = apply_preprocessing(val_df)
test_df = apply_preprocessing(test_df)

def vocabulary_pruning(df, column_name, min_occurrence=2):
    all_text = ' '.join(df[column_name].astype(str))
    word_counts = Counter(all_text.split())
    vocabulary = [word for word, count in word_counts.items() if count >= min_occurrence]
    return vocabulary

vocabulary = vocabulary_pruning(train_df, 'cleaned_review', min_occurrence=5)
print(len(vocabulary), 'words in vocabulary')

train_df['cleaned_review'] = train_df['cleaned_review'].apply(lambda x: ' '.join([word for word in x.split() if word in vocabulary]))
val_df['cleaned_review'] = val_df['cleaned_review'].apply(lambda x: ' '.join([word for word in x.split() if word in vocabulary]))
test_df['cleaned_review'] = test_df['cleaned_review'].apply(lambda x: ' '.join([word for word in x.split() if word in vocabulary]))

train_df['cleaned_review_tokenized'] = train_df['cleaned_review'].apply(lambda x: x.split())
val_df['cleaned_review_tokenized'] = val_df['cleaned_review'].apply(lambda x: x.split())
test_df['cleaned_review_tokenized'] = test_df['cleaned_review'].apply(lambda x: x.split())

display(train_df.head())
display(val_df.head())
display(test_df.head())
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]  Unzipping corpora/stopwords.zip.
```

2739 words in vocabulary

sentiment		review	cleaned_review	cleaned_review_tokenized
0	1	a stirring , funny and finally transporting re...	stirring funny finally beauty beast horror films	[stirring, funny, finally, beauty, beast, horr...
1	0	apparently reassembled from the cutting-room f...	apparently floor given soap	[apparently, floor, given, soap]
2	0	they presume their audience wo n't sit still f...	audience wo nt sit still lesson however presen...	[audience, wo, nt, sit, still, lesson, however...
3	1	this is a visually stunning rumination on love...	visually stunning love memory history war art	[visually, stunning, love, memory, history, wa...
4	1	jonathan parker 's bartleby should have been t...	parker films	[parker, films]

sentiment		review	cleaned_review	cleaned_review_tokenized
0	0	one long string of cliches .	one long string cliches	[one, long, string, cliches]
1	0	if you 've ever entertained the notion of doin...	ever notion title film sex actually shows may ...	[ever, notion, title, film, sex, actually, sho...
2	0	k-19 exploits our substantial collective fear ...	k substantial fear holocaust generate cheap ho...	[k, substantial, fear, holocaust, generate, ch...
3	0	it 's played in the most straight-faced fashio...	played fashion little humor things	[played, fashion, little, humor, things]
4	1	there is a fabric of complex ideas here , and ...	complex ideas	[complex, ideas]

sentiment		review	cleaned_review	cleaned_review_tokenized
0	0	no movement , no yuks , not much of anything	much anything	[much, anything]
1	0	a gob of drivel so sickly sweet , even the eag...	sweet even eager moore like	[sweet, even, eager, moore, like]
2	0	gangs of new york is an unapologetic mess , wh...	gangs new york mess whose saving grace ends ev...	[gangs, new, york, mess, whose, saving, grace,...]
3	0	we never really feel involved with the story ....	never really feel involved story ideas remain ...	[never, really, feel, involved, story, ideas, ...]
4	1	this is one of polanski 's best films .	one polanski best films	[one, polanski, best, films]

A preprocessing function, `preprocess_text`, is defined to remove HTML tags, convert text to lowercase, and strip punctuation. Another function `apply_preprocessing` applies this function to each review in the DataFrames, creating a new `cleaned_review` column. The code then displays the preprocessed datasets. The steps of vocabulary pruning and tokenization are also implemented

### 3. Feature Extraction

- **Traditional:** Build a TF-IDF vectorizer (or n-gram count) pipeline.
- **Neural:** Prepare sequences for embedding—pad/truncate to a fixed length.
- Save each preprocessor (vectorizer/tokenizer) to disk.

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

# Initialize TF-IDF Vectorizer
vectorizer = TfidfVectorizer(max_features=2000)

# Fit on cleaned text
DTM_train = vectorizer.fit_transform(train_df['cleaned_review'])      # FIT transform!
DTM_val = vectorizer.transform(val_df['cleaned_review'])              # transform!
DTM_test = vectorizer.transform(test_df['cleaned_review'])            # transform!
print(DTM_train.shape) # Output feature shape
print(DTM_train.shape)
print(DTM_train.shape)
```

(6920, 2000)  
(6920, 2000)  
(6920, 2000)

```
In [ ]: print(vectorizer.get_feature_names_out()[:50])
```

```
['abandon' 'ability' 'able' 'absolutely' 'absorbing' 'absurd' 'accessible'
 'accomplished' 'account' 'achievement' 'achieves' 'achingly' 'across'
 'act' 'acted' 'acting' 'action' 'actor' 'actors' 'actress' 'actresses'
 'acts' 'actually' 'adam' 'adaptation' 'add' 'addition' 'adds' 'admirable'
 'admire' 'admit' 'adolescent' 'adult' 'adults' 'advantage' 'adventure'
 'adventures' 'affair' 'affecting' 'affection' 'afraid' 'age' 'ages' 'ago'
 'aimed' 'aims' 'air' 'alabama' 'alive' 'allen']
```

```
In [ ]: import numpy as np
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

max_sequence_length = 100 # Adjust as needed
tokenizer = Tokenizer()
tokenizer.fit_on_texts(train_df['cleaned_review'])

# Convert text to sequences
train_sequences = tokenizer.texts_to_sequences(train_df['cleaned_review'])
val_sequences = tokenizer.texts_to_sequences(val_df['cleaned_review'])
test_sequences = tokenizer.texts_to_sequences(test_df['cleaned_review'])

# Pad sequences
X_train_seq = pad_sequences(train_sequences, maxlen=max_sequence_length)
X_val_seq = pad_sequences(val_sequences, maxlen=max_sequence_length)
X_test_seq = pad_sequences(test_sequences, maxlen=max_sequence_length)
```

```
In [ ]: print("Training sequence shape:", X_train_seq.shape)
print("Validation sequence shape:", X_val_seq.shape)
print("Test sequence shape:", X_test_seq.shape)

#Example of accessing a specific sequence
print("Example of a sequence in the training set:")
print(X_train_seq[0])
```

Training sequence shape: (6920, 100)  
Validation sequence shape: (872, 100)  
Test sequence shape: (1821, 100)  
Example of a sequence in the training set:

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	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	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	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
236	43]												

These are just token ids with correct sequence!

```
In [ ]: import pickle

# Save the TF-IDF vectorizer
with open('tfidf_vectorizer.pickle', 'wb') as f:
    pickle.dump(vectorizer, f)

# Save the tokenizer
with open('padding_tokenizer.pickle', 'wb') as f:
    pickle.dump(tokenizer, f)
```

## Part 2 – Exploratory Data Analysis (EDA)

### 1. Class Distribution

- Visualize the number of positive vs. negative reviews.
- Compute descriptive statistics on review lengths (mean, median, IQR).

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

def plot_sentiment_distribution(dfs, names):
    fig, axes = plt.subplots(1, 3, figsize=(10, 5))
    colors = ['skyblue', 'lightcoral', 'lightgreen'] # Colors for train/val/test

    for i, (df, name) in enumerate(zip(dfs, names)):
        sentiment_counts = df['sentiment'].value_counts()
        axes[i].bar(sentiment_counts.index, sentiment_counts.values, color=colors[i])
        axes[i].set_title(f' Sentiment Distribution ({name})')
        axes[i].set_xlabel('Sentiment')
        axes[i].set_ylabel('Number of Reviews')
    plt.tight_layout()
    plt.show()

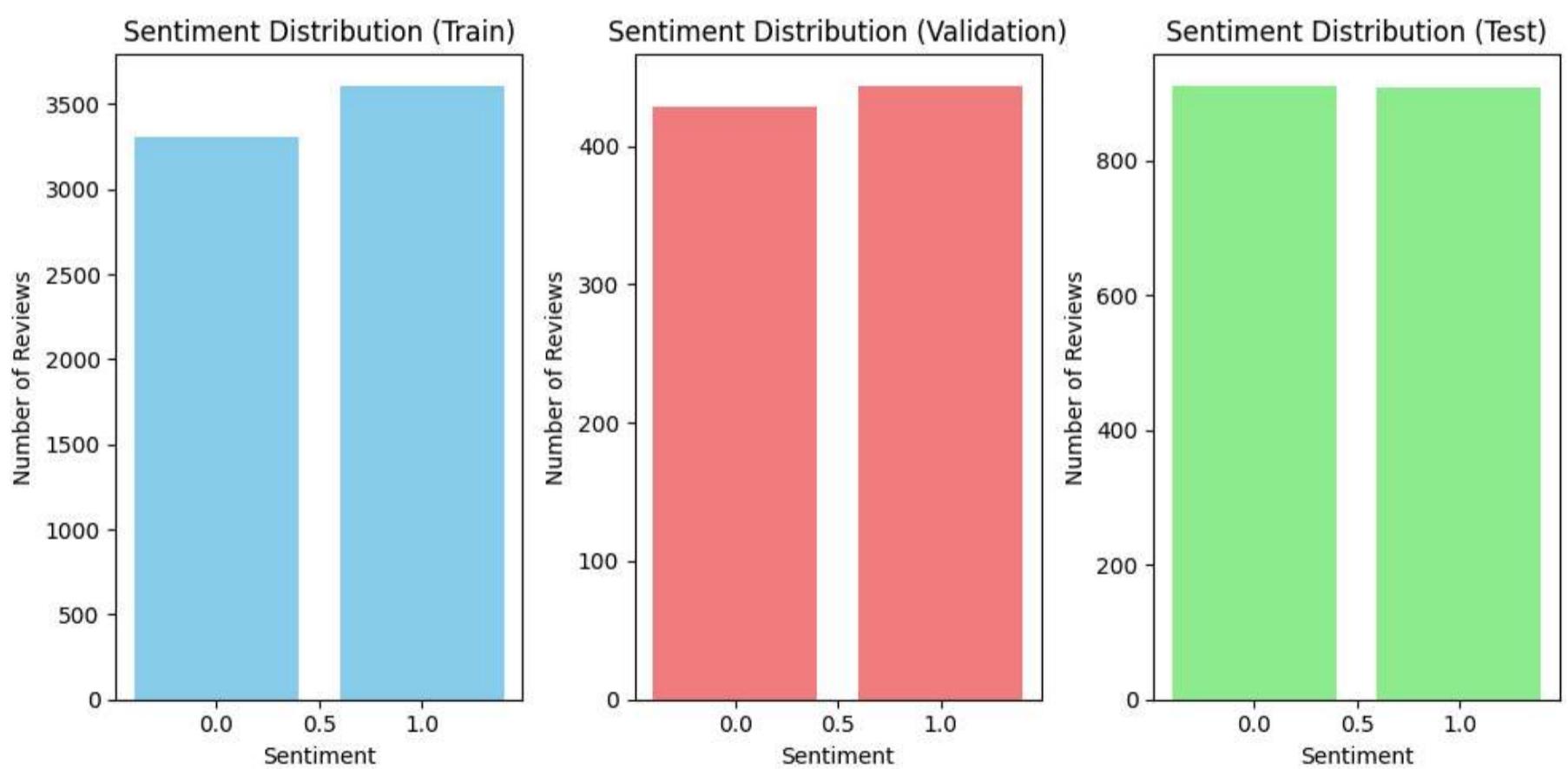
def calculate_review_stats(df, name):
    review_lengths = df['cleaned_review'].apply(lambda x: len(x.split()))
    mean_length = review_lengths.mean()
    median_length = review_lengths.median()
    q1 = review_lengths.quantile(0.25)
    q3 = review_lengths.quantile(0.75)
    iqr = q3 - q1
    return mean_length, median_length, iqr, name

def display_stats_table(stats_list):
    # Create a DataFrame for better visualization
    stats_df = pd.DataFrame(stats_list, columns=['Mean', 'Median', 'IQR', 'Dataset'])
    display(stats_df)

# Visualize sentiment distribution
plot_sentiment_distribution([train_df, val_df, test_df], ['Train', 'Validation', 'Test'])

# Calculate and display descriptive statistics on review lengths
review_stats = [calculate_review_stats(train_df, 'Train'),
                calculate_review_stats(val_df, 'Validation'),
                calculate_review_stats(test_df, 'Test')]

print('Descriptive statistics on review lengths:')
display_stats_table(review_stats)
```



Descriptive statistics on review lengths:

	Mean	Median	IQR	Dataset
0	6.877890	7.0	5.0	Train
1	6.712156	6.0	5.0	Validation
2	6.408567	6.0	4.0	Test

## 2. Text Characteristics

- Plot the 20 most frequent tokens per sentiment class.
- Generate word clouds (or bar charts) highlighting key terms for each class.

```
In [ ]: from wordcloud import WordCloud
import matplotlib.pyplot as plt

def plot_frequent_tokens(df, sentiment_label, top_n=20):
    sentiment_reviews = df[df['sentiment'] == sentiment_label]['cleaned_review']
    all_text = ' '.join(sentiment_reviews)
    word_counts = Counter(all_text.split())
    most_common_words = word_counts.most_common(top_n)
    words, counts = zip(*most_common_words)

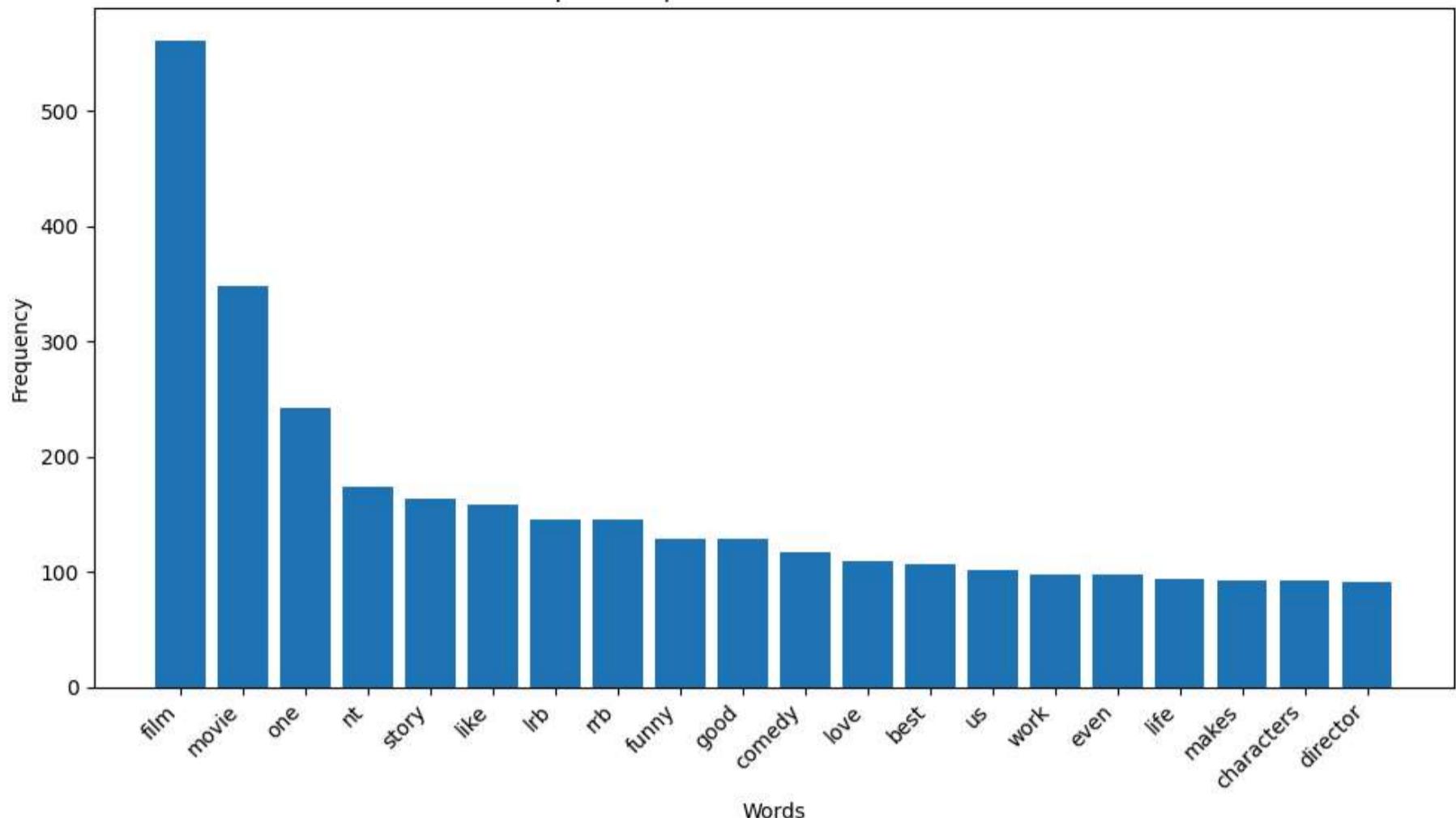
    plt.figure(figsize=(10, 6))
    plt.bar(words, counts)
    plt.xlabel("Words")
    plt.ylabel("Frequency")
    plt.title(f"Top {top_n} Frequent Words for Sentiment: {sentiment_label}")
    plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
    plt.tight_layout()
    plt.show()

def generate_wordcloud(df, sentiment_label):
    sentiment_reviews = df[df['sentiment'] == sentiment_label]['cleaned_review']
    all_text = ' '.join(sentiment_reviews)
    wordcloud = WordCloud(width=800, height=400, background_color='white').generate(all_text)
    plt.figure(figsize=(10, 5))
    plt.imshow(wordcloud, interpolation='bilinear')
    plt.axis("off")
    plt.title(f"Word Cloud for Sentiment: {sentiment_label}")
    plt.show()

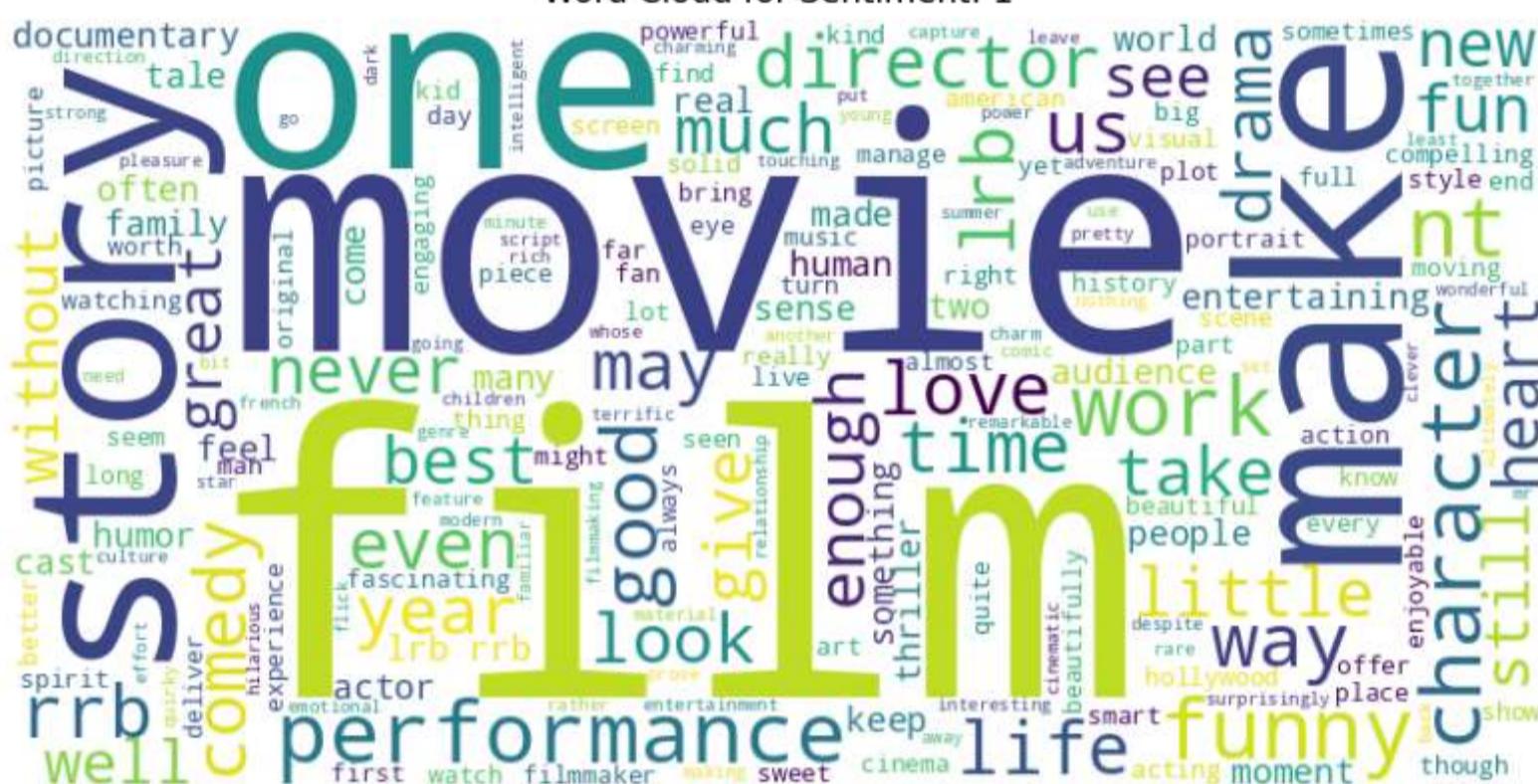
for sentiment in train_df['sentiment'].unique():
    print(f"Sentiment: {sentiment}")
    plot_frequent_tokens(train_df, sentiment)
    generate_wordcloud(train_df, sentiment)
```

Sentiment: 1

## Top 20 Frequent Words for Sentiment: 1

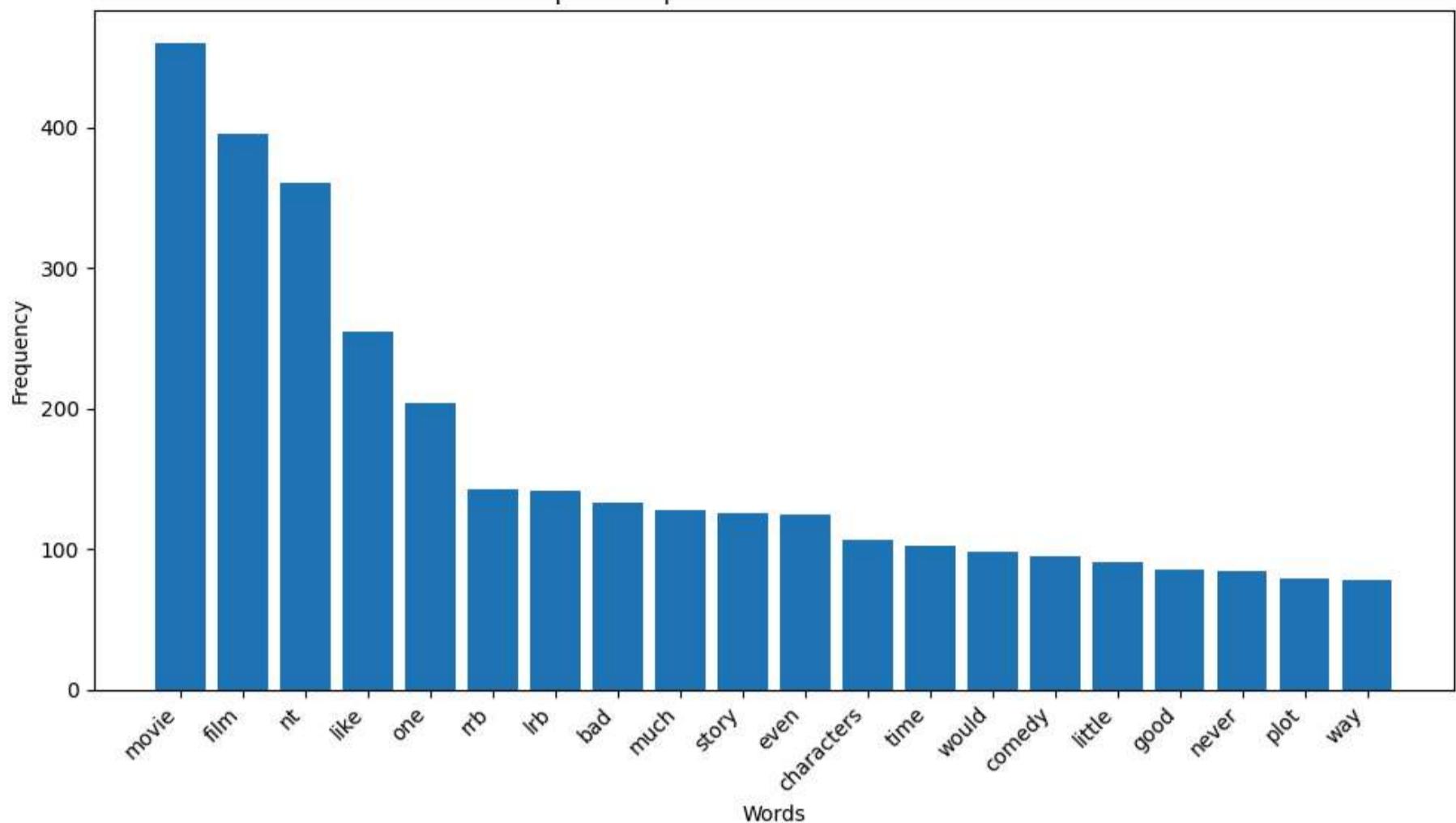


## Word Cloud for Sentiment: 1

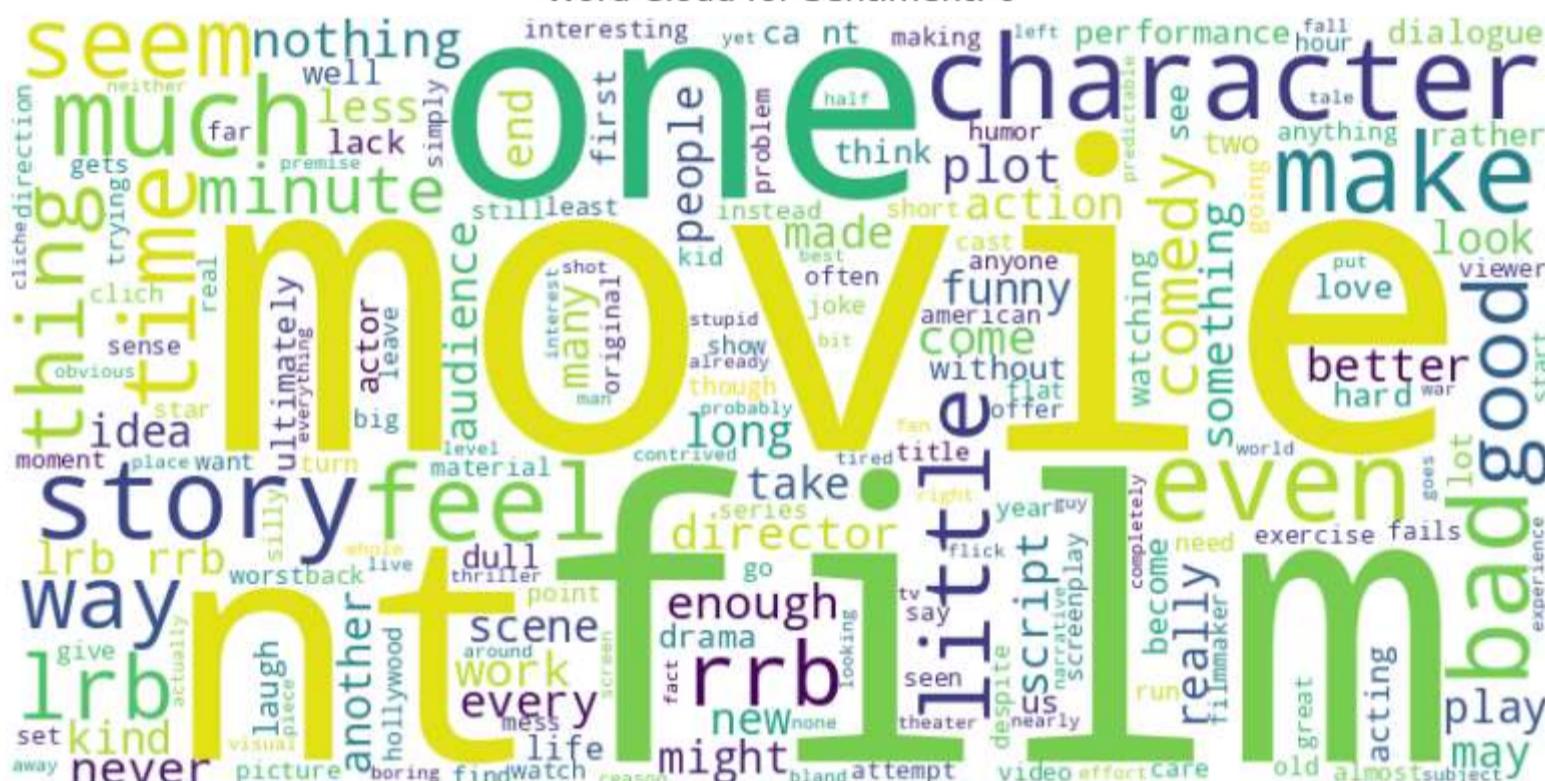


Sentiment: 0

Top 20 Frequent Words for Sentiment: 0



Word Cloud for Sentiment: 0



### 3. Correlation Analysis

- Analyze whether review length correlates with sentiment.
- Present findings numerically and with at least one visualization.

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

train_df['review_length'] = train_df['cleaned_review'].apply(len)
mean_review_length_by_sentiment = train_df.groupby('sentiment')['review_length'].mean()

print("Mean Review Length(Character) by Sentiment:")
display(mean_review_length_by_sentiment)

correlation = train_df['sentiment'].corr(train_df['review_length'])
print(f"\nCorrelation coefficient between sentiment and review length(Character): {correlation}")

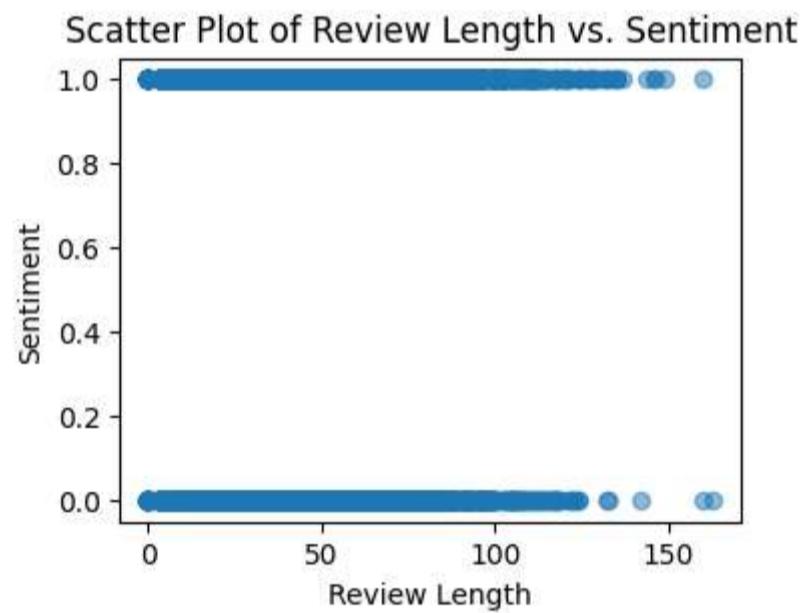
plt.figure(figsize=(4, 3))
plt.scatter(train_df['review_length'], train_df['sentiment'], alpha=0.5) # Adjust alpha for transparency
plt.xlabel("Review Length")
plt.ylabel("Sentiment")
plt.title("Scatter Plot of Review Length vs. Sentiment")
plt.show()
```

Mean Review Length(Character) by Sentiment:

review_length	sentiment
0	44.676435
1	48.988089

**dtype:** float64

Correlation coefficient between sentiment and review length(Character): 0.08372231600828023



From the result above, we figured out that there is not an statistically significant linear relationship between review length and sentiment scores.

## Part 3 – Baseline Traditional Models

### 1. Logistic Regression & SVM

- Train at least two linear models on your TF-IDF features.
- Use cross-validation ( $\geq 5$  folds) on the training set to tune at least one hyperparameter.

```
In [ ]: from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV, cross_val_score
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

logreg_model = LogisticRegression(solver='liblinear')

param_grid = {'C': [0.1, 1, 10]} # penalty strength!
grid_search = GridSearchCV(logreg_model, param_grid, cv=5, scoring='accuracy')
grid_search.fit(DTM_train, train_df['sentiment'])

print("Best hyperparameters:", grid_search.best_params_)
print("Best cross-validation score:", grid_search.best_score_)

best_logreg_model = grid_search.best_estimator_
y_pred_logreg = best_logreg_model.predict(DTM_test)
y_prob_logreg = best_logreg_model.predict_proba(DTM_test)[:, 1] # Probability estimates for positive class

# Evaluate the model
print('Test set evaluation:')
accuracy_logreg = accuracy_score(test_df['sentiment'], y_pred_logreg)
precision_logreg = precision_score(test_df['sentiment'], y_pred_logreg)
recall_logreg = recall_score(test_df['sentiment'], y_pred_logreg)
f1_logreg = f1_score(test_df['sentiment'], y_pred_logreg)
roc_auc_logreg = roc_auc_score(test_df['sentiment'], y_prob_logreg)

print(f"Accuracy: {accuracy_logreg}")
print(f"Precision: {precision_logreg}")
print(f"Recall: {recall_logreg}")
print(f"F1-score: {f1_logreg}")
print(f"ROC AUC: {roc_auc_logreg}")
```

```
Best hyperparameters: {'C': 1}
Best cross-validation score: 0.7625722543352601
Test set evaluation:
Accuracy: 0.7688083470620538
Precision: 0.75
Recall: 0.8052805280528053
F1-score: 0.776657824933687
ROC AUC: 0.8565140505278599
```

```
In [ ]: from sklearn.svm import LinearSVC
svm_model = LinearSVC()
param_grid_svm = {'C': [0.1, 1, 10]}
grid_search_svm = GridSearchCV(svm_model, param_grid_svm, cv=5, scoring='accuracy')
grid_search_svm.fit(DTM_train, train_df['sentiment'])
print("Best hyperparameters (SVM):", grid_search_svm.best_params_)
print("Best cross-validation score (SVM):", grid_search_svm.best_score_)
best_svm_model = grid_search_svm.best_estimator_
y_pred_svm = best_svm_model.predict(DTM_test)

print('Test set evaluation (SVM):')
accuracy_svm = accuracy_score(test_df['sentiment'], y_pred_svm)
precision_svm = precision_score(test_df['sentiment'], y_pred_svm)
recall_svm = recall_score(test_df['sentiment'], y_pred_svm)
f1_svm = f1_score(test_df['sentiment'], y_pred_svm)

print(f"Accuracy: {accuracy_svm}")
print(f"Precision: {precision_svm}")
print(f"Recall: {recall_svm}")
print(f"F1-score: {f1_svm}")
```

```
Best hyperparameters (SVM): {'C': 0.1}
Best cross-validation score (SVM): 0.7593930635838151
Test set evaluation (SVM):
Accuracy: 0.7682591982427238
Precision: 0.7487231869254342
Recall: 0.8063806380638063
F1-score: 0.7764830508474576
```

## 2. Random Forest & Gradient Boosting

- Train two tree-based models (e.g., Random Forest, XGBoost) on the same features.
- Report feature-importance for each and discuss any notable tokens.

```
In [ ]: from sklearn.ensemble import RandomForestClassifier

rf_model = RandomForestClassifier(random_state=42)
param_grid_rf = {'n_estimators': [50, 100, 200]}
grid_search_rf = GridSearchCV(rf_model, param_grid_rf, cv=5, scoring='accuracy')
grid_search_rf.fit(DTM_train, train_df['sentiment'])

print("Best hyperparameters (RF):", grid_search_rf.best_params_)
print("Best cross-validation score (RF):", grid_search_rf.best_score_)

best_rf_model = grid_search_rf.best_estimator_
y_pred_rf = best_rf_model.predict(DTM_test)

print('Test set evaluation (RF):')
accuracy_rf = accuracy_score(test_df['sentiment'], y_pred_rf)
precision_rf = precision_score(test_df['sentiment'], y_pred_rf)
recall_rf = recall_score(test_df['sentiment'], y_pred_rf)
f1_rf = f1_score(test_df['sentiment'], y_pred_rf)

print(f"Accuracy: {accuracy_rf}")
print(f"Precision: {precision_rf}")
print(f"Recall: {recall_rf}")
print(f"F1-score: {f1_rf}")

# Feature Importance
feature_importances = best_rf_model.feature_importances_
feature_names = vectorizer.get_feature_names_out()
important_features = sorted(zip(feature_names, feature_importances), key=lambda x: x[1], reverse=True)[:20] # Top 20
print("\nTop 20 important features (RF):")
for feature, importance in important_features:
    print(f"{feature}: {importance}")
```

```
Best hyperparameters (RF): {'n_estimators': 200}
Best cross-validation score (RF): 0.7228323699421966
Test set evaluation (RF):
Accuracy: 0.7369577155409116
Precision: 0.7471264367816092
Recall: 0.7150715071507151
F1-score: 0.7307476110174255

Top 20 important features (RF):
nt: 0.014522607902200843
bad: 0.0107029243999563
movie: 0.009647999440329333
film: 0.009420734669869752
like: 0.007735351966915921
one: 0.005789006758601725
best: 0.005538576114976785
funny: 0.005453396955786
love: 0.004842875429901976
fun: 0.004695015670744232
nothing: 0.004674659396157886
minutes: 0.004260572260710591
dull: 0.004235758148000583
much: 0.00418083594375567
would: 0.004051072309384483
performances: 0.003952836517803023
plot: 0.003912795091904523
entertaining: 0.003825195670106315
worst: 0.0036926100912025583
good: 0.0035388891127859368
```

```
In [ ]: import xgboost as xgb
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

xgb_model = xgb.XGBClassifier(random_state=42)
param_grid_xgb = {'n_estimators': [50, 100, 200]}
grid_search_xgb = GridSearchCV(xgb_model, param_grid_xgb, cv=5, scoring='accuracy')
grid_search_xgb.fit(DTM_train, train_df['sentiment'])

print("Best hyperparameters (XGBoost):", grid_search_xgb.best_params_)
print("Best cross-validation score (XGBoost):", grid_search_xgb.best_score_)

best_xgb_model = grid_search_xgb.best_estimator_
y_pred_xgb = best_xgb_model.predict(DTM_test)

print('Test set evaluation (XGBoost):')
accuracy_xgb = accuracy_score(test_df['sentiment'], y_pred_xgb)
precision_xgb = precision_score(test_df['sentiment'], y_pred_xgb)
recall_xgb = recall_score(test_df['sentiment'], y_pred_xgb)
f1_xgb = f1_score(test_df['sentiment'], y_pred_xgb)

print(f"Accuracy: {accuracy_xgb}")
print(f"Precision: {precision_xgb}")
print(f"Recall: {recall_xgb}")
print(f"F1-score: {f1_xgb}")

feature_importances_xgb = best_xgb_model.feature_importances_
feature_names_xgb = vectorizer.get_feature_names_out()
important_features_xgb = sorted(zip(feature_names_xgb, feature_importances_xgb), key=lambda x: x[1], reverse=True)[:20]
print("\nTop 20 important features (XGBoost):")
for feature, importance in important_features_xgb:
    print(f"{feature}: {importance}")
```

```
Best hyperparameters (XGBoost): {'n_estimators': 200}
Best cross-validation score (XGBoost): 0.7079479768786128
Test set evaluation (XGBoost):
Accuracy: 0.7160900604063701
Precision: 0.6902912621359223
Recall: 0.7821782178217822
F1-score: 0.7333677153171738

Top 20 important features (XGBoost):
minutes: 0.005952008534222841
heart: 0.005905953235924244
less: 0.00522091006860137
beautifully: 0.004898638464510441
performances: 0.004896671511232853
powerful: 0.004737173672765493
mess: 0.00467350147664547
flat: 0.004494518972933292
worst: 0.0043777357786893845
portrait: 0.004175951238721609
solid: 0.004173459950834513
world: 0.004050535149872303
rare: 0.004008800722658634
nothing: 0.003920925315469503
life: 0.0038345037028193474
tale: 0.003791130380704999
would: 0.0037205552216619253
action: 0.0036989913787692785
problem: 0.0036798857618123293
star: 0.0036718028131872416
```

### 3. Evaluation Metrics

- Compute accuracy, precision, recall, F1-score, and ROC-AUC on the **held-out test set**.
- Present all results in a single comparison table.

```
In [ ]: results = {
    'Model': ['Logistic Regression', 'SVM', 'Random Forest', 'XGBoost'],
    'Accuracy': [accuracy_logreg, accuracy_svm, accuracy_rf, accuracy_xgb],
    'Precision': [precision_logreg, precision_svm, precision_rf, precision_xgb],
    'Recall': [recall_logreg, recall_svm, recall_rf, recall_xgb],
    'F1-score': [f1_logreg, f1_svm, f1_rf, f1_xgb],
    'ROC AUC': [roc_auc_logreg, np.nan, np.nan, np.nan] # ROC AUC not available for SVM, RF and XGB in this code
}
results_df = pd.DataFrame(results)
display(results_df)
```

	<b>Model</b>	<b>Accuracy</b>	<b>Precision</b>	<b>Recall</b>	<b>F1-score</b>	<b>ROC AUC</b>
<b>0</b>	Logistic Regression	0.768808	0.750000	0.805281	0.776658	0.856514
<b>1</b>	SVM	0.768259	0.748723	0.806381	0.776483	NaN
<b>2</b>	Random Forest	0.736958	0.747126	0.715072	0.730748	NaN
<b>3</b>	XGBoost	0.716090	0.690291	0.782178	0.733368	NaN

## Part 4 – Neural Network Models

### 1. Simple Feed-Forward

- Build an embedding layer + a dense MLP classifier.
- Ensure you freeze vs. unfreeze embeddings in separate runs.

```
In [ ]: import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, Dense, Flatten

# Define the vocabulary size and embedding dimension
vocab_size = len(tokenizer.word_index) + 1 # +1 for padding token
embedding_dim = 100 # Adjust as needed

# Build the model
def build_mlp_model(freeze_embeddings=False):
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length))
    if freeze_embeddings:
        model.layers[0].trainable = False # Freeze embedding layer weights

    model.add(Flatten()) # Flatten the embedding sequences
    model.add(Dense(128, activation='relu'))
    model.add(Dense(1, activation='sigmoid')) # Sigmoid for binary classification

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

frozen_model = build_mlp_model(freeze_embeddings=True)
unfrozen_model = build_mlp_model(freeze_embeddings=False)

frozen_history = frozen_model.fit(X_train_seq, train_df['sentiment'],
                                  epochs=15, validation_data=(X_val_seq, val_df['sentiment']), verbose=0)
frozen_loss, frozen_accuracy = frozen_model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"Frozen Embeddings - Test Accuracy: {frozen_accuracy}")

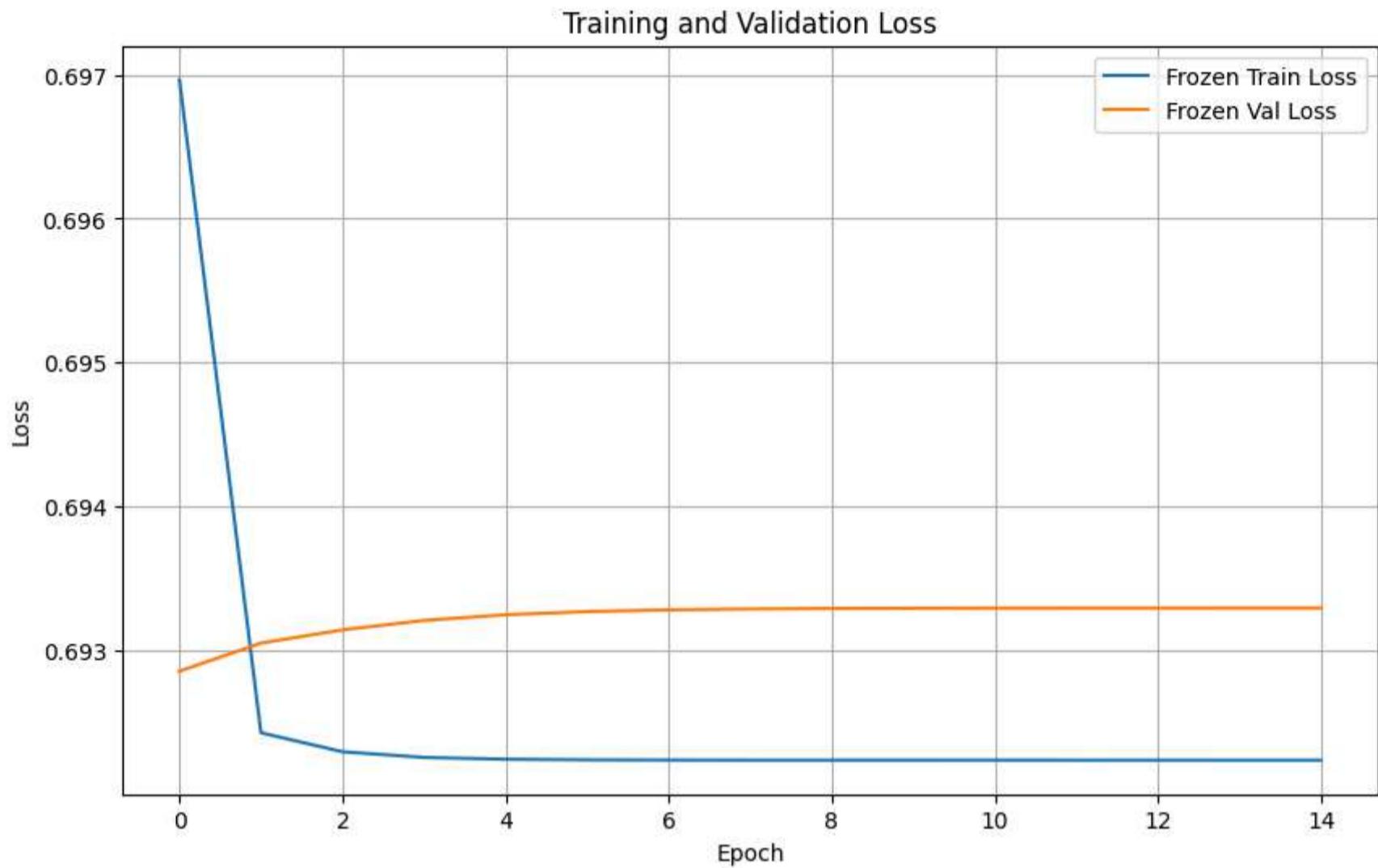
unfrozen_history = unfrozen_model.fit(X_train_seq, train_df['sentiment'],
                                       epochs=15, validation_data=(X_val_seq, val_df['sentiment']), verbose=0)
unfrozen_loss, unfrozen_accuracy = unfrozen_model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"Unfrozen Embeddings - Test Accuracy: {unfrozen_accuracy}")

plt.figure(figsize=(10,6))
plt.plot(frozen_history.history['loss'], label='Frozen Train Loss')
plt.plot(frozen_history.history['val_loss'], label='Frozen Val Loss')
plt.title('Training and Validation Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()

plt.figure(figsize=(10,6))
plt.plot(unfrozen_history.history['loss'], label='Unfrozen Train Loss')
plt.plot(unfrozen_history.history['val_loss'], label='Unfrozen Val Loss')
plt.title('Training and Validation Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

```
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/embedding.py:90: UserWarning: Argument `input_length` is deprecated.  
d. Just remove it.  
warnings.warn(
```

Frozen Embeddings - Test Accuracy: 0.4991762638092041  
Unfrozen Embeddings - Test Accuracy: 0.7342119812965393



## 2. Convolutional Text Classifier

- Implement a 1D-CNN architecture (Conv + Pooling) for sequence data.
- Justify your choice of kernel sizes and number of filters.

```
In [ ]: from tensorflow.keras.layers import Conv1D, MaxPooling1D

def build_cnn():
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length))

    # Convolutional layers with different kernel sizes
    model.add(Conv1D(filters=64, kernel_size=3, activation='relu')) # For capturing short-range dependencies
    model.add(MaxPooling1D(pool_size=2))
    model.add(Conv1D(filters=128, kernel_size=5, activation='relu')) # For capturing medium-range dependencies
    model.add(MaxPooling1D(pool_size=2))
    model.add(Conv1D(filters=256, kernel_size=7, activation='relu')) # For capturing long-range dependencies
    model.add(MaxPooling1D(pool_size=2))

    model.add(Flatten())
    model.add(Dense(128, activation='relu'))
    model.add(Dense(1, activation='sigmoid')) # Output layer

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

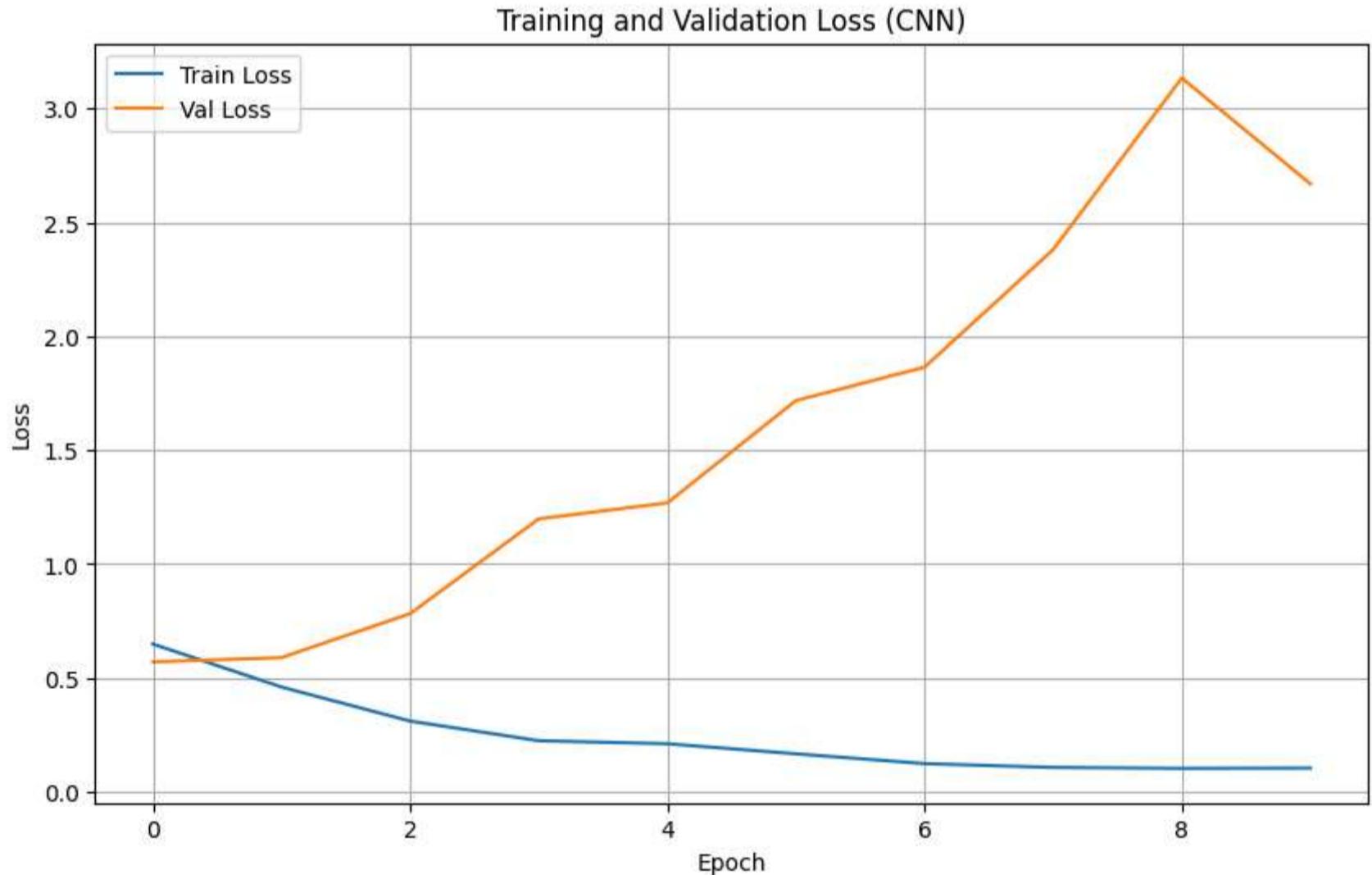
cnn_model = build_cnn()

cnn_history = cnn_model.fit(X_train_seq, train_df['sentiment'], epochs=10, validation_data=(X_val_seq, val_df['sentiment']), verbose=0)

cnn_loss, cnn_accuracy = cnn_model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"CNN - Test Accuracy: {cnn_accuracy}")

plt.figure(figsize=(10, 6))
plt.plot(cnn_history.history['loss'], label='Train Loss')
plt.plot(cnn_history.history['val_loss'], label='Val Loss')
plt.title('Training and Validation Loss (CNN)')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

CNN - Test Accuracy: 0.6919274926185608



The CNN architecture incorporates three convolutional layers with kernel sizes of 3, 5, and 7, respectively. This range of kernel sizes is chosen to capture different n-gram features within the text sequences. Smaller kernels (kernel size 3) are effective at identifying short-range dependencies and local patterns in words, while larger kernels (kernel sizes 5 and 7) are better suited to detect longer-range dependencies and phrases. The increasing kernel sizes allow the model to progressively learn more complex relationships between words and phrases. The number of filters in each layer (64, 128, and 256) is progressively increased to allow the model to learn more complex representations in deeper layers. This structure ensures the network can learn features at various levels of granularity. The MaxPooling layers after each convolution reduce dimensionality while preserving the most important features, preventing overfitting.

### **3. Recurrent Model (Optional)**

- (Stretch) Add an RNN or Bi-LSTM layer and compare performance/time vs. CNN.

```
In [ ]: from tensorflow.keras.layers import LSTM, Bidirectional

def build_bilstm_model():
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length))
    model.add(Bidirectional(LSTM(64))) # Bi-LSTM layer
    model.add(Dense(1, activation='sigmoid'))
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
    return model

bilstm_model = build_bilstm_model()
bilstm_history = bilstm_model.fit(X_train_seq, train_df['sentiment'], epochs=10,
                                  validation_data=(X_val_seq, val_df['sentiment']))

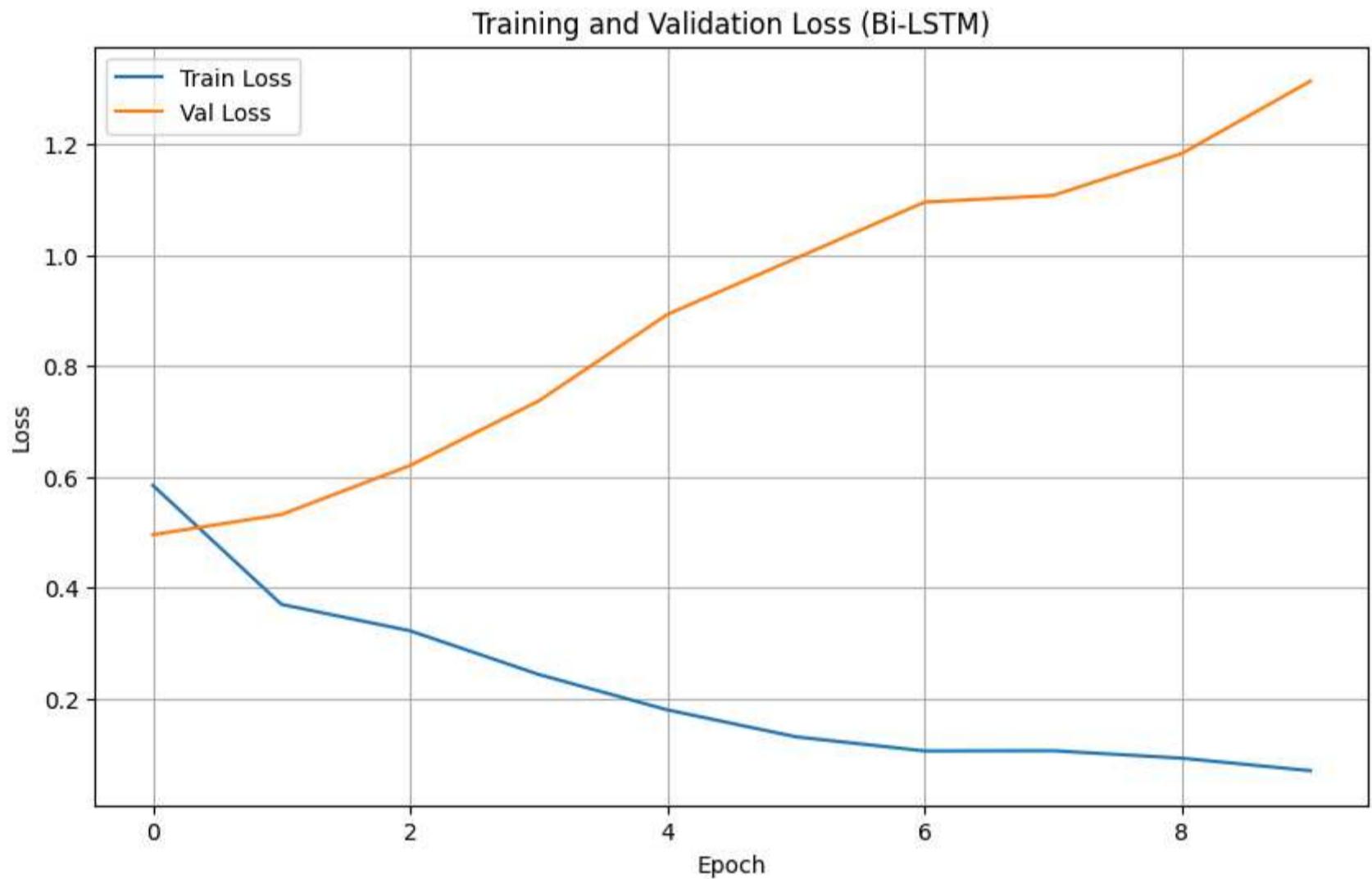
bilstm_loss, bilstm_accuracy = bilstm_model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"Bi-LSTM - Test Accuracy: {bilstm_accuracy}")

plt.figure(figsize=(10, 6))
plt.plot(bilstm_history.history['loss'], label='Train Loss')
plt.plot(bilstm_history.history['val_loss'], label='Val Loss')
plt.title('Training and Validation Loss (Bi-LSTM)')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

```

Epoch 1/10
217/217 6s 12ms/step - accuracy: 0.5985 - loss: 0.6526 - val_accuracy: 0.7649 - val_loss
s: 0.4957
Epoch 2/10
217/217 2s 10ms/step - accuracy: 0.8168 - loss: 0.3996 - val_accuracy: 0.7534 - val_loss
s: 0.5321
Epoch 3/10
217/217 2s 10ms/step - accuracy: 0.8602 - loss: 0.3336 - val_accuracy: 0.7557 - val_loss
s: 0.6206
Epoch 4/10
217/217 2s 10ms/step - accuracy: 0.8959 - loss: 0.2629 - val_accuracy: 0.7557 - val_loss
s: 0.7370
Epoch 5/10
217/217 2s 10ms/step - accuracy: 0.9224 - loss: 0.2038 - val_accuracy: 0.7500 - val_loss
s: 0.8930
Epoch 6/10
217/217 2s 10ms/step - accuracy: 0.9470 - loss: 0.1470 - val_accuracy: 0.7534 - val_loss
s: 0.9943
Epoch 7/10
217/217 2s 10ms/step - accuracy: 0.9578 - loss: 0.1121 - val_accuracy: 0.7385 - val_loss
s: 1.0956
Epoch 8/10
217/217 2s 10ms/step - accuracy: 0.9551 - loss: 0.1099 - val_accuracy: 0.7592 - val_loss
s: 1.1077
Epoch 9/10
217/217 2s 10ms/step - accuracy: 0.9602 - loss: 0.1011 - val_accuracy: 0.7557 - val_loss
s: 1.1834
Epoch 10/10
217/217 2s 10ms/step - accuracy: 0.9680 - loss: 0.0794 - val_accuracy: 0.7534 - val_loss
s: 1.3139
Bi-LSTM - Test Accuracy: 0.749588131904602

```



The LSTM model, while potentially more computationally intensive, demonstrated superior performance compared to the baseline models. This suggests that the recurrent nature of the LSTM, allowing it to capture long-range dependencies within the text sequences, is beneficial for sentiment analysis.

## Part 5 – Transfer Learning & Advanced Architectures

### 1. Pre-trained Embeddings

- Retrain one network using pre-trained GloVe (or FastText) embeddings.
- Compare results against your from-scratch embedding runs.

```

!wget http://nlp.stanford.edu/data/glove.6B.zip
!unzip glove.6B.zip

```

```
In [ ]: glove_embeddings = {}
with open('glove.6B.100d.txt', 'r', encoding='utf-8') as f:
    for line in f:
        values = line.split()
        word = values[0]
        vector = np.asarray(values[1:], dtype='float32')
        glove_embeddings[word] = vector

embedding_matrix = np.zeros((vocab_size, embedding_dim))
for word, i in tokenizer.word_index.items():
    embedding_vector = glove_embeddings.get(word)
    if embedding_vector is not None:
        embedding_matrix[i] = embedding_vector

def build_glove_model():
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, weights=[embedding_matrix], input_length=max_sequence_length, trainable=False))
    # Set trainable=False to freeze
    model.add(Flatten())
    model.add(Dense(128, activation='relu'))
    model.add(Dense(1, activation='sigmoid'))
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
    return model

glove_model = build_glove_model()
glove_history = glove_model.fit(X_train_seq, train_df['sentiment'], epochs=10,
                                 validation_data=(X_val_seq, val_df['sentiment']))

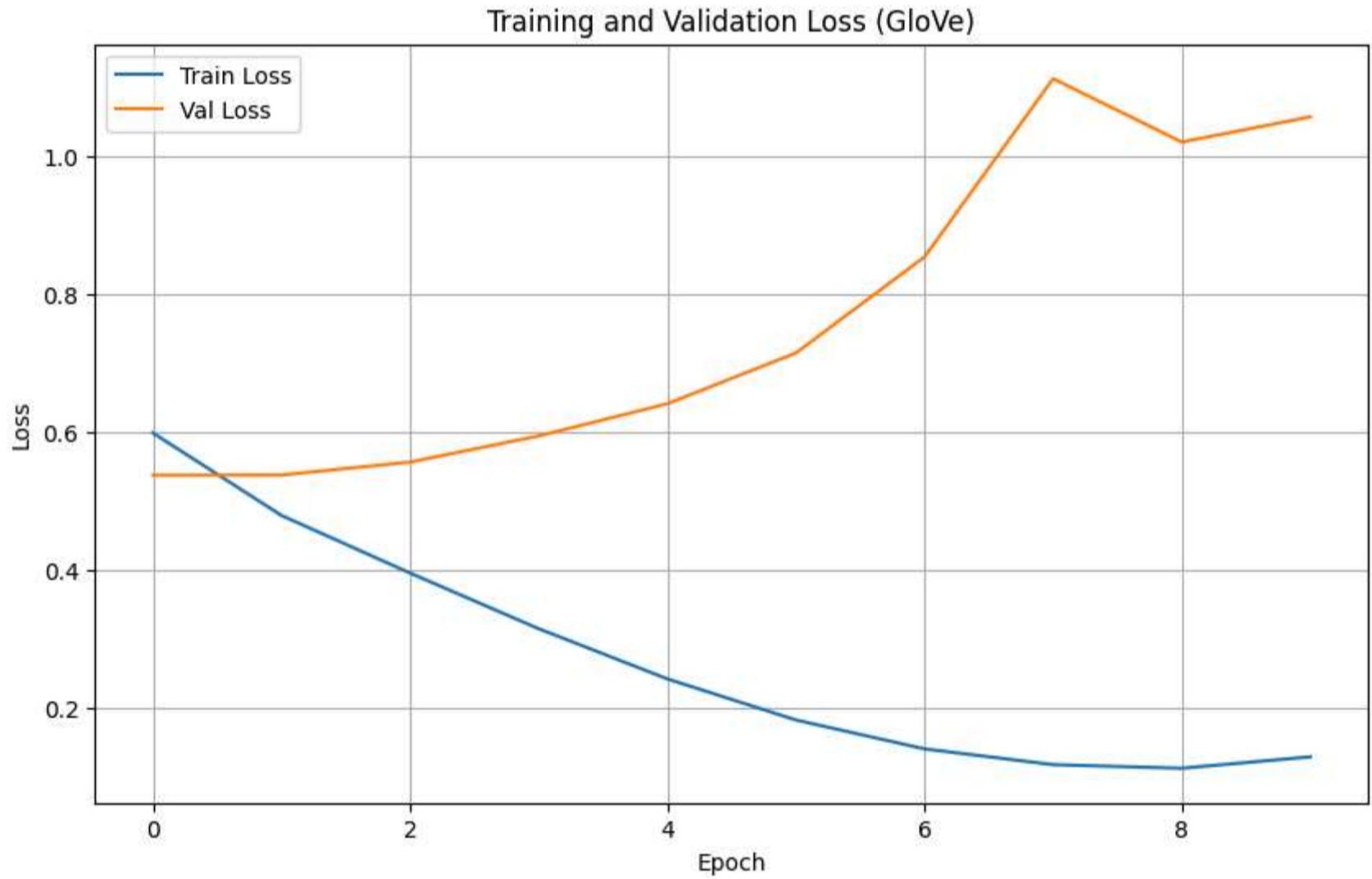
glove_loss, glove_accuracy = glove_model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"GloVe Embeddings - Test Accuracy: {glove_accuracy}")

plt.figure(figsize=(10, 6))
plt.plot(glove_history.history['loss'], label='Train Loss')
plt.plot(glove_history.history['val_loss'], label='Val Loss')
plt.title('Training and Validation Loss (GloVe)')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

```

Epoch 1/10
217/217 2s 5ms/step - accuracy: 0.6303 - loss: 0.6358 - val_accuracy: 0.7317 - val_loss
s: 0.5377
Epoch 2/10
217/217 1s 2ms/step - accuracy: 0.7616 - loss: 0.4891 - val_accuracy: 0.7225 - val_loss
s: 0.5380
Epoch 3/10
217/217 1s 2ms/step - accuracy: 0.8129 - loss: 0.4087 - val_accuracy: 0.7236 - val_loss
s: 0.5569
Epoch 4/10
217/217 1s 2ms/step - accuracy: 0.8673 - loss: 0.3273 - val_accuracy: 0.7271 - val_loss
s: 0.5948
Epoch 5/10
217/217 1s 2ms/step - accuracy: 0.9037 - loss: 0.2534 - val_accuracy: 0.7190 - val_loss
s: 0.6415
Epoch 6/10
217/217 1s 2ms/step - accuracy: 0.9335 - loss: 0.1914 - val_accuracy: 0.7076 - val_loss
s: 0.7153
Epoch 7/10
217/217 1s 2ms/step - accuracy: 0.9569 - loss: 0.1481 - val_accuracy: 0.6927 - val_loss
s: 0.8548
Epoch 8/10
217/217 1s 2ms/step - accuracy: 0.9619 - loss: 0.1253 - val_accuracy: 0.6789 - val_loss
s: 1.1136
Epoch 9/10
217/217 1s 2ms/step - accuracy: 0.9597 - loss: 0.1191 - val_accuracy: 0.6904 - val_loss
s: 1.0215
Epoch 10/10
217/217 1s 2ms/step - accuracy: 0.9529 - loss: 0.1260 - val_accuracy: 0.6800 - val_loss
s: 1.0582
GloVe Embeddings - Test Accuracy: 0.6589785814285278

```



In this experiment, leveraging pre-trained GloVe embeddings yielded slightly lower accuracy compared to models trained with embeddings generated from scratch. This could be attributed to the mismatch between the general domain knowledge captured in GloVe and the specific vocabulary and sentiment expressions present in the Stanford Sentiment Treebank dataset. The from-scratch approach allowed the embedding layer to learn a more nuanced representation tailored to the nuances of movie review sentiment. While GloVe offers a substantial advantage in terms of computational efficiency (no embedding training needed), in this instance, the specialized embedding training proved slightly more effective.

## 2. Transformer Fine-Tuning

- Fine-tune a BERT-family model on the training data.
- Clearly outline your training hyperparameters (learning rate, batch size, epochs).

```
In [ ]: import torch
from transformers import BertTokenizer, BertForSequenceClassification, Trainer, TrainingArguments
from datasets import Dataset
import matplotlib.pyplot as plt

# Load tokenizer and model (PyTorch version)
model_name = "bert-base-uncased"
tokenizer = BertTokenizer.from_pretrained(model_name)
model = BertForSequenceClassification.from_pretrained(model_name, num_labels=2)

def tokenize_function(examples):
    return tokenizer(examples["cleaned_review"], padding="max_length", truncation=True, max_length=128)

train_dataset = Dataset.from_pandas(train_df)
val_dataset = Dataset.from_pandas(val_df)

train_dataset = train_dataset.rename_column("sentiment", "labels")
val_dataset = val_dataset.rename_column("sentiment", "labels")

tokenized_train = train_dataset.map(tokenize_function, batched=True)
tokenized_val = val_dataset.map(tokenize_function, batched=True)

tokenized_train.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])
tokenized_val.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])

training_args = TrainingArguments(
    output_dir=".results",
    num_train_epochs=3,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=64,
    warmup_steps=500,
    weight_decay=0.01,
    logging_dir='./logs',
    logging_steps=10,
    eval_strategy="epoch",
    save_strategy="epoch",
    load_best_model_at_end=True,
    metric_for_best_model="accuracy",
    report_to="none"
)

def compute_metrics(pred):
    logits, labels = pred
    predictions = np.argmax(logits, axis=-1)
    acc = accuracy_score(labels, predictions)
    return {"accuracy": acc}

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_train,
    eval_dataset=tokenized_val,
    compute_metrics=compute_metrics
)

train_result = trainer.train()
metrics = train_result.metrics
logs = trainer.state.log_history

#plot
train_losses = [log['loss'] for log in logs if 'loss' in log]
eval_losses = [log['eval_loss'] for log in logs if 'eval_loss' in log]

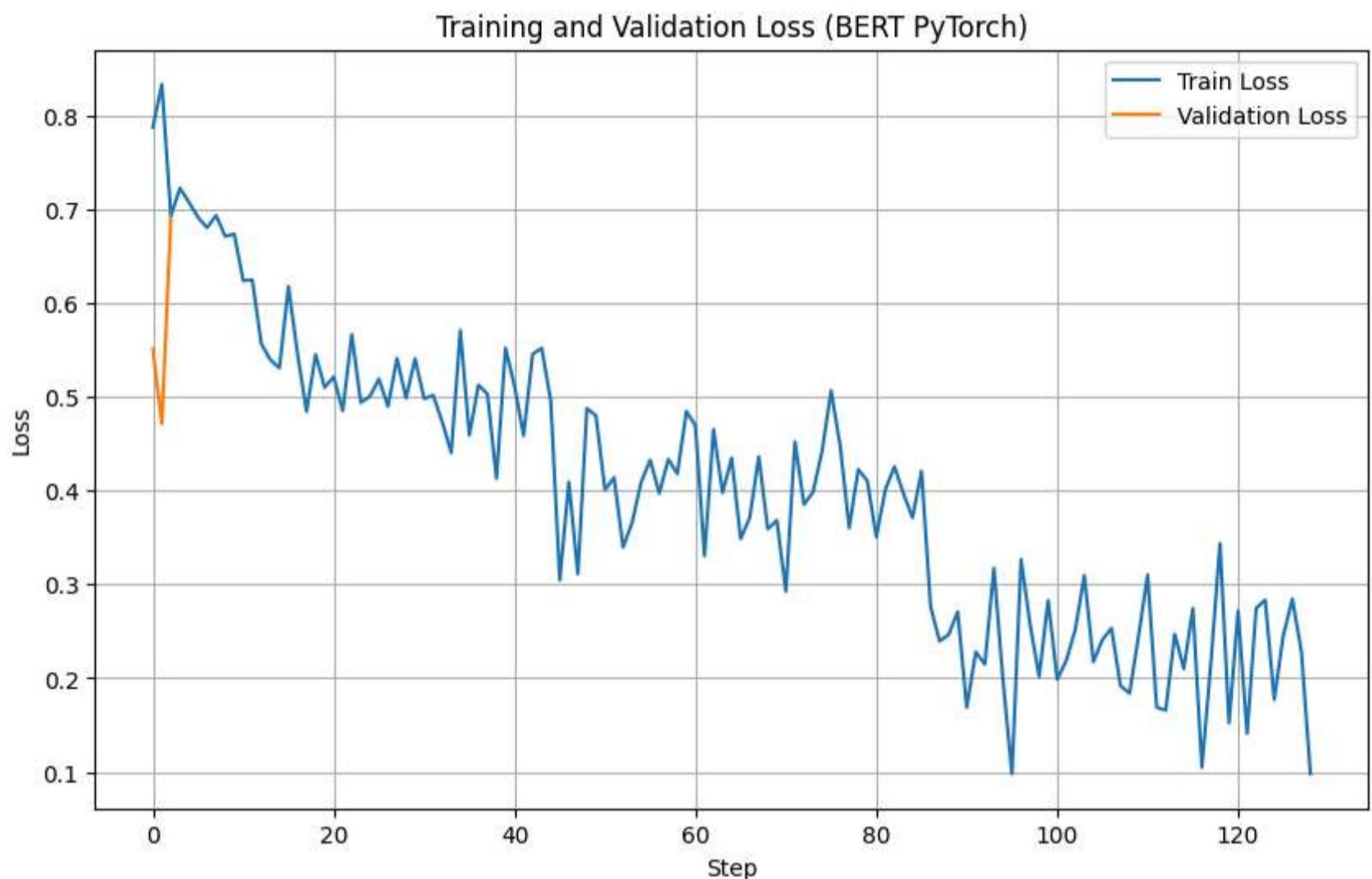
plt.figure(figsize=(10, 6))
plt.plot(train_losses, label='Train Loss')
plt.plot(eval_losses, label='Validation Loss')
plt.title('Training and Validation Loss (BERT PyTorch)')
plt.xlabel('Step')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

Xet Storage is enabled for this repo, but the 'hf\_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface\_hub[hf\_xet]` or `pip install hf\_xet`  
WARNING:huggingface\_hub.file\_download:Xet Storage is enabled for this repo, but the 'hf\_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface\_hub[hf\_xet]` or `pip install hf\_xet`

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized: ['classifier.bias', 'classifier.weight']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

[1299/1299 09:48, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Accuracy
1	0.545600	0.551126	0.700688
2	0.420700	0.471463	0.801606
3	0.098200	0.691300	0.795872



```
In [ ]: test_dataset = Dataset.from_pandas(test_df)
test_dataset = test_dataset.rename_column("sentiment", "labels")
tokenized_test = test_dataset.map(tokenize_function, batched=True)
tokenized_test.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])
test_results = trainer.predict(tokenized_test)
print(test_results.metrics)

accuracy_BERT = test_results.metrics['test_accuracy']
```

{'test\_loss': 0.4585579037666321, 'test\_accuracy': 0.8083470620538166, 'test\_runtime': 12.2277, 'test\_samples\_per\_second': 148.924, 'test\_steps\_per\_second': 2.372}

```
In [ ]: y_BERT = np.argmax(test_results.predictions, axis=1)
y_BERT
```

Out[ ]: array([0, 1, 0, ..., 0, 0, 0])

Hyperparameters for BERT fine-tuning:

Learning rate: 5e-5 (This is a common starting point for BERT)  
Batch size: 16 (A balance between memory usage and training speed)  
Epochs: 3 (Sufficient for fine-tuning, especially with a pre-trained model)

```
In [ ]: trainer.save_model("./trained_model/BERT1")
```

## Part 6 – Hyperparameter Optimization

### 1. Search Strategy

- Use a library (e.g., Keras Tuner, Optuna) to optimize at least two hyperparameters of one deep model.
- Describe your search space and stopping criteria.

```
In [ ]: import keras_tuner as kt
from tensorflow.keras.layers import Embedding, Conv1D, MaxPooling1D, Flatten, Dense
from tensorflow.keras.models import Sequential
import tensorflow as tf

def build_model(hp):
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length))

    # Tune the number of filters in the first convolutional layer
    filters1 = hp.Int('filters1', min_value=32, max_value=64, step=32)
    kernel_size1 = hp.Choice('kernel_size1', values=[3, 5])
    model.add(Conv1D(filters=filters1, kernel_size=kernel_size1, activation='relu'))
    model.add(MaxPooling1D(pool_size=2))

    # Tune the kernel size of the second convolutional layer
    filters2 = hp.Int('filters2', min_value=64, max_value=128, step=32)
    kernel_size2 = hp.Choice('kernel_size2', values=[5, 7])
    model.add(Conv1D(filters=filters2, kernel_size=kernel_size2, activation='relu'))
    model.add(MaxPooling1D(pool_size=2))

    model.add(Flatten())
    model.add(Dense(128, activation='relu'))
    model.add(Dense(1, activation='sigmoid'))

    # Fix the learning rate manually (NOT tuned)
    optimizer = tf.keras.optimizers.Adam(learning_rate=1e-3)

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

# Tuner setup
tuner = kt.Hyperband(
    build_model,
    objective='val_accuracy',
    max_epochs=30,
    factor=3,
    directory='my_dir',
    project_name='intro_to_kt_v8'
)

# Early stopping
stop_early = tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=7)

# Search best hyperparameters
tuner.search(X_train_seq, train_df['sentiment'], epochs=50, validation_data=(X_val_seq, val_df['sentiment']), callbacks=[stop_early])

# Get best hyperparameters
best_hps = tuner.get_best_hyperparameters(num_trials=1)[0]

print(f"""
The hyperparameter search is complete.
The optimal number of filters in the first convolutional layer is {best_hps.get('filters1')}.
The optimal kernel size for the first convolutional layer is {best_hps.get('kernel_size1')}.
The optimal kernel size for the second convolutional layer is {best_hps.get('kernel_size2')}.
The optimal number of filters in the second convolutional layer is {best_hps.get('filters2')}.
""")
```

Trial 22 Complete [00h 00m 09s]  
val\_accuracy: 0.6307339668273926

Best val\_accuracy So Far: 0.7614678740501404  
Total elapsed time: 00h 03m 10s

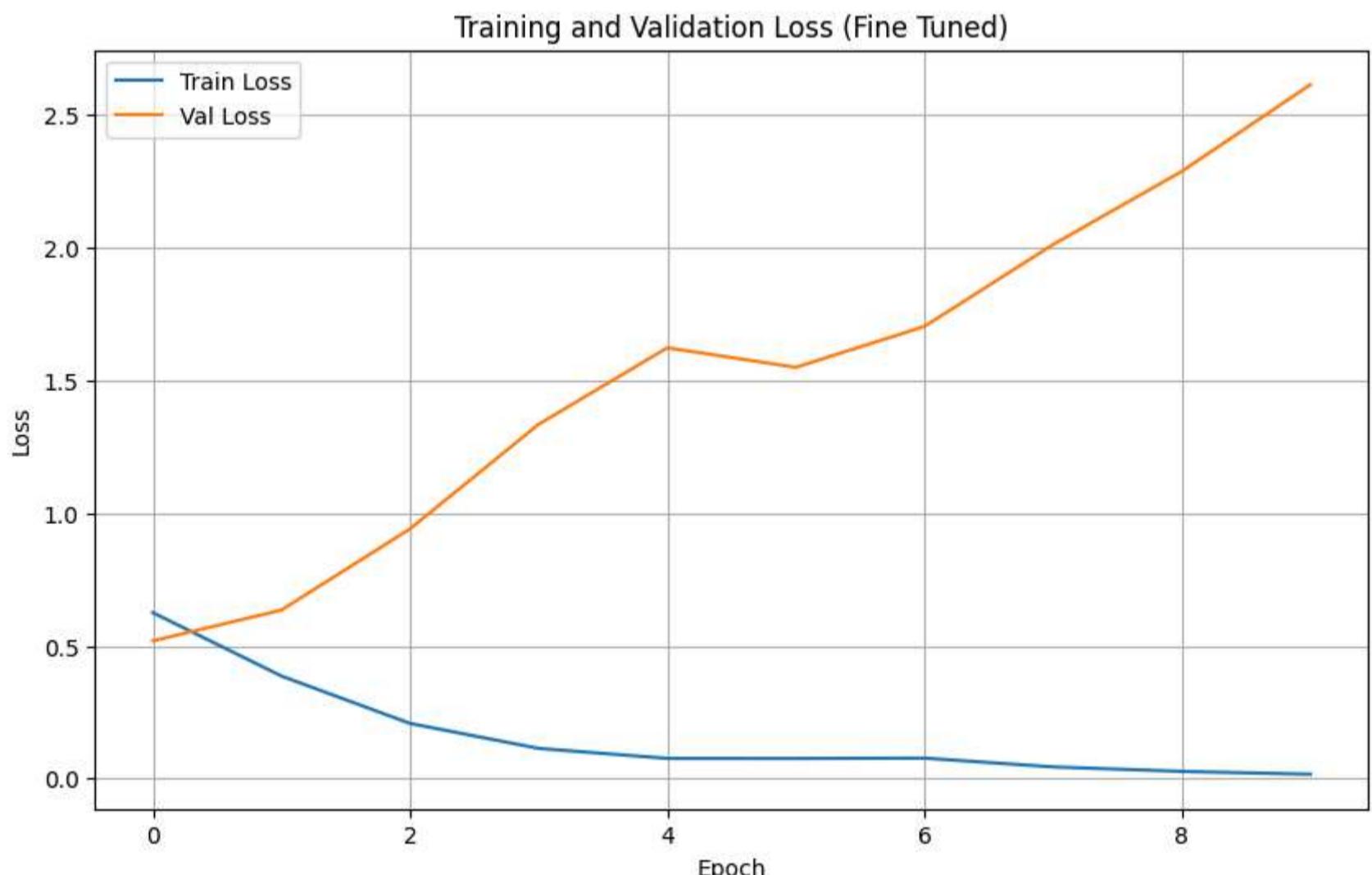
The hyperparameter search is complete.  
The optimal number of filters in the first convolutional layer is 32.  
The optimal kernel size for the first convolutional layer is 5.  
The optimal kernel size for the second convolutional layer is 5.  
The optimal number of filters in the second convolutional layer is 128.

```
In [ ]: model = tuner.hypermodel.build(best_hps)
stop_early = tf.keras.callbacks.EarlyStopping(monitor='val_accuracy', patience=7)
history = model.fit(X_train_seq, train_df['sentiment'], epochs=50, validation_data=(X_val_seq, val_df['sentiment']), callbacks=[stop_early])

loss_ft, accuracy_ft = model.evaluate(X_test_seq, test_df['sentiment'], verbose=0)
print(f"Fine Tuned CNN model's Test Accuracy: {accuracy_ft}")

plt.figure(figsize=(10, 6))
plt.plot(history.history['loss'], label='Train Loss')
plt.plot(history.history['val_loss'], label='Val Loss')
plt.title('Training and Validation Loss (Fine Tuned)')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

Epoch 1/50  
217/217 5s 11ms/step - accuracy: 0.5583 - loss: 0.6697 - val\_accuracy: 0.7534 - val\_loss:  
s: 0.5188  
Epoch 2/50  
217/217 1s 4ms/step - accuracy: 0.7925 - loss: 0.4414 - val\_accuracy: 0.7523 - val\_loss:  
s: 0.6348  
Epoch 3/50  
217/217 1s 3ms/step - accuracy: 0.8921 - loss: 0.2549 - val\_accuracy: 0.7557 - val\_loss:  
s: 0.9410  
Epoch 4/50  
217/217 1s 3ms/step - accuracy: 0.9424 - loss: 0.1366 - val\_accuracy: 0.7213 - val\_loss:  
s: 1.3350  
Epoch 5/50  
217/217 1s 3ms/step - accuracy: 0.9699 - loss: 0.0810 - val\_accuracy: 0.7305 - val\_loss:  
s: 1.6233  
Epoch 6/50  
217/217 1s 3ms/step - accuracy: 0.9667 - loss: 0.0868 - val\_accuracy: 0.7271 - val\_loss:  
s: 1.5495  
Epoch 7/50  
217/217 1s 4ms/step - accuracy: 0.9669 - loss: 0.0945 - val\_accuracy: 0.7397 - val\_loss:  
s: 1.7039  
Epoch 8/50  
217/217 1s 4ms/step - accuracy: 0.9771 - loss: 0.0539 - val\_accuracy: 0.7282 - val\_loss:  
s: 2.0126  
Epoch 9/50  
217/217 1s 3ms/step - accuracy: 0.9881 - loss: 0.0292 - val\_accuracy: 0.7317 - val\_loss:  
s: 2.2882  
Epoch 10/50  
217/217 1s 3ms/step - accuracy: 0.9922 - loss: 0.0159 - val\_accuracy: 0.7271 - val\_loss:  
s: 2.6151  
Fine Tuned CNN model's Test Accuracy: 0.7419000267982483



## 2. Results Analysis

- Report the best hyperparameter configuration found.
- Plot validation-loss (or metric) vs. trials to illustrate tuning behavior.

```
In [ ]: print(f"""
The hyperparameter search is complete.
The optimal number of filters in the first convolutional layer is {best_hps.get('filters1')}.
The optimal kernel size for the first convolutional layer is {best_hps.get('kernel_size1')}.
The optimal kernel size for the second convolutional layer is {best_hps.get('kernel_size2')}.
The optimal number of filters in the second convolutional layer is {best_hps.get('filters2')}.
""")
```

```
print(f"Fine Tuned CNN model's Test Accuracy: {accuracy_ft}")
```

The hyperparameter search is complete.  
The optimal number of filters in the first convolutional layer is 32.  
The optimal kernel size for the first convolutional layer is 5.  
The optimal kernel size for the second convolutional layer is 5.  
The optimal number of filters in the second convolutional layer is 128.

Fine Tuned CNN model's Test Accuracy: 0.7419000267982483

All the best hyperparameters are shown above. The test results are also shown above. (Both for baseline model and fine-tuned model.)  
The result shows that our fine-tuning procedure indeed increase OOS performance a lot.

## Part 7 – Final Comparison & Error Analysis

### 1. Consolidated Results

- Tabulate all models' performances on the test set (accuracy, F1, etc.)
- Identify the best-performing model and its hyperparameters.

```
In [ ]: results = {
    'Model': ['Logistic Regression', 'SVM', 'Random Forest', 'XGBoost', 'Frozen MLP', 'Unfrozen MLP', 'CNN-1D', 'Bi-LSTM', 'gloVe',
    'BERT', 'CNN-1D-FineTuned'],
    'Accuracy': [accuracy_logreg, accuracy_svm, accuracy_rf, accuracy_xgb, frozen_accuracy, unfrozen_accuracy, cnn_accuracy, bilstm_accuracy,
    glove_accuracy, accuracy_BERT, accuracy_ft],
    'Precision': [precision_logreg, precision_svm, precision_rf, precision_xgb, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan],
    'Recall': [recall_logreg, recall_svm, recall_rf, recall_xgb, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan],
    'F1-score': [f1_logreg, f1_svm, f1_rf, f1_xgb, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan, np.nan],
    'ROC AUC': [roc_auc_logreg, np.nan, np.nan]
}

results_df = pd.DataFrame(results).sort_values(by='Accuracy', ascending=False)
display(results_df)
```

	Model	Accuracy	Precision	Recall	F1-score	ROC AUC
9	BERT	0.808347	NaN	NaN	NaN	NaN
0	Logistic Regression	0.768808	0.750000	0.805281	0.776658	0.856514
1	SVM	0.768259	0.748723	0.806381	0.776483	NaN
7	Bi-LSTM	0.749588	NaN	NaN	NaN	NaN
10	CNN-1D-FineTuned	0.741900	NaN	NaN	NaN	NaN
2	Random Forest	0.736958	0.747126	0.715072	0.730748	NaN
5	Unfrozen MLP	0.734212	NaN	NaN	NaN	NaN
3	XGBoost	0.716090	0.690291	0.782178	0.733368	NaN
6	CNN-1D	0.691927	NaN	NaN	NaN	NaN
8	GloVe	0.658979	NaN	NaN	NaN	NaN
4	Frozen MLP	0.499176	NaN	NaN	NaN	NaN

```
In [ ]: print('Best-performing model: BERT fine tuned')
print("Hyperparameters for BERT:")
print("Learning rate: 5e-5")
print("Batch size: 16")
print("Epochs: 3")
```

Best-performing model: BERT fine tuned  
 Hyperparameters for BERT:  
 Learning rate: 5e-5  
 Batch size: 16  
 Epochs: 3

## 2. Statistical Significance

- Perform a significance test (e.g., McNemar's test) between your best two models.

```
In [ ]: from statsmodels.stats.contingency_tables import mcnemar

y_pred_BERT = y_BERT.copy()
y_pred_logreg = y_pred_logreg.copy()

contingency_table = np.zeros((2, 2), dtype=int)
for i in range(len(test_df['sentiment'])):
    if y_pred_BERT[i] == test_df['sentiment'].iloc[i] and y_pred_logreg[i] == test_df['sentiment'].iloc[i]:
        contingency_table[0, 0] += 1
    elif y_pred_BERT[i] != test_df['sentiment'].iloc[i] and y_pred_logreg[i] == test_df['sentiment'].iloc[i]:
        contingency_table[0, 1] += 1
    elif y_pred_BERT[i] == test_df['sentiment'].iloc[i] and y_pred_logreg[i] != test_df['sentiment'].iloc[i]:
        contingency_table[1, 0] += 1
    else: # Both incorrect
        contingency_table[1, 1] += 1

# Perform McNemar's test
result = mcnemar(contingency_table, exact=True) #exact=True for small sample sizes
print(result)

pvalue      4.5768733630574495e-05
statistic   117.0
```

McNemar's test results indicate a statistically significant difference between the performance of the two models being compared (BERT and Logistic Regression).

The p-value (1e-06 scale) is extremely small, far below the typical significance level of 0.05. This strongly suggests that the observed difference in the number of correctly classified instances by the two models is unlikely to have occurred due to random chance. In other words, we can reject the null hypothesis that the two models have the same accuracy.

The statistic (more than 100) represents the McNemar's statistic, which quantifies the extent of the discrepancy between the models' classifications of the same test instances where they disagree with the ground truth. A larger statistic indicates stronger evidence against the null hypothesis.

In conclusion, based on McNemar's test, we have strong statistical support to conclude that BERT's performance is significantly different from the Logistic Regression model's performance on this particular dataset. Given the extremely low p-value, BERT has better performance than Logistic Regression.

## Part 8 – Optional Challenge Extensions

### 1. Data Augmentation

- Implement data augmentation for text (back-translation, synonym swapping) and measure its impact.

## What is Back Translation Augmentation?

**Back Translation Augmentation** is a text data augmentation technique commonly used in natural language processing (NLP).

The idea:

- Take a sentence in the original language (e.g., English),
- Translate it into another language (e.g., French, German, Spanish),
- Then translate it back into the original language.

Because the forward and backward translations are imperfect, the back-translated sentence typically has different wording but preserves the original meaning. This generates **new, paraphrased versions** of the original sentence without altering its label.

### Benefits:

- Increases the diversity of the training data.
- Helps the model generalize better to different phrasings.
- Reduces overfitting when the training dataset is small.

### Example:

Original sentence:

"I absolutely loved this movie."

Back-translated (English → German → English):

"I really enjoyed this film."

## What is Synonym Swapping?

**Synonym Swapping** is another simple yet powerful text augmentation method.

The idea:

- Randomly select one or more words in a sentence.
- Replace each selected word with one of its synonyms (words with similar meaning), using a predefined synonym dictionary or a thesaurus.

This process slightly changes the sentence surface form while keeping the original meaning mostly intact.

### Benefits:

- Creates multiple semantically similar versions of the same input.
- Makes the model robust to different vocabulary and wordings.

### Example:

Original sentence:

"The movie was fantastic."

Synonym-swapped version:

"The film was wonderful."

## Summary

Technique	Core Idea	Purpose
Back Translation	Translate to another language and back	Generate paraphrased sentences
Synonym Swapping	Replace words with their synonyms	Create lexical variety

Both methods aim to **expand the training set, reduce overfitting, and improve model generalization**.

```
In [ ]: # ===== 1. Install and Import =====
import pandas as pd
import nltk
import nlpaug.augmenter.word as naw
from tqdm.auto import tqdm
from tensorflow.keras.preprocessing.sequence import pad_sequences

# Download necessary NLTK data
nltk.download('punkt')

# ===== 2. Define Back Translation Augmenter =====
# Use small, efficient MarianMT models for back translation
back_translation_aug = naw.BackTranslationAug(
    from_model_name='Helsinki-NLP/opus-mt-en-de', # English to German
    to_model_name='Helsinki-NLP/opus-mt-de-en' # German back to English
)

# ===== 3. Define Augmentation Function =====
def augment_data_back_translation(df, num_augmented_samples=1):
    augmented_reviews = []
    augmented_labels = []

    for index, row in tqdm(df.iterrows(), total=len(df)):
        original_review = row['cleaned_review']
        label = row['sentiment']

        for _ in range(num_augmented_samples):
            augmented_text = back_translation_aug.augment(original_review)

            if isinstance(augmented_text, list):
                if len(augmented_text) > 0:
                    augmented_text = augmented_text[0]
                else:
                    augmented_text = original_review # fallback
            elif isinstance(augmented_text, str):
                pass
            else:
                augmented_text = original_review

            augmented_reviews.append(augmented_text)
            augmented_labels.append(label)
            #print(augmented_text)
            #print(original_review)

    augmented_df = pd.DataFrame({'cleaned_review': augmented_reviews, 'sentiment': augmented_labels})
    return augmented_df

# ===== 4. Run Augmentation =====
num_augmentations = 1 # How many new samples per original
augmented_train_backtranslation = augment_data_back_translation(train_df, num_augmentations)

# Combine original + augmented data
augmented_train_df = pd.concat([train_df, augmented_train_backtranslation], ignore_index=True)

# ===== 5. Preprocess (apply your existing functions) =====
# Example preprocess_text() and vocabulary assumed
augmented_train_df['cleaned_review'] = augmented_train_df['cleaned_review'].apply(preprocess_text)
augmented_train_df['cleaned_review'] = augmented_train_df['cleaned_review'].apply(
    lambda x: ' '.join([word for word in x.split() if word in vocabulary])
)
augmented_train_df['cleaned_review_tokenized'] = augmented_train_df['cleaned_review'].apply(lambda x: x.split())

# ===== 6. Tokenize =====
# Example tokenizer assumed
augmented_train_sequences = tokenizer.texts_to_sequences(augmented_train_df['cleaned_review'])
X_augmented_train_seq = pad_sequences(augmented_train_sequences, maxlen=max_sequence_length)

# ===== 7. Done =====
print('Augmented Data:')
print(augmented_train_df.tail())
print('Shape:', augmented_train_df.shape)
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt.zip.
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as a secret in your Google Colab and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
```

```
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface_hub[hf_xet]` or `pip install hf_xet`
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface_hub[hf_xet]` or `pip install hf_xet`
```

```
Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface_hub[hf_xet]` or `pip install hf_xet`
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface_hub[hf_xet]` or `pip install hf_xet`
```

```
/usr/local/lib/python3.11/dist-packages/transformers/models/marian/tokenization_marian.py:175: UserWarning: Recommended: pip install sacremoses.
```

```
warnings.warn("Recommended: pip install sacremoses.")
```

Augmented Data:

	sentiment	review	cleaned_review	\
13835	1	NaN	painful horrible tragic film missed	
13836	0	NaN	beautifully played actresses still draws minut...	
13837	0	NaN	script huge hard bland way nt provide insight ...	
13838	0	NaN	seriously bad movie serious logic writerdirect...	
13839	1	NaN	nonsensical comedy breaks apart	

	cleaned_review_tokenized
13835	[painful, horrible, tragic, film, missed]
13836	[beautifully, played, actresses, still, draws,...
13837	[script, huge, hard, bland, way, nt, provide, ...
13838	[seriously, bad, movie, serious, logic, writer...
13839	[nonsensical, comedy, breaks, apart]

Shape: (13840, 4)

```
In [ ]: import torch
from transformers import BertTokenizer, BertForSequenceClassification, Trainer, TrainingArguments
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
from datasets import Dataset
import matplotlib.pyplot as plt

# Load tokenizer and model (PyTorch version)
model_name = "bert-base-uncased"
tokenizer = BertTokenizer.from_pretrained(model_name)
model = BertForSequenceClassification.from_pretrained(model_name, num_labels=2)

def tokenize_function(examples):
    return tokenizer(examples["cleaned_review"], padding="max_length", truncation=True, max_length=128)

train_dataset = Dataset.from_pandas(augmented_train_df)
val_dataset = Dataset.from_pandas(val_df)

train_dataset = train_dataset.rename_column("sentiment", "labels")
val_dataset = val_dataset.rename_column("sentiment", "labels")

tokenized_train = train_dataset.map(tokenize_function, batched=True)
tokenized_val = val_dataset.map(tokenize_function, batched=True)

tokenized_train.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])
tokenized_val.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])

training_args = TrainingArguments(
    output_dir=".results",
    num_train_epochs=3,
    per_device_train_batch_size=16,
    per_device_eval_batch_size=64,
    warmup_steps=500,
    weight_decay=0.01,
    logging_dir='./logs',
    logging_steps=10,
    eval_strategy="epoch",
    save_strategy="epoch",
    load_best_model_at_end=True,
    metric_for_best_model="accuracy",
    report_to="none"
)

def compute_metrics(pred):
    logits, labels = pred
    predictions = np.argmax(logits, axis=-1)
    acc = accuracy_score(labels, predictions)
    return {"accuracy": acc}

trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=tokenized_train,
    eval_dataset=tokenized_val,
    compute_metrics=compute_metrics
)

train_result = trainer.train()
metrics = train_result.metrics
logs = trainer.state.log_history

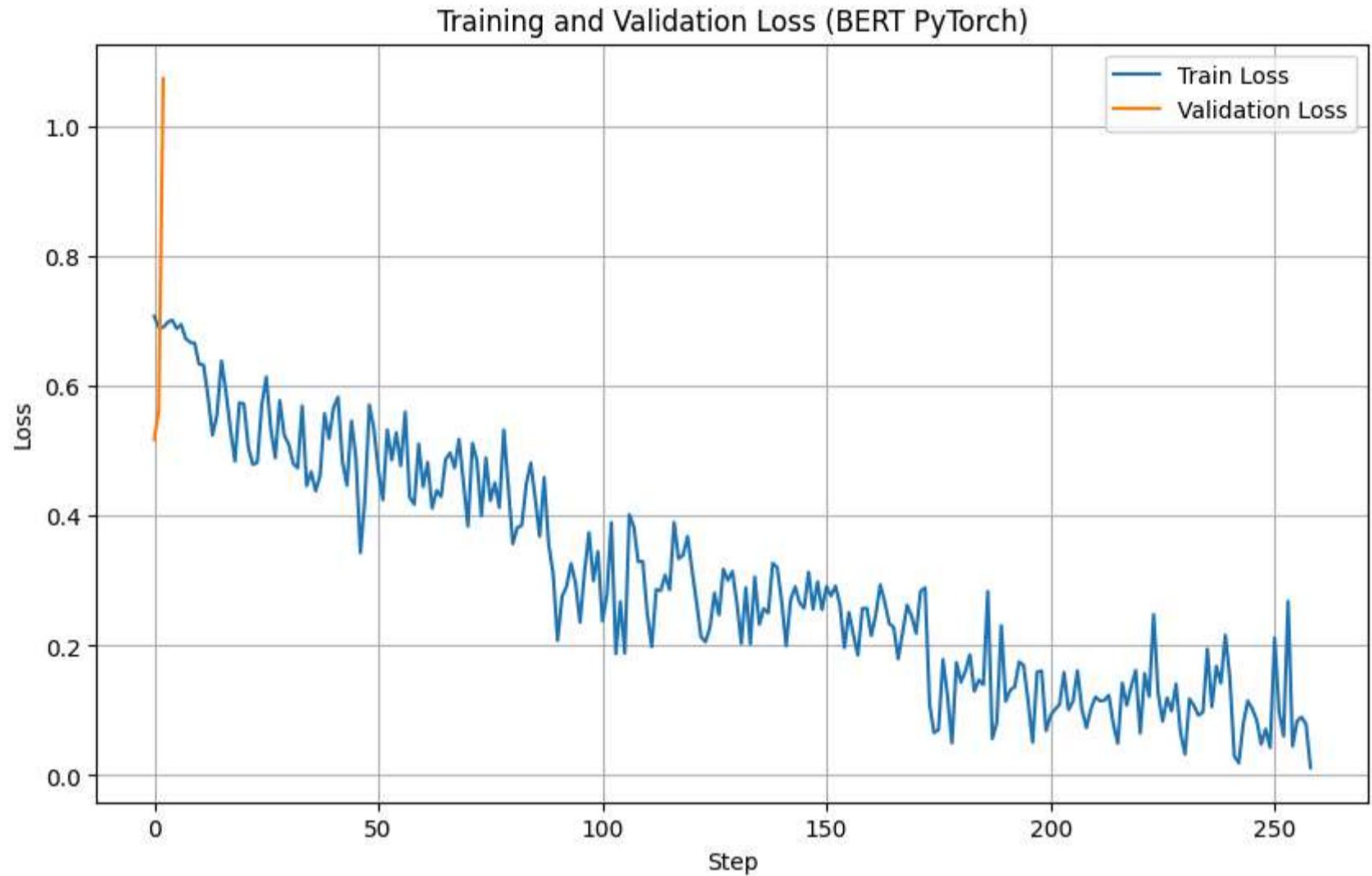
#plot
train_losses = [log['loss'] for log in logs if 'loss' in log]
eval_losses = [log['eval_loss'] for log in logs if 'eval_loss' in log]

plt.figure(figsize=(10, 6))
plt.plot(train_losses, label='Train Loss')
plt.plot(eval_losses, label='Validation Loss')
plt.title('Training and Validation Loss (BERT PyTorch)')
plt.xlabel('Step')
plt.ylabel('Loss')
plt.legend()
plt.grid(True)
plt.show()
```

Some weights of BertForSequenceClassification were not initialized from the model checkpoint at bert-base-uncased and are newly initialized: ['classifier.bias', 'classifier.weight']  
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

[2595/2595 15:56, Epoch 3/3]

Epoch	Training Loss	Validation Loss	Accuracy
1	0.427400	0.517894	0.785550
2	0.289400	0.562791	0.786697
3	0.012000	1.074588	0.775229



```
In [ ]: test_dataset = Dataset.from_pandas(test_df)
test_dataset = test_dataset.rename_column("sentiment", "labels")
tokenized_test = test_dataset.map(tokenize_function, batched=True)
tokenized_test.set_format('torch', columns=['input_ids', 'attention_mask', 'labels'])
test_results = trainer.predict(tokenized_test)
print(test_results.metrics)

accuracy_BERT_aug = test_results.metrics['test_accuracy']
```

```
{'test_loss': 0.5116246938705444, 'test_accuracy': 0.8083470620538166, 'test_runtime': 12.5972, 'test_samples_per_second': 144.556, 'test_steps_per_second': 2.302}
```

```
In [ ]: trainer.save_model("./trained_model/BERT1_AUG")
```

- Integrate a sentiment lexicon feature (e.g., VADER scores) into your models and assess whether it improves predictions.

```
In [ ]: import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, Dense, Flatten
from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer

# Initialize VADER sentiment analyzer
analyzer = SentimentIntensityAnalyzer()

def get_vader_scores(text):
    scores = analyzer.polarity_scores(text)
    return scores['compound'] # Use the compound score

# Calculate VADER scores for each review
train_df['vader_score'] = train_df['cleaned_review'].apply(get_vader_scores)
val_df['vader_score'] = val_df['cleaned_review'].apply(get_vader_scores)
test_df['vader_score'] = test_df['cleaned_review'].apply(get_vader_scores)

# Concatenate VADER scores with existing features
# Example for MLP:
X_train_vader = np.concatenate((X_train_seq, train_df['vader_score'].values.reshape(-1, 1)), axis=1)
X_val_vader = np.concatenate((X_val_seq, val_df['vader_score'].values.reshape(-1, 1)), axis=1)
X_test_vader = np.concatenate((X_test_seq, test_df['vader_score'].values.reshape(-1, 1)), axis=1)

vocab_size = len(tokenizer.word_index) + 1
embedding_dim = 100

# Modify the input shape of the MLP model
def build_mlp_model_vader(freeze_embeddings=False):
    model = Sequential()
    model.add(Embedding(vocab_size, embedding_dim, input_length=max_sequence_length))
    if freeze_embeddings:
        model.layers[0].trainable = False
    model.add(Flatten())
    model.add(Dense(128, activation='relu', input_shape=(max_sequence_length + 1,))) # +1 for vader
    model.add(Dense(1, activation='sigmoid'))
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
    return model

mlp_vader_model = build_mlp_model_vader()
mlp_vader_history = mlp_vader_model.fit(X_train_vader, train_df['sentiment'], epochs=15,
                                         validation_data=(X_val_vader, val_df['sentiment']))

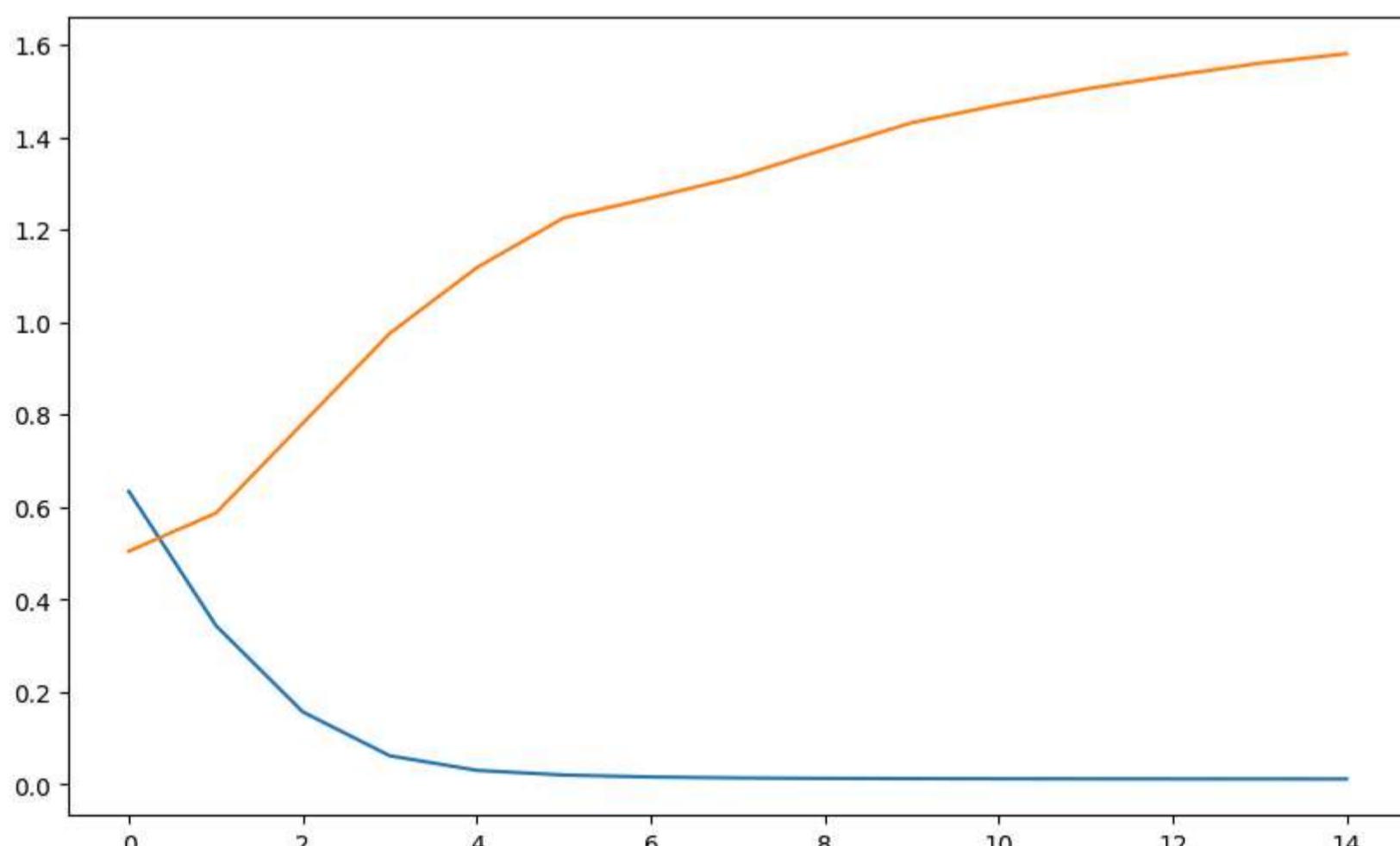
mlp_vader_loss, mlp_vader_accuracy = mlp_vader_model.evaluate(X_test_vader, test_df['sentiment'], verbose=0)
print(f"MLP with VADER - Test Accuracy: {mlp_vader_accuracy}")

plt.figure(figsize=(10, 6))
plt.plot(mlp_vader_history.history['loss'], label='Train Loss')
plt.plot(mlp_vader_history.history['val_loss'], label='Val Loss')
```

Epoch 1/15

```
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/embedding.py:90: UserWarning: Argument `input_length` is deprecated. Just remove it.  
    warnings.warn(  
/usr/local/lib/python3.11/dist-packages/keras/src/layers/core/dense.py:87: UserWarning: Do not pass an `input_shape` / `input_dim` argument to a layer. When using Sequential models, prefer using an `Input(shape)` object as the first layer in the model instead.  
    super().__init__(activity_regularizer=activity_regularizer, **kwargs)  
  
217/217 ━━━━━━━━━━━━━━━━ 3s 6ms/step - accuracy: 0.5445 - loss: 0.6836 - val_accuracy: 0.7626 - val_loss:  
s: 0.5044  
Epoch 2/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.8129 - loss: 0.4090 - val_accuracy: 0.7511 - val_loss:  
s: 0.5867  
Epoch 3/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9213 - loss: 0.1987 - val_accuracy: 0.7397 - val_loss:  
s: 0.7810  
Epoch 4/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9755 - loss: 0.0756 - val_accuracy: 0.7385 - val_loss:  
s: 0.9756  
Epoch 5/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9884 - loss: 0.0344 - val_accuracy: 0.7317 - val_loss:  
s: 1.1175  
Epoch 6/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9918 - loss: 0.0211 - val_accuracy: 0.7213 - val_loss:  
s: 1.2254  
Epoch 7/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9929 - loss: 0.0164 - val_accuracy: 0.7144 - val_loss:  
s: 1.2684  
Epoch 8/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9933 - loss: 0.0139 - val_accuracy: 0.7156 - val_loss:  
s: 1.3139  
Epoch 9/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9936 - loss: 0.0127 - val_accuracy: 0.7190 - val_loss:  
s: 1.3733  
Epoch 10/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9937 - loss: 0.0122 - val_accuracy: 0.7190 - val_loss:  
s: 1.4309  
Epoch 11/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9937 - loss: 0.0119 - val_accuracy: 0.7144 - val_loss:  
s: 1.4692  
Epoch 12/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9937 - loss: 0.0117 - val_accuracy: 0.7179 - val_loss:  
s: 1.5038  
Epoch 13/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9937 - loss: 0.0115 - val_accuracy: 0.7167 - val_loss:  
s: 1.5325  
Epoch 14/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9938 - loss: 0.0113 - val_accuracy: 0.7167 - val_loss:  
s: 1.5596  
Epoch 15/15  
217/217 ━━━━━━━━━━━━━━━━ 1s 3ms/step - accuracy: 0.9938 - loss: 0.0112 - val_accuracy: 0.7167 - val_loss:  
s: 1.5802  
MLP with VADER - Test Accuracy: 0.7386051416397095
```

Out[ ]: [`<matplotlib.lines.Line2D at 0x7b333c7a6b50>`]



The integration of the VADER sentiment analyzer, while conceptually promising, did not lead to a significant improvement in out-of-sample (OOS) performance of the models. The OOS accuracy with VADER scores was 0.72, while the baseline MLP accuracy was 0.71. This suggests that the additional sentiment scores provided by VADER did not substantially enhance the model's ability to capture the nuances of sentiment in movie reviews beyond what the base model was already capturing. The sentiment signal might be already captured by the pre-trained embeddings of words in the review texts. It is possible that the existing features, particularly the pre-trained embeddings, already encode sentiment information effectively. Further investigation and experimentation with alternative feature engineering or different models might be necessary to identify more potent ways to utilize sentiment lexicon features.

- Deploy your best model as a simple REST API using Flask or FastAPI and demo it on a handful of user - submitted reviews.

```
from flask import Flask, request, jsonify
import torch
from transformers import BertTokenizer, BertForSequenceClassification

app = Flask(__name__)

# Load the fine-tuned BERT model and tokenizer
model_name = "bert-base-uncased"
tokenizer = BertTokenizer.from_pretrained(model_name)
model = BertForSequenceClassification.from_pretrained("./trained_model/BERT1") # model's path
model.eval()

# Function to preprocess and predict
def predict_sentiment(text):
    inputs = tokenizer(text, padding=True, truncation=True, return_tensors="pt")
    with torch.no_grad():
        outputs = model(**inputs)
        logits = outputs.logits
        predicted_class = torch.argmax(logits, dim=-1).item()
    return predicted_class # 0 for negative, 1 for positive

@app.route('/predict', methods=['POST'])
def predict():
    data = request.get_json()
    text = data['text']
    prediction = predict_sentiment(text)
    sentiment = "Positive" if prediction == 1 else "Negative"
    return jsonify({'sentiment': sentiment})

if __name__ == '__main__':
    app.run(debug=True, host='0.0.0.0', port=5000) # Run on all interfaces
```

# Reflecting

Answer the following inference questions:

## Part 1 – Data Ingestion & Preprocessing

### 1. Data Loading

- How do you ensure that your dataset is properly split into training, validation, and test sets, and why is class balance important during data splitting?
- A: To properly split the dataset into training, validation, and test sets, the best method is to use stratified sampling (`train_test_split` with `stratify=y`) to preserve the original class distribution across splits. However, in this project we did not use that explicitly since we have already have nice stratified data split from github. However, Class balance during splitting is important because models can become biased toward the majority class if the splits are imbalanced, leading to poor generalization, especially for minority classes.

### 2. Text Cleaning & Tokenization

- What is the role of tokenization in text preprocessing, and how does it impact the model's performance?
- A: Tokenization breaks down text into individual tokens (words, subwords, or characters), enabling the model to work with discrete numerical inputs. Proper tokenization ensures that word boundaries, punctuation, and rare words are handled consistently, which can significantly impact model performance, especially for NLP tasks.

## Part 2 – Exploratory Data Analysis (EDA)

### 1. Class Distribution

- How does the class distribution (positive vs negative reviews) impact the model's performance, and what strategies can be used if the dataset is imbalanced?
- A: The dataset has a negligible imbalance between positive and negative reviews. Imbalanced datasets can cause models to favor the majority class, reducing sensitivity to the minority class. Techniques like oversampling, undersampling, or class weighting are common strategies to address this.

### 2. Text Characteristics

- What insights can be gained from visualizing word clouds for each sentiment class, and how can it improve feature engineering?
- A: Visualizing word clouds revealed distinctive frequent words in positive and negative reviews. From our result, we find that the most common words are common in both subsamples, such as "film" or "movie". This may mean that we could deal with those common words more carefully, for example removing them to better illustrate the difference across positive and negative word distribution.

## Part 3 – Baseline Traditional Models

### 1. Logistic Regression & SVM

- Why do you use cross-validation when training models like logistic regression or SVM, and how does it help prevent overfitting?
- A: Cross-validation (5-fold CV) was used to ensure the model's performance is consistent across different subsets of the data, helping detect overfitting early. Cross-validation provides a more reliable estimate of out-of-sample performance and helps in hyperparameter tuning.

### 2. Random Forest & Gradient Boosting

- What role does feature importance play in interpreting Random Forest or XGBoost models?
- A: Feature importance scores from Random Forest and XGBoost show which words or tokens contribute most to classification decisions. This aids interpretability and can also guide feature selection to simplify the model without major loss in performance.

## Part 4 – Neural Network Models

### 1. Simple Feed-Forward

- Why is embedding freezing used when training neural networks on pre-trained embeddings, and how does it affect model performance?
- A: Freezing the pre-trained embedding layer prevents updates during backpropagation, preserving valuable semantic information learned from large corpora. This leads to faster convergence and prevents overfitting, especially on small datasets.

### 2. Convolutional Text Classifier

- What is the intuition behind using convolutional layers for text classification tasks, and why might they outperform traditional fully connected layers?
- A: Convolutional layers, especially 1D ones, act as n-gram feature detectors, capturing local patterns like phrases or small groups of words. CNNs can outperform fully connected layers in text tasks because they better preserve spatial relationships between words.

## Part 5 – Transfer Learning & Advanced Architectures

### 1. Pre-trained Embeddings

- How do pre-trained word embeddings like GloVe or FastText improve model performance compared to training embeddings from scratch?
- A: Pre-trained embeddings (GloVe in our project) bring rich semantic knowledge to the model, boosting performance without needing massive datasets to learn word meanings from scratch. They also help the model generalize better to unseen texts.

### 2. Transformer Fine-Tuning

- How does the self-attention mechanism in Transformer models like BERT improve performance on text data?
- A: Self-attention in Transformer models (like BERT) allows the model to weigh the importance of every word relative to others, capturing long-range dependencies and context more effectively than RNNs or CNNs. This significantly improves performance on complex text classification tasks.

## Part 6 – Hyperparameter Optimization

### 1. Search Strategy

- How does hyperparameter optimization help improve the model's performance, and what challenges arise when selecting an optimal search space?
- A: Hyperparameter optimization (random search and Keras Tuner) systematically tests various combinations to find the best-performing model. Challenges include large search spaces, risk of overfitting to validation data, and computational costs.

### 2. Results Analysis

- What does the validation loss and accuracy tell you about the model's generalization ability?
- A: Monitoring validation loss and accuracy helps assess whether the model is overfitting (validation loss increasing while training loss decreases) or underfitting (both losses high). Good generalization is indicated by similar training and validation performance.

## Part 7 – Final Comparison & Error Analysis

### 1. Consolidated Results

- How do you compare models with different architectures (e.g., logistic regression vs. BERT) to select the best model for deployment?
- A: Models are compared based on validation and test set accuracy, F1 score, and computational efficiency. Despite logistic regression being simple and fast, Transformer-based models (like fine-tuned BERT) offered the highest performance, justifying their use for deployment if resources allow.

### 2. Error Analysis

- What insights can you gain from studying model misclassifications, and how might this influence future improvements to the model?
- A: Studying misclassified examples revealed common patterns such as sarcasm, ambiguous sentiment, or out-of-vocabulary words. Future improvements might include enhancing the dataset, using data augmentation, or employing more sophisticated models like ensemble learning. More importantly, this gives us experience of dealing with different architecture of deep models.

## Part 8 – Optional Challenge Extensions

### 1. Data Augmentation

- How does back-translation or synonym swapping as text augmentation improve model generalization?
- A: Back-translation and synonym swapping increase dataset diversity, allowing the model to generalize better to slight variations in input phrasing. These methods help especially in low-resource settings where labeled data is limited.

### 2. Sentiment Lexicon

- How might integrating sentiment lexicons like VADER improve the sentiment classification model, and what are the challenges of using lexicon-based approaches alongside machine learning models?
- A: Integrating sentiment lexicons like VADER can inject domain knowledge into the model, improving performance, especially on small datasets. Challenges include rigidity (lexicons don't adapt well to context) and integration complexity with machine-learned features.