



Customer Sentiment & Trend Analysis

Artificial Intelligence

Sentiment Analysis of
Product Reviews and Tweets
(English and Arabic)

Team members

- Ahmed Abdelhamed Hussein
 - Ahmed Awad Ata
 - Ahmed Mohamed Emad
- Ahmed Samy Farahat Abdelfatah
 - Amr Maher Abdallah Mohamed
 - Mohamed Wael Mohamed Khalifa

الفريق الثالث	الفريق الثاني	الفريق الأول
وائل + عبد الحميد	عمرو + عماد	فرحات + عوض
IMDB Movie Reviews	Amazon Reviews	Sentiment140

المهام	الأعضاء	الفريق
<ul style="list-style-type: none"> - استخدام MLflow لتبعد وإدارة نماذج التعلم الآلي. - إعداد التقارير الخاصة بإدارة النماذج. - ضمان تكامل النماذج مع البنية التحتية الحالية. 	Farahat - Wael - Emad -	MLOps فريق
<ul style="list-style-type: none"> - بناء وتنفيذ نموذج Generative Adversarial Network (GAN). - توليد بيانات اصطناعية لتحليل تعليقات العملاء. - اختبار وتحسين أداء النموذج. 	Amr - Abdelhamid - Awad -	GANs فريق
<ul style="list-style-type: none"> - إعداد التقرير النهائي للمشروع. - تحضير العرض التقديمي للعرض النهائي. - جمع المعلومات وتنسيق الشرائح. - مراجعة وتحرير المحتوى لضمان الجودة. 	Farahat - Amr - Wael - Abdelhamid - Awad - Emad -	فريق العرض النهائي والتقارير

TABLE OF CONTENTS

- Project Overview
- Data Sources
- Methodology
- Models & Techniques
- MLflow Integration
- Evaluation Metrics & Results
- Key Challenges & Solutions
- Future Work
- Conclusion



Project Overview

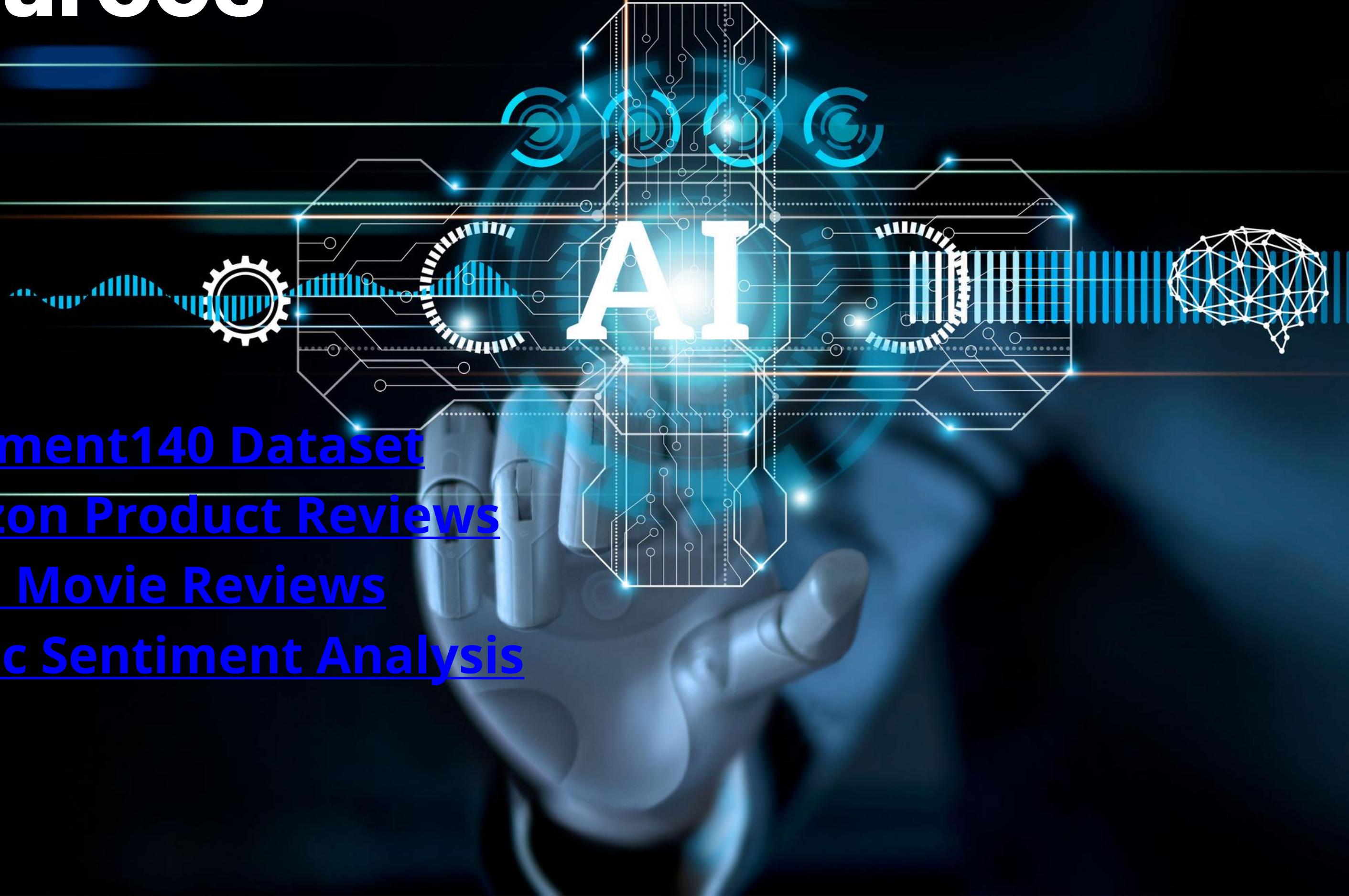
Objective: Analyze customer sentiment and uncover trends using datasets from product reviews and tweets in both English and Arabic.

Techniques: Apply NLP and machine learning techniques for sentiment classification (positive, negative, neutral).

Goal: Improve prediction accuracy using traditional machine learning and deep learning (e.g., transformers, LSTMs).

Data Sources

- Dataset 1: [Sentiment140 Dataset](#)
- Dataset 2: [Amazon Product Reviews](#)
- Dataset 3: [IMDB Movie Reviews](#)
- Dataset 4: [Arabic Sentiment Analysis](#)



METHODOLOGY

01

Data Preprocessing:
Cleaning, tokenization,
and padding to handle
both English and Arabic
text.

02

Feature Engineering: Using
word embeddings (GloVe
for English, custom
embeddings for Arabic) to
capture word semantics.

03

Model Training: Applying
ML algorithms and deep
learning models such as
LSTMs, transformers, and
RNNs.

Models & Techniques

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis tempor porttitor velit nec accumsan. Integer vehicula augue at purus maximus placerat.

In id nisl nec odio varius sodales. In pellentesque massa eu tempor congue. In in vestibulum nibh.



OUR PROJECTS

01

Sentiment Analysis of Twitter Data

02

Sentiment Analysis of IMDB Movie Reviews

03

Advanced Tracking and Prediction System

04

Sentiment Analysis of Arabic Tweets

05

Sentiment Analysis of Amazon reviews





Sentiment Analysis of *Twitter Data*

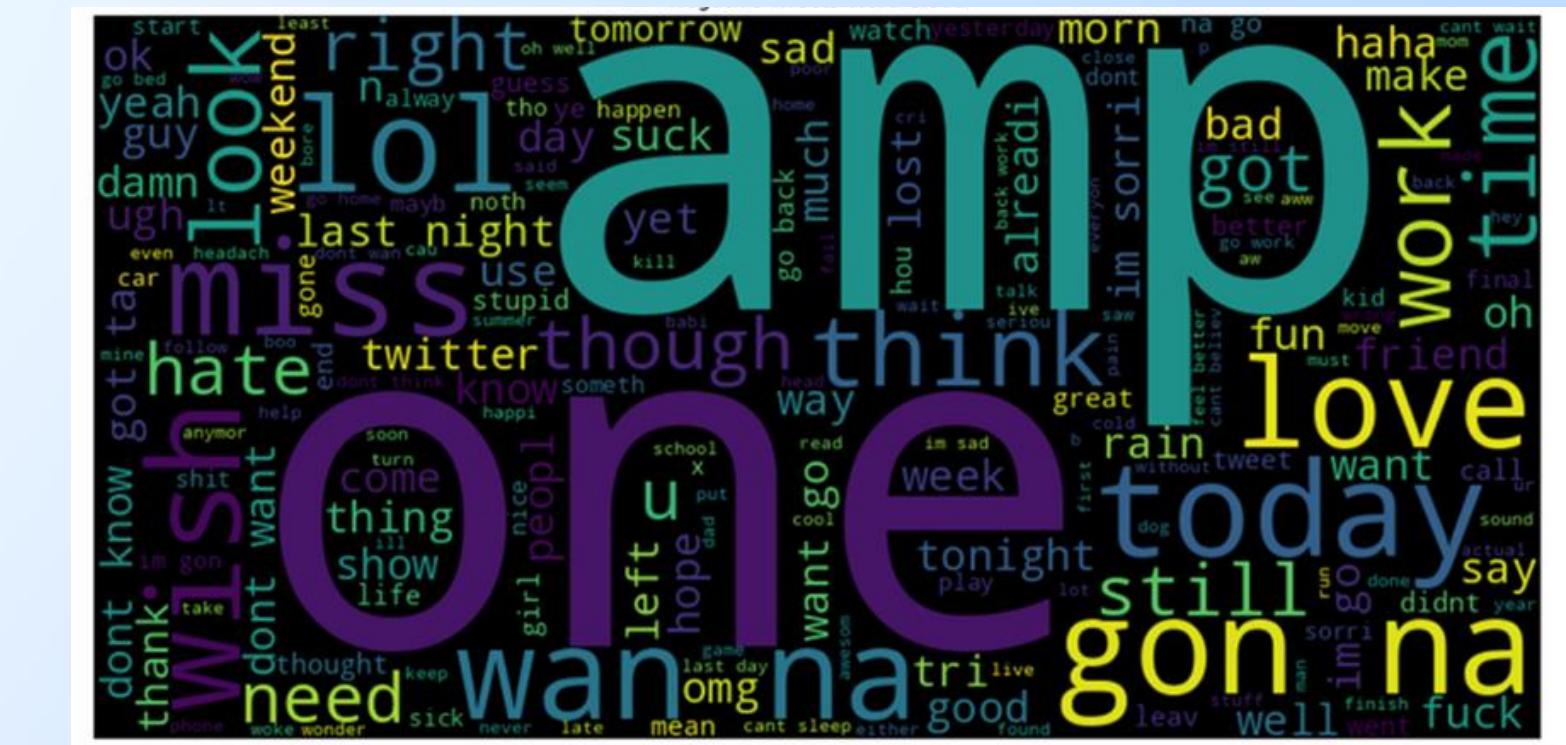
This project explores sentiment analysis using a dataset of 1.6 million tweets, categorized into positive and negative sentiments. The key steps include:

Data Preprocessing: Tweets were cleaned by removing URLs, mentions, numbers, and punctuation. Stopwords were also removed to improve model performance.

Positive Tweets Word Cloud



Negative Tweets Word Cloud

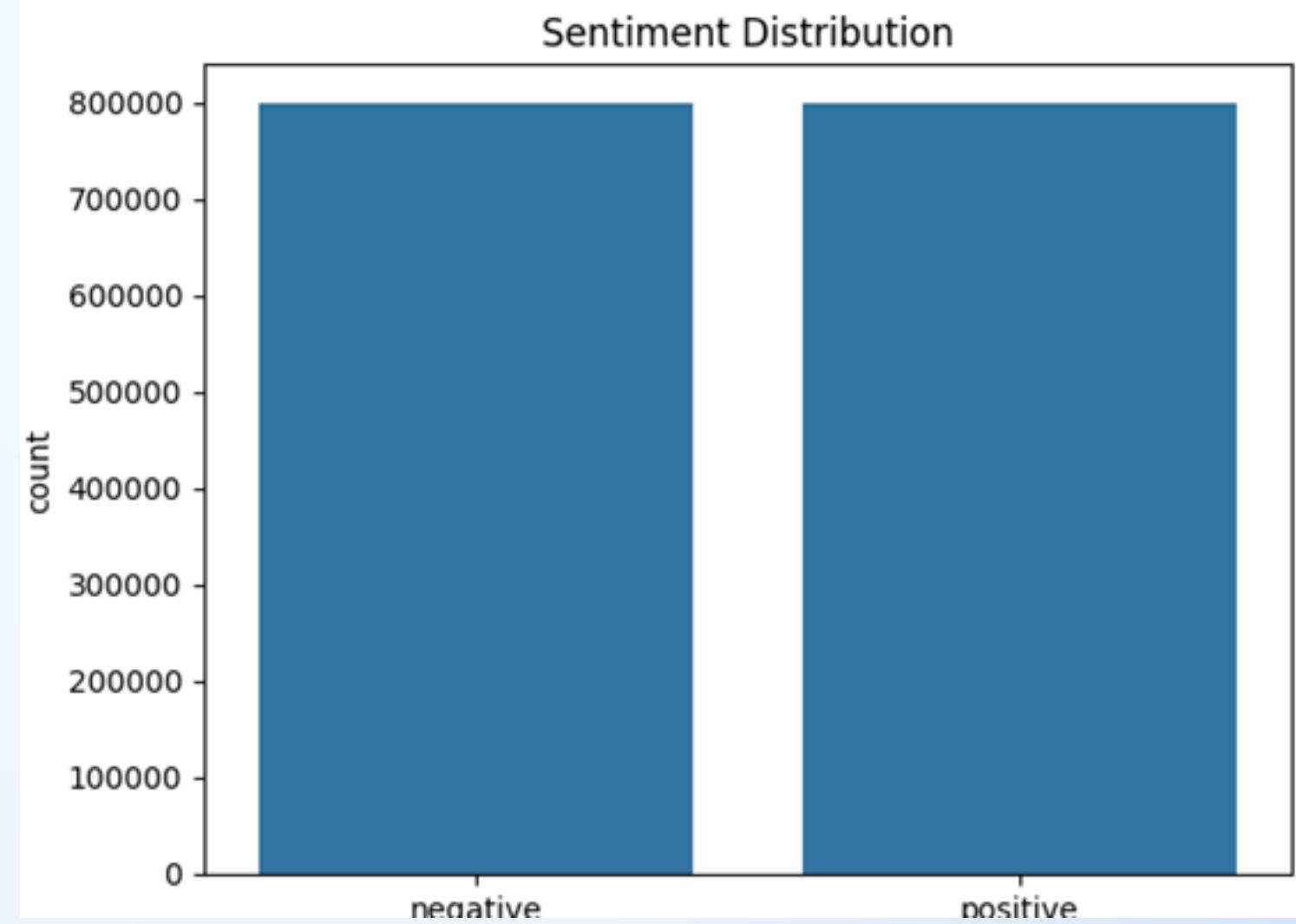


```
1      upset cant update facebook texting might cry r...
2      dived many times ball managed save rest go bounds
3                  whole body feels itchy like fire
4                  behaving im mad cant see
...
1599995          woke school best feeling ever
1599996      thewdbcom cool hear old walt interviews â« «
1599997          ready mojo makeover ask details
1599998      happy th birthday boo alll time tupac amaru sh...
1599999          happy charitytuesday
Name: cleaned_text, Length: 1600000, dtype: object

stemmer = PorterStemmer()
lemmatizer = WordNetLemmatizer()

df['cleaned_text'] = df['cleaned_text'].apply(lambda text: ' '.join([stemmer.stem(word) for word in text.split()]))
df['cleaned_text'] = df['cleaned_text'].apply(lambda text: ' '.join([lemmatizer.lemmatize(word) for word in text.split()]))

sns.countplot(x='sentiment', data=df)
plt.title('Sentiment Distribution')
plt.show()
```



Text Tokenization and Padding: The text data was tokenized and transformed into sequences, which were padded to ensure uniform input length for the model.

```
[10]: def get_most_common_words(texts, num=20):
    words = ' '.join(texts).split()
    counter = Counter(words)
    return counter.most_common(num)

print("Most common words in positive tweets:")
print(get_most_common_words(df[df['sentiment'] == 'positive']['cleaned_text']))

print("\nMost common words in negative tweets:")
print(get_most_common_words(df[df['sentiment'] == 'negative']['cleaned_text']))
```

Most common words in positive tweets:

```
[('im', 4590), ('love', 3893), ('good', 3859), ('go', 3698), ('day', 3485), ('than
k', 3159), ('get', 3072), ('like', 2509), ('u', 2322), ('lol', 2126), ('time', 210
7), ('got', 2019), ('work', 1844), ('today', 1818), ('know', 1758), ('see', 1747),
('one', 1744), ('watch', 1655), ('new', 1630), ('great', 1566)]
```

Most common words in negative tweets:

```
[('im', 6574), ('go', 4883), ('get', 3821), ('work', 3624), ('day', 3104), ('miss',
2978), ('cant', 2787), ('dont', 2748), ('like', 2578), ('got', 2446), ('want', 244
4), ('today', 2332), ('feel', 2257), ('back', 2000), ('realli', 1944), ('time', 189
7), ('u', 1822), ('one', 1796), ('sad', 1781), ('still', 1771)]
```

Feature Engineering: Both stemming and lemmatization techniques were applied to reduce words to their root form. Additionally, emoticons were replaced with corresponding textual representations.

```
: le = LabelEncoder()
y = le.fit_transform(y) # Encode positive: 1, negative: 0

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

: max_words = 10000 # Maximum number of words in the vocabulary
max_len = 100 # Maximum sequence Length

tokenizer = Tokenizer(num_words=max_words, oov_token='<OOV>')
tokenizer.fit_on_texts(X_train)

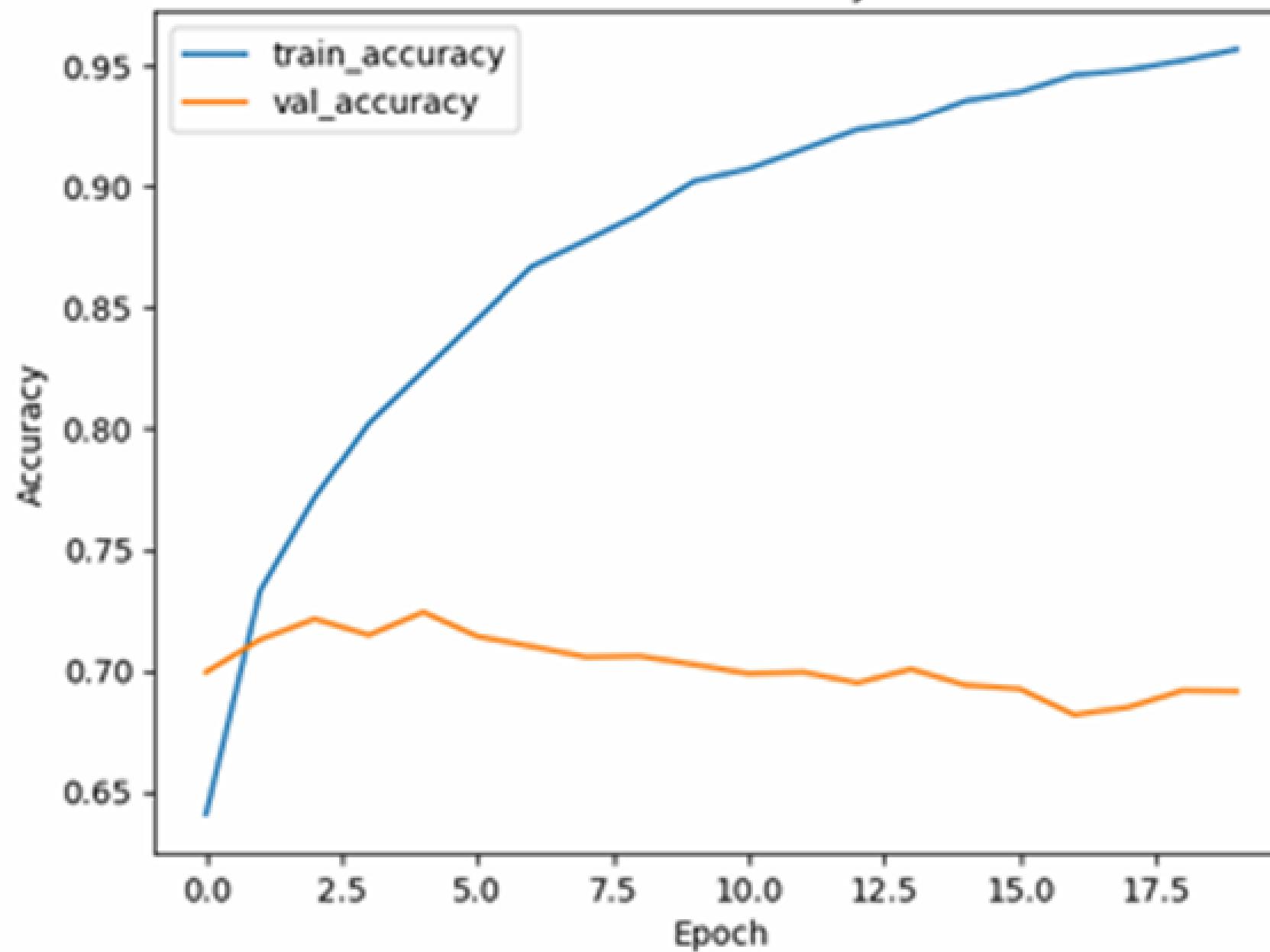
X_train_seq = tokenizer.texts_to_sequences(X_train)
X_test_seq = tokenizer.texts_to_sequences(X_test)

# Padding sequences to ensure uniform input size
X_train_padded = pad_sequences(X_train_seq, maxlen=max_len, padding='post', truncating='post')
X_test_padded = pad_sequences(X_test_seq, maxlen=max_len, padding='post', truncating='post')
```

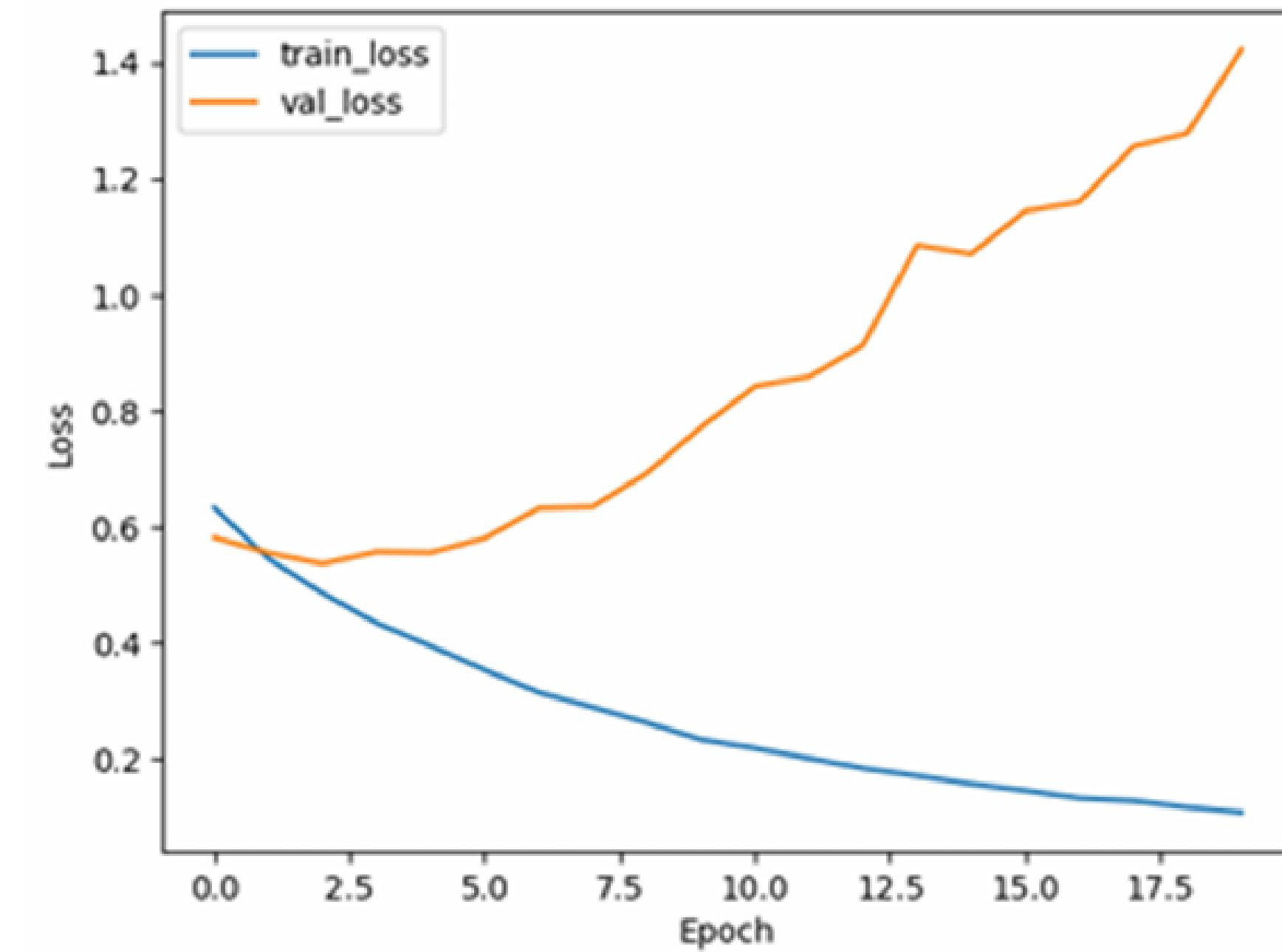
Model Training: A Bidirectional LSTM model was used for sentiment classification. The model leveraged GloVe pre-trained word embeddings for enhanced accuracy in understanding the semantic meaning of words.

Layer (type)	Output Shape	Param #
<hr/>		
embedding_1 (Embedding)	(None, 100, 100)	1000000
bidirectional (Bidirectional) 1)	(None, 256)	234496
dense_1 (Dense)	(None, 1)	257
<hr/>		
Total params: 1,234,753		
Trainable params: 1,234,753		
Non-trainable params: 0		

Model Accuracy



Model Loss





Sentiment Analysis of *IMDB Movie Reviews*

Project Description:

This project focuses on sentiment analysis of the IMDB movie reviews dataset, which is a binary classification task. The objective is to classify reviews as either positive or negative using deep learning models. The notebook leverages multiple recurrent neural network (RNN) architectures to build, train, and evaluate the performance of these models.

Key Steps and Features

Data Preprocessing:
Dataset: The dataset consists of IMDB movie reviews labeled as either positive or negative.

Text Cleaning and Tokenization:
The reviews are first cleaned to remove unnecessary characters, such as punctuation and special symbols.
Text is tokenized using a tokenizer, converting the words into sequences of integers.
The sequences are then padded to ensure a uniform length across the dataset. This step is crucial since RNNs require fixed-length inputs.

Feature Engineering:
The text data is transformed into a format suitable for neural network models through embeddings. This allows the model to capture semantic relationships between words.

Restricted Mode is intended for safe code browsing. Trust this window to enable all features. [Manage](#) [Learn More](#)

Untitled (10).ipynb X

X

D: > Downloads > Downloads > Untitled (10).ipynb > import nltk

...

+ Code + Markdown | ▶ Run All ⌘ Clear All Outputs | ⌘ Outline ...

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

def preprocess_text(sentence):

    sentence = re.sub(r"https?://\S+|www\S+", " ", sentence)

    sentence = re.sub(r"<.*?>|&([a-z0-9]+|[#][0-9]{1,6}|#[x][0-9a-f]{1,6});", " ", sentence)

    sentence = re.sub(r"^\w\s]", " ", sentence)

    sentence = re.sub(r"\w*\d\w*", " ", sentence)

    sentence = re.sub(r"[0-9]+", " ", sentence)

    sentence = re.sub(r"\s+", " ", sentence).strip()

    sentence = sentence.lower()

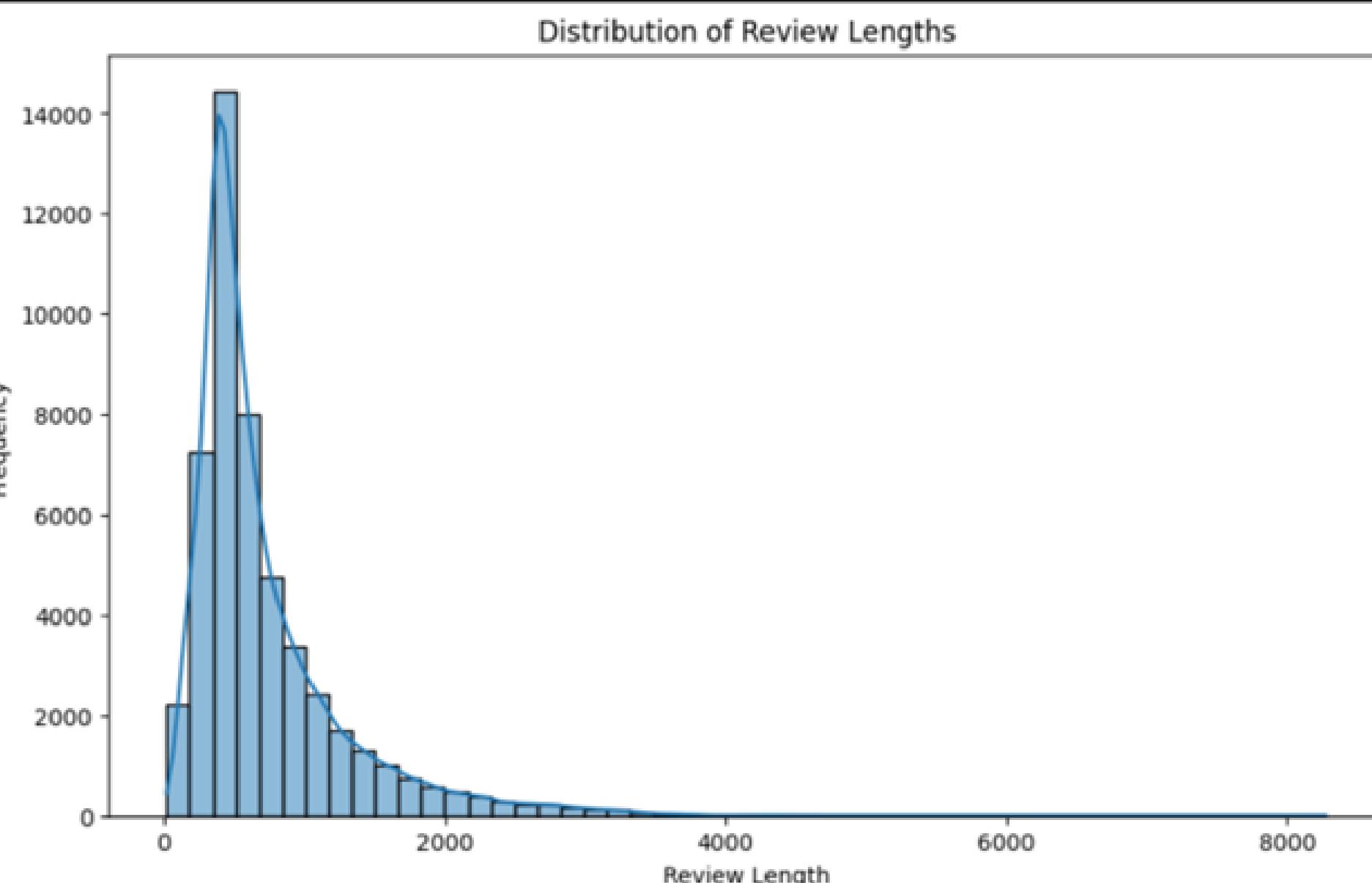
    tokens = []
    for token in sentence.split():
        if token not in stop_words:
            lemmatized_token = lemmatizer.lemmatize(token)
            stemmed_token = stemmer.stem(lemmatized_token)
            tokens.append(stemmed_token)

    return " ".join(tokens)
```

[12]

Restricted Mode ⌘ 7 ▲ 56 ⌘ 0

Select Kernel

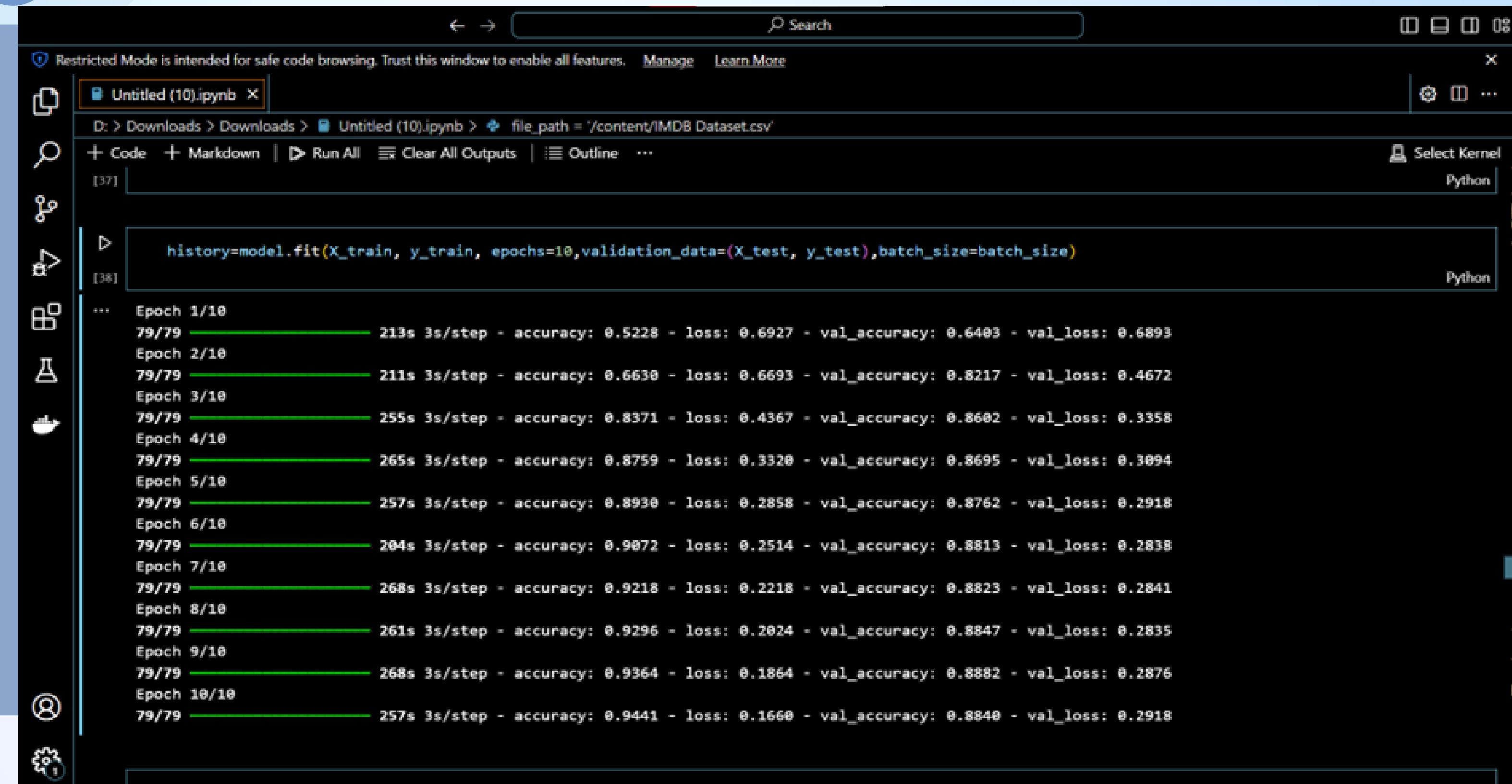


LSTM Model:

The LSTM model is designed to capture long-term dependencies in the text, which is essential for understanding the context in movie reviews.

Model Architectures: The notebook explores three types of RNN-based models to classify the reviews:

The model architecture includes an embedding layer, an LSTM layer, dropout layers to prevent overfitting, and a dense layer with a sigmoid activation function for binary classification.



A screenshot of a Jupyter Notebook interface. The top bar shows the file "Untitled (10).ipynb" is selected. The toolbar includes icons for Restricted Mode, Search, and various notebook operations. The menu bar has "File", "Edit", "Cell", "Kernel", "Help", and "About". The main area shows a code cell with the command:

```
history=model.fit(X_train, y_train, epochs=10,validation_data=(X_test, y_test),batch_size=batch_size)
```

Below the code cell, the output displays the training progress for 10 epochs:

```
... Epoch 1/10  
79/79 213s 3s/step - accuracy: 0.5228 - loss: 0.6927 - val_accuracy: 0.6489 - val_loss: 0.6893  
Epoch 2/10  
79/79 211s 3s/step - accuracy: 0.6638 - loss: 0.6693 - val_accuracy: 0.8217 - val_loss: 0.4672  
Epoch 3/10  
79/79 255s 3s/step - accuracy: 0.8371 - loss: 0.4367 - val_accuracy: 0.8682 - val_loss: 0.3358  
Epoch 4/10  
79/79 265s 3s/step - accuracy: 0.8759 - loss: 0.3320 - val_accuracy: 0.8695 - val_loss: 0.3094  
Epoch 5/10  
79/79 257s 3s/step - accuracy: 0.8938 - loss: 0.2858 - val_accuracy: 0.8762 - val_loss: 0.2918  
Epoch 6/10  
79/79 204s 3s/step - accuracy: 0.9072 - loss: 0.2514 - val_accuracy: 0.8813 - val_loss: 0.2838  
Epoch 7/10  
79/79 268s 3s/step - accuracy: 0.9218 - loss: 0.2218 - val_accuracy: 0.8823 - val_loss: 0.2841  
Epoch 8/10  
79/79 261s 3s/step - accuracy: 0.9296 - loss: 0.2024 - val_accuracy: 0.8847 - val_loss: 0.2835  
Epoch 9/10  
79/79 268s 3s/step - accuracy: 0.9364 - loss: 0.1864 - val_accuracy: 0.8882 - val_loss: 0.2876  
Epoch 10/10  
79/79 257s 3s/step - accuracy: 0.9441 - loss: 0.1660 - val_accuracy: 0.8848 - val_loss: 0.2918
```

Bidirectional LSTM Model:

This model enhances the LSTM by reading the sequence in both directions (forward and backward), improving the ability to understand context.

Similar to the LSTM model but with additional power due to bidirectional processing.

The screenshot shows a Jupyter Notebook interface with a dark theme. The top bar includes File, Edit, Selection, View, Go, Run, and other navigation options. A search bar is present, along with window control buttons. A message about Restricted Mode is displayed. The left sidebar features a file tree showing 'Untitled (10).ipynb' and various notebook icons. The main area contains a code cell with the following Python code:

```
model_2 = Sequential()
model_2.add(Embedding(input_dim=vocabulary_size, output_dim=128, input_length=X_train.shape[1]))
model_2.add(Bidirectional(LSTM(units=128, dropout=0.2, recurrent_dropout=0.2)))
model_2.add(Dropout(0.5))
model_2.add(Dense(128, activation='relu'))
model_2.add(Dropout(0.5))
model_2.add(Dense(1, activation='sigmoid'))

model_2.compile(loss='binary_crossentropy', optimizer=Adam(learning_rate=0.0001), metrics=['accuracy'])

history_2 = model_2.fit(X_train,y_train, validation_data=(X_test, y_test), epochs=10, batch_size=64)
```

Below the code cell, the output of the execution shows training progress for 10 epochs:

```
[41] Python
...
Epoch 1/10
/usr/local/lib/python3.10/dist-packages/keras/src/layers/core/embedding.py:90: UserWarning: Argument `input_length` is deprecated. Just remove it.
  warnings.warn(
625/625 646s 1s/step - accuracy: 0.5755 - loss: 0.6521 - val_accuracy: 0.8531 - val_loss: 0.3420
Epoch 2/10
625/625 676s 1s/step - accuracy: 0.8776 - loss: 0.3118 - val_accuracy: 0.8779 - val_loss: 0.2936
Epoch 3/10
625/625 682s 1s/step - accuracy: 0.9134 - loss: 0.2395 - val_accuracy: 0.8800 - val_loss: 0.2902
Epoch 4/10
625/625 629s 1s/step - accuracy: 0.9308 - loss: 0.1989 - val_accuracy: 0.8780 - val_loss: 0.2983
Epoch 5/10
625/625 671s 1s/step - accuracy: 0.9453 - loss: 0.1637 - val_accuracy: 0.8794 - val_loss: 0.3094
```

The bottom of the screen shows the Windows taskbar with various application icons and the system tray with battery, signal, and clock information.

GRU Model

The Gated Recurrent Unit (GRU) model is another variant of RNN, offering a simpler structure compared to LSTMs, while still capturing temporal dependencies in the data. GRU models are often faster to train and require fewer computational resources.

Training and Evaluation:

Training: Each model is trained on the preprocessed IMDB dataset using the Adam optimizer with binary cross-entropy loss. The models are trained for 10 epochs.

Evaluation Metrics: The performance is tracked using accuracy and loss. The training and validation accuracy are plotted to assess model performance.

Batch Size: The models are trained with a batch size of 64, balancing between computational efficiency and model convergence.

GRU Model

The image shows a Jupyter Notebook interface with the title "GRU Model". The notebook has a single open cell titled "Untitled (10).ipynb". The cell contains Python code for creating a GRU model and training it on an IMDB dataset. The code imports TensorFlow Keras layers, defines a sequential model with an embedding layer, two GRU layers, and two dense layers, then compiles it with binary crossentropy loss, Adam optimizer, and accuracy metric. It fits the model on training data and validates it on test data over 10 epochs with a batch size of 64. The output of the cell shows the training progress for each epoch, including the number of steps, duration, accuracy, loss, validation accuracy, and validation loss.

```
from tensorflow.keras.layers import GRU
model_3 = Sequential()
model_3.add(Embedding(input_dim=vocabulary_size, output_dim=128, input_length=X_train.shape[1]))
model_3.add(GRU(units=128, dropout=0.2, recurrent_dropout=0.2))
model_3.add(Dropout(0.5))
model_3.add(Dense(128, activation='relu'))
model_3.add(Dropout(0.5))
model_3.add(Dense(1, activation='sigmoid'))

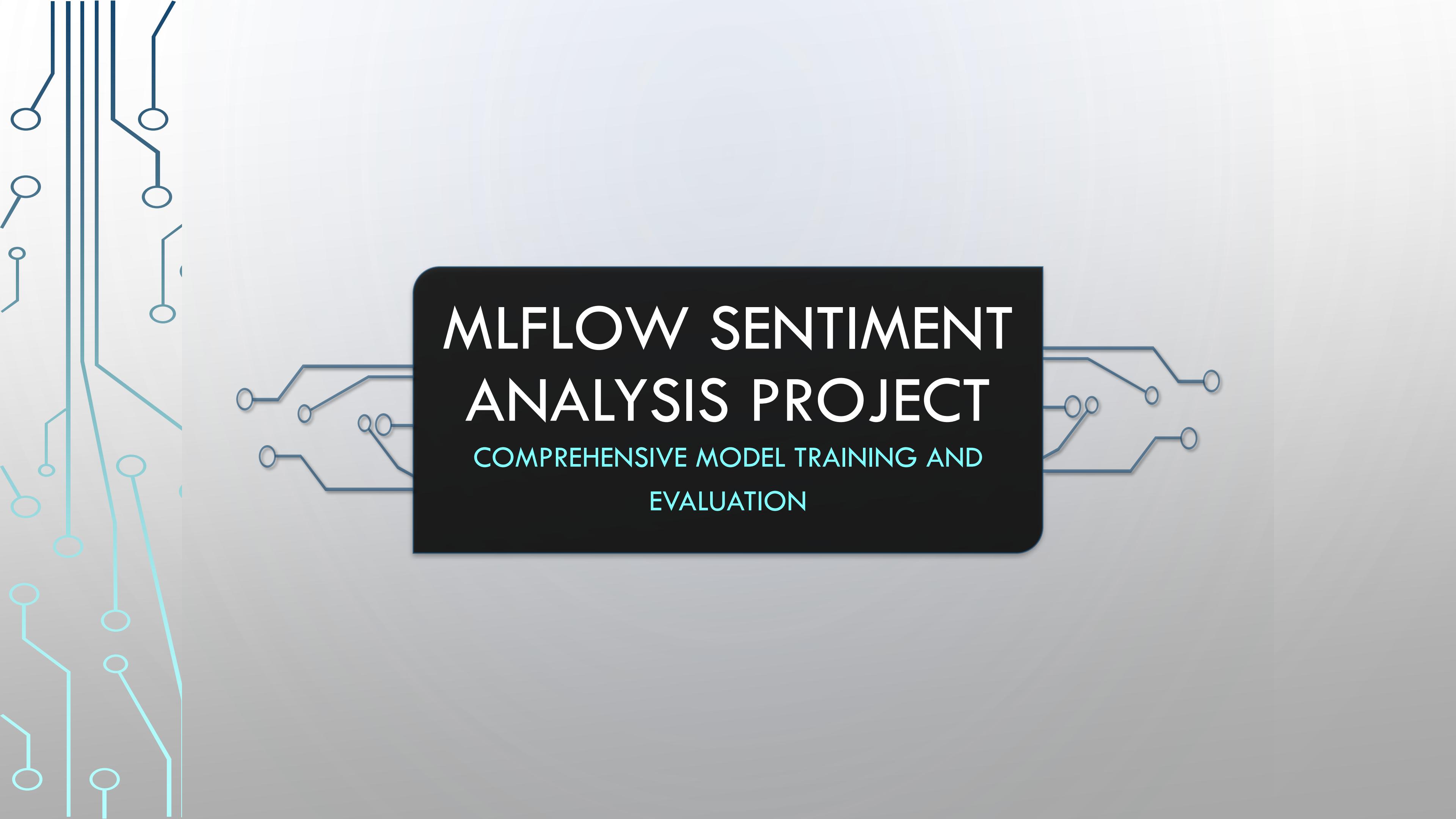
model_3.compile(loss='binary_crossentropy', optimizer=Adam(learning_rate=0.0001), metrics=['accuracy'])

history_3 = model_3.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=10, batch_size=64)
```

[43]

```
... Epoch 1/10  
625/625 455s 721ms/step - accuracy: 0.5699 - loss: 2.6377 - val_accuracy: 0.7545 - val_loss: 0.5430  
Epoch 2/10  
625/625 434s 693ms/step - accuracy: 0.8436 - loss: 0.3843 - val_accuracy: 0.7694 - val_loss: 0.4892  
Epoch 3/10  
625/625 437s 686ms/step - accuracy: 0.8811 - loss: 0.3087 - val_accuracy: 0.7710 - val_loss: 0.4775  
Epoch 4/10  
625/625 444s 691ms/step - accuracy: 0.8943 - loss: 0.2755 - val_accuracy: 0.7728 - val_loss: 0.4724  
Epoch 5/10  
625/625 416s 649ms/step - accuracy: 0.9182 - loss: 0.2495 - val_accuracy: 0.7728 - val_loss: 0.4713  
Epoch 6/10  
625/625 408s 652ms/step - accuracy: 0.9225 - loss: 0.2242 - val_accuracy: 0.7692 - val_loss: 0.4736  
Epoch 7/10
```

x Restricted Mode ⑧ 7 ▲ 56 ⌂ 0 ⌂ Prettier Q

A faint watermark of the MLflow logo is visible across the background of the slide. It consists of a dark blue rounded rectangle containing the word "MLFLOW" in white, with a smaller "mlflow" logo above it. The background of the slide is a light gray.

MLFLOW SENTIMENT ANALYSIS PROJECT

COMPREHENSIVE MODEL TRAINING AND
EVALUATION

DATASETS USED

- We utilized well-known datasets for sentiment analysis:
- Sentiment140 Dataset: 1.6M tweets labeled with sentiment.
- Amazon Product Reviews: Amazon product reviews with sentiment labels.
- IMDB Movie Reviews: 50k movie reviews labeled as positive or negative.
- Arabic Sentiment Analysis: Arabic tweets with sentiment categories.

The screenshot shows the mlflow interface with the title "mlflow 2.11.3". The "Experiments" tab is selected, indicated by a blue underline. The page has a dark theme. At the top right are a plus sign icon and a back arrow icon. Below the tabs is a search bar labeled "Search Experiments". The main area displays a list of experiments with the following details:

Experiment Name	Last Run Status	Edit Icon	Delete Icon
Sentiment Analysis - Bidirectional LSTM	Success		
IMDB Models	Not Started		
Sentiment Analysis Arabic tweets	Not Started		
Amazon-Reviews	Not Started		

MODELS TRAINED

- Multiple models were trained and evaluated:
 - Logistic Regression
 - Support Vector Machine (SVM)
 - Random Forest
 - Bidirectional LSTM
 - Pre-trained Transformer (BERT)

Run Name	Created	Dataset
GRU_sequential	8 hours ago	-
GRU_sequential	8 hours ago	-
Bidirectional_sequential	9 hours ago	-
Bidirectional_sequential	11 hours ago	-
Sequential_LSTM_Model	11 hours ago	-

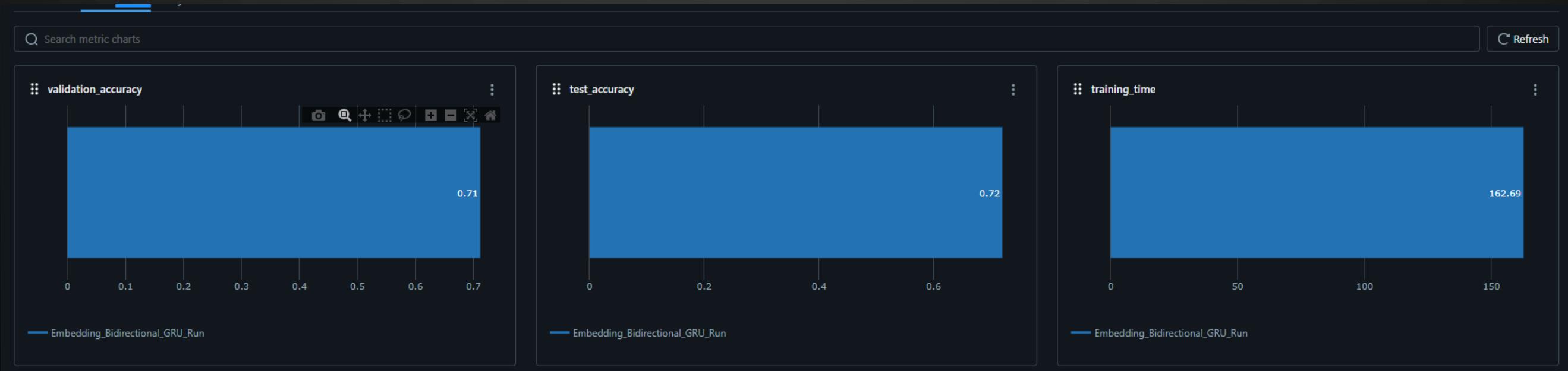
Run Name	Created	Dataset	Duration	Source	Models
Embedding_Bidirectional...	3 hours ago	-	6.1min	CUDNN...	keras
Embedding_Layer_LSTM...	4 hours ago	-	6.5min	CUDNN...	keras
Embedding_Layer_MLP_R...	4 hours ago	-	6.2min	CUDNN...	keras
Embedding_Layer_MLP_R...	4 hours ago	-	28.4min	CUDNN...	keras
Text_MLP_Run	5 hours ago	-	28.3min	CUDNN...	keras
Overfitter_MLP_Run	5 hours ago	-	30.3min	CUDNN...	keras
Text_MLP_Run	5 hours ago	-	30.3min	CUDNN...	keras
Overfitter_MLP_Run	6 hours ago	-	150min	CUDNN...	keras
Logistic_Regression_CHE...	1 day ago	-	5.5s	colab...	sklearn
Bi-directional_LSTM_Run	1 day ago	-	1.0s	CUDNN...	tensorflow

Run Name	Created	Dataset
CNN Model	10 hours ago	-
glamorous-hawk-503	10 hours ago	-

Run Name	Created	Dataset
LSTM with Tokenization	7 hours ago	-
Random Forest with Cou...	8 hours ago	-
Logistic Regression with T...	10 hours ago	-

MLFLOW FOR MODEL TRACKING

- MLFlow was used for tracking the performance of different models:
 - Track metrics such as accuracy, precision, recall, and F1-score.
 - Compare model versions and hyperparameter tuning results.

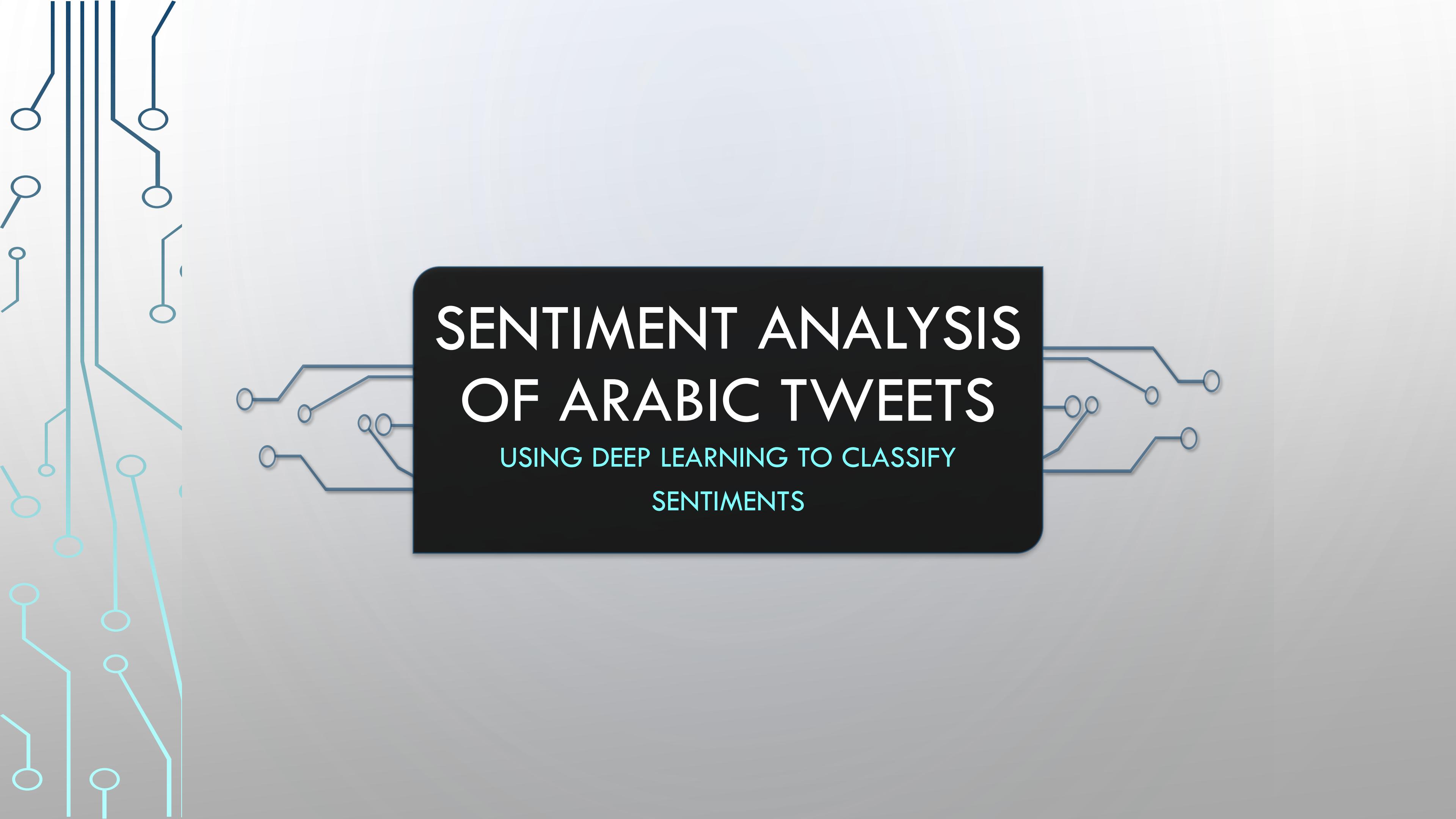


LINK TO MLFLOW EXPERIMENT

- Access the full experiment details here:
- MLFlow Link:
https://dagshub.com/asamy4194/DEPI_Final_Project.mlflow/#/experiments/0?viewStateShareKey=d537fdc0c959fa49a2e0b67727b54010fcd4ffd194d195a2598ac3fdf5c6e596

MLflow Integration

- **MLflow for Collaboration:** Linking team notebooks and tracking experiments to ensure all models and results are centrally managed.
- **Experiment Tracking:** Tracking model parameters, accuracy, and results for different datasets.



SENTIMENT ANALYSIS OF ARABIC TWEETS

USING DEEP LEARNING TO CLASSIFY
SENTIMENTS

OVERVIEW OF THE PROJECT



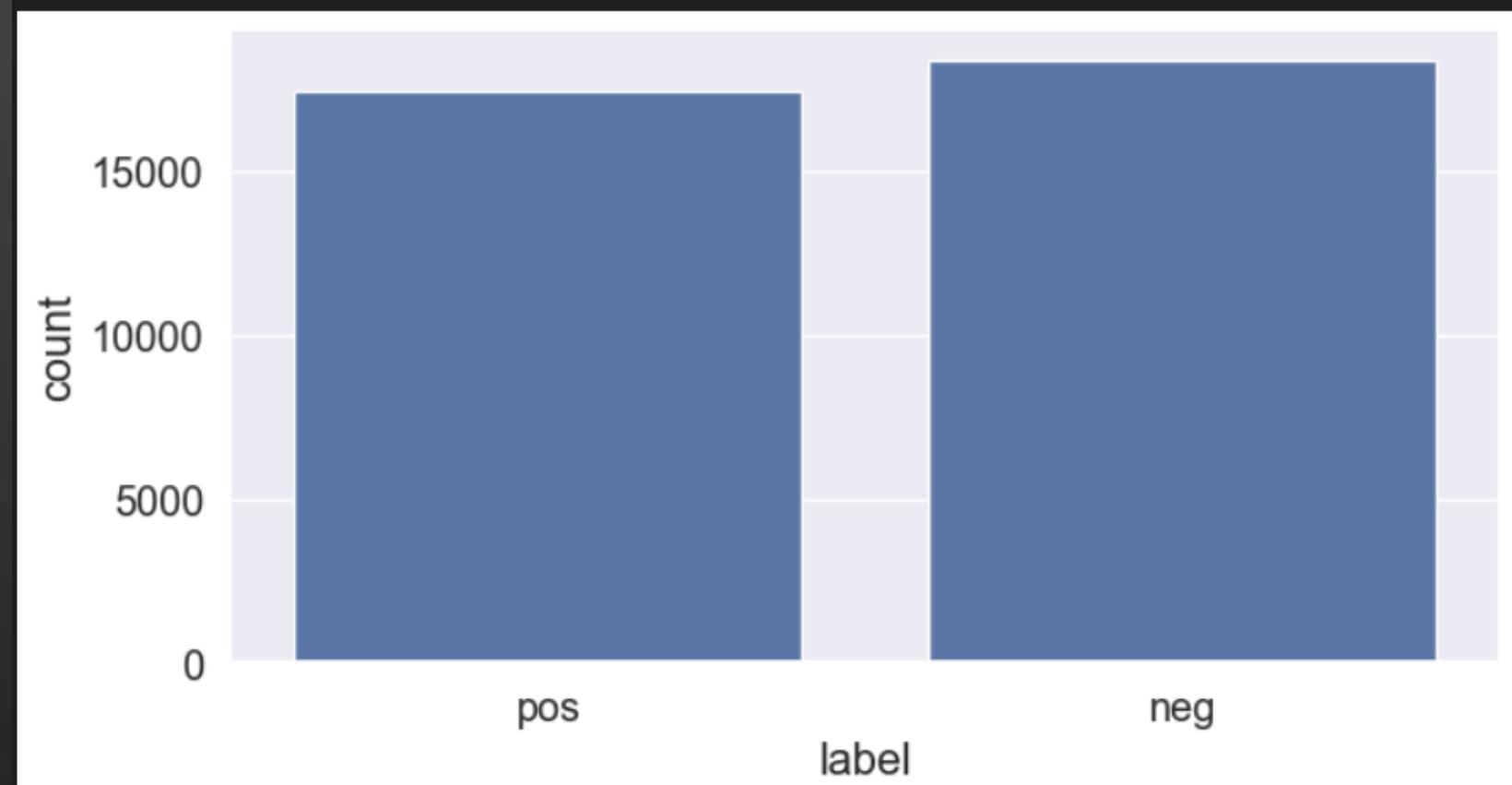
Objective: Build a model to classify Arabic tweets into positive and negative sentiments.



Scope: Utilize deep learning and natural language processing (NLP) for sentiment analysis.

DATA COLLECTIO N

- Datasets: Merged two datasets of positive and negative Arabic tweets.
- Data Cleaning: Removed duplicates, handled missing values.



PREPROCESSING STEPS



Stop Words Removal: Custom and existing stop words were used.



Emoji/Text Handling: Converted emojis, cleaned URLs, mentions, and punctuation.



Text Normalization: ISRI Stemmer and Arabic Light Stemmer were used.

FEATURE ENGINEERIN G

- Tokenization & Padding: Keras tokenizer applied to standardize input length.
- Vocabulary Size Limitation: Most frequent words were selected to reduce complexity.

```
print(X[:5])  
  
print('#' * 50)  
  
print(y[:5])
```

```
... مناسب وز هلال حب سحب على يفو ...  
0 ... اتحاد نصر اتحسي سى يالطواق انب وقف مع ماج حاجه  
1 ... مفروض كو حساب موبایل ويتر يافيس لي عند يس مينز  
2 ...  
3 ...  
4 ...  
Name: text, dtype: object  
#####  
0 pos  
1 pos  
2 neg  
3 neg  
4 neg  
Name: label, dtype: object
```

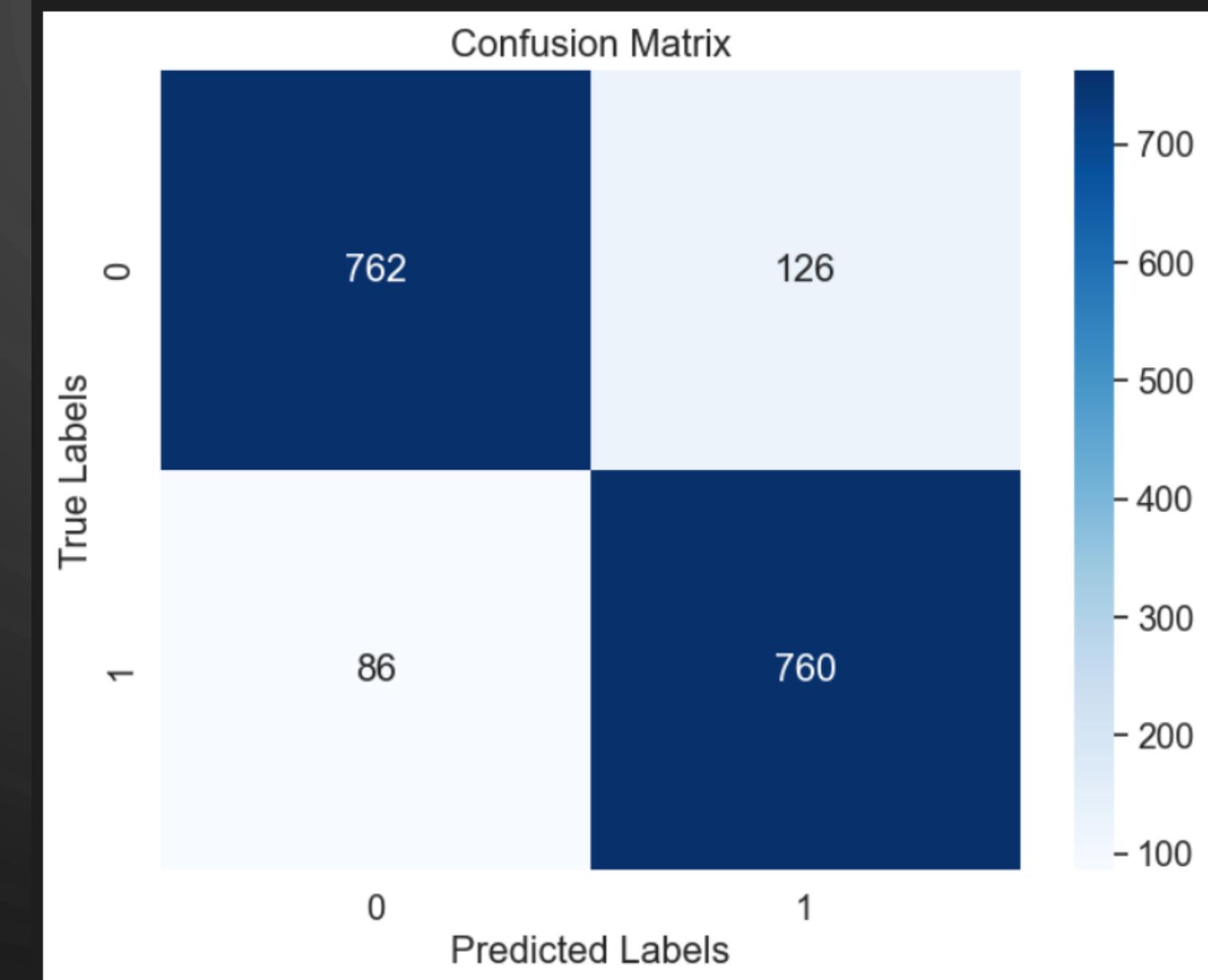
MODEL BUILDING

- Model: Bidirectional LSTM to capture sequential dependencies.
- Training & Validation: Trained on 30,000 tweets, validated with a test set.
- Overfitting Prevention: Applied dropout and batch normalization.

Model: "sequential"		
Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, None, 64)	1754048
bidirectional (Bidirectiona l1)	(None, None, 128)	66048
dropout (Dropout)	(None, None, 128)	0
batch_normalization (BatchN ormalization)	(None, None, 128)	512
bidirectional_1 (Bidirectio nal)	(None, 64)	41216
dropout_1 (Dropout)	(None, 64)	0
batch_normalization_1 (BathchNormalization)	(None, 64)	256
dense (Dense)	(None, 32)	2080
dropout_2 (Dropout)	(None, 32)	0
...		
Total params: 1,864,226		
Trainable params: 1,863,842		
Non-trainable params: 384		

RESULTS AND FUTURE WORK

- Accuracy: 96.2% on training, 87.7% on validation.
- Confusion Matrix & Accuracy Graph: Used to visualize performance.
- Future Work: Expand dataset, fine-tune hyperparameters, and incorporate more advanced techniques.





Sentiment Analysis of *Amazon reviews*

Taking sample from train data and test data:

In [3]:

```
train.columns = ["polarity", "title", "text"]
train = train.sample(100000, random_state = 99)
train.head()
```

Out[3]:

	polarity	title	text
2136295	2	Making a suit of human skin is hard work. Real...	Yes, go ahead and laugh. But some of the other...
2620775	2	If Only Amazon allowed my 6 "s.	Wow, This movie is a classic. Family movie by a...
1588872	1	I THOUGHT APPLE WAS USER FRIENDLY	I THOUGHT APPLE WAS USER FRIENDLY????????? I GU...
3125694	2	Very Helpful	I bought this book after joining Curves. It he...
1668956	1	Academically Repugnant	The first ten minutes of the film passed with ...

In [4]:

```
test.columns = ["polarity", "title", "text"]
test = test.sample(50000, random_state = 99)
test.head()
```

Out[4]:

	polarity	title	text
110227	1	A Useless Tool for Novel Writing	I met the author of this book/software at a wr...
118016	1	A load of historical BS	I bought this book several years ago and it wa...
180926	1	Don't do it	I have had two of these units and both worked ...
121412	1	Cute, but...	Really liked the Mobiblue 1500 cube at first -...
128565	1	Totally disappointed	This trilogy was recommended to me and I was s...

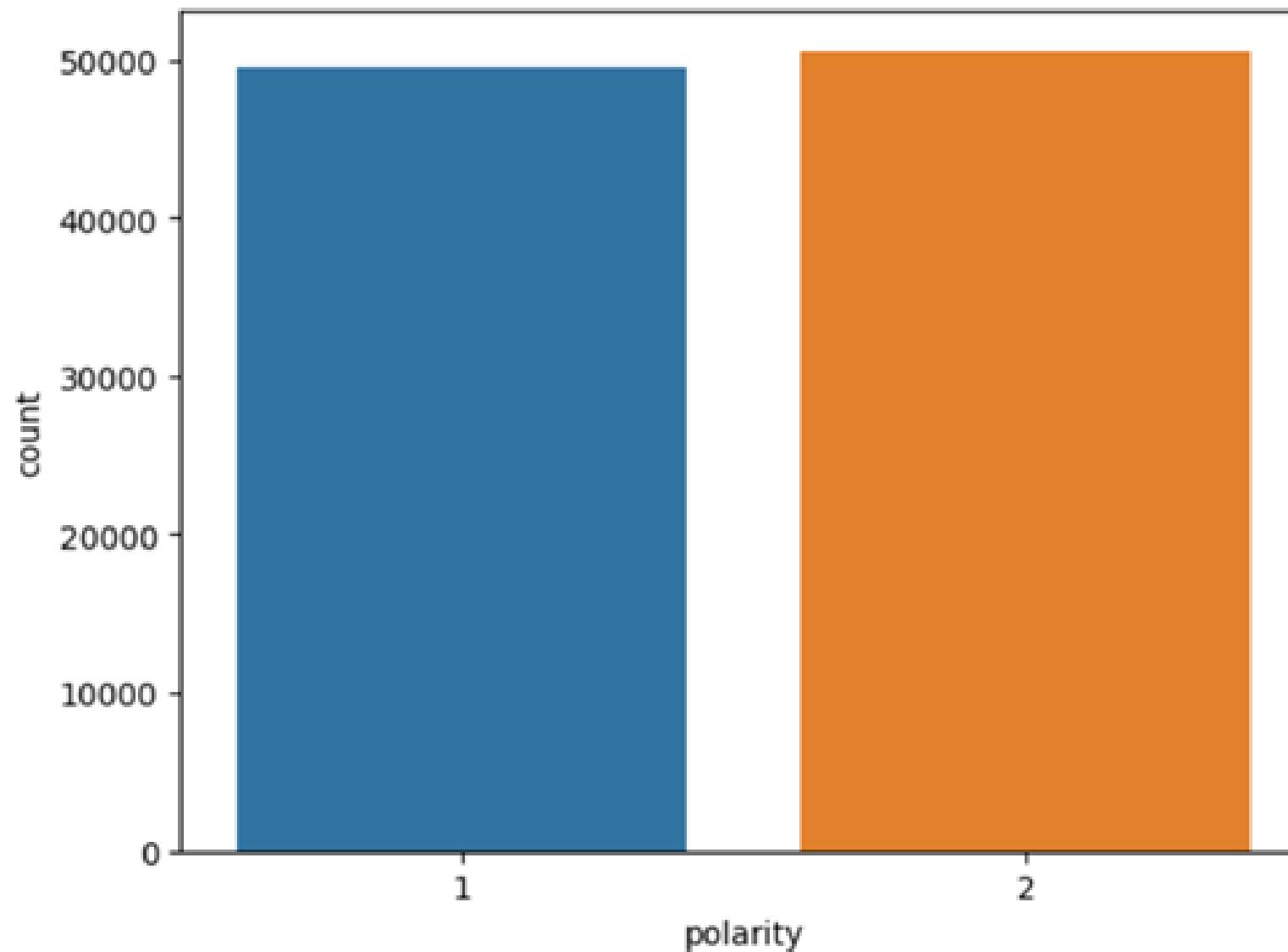
Count plot between negative (2) and positive (1) reviews in train and test data:

In [9]:

```
sns.countplot(x ='polarity', data = train)
```

Out[9]:

```
<Axes: xlabel='polarity', ylabel='count'>
```

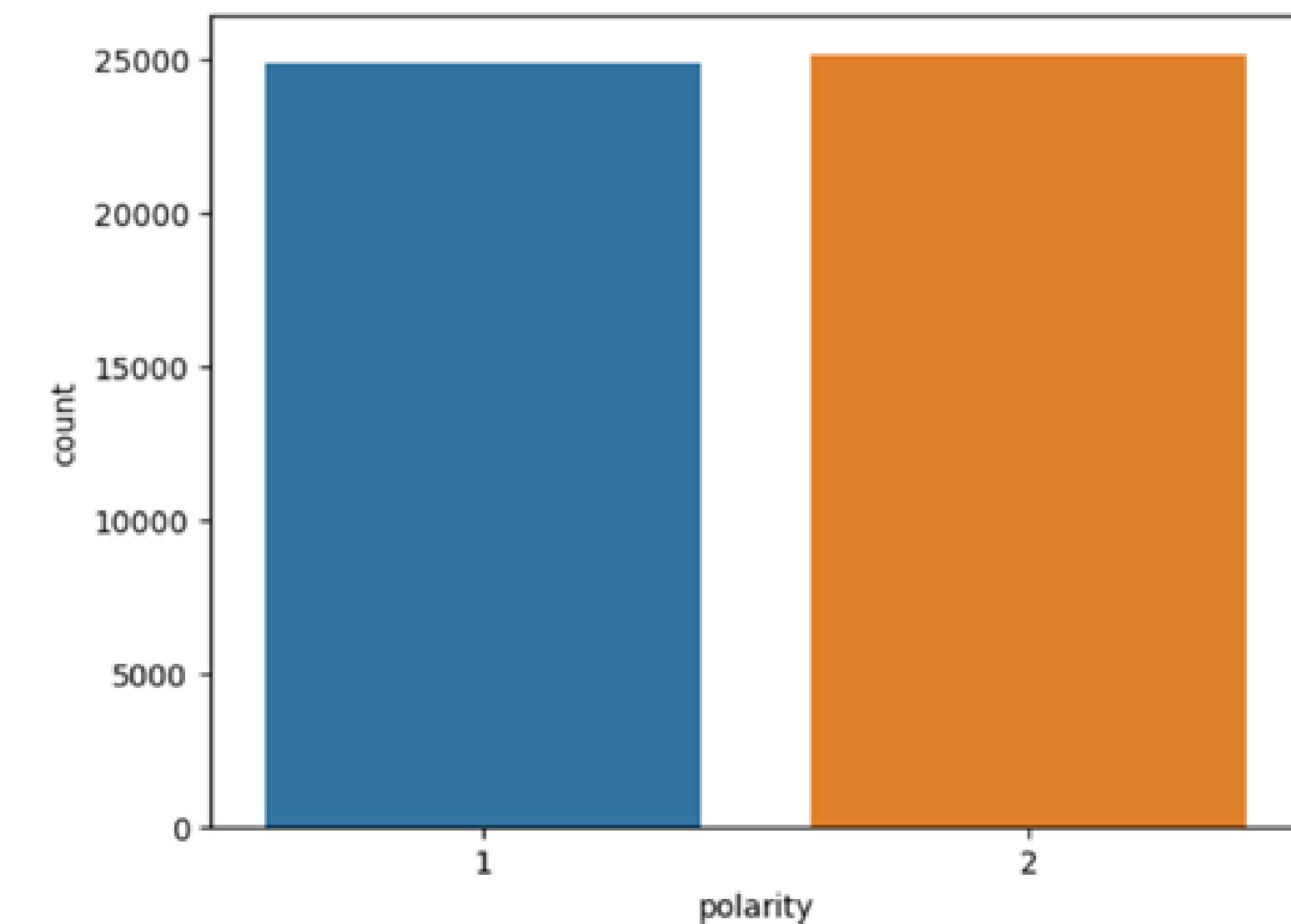


In [10]:

```
sns.countplot(x ='polarity', data = test)
```

Out[10]:

```
<Axes: xlabel='polarity', ylabel='count'>
```



Text preprocessing

```
pr=PorterStemmer()
def clean(text):
    text=text.lower()
    without_punctuation = ''.join(i for i in text if i not in string.punctuation)
    Without_punctuation= ''.join(c for c in without_punctuation if not ud.category(c).startswith('P'))
    nums = re.compile(r'\d+')
    clean_text = nums.sub(r'', Without_punctuation)
    clean_text = ' '.join([word for word in clean_text.split() if word not in stop_words])
    emoji_pattern = re.compile("["
                                u"\U0001F600-\U0001F64F"
                                u"\U0001F300-\U0001F5FF"
                                u"\U0001F680-\U0001F6FF"
                                u"\U0001F1E0-\U0001F1FF"
                                u"\U00002500-\U00002BEF"
                                u"\U00002702-\U000027B0"
                                u"\U00002702-\U000027B0"
                                u"\U000024C2-\U0001F251"
                                u"\U0001f926-\U0001f937"
                                u"\U00010000-\U0010ffff"
                                u"\u2640-\u2642"
                                u"\u2600-\u2B55"
                                u"\u200d"
                                u"\u23cf"
                                u"\u23e9"
                                u"\u231a"
                                u"\ufe0f"
                                u"\u3030"
                                "]+", flags=re.UNICODE)

    clean = emoji_pattern.sub(r'', clean_text)
    words = word_tokenize(clean)
    stemmed_words=[pr.stem(word) for word in words]
    return TreebankWordDetokenizer().detokenize(stemmed_words)
```

Splitting data and applying vectorization using TFIDF:

```
In [18]:  
x_train = train["text"]  
y_train = train["polarity"]  
x_test = test["text"]  
y_test = test["polarity"]
```

```
In [19]:  
vectorizer = TfidfVectorizer().fit(x_train)  
x_train_tfidf = vectorizer.transform(x_train)  
x_test_tfidf = vectorizer.transform(x_test)
```

Machine learning section: 1- logistic model with TFIDF

In [20]:

```
model = LogisticRegression()
model.fit(x_train_tfidf,y_train)
y_train_predict = model.predict(x_train_tfidf)
y_test_predict = model.predict(x_test_tfidf)
```

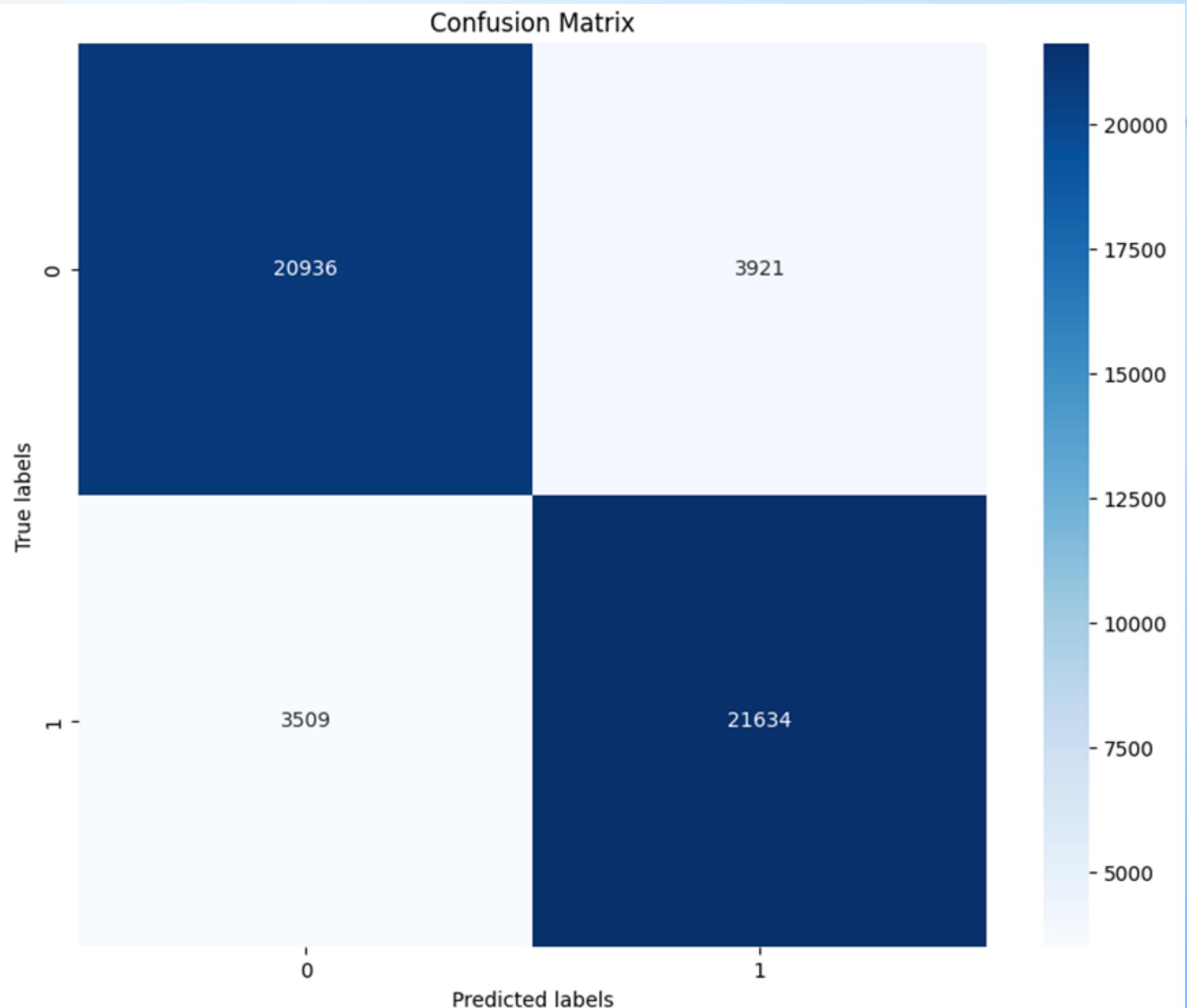
In [21]:

```
train_mse = mean_squared_error(y_train, y_train_predict)
train_r2 = r2_score(y_train, y_train_predict)
train_accuracy=accuracy_score(y_train, y_train_predict)

test_mse = mean_squared_error(y_test, y_test_predict)
test_r2 = r2_score(y_test, y_test_predict)
test_accuracy=accuracy_score(y_test, y_test_predict)

print(f"Training MSE: {train_mse}")
print(f"Training R2 Score: {train_r2}")
print(f"training accuracy: {train_accuracy}")
print(f"Testing MSE: {test_mse}")
print(f"Testing R2 Score: {test_r2}")
print(f"testing accuracy: {test_accuracy}")
```

```
Training MSE: 0.11279
Training R2 Score: 0.5487854030337671
training accuracy: 0.88721
Testing MSE: 0.1486
Testing R2 Score: 0.4055805515467178
testing accuracy: 0.8514
```



Machine learning section: Logistic with N-Gram

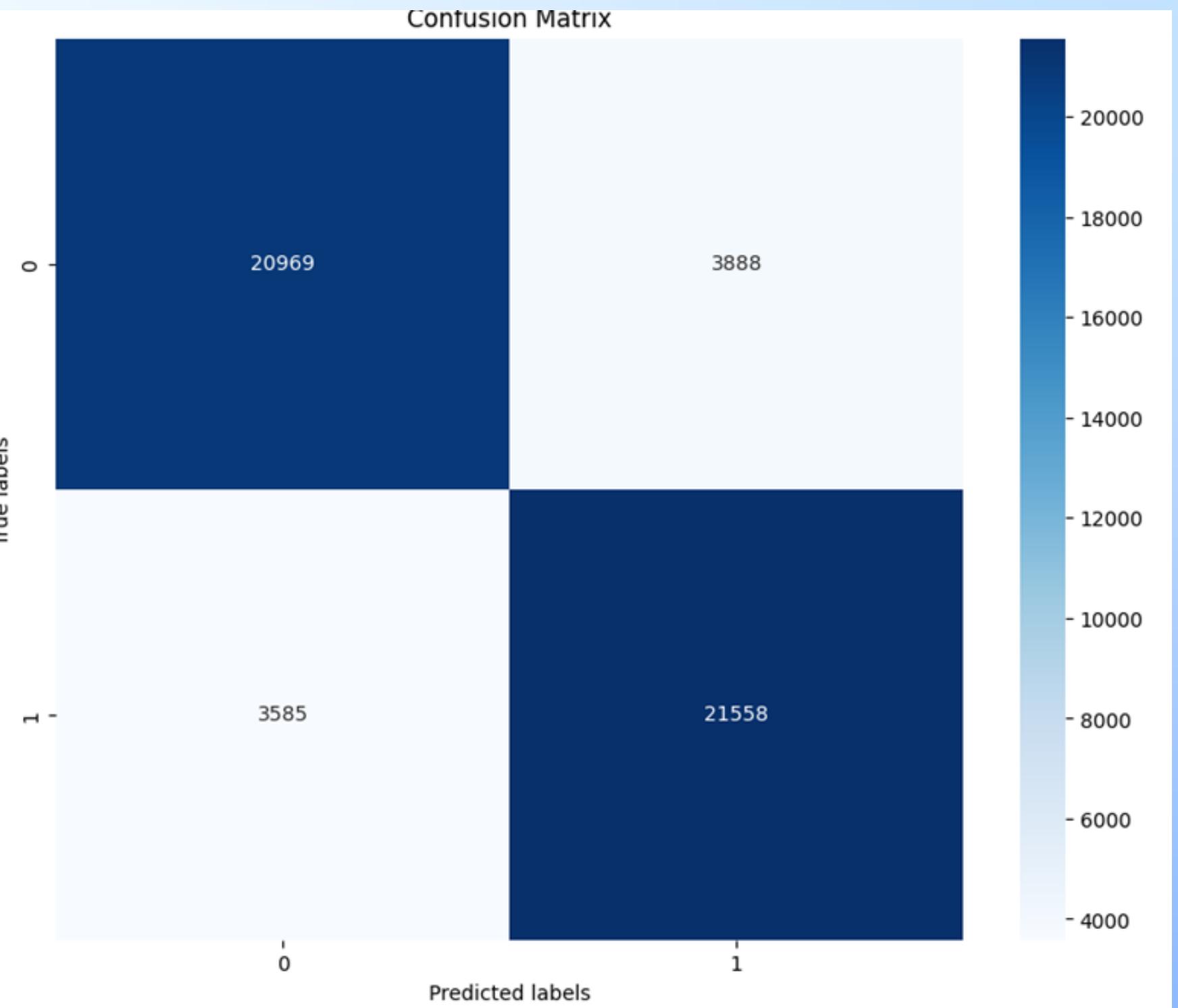
```
model = LogisticRegression()
model.fit(X_train_ngrams, y_train)
y_train_predictlo = model.predict(X_train_ngrams)
y_test_predictlo = model.predict(X_test_ngrams)
```

```
train_mse1 = mean_squared_error(y_train, y_train_predictlo)
train_r21 = r2_score(y_train, y_train_predictlo)
train_accuracy=accuracy_score(y_train, y_train_predictlo)

test_mse1 = mean_squared_error(y_test, y_test_predictlo)
test_r21 = r2_score(y_test, y_test_predictlo)
test_accuracy=accuracy_score(y_test, y_test_predictlo)

print(f"Training MSE: {train_mse1}")
print(f"Training R2 Score: {train_r21}")
print(f"training accuracy: {train_accuracy}")
print(f"Testing MSE: {test_mse1}")
print(f"Testing R2 Score: {test_r21}")
print(f"testing accuracy: {test_accuracy}")
```

```
Training MSE: 0.00033
Training R2 Score: 0.9986798402606716
training accuracy: 0.99967
Testing MSE: 0.1306
Testing R2 Score: 0.4775829073485959
testing accuracy: 0.8694
```



Deep learning model

In [33]:

```
max_words = 2000
max_len = 100

tokenizer = Tokenizer(num_words=max_words)
tokenizer.fit_on_texts(train['text'])

X_train = tokenizer.texts_to_sequences(train['text'])
X_train = pad_sequences(X_train, maxlen=max_len)
X_test = tokenizer.texts_to_sequences(test['text'])
X_test = pad_sequences(X_test, maxlen=max_len)

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(input_dim=max_words, output_dim=128, input_length=max_len),
    tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(64, return_sequences=True)),
    tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(64, return_sequences=False)),
    tf.keras.layers.Dropout(0.5),
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.BatchNormalization(),
    tf.keras.layers.Dense(1, activation='sigmoid')
])
```

In [34]:

```
Lr_schedule = tf.keras.callbacks.ReduceLROnPlateau(monitor='val_loss', factor=0.1, patience=3, verbose=1)
optimizer = tf.keras.optimizers.Adam(learning_rate=.0001)
model.compile(optimizer=optimizer, loss='binary_crossentropy', metrics=['accuracy'])
history = model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=10, batch_size=64, callbacks=[Lr_schedule])
```

Deep learning model

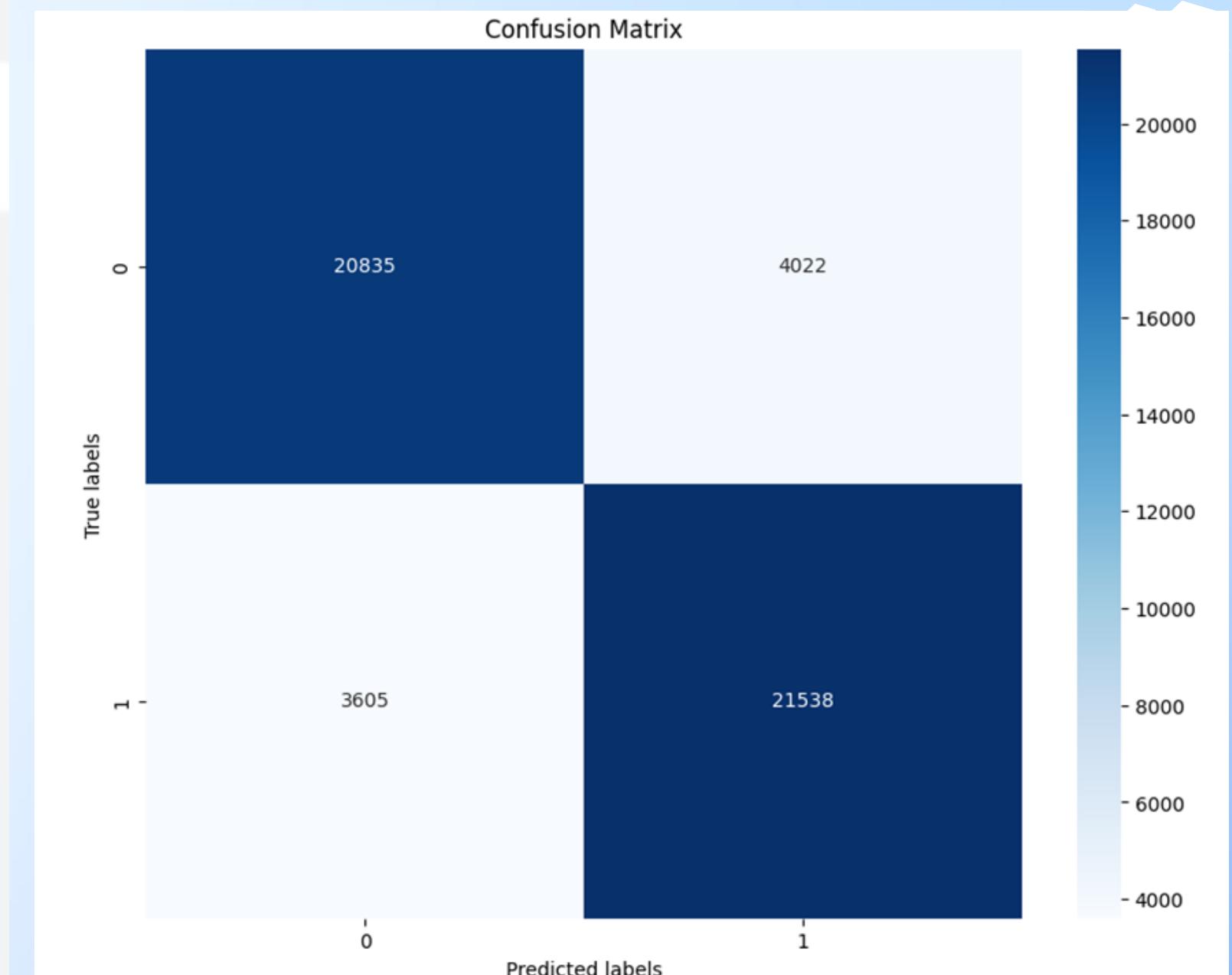
```
Epoch 2/10
1563/1563 49s 31ms/step - accuracy: 0.8394 - loss: 0.3657 - val_accuracy: 0.8377 - val_loss: 0.3616 - learning_rate: 1.0000e-04
Epoch 3/10
1563/1563 49s 31ms/step - accuracy: 0.8472 - loss: 0.3476 - val_accuracy: 0.8295 - val_loss: 0.3807 - learning_rate: 1.0000e-04
Epoch 4/10
1563/1563 49s 31ms/step - accuracy: 0.8497 - loss: 0.3393 - val_accuracy: 0.8376 - val_loss: 0.3698 - learning_rate: 1.0000e-04
Epoch 5/10
1561/1563 0s 26ms/step - accuracy: 0.8587 - loss: 0.3266
Epoch 5: ReduceLROnPlateau reducing learning rate to 9.999999747378752e-06.
1563/1563 49s 31ms/step - accuracy: 0.8587 - loss: 0.3266 - val_accuracy: 0.8400 - val_loss: 0.3616 - learning_rate: 1.0000e-04
Epoch 6/10
1563/1563 49s 32ms/step - accuracy: 0.8663 - loss: 0.3099 - val_accuracy: 0.8475 - val_loss: 0.3510 - learning_rate: 1.0000e-05
Epoch 7/10
1563/1563 49s 31ms/step - accuracy: 0.8651 - loss: 0.3129 - val_accuracy: 0.8479 - val_loss: 0.3516 - learning_rate: 1.0000e-05
Epoch 8/10
1563/1563 49s 31ms/step - accuracy: 0.8661 - loss: 0.3075 - val_accuracy: 0.8481 - val_loss: 0.3544 - learning_rate: 1.0000e-05
Epoch 9/10
1563/1563 0s 27ms/step - accuracy: 0.8675 - loss: 0.3061
Epoch 9: ReduceLROnPlateau reducing learning rate to 9.999999747378752e-07.
1563/1563 49s 31ms/step - accuracy: 0.8675 - loss: 0.3061 - val_accuracy: 0.8480 - val_loss: 0.3537 - learning_rate: 1.0000e-05
Epoch 10/10
1563/1563 49s 31ms/step - accuracy: 0.8683 - loss: 0.3038 - val_accuracy: 0.8475 - val_loss: 0.3538 - learning_rate: 1.0000e-06
```

Model summary

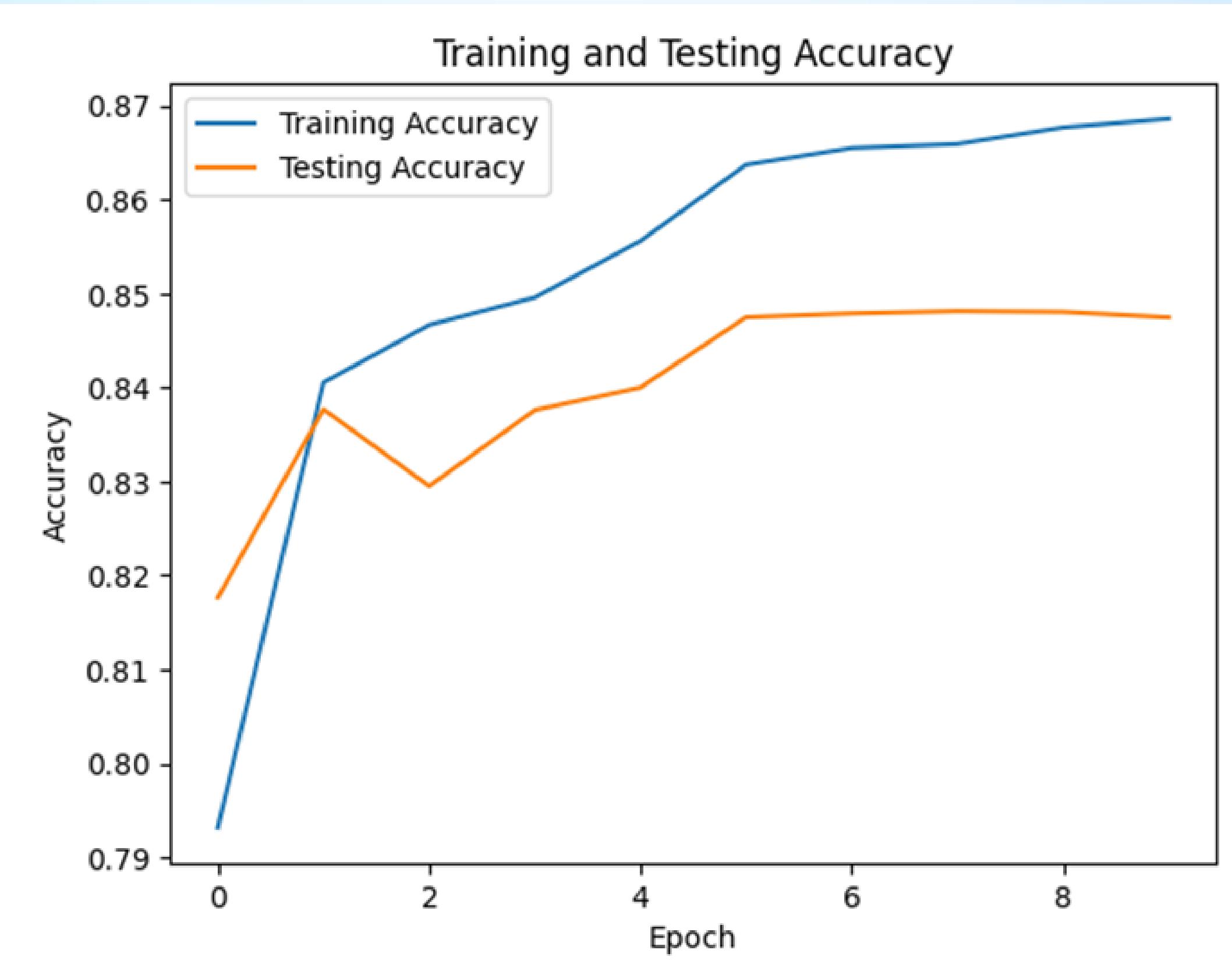
```
model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 100, 128)	256,000
bidirectional (Bidirectional)	(None, 100, 128)	98,816
bidirectional_1 (Bidirectional)	(None, 128)	98,816
dropout (Dropout)	(None, 128)	0
dense (Dense)	(None, 64)	8,256
dense_1 (Dense)	(None, 32)	2,080
dense_2 (Dense)	(None, 32)	1,056
batch_normalization (BatchNormalization)	(None, 32)	128
dense_3 (Dense)	(None, 1)	33



The compression between training and testing accuracy:



Thank You

