## Imports

```
 1   import pickle
 2   import numpy as np
 3   import pandas as pd
 4   import matplotlib.pyplot as plt
 5   import matplotlib.colors as mcolors
 6
 7   from tensorflow import keras
 8   from tensorflow.keras.models import Sequential
 9   from tensorflow.keras.preprocessing.text import Tokenizer
10   from tensorflow.keras.preprocessing.sequence import pad_sequences
11   from tensorflow.keras.optimizers import Adam
12   from tensorflow.keras.layers import (
13       Dense, Flatten, Dropout, Embedding,
14       Conv1D, GlobalMaxPooling1D, LSTM,
15       Bidirectional, BatchNormalization,
16       SimpleRNN
17   )
18
19
20   colors = mcolors.TABLEAU_COLORS.keys()
21   dir = '/content/drive/MyDrive/NLP/textClassification2/'
```

## Data Load

## Mounting google drive

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

## Loading and spliting data

```
1 # Entire Dataset
2
3 train_data = pd.read_csv(
4     '/content/drive/MyDrive/NLP/Corona_NLP_train.csv',
5     encoding='latin1'
6     )
7
8 test_data = pd.read_csv(
9     '/content/drive/MyDrive/NLP/Corona_NLP_test.csv',
10     encoding='latin1'
11     )
```

```
1 # Small Dataset with 80–20 split
2
3 # data = pd.read_csv(
4 #     '/content/drive/MyDrive/NLP/Corona_NLP_test.csv',
5 #      encoding='latin1'
6 #        )
7 # train_data = data.sample(frac=0.8, random_state=42)
8 # test_data = data.drop(train_data.index)
```

```
1 print(f"Training Data Shape: {train_data.shape}")
2 print(f"Test Data Shape: {test_data.shape}")
```

```
Training Data Shape: (41157, 6)
Test Data Shape: (3798, 6)
```

▾ Data Visualization and Analysis

# About Dataset

## Coronavirus Tweets

**Dataset:** https://www.kaggle.com/datasets/datatattle/covid-19-nlp-text-classification

This dataset is about **Coronavirus Tweets** and has **sentiments** for each tweet. The tweets have been pulled from Twitter and manual tagging has been done then.

It has the following columns:

1. Location
2. Tweet At
3. Original Tweet
4. Label (Sentiment: ['Positive', 'Negative', 'Neutral', 'Extremely Positive', 'Extremely Negative'])
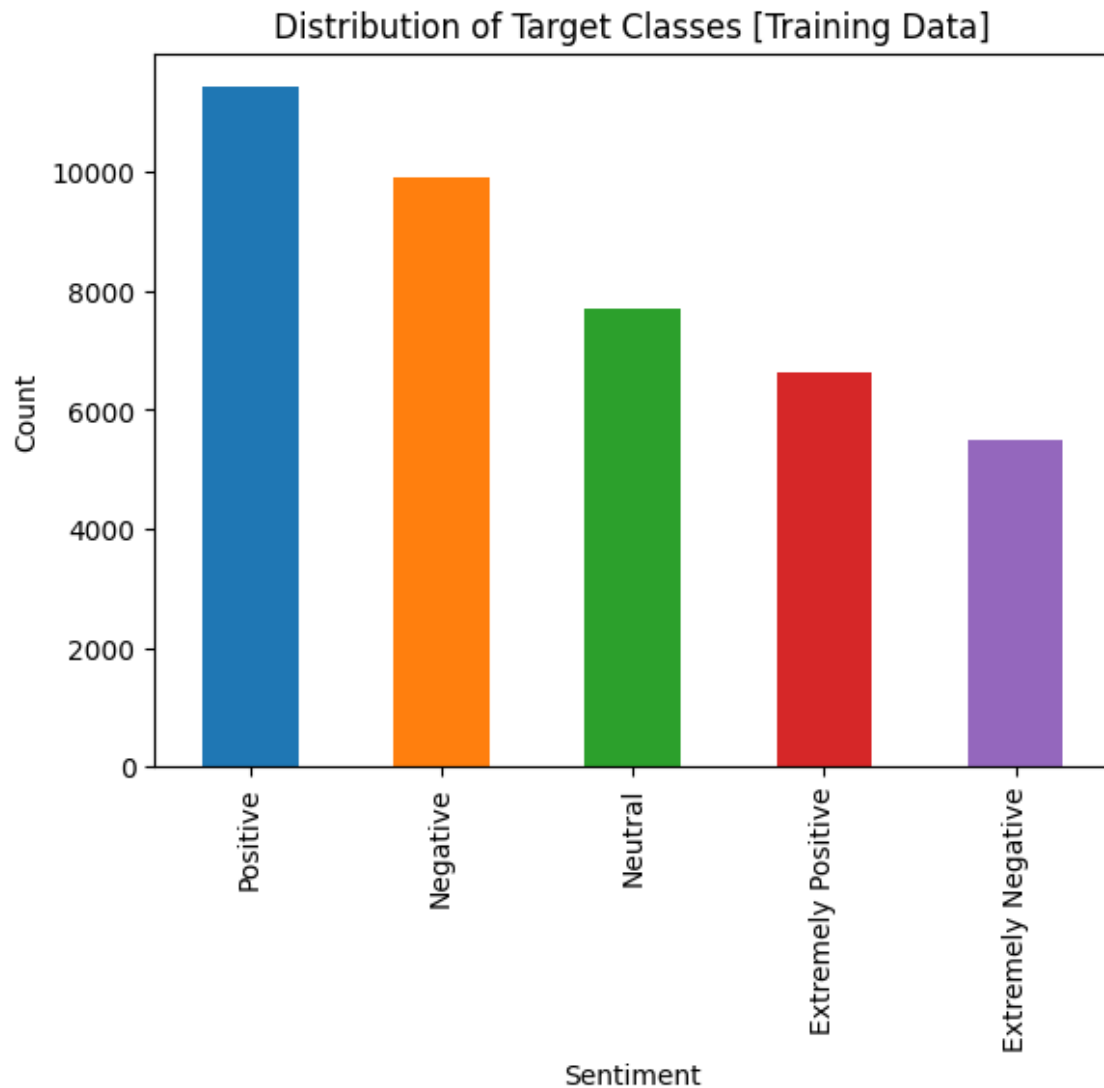
## Model Objective

**Model should be able to predict the sentiment of the tweet by examining it**.

## Distribution of Classes

Below graph shows class distribution

# Classes Distribution [Training Data]

```
1
2  fig = plt.figure()
3  train_data.value_counts('Sentiment').plot(kind='bar', color=colors
4  plt.ylabel('Count')
5  plt.xlabel('Sentiment')
6  plt.title('Distribution of Target Classes [Training Data]')
7  plt.show()
```



Distribution of Target Classes [Training Data]

▼ Classes Distribution [Test Data]

```
1 fig = plt.figure()
2 test_data.value_counts('Sentiment').plot(kind='bar', color=colors)
3 plt.ylabel('Count')
4 plt.xlabel('Sentiment')
5 plt.title('Distribution of Target Classes [Test Data]')
6 plt.show()
```
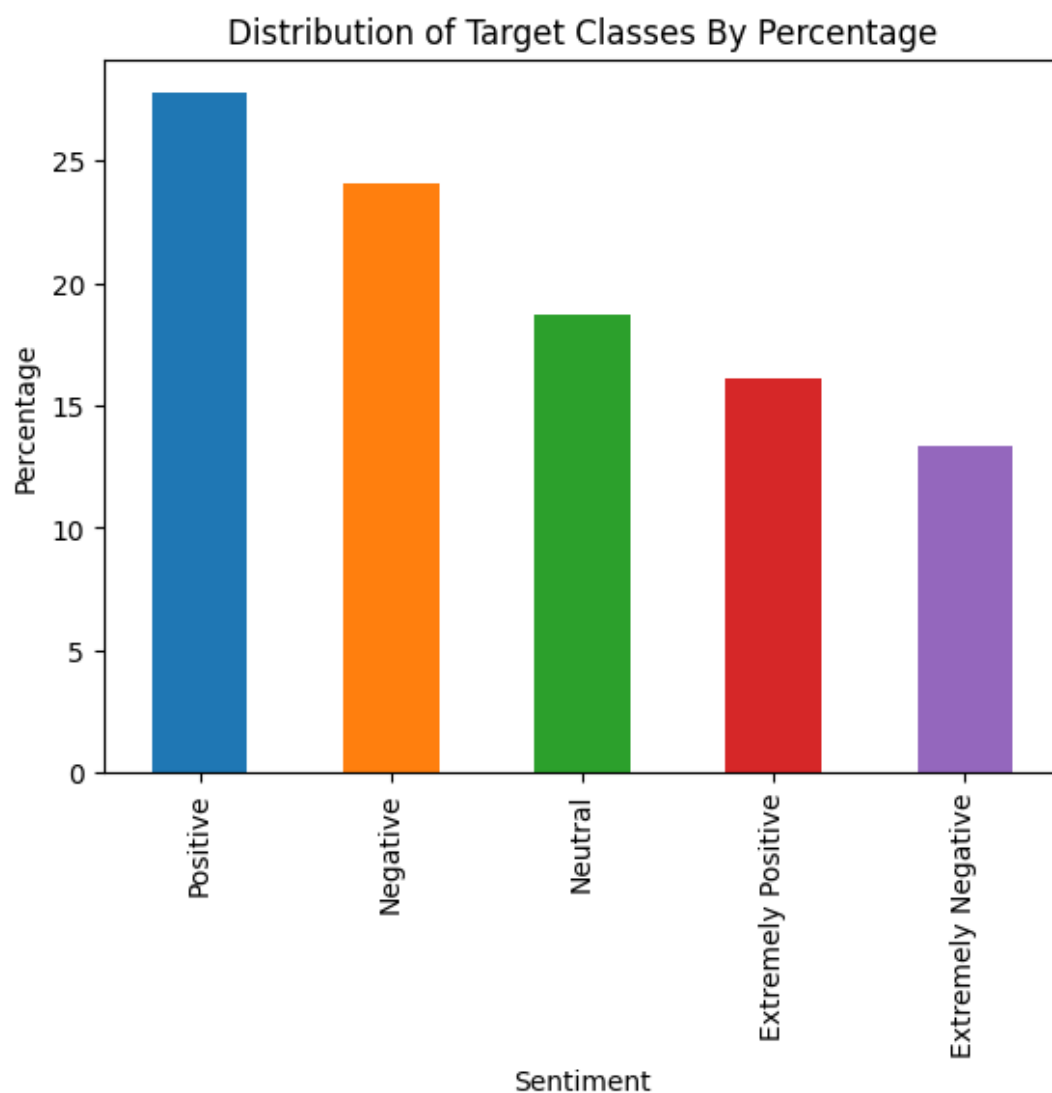


Distribution of Target Classes [Test Data]

▼ Classes Distribution in percentage [Training Data]

```
 1 perct_df = (
 2     train_data.groupby('Sentiment')
 3     .size().sort_values(ascending=False) /
 4     train_data.groupby('Sentiment')
 5     .size().sort_values(ascending=False).sum()
 6     )*100
 7
 8 fig = plt.figure()
 9 perct_df.plot(kind='bar', color=colors)
10 plt.ylabel('Percentage')
11 plt.xlabel('Sentiment')
12 plt.title('Distribution of Target Classes By Percentage [Training
13 plt.show()
```
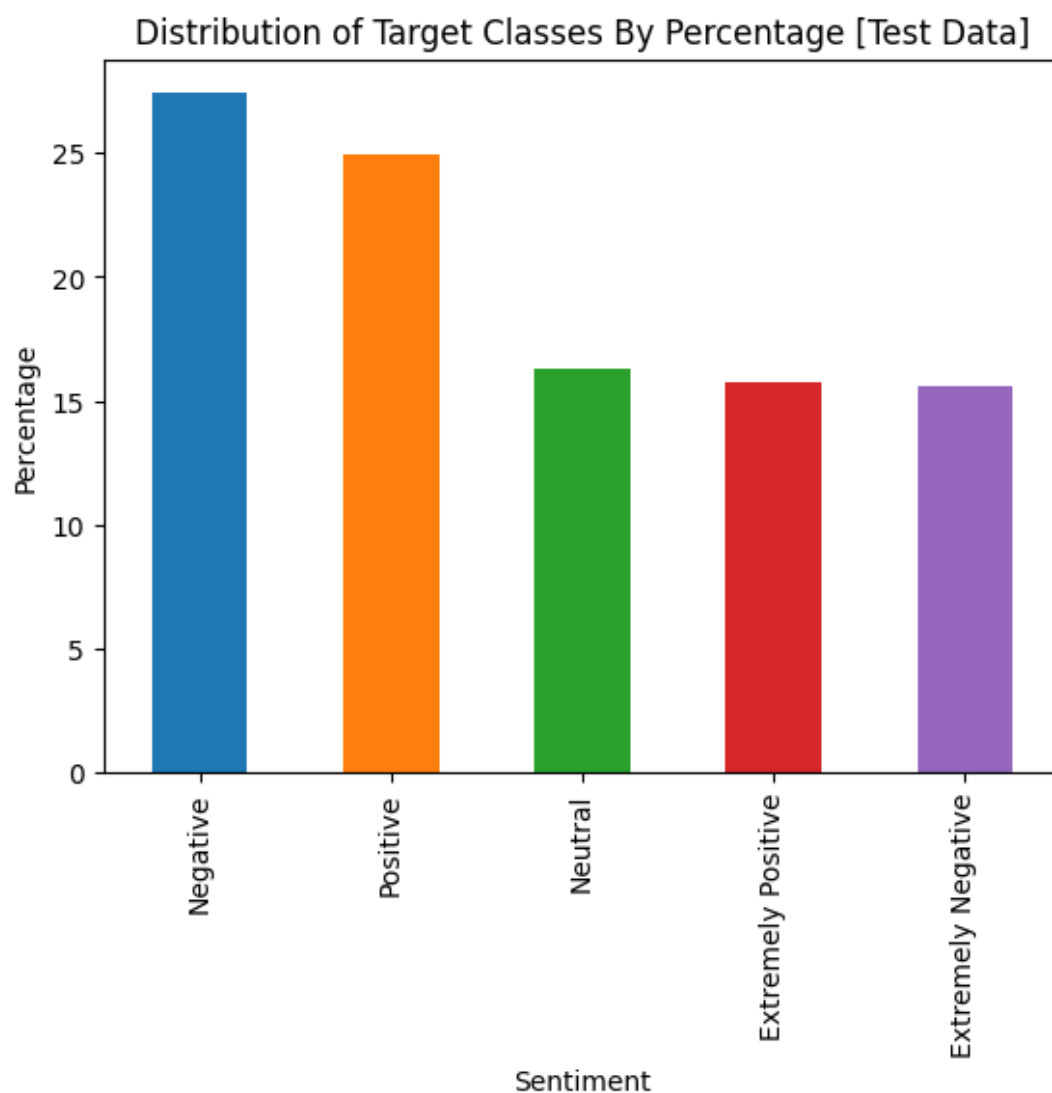


Distribution of Target Classes By Percentage

Classes Distribution in percentage [Test Data]

```
 1 perct_df = (
 2     test_data.groupby('Sentiment')
 3     .size().sort_values(ascending=False) /
 4     test_data.groupby('Sentiment')
 5     .size().sort_values(ascending=False).sum()
 6     )*100
 7
 8 fig = plt.figure()
 9 perct_df.plot(kind='bar', color=colors)
10 plt.ylabel('Percentage')
11 plt.xlabel('Sentiment')
12 plt.title('Distribution of Target Classes By Percentage [Test Data
13 plt.show()
```



Distribution of Target Classes By Percentage [Test Data]

▾ Preprocessing

## Dropping unwanted columns

```
1  train_data.columns
```

```
Index(['UserName', 'ScreenName', 'Location', 'TweetAt', 'OriginalTweet',
       'Sentiment'],
      dtype='object')
```

```
1  train_data.drop(
2      ['UserName', 'ScreenName','Location', 'TweetAt' ],
3       axis=1, inplace=True
4       )
5  test_data.drop(
6      ['UserName', 'ScreenName','Location', 'TweetAt' ],
7      axis=1, inplace=True
8      )
9  train_data.head()
```

|   | OriginalTweet | Sentiment |
|---|---|---|
| **0** | @MeNyrbie @Phil_Gahan @Chrisitv https://t.co/i... | Neutral |
| **1** | advice Talk to your neighbours family to excha... | Positive |
| **2** | Coronavirus Australia: Woolworths to give elde... | Positive |
| **3** | My food stock is not the only one which is emp... | Positive |
| **4** | Me, ready to go at supermarket during the #COV... | Extremely Negative |

## Converting Sentiments to numeric values

```
 1
 2  encoding = {
 3      'Positive':0, 'Negative':1, 'Neutral':2,
 4      'Extremely Positive':3, 'Extremely Negative':4
 5      }
 6  train_data = train_data.replace({'Sentiment':encoding})
 7  test_data = test_data.replace({'Sentiment':encoding})
 8
 9  train_labels = train_data['Sentiment']
10  test_labels = test_data['Sentiment']
11
12  print(train_data.Sentiment.value_counts())
13
14  train_data.head()
```

```
0    11422
1     9917
2     7713
3     6624
4     5481
Name: Sentiment, dtype: int64
```

| | OriginalTweet | Sentiment |
|---|---|---|
| **0** | @MeNyrbie @Phil_Gahan @Chrisitv https://t.co/i... | 2 |
| **1** | advice Talk to your neighbours family to excha... | 0 |
| **2** | Coronavirus Australia: Woolworths to give elde... | 0 |
| **3** | My food stock is not the only one which is emp... | 0 |
| **4** | Me, ready to go at supermarket during the #COV... | 4 |

## ▾ Models

## ▾ Model Utils

```
1 # Tokenize the tweets
2 tokenizer = Tokenizer(num_words=5000, lower=True)
3 tokenizer.fit_on_texts(train_data['OriginalTweet'])
4 train_sequences = tokenizer.texts_to_sequences(
5     train_data['OriginalTweet'])
6 test_sequences = tokenizer.texts_to_sequences(
7     test_data['OriginalTweet'])
```

```
1 vocab_size = len(tokenizer.word_index)
2 # vocab_size = len(train_sequences)
3 vocab_size
```

    85198

```
1 embedding_dim = min(16, round(vocab_size ** 0.25))
2 print("Recommended embedding dimension:", embedding_dim)
3
4 # embedding_dim = 16
```

    Recommended embedding dimension: 16

```
1 num_classes = len(train_data['Sentiment'].unique())
2 num_classes
```

    5

```
1 maxlen = max(len(x_tr_sqe) for x_tr_sqe in train_sequences)
2 maxlen
```

    64

```
1 # Pad the sequences
2 train_padded = pad_sequences(
3     train_sequences, padding='post', maxlen=maxlen)
4 test_padded = pad_sequences(
5     test_sequences, padding='post', maxlen=maxlen)
```

```
1 from tensorflow.keras.utils import to_categorical
2
3 train_labels = to_categorical(train_labels)
4 test_labels = to_categorical(test_labels)
```

```
1 from tensorflow.keras.callbacks import EarlyStopping
2
3 early_stopping = EarlyStopping(
4     min_delta=0.001,
5     patience=4,
6     restore_best_weights=True,
7 )
```

## Sequential

```python
# Define the model
model_seq = Sequential()
model_seq.add(
    Embedding(
        vocab_size, embedding_dim, input_length=maxlen
        ))
model_seq.add(Flatten())
model_seq.add(Dense(16, activation='relu'))
model_seq.add(Dense(num_classes, activation='softmax'))

model_seq.compile(
    optimizer='adam',
    loss='categorical_crossentropy',
    metrics=['accuracy']
    )

# Train the model
history_seq = model_seq.fit(
    train_padded, train_labels,
    callbacks=[early_stopping], epochs=10,
    batch_size=32, validation_split=0.2
    )
```

```
Epoch 1/10
1029/1029 [==============================] – 25s 23ms/step – loss: 1.4144 –
Epoch 2/10
1029/1029 [==============================] – 25s 24ms/step – loss: 0.9268 –
Epoch 3/10
1029/1029 [==============================] – 24s 23ms/step – loss: 0.6900 –
Epoch 4/10
1029/1029 [==============================] – 25s 24ms/step – loss: 0.5331 –
Epoch 5/10
1029/1029 [==============================] – 22s 21ms/step – loss: 0.3988 –
Epoch 6/10
1029/1029 [==============================] – 25s 24ms/step – loss: 0.2870 –
Epoch 7/10
1029/1029 [==============================] – 26s 25ms/step – loss: 0.1988 –
```

```
1 # Evaluate the model
2 loss_seq, accuracy_seq = model_seq.evaluate(
3     test_padded, test_labels)
4 print(f'Test Accuracy with Sequential: {accuracy_seq}')
5 print(f'Test Loss with Sequential: {loss_seq}')
```

```
119/119 [==============================] – 0s 2ms/step – loss: 0.9686 – acc
Test Accuracy with Sequential: 0.6327013969421387
Test Loss with Sequential: 0.9686248898506165
```
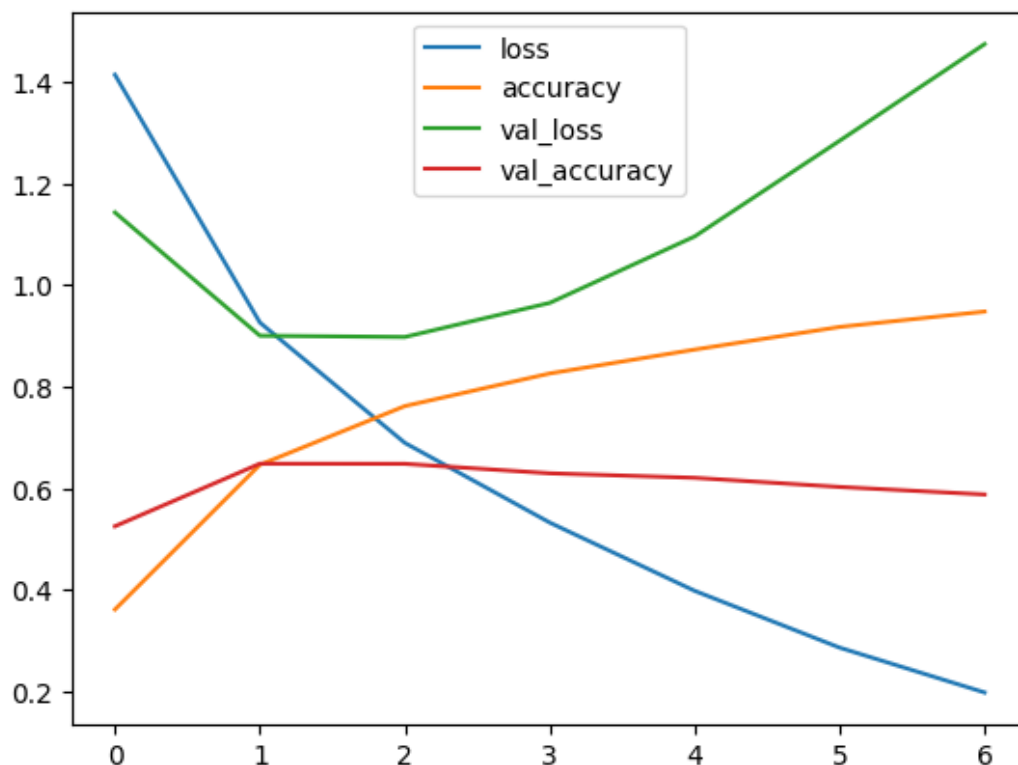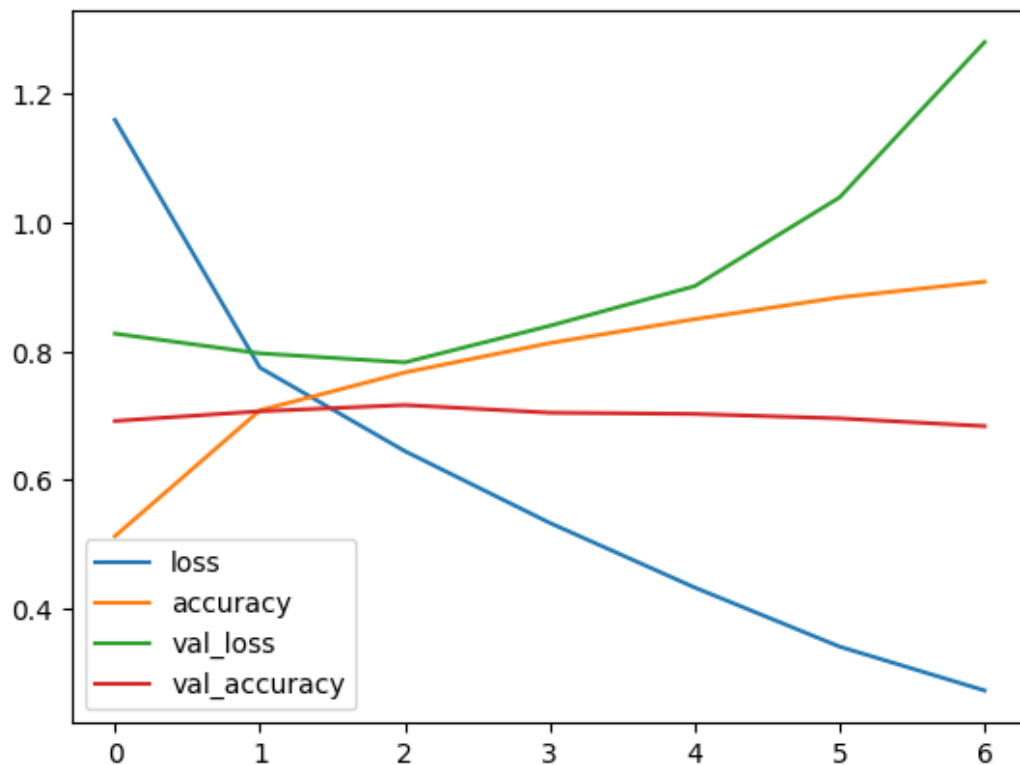
```
1 pd.DataFrame(history_seq.history).plot()
```

```
<Axes: >
```



```
1 # Save and delete the model and history
2 model_seq.save(dir+'model_seq')
3 del model_seq
4
5 with open(dir+'history_seq', 'wb') as file_pi:
6     pickle.dump(history_seq, file_pi)
7 del history_seq
```

```
WARNING:absl:Found untraced functions such as _update_step_xla while saving
```

## CNN

```python
1  filters = 128
2  kernel_size = 5
3
4  # Define the CNN model
5  model_cnn = Sequential()
6  model_cnn.add(
7      Embedding(vocab_size , embedding_dim))
8  model_cnn.add(
9      Conv1D(filters , kernel_size , activation = 'relu'
10     ))
11 model_cnn.add(GlobalMaxPooling1D())
12 model_cnn.add(Dense(512 , activation = 'relu'))
13 model_cnn.add(Dropout(0.2))
14 model_cnn.add(Dense(num_classes , activation = 'softmax'))
15
16 # Compile the model
17 model_cnn.compile(
18     optimizer='adam',
19     loss='categorical_crossentropy',
20     metrics=['accuracy'])
21
22
23 # Train the model
24 history_cnn = model_cnn.fit(
25     train_padded, train_labels,
26     callbacks=[early_stopping], epochs=10,
27     batch_size=32, validation_split=0.2)
```

```
Epoch 1/10
1029/1029 [==============================] – 32s 31ms/step – loss: 1.1594 –
Epoch 2/10
1029/1029 [==============================] – 30s 30ms/step – loss: 0.7746 –
Epoch 3/10
1029/1029 [==============================] – 30s 29ms/step – loss: 0.6449 –
Epoch 4/10
1029/1029 [==============================] – 30s 29ms/step – loss: 0.5333 –
Epoch 5/10
1029/1029 [==============================] – 30s 29ms/step – loss: 0.4329 –
Epoch 6/10
1029/1029 [==============================] – 31s 30ms/step – loss: 0.3411 –
Epoch 7/10
1029/1029 [==============================] – 31s 30ms/step – loss: 0.2730 –
```

```
1 # Evaluate the model
2 loss_cnn, accuracy_cnn = model_cnn.evaluate(
3     test_padded, test_labels)
4 print(f'Test Accuracy with CNN: {accuracy_cnn}')
5 print(f'Test Loss with CNN: {loss_cnn}')
```

```
119/119 [==============================] – 0s 3ms/step – loss: 0.8610 – acc
Test Accuracy with CNN: 0.6750921607017517
Test Loss with CNN: 0.8609734177589417
```

```
1 pd.DataFrame(history_cnn.history).plot()
```

```
<Axes: >
```



```
1 # Save and delete the model and history
2 model_cnn.save(dir+'model_cnn')
3 del model_cnn
4
5 with open(dir+'history_cnn', 'wb') as file_pi:
6     pickle.dump(history_cnn, file_pi)
7 del history_cnn
```

```
WARNING:absl:Found untraced functions such as _jit_compiled_convolution_op,
```

▾ Simple RNN

```python
# Define the RNN model
model_rnn = Sequential()
model_rnn.add(
    Embedding(vocab_size, embedding_dim,
              input_length=maxlen))
model_rnn.add(
    Bidirectional(SimpleRNN(32 , return_sequences = True)))
model_rnn.add(BatchNormalization())
model_rnn.add(Bidirectional(SimpleRNN(64)))
model_rnn.add(BatchNormalization())
model_rnn.add(Dense(512 , activation = 'relu'))
model_rnn.add(Dropout(0.2))
model_rnn.add(Dense(5 , activation = 'softmax'))

# Compile the model
model_rnn.compile(
    optimizer='adam', loss='categorical_crossentropy',
    metrics=['accuracy'])

# Train the model
history_rnn = model_rnn.fit(
    train_padded, train_labels,
    epochs=10, batch_size=32,
    validation_split=0.2)

```

```
Epoch 1/10
1029/1029 [==============================] – 113s 104ms/step – loss: 1.5177
Epoch 2/10
1029/1029 [==============================] – 102s 100ms/step – loss: 1.1001
Epoch 3/10
1029/1029 [==============================] – 104s 101ms/step – loss: 0.8624
Epoch 4/10
1029/1029 [==============================] – 105s 102ms/step – loss: 0.7748
Epoch 5/10
1029/1029 [==============================] – 101s 98ms/step – loss: 0.7151
Epoch 6/10
1029/1029 [==============================] – 102s 99ms/step – loss: 0.6473
Epoch 7/10
1029/1029 [==============================] – 100s 97ms/step – loss: 0.6147
Epoch 8/10
1029/1029 [==============================] – 100s 97ms/step – loss: 0.5866
Epoch 9/10
1029/1029 [==============================] – 101s 98ms/step – loss: 0.5035
Epoch 10/10
1029/1029 [==============================] – 105s 102ms/step – loss: 0.4597
```

```
1 # Evaluate the model
2 loss_rnn, accuracy_rnn = model_rnn.evaluate(
3     test_padded, test_labels)
4 print(f'Test Accuracy with Simple RNN: {accuracy_rnn}')
5 print(f'Test Loss with Simple RNN: {loss_rnn}')
```

```
119/119 [==============================] - 2s 16ms/step - loss: 2.8046 - ac
Test Accuracy with Simple RNN: 0.5589784383773804
Test Loss with Simple RNN: 2.804551839828491
```
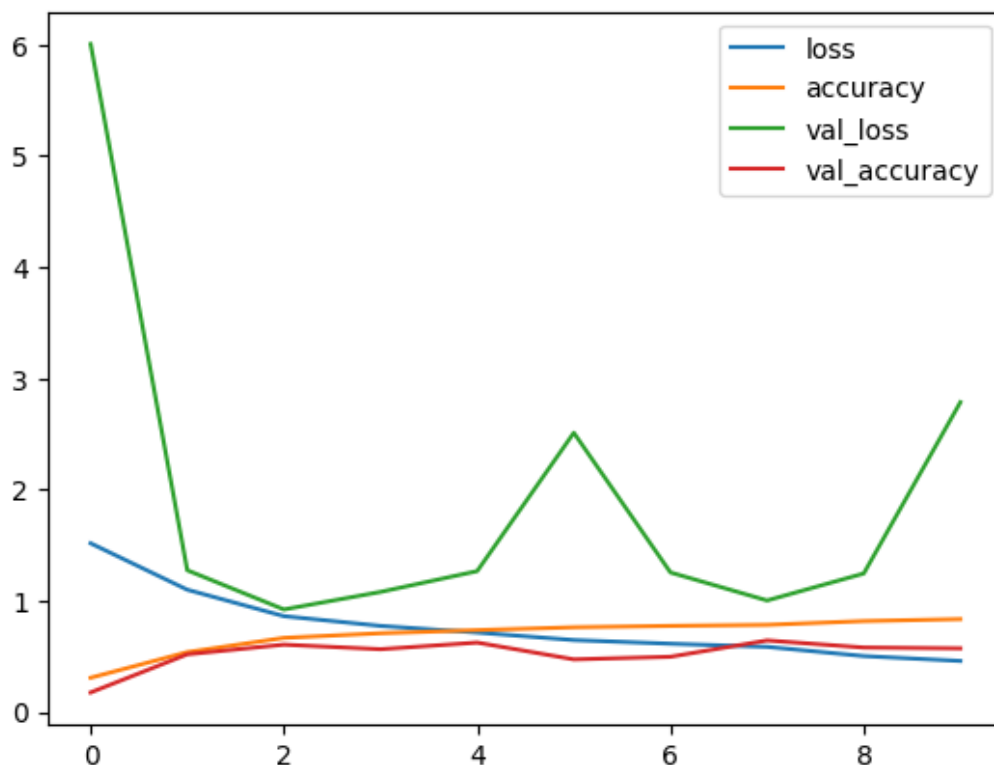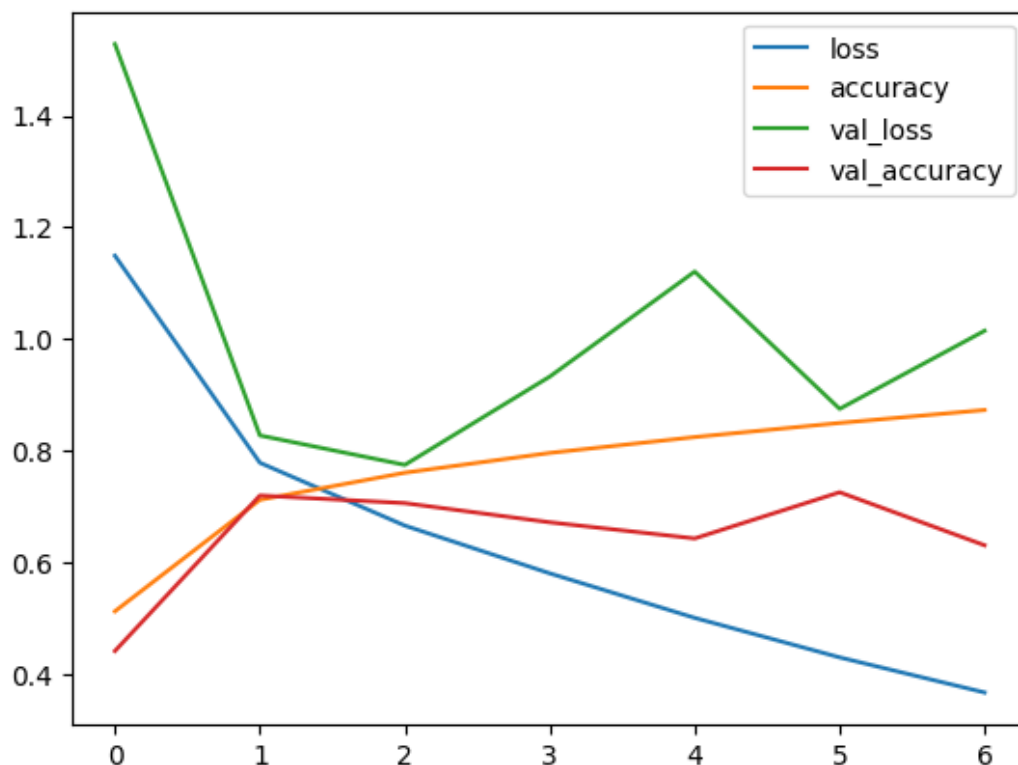
```
1 pd.DataFrame(history_rnn.history).plot()
```

```
<Axes: >
```



```
1 # Save and delete the model and history
2 model_rnn.save(dir+'model_rnn')
3 del model_rnn
4
5 with open(dir+'history_rnn', 'wb') as file_pi:
6     pickle.dump(history_rnn, file_pi)
7 del history_rnn
```

```
WARNING:absl:Found untraced functions such as _update_step_xla while saving
```

## LSTM

```python
1  # Define the model
2  model_lstm = Sequential()
3  model_lstm.add(
4      Embedding(vocab_size, embedding_dim,
5                  input_length=maxlen))
6  model_lstm.add(
7      Bidirectional(LSTM(32 , return_sequences = True)))
8  model_lstm.add(BatchNormalization())
9  model_lstm.add(Bidirectional(LSTM(64)))
10 model_lstm.add(BatchNormalization())
11 model_lstm.add(Dense(512 , activation = 'relu'))
12 model_lstm.add(Dropout(0.2))
13 model_lstm.add(Dense(num_classes , activation = 'softmax'))
14
15 model_lstm.compile(
16     optimizer='adam',
17     loss='categorical_crossentropy',
18     metrics=['accuracy'])
19
20 # Train the model
21 history_lstm = model_lstm.fit(
22     train_padded, train_labels,
23     callbacks=[early_stopping],
24     epochs=10, batch_size=32,
25     validation_split=0.2)
26
```

```
Epoch 1/10
1029/1029 [==============================] - 204s 188ms/step - loss: 1.1492
Epoch 2/10
1029/1029 [==============================] - 202s 196ms/step - loss: 0.7783
Epoch 3/10
1029/1029 [==============================] - 186s 181ms/step - loss: 0.6656
Epoch 4/10
1029/1029 [==============================] - 189s 184ms/step - loss: 0.5801
Epoch 5/10
1029/1029 [==============================] - 201s 196ms/step - loss: 0.5003
Epoch 6/10
1029/1029 [==============================] - 204s 198ms/step - loss: 0.4295
Epoch 7/10
1029/1029 [==============================] - 201s 196ms/step - loss: 0.3666
```

```
1 # Evaluate the model
2 loss_lstm, accuracy_lstm = model_lstm.evaluate(
3     test_padded, test_labels)
4 print(f'Test Accuracy with LSTM: {accuracy_lstm}')
5 print(f'Test Loss with LSTM: {loss_lstm}')
```

```
119/119 [==============================] – 4s 33ms/step – loss: 0.8378 – ac
Test Accuracy with LSTM: 0.6829910278320312
Test Loss with LSTM: 0.837775468826294
```

```
1 pd.DataFrame(history_lstm.history).plot()
```

```
<Axes: >
```



```
1 # Save and delete the model and history
2 model_lstm.save(dir+'model_lstm')
3 del model_lstm
4
5 with open(dir+'history_lstm', 'wb') as file_pi:
6     pickle.dump(history_lstm, file_pi)
7 del history_lstm
```

```
WARNING:absl:Found untraced functions such as _update_step_xla, lstm_cell_1
```

## Pretrained LSTM - Pretrained Word Embeddings

```
1 # Download GloVe embeddings
2 !wget http://nlp.stanford.edu/data/glove.6B.zip
3 !unzip -q glove.6B.zip
```

```
--2023-04-20 06:54:55--  http://nlp.stanford.edu/data/glove.6B.zip
Resolving nlp.stanford.edu (nlp.stanford.edu)... 171.64.67.140
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:80... conn
HTTP request sent, awaiting response... 302 Found
Location: https://nlp.stanford.edu/data/glove.6B.zip [following]
--2023-04-20 06:54:55--  https://nlp.stanford.edu/data/glove.6B.zip
Connecting to nlp.stanford.edu (nlp.stanford.edu)|171.64.67.140|:443... con
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://downloads.cs.stanford.edu/nlp/data/glove.6B.zip [followin
--2023-04-20 06:54:55--  https://downloads.cs.stanford.edu/nlp/data/glove.6
Resolving downloads.cs.stanford.edu (downloads.cs.stanford.edu)... 171.64.6
Connecting to downloads.cs.stanford.edu (downloads.cs.stanford.edu)|171.64.
HTTP request sent, awaiting response... 200 OK
Length: 862182613 (822M) [application/zip]
Saving to: 'glove.6B.zip'

glove.6B.zip          100%[===================>] 822.24M  5.04MB/s    in 2m 4

2023-04-20 06:57:35 (5.14 MB/s) - 'glove.6B.zip' saved [862182613/862182613
```

```
1 # Load the pre-trained embeddings into a dictionary
2 embeddings_index = {}
3 with open('glove.6B.100d.txt', 'r') as f:
4     for line in f:
5         values = line.split()
6         word = values[0]
7         coefs = np.asarray(values[1:], dtype='float32')
8         embeddings_index[word] = coefs
```

```
1 # Create an embedding matrix for the words in the training data
2 embedding_dim = 100
3 embedding_matrix = np.zeros((vocab_size, embedding_dim))
4 for word, i in tokenizer.word_index.items():
5     embedding_vector = embeddings_index.get(word)
6     if embedding_vector is not None:
7         embedding_matrix[i] = embedding_vector
```

```
1 # Define the model with pre-trained embeddings
2 model_pretrained = Sequential()
3 model_pretrained.add(
```

```
 4      Embedding(
 5          vocab_size, embedding_dim,
 6          input_length=maxlen,
 7          weights=[embedding_matrix],
 8          trainable=False))
 9 model_pretrained.add(
10     Bidirectional(LSTM(32 , return_sequences = True)))
11 model_pretrained.add(BatchNormalization())
12 model_pretrained.add(Bidirectional(LSTM(64)))
13 model_pretrained.add(BatchNormalization())
14 model_pretrained.add(Dense(512 , activation = 'relu'))
15 model_pretrained.add(Dropout(0.2))
16 model_pretrained.add(
17     Dense(num_classes , activation = 'softmax'))
18
19 # Compile the model
20 model_pretrained.compile(
21     optimizer='adam',
22     loss='categorical_crossentropy',
23     metrics=['accuracy'])
24
25 # Train the model
26 history_pretrained = model_pretrained.fit(
27     train_padded, train_labels,
28     callbacks=[early_stopping],
29     epochs=10, batch_size=32,
30     validation_split=0.2)
```

```
Epoch 1/10
1029/1029 [==============================] – 214s 197ms/step – loss: 1.2831
Epoch 2/10
1029/1029 [==============================] – 189s 183ms/step – loss: 1.0389
Epoch 3/10
1029/1029 [==============================] – 188s 183ms/step – loss: 0.8872
Epoch 4/10
1029/1029 [==============================] – 196s 190ms/step – loss: 0.8013
Epoch 5/10
1029/1029 [==============================] – 186s 181ms/step – loss: 0.7358
Epoch 6/10
1029/1029 [==============================] – 202s 196ms/step – loss: 0.6900
Epoch 7/10
1029/1029 [==============================] – 199s 193ms/step – loss: 0.6428
Epoch 8/10
1029/1029 [==============================] – 201s 195ms/step – loss: 0.6001
Epoch 9/10
1029/1029 [==============================] – 202s 196ms/step – loss: 0.5553
```
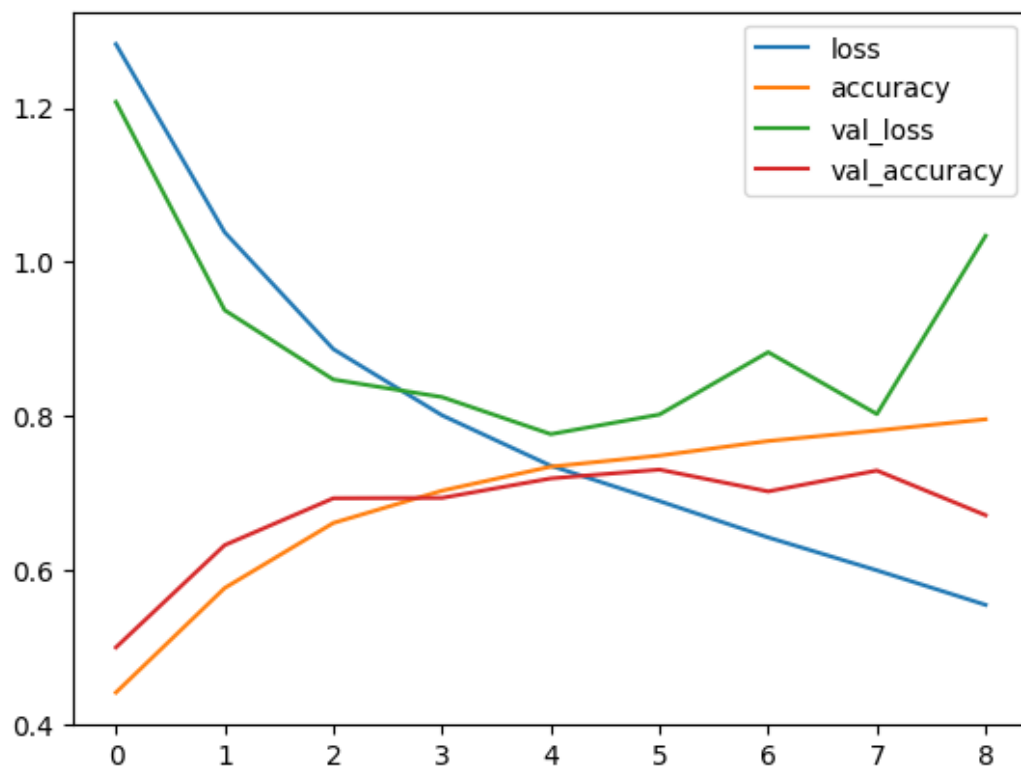
```
1 # Evaluate the model
2 loss_pretrained, accuracy_pretrained = model_pretrained.evaluate(
3     test_padded, test_labels)
4 print(
5     f'Test Accuracy with LSTM Pretrained: {accuracy_pretrained}'
6     )
7 print(
8     f'Test Loss with LSTM Pretrained: {loss_pretrained}'
9     )
```

```
119/119 [==============================] - 7s 59ms/step - loss: 0.8615 - ac
Test Accuracy with LSTM Pretrained: 0.6803581118583679
Test Loss with LSTM Pretrained: 0.8615074753761292
```

```
1 pd.DataFrame(history_pretrained.history).plot()
```

<Axes: >

```
1 # Save and delete the model and history
2 model_pretrained.save(dir+'model_pretrained')
3 del model_pretrained
4
5 with open(dir+'history_pretrained', 'wb') as file_pi:
6     pickle.dump(history_pretrained, file_pi)
7 del history_pretrained
```

WARNING:absl:Found untraced functions such as _update_step_xla, lstm_cell_7

## Character Embeddings

```
1 # Convert text to lowercase
2 train_text = train_data["OriginalTweet"].str.lower()
3 test_text = test_data["OriginalTweet"].str.lower()
4
5 # Convert text to sequences of character indices
6 tokenizer = Tokenizer(char_level=True)
7 tokenizer.fit_on_texts(train_text)
8 train_seq = tokenizer.texts_to_sequences(train_text)
9 test_seq = tokenizer.texts_to_sequences(test_text)
10
11 # Pad sequences to a fixed length
12 # maxlen = 140  # Example sequence length
13 train_seq = pad_sequences(train_seq, maxlen=maxlen)
14 test_seq = pad_sequences(test_seq, maxlen=maxlen)
15
16 # Define the model architecture
17 model_char = Sequential()
18 model_char.add(
19     Embedding(
20         input_dim=len(tokenizer.word_index)+1,
21         output_dim=100, input_length=maxlen))
22 model_char.add(
23     Bidirectional(LSTM(32 , return_sequences = True)))
24 model_char.add(BatchNormalization())
25 model_char.add(Bidirectional(LSTM(64)))
26 model_char.add(BatchNormalization())
27 model_char.add(Dense(512 , activation = 'relu'))
28 model_char.add(Dropout(0.2))
29 model_char.add(
30     Dense(num_classes , activation = 'softmax'))
31
32 # Compile the model
```

```
33 optimizer = Adam(learning_rate=1e-4)
34 model_char.compile(
35     loss='categorical_crossentropy',
36     optimizer=optimizer,
37     metrics=['accuracy'])
38
39 # Train the model
40 y_train = pd.get_dummies(
41     train_data["Sentiment"]).values
42 y_test = pd.get_dummies(
43     test_data["Sentiment"]).values
44 history_char = model_char.fit(
45     train_seq, y_train,
46     validation_data=(test_seq, y_test),
47     callbacks=[early_stopping],
48     epochs=10, batch_size=32)
```

```
Epoch 1/10
1287/1287 [==============================] - 267s 198ms/step - loss: 1.6165
Epoch 2/10
1287/1287 [==============================] - 249s 194ms/step - loss: 1.5773
Epoch 3/10
1287/1287 [==============================] - 247s 192ms/step - loss: 1.5615
Epoch 4/10
1287/1287 [==============================] - 249s 194ms/step - loss: 1.5480
Epoch 5/10
1287/1287 [==============================] - 242s 188ms/step - loss: 1.5363
Epoch 6/10
1287/1287 [==============================] - 250s 195ms/step - loss: 1.5280
Epoch 7/10
1287/1287 [==============================] - 248s 193ms/step - loss: 1.5183
Epoch 8/10
1287/1287 [==============================] - 242s 188ms/step - loss: 1.5104
Epoch 9/10
1287/1287 [==============================] - 243s 189ms/step - loss: 1.5025
```
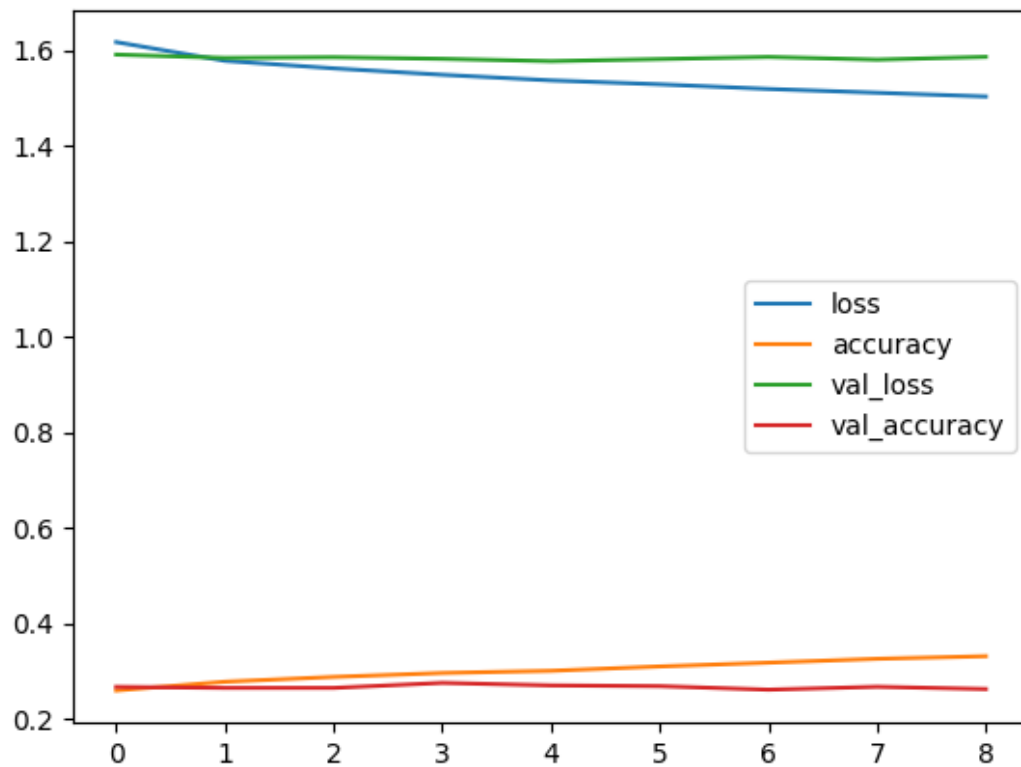
```
1 # Evaluate the model
2 loss_char, accuracy_char = model_char.evaluate(
3     test_seq, y_test, verbose=0)
4 print(
5     f'Test Accuracy with Char Embeddings: {accuracy_char}')
6 print(
7     f'Test Loss with Char Embeddings: {loss_char}')
```

```
Test Accuracy with Char Embeddings: 0.2693522870540619
Test Loss with Char Embeddings: 1.5764895677566528
```

```
1  pd.DataFrame(history_char.history).plot()
```

<Axes: >



```
1 # Save and delete the model and history
2 model_char.save(dir+'model_char')
3 del model_char
4
5 with open(dir+'history_char', 'wb') as file_pi:
6     pickle.dump(history_char, file_pi)
7 del history_char
```

WARNING:absl:Found untraced functions such as _update_step_xla, lstm_cell_1

▾ Model Comparison

▾ Load model history

```
 1 with open(
 2     dir+'history_seq',
 3     "rb"
 4     ) as f:
 5     history_seq = pickle.load(f)
 6
 7 with open(
 8     dir+'history_cnn',
 9     "rb"
10     ) as f:
11     history_cnn = pickle.load(f)
12
13 with open(
14     dir+'history_rnn',
15     "rb"
16     ) as f:
17     history_rnn = pickle.load(f)
18
19 with open(
20     dir+'history_lstm',
21     "rb"
22     ) as f:
23     history_lstm = pickle.load(f)
24
25 with open(
26     dir+'history_pretrained',
27     "rb"
28     ) as f:
29     history_pretrained = pickle.load(f)
30
31 with open(
32     dir+'history_char',
33     "rb"
34     ) as f:
35     history_char = pickle.load(f)
36
```
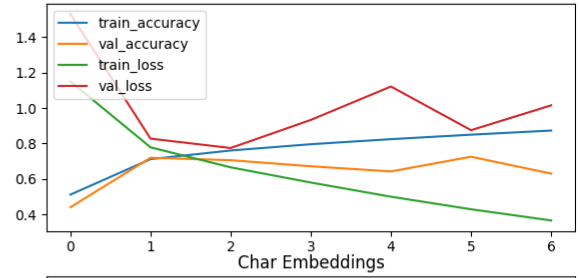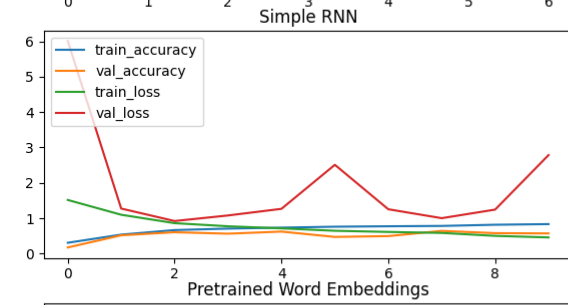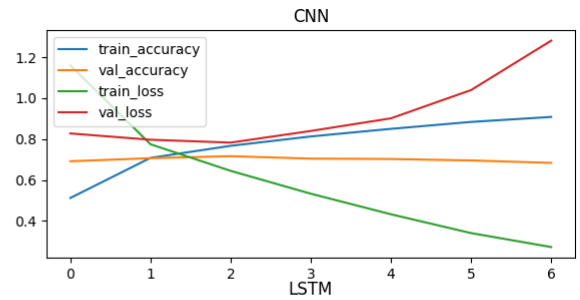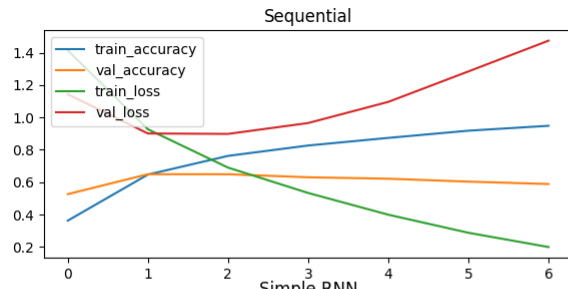
Plot history

```
1 figure, axis = plt.subplots(3, 2,  figsize=(15, 10))
2
3 # Sequential
4 axis[0,0].plot(history_seq.history['accuracy'])
5 axis[0,0].plot(history_seq.history['val_accuracy'])
```

```python
 6 axis[0,0].plot(history_seq.history['loss'])
 7 axis[0,0].plot(history_seq.history['val_loss'])
 8 axis[0,0].legend(
 9     ['train_accuracy', 'val_accuracy',
10     'train_loss', 'val_loss'], loc='upper left')
11 axis[0,0].set_title("Sequential")
12
13 # CNN
14 axis[0,1].plot(history_cnn.history['accuracy'])
15 axis[0,1].plot(history_cnn.history['val_accuracy'])
16 axis[0,1].plot(history_cnn.history['loss'])
17 axis[0,1].plot(history_cnn.history['val_loss'])
18 axis[0,1].legend(
19     ['train_accuracy', 'val_accuracy',
20     'train_loss', 'val_loss'], loc='upper left')
21 axis[0,1].set_title("CNN")
22
23 # Simple RNN
24 axis[1,0].plot(history_rnn.history['accuracy'])
25 axis[1,0].plot(history_rnn.history['val_accuracy'])
26 axis[1,0].plot(history_rnn.history['loss'])
27 axis[1,0].plot(history_rnn.history['val_loss'])
28 axis[1,0].legend(
29     ['train_accuracy', 'val_accuracy',
30     'train_loss', 'val_loss'], loc='upper left')
31 axis[1,0].set_title("Simple RNN")
32
33 # LSTM
34 axis[1,1].plot(history_lstm.history['accuracy'])
35 axis[1,1].plot(history_lstm.history['val_accuracy'])
36 axis[1,1].plot(history_lstm.history['loss'])
37 axis[1,1].plot(history_lstm.history['val_loss'])
38 axis[1,1].legend(
39     ['train_accuracy', 'val_accuracy',
40     'train_loss', 'val_loss'], loc='upper left')
41 axis[1,1].set_title("LSTM")
42
43
44 # Pretrained Word Embeddings
45 axis[2,0].plot(history_pretrained.history['accuracy'])
46 axis[2,0].plot(history_pretrained.history['val_accuracy'])
47 axis[2,0].plot(history_pretrained.history['loss'])
48 axis[2,0].plot(history_pretrained.history['val_loss'])
49 axis[2,0].legend(
50     ['train_accuracy', 'val_accuracy',
51     'train_loss', 'val_loss'], loc='upper left')
```

```
52 axis[2,0].set_title("Pretrained Word Embeddings")
53
54 # Char Embeddings
55 axis[2,1].plot(history_char.history['accuracy'])
56 axis[2,1].plot(history_char.history['val_accuracy'])
57 axis[2,1].plot(history_char.history['loss'])
58 axis[2,1].plot(history_char.history['val_loss'])
59 axis[2,1].legend(
60     ['train_accuracy', 'val_accuracy',
61     'train_loss', 'val_loss'], loc='upper left')
62 axis[2,1].set_title("Char Embeddings")
```

Text(0.5, 1.0, 'Char Embeddings')



Analysis of the performance of various approaches

# Analysis of the performance of the various approaches:

## Simple Sequential Model:

The simple sequential model achieved a test accuracy of around 63%, which is not bad for a simple model. However, the model may not be able to capture complex relationships between words and may suffer from overfitting.

The Sequential Model plot above shows that the training loss continues to decrease while the validation loss continues to increase from epoch 1, this indicates overfitting. The model is optimizing the training data too well and is starting to overfit to noise and outliers in the data.

## CNN Model:

The CNN model achieved a test accuracy of around 67%, which is comparable to the LSTM model. The CNN model uses convolutions to extract local features from the text data, which can be useful for tasks like text classification.

The CNN Model plot above shows that training accuracy continues to increase while the validation accuracy is plateau at 0.7, this indicates overfitting. The model is starting to memorize the training data and is no longer able to generalize well to new data.

## Simple RNN Model:

The Simple RNN model achieved a test accuracy of around 55%, which is less as compared to the CNN and the LSTM models.

## LSTM Model:

The LSTM model performed better than the simple sequential model, achieving a test accuracy of around 68.2%. This is because the LSTM model can capture the sequential nature of text data and remember important information from earlier time steps.

The LSTM Model plot above shows that the training loss continues to decrease but the validation loss increases from epoch 2, this indicates overfitting. The model is optimizing the training data too well and is starting to overfit to noise and outliers in the data.

## Pretrained Word Embeddings LSTM Model:

Using pretrained word embeddings like GloVe improved the performance of the model, achieving a test accuracy of 68%. This is because pre-trained embeddings can capture the semantic meaning of words and improve the model's ability to understand the text data.

The Pretrained Word Embeddings LSTM Model shows the train and validation accuracy increaes and the train and validation loss decreases. So there is no overfitting and the

accuracy of the model is best in all the models.

## Character Embeddings:

Using character-level embeddings the model performed poorly, only achieving a test accuracy of 26%. By far the lowest accuracy.

---

## Conslusion:

**The LSTM model with pre-trained GloVe word embeddings is the best performing model** among the ones tested. However, further improvements can be made by experimenting with different hyperparameters, using more advanced architectures like transformers, and incorporating additional features.