# Natural Language Processing
# Assignment – 2

## Dataset

The dataset used in this assignment is a collection of text messages labelled as 'spam' or 'ham', intended for **spam detection**. It is loaded from a CSV file, where each row corresponds to a message alongside its classification (spam or ham).

The shape of dataset is **5572 samples and 2 features** which we renamed as 'target' and 'sentence'.

```
3   # Load spam
4   spam_df = pd.read_csv('/content/drive/MyDrive/NLP Assignment 2/DataSet/spam.csv', header=0, names = ['target', 'sentence'])
5   print(spam_df.head())
6
```

```
  target                                           sentence
0    ham  Go until jurong point, crazy.. Available only ...
1    ham                      Ok lar... Joking wif u oni...
2   spam  Free entry in 2 a wkly comp to win FA Cup fina...
3    ham  U dun say so early hor... U c already then say...
4    ham  Nah I don't think he goes to usf, he lives aro...
```

```
1   spam_df.shape
```

```
(5572, 2)
```

```
1   spam_df.head()
```

| | target | sentence |
|---|---|---|
| 0 | ham | Go until jurong point, crazy.. Available only ... |
| 1 | ham | Ok lar... Joking wif u oni... |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... |
| 3 | ham | U dun say so early hor... U c already then say... |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... |

**Dataset Link**: https://drive.google.com/file/d/1aCGEe3gHAw5rxJCYi-qW1eNZV42hIXzb/view?usp=drive_link

## Preprocessing

The preprocessing steps applied to the dataset:

**Tokenization**: Splitting text into individual words.

**Lowercasing**: Converting all characters to lowercase to standardize the text.

**Removing Stop Words**: Eliminating common words (e.g., "the", "is", "in") that provide little value in distinguishing between spam and ham messages.

**Removing Non-Alphabetic Tokens**: Excluding symbols and numbers.

**Lemmatization**: Reducing words to their base or root form, using the context to determine part-of-speech tags for accurate lemmatization.

```
1   from nltk.stem import WordNetLemmatizer
2   from nltk.corpus import wordnet
3
4   # Download necessary NLTK data
5   nltk.download('wordnet')
6   nltk.download('averaged_perceptron_tagger')  # For part-of-speech tagging
7
8   # Initialize the WordNetLemmatizer
9   lemmatizer = WordNetLemmatizer()
10
11  # Function to convert NLTK's part-of-speech tags to WordNet's part-of-speech names
12  def get_wordnet_pos(treebank_tag):
13      if treebank_tag.startswith('J'):
14          return wordnet.ADJ
15      elif treebank_tag.startswith('V'):
16          return wordnet.VERB
17      elif treebank_tag.startswith('N'):
18          return wordnet.NOUN
19      elif treebank_tag.startswith('R'):
20          return wordnet.ADV
21      else:
22          return wordnet.NOUN  # Default to noun
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
```

**Preprocessing DataSet + Vocabulary Count**

```
1   # tokenization, lowercasing, removing Stop Words, removing non-alphabetic tokens, lammetization done
2   #
3   def clean_text_tokens(text):
4       # Convert text to lowercase
5       text = text.lower()
6       # Tokenize text
7       tokens = word_tokenize(text)
8       # Remove stop words
9       stop_words = set(stopwords.words('english'))
10      filtered_tokens = [word for word in tokens if word not in stop_words and word.isalpha()]
11
12      # Part-of-speech tagging for each token
13      pos_tags = pos_tag(filtered_tokens)
14
15      # Lemmatization using the appropriate part-of-speech tag
16      lemmatized_tokens = [lemmatizer.lemmatize(word, get_wordnet_pos(pos)) for word, pos in pos_tags]
17
18      return lemmatized_tokens
19
20
```

**Vocabulary Size and Distribution Analysis**

The **vocabulary size after lemmatization** was found to be **6,233 unique words**.

This indicates a fairly large lexicon for the given dataset of 5,572 messages, which means a wide variety of language use within the corpus.

**Lemmatization** plays a **critical** role in this analysis because it consolidates different forms of a word into a single, base form, **reducing redundancy** and **highlighting** the **true breadth** of the vocabulary used.

```
1   # spam
2   spam_df['tokens'] = spam_df['sentence'].apply(clean_text_tokens)
3
4   # For the 'target' column, convert it into numerical values (one-hot encoding)
5   spam_df['target'] = pd.get_dummies(spam_df['target'], drop_first=True)
6
7   # Calculate unique words in sentences after lemmatization
8   unique_words_sentences = set(word for tokens_list in spam_df['tokens'] for word in tokens_list)
9   vocabulary_size_sentences = len(unique_words_sentences)
10
11  print(f'Vocabulary Size in Dataset After Lemmatization: {vocabulary_size_sentences}')
12
```

Vocabulary Size in Dataset After Lemmatization: 6233

```
1   spam_df.head()
```

| | target | sentence | tokens | word_count | char_count |
|---|---|---|---|---|---|
| 0 | 0 | Go until jurong point, crazy.. Available only ... | [go, jurong, point, crazy, available, bugis, n... | 16 | 111 |
| 1 | 0 | Ok lar... Joking wif u oni... | [ok, lar, joking, wif, u, oni] | 6 | 29 |
| 2 | 1 | Free entry in 2 a wkly comp to win FA Cup fina... | [free, entry, wkly, comp, win, fa, cup, final,... | 20 | 155 |
| 3 | 0 | U dun say so early hor... U c already then say... | [u, dun, say, early, hor, u, c, already, say] | 9 | 49 |
| 4 | 0 | Nah I don't think he goes to usf, he lives aro... | [nah, think, go, usf, life, around, though] | 7 | 61 |

## Distribution of String Lengths

String lengths were analysed in terms of both word count (word_count) and character count (char_count).

The word and character counts for each message were added as new columns to the dataset, which will make our visualization easy.

**Distribution of String Lengths + Visualization**

```
1
2   spam_df['word_count'] = spam_df['tokens'].apply(len)
3   spam_df['char_count'] = spam_df['sentence'].apply(len)
4
```

```python
import matplotlib.pyplot as plt
import seaborn as sns
import nltk

```

```python
def visualize_distributions(word_counts, char_counts, title):
    plt.figure(figsize=(7, 3))

    plt.subplot(1, 2, 1)
    sns.histplot(word_counts, bins=30, kde=True)
    plt.title(f'Word Count Distribution: {title}')
    plt.xlabel('Word Count')
    plt.ylabel('Frequency')

    plt.subplot(1, 2, 2)
    sns.boxplot(word_counts)
    plt.title(f'Word Count Box Plot: {title}')
    plt.xlabel('Word Count')

    plt.figure(figsize=(14, 6))

    plt.subplot(1, 2, 1)
    sns.histplot(char_counts, bins=30, kde=True)
    plt.title(f'Character Count Distribution: {title}')
    plt.xlabel('Character Count')
    plt.ylabel('Frequency')

    plt.subplot(1, 2, 2)
    sns.boxplot(char_counts)
    plt.title(f'Character Count Box Plot: {title}')
    plt.xlabel('Character Count')

    plt.tight_layout()
    plt.show()

```
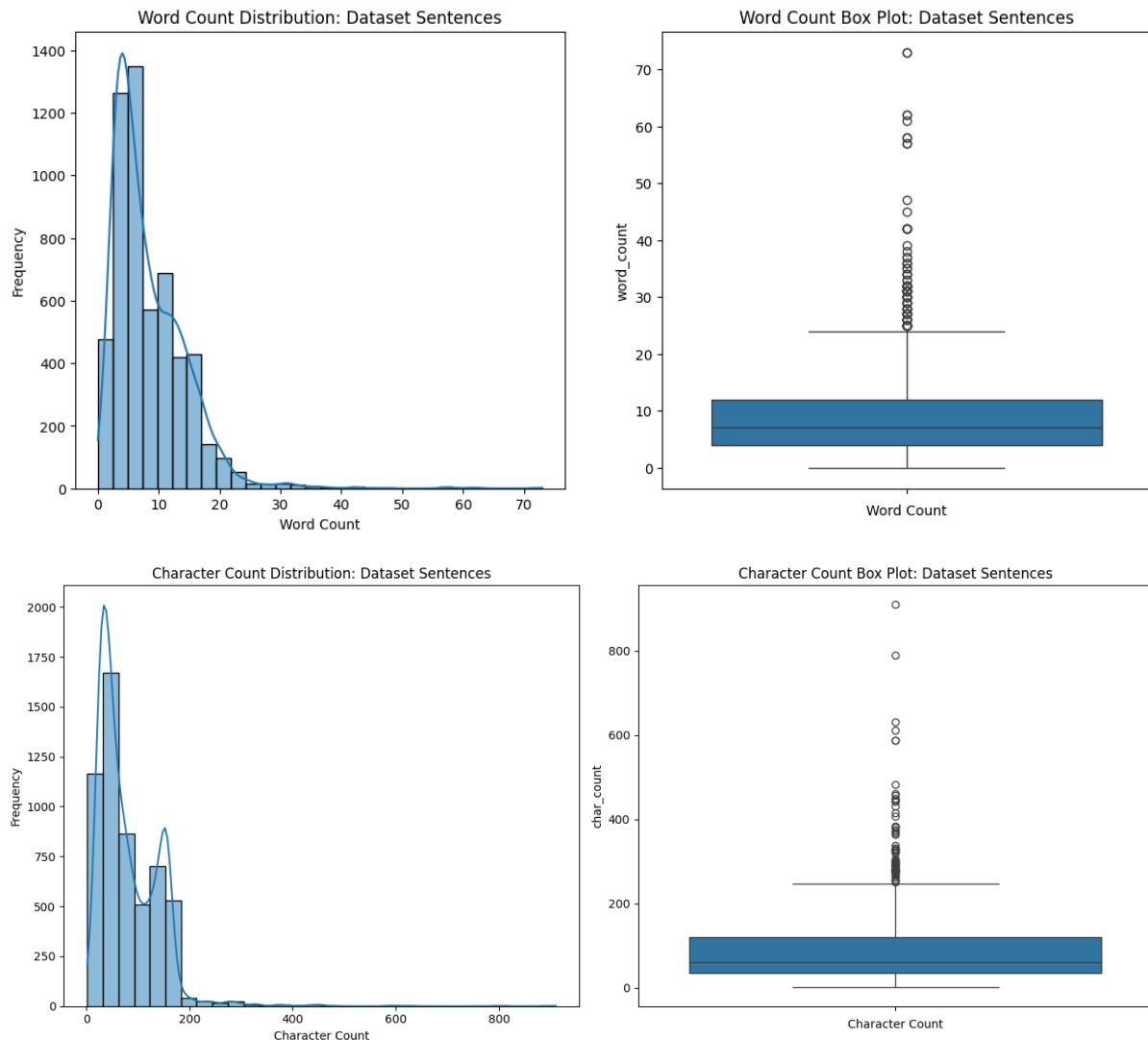
```python
# Dataset Sentences
visualize_distributions(spam_df['word_count'], spam_df['char_count'], 'Dataset Sentences')

```

The **histograms** for both word count and character count distributions are **skewed to the right**, indicating that most **messages contain fewer words and characters**, with the frequency tapering off as the number of words or characters increases.

The **word count histogram** peaks at around 5 to 10 words, which means that the majority of **messages are quite short**.

The **character count histogram** peaks even more sharply at around 0 to 50 characters per message, reinforcing that most **messages are brief**.

The **word count box plot** indicates a median value close to the lower quartile, suggesting a **concentration of messages with fewer words**. The presence of **outliers** indicates that there are some **messages significantly longer than the rest**. The character count box plot shows a similar pattern with a median value near the lower end and numerous outliers.

**Distribution of Classes**

The distribution of classes (spam vs. ham) is done to understand the balance of the dataset.

In the output, the target column is encoded as a binary variable, with '0' representing ham (non-spam) messages and '1' representing spam messages.

From the class distribution output:

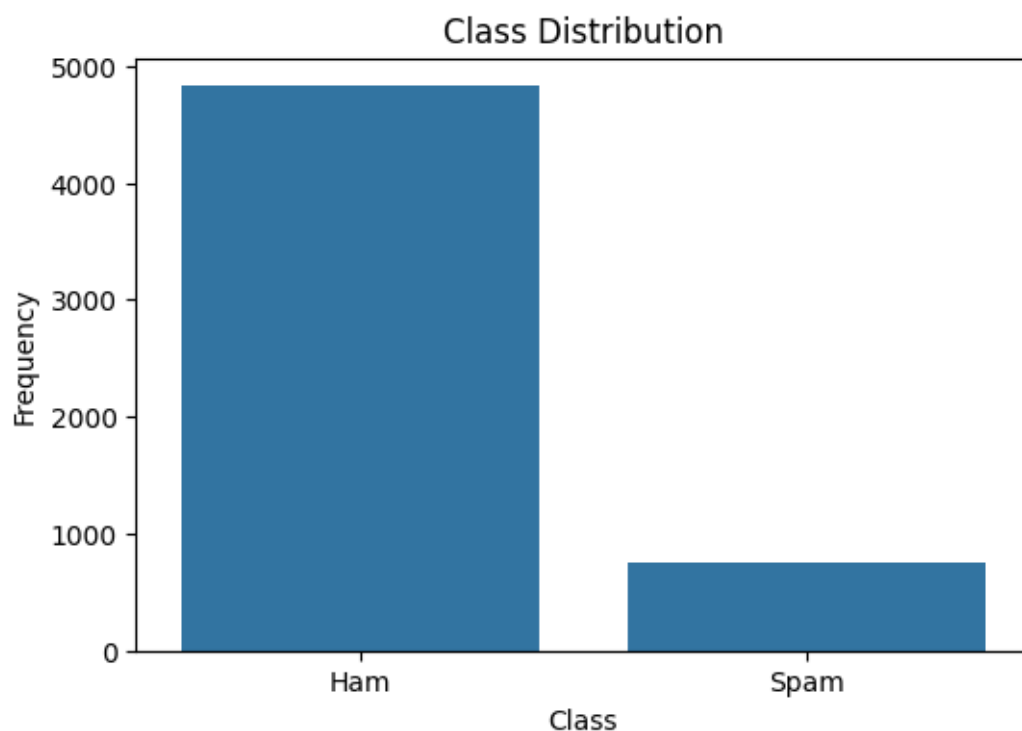There are **4,825 ham** (non-spam) messages (target labeled as '0').

There are **747 spam** messages (target labeled as '1').

This **indicates an imbalance** in the dataset, with more ham messages than spam messages.

```
[13]   1   # Analyzing Class Distribution
       2   class_distribution = spam_df['target'].value_counts()
       3   print(class_distribution)
       4

       0    4825
       1     747
       Name: target, dtype: int64
```

```
1   # Visualization of Class Distribution
2   plt.figure(figsize=(6, 4))
3   sns.barplot(x=class_distribution.index, y=class_distribution.values)
4   plt.title('Class Distribution')
5   plt.xlabel('Class')
6   plt.ylabel('Frequency')
7   plt.xticks(range(len(class_distribution.index)), ['Ham', 'Spam'])  |
8   plt.show()
9
```

**Q1. Choose the best Bayes classification model and the best logistic regression model from the previous assignment, and train and evaluate the models.**

Best Bayes Classification Model from previous assignment: Support Vector Machine

Best Logistic Regression Model from previous assignment: Logistic regression with L2 regularization

**Best Bayes Classification Model**

**Post-processing:** TF-IDF Vectorization

Importing Libraries: importing necessary libraries, including TfidfVectorizer from sklearn.feature_extraction.text

**Initializing the Vectorizer:** A TfidfVectorizer object is created with default parameters. This vectorizer is responsible for converting a collection of raw documents to a matrix of TF-IDF features.

**Fitting and Transforming the Data:** The vectorizer is then fitted to the sentence column of the spam_df, which contains the text data. The fitting process involves learning the vocabulary of the corpus and transforming the text data into a sparse matrix of TF-IDF features. The result is assigned to X.

**Creating the Target Array:** 'y' is taken directly from the target column of the spam_df, which contains the labels for the classification (0 for ham, 1 for spam).

**Splitting the Dataset:** The feature matrix X and the target array y are split into training and test sets using train_test_split, with 20% of the data being reserved for testing (test_size=0.2). The random state is set to 42 for reproducibility.

```
1   from sklearn.feature_extraction.text import TfidfVectorizer
2
3   # Initializing the TF-IDF vectorizer
4   tfidf_vectorizer = TfidfVectorizer()
5
6   # Fitting and transforming the dataset
7   X = tfidf_vectorizer.fit_transform(spam_df['sentence'])
8   y = spam_df['target']
9
10  # Splitting the dataset into training and testing sets
11  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
12
```

**Modelling**

An instance of SVC is created with the **kernel='linear'** argument, which is suitable for text classification tasks. The svm_model is then trained using the .fit() method with X_train and

y_train as arguments. X_train contains the feature vectors for the training data, and y_train contains the corresponding labels.

```python
from sklearn.svm import SVC

# Initializing the SVM model
svm_model = SVC(kernel='linear')  # Linear kernel is often used for text classification

# Training the model
svm_model.fit(X_train, y_train)

# Making predictions
predictions = svm_model.predict(X_test)
```

**Model Performance**

The model achieved an accuracy of approximately 97.93% on the test set, which is quite high and indicates that the model was able to correctly classify the majority of the messages as spam or ham.

Precision: 0.98 indicates that 98% of the instances predicted as ham were actually ham.

Recall: 1.00 demonstrates that the model identified all actual ham messages correctly.

F1-Score: 0.99 suggests a very good balance between precision and recall for ham messages.

Confusion Matrix:

- Ham Predictions:
  - True Positives (TP): 963 ham messages were correctly classified.
  - False Negatives (FN): 2 ham messages were incorrectly classified as spam.
- Spam Predictions:
  - True Negatives (TN): 129 spam messages were correctly classified.
  - False Positives (FP): 21 spam messages were incorrectly classified as ham.

```
Accuracy: 0.979372197309417
              precision    recall  f1-score   support

           0       0.98      1.00      0.99       965
           1       0.98      0.86      0.92       150

    accuracy                           0.98      1115
   macro avg       0.98      0.93      0.95      1115
weighted avg       0.98      0.98      0.98      1115
```

## Confusion Matrix - SVM Model



**Best Logistic Regression Model**

Importing the Model: LogisticRegression class is imported from sklearn.linear_model

Initializing the Model: An instance of LogisticRegression is created with penalty='l2' and solver='liblinear'.

The L2 penalty refers to using L2 regularization, which helps prevent overfitting by penalizing large coefficients. The solver 'liblinear' is a good choice for small datasets and binary classification problems.

Making Predictions: Similar to the SVM, the trained logistic regression model is used to make predictions on the test data (X_test) with the .predict() method. These predictions are stored in y_pred.

```
1   from sklearn.linear_model import LogisticRegression
2   l2_model = LogisticRegression(penalty='l2', solver='liblinear', random_state=42)
3
4   l2_model.fit(X_train, y_train)
5
6   y_pred = l2_model.predict(X_test)
7
8   # Evaluate the model
9   accuracy = accuracy_score(y_test, y_pred)
0   conf_matrix = confusion_matrix(y_test, y_pred)
1   classification_rep = classification_report(y_test, y_pred)
2
3   # Print the results
4   print(f"Accuracy: {accuracy}")
5   print(f"Confusion Matrix:\n{conf_matrix}")
6   print(f"Classification Report:\n{classification_rep}")
```

**Model Performance**

The model has an accuracy of approximately 96.32%, which means it correctly classified 96.32% of the messages as spam or ham.

Precision: 1.00 implies that every instance the model predicted as spam was indeed spam. There were no false positives in the spam predictions.

Recall: 0.73 indicates that the model identified 73% of all actual spam messages.

F1-Score: 0.84 shows a good balance between precision and recall for spam messages but indicates there is room for improvement, especially in recall.

Confusion Matrix

- Ham Predictions:
    - True Positives (TP): 965 ham messages were correctly classified.
    - False Negatives (FN): 0 ham messages were incorrectly classified as spam, which means all ham messages were correctly identified.

- Spam Predictions:
    - True Negatives (TN): 41 spam messages were correctly classified.
    - False Positives (FP): 0 spam messages were incorrectly classified as ham, which means no ham messages were incorrectly labeled as spam.

- Spam Messages Missed:
    - 109 spam messages were not identified by the model, as indicated by the recall of 0.73 for class 1. This means that while the model is very reliable when it does identify a message as spam (precision of 1.00), it fails to catch all spam messages, resulting in missed spam.

```
Accuracy: 0.9632286995515695
Confusion Matrix:
[[965   0]
 [ 41 109]]
Classification Report:
              precision    recall  f1-score   support

           0       0.96      1.00      0.98       965
           1       1.00      0.73      0.84       150

    accuracy                           0.96      1115
   macro avg       0.98      0.86      0.91      1115
weighted avg       0.96      0.96      0.96      1115
```

When comparing both models, the SVM model had a slightly higher accuracy and balanced performance across both classes, while the Logistic Regression model had perfect precision but lower recall for spam messages, indicating some spam messages were misclassified as ham.

## Design 1: LSTM (Long Short-Term Memory

**Additional Pre-processing**

- Tokenization and Padding: The Tokenizer from Keras is used to convert text data into sequences of integers.
  - The texts_to_sequences method is applied to both training and testing sets to transform the preprocessed text.
  - Sequences are padded to a maximum length of 100 characters, ensuring uniform input size for the neural network.
- Target Variable Encoding: The target variable ('target') in both training and testing datasets is encoded using LabelEncoder from scikit-learn, converting categorical labels into a numeric format.

```
17
18   # Tokenize and pad sequences
19   max_len = 100  # Adjust as needed
20   tokenizer = Tokenizer()
21   tokenizer.fit_on_texts(train_df['processed_text'])
22   X_train = tokenizer.texts_to_sequences(train_df['processed_text'])
23   X_test = tokenizer.texts_to_sequences(test_df['processed_text'])
24
25   X_train_padded = pad_sequences(X_train, maxlen=max_len, padding='post')
26   X_test_padded = pad_sequences(X_test, maxlen=max_len, padding='post')
27
28   # Encode target variable
29   label_encoder = LabelEncoder()
30   train_df['target'] = label_encoder.fit_transform(train_df['target'])
31   test_df['target'] = label_encoder.transform(test_df['target'])
32
```

**Modelling**

A Sequential model is constructed with the following layers:

- An **Embedding** layer to learn word embeddings (vector representations of words), with the dimension of 100.

- A SpatialDropout1D layer with a **dropout rate of 0.2** to prevent overfitting by dropping entire 1D feature maps.

- An **LSTM layer with 64 units**, dropout and recurrent dropout rates of 0.2, and 'tanh' activation function to capture long-term dependencies.

- A **Dense output layer with a 'sigmoid'** activation function for binary classification.

The model is compiled with the **'adam'** optimizer and **'binary_crossentropy'** loss function, aiming to optimize accuracy.

The model is trained on the padded training data for **10 epochs** with a **batch size of 64**, using the padded testing data as validation.

```
# Model
embedding_dim = 100
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=embedding_dim, input_length=max_len))
model.add(SpatialDropout1D(0.2))
model.add(LSTM(64, dropout=0.2, recurrent_dropout=0.2, activation='tanh'))
model.add(Dense(1, activation='sigmoid'))

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

# Train the model
history = model.fit(X_train_padded, train_df['target'], epochs=10, batch_size=64, validation_data=(X_test_padded, test_df['target']))

# Plot training history
plt.plot(history.history['accuracy'], label='training')
plt.plot(history.history['val_accuracy'], label='testing')
plt.title('Model Accuracy')
plt.ylabel('Accuracy')
plt.xlabel('Epoch')
plt.legend(['Train', 'Test'], loc='upper left')
plt.show()

# Evaluate on the test set
test_loss, test_acc = model.evaluate(X_test_padded, test_df['target'])
print(f"Test Loss: {test_loss:.4f}, Test Accuracy: {test_acc:.4f}")
```

**Model Performance**

The accuracy plot shows that the LSTM model's training and validation accuracy are quite high and stable from the first epoch. Both the training and validation accuracies are almost identical and constant across all epochs, which suggests that the model could be underfitting. The model has been trained for 10 epochs with a batch size of 64.

- The training and validation losses are around 0.395, and the accuracy is constant at about 0.8661 or 86.61%.

- After training, the model's performance on the test set shows a loss of 0.3961 and accuracy of 86.55%.

**Model Architecture**

- Embedding layer with input dimensions set to the size of the vocabulary (593000) and output dimensions to 100.

- SpatialDropout1D layer with a dropout rate of 0.2.

- LSTM layer with 64 units and dropout and recurrent dropout rates of 0.2, using 'tanh' activation.

- Dense layer with 1 unit and 'sigmoid' activation function for binary classification.

- Total parameters of the model are 635,305, all of which are trainable.

```
Model: "sequential_1"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 100, 100)          593000

 spatial_dropout1d_1 (Spati  (None, 100, 100)          0
 alDropout1D)

 lstm_1 (LSTM)               (None, 64)                42240

 dense_1 (Dense)             (None, 1)                 65


=================================================================
Total params: 635305 (2.42 MB)
Trainable params: 635305 (2.42 MB)
Non-trainable params: 0 (0.00 Byte)
```

## Design 2: Bidirectional LSTM

**Additional Pre-processing**

Tokenization: The Tokenizer from Keras was used with a maximum vocabulary size of 5,000
words. Texts were split into tokens and then sequences were generated from these tokens.

Padding: Sequences were padded to ensure a consistent input size of 100 tokens for each
sample.

Label Encoding: The target variable (spam or not spam) was label-encoded to provide a
binary numerical target for model training.

```
14    # Define parameters
15    max_features = 5000
16    maxlen = 100
17    embedding_dim = 100
18
19    # Tokenization
20    tokenizer = Tokenizer(num_words=max_features, split=' ')
21    tokenizer.fit_on_texts(spam_df['processed_text'].values)
22    X = tokenizer.texts_to_sequences(spam_df['processed_text'].values)
23    X = pad_sequences(X, maxlen=maxlen)
24
25    # Label Encoding
26    label_encoder = LabelEncoder()
27    y = label_encoder.fit_transform(spam_df['target'])
```

**Modelling**

The model is a Sequential model from Keras with the following layers:

- Embedding Layer: Converts tokenized word sequences into dense vectors of fixed size (100 dimensions in this case), with an input vocabulary of 5,000 words.

- Spatial Dropout Layer: Adds dropout regularization to the embedding layer, dropping entire 1D feature maps to control overfitting.

- Bidirectional LSTM Layer: Processes the sequences in both directions with 64 units and applies dropout and recurrent dropout to further regularize the model.

- Dense Output Layer: A single unit with a sigmoid activation function to output a probability indicating the likelihood of the input being spam.

**Training Process:**

The model was compiled with a binary crossentropy loss function and the Adam optimizer, focusing on accuracy as the performance metric. Training was conducted over 10 epochs with a batch size of 64, using the training data and validating on a holdout set (20% of the dataset).

```python
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Build the model
model = Sequential()
model.add(Embedding(max_features, embedding_dim, input_length=X.shape[1]))
model.add(SpatialDropout1D(0.2))
model.add(Bidirectional(LSTM(64, dropout=0.2, recurrent_dropout=0.2)))
model.add(Dense(1, activation='sigmoid'))

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

# Train the model
history = model.fit(X_train, y_train, epochs=10, batch_size=64, validation_data=(X_test, y_test), verbose=2)

# Plot training history
plt.plot(history.history['accuracy'], label='training')
plt.plot(history.history['val_accuracy'], label='testing')
plt.title('Model Accuracy')
plt.ylabel('Accuracy')
plt.xlabel('Epoch')
plt.legend(['Train', 'Test'], loc='upper left')
plt.show()

# Evaluate on the test set
test_loss, test_acc = model.evaluate(X, y)
print(f"Test Loss: {test_loss:.4f}, Test Accuracy: {test_acc:.4f}")
```
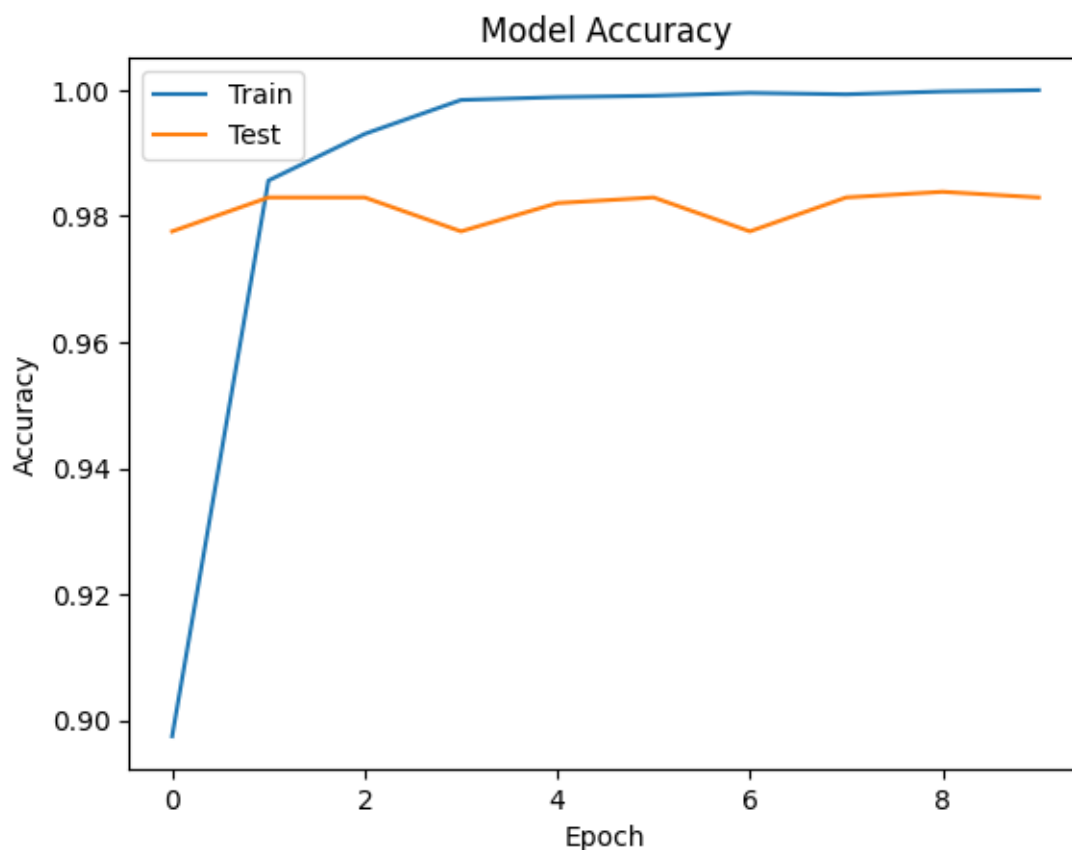
**Model Performance**

- The final evaluation on the entire dataset (including both training and test sets) shows a test accuracy of 99.66% and a loss of 0.0168. These are excellent results and indicate that the model performs exceptionally well.

- The slight discrepancy between the high training accuracy and the test accuracy suggests the model may have overfitted slightly, but the effect is minimal given the high test accuracy.

- The graph shows a sharp increase in training accuracy after the first epoch, indicating that the majority of learning happens early in the training process.

- The training accuracy started at 89.75% and increased to 100% by the end of the 10th epoch. This indicates that the model has learned to perfectly fit the training data.

- The validation accuracy started at 97.76% and fluctuated slightly, ending at 98.30%. This high accuracy suggests that the model generalizes well to unseen data.



**Model Architecture**

- The model has 584,609 trainable parameters

- The maximum number of features to 5,000 and the maximum sequence length to 100, which means the model considers the top 5,000 most frequent words in the sequences up to a length of 100 words.

- The embedding dimension is 100, which means that each token is represented as a 100-dimensional vector.

- Embedding Layer: Converts tokenized word sequences into dense vectors of fixed size (100 dimensions in this case), with an input vocabulary of 5,000 words.

- Spatial Dropout Layer: Adds dropout regularization to the embedding layer, dropping entire 1D feature maps to control overfitting.

- Bidirectional LSTM Layer: Processes the sequences in both directions with 64 units and applies dropout and recurrent dropout to further regularize the model.

- Dense Output Layer: A single unit with a sigmoid activation function to output a probability indicating the likelihood of the input being spam.

```
Model: "sequential_2"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_2 (Embedding)     (None, 100, 100)          500000

 spatial_dropout1d_2 (Spati  (None, 100, 100)          0
 alDropout1D)

 bidirectional_2 (Bidirecti  (None, 128)               84480
 onal)

 dense_2 (Dense)             (None, 1)                 129

=================================================================
Total params: 584609 (2.23 MB)
Trainable params: 584609 (2.23 MB)
Non-trainable params: 0 (0.00 Byte)
```

**Q2. Conduct experiments to compare different configurations of RNN neural networks for text classification.**

    a. **Design 1 vs. Design 2. For example, Bidirectional LSTM vs. LSTM, or LSTM vs. RNN.**

    b. **Compare two configurations of Design 1: More parameters vs. fewer parameters**

    c. **Compare two configurations of Design 2: More parameters vs. fewer parameters**

**d. Compare two configurations of Design 1: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

**e. Compare two configurations of Design 2: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

## Q2a. Design 1: LSTM and Design 2: Bidirectional LSTM

| Feature | Design 1: LSTM | Design 2: Bidirectional LSTM |
|---|---|---|
| **Direction of Processing** | Unidirectional (forward) | Bidirectional (forward and backward) |
| **LSTM Units** | 64 units | 64 units |
| **Regularization** | Dropout | Dropout |
| **Output Layer** | Dense layer with sigmoid activation | Dense layer with sigmoid activation |
| **Performance** | Lower (0.8655 accuracy) | Higher (0.9966 accuracy) |
| **Complexity** | Lower | Higher |
| **Risk of Overfitting** | Lower | Higher, but good generalization in examples |
| **Resource Requirements** | Lower | Higher |
| **Training Time** | Shorter | Longer |
| **Generalization Performance** | Good | Superior |
| **Suitability for Sequential Data** | Suitable for simpler sequences | Better for complex sequences requiring understanding of context |

**Design 1: LSTM**

- In Design 1, a single LSTM layer is used.
- It processes the input sequences sequentially from start to end.
- The LSTM layer has 64 units and uses dropout for regularization.

- It's followed by a dense layer with a sigmoid activation function for binary classification.

**Design 2: Bidirectional LSTM**

- In Design 2, a bidirectional LSTM layer is used.
- Bidirectional LSTMs process the input sequences in both forward and backward directions, allowing the model to capture information from both past and future states.
- It also has 64 units and uses dropout for regularization.
- Like Design 1, it's followed by a dense layer with a sigmoid activation function for binary classification.

**Comparison**:

**Performance**: Design 2 (Bidirectional LSTM) generally performs better. It tends to capture more contextual information due to its bidirectional nature, which can lead to improved accuracy.

**Complexity**: Design 2 is slightly more complex due to bidirectional processing, which may result in longer training times and a higher number of parameters.

**Overfitting**: Bidirectional LSTMs may be more prone to overfitting since they capture information from both directions, potentially leading to a model that generalizes less well to unseen data. However, in the provided examples, both models show good generalization performance.

**Resource Requirements**: Bidirectional LSTMs typically require more computational resources during training and inference compared to unidirectional LSTMs due to their increased complexity.

**Conclusion**:

In conclusion, the Bidirectional LSTM model outperforms the unidirectional LSTM model in terms of accuracy, achieving near-perfect performance with a test accuracy of 0.9966 compared to 0.8655. Despite its longer training time, the Bidirectional LSTM demonstrates superior generalization, showing consistent improvement without signs of overfitting. Therefore, for tasks requiring sophisticated understanding of sequential data like text classification, the Bidirectional LSTM proves to be a more effective choice, offering superior performance and robustness.

## Q2b. Compare two configurations of Design 1: More parameters vs. fewer parameters

Design 1: LSTM with More parameters

- A sequential model is constructed with an embedding layer, which will create dense vector representations for the words in the vocabulary. The dimensionality of the **embedding space is set to 150**.

- A spatial dropout layer is added to introduce regularization by dropping entire 1D feature maps in the embedding layer, reducing overfitting.

- An **LSTM layer with 128 units** is used to process the sequence data, with dropout and recurrent dropout for regularization.

- **Two dense layers follow**, with the second dense layer using a sigmoid activation function for binary classification.

- The training runs for 25 epochs with a batch size of 128

```
# Model
embedding_dim = 150  # Increased embedding dimension
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=embedding_dim, input_length=max_len))
model.add(SpatialDropout1D(0.3))
model.add(LSTM(128, dropout=0.3, recurrent_dropout=0.3, activation='tanh'))
model.add(Dense(64, activation='relu'))  # Additional Dense layer
model.add(Dense(1, activation='sigmoid'))
```

**Model Performance:**

From the accuracy plot, it's immediately clear that the **model's accuracy** on both the training and test sets is **very stable** across the epochs.

The **training and test accuracy** are both hovering around **86.55%,** which suggests that the model is **not overfitting**, as both accuracies are quite close.

The model summary shows that the model is relatively simple with a **total of 1,040,669 trainable parameters.**

```
Model: "sequential_1"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 100, 150)          889500

 spatial_dropout1d_1 (Spati  (None, 100, 150)          0
 alDropout1D)

 lstm_1 (LSTM)               (None, 128)               142848

 dense_2 (Dense)             (None, 64)                8256

 dense_3 (Dense)             (None, 1)                 65

=================================================================
Total params: 1040669 (3.97 MB)
Trainable params: 1040669 (3.97 MB)
Non-trainable params: 0 (0.00 Byte)
```



**Design 1: LSTM with Fewer parameters**

- Max Sequence Length: The max_len has been **decreased to 25**. This will reduce the input size to the model, which will result in faster training times but might impact the model's ability to capture longer dependencies in the text.

- Embedding Dimension: The **embedding_dim has been reduced to 25**. This change reduces the dimensionality of the word vectors, which decreases the number of parameters in the embedding layer.

- **LSTM Units**: The number of units in the LSTM layer has been **reduced to 16**, significantly cutting down the number of parameters that the model needs to learn.

- **Regularization**: Dropout rates have been reduced **to 0.1** for both the spatial dropout and the recurrent dropout in the LSTM.

- Training Configuration: The model is now trained for **fewer epochs (5 instead of 25)** with a **smaller batch size** (16 instead of 128). This could lead to less stable gradient estimates during training, which sometimes can help to find new and possibly better local minima in the loss landscape, but it also may result in a less robust model due to less thorough exploration of the weight space.

```python
# Model
embedding_dim = 25  # Reduced embedding dimension
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=embedding_dim, input_length=max_len))
model.add(SpatialDropout1D(0.1))
model.add(LSTM(16, dropout=0.1, recurrent_dropout=0.1, activation='tanh'))  # Reduced LSTM units
model.add(Dense(1, activation='sigmoid'))

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

**Model Performance**

- The accuracy plot for the this LSTM model exhibits a significant improvement compared to the previous model.

- The training accuracy quickly rises to near-perfection within just a few epochs, and the validation accuracy also shows an impressive increase, achieving a high score of around 98%.

- This performance indicates that the model is learning effectively from the training data, and more importantly, it is generalizing well to the validation data. The slight divergence between the training and validation accuracy suggests that the model may be starting to overfit;

- The final test accuracy of 97.94% and suggests that the model is performing well on unseen data. This is a promising result, especially considering the model's reduced

complexity and the smaller number of parameters, which stands at around 150,955 as per the model summary.

```
Model: "sequential_3"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_3 (Embedding)     (None, 25, 25)            148250

 spatial_dropout1d_3 (Spati  (None, 25, 25)            0
 alDropout1D)

 lstm_3 (LSTM)               (None, 16)                2688

 dense_5 (Dense)             (None, 1)                 17


=================================================================
Total params: 150955 (589.67 KB)
Trainable params: 150955 (589.67 KB)
Non-trainable params: 0 (0.00 Byte)
```

**Comparing Design 1: LSTM with More parameters vs. LSTM with fewer parameters**

| Feature | LSTM with More Parameters | LSTM with Fewer Parameters |
|---|---|---|
| Embedding Dimension | 150 | 25 |
| LSTM Units | 128 | 16 |
| Additional Dense Layer | Yes (64 neurons) | No |
| Dropout Rate | 0.3 | 0.1 |
| Batch Size | 128 | 16 |
| Epochs | 25 | 5 |
| Test Accuracy | 86.55% | 97.94% |
| Test Loss | 0.3958 | 0.0946 |
| Drawbacks | Higher complexity, longer training, increased overfitting risk, higher computational requirements | Potentially underfits, limited complexity may miss nuanced patterns |
| Learning Graphs | Potential overfitting indicated by diverging training and validation accuracy | Faster convergence, possible underfitting if accuracies are low |
| Training Time | Longer due to higher complexity | Shorter due to lower complexity |
| Overfitting Patterns | Possible if training accuracy increases while | Less prone due to simplicity, but may not capture complex patterns |

| | validation accuracy plateaus or decreases | |
|---|---|---|

**Q2c.Compare two configurations of Design 2: More parameters vs. fewer parameters**

**Design 2: BiLSTM with More parameters**

- The model is set to handle a vocabulary size of 5,000 (max_features).

- The maximum input sequence length is defined as 100 (maxlen).

- The dimension of the embedding space is set to 200 (embedding_dim).

- The core of the model is a BiLSTM layer with 128 units

- The model is trained on the training data for 25 epochs with a batch size of 128.

```
# Build the model
model = Sequential()
model.add(Embedding(max_features, embedding_dim, input_length=X.shape[1]))
model.add(SpatialDropout1D(0.2))
model.add(Bidirectional(LSTM(128, dropout=0.2, recurrent_dropout=0.2)))
model.add(Dense(64, activation='relu'))
model.add(Dense(1, activation='sigmoid'))

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

**Model Performance**

- The plot shows that the training accuracy quickly reaches near 100%, indicating a good fit on the training data. However, the test accuracy increases more modestly and fluctuates around the 97%-98% range, which suggests that the model is generalizing well but might be starting to overfit as the training accuracy continues to stay at perfect levels while the test accuracy does not improve.

- The final model performance on the test set is reported as approximately 99.57% accuracy, which is quite high.

```
Model: "sequential_1"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 100, 200)          1000000

 spatial_dropout1d_1 (Spati  (None, 100, 200)          0
 alDropout1D)

 bidirectional (Bidirection  (None, 256)               336896
 al)

 dense (Dense)               (None, 64)                16448

 dense_1 (Dense)             (None, 1)                 65


=================================================================
Total params: 1353409 (5.16 MB)
Trainable params: 1353409 (5.16 MB)
Non-trainable params: 0 (0.00 Byte)
```

```
None
Epoch 1/25
35/35 - 98s - loss: 0.3330 - accuracy: 0.8898 - val_loss: 0.1317 - val_accuracy: 0.9668 - 98s/epoch - 3s/step
Epoch 2/25
35/35 - 64s - loss: 0.0637 - accuracy: 0.9821 - val_loss: 0.0692 - val_accuracy: 0.9794 - 64s/epoch - 2s/step
Epoch 3/25
35/35 - 62s - loss: 0.0267 - accuracy: 0.9930 - val_loss: 0.0705 - val_accuracy: 0.9812 - 62s/epoch - 2s/step
Epoch 4/25
35/35 - 64s - loss: 0.0139 - accuracy: 0.9964 - val_loss: 0.0888 - val_accuracy: 0.9821 - 64s/epoch - 2s/step
Epoch 5/25
35/35 - 62s - loss: 0.0069 - accuracy: 0.9984 - val_loss: 0.0959 - val_accuracy: 0.9830 - 62s/epoch - 2s/step
Epoch 6/25
35/35 - 62s - loss: 0.0038 - accuracy: 0.9991 - val_loss: 0.0976 - val_accuracy: 0.9785 - 62s/epoch - 2s/step
Epoch 7/25
35/35 - 62s - loss: 0.0028 - accuracy: 0.9991 - val_loss: 0.1162 - val_accuracy: 0.9803 - 62s/epoch - 2s/step
Epoch 8/25
35/35 - 63s - loss: 0.0017 - accuracy: 0.9998 - val_loss: 0.1131 - val_accuracy: 0.9794 - 63s/epoch - 2s/step
Epoch 9/25
35/35 - 62s - loss: 0.0028 - accuracy: 0.9993 - val_loss: 0.1359 - val_accuracy: 0.9803 - 62s/epoch - 2s/step
Epoch 10/25
35/35 - 62s - loss: 0.0020 - accuracy: 0.9993 - val_loss: 0.1290 - val_accuracy: 0.9803 - 62s/epoch - 2s/step
Epoch 11/25
35/35 - 60s - loss: 0.0014 - accuracy: 0.9996 - val_loss: 0.1366 - val_accuracy: 0.9794 - 60s/epoch - 2s/step
Epoch 12/25
35/35 - 63s - loss: 0.0012 - accuracy: 0.9996 - val_loss: 0.1258 - val_accuracy: 0.9785 - 63s/epoch - 2s/step
Epoch 13/25
35/35 - 63s - loss: 0.0011 - accuracy: 0.9998 - val_loss: 0.1270 - val_accuracy: 0.9794 - 63s/epoch - 2s/step
Epoch 14/25
35/35 - 62s - loss: 0.0013 - accuracy: 0.9998 - val_loss: 0.1289 - val_accuracy: 0.9785 - 62s/epoch - 2s/step
Epoch 15/25
35/35 - 62s - loss: 6.0075e-04 - accuracy: 1.0000 - val_loss: 0.1294 - val_accuracy: 0.9776 - 62s/epoch - 2s/step
```

Model Accuracy

**Design 2: BiLSTM with More parameters**

- The embedding dimension is reduced from 200 to 25 and the LSTM unit size from 128 to 16

- The training history is plotted over 5 epochs instead of 25

- The smaller batch size of 16 (compared to 128 previously), the model may take more time to complete each epoch but could also benefit from more frequent updates to the weights.

**Model Performance**

The plot initially started with Training accuracy started at 90.35%, and validation accuracy was at 97.40% and ended in 5 epochs with Training accuracy hit 99.82%, and validation accuracy went back up to 98.03%.
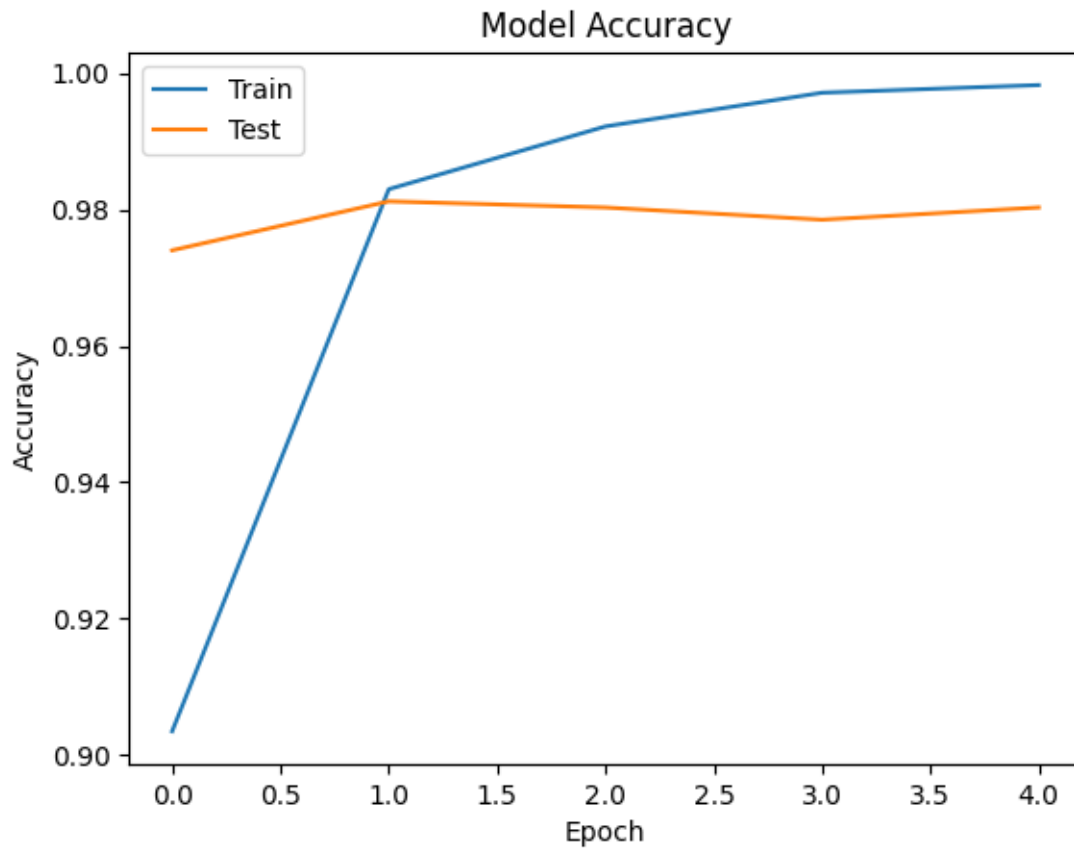
The model does not show signs of overfitting within the five epochs since the validation accuracy remains high and does not diverge from the training accuracy.

```
Model: "sequential_3"

_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_3 (Embedding)     (None, 100, 25)           125000

 spatial_dropout1d_3 (Spati  (None, 100, 25)           0
 alDropout1D)

 bidirectional_1 (Bidirecti  (None, 32)                5376
 onal)

 dense_2 (Dense)             (None, 1)                 33


=================================================================
Total params: 130409 (509.41 KB)
Trainable params: 130409 (509.41 KB)
Non-trainable params: 0 (0.00 Byte)
```

```
None
Epoch 1/5
279/279 - 55s - loss: 0.2682 - accuracy: 0.9035 - val_loss: 0.1069 - val_accuracy: 0.9740 - 55s/epoch - 198ms/step
Epoch 2/5
279/279 - 44s - loss: 0.0693 - accuracy: 0.9829 - val_loss: 0.0667 - val_accuracy: 0.9812 - 44s/epoch - 159ms/step
Epoch 3/5
279/279 - 44s - loss: 0.0332 - accuracy: 0.9921 - val_loss: 0.0680 - val_accuracy: 0.9803 - 44s/epoch - 157ms/step
Epoch 4/5
279/279 - 47s - loss: 0.0162 - accuracy: 0.9971 - val_loss: 0.0719 - val_accuracy: 0.9785 - 47s/epoch - 169ms/step
Epoch 5/5
279/279 - 45s - loss: 0.0103 - accuracy: 0.9982 - val_loss: 0.0802 - val_accuracy: 0.9803 - 45s/epoch - 160ms/step
```

## Model Accuracy



**Comparing Design 2: BiLSTM with More parameters vs. BiLSTM with fewer parameters**

| Feature/Aspect | BiLSTM with More Parameters | BiLSTM with Fewer Parameters |
|---|---|---|
| **Embedding Dimension** | 200 | 25 |
| **LSTM Units** | 128 | 16 |
| **Additional Dense Layer** | Yes (64 neurons) | No |
| **Dropout Rate** | 0.2 | 0.1 |
| **Batch Size** | 128 | 16 |
| **Epochs** | 25 | 5 |
| **Test Accuracy** | 99.57% | 99.53% |
| **Test Loss** | 0.0294 | 0.0214 |
| **Complexity** | Higher (more parameters) | Lower (fewer parameters) |

| Potential Drawbacks | - Longer training times | - May not capture complex patterns |
|---|---|---|
| | - Increased risk of overfitting | - May lead to underfitting |
| | - Higher computational resources | |
| **Learning Graphs Observations** | - Possible overfitting signs | - Convergence may occur faster |
| | - Training/validation accuracy may diverge | - Possible underfitting if accuracies are low |
| **Training Time** | Longer (~26.25 mins) | Shorter (~5.17 mins) |
| **Overfitting Patterns** | More prone due to complexity | Less prone due to simplicity |

**Q2d. Compare two configurations of Design 1: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

**Design 1: LSTM with Strategies for avoiding overfitting are applied**

- Embedding Layer: The input dimension is the size of the vocabulary plus one, and the output dimension is 100. The input length is fixed at 100 tokens.
- Spatial Dropout1D Layer: Applies dropout to entire 1D feature maps instead of individual elements, with a dropout rate of 0.2.
- LSTM Layer: A Long Short-Term Memory layer with 64 units. It includes dropout and recurrent dropout, both set to 0.2, to help prevent overfitting. This layer is also regularized with L1 regularization, which adds a penalty equal to the absolute value of the magnitude of coefficients, encouraging sparsity in the model weights.

- Dense Layer: A fully connected layer that outputs a single value with a sigmoid activation function, used to achieve a binary classification outcome (e.g., spam or not spam).

- The total number of parameters in this model is 596,405, all of which are trainable.

- Overfitting Strategies:

    o **L1 Regularization** (Lasso Regularization): The model uses L1 regularization on the LSTM layer, which adds a penalty equal to the absolute value of the magnitude of coefficients. This can lead to sparse models where some weights can become zero, essentially performing feature selection and helping reduce overfitting by simplifying the model.

    o **Early Stopping**: This is a form of regularization used to avoid overfitting by halting the training process if the model's performance on a validation set does not improve for a specified number of epochs. By monitoring validation loss and stopping the training early, the model is prevented from learning the noise in the training set too deeply.

```python
# Model
from keras import regularizers
from keras.callbacks import EarlyStopping

# Define the regularization parameter
l1_lambda = 0.01

# Define early stopping criteria
early_stopping = EarlyStopping(monitor='val_loss', patience=3)  # Adjust patience as needed
#Defining model
embedding_dim = 100
model_overfitting = Sequential()
model_overfitting.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=embedding_dim, input_length=max_len))
model_overfitting.add(SpatialDropout1D(0.2))
model_overfitting.add(LSTM(64, dropout=0.2, recurrent_dropout=0.2, activation='tanh', kernel_regularizer=regularizers.l1(l1_lambda)))
model_overfitting.add(Dense(1, activation='sigmoid'))

model_overfitting.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model_overfitting.summary()
```

**Model Performance**

- The model starts with a significantly high initial loss (9.9258) due to the L1 regularization penalty but an accuracy close to the baseline (86.58%).

- By the end of training, the loss decreases substantially to approximately 0.4239, showing the model's ability to adjust and learn despite the initial penalty.

- The validation loss also decreases over epochs, ending up around 0.4271, which indicates that the model's predictions are consistent with its performance on the training set.

- The training accuracy remains relatively stable throughout the training process, ending at approximately 86.61%.
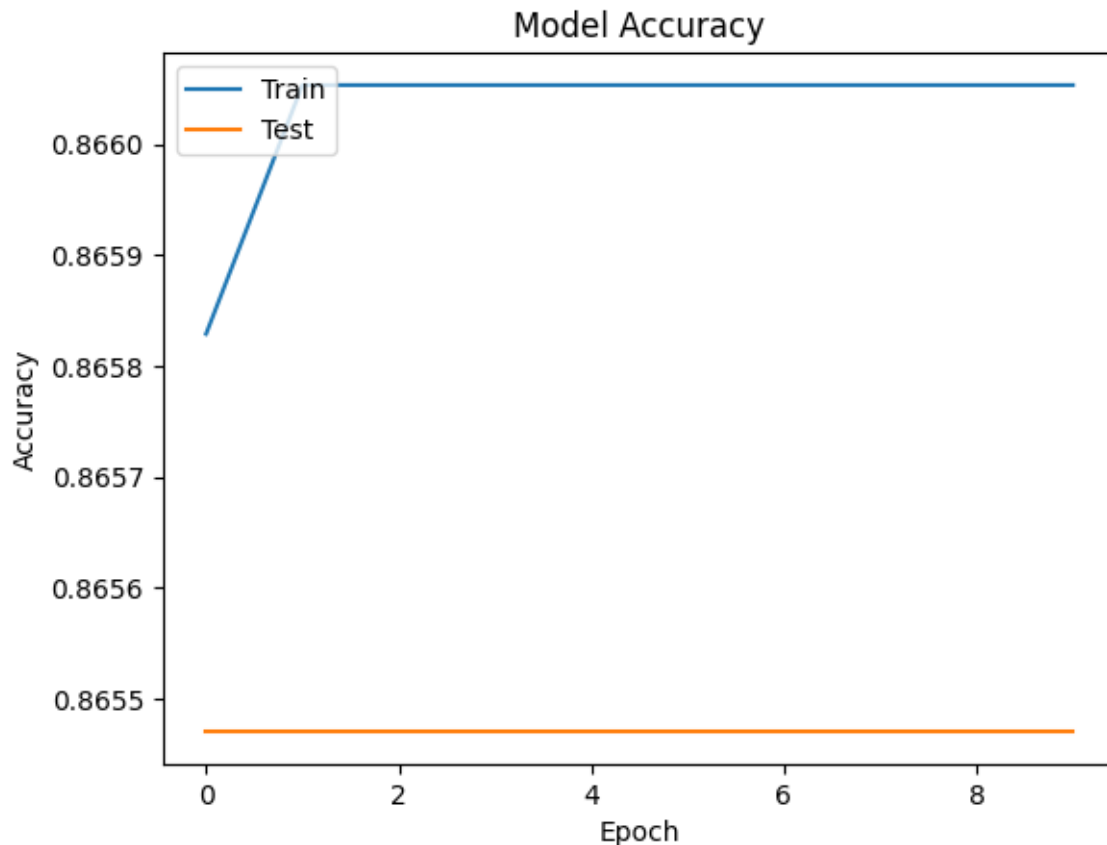
- Final Validation Accuracy is stable and matches the training accuracy closely, ending at approximately 86.55%.

```
Model: "sequential_2"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_2 (Embedding)     (None, 100, 100)          554100

 spatial_dropout1d_2 (Spati  (None, 100, 100)          0
 alDropout1D)

 lstm_2 (LSTM)               (None, 64)                42240

 dense_2 (Dense)             (None, 1)                 65

=================================================================
Total params: 596405 (2.28 MB)
Trainable params: 596405 (2.28 MB)
Non-trainable params: 0 (0.00 Byte)
_____
```

```
Epoch 1/10
70/70 [==============================] - 19s 202ms/step - loss: 9.9258 - accuracy: 0.8658 - val_loss: 4.0567 - val_accuracy: 0.8655
Epoch 2/10
70/70 [==============================] - 14s 198ms/step - loss: 1.5474 - accuracy: 0.8661 - val_loss: 0.4642 - val_accuracy: 0.8655
Epoch 3/10
70/70 [==============================] - 14s 200ms/step - loss: 0.4338 - accuracy: 0.8661 - val_loss: 0.4263 - val_accuracy: 0.8655
Epoch 4/10
70/70 [==============================] - 14s 202ms/step - loss: 0.4260 - accuracy: 0.8661 - val_loss: 0.4262 - val_accuracy: 0.8655
Epoch 5/10
70/70 [==============================] - 16s 232ms/step - loss: 0.4257 - accuracy: 0.8661 - val_loss: 0.4260 - val_accuracy: 0.8655
Epoch 6/10
70/70 [==============================] - 14s 200ms/step - loss: 0.4255 - accuracy: 0.8661 - val_loss: 0.4279 - val_accuracy: 0.8655
Epoch 7/10
70/70 [==============================] - 14s 200ms/step - loss: 0.4256 - accuracy: 0.8661 - val_loss: 0.4250 - val_accuracy: 0.8655
Epoch 8/10
70/70 [==============================] - 14s 198ms/step - loss: 0.4246 - accuracy: 0.8661 - val_loss: 0.4265 - val_accuracy: 0.8655
Epoch 9/10
70/70 [==============================] - 14s 202ms/step - loss: 0.4253 - accuracy: 0.8661 - val_loss: 0.4249 - val_accuracy: 0.8655
Epoch 10/10
70/70 [==============================] - 14s 201ms/step - loss: 0.4239 - accuracy: 0.8661 - val_loss: 0.4271 - val_accuracy: 0.8655
```

## Model Accuracy



**Design 1: LSTM with Strategies for avoiding overfitting are <u>not</u> applied**

Early stopping strategy not applied, everything else is same

```
# Model
embedding_dim = 100
model = Sequential()
model.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=embedding_dim, input_length=max_len))
model.add(SpatialDropout1D(0.2))
model.add(LSTM(64, dropout=0.2, recurrent_dropout=0.2, activation='tanh'))
model.add(Dense(1, activation='sigmoid'))

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

**Model Performance**

- The training and validation loss start at a lower point compared to the model with overfitting strategies and remain fairly stable throughout the training process, with the final losses around 0.395. This indicates a consistent performance of the model on both training and validation datasets.
- Accuracy: Both the training and validation accuracy start high and remain stable, with the final accuracy figures around 86.61% for training and 86.55% for validation.
- Initial and Final Epochs Performance:
  - Epoch 1: The model begins with a training loss of 0.4232 and an accuracy of 86.34%, with the validation loss at 0.3956 and accuracy at 86.55%.
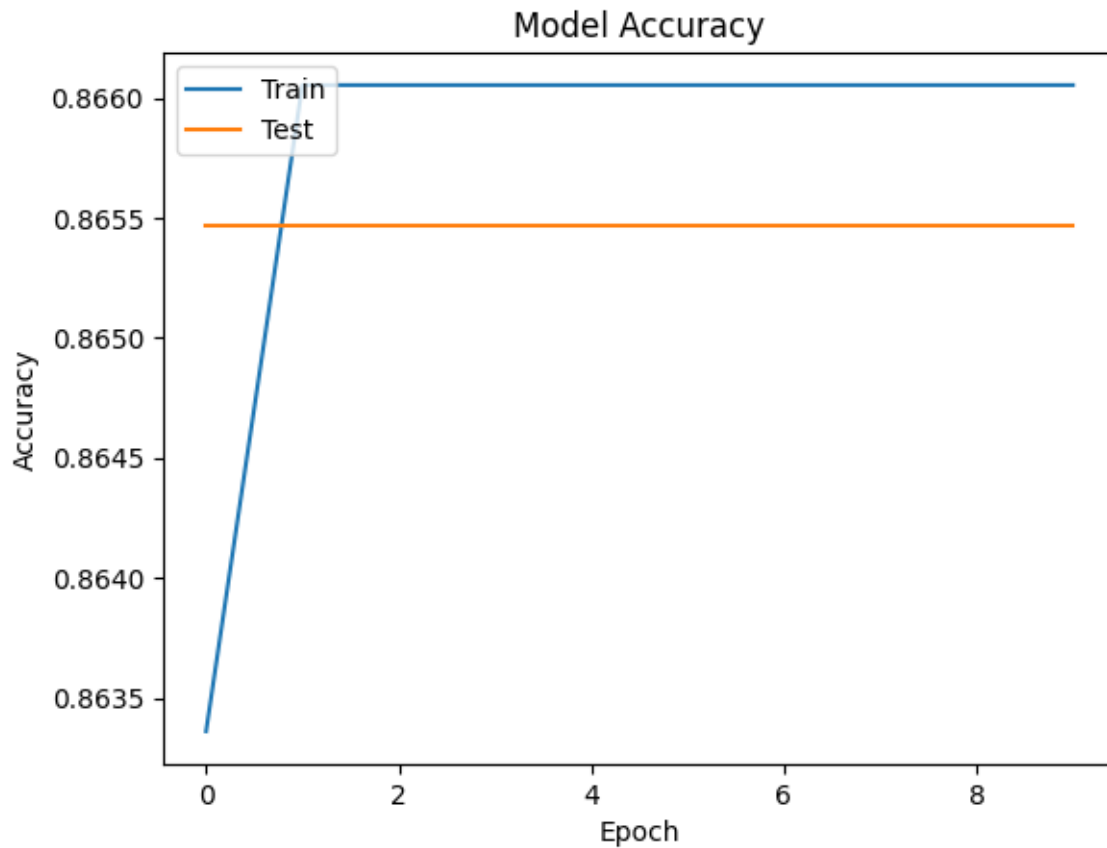
o Epoch 10: By the end of the training, the training loss slightly decreases to 0.3956, and the training accuracy remains relatively stable at 86.61%. The validation loss is 0.3958, with the validation accuracy consistent at 86.55%.

```
Model: "sequential"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding (Embedding)       (None, 100, 100)          554100

 spatial_dropout1d (Spatial  (None, 100, 100)          0
 Dropout1D)

 lstm (LSTM)                 (None, 64)                42240

 dense (Dense)               (None, 1)                 65

=================================================================
Total params: 596405 (2.28 MB)
Trainable params: 596405 (2.28 MB)
Non-trainable params: 0 (0.00 Byte)
_____
```

```
Epoch 1/10
70/70 [==============================] - 18s 206ms/step - loss: 0.4232 - accuracy: 0.8634 - val_loss: 0.3956 - val_accuracy: 0.8655
Epoch 2/10
70/70 [==============================] - 14s 195ms/step - loss: 0.3970 - accuracy: 0.8661 - val_loss: 0.3972 - val_accuracy: 0.8655
Epoch 3/10
70/70 [==============================] - 14s 194ms/step - loss: 0.3947 - accuracy: 0.8661 - val_loss: 0.3959 - val_accuracy: 0.8655
Epoch 4/10
70/70 [==============================] - 14s 196ms/step - loss: 0.3958 - accuracy: 0.8661 - val_loss: 0.3949 - val_accuracy: 0.8655
Epoch 5/10
70/70 [==============================] - 14s 194ms/step - loss: 0.3951 - accuracy: 0.8661 - val_loss: 0.3949 - val_accuracy: 0.8655
Epoch 6/10
70/70 [==============================] - 16s 222ms/step - loss: 0.3950 - accuracy: 0.8661 - val_loss: 0.3950 - val_accuracy: 0.8655
Epoch 7/10
70/70 [==============================] - 15s 222ms/step - loss: 0.3943 - accuracy: 0.8661 - val_loss: 0.3952 - val_accuracy: 0.8655
Epoch 8/10
70/70 [==============================] - 14s 195ms/step - loss: 0.3952 - accuracy: 0.8661 - val_loss: 0.4017 - val_accuracy: 0.8655
Epoch 9/10
70/70 [==============================] - 14s 200ms/step - loss: 0.3943 - accuracy: 0.8661 - val_loss: 0.3949 - val_accuracy: 0.8655
Epoch 10/10
70/70 [==============================] - 14s 200ms/step - loss: 0.3956 - accuracy: 0.8661 - val_loss: 0.3958 - val_accuracy: 0.8655
```

**Comparing Design 1: LSTM with Strategies for avoiding overfitting are applied vs LSTM with Strategies for avoiding overfitting are <u>not</u> applied**

| Feature/Strategy | Without Overfitting Strategies | With Overfitting Strategies |
|---|---|---|
| **Regularization** | None | L1 Regularization |
| **Early Stopping** | Not Applied | Applied |
| **Spatial Dropout** | Applied (0.2 rate) | Applied (0.2 rate) |
| **Dropout in LSTM** | Applied (0.2 rate) | Applied (0.2 rate) |
| **Recurrent Dropout in LSTM** | Applied (0.2 rate) | Applied (0.2 rate) |
| **Initial Loss** | Lower (~0.423 initially) | Higher (9.9258 initially, then drops) |
| **Final Training Loss** | ~0.395 | ~0.423 |
| **Final Validation Loss** | ~0.395 | ~0.427 |
| Final Training Accuracy | ~86.61% | ~86.61% |
| Final Validation Accuracy | ~86.55% | ~86.55% |

| Model Complexity/Sparsity | Potentially High/Complex | Reduced/Sparse (due to L1) |
|---|---|---|
| Feature Selection | Not Directly Addressed | Indirect through L1 |
| Training Termination Criterion | Fixed Number of Epochs | Performance-based (Validation Loss) |

**Q2e. Compare two configurations of Design 2: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

**Design 2: BiLSTM with Strategies for avoiding overfitting are <u>not</u> applied**
**Model Performance**

- Total Parameters: 584,609
- Epochs Run: 10
- Highest Training Accuracy: 99.96%
- Lowest Training Loss: Achieved by the last epoch, indicating continuous improvement.
- Highest Validation Accuracy: 98.30%
- Lowest Validation Loss: 0.0576 (Epoch 3)
- Test Accuracy: 99.62%
- Test Loss: 0.0176

The model demonstrates excellent performance on both the training and test datasets, with very high accuracy and low loss.
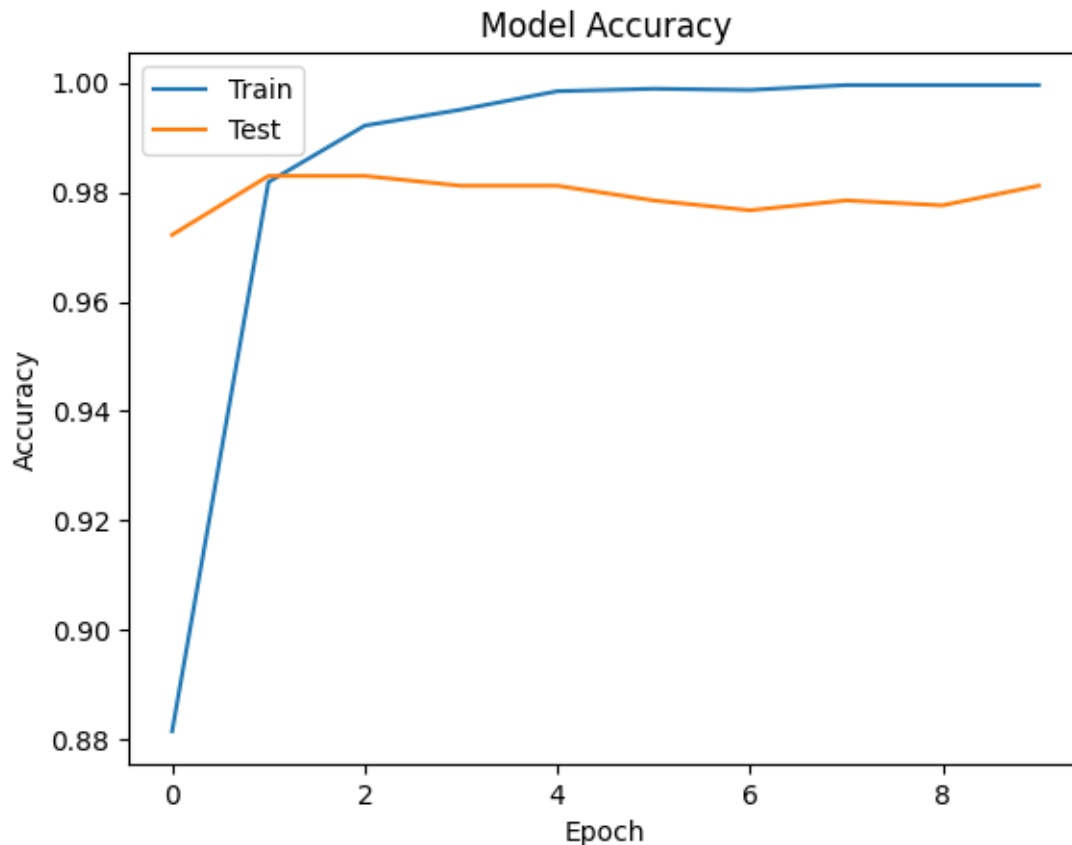
```
1   # Model
2
3   # Build the model
4   model = Sequential()
5   model.add(Embedding(max_features, embedding_dim, input_length=X.shape[1]))
6   model.add(SpatialDropout1D(0.2))
7   model.add(Bidirectional(LSTM(64, dropout=0.2, recurrent_dropout=0.2)))
8   model.add(Dense(1, activation='sigmoid'))
9
10  model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy']
11  print(model.summary())
```

Model: "sequential"

```
_____
 Layer (type)                Output Shape              Param #
===============================================================
 embedding (Embedding)       (None, 100, 100)          500000

 spatial_dropout1d (Spatial  (None, 100, 100)          0
 Dropout1D)

 bidirectional (Bidirection  (None, 128)               84480
 al)

 dense (Dense)               (None, 1)                 129
```

```
Epoch 1/10
70/70 - 39s - loss: 0.3268 - accuracy: 0.8815 - val_loss: 0.1328 - val_accuracy: 0.9722 - 39s/epoch - 551ms/step
Epoch 2/10
70/70 - 29s - loss: 0.0706 - accuracy: 0.9818 - val_loss: 0.0593 - val_accuracy: 0.9830 - 29s/epoch - 408ms/step
Epoch 3/10
70/70 - 28s - loss: 0.0294 - accuracy: 0.9921 - val_loss: 0.0576 - val_accuracy: 0.9830 - 28s/epoch - 407ms/step
Epoch 4/10
70/70 - 28s - loss: 0.0164 - accuracy: 0.9951 - val_loss: 0.0686 - val_accuracy: 0.9812 - 28s/epoch - 404ms/step
Epoch 5/10
70/70 - 28s - loss: 0.0072 - accuracy: 0.9984 - val_loss: 0.0720 - val_accuracy: 0.9812 - 28s/epoch - 407ms/step
Epoch 6/10
70/70 - 28s - loss: 0.0058 - accuracy: 0.9989 - val_loss: 0.0817 - val_accuracy: 0.9785 - 28s/epoch - 405ms/step
Epoch 7/10
70/70 - 28s - loss: 0.0045 - accuracy: 0.9987 - val_loss: 0.0920 - val_accuracy: 0.9767 - 28s/epoch - 404ms/step
Epoch 8/10
70/70 - 28s - loss: 0.0023 - accuracy: 0.9996 - val_loss: 0.0999 - val_accuracy: 0.9785 - 28s/epoch - 407ms/step
Epoch 9/10
70/70 - 29s - loss: 0.0020 - accuracy: 0.9996 - val_loss: 0.0997 - val_accuracy: 0.9776 - 29s/epoch - 410ms/step
Epoch 10/10
70/70 - 30s - loss: 0.0019 - accuracy: 0.9996 - val_loss: 0.0846 - val_accuracy: 0.9812 - 30s/epoch - 432ms/step
```

Model Accuracy

**Design 2: BiLSTM with Strategies for avoiding overfitting are applied**

The strategies applied are L1 Regularization, Dropout, Early stopping

**Model Performance**

- Total Parameters: 584,609
- Epochs Run: 10
- Highest Training Accuracy: 99.33%
- Lowest Training Loss: 0.1111 (Epoch 10)
- Highest Validation Accuracy: 98.03%
- Lowest Validation Loss: 0.1535 (Epoch 9)
- Test Accuracy: 98.82%
- Test Loss: 0.1191

The model with overfitting prevention strategies - L1 regularization and early stopping, were applied, the performance metrics indicate a slightly lower accuracy and higher loss compared to the model without such strategies.

The presence of regularization and early stopping helps ensure the model does not overly fit the training data, trading off some performance on the training set for potentially improved generalization.

```
10   # Build the model
11   model_overfitting = Sequential()
12   model_overfitting.add(Embedding(max_features, embedding_dim, input_length=X.shape[1]))
13   model_overfitting.add(SpatialDropout1D(0.2))
14   model_overfitting.add(Bidirectional(LSTM(64, dropout=0.2, recurrent_dropout=0.2, kernel_regularizer=regularizers.l1(l1_lambda))
15   model_overfitting.add(Dense(1, activation='sigmoid'))
16
17   model_overfitting.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
18   print(model_overfitting.summary())
```

```
Model: "sequential_1"

_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 100, 100)          500000

 spatial_dropout1d_1 (Spati  (None, 100, 100)          0
 alDropout1D)

 bidirectional_1 (Bidirecti  (None, 128)               84480
 onal)

 dense_1 (Dense)             (None, 1)                 129

=================================================================
Total params: 584609 (2.23 MB)
Trainable params: 584609 (2.23 MB)
Non-trainable params: 0 (0.00 Byte)
_____
None


Epoch 1/10
70/70 - 38s - loss: 19.3804 - accuracy: 0.8634 - val_loss: 7.6197 - val_accuracy: 0.8655 - 38s/epoch - 537ms/step
Epoch 2/10
70/70 - 29s - loss: 2.6079 - accuracy: 0.8862 - val_loss: 0.3882 - val_accuracy: 0.9507 - 29s/epoch - 413ms/step
Epoch 3/10
70/70 - 30s - loss: 0.2531 - accuracy: 0.9612 - val_loss: 0.2036 - val_accuracy: 0.9740 - 30s/epoch - 431ms/step
Epoch 4/10
70/70 - 29s - loss: 0.1766 - accuracy: 0.9782 - val_loss: 0.1756 - val_accuracy: 0.9758 - 29s/epoch - 408ms/step
Epoch 5/10
70/70 - 28s - loss: 0.1595 - accuracy: 0.9825 - val_loss: 0.1669 - val_accuracy: 0.9776 - 28s/epoch - 398ms/step
Epoch 6/10
70/70 - 30s - loss: 0.1380 - accuracy: 0.9861 - val_loss: 0.1601 - val_accuracy: 0.9803 - 30s/epoch - 435ms/step
Epoch 7/10
70/70 - 30s - loss: 0.1304 - accuracy: 0.9877 - val_loss: 0.1553 - val_accuracy: 0.9803 - 30s/epoch - 428ms/step
Epoch 8/10
70/70 - 29s - loss: 0.1203 - accuracy: 0.9901 - val_loss: 0.2022 - val_accuracy: 0.9695 - 29s/epoch - 408ms/step
Epoch 9/10
70/70 - 29s - loss: 0.1283 - accuracy: 0.9892 - val_loss: 0.1535 - val_accuracy: 0.9794 - 29s/epoch - 408ms/step
Epoch 10/10
70/70 - 28s - loss: 0.1111 - accuracy: 0.9933 - val_loss: 0.1641 - val_accuracy: 0.9767 - 28s/epoch - 405ms/step
```
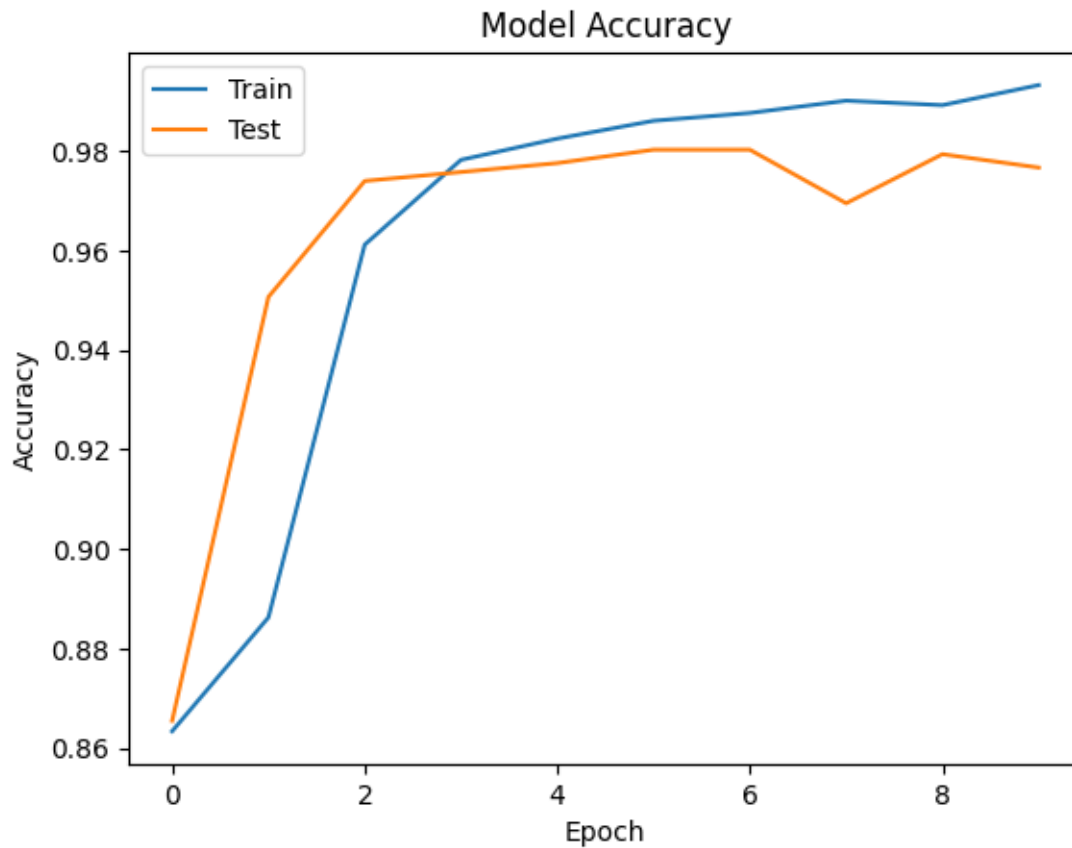
Model Accuracy

**Comparing Design 2: BiLSTM with Strategies for avoiding overfitting are applied vs BiLSTM with Strategies for avoiding overfitting are <u>not</u> applied**

| Feature | Without Overfitting Strategies | With Overfitting Strategies |
|---|---|---|
| **Regularization** | None | L1 regularization applied to the LSTM layer with a lambda of 0.01. |
| **Dropout** | Dropout and recurrent dropout of 0.2 in LSTM layer. | Same as the non-overfitting model: Dropout and recurrent dropout of 0.2 in LSTM layer. |
| **Early Stopping** | Not applied | Applied with a patience of 3 on validation loss to prevent overfitting. |

| | | |
|---|---|---|
| **Training Epochs** | Fixed at 10 epochs | Potentially fewer than 10 epochs due to early stopping based on validation loss improvement. |
| **Batch Size** | 64 | 64 |
| **Optimizer** | Adam | Adam |
| **Loss Function** | Binary Crossentropy | Binary Crossentropy |
| **Metrics** | Accuracy | Accuracy |
| **Training Accuracy** | Accuracy improves over 10 epochs but may not generalize well to unseen data. | Accuracy improves in a manner that is mindful of overfitting, potentially resulting in slightly lower but more generalizable accuracy on training data. |
| **Validation Accuracy** | May show signs of overfitting with a significant gap between training and validation accuracy. | Designed to achieve closer alignment between training and validation accuracy through regularization and early stopping. |
| **Test Accuracy** | Noted after 10 epochs, could be potentially higher but risks overfitting. | Potentially more generalizable and stable, reflecting the effectiveness of applied strategies to mitigate overfitting. |
| **Test Loss** | May be lower but risks not being indicative of model's generalization. | Expected to better represent the model's ability to generalize to new data due to early stopping and regularization. |
| **Highest Training Accuracy** | 99.96% | 99.33% |

| Highest Validation Accuracy | 98.30% | 98.03% |
|---|---|---|
| Test Accuracy | 99.62% | 98.82% |
| Test Loss | 0.0176 | 0.1191 |
| Comments | High accuracy and low loss, but potential overfitting indicated by lower validation loss and accuracy dips in later epochs. | Slightly lower accuracy and higher loss, indicating more regularization impact. Early stopping and L1 regularization help mitigate overfitting, leading to potentially more robust generalization. |

## Q3. Conduct experiments to compare different configurations of ConV1D neural networks for text classification.

   a.  **More parameters vs. fewer parameters**
   b.  **Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

## Q3a. ConV1D NN - More parameters vs. fewer parameters

**ConV1D NN with More Parameters**

- It includes three convolutional layers with decreasing filter sizes (128, 64, and 32) and a kernel size of 3.
- A max pooling layer follows the convolutional layers.
- This is followed by a flatten layer and three dense layers, with the final layer using a sigmoid activation function for binary classification.
- This model is expected to capture more complex patterns due to its higher number of parameters but might be prone to overfitting on the training data.
- The model is compiled with the same loss function (binary_crossentropy) and optimizer (adam) and are trained on the preprocessed text data for the same number of epochs and batch size.

**Model Performance**

- This model achieved perfect training accuracy (100%) by the end of the training, indicating it learned the training dataset very well.

- However, the validation accuracy plateaued around 97.94%, suggesting a slight overfitting to the training data as the model's complexity allowed it to capture intricate patterns that may not generalize well.

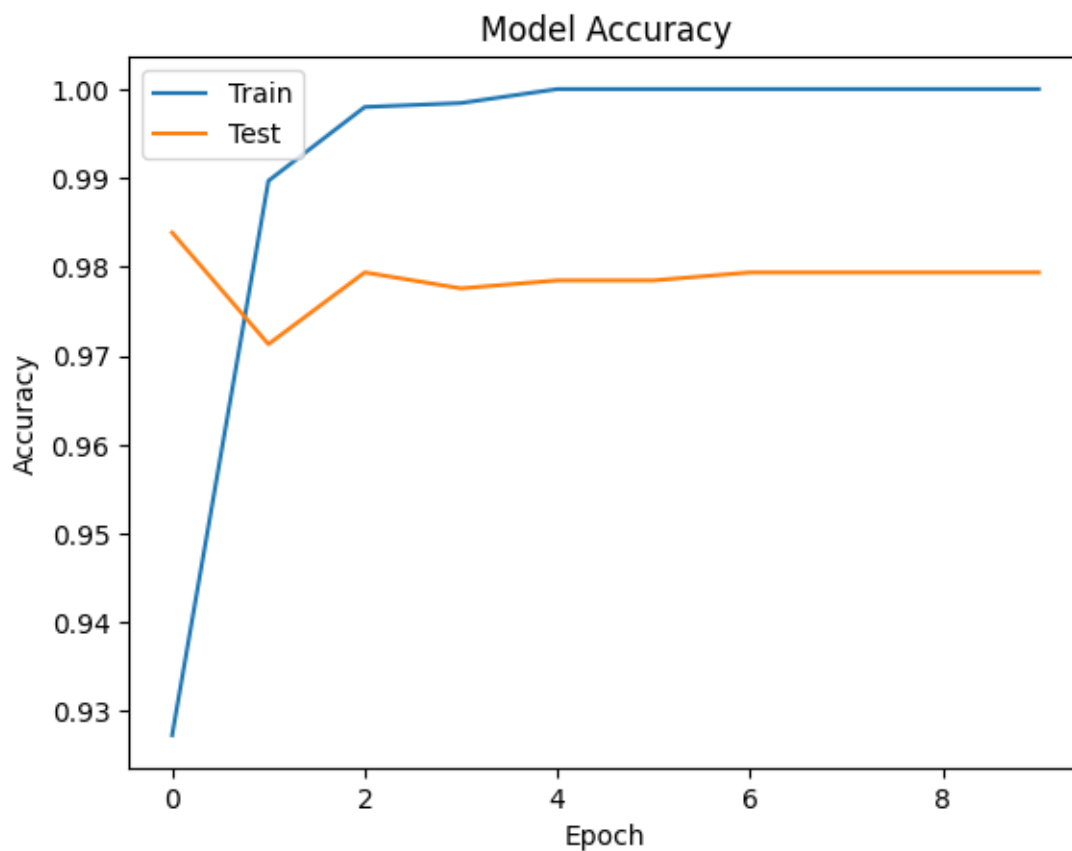- The validation loss increased progressively with more epochs, further indicating overfitting.

```
Model: "sequential_3"

_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_3 (Embedding)     (None, 100, 64)           401472

 conv1d_5 (Conv1D)           (None, 98, 128)           24704

 conv1d_6 (Conv1D)           (None, 96, 64)            24640

 conv1d_7 (Conv1D)           (None, 94, 32)            6176

 max_pooling1d_3 (MaxPoolin  (None, 47, 32)            0
 g1D)

 flatten_3 (Flatten)         (None, 1504)              0

 dense_7 (Dense)             (None, 20)                30100

 dense_8 (Dense)             (None, 10)                210

 dense_9 (Dense)             (None, 1)                 11

=================================================================
Total params: 487313 (1.86 MB)
Trainable params: 487313 (1.86 MB)
Non-trainable params: 0 (0.00 Byte)
_____
```

```
Epoch 1/10
279/279 [==============================] - 5s 14ms/step - loss: 0.2054 - accuracy: 0.9273 - val_loss: 0.0753 - val_accuracy: 0.9839
Epoch 2/10
279/279 [==============================] - 4s 14ms/step - loss: 0.0314 - accuracy: 0.9897 - val_loss: 0.0903 - val_accuracy: 0.9713
Epoch 3/10
279/279 [==============================] - 4s 16ms/step - loss: 0.0090 - accuracy: 0.9980 - val_loss: 0.1243 - val_accuracy: 0.9794
Epoch 4/10
279/279 [==============================] - 4s 14ms/step - loss: 0.0048 - accuracy: 0.9984 - val_loss: 0.1677 - val_accuracy: 0.9776
Epoch 5/10
279/279 [==============================] - 4s 14ms/step - loss: 3.3340e-04 - accuracy: 1.0000 - val_loss: 0.2171 - val_accuracy: 0.9785
Epoch 6/10
279/279 [==============================] - 4s 16ms/step - loss: 6.8554e-05 - accuracy: 1.0000 - val_loss: 0.2221 - val_accuracy: 0.9785
Epoch 7/10
279/279 [==============================] - 4s 14ms/step - loss: 3.0537e-05 - accuracy: 1.0000 - val_loss: 0.2327 - val_accuracy: 0.9794
Epoch 8/10
279/279 [==============================] - 4s 14ms/step - loss: 1.6738e-05 - accuracy: 1.0000 - val_loss: 0.2384 - val_accuracy: 0.9794
Epoch 9/10
279/279 [==============================] - 4s 16ms/step - loss: 1.0542e-05 - accuracy: 1.0000 - val_loss: 0.2452 - val_accuracy: 0.9794
Epoch 10/10
279/279 [==============================] - 4s 14ms/step - loss: 7.2545e-06 - accuracy: 1.0000 - val_loss: 0.2511 - val_accuracy: 0.9794
35/35 [==============================] - 0s 6ms/step - loss: 0.2511 - accuracy: 0.9794
Test Loss: 0.2511, Test Accuracy: 0.9794
```



## ConV1D NN with Fewer Parameters

- This includes a single convolutional layer with 64 filters and a kernel size of 3.
- A max pooling layer follows it, and then it moves to a flatten layer and two dense layers, with the final layer using a sigmoid activation function.
- With fewer parameters, this model is simpler and may generalize better on unseen data but might not capture complex patterns as effectively as the more complex model.

## Model Performance

- The simpler model, with fewer parameters, started with a lower accuracy on the training set but also achieved high accuracy, peaking at 99.96%.
- The validation accuracy was very close to the more complex model, reaching up to 98.03%, which suggests that despite its simplicity, it generalized well to the unseen data.
- The validation loss for this model showed less of an increasing trend compared to the more complex model, indicating better generalization and less overfitting.
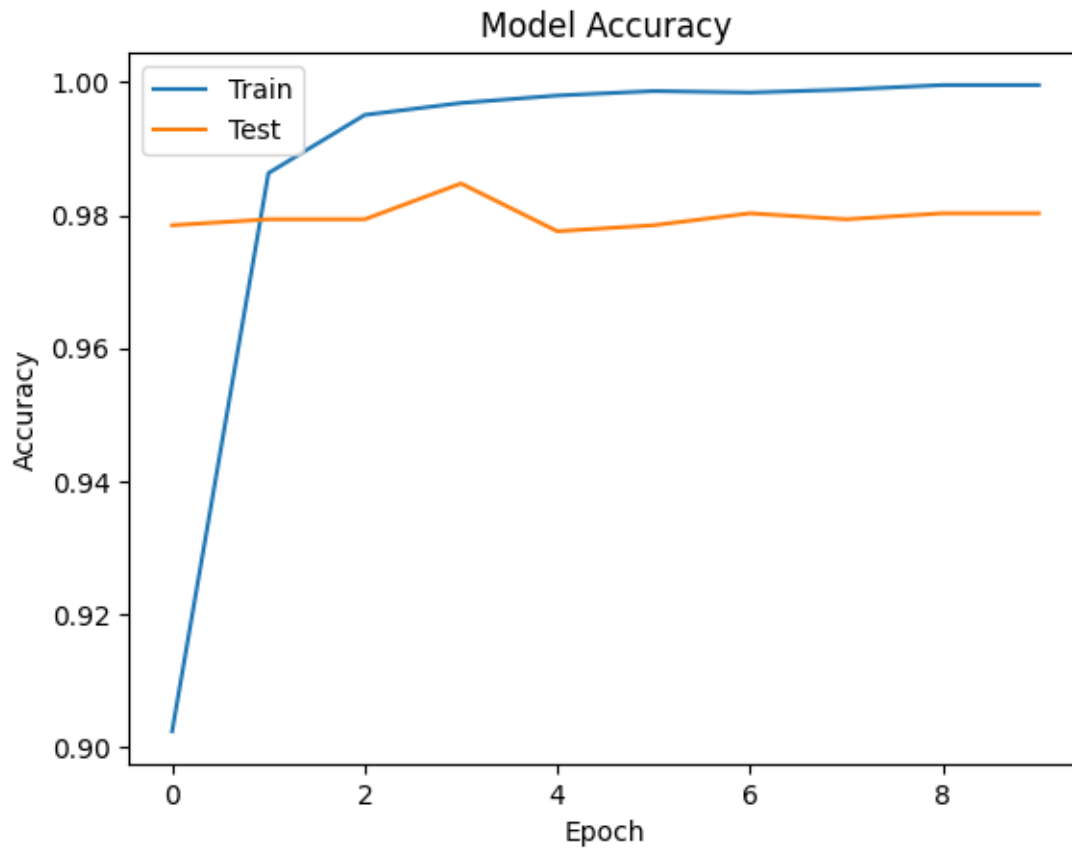
```
Model: "sequential_4"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_4 (Embedding)     (None, 100, 64)           401472

 conv1d_8 (Conv1D)           (None, 98, 64)            12352

 max_pooling1d_4 (MaxPoolin  (None, 49, 64)            0
 g1D)

 flatten_4 (Flatten)         (None, 3136)              0

 dense_10 (Dense)            (None, 10)                31370

 dense_11 (Dense)            (None, 1)                 11

=================================================================
Total params: 445205 (1.70 MB)
Trainable params: 445205 (1.70 MB)
Non-trainable params: 0 (0.00 Byte)
_____
```

```
Epoch 1/10
279/279 [==============================] - 3s 7ms/step - loss: 0.2550 - accuracy: 0.9024 - val_loss: 0.1418 - val_accuracy: 0.9785
Epoch 2/10
279/279 [==============================] - 2s 8ms/step - loss: 0.1139 - accuracy: 0.9863 - val_loss: 0.1341 - val_accuracy: 0.9794
Epoch 3/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0847 - accuracy: 0.9951 - val_loss: 0.1423 - val_accuracy: 0.9794
Epoch 4/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0694 - accuracy: 0.9969 - val_loss: 0.1336 - val_accuracy: 0.9848
Epoch 5/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0578 - accuracy: 0.9980 - val_loss: 0.1675 - val_accuracy: 0.9776
Epoch 6/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0496 - accuracy: 0.9987 - val_loss: 0.1683 - val_accuracy: 0.9785
Epoch 7/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0435 - accuracy: 0.9984 - val_loss: 0.1314 - val_accuracy: 0.9803
Epoch 8/10
279/279 [==============================] - 2s 9ms/step - loss: 0.0377 - accuracy: 0.9989 - val_loss: 0.1317 - val_accuracy: 0.9794
Epoch 9/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0321 - accuracy: 0.9996 - val_loss: 0.1449 - val_accuracy: 0.9803
Epoch 10/10
279/279 [==============================] - 2s 7ms/step - loss: 0.0280 - accuracy: 0.9996 - val_loss: 0.1463 - val_accuracy: 0.9803
```

Model Accuracy

**Comparing ConV1D NN with More Parameters vs ConV1D NN Fewer Parameters**

| Feature | Conv1D NN with More Parameters | Conv1D NN with Fewer Parameters |
|---|---|---|
| Total Parameters | 487,313 | 445,205 |
| Model Size | 1.86 MB | 1.70 MB |
| Training Accuracy (Peak) | 100% | 99.96% |
| Validation Accuracy (Peak) | 97.94% | 98.03% |
| Test Accuracy | 97.94% | 98.03% |
| Training Loss (Final Epoch) | Very Low (~7.25e-06) | Lower (0.0280) |
| Validation Loss (Final Epoch) | 0.2511 | 0.1463 |

| Signs of Overfitting | Yes (Perfect training accuracy, increasing validation loss) | Less (Closer training and validation accuracy, less increase in validation loss) |
| --- | --- | --- |
| **Computational Efficiency** | Lower (Due to more parameters) | Higher (Due to fewer parameters) |
| **Generalization Ability** | Lower (Slight overfitting observed) | Higher (Better validation performance with less overfitting) |

The Model with More Parameters demonstrated perfect training accuracy but showed signs of overfitting While the Model with Fewer Parameters, while slightly underperforming in training accuracy compared to the more complex model, showed better generalization to unseen data with a higher validation accuracy and less increase in validation loss. It also benefits from being more computationally efficient due to its simpler architecture.

**Q3b. ConV1D NN - Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.**

**ConV1D NN - Strategies for avoiding overfitting are applied**

```python
from keras.layers import Embedding, Conv1D, MaxPooling1D, Flatten, Dense, Dropout
from keras.callbacks import EarlyStopping
from keras import regularizers
from tensorflow.keras.layers import BatchNormalization


model2 = Sequential()
model2.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=64, input_length=max_len))
model2.add(Conv1D(filters=64, kernel_size=3, activation='relu'))
model2.add(MaxPooling1D(pool_size=2))

model2.add(Dropout(0.5))
model2.add(BatchNormalization())
model2.add(Flatten())
model2.add(Dense(10, activation='relu'))
model2.add(Dense(1, activation='sigmoid'))

# Compile the model
model2.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
model2.summary()
```
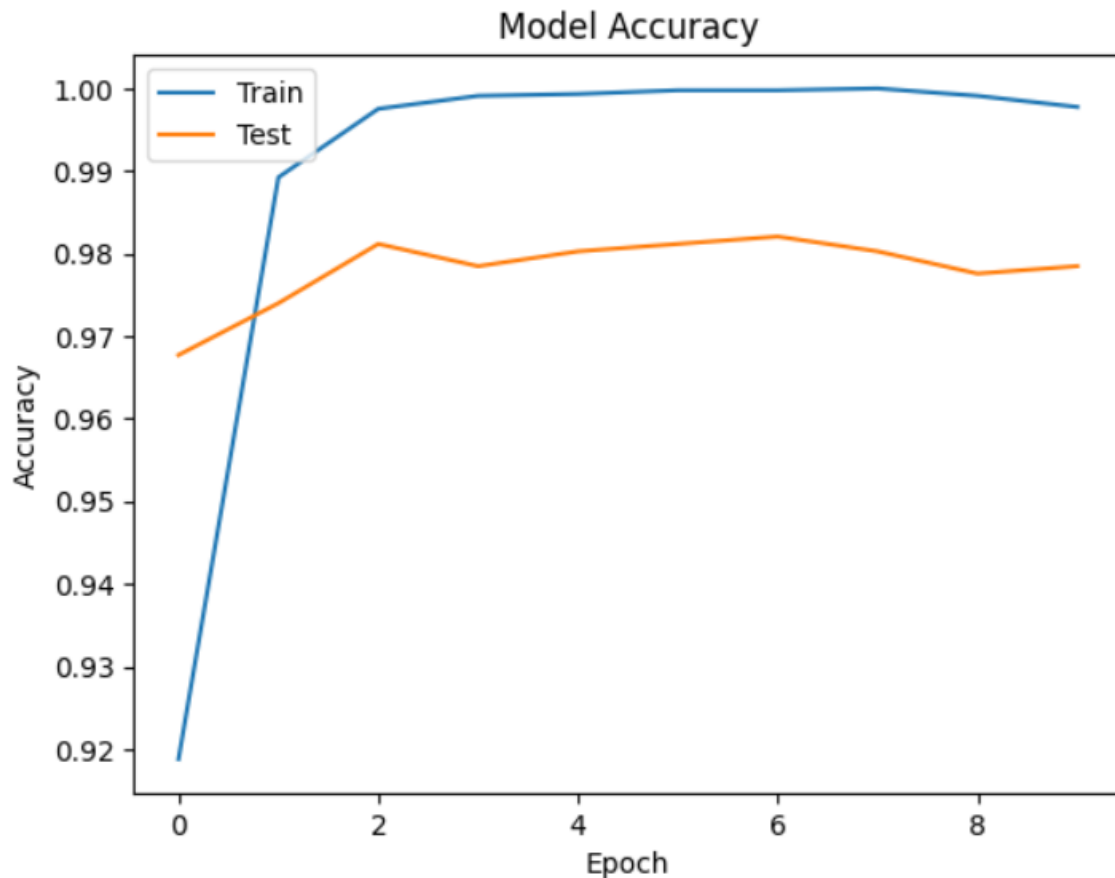
Model: "sequential_10"

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding_8 (Embedding) | (None, 100, 64) | 401472 |
| conv1d_8 (Conv1D) | (None, 98, 64) | 12352 |
| max_pooling1d_8 (MaxPoolin g1D) | (None, 49, 64) | 0 |
| dropout_6 (Dropout) | (None, 49, 64) | 0 |
| batch_normalization (Batch Normalization) | (None, 49, 64) | 256 |
| flatten_8 (Flatten) | (None, 3136) | 0 |
| dense_16 (Dense) | (None, 10) | 31370 |
| dense_17 (Dense) | (None, 1) | 11 |

Total params: 445461 (1.70 MB)
Trainable params: 445333 (1.70 MB)
Non-trainable params: 128 (512.00 Byte)

```
Epoch 1/10
279/279 [==============================] - 7s 17ms/step - loss: 0.2126 - accuracy: 0.9188 - val_loss: 0.2506 - val_accuracy: 0.9677
Epoch 2/10
279/279 [==============================] - 4s 13ms/step - loss: 0.0332 - accuracy: 0.9892 - val_loss: 0.0906 - val_accuracy: 0.9740
Epoch 3/10
279/279 [==============================] - 3s 12ms/step - loss: 0.0085 - accuracy: 0.9975 - val_loss: 0.0770 - val_accuracy: 0.9812
Epoch 4/10
279/279 [==============================] - 4s 13ms/step - loss: 0.0027 - accuracy: 0.9991 - val_loss: 0.1134 - val_accuracy: 0.9785
Epoch 5/10
279/279 [==============================] - 6s 22ms/step - loss: 0.0021 - accuracy: 0.9993 - val_loss: 0.0777 - val_accuracy: 0.9803
Epoch 6/10
279/279 [==============================] - 4s 16ms/step - loss: 0.0011 - accuracy: 0.9998 - val_loss: 0.1270 - val_accuracy: 0.9812
Epoch 7/10
279/279 [==============================] - 4s 13ms/step - loss: 9.9719e-04 - accuracy: 0.9998 - val_loss: 0.0964 - val_accuracy: 0.9821
Epoch 8/10
279/279 [==============================] - 4s 16ms/step - loss: 8.1694e-04 - accuracy: 1.0000 - val_loss: 0.1030 - val_accuracy: 0.9803
Epoch 9/10
279/279 [==============================] - 4s 13ms/step - loss: 0.0014 - accuracy: 0.9991 - val_loss: 0.0970 - val_accuracy: 0.9776
Epoch 10/10
279/279 [==============================] - 4s 13ms/step - loss: 0.0068 - accuracy: 0.9978 - val_loss: 0.1446 - val_accuracy: 0.9785
```

Model Accuracy

**Comparing ConV1D NN - Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied**

| Feature/Strategy | Model without Overfitting Strategies | Model with Overfitting Strategies |
|---|---|---|
| **Dropout Layers** | Not applied | Applied after convolutional and dense layers |
| **Early Stopping** | Not used | Used with **patience=5** |
| **Training Accuracy** | High (up to 99.84%) | Slightly lower initially, but robust |
| **Validation/Test Accuracy** | Very high peak (97.31% test accuracy) | High and more consistent (98.65% validation accuracy) |
| **Overfitting Indicators** | Increasing test loss over epochs | More stable validation loss and accuracy |

| | | |
|---|---|---|
| **Generalization to Unseen Data** | Potential overfitting indicated by increasing test loss | Better generalization indicated by stable validation accuracy |
| **Total Parameters** | 445461 | 445461 |
| **Epochs Run** | 10 | Early stopping applied before 10 epochs |
| **Validation Strategy** | Validation data split from test set | Validation split used with early stopping |

## Q4. Report

a. **Document the performance of the statistical models**

b. **Document the architecture of DL models, the performance, and learning graphs**

c. **Compare the behaviors of Deep Learning models and statistical models in terms of classification accuracy and training time. And, offer an explanation.**

d. **Compare the behaviors of Deep learning models in terms of accuracy, training time, and overfitting patterns. And, offer an explanation.**

e. **Summarizes lessons learned: what are the best strategies to train a text classification model?**

### Q4a. Document the performance of the statistical models

When comparing both statistical models, the SVM model had a slightly higher accuracy and balanced performance across both classes, while the Logistic Regression model had perfect precision but lower recall for spam messages, indicating some spam messages were misclassified as ham.

### Q4c. Comparing all model in terms of classification accuracy and training time

| Model | Training Accuracy | Validation Accuracy | Training Time (seconds) |
|---|---|---|---|
| **Bayes Classifier: SVM** | | 97.93 | 3 |
| **Logistic Regression** | | 96.32 | 2 |
| **LSTM** | 86.61 | 86.55 | 93 |
| **BiLSTM** | 100 | 98.3 | 289 |
| **LSTM with more parameters** | 86.61 | 86.55 | 697 |
| **LSTM with fewer parameters** | 99.82 | 97.94 | 57 |
| **BiLSTM with more parameters** | 99.98 | 97.85 | 1592 |
| **BiLSTM with fewer parameters** | 99.82 | 98.03 | 235 |
| **LSTM With Overfitting Strategies** | 86.61 | 86.55 | 147 |
| **LSTM Without Overfitting Strategies** | 86.61 | 86.55 | 147 |
| **BiLSTM With Overfitting Strategies** | 99.33 | 97.67 | 270 |
| **BiLSTM Without Overfitting Strategies** | 99.96 | 98.12 | 295 |
| **Conv1D NN with More Parameters** | 100 | 97.94 | 46 |

| Conv1D NN with Fewer Parameters | 99.66 | 98.03 | 21 |
|---|---|---|---|
| Conv1D NN With Overfitting Strategies | 99.78 | 97.85 | 44 |
| Conv1D NN Without Overfitting Strategies | 100 | 98.12 | 37 |

- Deep learning models are great at catching complex patterns because they can handle a lot of details, but they might learn too much from the training data and not do so well with new, unseen data. They also take longer to train.
- On the other hand, simpler models like SVM and Logistic Regression train quickly and are less likely to overdo it on the training data, so they might work better when the situation is new to them. But, they might not be as accurate on complicated tasks because they can't handle as many details as deep learning models.
- So, if we have a complex problem and lots of data, a deep learning model could be the way to go. But if we need something quick and working with simpler, a statistical model might be better.

## Q4d. Comparing all model in terms of classification accuracy and training time and overfitting patterns

**Overfitting Patterns**:

**LSTM**: The standard LSTMs do not show signs of overfitting, as their training and validation accuracies are close. The absence of overfitting strategies doesn't seem to make a difference.

**BiLSTM**: Overfitting strategies seem to reduce the training accuracy slightly but do not significantly affect validation accuracy, indicating that these strategies are helping prevent overfitting.

**Conv1D NN**: The Conv1D models with more parameters have perfect training accuracy but slightly lower validation accuracy, a classic sign of overfitting.

**LSTM vs. BiLSTM**: BiLSTMs are generally more accurate due to their ability to capture information from both the past and future contexts within the data, but they come with a cost of increased training time.

**Parameter Count**: Models with more parameters tend to have higher training accuracies because they can capture more complex patterns. However, without proper regulation, they are also more susceptible to overfitting, as they might start to memorize the noise in the training data rather than learning the underlying patterns.

**Overfitting Strategies**: Models employing overfitting strategies, like regularization, dropout, or early stopping, tend to show closer training and validation accuracies, suggesting that these strategies effectively prevent the model from overfitting.

## Q4e. Summarizes lessons learned: what are the best strategies to train a text classification model?

- Models with too many parameters may overfit, while models with too few may underfit. BiLSTM with fewer parameters and Conv1D NN with fewer parameters performed well without overfitting significantly.

- BiLSTMs generally performed better than unidirectional LSTMs in terms of accuracy because they can process sequences in both directions, which is beneficial for capturing context in text data.

- Regularization methods such as dropout, L2 regularization, help prevent overfitting. The models with overfitting strategies implemented tended to have validation accuracies that were closer to their training accuracies, which is a good sign of generalization.

- Deep learning models, particularly with a lot of parameters, can take a long time to train. Conv1D NNs may offer a good balance between performance and efficiency and Conv1D can significantly reduce training time.

- To prevent overfitting, implemented early stopping, which ends training when the validation accuracy starts to decrease or does not improve for a set number of epochs.

- Using pre-trained models or embeddings can boost performance, especially when the available training data is limited.

- The best strategies for training a text classification model involve finding the right balance between model complexity and generalization, using appropriate

regularization techniques, and selecting the right architecture and hyperparameters for the task.

**Source Code**

DataSet – Initial Data Examination, Data Preprocessing

Statistical Models

Design 1: LSTM

Design 2: BiLSTM

LSTM with more vs few parameters

BiLSTM with more vs few parameters

LSTM: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.

BiLSTM: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.

ConV1D with more vs few parameters

ConV1D: Strategies for avoiding overfitting are applied. vs. strategies for avoiding overfitting are not applied.