# SMS SPAM DETECTION

Pavani Billapati(<a href="Pbill1@unh.newhaven.edu">Pbill1@unh.newhaven.edu</a>)
Trilok kumar pidikiti(<a href="tpidi1@unh.newhaven.edu">tpidi1@unh.newhaven.edu</a>)
Harianth Kancherla(<a href="https://hkanc1@unh.newhaven.edu">hkanc1@unh.newhaven.edu</a>)

#### **Abstract:**

In the era of pervasive mobile communication, the surge in SMS spam presents a formidable threat, with malicious actors exploiting unsuspecting individuals through deceptive messages for extracting sensitive information. This research addresses the urgent need for effective SMS spam detection, aiming to curtail the risk of individuals falling victim to phishing attacks and incurring financial losses. We propose an innovative SMS spam detection system, leveraging advanced Natural Language Processing (NLP) techniques and employing Deep Learning architectures such as Long Short-Term Memory (LSTM), Bidirectional LSTM, and Dense Neural Networks.

Concentrating on English language SMS spams, our methodology encompasses meticulous data preparation involving word embedding, tokenization, padding, and truncation to enhance the dimensionality of the data. The model development phase incorporates a Dense Neural Network for baseline understanding and compares it with the sequence modeling capabilities of LSTM and the bidirectional nature of Bi-LSTM.

This research contributes a robust solution to SMS spam detection, emphasizing the efficacy of advanced NLP and Deep Learning techniques in fortifying individuals against phishing attempts and preserving their financial well-being. The diverse model architectures considered in our approach signify a comprehensive exploration of strategies to address the evolving landscape of SMS spam threats.

#### **Introduction:**

The state of communication technology is currently quickly advancing. This implies that with the use of email, text messaging, and web surfing, anyone may now obtain or receive information more easily than in the past. By sending a false notice of a bank transaction or a misleading advertisement, a hacker can use these technologies' benefits to obtain sensitive information from users, including phone numbers, credit card numbers, and online banking account information [1]. Furthermore, the unsuspecting individuals who received the message in a panic then obeyed the hacker's instructions, providing them with their private information. The hacker then utilizes this data to obtain assets from individuals.

Spam that is sent via mobile devices and targets users' sensitive information is known as short messaging service (SMS) spam [2], [3]. These days, a variety of people who follow the hacker's instructions and believe the information in the message are severely impacted by an SMS spam attack. Having a technology that can efficiently identify spam messages can help to mitigate this issue.

Deep Learning is currently a widely utilized data analysis technique since it often yields excellent prediction accuracy [4], [5]. Recurrent Neural Networks (RNNs) are the method that is commonly applied when assessing sequential data, like text data [5]. The LSTM, or long short-term memory, Bi-LSTM and Dense Model algorithms are adjusted for RNNs. Many published studies assert that LSTM Bi-LSTM and Dense

Model perform better in deep learning than other methods, especially for the analysis of sequential data. Taking Kraus and Feuerriegel as examples developed a decision support system to help investors feel more safe about their investments by using economic news and long short-term memory (LSTM). Moreover, they employ word tokenization, transfer learning, utilizing word embedding to get data ready.

In this study, we make use of Keras, a Python toolkit for developing deep learning models with a Tensorflow backend. Its data preparation toolset includes word tokenization, data padding and truncation, and word embedding. Using the word tokenization process, text inputs are converted into sequential data using the words' index values. While the word embedding approach is used to add extra dimensions to the sequence into the vector, To ensure that each sequence is the same length, the techniques of padding and truncating data are applied. After the data preparation step, we use the LSTM, Bi-LSTM and Dense Model algorithms to train the model. Next, we assess the models' performance.

#### **Dataset:**

We use an SMS spam dataset that was downloaded from UCI datasets. This collection has roughly 5,572 records. It includes English-language SMS text messaging chats with text and numbers in varying sentence lengths. Every record in this dataset has a label already. The label for the spam messages is 'spam' (747 records), while the label for the regular mails is 'ham' (4,825 records).

[5]	messages.describe()			
		label	message	
	count	5572	5572	
	unique	2	5169	
	top	ham	Sorry, I'll call later	
	freq	4825	30	

## **Preparing Data:**

Natural Language Process (NLP) is used in this procedure to pre-process natural language data. The goal of natural language processing (NLP) is to enable computers to comprehend natural language in the same way that humans do. It offers a plethora of methods for preprocessing data into a machine-readable format. In this study, we text input into sequential data, which we then use to build LSTM, Bidirectional LSTM, and Dense Model algorithms are used to create SMS classification models. For pre-processing data, we additionally employ word tokenization, padding, truncating, and word embedding algorithms. The following is a description of each method we employ in the data pre-processing step.

#### **Word Tokenization:**

The process of converting sentences words into index values represented by a number is called word tokenization. To develop a word tokenizer, we established a lot of intriguing vocabulary terms in this process. We utilize word token to convert the words in the sentence to serial data after generating the token. If the word is unknown, the tokenizer sets the index to 0 and converts the word to index. Furthermore, we

generate a tokenizer by assigning a fixed 10,000 vocabulary items. Additionally, we employ the tokenizer is used to convert text data into an index number word sequence.

```
word_index = tokenizer.word_index
word index
{'<00V>': 1,
 'to': 2,
 'you': 3,
 'a': 4,
 'i': 5,
 'call': 6,
 'the': 7,
 'your': 8,
 'u': 9,
 'for': 10,
 'is': 11,
 'and': 12,
 'free': 13,
 'now': 14,
 'have': 15,
 '2': 16,
 'or': 17,
 'in': 18,
 'on': 19,
 'ur': 20,
 'of': 21,
 '4': 22,
 'txt': 23,
 'are': 24,
 'from': 25,
```

## **Padding and Truncating Data:**

Data Truncation and Padding During this procedure, we use the LSTM, BI-LSTM and Dense Model algorithms to ensure that every sequence in the dataset has the same length for training. Based on (1), we determine the optimal message length. Once the message's length has been optimized, we pad any data that is shorter than the ideal length by adding 0 to the start of the sequence until the data length corresponds to the ideal length. We start at the beginning and truncate the data until its length equals the ideal length. if the data's length exceeds that of the sample. Optimize is equal to mean  $(len(xi)) + 2 \times std(len(xi))$ .

where xi represents the dataset's records, mean(x) represents the mean of the data, std(x) represents the message length function. Based on (1), we calculate the ideal message length. function. Based on (1), we determine the optimal message length. After computation, the ideal length for cover data is 200, which accounts for 97.95% of datasets.

Consequently, we regulate padding and truncating in sequence using this ideal length.

## **Word Embedding:**

The word embedding technique developed by Pennington et al. is employed in this study. Using this method, a preprocessed word sequence is transformed into a vector representation known as embedding space, which has It has more dimensions than traditional word data and is used to train LSTM, BI-LSTM, and Dense Model algorithms. We utilize the word embedding approach to create additional dimensions for the data in sequence after padding and truncating it. We set the embedding size to 32.

## **Modeling**

In this study, we are developing SMS spam classification models using the three deep learning algorithms listed below: Bidirectional LSTM, LSTM and Dense model. Detailed description of each algorithm below.

## **Long Short-Term Memory (LSTM):**

Hochreiter and Schmidhuber, a German team, first introduced LSTM in 1997. To solve the vanishing problem, LSTM enhances the base RNN algorithm. LSTM incorporates cell states to remember or forget data. These cell states have a structure known as a cell gate. The cell gate is made up of four pieces. The input gate, the forget gate, the memory cell state gate, and the output gate.

The gate that controls the input data is known as the input. The forgotten gate is in charge of the previous hidden state. The previous hidden state is to be saved in the current hidden state's memory cell.

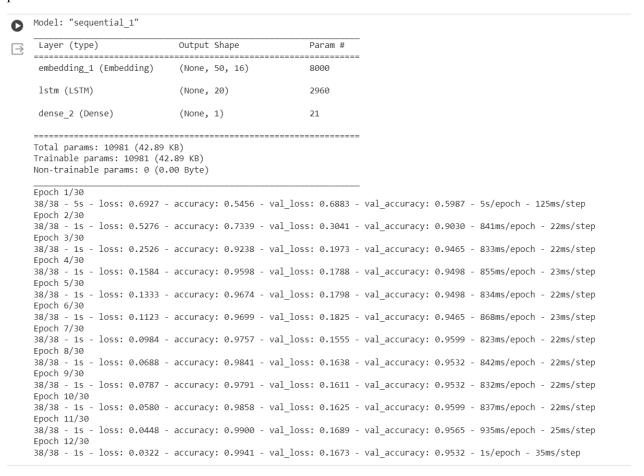
The memory cell state gate updates the data based on input and forget gate information. The output gate computes the network's output data based on memory-cell state.

```
#LSTM hyperparameters
    n lstm = 20
    drop lstm =0.2
    # Define LSTM Spam detection architecture
    lstm model = Sequential()
    lstm model.add(Embedding(vocab size, embedding dim, input length=max len))
    lstm model.add(LSTM(n lstm, dropout=drop lstm))
    lstm_model.add(Dense(1, activation='sigmoid'))
    lstm model.summary()
    # Compile the model
    lstm model.compile(loss='binary crossentropy', optimizer='adam', metrics=['accuracy'])
    # Fitting the LSTM spam detection model
    num epochs = 30
    early_stop = EarlyStopping(monitor='val_loss', patience=5)
    history = lstm_model.fit(training_padded, train_labels, epochs=num_epochs, validation_data=(testing_padded, test_labels), callbacks=[early_stop], verbose=2)
    # Save the model
    lstm model.save('LSTM Spam Detection.h5')
```

The architecture of the LSTM is very good at capturing sequence patterns and dependencies. This makes it ideal for natural language understanding (NLU) tasks. In the SMS spam detection context, it allows the model to understand the nuances of the spam messages and distinguish them from the legitimate ones.

By adjusting the hyperparameters (e.g., LSTM number of units), the model can capture long-lasting dependencies and improve its performance.

This code implements an LSTM-based neural network for SMS spam detection. Within the architecture, words are represented by an embedding layer, sequential patterns are captured by an LSTM layer, and binary classification is performed by a dense layer. The Adam optimizer and binary cross-entropy loss are used to train the model with early stopping, to prevent overfitting. The trained model is then saved for use later. Experimenting with model architecture and hyperparameter adjustments can help further optimize performance.



#### **BIDIRECTIONAL LSTM:**

A variation of the classic Long Short-Term Memory (LSTM) neural network architecture, Bidirectional Long Short-Term Memory (Bi-LSTM) anticipated to extract contextual information from both previous and subsequent elements in a series. Information processing usually happens in a unidirectional fashion, starting at the beginning and ending at the end of a sequence, in many sequential data tasks. The model's

comprehension of the sequence's dependencies and context may be limited by this unidirectional method.

```
# Biderectional LSTM Spam detection architecture

bidirectional_lstm_model = Sequential()

bidirectional_lstm_model.add(Embedding(vocab_size, embedding_dim, input_length=max_len))

bidirectional_lstm_model.add(Bidirectional(LSTM(n_lstm, dropout=drop_lstm)))

bidirectional_lstm_model.add(Dense(1, activation='sigmoid'))

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

# Training

num_epochs = 30

early_stop = EarlyStopping(monitor='val_loss', patience=4)

history = bidirectional_lstm_model.fit(training_padded, train_labels, epochs=num_epochs,

validation_data=(testing_padded, test_labels),callbacks =[early_stop], verbose=2)

bidirectional_lstm_model.save('BiLSTM_Spam_Detection.h5')
```

```
→ Epoch 1/30
    38/38 - 7s - loss: 0.6863 - accuracy: 0.6025 - val loss: 0.6621 - val accuracy: 0.8629 - 7s/epoch - 191ms/step
    Epoch 2/30
    38/38 - 1s - loss: 0.5354 - accuracy: 0.8360 - val loss: 0.4017 - val accuracy: 0.8863 - 1s/epoch - 35ms/step
    Epoch 3/30
    38/38 - 1s - loss: 0.3001 - accuracy: 0.9063 - val loss: 0.1975 - val accuracy: 0.9264 - 1s/epoch - 35ms/step
    Epoch 4/30
    38/38 - 1s - loss: 0.1194 - accuracy: 0.9640 - val loss: 0.1671 - val accuracy: 0.9365 - 1s/epoch - 33ms/step
    Epoch 5/30
    38/38 - 2s - loss: 0.0932 - accuracy: 0.9732 - val loss: 0.1433 - val accuracy: 0.9565 - 2s/epoch - 46ms/step
    Epoch 6/30
    38/38 - 2s - loss: 0.0735 - accuracy: 0.9749 - val loss: 0.1536 - val accuracy: 0.9498 - 2s/epoch - 62ms/step
    Epoch 7/30
    38/38 - 2s - loss: 0.0741 - accuracy: 0.9782 - val loss: 0.1508 - val accuracy: 0.9532 - 2s/epoch - 52ms/step
    Epoch 8/30
    38/38 - 1s - loss: 0.0406 - accuracy: 0.9908 - val loss: 0.1526 - val accuracy: 0.9565 - 1s/epoch - 33ms/step
    Epoch 9/30
    38/38 - 1s - loss: 0.0427 - accuracy: 0.9883 - val loss: 0.1813 - val accuracy: 0.9465 - 1s/epoch - 33ms/step
```

To overcome this constraint, Bi-LSTM combines two independent LSTM layers: one that processes the sequence forward (from the start) and the other that processes it backward (from the finish). The model can consider each element in the sequence's past and future context thanks to this bidirectional processing.

While the backward LSTM layer gathers data from the sequence's future elements, the forward LSTM layer gathers data from the sequence's past elements. The model obtains a more thorough grasp of the context and dependencies within the sequence by integrating the two directions.

Before the outputs from both LSTM layers are sent to the neural network's later layers, they are usually combined, frequently by concatenation. Information from both directions is combined in this concatenated representation.

#### **Dense Model:**

Every node (or neuron) in one layer of an artificial neural network is connected to every other layer's node through a dense model architecture, also known as a dense neural network or fully connected neural network. Put differently, there is no layer skip in the information flow, both forward and backward. The layers in this architecture are highly connected to one another.

```
# Define Dense model architecture
    dense model = Sequential()
    dense model.add(Embedding(vocab size, embedding dim, input length=max len))
    dense model.add(GlobalAveragePooling1D())
    dense model.add(Dense(n dense, activation='relu'))
    dense model.add(Dropout(drop value))
    dense model.add(Dense(1, activation='sigmoid'))
    dense model.summary()
    # Compile the model
    dense model.compile(loss='binary crossentropy', optimizer='adam', metrics=['accuracy'])
    # Fitting the dense spam detector model
    num epochs = 50
    early stop = EarlyStopping(monitor='val loss', patience=3)
    history = dense model.fit(training padded, train labels, epochs=num epochs, validation data=(testing padded, test labels), callbacks=[early stop], verbose=2)
    # Save the model
    dense model.save('Dense Spam Detection.h5')
```

In a dense model, weighted connections and activation functions are used by later hidden layers to process the input layer's received data. The output layer generates the model's result. The dense connections between neurons in neighboring layers are what give rise to the word "dense." The weights corresponding to every connection and the biases of every neuron are among the

parameters of the model. During training, these weights and biases are modified using optimization algorithms such as gradient descent to minimize the discrepancy between the target values and the predicted output.

Regression, feature learning, and classification are just a few of the applications for which dense architectures are useful. On complicated tasks, though, they might become computationally costly and prone to overfitting, which would prompt the creation of more advanced architectures like convolutional neural networks (CNNs) and recurrent neural networks (RNNs) for uses.

```
Model: "sequential"
                                                     Param #
Layer (type)
                            Output Shape
 embedding (Embedding)
                            (None, 50, 16)
                                                     8000
 global_average_pooling1d ( (None, 16)
                                                     0
 GlobalAveragePooling1D)
 dense (Dense)
                            (None, 24)
                                                     408
 dropout (Dropout)
                            (None, 24)
                                                     0
 dense 1 (Dense)
                            (None, 1)
                                                     25
______
Total params: 8433 (32.94 KB)
Trainable params: 8433 (32.94 KB)
Non-trainable params: 0 (0.00 Byte)
Epoch 1/50
38/38 - 4s - loss: 0.6833 - accuracy: 0.5958 - val loss: 0.6710 - val accuracy: 0.6890 - 4s/epoch - 117ms/step
Fnoch 2/50
38/38 - 0s - loss: 0.6497 - accuracy: 0.7481 - val loss: 0.6211 - val accuracy: 0.8328 - 403ms/epoch - 11ms/step
Epoch 3/50
38/38 - 0s - loss: 0.5844 - accuracy: 0.8510 - val loss: 0.5446 - val accuracy: 0.8562 - 331ms/epoch - 9ms/step
Epoch 4/50
38/38 - 0s - loss: 0.5010 - accuracy: 0.8711 - val loss: 0.4588 - val accuracy: 0.8896 - 224ms/epoch - 6ms/step
Epoch 5/50
38/38 - 0s - loss: 0.4189 - accuracy: 0.8837 - val_loss: 0.3820 - val_accuracy: 0.9064 - 320ms/epoch - 8ms/step
Epoch 6/50
38/38 - 0s - loss: 0.3556 - accuracy: 0.8996 - val loss: 0.3193 - val accuracy: 0.9197 - 259ms/epoch - 7ms/step
Epoch 7/50
38/38 - 0s - loss: 0.2988 - accuracy: 0.9155 - val loss: 0.2728 - val accuracy: 0.9264 - 217ms/epoch - 6ms/step
38/38 - 0s - loss: 0.2491 - accuracy: 0.9272 - val loss: 0.2385 - val accuracy: 0.9365 - 219ms/epoch - 6ms/step
```

An embedding layer for word mapping and a global average pooling 1D layer for reducing spatial dimensions define the Dense model for SMS spam detection. There are two dense layers: a dropout layer for regularization and a first layer with 24 units and ReLU activation. Binary predictions are generated by the final Dense layer, which has a sigmoid activation. The accuracy metric, Adam optimizer, and binary cross-entropy loss are used to compile the model. Training lasts 50 epochs, with an early termination occurring after 3 epochs in which the validation loss does not improve. The trained model is stored for later use, and the model successfully detects spam. Optimization can be achieved by experimenting with model architecture and hyperparameter adjustments.

#### **Evaluation:**

For loss and accuracy, all Dense, LSTM, and Bi-LSTM models are similar. These three models have validation losses of 0.14, 0.16, and 0.18, respectively. Furthermore, the respective validation accuracy rates are 95.99%, 95.32%, and 94%.

As a final model, we decide to classify the text messages as spam or ham using Dense architecture. Compared to LSTM and Bi-LSTM, the accuracy and loss over epochs are more stable in the dense classifier's straightforward structure.

### **Predictions:**

## First scenario: Using our data's raw text.

The third message below is spam, while the first and second are ham. To turn them into sequences, we employed the same tokenizer that we had previously written in the code. By doing this, the training sets token and the new words' token are guaranteed to match. After tokenization, we apply padding in the same way as previously and supply the same dimension as the training set.



The model correctly predicts that the third sentence is spam, while the first two sentences are not. The third sentence has a 99% probability of being spam.

## Scenario 2: Use newly created text messages to see how the model classifies them.

The first sentence below is more like spam, while the remaining two sentences are more like ham.

According to our model, the first message has a 78% probability of being spam, while the remaining messages are classified as ham.

#### **REFERENCES:**

- 1. "Review spam detection," N. Jindal and B. Liu, Proceedings of the 16th International Conference on World Wide Web, 2007, pp. 1189–1190.
- 2. "Towards SMS spam filtering: Results under a new dataset," T. A. Almeida, J. M. G. Hidalgo, and T. P. Silva, International Journal of Information Security Science, vol. 2, No. 1, pp. 1-18, 2013.
- 3. "Text normalization and semantic indexing to enhance instant messaging and SMS spam filtering," Knowledge-Based Systems, vol. 108, pp. 25–32, 2016, T. A. Almeida, T. P. Silva, I. Santos, and J. M. G. Hidalgo.

- 4. J. Schmidhuber, "An overview of deep learning in neural networks," Neural Networks, vol. 61, 2015, pp. 85–117.
- 5. "Deep learning," Y. LeCun, Y. Bengio, and G. Hinton, Nature, vol. 521, no. 7553, pp. 436-444,