# Ungraded Lab: Training a Sarcasm Detection Model using Bidirectional LSTMs

In this lab, you will revisit the News Headlines Dataset for Sarcasm Detection dataset and use it to train a Bi-LSTM Model.

## Imports

```
In [1]:  import json
         import matplotlib.pyplot as plt
         import tensorflow as tf
```

## Load the Dataset

First, you will download the JSON file and extract the contents into lists.

```
In [2]:  # The dataset is already downloaded for you. For downloading you can use the code below.
         # !wget https://storage.googleapis.com/tensorflow-1-public/course3/sarcasm.json
```

```
In [3]:  # Load the JSON file
         with open("./sarcasm.json", 'r') as f:
             datastore = json.load(f)

         # Initialize the lists
         sentences = []
         labels = []

         # Collect sentences and labels into the lists
         for item in datastore:
             sentences.append(item['headline'])
             labels.append(item['is_sarcastic'])
```

## Parameters

We placed the constant parameters in the cell below so you can easily tweak it later:

```
In [4]:  # Number of examples to use for training
         TRAINING_SIZE = 20000

         # Vocabulary size of the tokenizer
         VOCAB_SIZE = 10000

         # Maximum length of the padded sequences
         MAX_LENGTH = 32

         # Type of padding
         PADDING_TYPE = 'pre'

         # Specifies how to truncate the sequences
         TRUNC_TYPE = 'post'
```

## Split the Dataset

You will then split the lists into train and test sets.

```
In [5]:  # Split the sentences
         train_sentences = sentences[0:TRAINING_SIZE]
         test_sentences = sentences[TRAINING_SIZE:]

         # Split the labels
         train_labels = labels[0:TRAINING_SIZE]
         test_labels = labels[TRAINING_SIZE:]
```

## Data preprocessing

Next, you will generate the vocabulary and padded sequences.

```
In [6]:  # Instantiate the vectorization layer
         vectorize_layer = tf.keras.layers.TextVectorization(max_tokens=VOCAB_SIZE)

         # Generate the vocabulary based on the training inputs
         vectorize_layer.adapt(train_sentences)
```

You will combine the sentences and labels, then put them in a `tf.data.Dataset`. This will let you leverage the `tf.data` pipeline methods you've been using to preprocess the dataset.

```
In [7]:  # Put the sentences and labels in a tf.data.Dataset
         train_dataset = tf.data.Dataset.from_tensor_slices((train_sentences,train_labels))
         test_dataset = tf.data.Dataset.from_tensor_slices((test_sentences,test_labels))
```

You will use the same preprocessing function from the previous lab to generate the padded sequences.

```
In [8]:  def preprocessing_fn(dataset):
             '''Generates padded sequences from a tf.data.Dataset'''

             # Apply the vectorization layer to the string features
             dataset_sequences = dataset.map(
                 lambda text, label: (vectorize_layer(text), label)
                 )

             # Put all elements in a single ragged batch
             dataset_sequences = dataset_sequences.ragged_batch(
                 batch_size=dataset_sequences.cardinality()
                 )

             # Output a tensor from the single batch. Extract the sequences and labels.
             sequences, labels = dataset_sequences.get_single_element()

             # Pad the sequences
             padded_sequences = tf.keras.utils.pad_sequences(
                 sequences.numpy(),
                 maxlen=MAX_LENGTH,
                 truncating=TRUNC_TYPE,
                 padding=PADDING_TYPE
                 )

             # Convert back to a tf.data.Dataset
             padded_sequences = tf.data.Dataset.from_tensor_slices(padded_sequences)
             labels = tf.data.Dataset.from_tensor_slices(labels)

             # Combine the padded sequences and labels
             dataset_vectorized = tf.data.Dataset.zip(padded_sequences, labels)

             return dataset_vectorized
```

```
In [9]:  # Preprocess the train and test data
         train_dataset_vectorized = train_dataset.apply(preprocessing_fn)
         test_dataset_vectorized = test_dataset.apply(preprocessing_fn)
```

It's always good to check a few examples to see if the transformation works as expected.

```
In [10]: # View 2 training sequences and its labels
         for example in train_dataset_vectorized.take(2):
             print(example)
             print()
```

```
(<tf.Tensor: shape=(32,), dtype=int32, numpy=
array([   0,    0,    0,    0,    0,    0,    0,    0,    0,    0,    0,
          0,    0,    0,    0,    0,    0,    0,    0,  319,    1,
        943, 4079, 2366,   47,  366,   94, 2026,    6, 2653, 9470],
      dtype=int32)>, <tf.Tensor: shape=(), dtype=int32, numpy=0>)

(<tf.Tensor: shape=(32,), dtype=int32, numpy=
array([   0,    0,    0,    0,    0,    0,    0,    0,    0,    0,    0,
          0,    0,    0,    0,    0,    0,    4, 7185, 3128, 3305,
         28,    2,  152,    1,  358, 2902,    6,  236,    9,  844],
      dtype=int32)>, <tf.Tensor: shape=(), dtype=int32, numpy=0>)
```

Then, you will optimize and batch the dataset.

```
In [11]: SHUFFLE_BUFFER_SIZE = 1000
         PREFETCH_BUFFER_SIZE = tf.data.AUTOTUNE
         BATCH_SIZE = 32
```

```
# Optimize and batch the datasets for training
train_dataset_final = (train_dataset_vectorized
                        .cache()
                        .shuffle(SHUFFLE_BUFFER_SIZE)
                        .prefetch(PREFETCH_BUFFER_SIZE)
                        .batch(BATCH_SIZE)
                        )

test_dataset_final = (test_dataset_vectorized
                       .cache()
                       .prefetch(PREFETCH_BUFFER_SIZE)
                       .batch(BATCH_SIZE)
                       )
```

## Plot Utility

In [12]:
```python
def plot_loss_acc(history):
    '''Plots the training and validation loss and accuracy from a history object'''
    acc = history.history['accuracy']
    val_acc = history.history['val_accuracy']
    loss = history.history['loss']
    val_loss = history.history['val_loss']

    epochs = range(len(acc))

    fig, ax = plt.subplots(1,2, figsize=(12, 6))
    ax[0].plot(epochs, acc, 'bo', label='Training accuracy')
    ax[0].plot(epochs, val_acc, 'b', label='Validation accuracy')
    ax[0].set_title('Training and validation accuracy')
    ax[0].set_xlabel('epochs')
    ax[0].set_ylabel('accuracy')
    ax[0].legend()

    ax[1].plot(epochs, loss, 'bo', label='Training Loss')
    ax[1].plot(epochs, val_loss, 'b', label='Validation Loss')
    ax[1].set_title('Training and validation loss')
    ax[1].set_xlabel('epochs')
    ax[1].set_ylabel('loss')
    ax[1].legend()

    plt.show()
```

## Build and Compile the Model

The architecture here is almost identical to the one you used in the previous lab with the IMDB Reviews. Try to tweak the parameters and see how it affects the training time and accuracy (both training and validation).

In [13]:
```python
# Parameters
EMBEDDING_DIM = 16
LSTM_DIM = 32
DENSE_DIM = 24

# Model Definition with LSTM
model_lstm = tf.keras.Sequential([
    tf.keras.Input(shape=(MAX_LENGTH,)),
    tf.keras.layers.Embedding(input_dim=VOCAB_SIZE, output_dim=EMBEDDING_DIM),
    tf.keras.layers.Bidirectional(tf.keras.layers.LSTM(LSTM_DIM)),
    tf.keras.layers.Dense(DENSE_DIM, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

# Set the training parameters
model_lstm.compile(loss='binary_crossentropy',optimizer='adam',metrics=['accuracy'])

# Print the model summary
model_lstm.summary()
```
**Model: "sequential"**

| Layer (type) | Output Shape | Param # |
|---|---|---|
| embedding (Embedding) | (None, 32, 16) | 160,000 |
| bidirectional (Bidirectional) | (None, 64) | 12,544 |
| dense (Dense) | (None, 24) | 1,560 |
| dense_1 (Dense) | (None, 1) | 25 |

**Total params:** 174,129 (680.19 KB)
**Trainable params:** 174,129 (680.19 KB)
**Non-trainable params:** 0 (0.00 B)
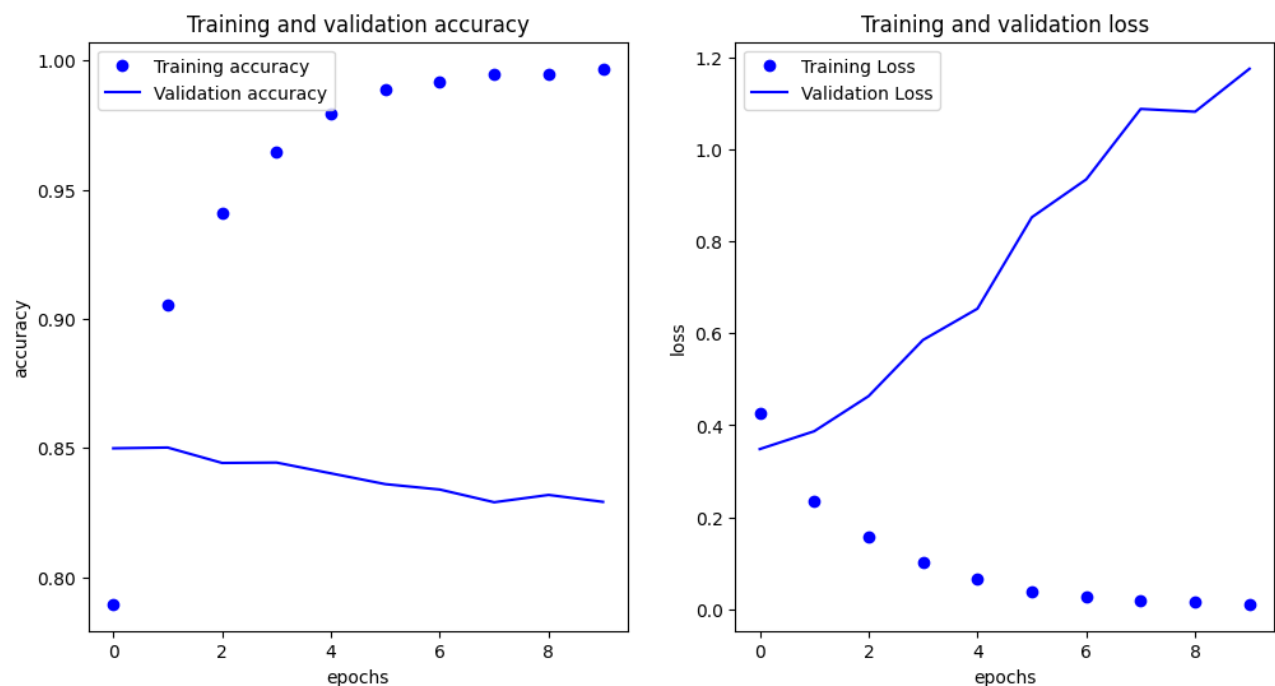
## Train the Model

```
In [14]: NUM_EPOCHS = 10

         # Train the model
         history_lstm = model_lstm.fit(train_dataset_final, epochs=NUM_EPOCHS, validation_data=test_dataset_final)
```

```
Epoch 1/10
625/625 ──────────────── 7s 6ms/step - accuracy: 0.7031 - loss: 0.5345 - val_accuracy: 0.8499 - val_loss: 0.3483
Epoch 2/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.8959 - loss: 0.2573 - val_accuracy: 0.8502 - val_loss: 0.3871
Epoch 3/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9344 - loss: 0.1768 - val_accuracy: 0.8442 - val_loss: 0.4635
Epoch 4/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9601 - loss: 0.1185 - val_accuracy: 0.8444 - val_loss: 0.5856
Epoch 5/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9763 - loss: 0.0756 - val_accuracy: 0.8402 - val_loss: 0.6536
Epoch 6/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9862 - loss: 0.0493 - val_accuracy: 0.8360 - val_loss: 0.8525
Epoch 7/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9893 - loss: 0.0356 - val_accuracy: 0.8340 - val_loss: 0.9349
Epoch 8/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9942 - loss: 0.0202 - val_accuracy: 0.8290 - val_loss: 1.0880
Epoch 9/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9943 - loss: 0.0205 - val_accuracy: 0.8319 - val_loss: 1.0820
Epoch 10/10
625/625 ──────────────── 4s 6ms/step - accuracy: 0.9964 - loss: 0.0122 - val_accuracy: 0.8292 - val_loss: 1.1754
```

```
In [15]: # Plot the accuracy and loss
         plot_loss_acc(history_lstm)
```



## Wrap Up

This concludes this lab on using LSTMs for the Sarcasm dataset. You will explore another architecture in the next lab. Before doing so, run the cell below to free up resources.

```
In [16]:  # Shutdown the kernel to free up resources.
          # Note: You can expect a pop-up when you run this cell. You can safely ignore that and just press `Ok`.

          from IPython import get_ipython

          k = get_ipython().kernel

          k.do_shutdown(restart=False)
```

```
Out[16]:  {'status': 'ok', 'restart': False}
```