

```
In [2]: #Team Members  
#Ayan Agrawal - aka398  
#Joshua Abraham - jma672
```

```
In [3]: #Submitted By: Ayan Agrawal - aka398
```

Sentiment classification with movie reviews

This notebook classifies movie reviews as *positive* or *negative* using the text of the review. This is an example of *binary*—or two-class—classification, an important and widely applicable kind of machine learning problem.

We'll use the [IMDB dataset](https://www.tensorflow.org/api_docs/python/tf/keras/datasets/imdb) (https://www.tensorflow.org/api_docs/python/tf/keras/datasets/imdb) that contains the text of 50,000 movie reviews from the [Internet Movie Database](https://www.imdb.com/) (<https://www.imdb.com/>). These are split into 25,000 reviews for training and 25,000 reviews for testing. The training and testing sets are *balanced*, meaning they contain an equal number of positive and negative reviews.

This notebook uses [tf.keras](https://www.tensorflow.org/guide/keras) (<https://www.tensorflow.org/guide/keras>), a high-level API to build and train models in TensorFlow. For a more advanced text classification tutorial using `tf.keras`, see the [MLCC Text Classification Guide](https://developers.google.com/machine-learning/guides/text-classification/) (<https://developers.google.com/machine-learning/guides/text-classification/>).

```
In [1]: import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras import regularizers  
  
import numpy as np  
  
print(tf.__version__)  
  
1.11.0
```

```
In [2]: #Setting up tensoBoard
!wget https://bin.equinox.io/c/4VmDzA7iaHb/ngrok-stable-linux-amd64.zip
!unzip ngrok-stable-linux-amd64.zip

--2018-10-11 19:35:27-- https://bin.equinox.io/c/4VmDzA7iaHb/ngrok-stable-linux-amd64.zip
Resolving bin.equinox.io (bin.equinox.io)... 34.232.181.106, 34.231.75.48, 34.235.97.255, ...
Connecting to bin.equinox.io (bin.equinox.io)|34.232.181.106|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 5363700 (5.1M) [application/octet-stream]
Saving to: 'ngrok-stable-linux-amd64.zip'

100%[=====>] 5,363,700  28.0MB/s  in 0.2s

2018-10-11 19:35:28 (28.0 MB/s) - 'ngrok-stable-linux-amd64.zip' saved [5363700/5363700]

Archive:  ngrok-stable-linux-amd64.zip
  inflating: ngrok
```

```
In [22]: LOG_DIR = './log'
get_ipython().system_raw(
    'tensorboard --logdir {} --host 0.0.0.0 --port 6006 &'
    .format(LOG_DIR)
)

get_ipython().system_raw('./ngrok http 6006 &')

! curl -s http://localhost:4040/api/tunnels | python3 -c '\
    "import sys, json; print(json.load(sys.stdin)[\'tunnels\'][0][\'public_url\'])"'

http://428f32bc.ngrok.io
```

Download the IMDB dataset

The IMDB dataset comes packaged with TensorFlow. It has already been preprocessed such that the reviews (sequences of words) have been converted to sequences of integers, where each integer represents a specific word in a dictionary.

The following code downloads the IMDB dataset to your machine (or uses a cached copy if you've already downloaded it):

```
In [4]: imdb = keras.datasets.imdb

(train_data, train_labels), (test_data, test_labels) = imdb.load_data(num_words=10000)
```

```
In [5]: set(list(train_labels))
```

```
Out[5]: {0, 1}
```

The argument `num_words=10000` keeps the top 10,000 most frequently occurring words in the training data. The rare words are discarded to keep the size of the data manageable.

Explore the data

Let's take a moment to understand the format of the data. The dataset comes preprocessed: each example is an array of integers representing the words of the movie review. Each label is an integer value of either 0 or 1, where 0 is a negative review, and 1 is a positive review.

```
In [6]: print("Training entries: {}, labels: {}".format(len(train_data), len(train_labels)))
```

```
Training entries: 25000, labels: 25000
```

The text of reviews have been converted to integers, where each integer represents a specific word in a dictionary. Here's what the first review looks like:

```
In [7]: print(train_data[0])
```

```
[1, 14, 22, 16, 43, 530, 973, 1622, 1385, 65, 458, 4468, 66, 3941, 4, 173, 3
6, 256, 5, 25, 100, 43, 838, 112, 50, 670, 2, 9, 35, 480, 284, 5, 150, 4, 17
2, 112, 167, 2, 336, 385, 39, 4, 172, 4536, 1111, 17, 546, 38, 13, 447, 4, 19
2, 50, 16, 6, 147, 2025, 19, 14, 22, 4, 1920, 4613, 469, 4, 22, 71, 87, 12, 1
6, 43, 530, 38, 76, 15, 13, 1247, 4, 22, 17, 515, 17, 12, 16, 626, 18, 2, 5,
62, 386, 12, 8, 316, 8, 106, 5, 4, 2223, 5244, 16, 480, 66, 3785, 33, 4, 130,
12, 16, 38, 619, 5, 25, 124, 51, 36, 135, 48, 25, 1415, 33, 6, 22, 12, 215, 2
8, 77, 52, 5, 14, 407, 16, 82, 2, 8, 4, 107, 117, 5952, 15, 256, 4, 2, 7, 376
6, 5, 723, 36, 71, 43, 530, 476, 26, 400, 317, 46, 7, 4, 2, 1029, 13, 104, 8
8, 4, 381, 15, 297, 98, 32, 2071, 56, 26, 141, 6, 194, 7486, 18, 4, 226, 22,
21, 134, 476, 26, 480, 5, 144, 30, 5535, 18, 51, 36, 28, 224, 92, 25, 104, 4,
226, 65, 16, 38, 1334, 88, 12, 16, 283, 5, 16, 4472, 113, 103, 32, 15, 16, 53
45, 19, 178, 32]
```

Movie reviews may be different lengths. The below code shows the number of words in the first and second reviews. Since inputs to a neural network must be the same length, we'll need to resolve this later.

```
In [8]: len(train_data[0]), len(train_data[1])
# print(len(train_data[0]))
```

```
Out[8]: (218, 189)
```

Convert the integers back to words

It may be useful to know how to convert integers back to text. Here, we'll create a helper function to query a dictionary object that contains the integer to string mapping:

```
In [9]: # A dictionary mapping words to an integer index
word_index = imdb.get_word_index()

# The first indices are reserved
word_index = {k:(v+3) for k,v in word_index.items()}
word_index["<PAD>"] = 0
word_index["<START>"] = 1
word_index["<UNK>"] = 2 # unknown
word_index["<UNUSED>"] = 3

reverse_word_index = dict([(value, key) for (key, value) in word_index.items()])

def decode_review(text):
    return ' '.join([reverse_word_index.get(i, '?') for i in text])
```

Now we can use the decode_review function to display the text for the first review:

```
In [10]: decode_review(train_data[0])
```

```
Out[10]: "<START> this film was just brilliant casting location scenery story direction everyone's really suited the part they played and you could just imagine being there robert <UNK> is an amazing actor and now the same being director <UNK> father came from the same scottish island as myself so i loved the fact there was a real connection with this film the witty remarks throughout the film were great it was just brilliant so much that i bought the film as soon as it was released for <UNK> and would recommend it to everyone to watch and the fly fishing was amazing really cried at the end it was so sad and you know what they say if you cry at a film it must have been good and this definitely was also <UNK> to the two little boy's that played the <UNK> of norman and paul they were just brilliant children are often left out of the <UNK> list i think because the stars that play them all grown up are such a big profile for the whole film but these children are amazing and should be praised for what they have done don't you think the whole story was so lovely because it was true and was someone's life after all that was shared with us all"
```

```
In [11]: decode_review(train_data[2])
```

```
Out[11]: "<START> this has to be one of the worst films of the 1990s when my friends i
were watching this film being the target audience it was aimed at we just sat
watched the first half an hour with our jaws touching the floor at how bad it
really was the rest of the time everyone else in the theatre just started tal
king to each other leaving or generally crying into their popcorn that they a
ctually paid money they had <UNK> working to watch this feeble excuse for a f
ilm it must have looked like a great idea on paper but on film it looks like
no one in the film has a clue what is going on crap acting crap costumes i ca
n't get across how <UNK> this is to watch save yourself an hour a bit of your
life"
```

Prepare the data

The reviews—the arrays of integers—must be converted to tensors before fed into the neural network. This conversion can be done a couple of ways:

- One-hot-encode the arrays to convert them into vectors of 0s and 1s. For example, the sequence [3, 5] would become a 10,000-dimensional vector that is all zeros except for indices 3 and 5, which are ones. Then, make this the first layer in our network—a Dense layer—that can handle floating point vector data. This approach is memory intensive, though, requiring a `num_words * num_reviews` size matrix.
- Alternatively, we can pad the arrays so they all have the same length, then create an integer tensor of shape `max_length * num_reviews`. We can use an embedding layer capable of handling this shape as the first layer in our network.

In this assignment, we will use the second approach.

Since the movie reviews must be the same length, we will use the [pad_sequences](https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences) (https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences) function to standardize the lengths:

```
In [12]: print(len(train_data[0]))
num_tokens = [len(words) for words in train_data + test_data]
num_tokens = np.array(num_tokens)
print(num_tokens)
np.mean(num_tokens)
np.max(num_tokens)
```

```
218
[286 449 744 ..., 259 249 325]
```

```
Out[12]: 2697
```

```
In [13]: max_tokens = np.mean(num_tokens) + 2 * np.std(num_tokens)
max_tokens = int(max_tokens)
max_tokens
#np.sum(num_tokens < max_tokens) / len(num_tokens)
```

```
Out[13]: 960
```

```
In [14]: train_data = keras.preprocessing.sequence.pad_sequences(train_data,
                                                                value=word_index["<PAD
>"],
                                                                padding='post',
                                                                maxlen=960)

test_data = keras.preprocessing.sequence.pad_sequences(test_data,
                                                        value=word_index["<PAD
>"],
                                                        padding='post',
                                                        maxlen=960)
```

Let's look at the length of the examples now:

```
In [15]: len(train_data[0]), len(train_data[1])
```

```
Out[15]: (960, 960)
```

And inspect the (now padded) first review:

```
In [16]: print(train_data[0])
```

8/15

0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0]

Build the model

The neural network is created by stacking layers—this requires two main architectural decisions:

- How many layers to use in the model?
- How many *hidden units* to use for each layer?

In this example, the input data consists of an array of word-indices. The labels to predict are either 0 or 1. Let's build a model for this problem:

```
In [17]: # input shape is the vocabulary count used for the movie reviews (10,000 words)
vocab_size = 10000

model = keras.Sequential()
model.add(keras.layers.Embedding(vocab_size, 126))
model.add(keras.layers.GRU(512, dropout=0.2, return_sequences = True))
model.add(keras.layers.GRU(256, return_sequences = True))
model.add(keras.layers.Dropout(0.2, seed=1))
model.add(keras.layers.GlobalAveragePooling1D())
model.add(keras.layers.Dense(64, activation=tf.nn.leaky_relu))
model.add(keras.layers.Dense(32, activation=tf.nn.leaky_relu))
model.add(keras.layers.Dense(16, activation=tf.nn.leaky_relu))
model.add(keras.layers.Dense(1, activation=tf.nn.sigmoid))

model.summary()
```

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, None, 126)	1260000

gru (GRU)	(None, None, 512)	981504

gru_1 (GRU)	(None, None, 256)	590592

dropout (Dropout)	(None, None, 256)	0

global_average_pooling1d (GlobalAveragePooling1D)	(None, 256)	0

dense (Dense)	(None, 64)	16448

dense_1 (Dense)	(None, 32)	2080

dense_2 (Dense)	(None, 16)	528

dense_3 (Dense)	(None, 1)	17
=====		
Total params: 2,851,169		
Trainable params: 2,851,169		
Non-trainable params: 0		

The layers are stacked sequentially to build the classifier:

1. The first layer is an Embedding layer. This layer takes the integer-encoded vocabulary and looks up the embedding vector for each word-index. These vectors are learned as the model trains. The vectors add a dimension to the output array. The resulting dimensions are: (batch, sequence, embedding).
2. Next, a GlobalAveragePooling1D layer returns a fixed-length output vector for each example by averaging over the sequence dimension. This allows the model can handle input of variable length, in the simplest way possible.
3. This fixed-length output vector is piped through a fully-connected (Dense) layer with 16 hidden units.
4. The last layer is densely connected with a single output node. Using the sigmoid activation function, this value is a float between 0 and 1, representing a probability, or confidence level.

Hidden units

The above model has two intermediate or "hidden" layers, between the input and output. The number of outputs (units, nodes, or neurons) is the dimension of the representational space for the layer. In other words, the amount of freedom the network is allowed when learning an internal representation.

If a model has more hidden units (a higher-dimensional representation space), and/or more layers, then the network can learn more complex representations. However, it makes the network more computationally expensive and may lead to learning unwanted patterns—patterns that improve performance on training data but not on the test data. This is called *overfitting*, and we'll explore it later.

Loss function and optimizer

A model needs a loss function and an optimizer for training. Since this is a binary classification problem and the model outputs of a probability (a single-unit layer with a sigmoid activation), we'll use the `binary_crossentropy` loss function.

This isn't the only choice for a loss function, you could, for instance, choose `mean_squared_error`. But, generally, `binary_crossentropy` is better for dealing with probabilities—it measures the "distance" between probability distributions, or in our case, between the ground-truth distribution and the predictions.

Now, configure the model to use an optimizer and a loss function:

```
In [23]: model.compile(optimizer=tf.train.AdamOptimizer(),  
                      loss='binary_crossentropy',  
                      metrics=['accuracy'])
```

Create a validation set

When training, we want to check the accuracy of the model on data it hasn't seen before. Create a *validation set* by setting apart 10,000 examples from the original training data. (Why not use the testing set now? Our goal is to develop and tune our model using only the training data, then use the test data just once to evaluate our accuracy).

```
In [24]: x_val = train_data[:5000]
        partial_x_train = train_data[5000:]

        y_val = train_labels[:5000]
        partial_y_train = train_labels[5000:]
```

```
In [25]: from tensorflow.keras.callbacks import TensorBoard
        from tensorflow.keras.callbacks import EarlyStopping

        tbCallBack = TensorBoard(log_dir='./log', histogram_freq=1,
                                write_graph=True,
                                batch_size=126)

        stopCallBack = EarlyStopping(monitor='val_loss', patience=2)
```

Train the model

Training the model for 40 epochs in mini-batches of 512 samples (Please decide the appropriate epochs and batch size for the sequence models). This is 10 iterations over all samples in the `x_train` and `y_train` tensors. While training, monitor the model's loss and accuracy on the 10,000 samples from the validation set:

```
In [26]: history = model.fit(partial_x_train,
                             partial_y_train,
                             epochs=30,
                             batch_size=126,
                             validation_data=(x_val, y_val),
                             verbose=1,
                             callbacks=[tbCallback, stopCallback])
```

Train on 20000 samples, validate on 5000 samples

Epoch 1/30

20000/20000 [=====] - 1241s 62ms/step - loss: 0.5899
- acc: 0.6423 - val_loss: 0.3045 - val_acc: 0.8798

Epoch 2/30

20000/20000 [=====] - 1249s 62ms/step - loss: 0.2543
- acc: 0.8987 - val_loss: 0.2693 - val_acc: 0.8932

Epoch 3/30

20000/20000 [=====] - 1248s 62ms/step - loss: 0.1816
- acc: 0.9323 - val_loss: 0.3197 - val_acc: 0.8838

Epoch 4/30

20000/20000 [=====] - 1250s 63ms/step - loss: 0.1347
- acc: 0.9515 - val_loss: 0.2929 - val_acc: 0.8890

Evaluate the model

And let's see how the model performs. Two values will be returned. Loss (a number which represents our error, lower values are better), and accuracy.

```
In [27]: results = model.evaluate(test_data, test_labels)

print(results)
```

25000/25000 [=====] - 498s 20ms/step
[0.32112948688983917, 0.8807599999999999]

This fairly naive approach achieves an accuracy of about 87%. With more advanced approaches, the model should get closer to 95%.

Create a graph of accuracy and loss over time

`model.fit()` returns a `History` object that contains a dictionary with everything that happened during training:

```
In [28]: history_dict = history.history
         history_dict.keys()
```

```
Out[28]: dict_keys(['val_loss', 'val_acc', 'loss', 'acc'])
```

There are four entries: one for each monitored metric during training and validation. We can use these to plot the training and validation loss for comparison, as well as the training and validation accuracy:

```
In [29]: import matplotlib.pyplot as plt

acc = history.history['acc']
val_acc = history.history['val_acc']
loss = history.history['loss']
val_loss = history.history['val_loss']

epochs = range(1, len(acc) + 1)

# "bo" is for "blue dot"
plt.plot(epochs, loss, 'bo', label='Training loss')
# b is for "solid blue line"
plt.plot(epochs, val_loss, 'b', label='Validation loss')
plt.title('Training and validation loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()

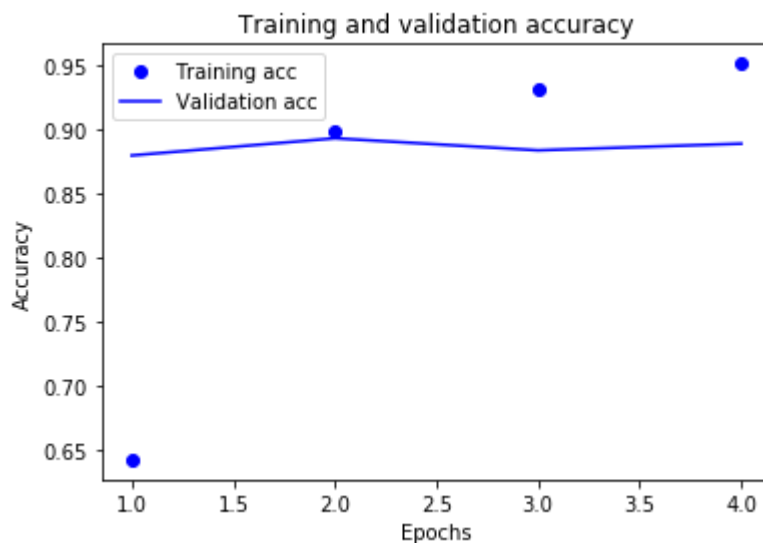
plt.show()
```

<matplotlib.figure.Figure at 0x2ae60ba44518>

```
In [30]: plt.clf() # clear figure
acc_values = history_dict['acc']
val_acc_values = history_dict['val_acc']

plt.plot(epochs, acc, 'bo', label='Training acc')
plt.plot(epochs, val_acc, 'b', label='Validation acc')
plt.title('Training and validation accuracy')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.legend()

plt.show()
```



In this plot, the dots represent the training loss and accuracy, and the solid lines are the validation loss and accuracy.

Notice the training loss *decreases* with each epoch and the training accuracy *increases* with each epoch. This is expected when using a gradient descent optimization—it should minimize the desired quantity on every iteration.

This isn't the case for the validation loss and accuracy—they seem to peak after about twenty epochs. This is an example of overfitting: the model performs better on the training data than it does on data it has never seen before. After this point, the model over-optimizes and learns representations *specific* to the training data that do not *generalize* to test data.

For this particular case, we could prevent overfitting by simply stopping the training after twenty or so epochs.