

```

audio
  "nsynth"
image
  "abstract_reasoning"
  "caltech101"
  "cats_vs_dogs"
  "celeb_a"
  "celeb_a_hq"
  "cifar10"
  "cifar100"
  "cifar10_corrupted"
  "coco2014"
  "colorectal_histology"
  "cycle_gan"
  "diabetic_retinopathy..."
  "dsprites"
  "dtd"
  "emnist"
  "fashion_mnist"
  "horses_or_humans"
  "image_label_folder"
  "imagenet2012"
  "imagenet2012_corrupted"
  "kmnist"
  "lsun"
  "mnist"
  "omniglot"
  "open_images_v4"
  "oxford_iiit_pet"
  "quickdraw_bitmap"
  "rock_paper_scissors"
  "shapes3d"
  "smallnorb"
  "sun397"
  "svhn_cropped"
  "tf_flowers"
structured
  "higgs"
  "iris"
  "titanic"
text
  "cnn_dailymail"
  "glue"
  "imdb_reviews"
  "lm1b"
  "multi_nli"
  "squad"
  "wikipedia"
  "xnli"
translate
  "flores"
  "para_crawl"
  "ted_hrlr_translate"
  "ted_multi_translate"
  "wmt15_translate"
  "wmt16_translate"
  "wmt17_translate"
  "wmt18_translate"
  "wmt19_translate"

```

There's a library called TensorFlow Data Services or TDFS for short, and that contains many data sets and lots of different categories.

<http://ai.stanford.edu/~amaas/data/sentiment/>

```

@InProceedings{maas-EtAl:2011:ACL-HLT2011,
  author    = {Maas, Andrew L. and Daly, Raymond E. and Pham, Peter T. and Huang, Dan and Ng, Andrew Y. and Potts, Christopher},
  title     = {Learning Word Vectors for Sentiment Analysis},
  booktitle = {Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies},
  month     = {June},
  year      = {2011},
  address   = {Portland, Oregon, USA},
  publisher = {Association for Computational Linguistics},
  pages     = {142--150},
  url       = {http://www.aclweb.org/anthology/P11-1015}
}

```

```

import tensorflow as tf
print(tf.__version__)

```

```
tf.enable_eager_execution()
```

```
!pip install -q tensorflow-datasets
```

Installs TensorFlow datasets (If you're using Google Colab, then you should have TensorFlow datasets already installed)

```

import tensorflow_datasets as tfds
imdb, info = tfds.load("imdb_reviews", with_info=True, as_supervised=True)

```

```
import numpy as np
```

```
train_data, test_data = imdb['train'], imdb['test']
```

Split

```
training_sentences = []
training_labels = []

testing_sentences = []
testing_labels = []
```

Define the lists containing the sentences and labels for both training and testing data.

```
# str(s.tonumpy()) is needed in Python3 instead of just s.numpy()
for s,l in train_data:
    training_sentences.append(str(s.numpy()))
    training_labels.append(l.numpy())

for s,l in test_data:
    testing_sentences.append(str(s.numpy()))
    testing_labels.append(l.numpy())
```

```
training_sentences = []
training_labels = []

testing_sentences = []
testing_labels = []
```

```
# str(s.tonumpy()) is needed in Python3 instead of just s.numpy()
```

```
for s,l in train_data:
    training_sentences.append(str(s.numpy()))
    training_labels.append(l.numpy())
```

Iterate over training data extracting the sentences and the labels. The values for S and l are tensors, so by calling their NumPy method, I'll actually extract their value.

```
for s,l in test_data:
    testing_sentences.append(str(s.numpy()))
    testing_labels.append(l.numpy())
```

```
training_sentences = []
training_labels = []
```

```
testing_sentences = []
testing_labels = []
```

```
# str(s.tonumpy()) is needed in Python3 instead of just s.numpy()
```

```
for s,l in train_data:
    training_sentences.append(str(s.numpy()))
    training_labels.append(l.numpy())
```

```
for s,l in test_data:
    testing_sentences.append(str(s.numpy()))
    testing_labels.append(l.numpy())
```

Do the same for the test set

```
tf.Tensor(b"As a lifelong fan of Dickens, I have invariably been disappointed  
by adaptations of his novels.<br /><br />Although his works presented an  
extremely accurate re-telling of human life at every level in Victorian Britain,  
throughout them all was a pervasive thread of humour that could be both playful  
or sarcastic as the narrative dictated. In a way, he was a literary  
caricaturist and cartoonist. He could be serious and hilarious in the same  
sentence. He pricked pride, lampooned arrogance, celebrated modesty,  
and empathised with loneliness and poverty. It may be a cliché, but  
he was a people's writer.<br /><br />And it is the comedy that is so often  
missing from his interpretations. At the time of writing, Oliver Twist  
is being dramatised in serial form on BBC television. All of the misery  
and cruelty is their, but non of the humour, irony, and savage lampoonery.",  
shape=(), dtype=string)
```

Here's an example of a review. I've truncated it to fit it on this slide, but you can see how it is stored as a tf.tensor.

```
tf.Tensor(1, shape=(), dtype=int64)  
tf.Tensor(1, shape=(), dtype=int64)  
tf.Tensor(1, shape=(), dtype=int64)  
tf.Tensor(0, shape=(), dtype=int64)  
tf.Tensor(0, shape=(), dtype=int64)  
tf.Tensor(1, shape=(), dtype=int64)
```

Here's a bunch of labels also stored as tensors. The value 1 indicates a positive review and zero a negative one.

```
training_labels_final = np.array(training_labels)  
testing_labels_final = np.array(testing_labels)
```

When training, my labels are expected to be NumPy arrays. So I'll turn the list of labels that I've just created into NumPy arrays

```
vocab_size = 10000  
embedding_dim = 16  
max_length = 120  
trunc_type='post'  
oov_tok = "<OOV>"
```

I've put the hyperparameters at the top like this for the reason that it makes it easier to change and edit them, instead of phishing through function sequences for the literals and then changing those.

```
from tensorflow.keras.preprocessing.text import Tokenizer  
from tensorflow.keras.preprocessing.sequence import pad_sequences  
  
tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)  
tokenizer.fit_on_texts(training_sentences)  
word_index = tokenizer.word_index  
sequences = tokenizer.texts_to_sequences(training_sentences)  
padded = pad_sequences(sequences,maxlen=max_length, truncating=trunc_type)  
  
testing_sequences = tokenizer.texts_to_sequences(testing_sentences)  
testing_padded = pad_sequences(testing_sequences,maxlen=max_length)
```

```
vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type='post'
oov_tok = "<OOV>"

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(training_sentences)
padded = pad_sequences(sequences,maxlen=max_length, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences,maxlen=max_length)
```

Import the tokenizer and the pad sequences.

```
vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type='post'
oov_tok = "<OOV>"

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(training_sentences)
padded = pad_sequences(sequences,maxlen=max_length, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences,maxlen=max_length)
```

Create an instance of tokenizer, giving it our vocab size and our desired out of vocabulary token

```
vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type='post'
oov_tok = "<OOV>"

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(training_sentences)
padded = pad_sequences(sequences,maxlen=max_length, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences,maxlen=max_length)
```

Fit the tokenizer on our training set of data.


```

vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type='post'
oov_tok = "<OOV>"

```

Once we have our word index, we can now replace the strings containing the words with the token value we created for them. This will be the list called sequences. As before, the sentences will have variant length. So we'll pad and or truncate the sequenced sentences until they're all the same length, determined by the maxlen parameter.

```

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(training_sentences)
padded = pad_sequences(sequences, maxlen=max_length, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences, maxlen=max_length)

```

```

vocab_size = 10000
embedding_dim = 16
max_length = 120
trunc_type='post'
oov_tok = "<OOV>"

```

Then we'll do the same for the testing sequences. Note that the word index is words that are derived from the training set, so you should expect to see a lot more out of vocabulary tokens in the test exam.

```

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences

tokenizer = Tokenizer(num_words = vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)
word_index = tokenizer.word_index
sequences = tokenizer.texts_to_sequences(training_sentences)
padded = pad_sequences(sequences, maxlen=max_length, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences, maxlen=max_length)

```

```

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(vocab_size, embedding_dim, input_length=max_length),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(6, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

```

The results of the embedding will be a 2D array with the length of the sentence and the embedding dimension for example 16 as its size. So we need to flatten it out in much the same way as we needed to flatten out our images

```

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(vocab_size, embedding_dim, input_length=max_length),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(6, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

```

```

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(vocab_size, embedding_dim, input_length=max_length),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(6, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

```

Layer (type)	Output Shape	Param #
embedding_9 (Embedding)	(None, 120, 16)	160000
flatten_3 (Flatten)	(None, 1920)	0
dense_14 (Dense)	(None, 6)	11526
dense_15 (Dense)	(None, 1)	7

Total params: 171,533
 Trainable params: 171,533
 Non-trainable params: 0

Often in natural language processing, a different layer type than a flatten is used, and this is a global average pooling 1D. The reason for this is the size of the output vector being fed into the dense. So for example, if I show the summary of the model with the flatten that we just saw, it will look like this

```

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(vocab_size, embedding_dim, input_length=max_length),
    tf.keras.layers.GlobalAveragePooling1D(),
    tf.keras.layers.Dense(6, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

```

Or alternatively, you can use a Global Average Pooling 1D like this, which averages across the vector to flatten it out.

Layer (type)	Output Shape	Param #
embedding_11 (Embedding)	(None, 120, 16)	160000
global_average_pooling1d_3 (GlobalAveragePooling1D)	(None, 16)	0
dense_16 (Dense)	(None, 6)	102
dense_17 (Dense)	(None, 1)	7

Total params: 160,109
 Trainable params: 160,109
 Non-trainable params: 0

```

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

```

```

num_epochs = 10
model.fit(padded,
          training_labels_final,
          epochs=num_epochs,
          validation_data=(testing_padded, testing_labels_final))

```

Now training is the simplest passing padded and your training labels final as your training set, specifying the number of epochs, and passing the testing padded and testing labels final as your test set.

```

Epoch 8/10
25000/25000 [=====] -
6s 256us/sample - loss: 5.2086e-04 - acc: 1.0000 - val_loss: 0.7252 - val_acc: 0.8270
Epoch 9/10
25000/25000 [=====] -
6s 222us/sample - loss: 3.0199e-04 - acc: 1.0000 - val_loss: 0.7628 - val_acc: 0.8269
Epoch 10/10
25000/25000 [=====] -
6s 224us/sample - loss: 1.7872e-04 - acc: 1.0000 - val_loss: 0.7997 - val_acc: 0.8259

```

```

e = model.layers[0]
weights = e.get_weights()[0]
print(weights.shape) # shape: (vocab_size, embedding_dim)

(10000, 16)

```

Now we need to talk about and demonstrate the embeddings, so you can visualize them like you did right back at the beginning of this lesson. We'll start by getting the results of the embeddings layer, which is layer zero. We can get the weights, and print out their shape like this. We can see that this is a 10,000 by 16 array, we have 10,000 words in our corpus, and we're working in a 16 dimensional array, so our embedding will have that shape.

```

Hello : 1
World : 2
How : 3
Are : 4
You : 5

```

To be able to plot it, we need a helper function to reverse our word index. As it currently stands, our word index has the key being the word, and the value being the token for the word. We'll need to flip this around, to look through the padded list to decode the tokens back into the words, so we've written this helper function.

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

```

1 : Hello
2 : World
3 : How
4 : Are
5 : You

```

```
import io
```

```

out_v = io.open('vecs.tsv', 'w', encoding='utf-8')
out_m = io.open('meta.tsv', 'w', encoding='utf-8')
for word_num in range(1, vocab_size):
    word = reverse_word_index[word_num]
    embeddings = weights[word_num]
    out_m.write(word + "\n")
    out_v.write('\t'.join([str(x) for x in embeddings]) + "\n")
out_v.close()
out_m.close()

```

Now it's time to write the vectors and their metadata auto files. The TensorFlow Projector reads this file type and uses it to plot the vectors in 3D space so we can visualize them. To the vectors file, we simply write out the value of each of the items in the array of embeddings, i.e., the co-efficient of each dimension on the vector for this word.

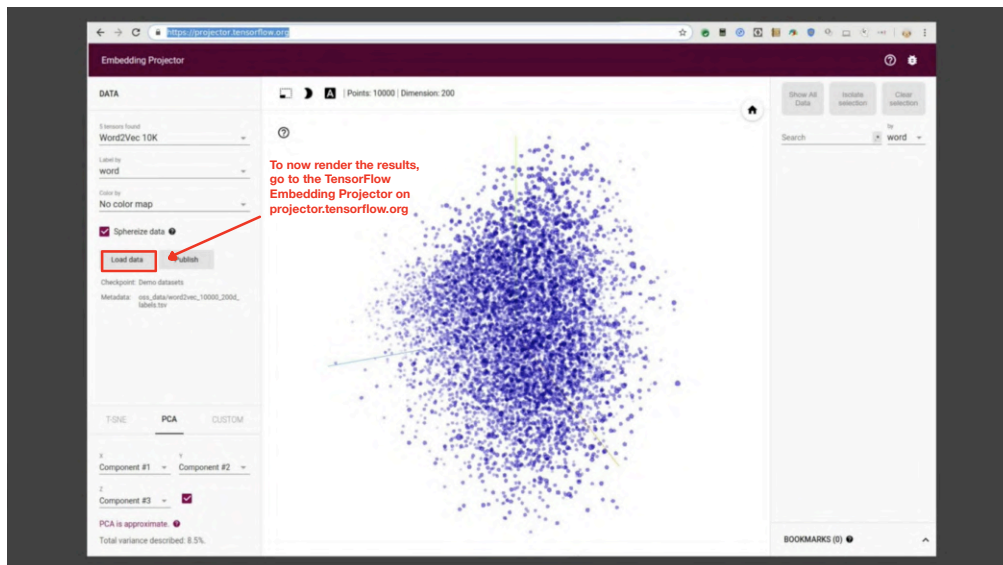
```
import io
```

```

out_v = io.open('vecs.tsv', 'w', encoding='utf-8')
out_m = io.open('meta.tsv', 'w', encoding='utf-8')
for word_num in range(1, vocab_size):
    word = reverse_word_index[word_num]
    embeddings = weights[word_num]
    out_m.write(word + "\n")
    out_v.write('\t'.join([str(x) for x in embeddings]) + "\n")
out_v.close()
out_m.close()

```

To the metadata array, we just write out the words.



Load data from your computer

Step 1: Load a TSV file of vectors.
Example of 3 vectors with dimension 4:

```
0.1\t0.2\t0.5\t0.9  
0.2\t0.1\t5.0\t0.2  
0.4\t0.1\t7.0\t0.8
```

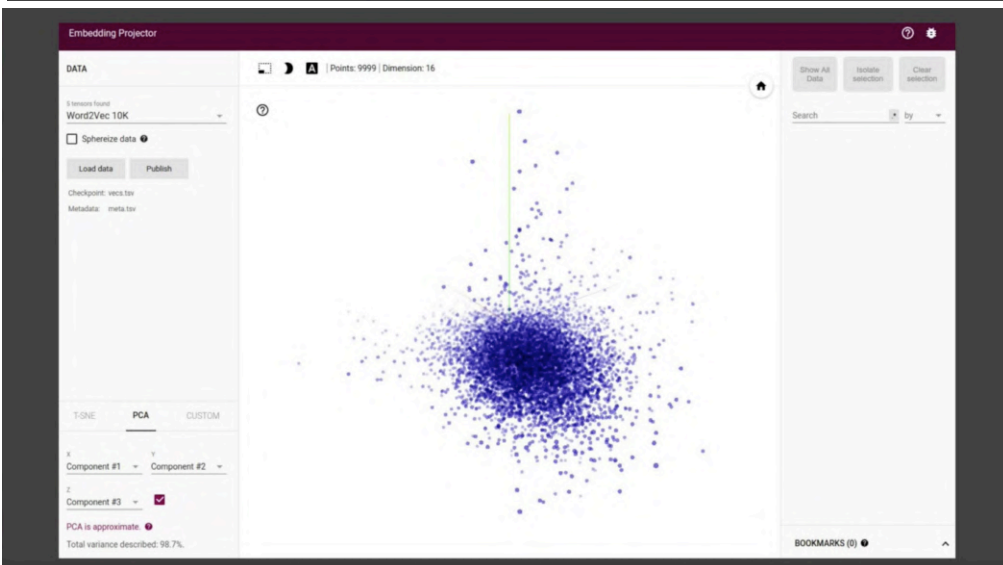
Choose file

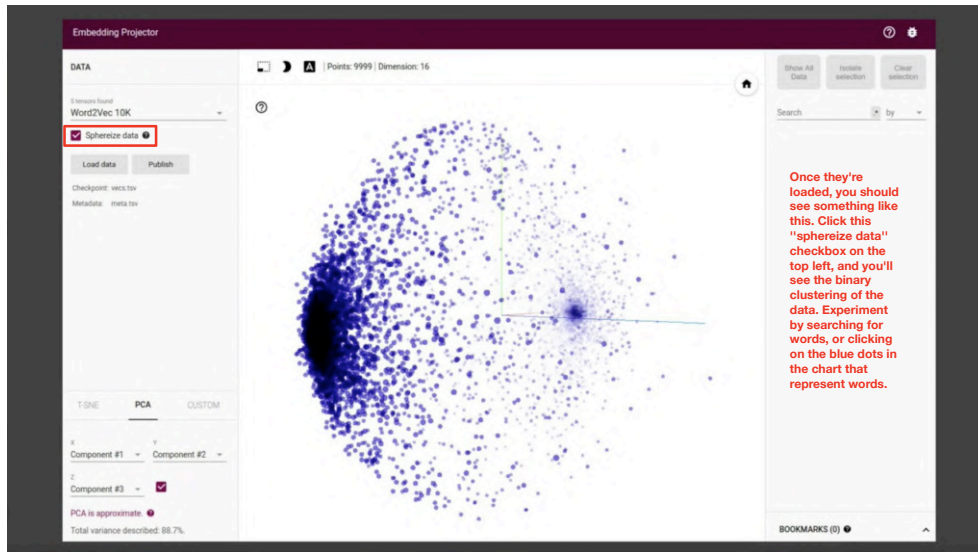
Step 2 (optional): Load a TSV file of metadata.
Example of 3 data points and 2 columns.
Note: If there is more than one column, the first row will be parsed as column labels.

```
Pokémon\tSpecies  
Wartortle\tTurtle  
Venusaur\tSeed  
Charmeleon\tFlame
```

Choose file

Click outside to dismiss.





```
import json
import tensorflow as tf

from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
```

```
vocab_size = 10000
embedding_dim = 16
max_length = 32
trunc_type='post'
padding_type='post'
oov_tok = "<OOV>"
training_size = 20000
```

```
!wget --no-check-certificate \
  https://storage.googleapis.com/laurencemoroney-blog.appspot.com/sarcasm.json \
  -O /tmp/sarcasm.json
```

The sarcasm data is stored at this URL, so you can download it to /tmp/sarcasm.json with this code.

```
with open("/tmp/sarcasm.json", 'r') as f:
    datastore = json.load(f)
```

```
sentences = []
labels = []
```

```
for item in datastore:
    sentences.append(item['headline'])
    labels.append(item['is_sarcastic'])
```

Now that you have the data set, you can open it and load it as an iterable with this code. You can create an array for sentences, and another for labels, and then iterate through the datastore, loading each headline as a sentence, and each is_sarcastic field, as your label.

```
training_sentences = sentences[0:training_size]
testing_sentences = sentences[training_size:]
training_labels = labels[0:training_size]
testing_labels = labels[training_size:]
```

To get the training set, you take array items from zero to the training size

```
training_sentences = sentences[0:training_size]
testing_sentences = sentences[training_size:]
training_labels = labels[0:training_size]
testing_labels = labels[training_size:]
```

To get the testing set, you can go from training size to the end of the array with code like this.

```

training_sentences = sentences[0:training_size]
testing_sentences = sentences[training_size:]
training_labels = labels[0:training_size]
testing_labels = labels[training_size:]

```

To get the training and testing labels, you'll use similar codes to slice the labels array.

```

tokenizer = Tokenizer(num_words=vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)

```

Start with a tokenizer, passing it the number of words you want to tokenize on and the desired out of vocabulary token

Now that we have training and test sets of sequences and labels, it's time to sequence them. To pad those sequences, you'll do that with this code.

```

word_index = tokenizer.word_index

training_sequences = tokenizer.texts_to_sequences(training_sentences)
training_padded = pad_sequences(training_sequences, maxlen=max_length,
                                padding=padding_type, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences, maxlen=max_length,
                               padding=padding_type, truncating=trunc_type)

```

```

tokenizer = Tokenizer(num_words=vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)

```

Then fit that on the training set by calling fit on texts, passing it the training sentences array.

```

word_index = tokenizer.word_index

training_sequences = tokenizer.texts_to_sequences(training_sentences)
training_padded = pad_sequences(training_sequences, maxlen=max_length,
                                padding=padding_type, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences, maxlen=max_length,
                               padding=padding_type, truncating=trunc_type)

```

```

tokenizer = Tokenizer(num_words=vocab_size, oov_token=oov_tok)
tokenizer.fit_on_texts(training_sentences)

```

Then you can use text to sequences to create the training sequence, replacing the words with their tokens.

```

word_index = tokenizer.word_index

training_sequences = tokenizer.texts_to_sequences(training_sentences)
training_padded = pad_sequences(training_sequences, maxlen=max_length,
                                padding=padding_type, truncating=trunc_type)

testing_sequences = tokenizer.texts_to_sequences(testing_sentences)
testing_padded = pad_sequences(testing_sequences, maxlen=max_length,
                               padding=padding_type, truncating=trunc_type)

```

Then you can pad the training sequences to the desired length or truncate if they're too long

```

model = tf.keras.Sequential([
    tf.keras.layers.Embedding(vocab_size, embedding_dim, input_length=max_length),
    tf.keras.layers.GlobalAveragePooling1D(),
    tf.keras.layers.Dense(24, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])

```

```
model.summary()
```

Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, 32, 16)	160000
global_average_pooling1d_2 ((None, 16)		0
dense_4 (Dense)	(None, 24)	408
dense_5 (Dense)	(None, 1)	25
Total params: 160,433		
Trainable params: 160,433		
Non-trainable params: 0		

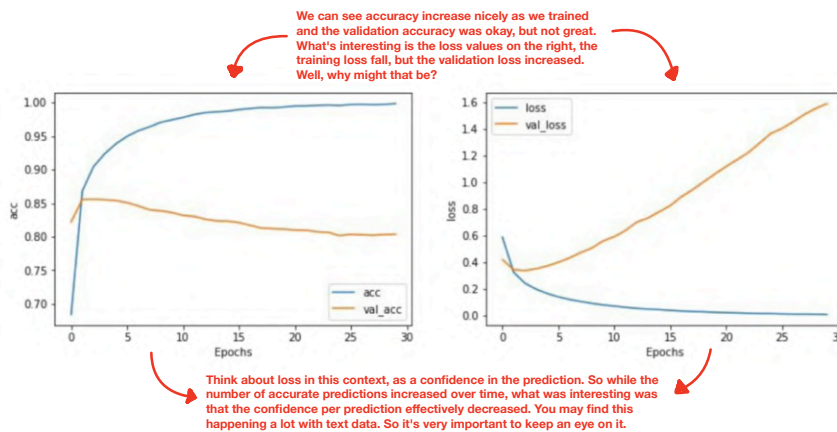
```
num_epochs = 30
history = model.fit(training_padded, training_labels, epochs=num_epochs,
                    validation_data=(testing_padded, testing_labels), verbose=2)
```

To train for 30 epochs, you pass in the padded data and labels. If you want to validate, you'll give the testing padded and labels to.

```
import matplotlib.pyplot as plt

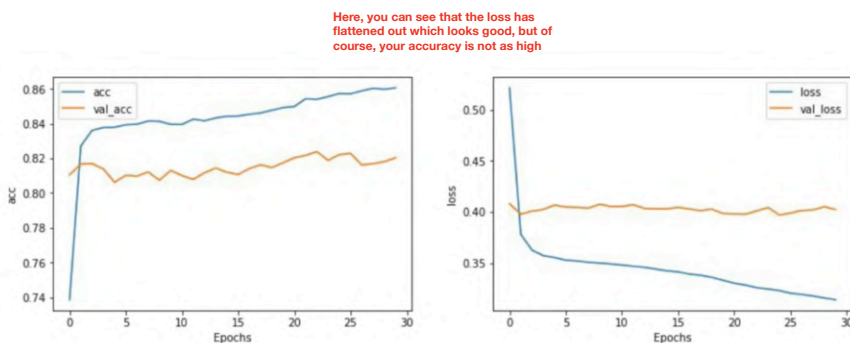
def plot_graphs(history, string):
    plt.plot(history.history[string])
    plt.plot(history.history['val_'+string])
    plt.xlabel("Epochs")
    plt.ylabel(string)
    plt.legend([string, 'val_'+string])
    plt.show()

plot_graphs(history, "acc")
plot_graphs(history, "loss")
```



```
vocab_size = 1000 (was 10,000)
embedding_dim = 16
max_length = 16 (was 32)
trunc_type='post'
padding_type='post'
oov_tok = "<OOV>"
training_size = 20000
```

If you consider these changes, a decrease in vocabulary size, and taking shorter sentences, reducing the likelihood of padding, and then rerun, you may see results like this.

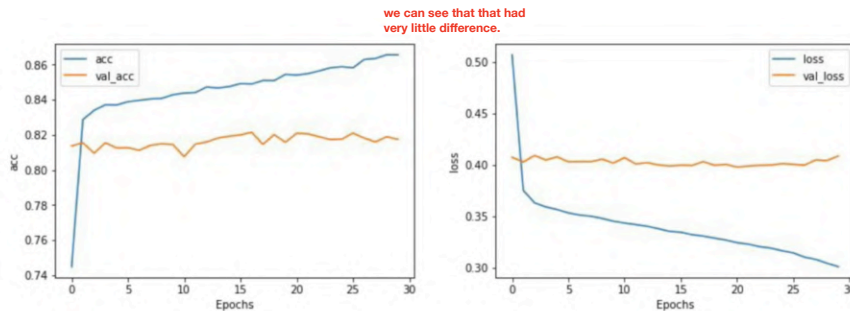


```

vocab_size = 1000 (was 10,000)
embedding_dim = 32 (was 16)
max_length = 16 (was 32)
trunc_type='post'
padding_type='post'
oov_tok = "<OOV>"
training_size = 20000

```

Another tweak. Changing the number of dimensions using the embedding was also tried.



https://github.com/tensorflow/datasets/blob/master/docs/datasets.md#imdb_reviews

"imdb_reviews"

Large Movie Review Dataset. This is a dataset for binary sentiment classification containing substantially more data than previous benchmark datasets. We provide a set of 25,000 highly polar movie reviews for training, and 25,000 for testing.

- URL: <http://ai.stanford.edu/~amaas/data/sentiment/>
- DatasetBuilder: `tfds.text.imdb.IMDBReviews`

`imdb_reviews` is configured with `tfds.text.imdb.IMDBReviewsConfig` and has the following configurations predefined (defaults to the first one):

- "plain_text" (v0.0.1) (Size: 80.23 MiB): Plain text
- "bytes" (v0.0.1) (Size: 80.23 MiB): Uses byte-level text encoding with `tfds.features.text.ByteTextEncoder`
- "subwords8k" (v0.0.1) (Size: 80.23 MiB): Uses `tfds.features.text.SubwordTextEncoder` with 8k vocab size
- "subwords32k" (v0.0.1) (Size: 80.23 MiB): Uses `tfds.features.text.SubwordTextEncoder` with 32k vocab size

https://github.com/tensorflow/datasets/blob/master/docs/datasets.md#imdb_reviews

"imdb_reviews"

Large Movie Review Dataset. This is a dataset for binary sentiment classification containing substantially more data than previous benchmark datasets. We provide a set of 25,000 highly polar movie reviews for training, and 25,000 for testing.

- URL: <http://ai.stanford.edu/~amaas/data/sentiment/>
- DatasetBuilder: `tfds.text.imdb.IMDBReviews`

`imdb_reviews` is configured with `tfds.text.imdb.IMDBReviewsConfig` and has the following configurations predefined (defaults to the first one):

- "plain_text" (v0.0.1) (Size: 80.23 MiB): Plain text
- "bytes" (v0.0.1) (Size: 80.23 MiB): Uses byte-level text encoding with `tfds.features.text.ByteTextEncoder`
- "subwords8k" (v0.0.1) (Size: 80.23 MiB): Uses `tfds.features.text.SubwordTextEncoder` with 8k vocab size
- "subwords32k" (v0.0.1) (Size: 80.23 MiB): Uses `tfds.features.text.SubwordTextEncoder` with 32k vocab size

<https://github.com/tensorflow/datasets/tree/master/docs/catalog>

https://github.com/tensorflow/datasets/blob/master/docs/catalog/imdb_reviews.md

```

import tensorflow as tf
print(tf.__version__)

```

```
!pip install tensorflow==2.0.0-alpha0
```

```

import tensorflow_datasets as tfds
imdb, info = tfds.load("imdb_reviews/subwords8k", with_info=True, as_supervised=True)

```

```
train_data, test_data = imdb['train'], imdb['test']
```



```
tokenizer = info.features['text'].encoder
```

If you want to access the sub words tokenizer, you can do it with this code. You can learn all about the sub-words texts encoder at this URL.

https://www.tensorflow.org/datasets/api_docs/python/tfds/features/text/SubwordTextEncoder

```
tensorflow.org/datasets/api_docs/python/tfds/features/text/SubwordTextEncoder
```

```
print(tokenizer.subwords)
```

```
['the_', ', ', '. ', 'a_', 'and_', 'of_', 'to_', 's_', 'is_',  
'br', 'in_', 'I_', 'that_', 'this_', 'it_', ... ]
```

```
sample_string = 'TensorFlow, from basics to mastery'
```

```
tokenized_string = tokenizer.encode(sample_string)  
print('Tokenized string is {}'.format(tokenized_string))
```

we can encode simply by calling the encode method passing it the string

```
original_string = tokenizer.decode(tokenized_string)  
print('The original string: {}'.format(original_string))
```

```
Tokenized string is [6307, 2327, 4043, 2120, 2, 48, 4249, 4429,  
7, 2652, 8050]
```

```
The original string: TensorFlow, from basics to mastery
```

```
sample_string = 'TensorFlow, from basics to mastery'
```

```
tokenized_string = tokenizer.encode(sample_string)  
print('Tokenized string is {}'.format(tokenized_string))
```

```
original_string = tokenizer.decode(tokenized_string)  
print('The original string: {}'.format(original_string))
```

decode by calling the decode method.

```
Tokenized string is [6307, 2327, 4043, 2120, 2, 48, 4249, 4429,  
7, 2652, 8050]
```

```
The original string: TensorFlow, from basics to mastery
```

```

sample_string = 'TensorFlow, from basics to mastery'

tokenized_string = tokenizer.encode(sample_string)
print ('Tokenized string is {}'.format(tokenized_string))

original_string = tokenizer.decode(tokenized_string)
print ('The original string: {}'.format(original_string))

```

We can see the results of the tokenization when we print out the encoded and decoded strings.

```

Tokenized string is [6307, 2327, 4043, 2120, 2, 48, 4249, 4429,
7, 2652, 8050]

```

```

The original string: TensorFlow, from basics to mastery

```

```

for ts in tokenized_string:
    print ('{} ----> {}'.format(ts, tokenizer.decode([ts])))

```

```

6307 ----> Ten
2327 ----> sor
4043 ----> Fl
2120 ----> ow
2 ----> ,
48 ----> from
4249 ----> basi
4429 ----> cs
7 ----> to
2652 ----> master
8050 ----> y

```

If we want to see the tokens themselves, we can take each element and decode that, showing the value to token. Note that this is case sensitive and punctuation is maintained

```

embedding_dim = 64
model = tf.keras.Sequential([
    tf.keras.layers.Embedding(tokenizer.vocab_size, embedding_dim),
    tf.keras.layers.GlobalAveragePooling1D(),
    tf.keras.layers.Dense(6, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

model.summary()

```

One thing to take into account though, is the shape of the vectors coming from the tokenizer through the embedding, and it's not easily flattened. So we'll use Global Average Pooling 1D instead. Trying to flatten them, will cause a TensorFlow crash.

Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, None, 64)	523840
global_average_pooling1d_1 (GlobalAveragePooling1D)	(None, 64)	0
dense_4 (Dense)	(None, 6)	390
dense_5 (Dense)	(None, 1)	7

Total params: 524,237
 Trainable params: 524,237
 Non-trainable params: 0

```

num_epochs = 10

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

history = model.fit(train_data,
                    epochs=num_epochs,
                    validation_data=test_data)

```

```
import matplotlib.pyplot as plt
```

```

def plot_graphs(history, string):
    plt.plot(history.history[string])
    plt.plot(history.history['val_'+string])
    plt.xlabel("Epochs")
    plt.ylabel(string)
    plt.legend([string, 'val_'+string])
    plt.show()

```

```

plot_graphs(history, "acc")
plot_graphs(history, "loss")

```

You can graph the results with this code, and your graphs will probably look something like this. In my case, the accuracy was barely about 50 percent, which you could get with a random guess. While losses decreasing, it's decreasing in a very small way. So why do you think that might be? Well, the keys in the fact that we're using sub-words and not for-words, sub-word meanings are often nonsensical and it's only when we put them together in sequences that they have meaningful semantics.

