

7. Workshop 2-1 : IMDB - NLP - One-Hot-Encoding

[Reference] :

François Chollet, **Deep Learning with Python**, Chapter 3, Section 5, Manning, 2018. http://www.deeplearningitalia.com/wp-content/uploads/2017/12/Dropbox_Chollet.pdf
(http://www.deeplearningitalia.com/wp-content/uploads/2017/12/Dropbox_Chollet.pdf)

In [1]:

```
1 import keras
2 keras.__version__
```

```
/Users/macmini1/anaconda3/lib/python3.6/site-packages/h5py
y/__init__.py:36: FutureWarning: Conversion of the second
argument of issubdtype from `float` to `np.floating` is de
precated. In future, it will be treated as `np.float64 ==
np.dtype(float).type`.
```

```
from ._conv import register_converters as _register_conv
eters
Using TensorFlow backend.
```

Out[1]:

```
'2.2.4'
```

The IMDB dataset

We'll be working with "IMDB dataset", a set of 50,000 highly-polarized reviews from the Internet Movie Database. They are split into 25,000 reviews for training and 25,000 reviews for testing, each set consisting in 50% negative and 50% positive reviews.

The following code will load the dataset (when you run it for the first time, about 80MB of data will be downloaded to your machine):

In [2]:

```
1 from keras.datasets import imdb
2
3 (train_data, train_labels), (test_data, test_labels) = imdb.load_data
```

The argument `num_words=10000` means that we will only keep the top 10,000 most frequently occurring words in the training data. Rare words will be discarded. This allows us to work with vector data of manageable size.

The variables `train_data` and `test_data` are lists of reviews, each review being a list of word indices (encoding a sequence of words). `train_labels` and `test_labels` are lists of 0s and 1s, where 0 stands for "negative" and 1 stands for "positive":

In [3]:

```
1 print(train_data[0])
```

```
[1, 14, 22, 16, 43, 530, 973, 1622, 1385, 65, 458, 4468, 6
6, 3941, 4, 173, 36, 256, 5, 25, 100, 43, 838, 112, 50, 67
0, 2, 9, 35, 480, 284, 5, 150, 4, 172, 112, 167, 2, 336, 3
85, 39, 4, 172, 4536, 1111, 17, 546, 38, 13, 447, 4, 192,
50, 16, 6, 147, 2025, 19, 14, 22, 4, 1920, 4613, 469, 4, 2
2, 71, 87, 12, 16, 43, 530, 38, 76, 15, 13, 1247, 4, 22, 1
7, 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, 28, 77, 52, 5, 14, 407, 16, 82, 2, 8, 4, 107,
117, 5952, 15, 256, 4, 2, 7, 3766, 5, 723, 36, 71, 43, 53
0, 476, 26, 400, 317, 46, 7, 4, 2, 1029, 13, 104, 88, 4, 3
81, 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, 5
1, 36, 28, 224, 92, 25, 104, 4, 226, 65, 16, 38, 1334, 88,
12, 16, 283, 5, 16, 4472, 113, 103, 32, 15, 16, 5345, 19,
178, 32]
```

In [4]:

```
1 train_labels[0]
```

Out[4]:

```
1
```

Since we restricted ourselves to the top 10,000 most frequent words, no word index will exceed 10,000:

In [5]:

```
1 max([max(sequence) for sequence in train_data])
```

Out[5]:

```
9999
```

For kicks, here's how you can quickly decode one of these reviews back to English words:

In [6]:

```
1 # word_index is a dictionary mapping words to an integer index
2 word_index = imdb.get_word_index()
3 # We reverse it, mapping integer indices to words
4 reverse_word_index = dict([(value, key) for (key, value) in word_in
5 # We decode the review; note that our indices were offset by 3
6 # because 0, 1 and 2 are reserved indices for "padding", "start of
7 decoded_review = ' '.join([reverse_word_index.get(i - 3, '?') for i
```

In [7]:

1	decoded_review
---	----------------

Out[7]:

```
"? this film was just brilliant casting location scenery s
tory direction everyone's really suited the part they play
ed and you could just imagine being there robert ? is an a
mazing actor and now the same being director ? father came
from the same scottish island as myself so i loved the fac
t there was a real connection with this film the witty rem
arks throughout the film were great it was just brilliant
so much that i bought the film as soon as it was released
for ? 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 mus
t have been good and this definitely was also ? to the two
little boy's that played the ? of norman and paul they wer
e just brilliant children are often left out of the ? list
i think because the stars that play them all grown up are
such a big profile for the whole film but these children a
re amazing and should be praised for what they have done d
on't you think the whole story was so lovely because it wa
s true and was someone's life after all that was shared wi
th us all"
```

Preparing the data - One-hot-encoding

We cannot feed lists of integers into a neural network. We have to turn our lists into tensors. There are two ways we could do that:

- We could pad our lists so that they all have the same length, and turn them into an integer tensor of shape `(samples, word_indices)`, then use as first layer in our network a layer capable of handling such integer tensors (the `Embedding` layer, which we will cover in detail later in the book).
- We could one-hot-encode our lists to turn them into vectors of 0s and 1s. Concretely, this would mean for instance turning the sequence `[3, 5]` into a 10,000-dimensional vector that would be all-zeros except for indices 3 and 5, which would be ones. Then we could use as first layer in our network a `Dense` layer, capable of handling floating point vector data.

We will go with the latter solution. Let's vectorize our data, which we will do manually for maximum clarity:

In [8]:

```
1 import numpy as np
2
3 def vectorize_sequences(sequences, dimension=10000):
4     # Create an all-zero matrix of shape (len(sequences), dimension)
5     results = np.zeros((len(sequences), dimension))
6     for i, sequence in enumerate(sequences):
7         results[i, sequence] = 1. # set specific indices of result
8     return results
9
10 # Our vectorized training data
11 x_train = vectorize_sequences(train_data)
12 # Our vectorized test data
13 x_test = vectorize_sequences(test_data)
```

Here's what our samples look like now:

In [9]:

```
1 x_train[0]
```

Out[9]:

```
array([0., 1., 1., ..., 0., 0., 0.])
```

We should also vectorize our labels, which is straightforward:

In [10]:

```
1 # Our vectorized labels
2 y_train = np.asarray(train_labels).astype('float32')
3 y_test = np.asarray(test_labels).astype('float32')
```

Now our data is ready to be fed into a neural network.

Building our network

Our input data is simply vectors, and our labels are scalars (1s and 0s): this is the easiest setup you will ever encounter. A type of network that performs well on such a problem would be a simple stack of fully-connected (`Dense`) layers with `relu` activations:

```
Dense(16, activation='relu')
```

The argument being passed to each `Dense` layer (16) is the number of "hidden units" of the layer. What's a hidden unit? It's a dimension in the representation space of the layer. You may remember from the previous chapter that each such `Dense` layer with a `relu` activation implements the following chain of tensor operations:

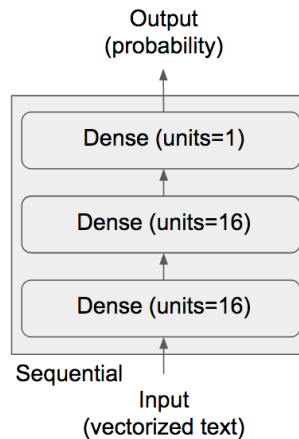
```
output = relu(dot(W, input) + b)
```

Having 16 hidden units means that the weight matrix `w` will have shape `(input_dimension, 16)`, i.e. the dot product with `w` will project the input data onto a 16-dimensional representation space (and then we would add the bias vector `b` and apply

the `relu` operation).

A `relu` (rectified linear unit) is a function meant to zero-out negative values, while a sigmoid "squashes" arbitrary values into the $[0, 1]$ interval, thus outputting something that can be interpreted as a probability.

Here's what our network looks like:



In [11]:

```
1 from keras import models
2 from keras import layers
3
4 model = models.Sequential()
5 model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
6 model.add(layers.Dense(16, activation='relu'))
7 model.add(layers.Dense(1, activation='sigmoid'))
```

Lastly, we need to pick a loss function and an optimizer.

- Since we are facing a binary classification problem and the output of our network is a probability (we end our network with a single-unit layer with a sigmoid activation), is it best to use the `binary_crossentropy` loss.
- Crossentropy is a quantity from the field of Information Theory, that measures the "distance" between probability distributions, or in our case, between the ground-truth distribution and our predictions.

We are passing our optimizer, loss function and metrics as strings, which is possible because `rmsprop`, `binary_crossentropy` and `accuracy` are packaged as part of Keras. Sometimes you may want to configure the parameters of your optimizer, or pass a custom loss function or metric function. This former can be done by passing an optimizer class instance as the `optimizer` argument:

Here's the step where we configure our model with the `rmsprop` optimizer and the `binary_crossentropy` loss function. Note that we will also monitor accuracy during training.

In [12]:

```
1 from keras import optimizers
2
3 model.compile(optimizer=optimizers.RMSprop(lr=0.001),
4               loss='binary_crossentropy',
5               metrics=['accuracy'])
```

The latter can be done by passing function objects as the `loss` or `metrics` arguments:

In [13]:

```
1 from keras import losses
2 from keras import metrics
3
4 model.compile(optimizer=optimizers.RMSprop(lr=0.001),
5               loss=losses.binary_crossentropy,
6               metrics=[metrics.binary_accuracy])
```

Validating our approach - 15000 Training Samples + 10000 Validating Samples

In order to monitor during training the accuracy of the model on data that it has never seen before, we will create a "validation set" by setting apart 10,000 samples from the original training data:

In [14]:

```
1 x_val = x_train[:10000]
2 partial_x_train = x_train[10000:]
3
4 y_val = y_train[:10000]
5 partial_y_train = y_train[10000:]
```

We will now train our model for 20 epochs (20 iterations over all samples in the `x_train` and `y_train` tensors), in mini-batches of 512 samples. At this same time we will monitor loss and accuracy on the 10,000 samples that we set apart. This is done by passing the validation data as the `validation_data` argument:

In [15]:

```
1 history = model.fit(partial_x_train,
2                     partial_y_train,
3                     epochs=20,
4                     batch_size=512,
5                     validation_data=(x_val, y_val))
```

Train on 15000 samples, validate on 10000 samples

Epoch 1/20

15000/15000 [=====] - 5s 363us/st
ep - loss: 0.5084 - binary_accuracy: 0.7813 - val_loss: 0.
3797 - val_binary_accuracy: 0.8684

Epoch 2/20

15000/15000 [=====] - 3s 212us/st
ep - loss: 0.3004 - binary_accuracy: 0.9047 - val_loss: 0.
3004 - val_binary_accuracy: 0.8897

Epoch 3/20

15000/15000 [=====] - 2s 164us/st
ep - loss: 0.2179 - binary_accuracy: 0.9285 - val_loss: 0.
3087 - val_binary_accuracy: 0.8711

Epoch 4/20

15000/15000 [=====] - 2s 150us/st
ep - loss: 0.1750 - binary_accuracy: 0.9438 - val_loss: 0.
2840 - val_binary_accuracy: 0.8832

Epoch 5/20

15000/15000 [=====] - 2s 138us/st
ep - loss: 0.1427 - binary_accuracy: 0.9543 - val_loss: 0.
2841 - val_binary_accuracy: 0.8872

Epoch 6/20

15000/15000 [=====] - 2s 136us/st
ep - loss: 0.1150 - binary_accuracy: 0.9650 - val_loss: 0.
3162 - val_binary_accuracy: 0.8770

Epoch 7/20

15000/15000 [=====] - 2s 134us/st
ep - loss: 0.0980 - binary_accuracy: 0.9707 - val_loss: 0.
3127 - val_binary_accuracy: 0.8846

Epoch 8/20

15000/15000 [=====] - 2s 131us/st
ep - loss: 0.0807 - binary_accuracy: 0.9763 - val_loss: 0.
3859 - val_binary_accuracy: 0.8650

Epoch 9/20

15000/15000 [=====] - 2s 132us/st
ep - loss: 0.0661 - binary_accuracy: 0.9821 - val_loss: 0.
3634 - val_binary_accuracy: 0.8782

Epoch 10/20

15000/15000 [=====] - 2s 135us/st
ep - loss: 0.0560 - binary_accuracy: 0.9853 - val_loss: 0.
3842 - val_binary_accuracy: 0.8793

Epoch 11/20

15000/15000 [=====] - 2s 134us/st
ep - loss: 0.0442 - binary_accuracy: 0.9890 - val_loss: 0.
4154 - val_binary_accuracy: 0.8778

Epoch 12/20

15000/15000 [=====] - 2s 130us/st
ep - loss: 0.0382 - binary_accuracy: 0.9917 - val_loss: 0.
4518 - val_binary_accuracy: 0.8688

Epoch 13/20

15000/15000 [=====] - 2s 129us/st

```

ep - loss: 0.0300 - binary_accuracy: 0.9929 - val_loss: 0.
4696 - val_binary_accuracy: 0.8731
Epoch 14/20
15000/15000 [=====] - 2s 127us/st
ep - loss: 0.0247 - binary_accuracy: 0.9947 - val_loss: 0.
5020 - val_binary_accuracy: 0.8721
Epoch 15/20
15000/15000 [=====] - 2s 147us/st
ep - loss: 0.0172 - binary_accuracy: 0.9983 - val_loss: 0.
5545 - val_binary_accuracy: 0.8665
Epoch 16/20
15000/15000 [=====] - 2s 127us/st
ep - loss: 0.0138 - binary_accuracy: 0.9981 - val_loss: 0.
5887 - val_binary_accuracy: 0.8686
Epoch 17/20
15000/15000 [=====] - 2s 127us/st
ep - loss: 0.0129 - binary_accuracy: 0.9979 - val_loss: 0.
6142 - val_binary_accuracy: 0.8678
Epoch 18/20
15000/15000 [=====] - 2s 127us/st
ep - loss: 0.0119 - binary_accuracy: 0.9971 - val_loss: 0.
6437 - val_binary_accuracy: 0.8699
Epoch 19/20
15000/15000 [=====] - 2s 127us/st
ep - loss: 0.0055 - binary_accuracy: 0.9997 - val_loss: 0.
7664 - val_binary_accuracy: 0.8519
Epoch 20/20
15000/15000 [=====] - 2s 131us/st
ep - loss: 0.0065 - binary_accuracy: 0.9994 - val_loss: 0.
7003 - val_binary_accuracy: 0.8677

```

On CPU, this will take less than two seconds per epoch -- training is over in 20 seconds. At the end of every epoch, there is a slight pause as the model computes its loss and accuracy on the 10,000 samples of the validation data.

Note that the call to `model.fit()` returns a `History` object. This object has a member `history`, which is a dictionary containing data about everything that happened during training. Let's take a look at it:

In [16]:

```

1 history_dict = history.history
2 history_dict.keys()

```

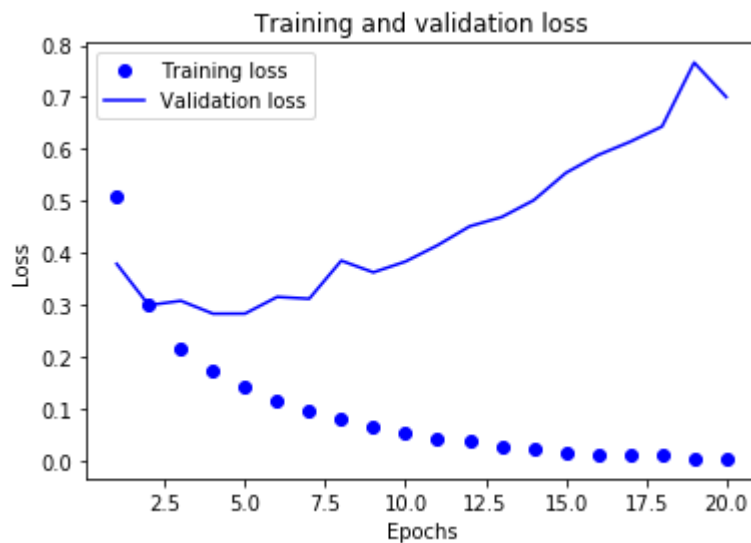
Out[16]:

```
dict_keys(['val_loss', 'val_binary_accuracy', 'loss', 'binary_accuracy'])
```

It contains 4 entries: one per metric that was being monitored, during training and during validation. Let's use Matplotlib to plot the training and validation loss side by side, as well as the training and validation accuracy:

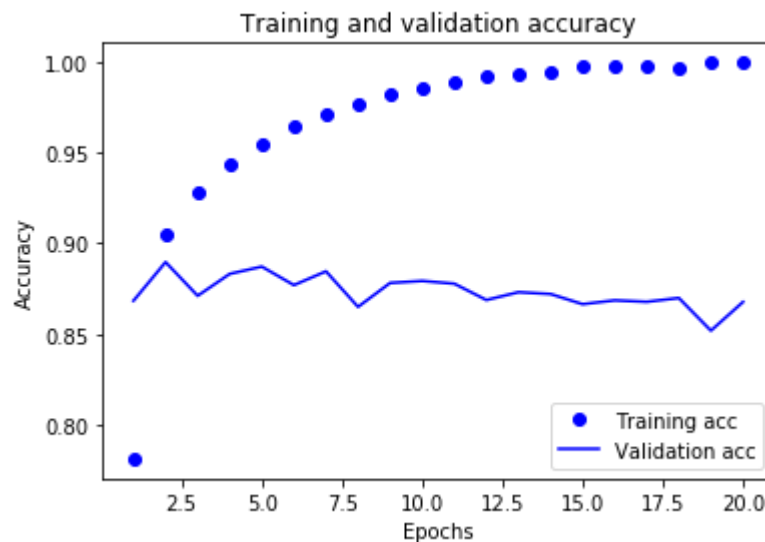
In [17]:

```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3
4 acc = history.history['binary_accuracy'] # 原始程式：此處為
5 val_acc = history.history['val_binary_accuracy'] # 原始程式：此處為
6 loss = history.history['loss']
7 val_loss = history.history['val_loss']
8
9 epochs = range(1, len(acc) + 1)
10
11 # "bo" is for "blue dot"
12 plt.plot(epochs, loss, 'bo', label='Training loss')
13 # b is for "solid blue line"
14 plt.plot(epochs, val_loss, 'b', label='Validation loss')
15 plt.title('Training and validation loss')
16 plt.xlabel('Epochs')
17 plt.ylabel('Loss')
18 plt.legend()
19
20 plt.show()
```



In [18]:

```
1 plt.clf() # clear figure
2 acc_values = history_dict['binary_accuracy'] # 原始程式：此處
3 val_acc_values = history_dict['val_binary_accuracy'] # 原始程式：此處
4
5 plt.plot(epochs, acc, 'bo', label='Training acc')
6 plt.plot(epochs, val_acc, 'b', label='Validation acc')
7 plt.title('Training and validation accuracy')
8 plt.xlabel('Epochs')
9 plt.ylabel('Accuracy')
10 plt.legend()
11
12 plt.show()
```



The dots are the training loss and accuracy, while the solid lines are the validation loss and accuracy. Note that your own results may vary slightly due to a different random initialization of your network.

As you can see, the training loss decreases with every epoch and the training accuracy increases with every epoch. That's what you would expect when running gradient descent optimization -- the quantity you are trying to minimize should get lower with every iteration. But that isn't the case for the validation loss and accuracy: they seem to peak at the fourth epoch. This is an example of what we were warning against earlier: a model that performs better on the training data isn't necessarily a model that will do better on data it has never seen before. **In precise terms, what you are seeing is "overfitting": after the second epoch, we are over-optimizing on the training data, and we ended up learning representations that are specific to the training data and do not generalize to data outside of the training set.**

- In this case, to prevent overfitting, we could simply stop training after three epochs. In general, there is a range of techniques you can leverage to mitigate overfitting, which we will cover in the next chapter.

Let's train a new network from scratch for four epochs, then evaluate it on our test data:

In [19]:

```
1 model = models.Sequential()
2 model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
3 model.add(layers.Dense(16, activation='relu'))
4 model.add(layers.Dense(1, activation='sigmoid'))
5
6 model.compile(optimizer='rmsprop',
7               loss='binary_crossentropy',
8               metrics=['accuracy'])
9
10 model.fit(x_train, y_train, epochs=4, batch_size=512)
11 results = model.evaluate(x_test, y_test)
```

```
Epoch 1/4
25000/25000 [=====] - 3s 137us/st
ep - loss: 0.4749 - acc: 0.8217
Epoch 2/4
25000/25000 [=====] - 3s 102us/st
ep - loss: 0.2658 - acc: 0.9098
Epoch 3/4
25000/25000 [=====] - 2s 89us/st
p - loss: 0.1983 - acc: 0.9297
Epoch 4/4
25000/25000 [=====] - 2s 96us/st
p - loss: 0.1678 - acc: 0.9402
25000/25000 [=====] - 5s 182us/st
ep
```

In [20]:

```
1 results
```

Out[20]:

```
[0.3235799854755402, 0.87348]
```

Our fairly naive approach achieves an accuracy of 88%. With state-of-the-art approaches, one should be able to get close to 95%.

Using a trained network to generate predictions on new data

After having trained a network, you will want to use it in a practical setting. You can generate the likelihood of reviews being positive by using the `predict` method:

In [21]:

```
1 model.predict(x_test)
```

Out[21]:

```
array([[0.13985816],
       [0.9997067 ],
       [0.29377815],
       ...,
       [0.07164008],
       [0.04311517],
       [0.4778198 ]], dtype=float32)
```

As you can see, the network is very confident for some samples (0.99 or more, or 0.01 or less) but less confident for others (0.6, 0.4).

Further experiments

- We were using 2 hidden layers. Try to use 1 or 3 hidden layers and see how it affects validation and test accuracy.
- Try to use layers with more hidden units or less hidden units: 32 units, 64 units...
- Try to use the `mse` loss function instead of `binary_crossentropy`.
- Try to use the `tanh` activation (an activation that was popular in the early days of neural networks) instead of `relu`.

These experiments will help convince you that the architecture choices we have made are all fairly reasonable, although they can still be improved!

[Further Reading] :

Deep Learning with Python, Chapter 4, Section 4 : Overfitting and Underfitting

(4.4-overfitting-and-underfitting.ipynb)