

## assignment 5.1 Implement the movie review classifier

October 3, 2020

### 0.0.1 5.1 Assignment pg 68

Implement the movie review classifier found in section 3.4 of Deep Learning with Python as a Luigi workflow. Example code and results can be found in [dsc650/assignments/assignment05/](#).

Note: `#pip install keras`

### 0.0.2 Import libraries

```
[40]: from tensorflow import keras
      from tensorflow.keras.datasets import mnist
      from tensorflow.keras.models import Sequential
      from tensorflow.keras.layers import Dense, Dropout
      from tensorflow.keras.optimizers import RMSprop
```

```
[41]: import keras
      keras.__version__
```

```
[41]: '2.4.3'
```

### 0.0.3 Load the dataset

when you run it for the first time, about 80MB of data will be downloaded to your machine:

```
[42]: from keras.datasets import imdb

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

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”:

```
[43]: train_data[0]
```

[43] : [1,  
14,  
22,  
16,  
43,  
530,  
973,  
1622,  
1385,  
65,  
458,  
4468,  
66,  
3941,  
4,  
173,  
36,  
256,  
5,  
25,  
100,  
43,  
838,  
112,  
50,  
670,  
2,  
9,  
35,  
480,  
284,  
5,  
150,  
4,  
172,  
112,  
167,  
2,  
336,  
385,  
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,  
22,  
71,  
87,  
12,  
16,  
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,  
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,  
530,  
476,  
26,  
400,  
317,  
46,  
7,  
4,  
2,  
1029,  
13,  
104,  
88,  
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,  
5345,  
19,  
178,  
32]
```

```
[44]: train_labels[0]
```

```
[44]: 1
```

```
[45]: #Since we restricted ourselves to the top 10,000 most frequent words, no word_  
↪index will exceed 10,000:  
max([max(sequence) for sequence in train_data])
```

```
[45]: 9999
```

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. It isn't the only viable choice: you could use, for instance, `mean_squared_error`. But crossentropy is usually the best choice when you are dealing with models that output probabilities. 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.

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.

```
[46]: # word_index is a dictionary mapping words to an integer index
word_index = imdb.get_word_index()

# We reverse it, mapping integer indices to words
reverse_word_index = dict([(value, key) for (key, value) in word_index.items()])

# We decode the review; note that our indices were offset by 3
# because 0, 1 and 2 are reserved indices for "padding", "start of sequence",
# and "unknown".
decoded_review = ' '.join([reverse_word_index.get(i - 3, '?') for i in
    train_data[0]])
decoded_review
```

```
[46]: "? 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 ? is an amazing actor and now the same being director ? 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
? 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 ? to the two little
boy's that played the ? of norman and paul they were 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 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"
```

#### 0.0.4 Preparing the data

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.

We could one-hot-encode our lists to turn them into vectors of 0s and 1s. Concretely, this would be a matrix of shape (number of reviews, number of words in vocabulary).

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

```
[47]: import numpy as np

def vectorize_sequences(sequences, dimension=10000):
```

```

# Create an all-zero matrix of shape (len(sequences), dimension)
results = np.zeros((len(sequences), dimension))
for i, sequence in enumerate(sequences):
    results[i, sequence] = 1. # set specific indices of results[i] to 1s
return results

# Our vectorized training data
x_train = vectorize_sequences(train_data)
# Our vectorized test data
x_test = vectorize_sequences(test_data)

```

Here's what our samples look like now:

```
[48]: x_train[0]
```

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

We should also vectorize our labels, which is straightforward:

```
[49]: # Our vectorized labels
y_train = np.asarray(train_labels).astype('float32')
y_test = np.asarray(test_labels).astype('float32')
```

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

**Building our network** The input data is simply vectors, and the 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:

$$\text{output} = \text{relu}(\text{dot}(W, \text{input}) + b)$$

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). You can intuitively understand the dimensionality of your representation space as “how much freedom you are allowing the network to have when learning internal representations”. Having more hidden units (a higher-dimensional representation space) allows your network to learn more complex representations, but it makes your network more computationally expensive and may lead to learning unwanted patterns (patterns that will improve performance on the training data but not on the test data).

There are two key architecture decisions to be made about such stack of dense layers:

How many layers to use.

How many "hidden units" to chose for each layer.



In the next chapter, you will learn formal principles to guide you in making these choices. For the time being, you will have to trust us with the following architecture choice: two intermediate layers with 16 hidden units each, and a third layer which will output the scalar prediction regarding the sentiment of the current review. The intermediate layers will use `relu` as their “activation function”, and the final layer will use a `sigmoid` activation so as to output a probability (a score between 0 and 1, indicating how likely the sample is to have the target “1”, i.e. how likely the review is to be positive). 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.

### 0.0.5 The Keras implementation, very similar to the MNIST example you saw previously:

```
[50]: #The Keras implementation

from keras import models
from keras import layers

model = models.Sequential()
model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(16, activation='relu'))
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. It isn’t the only viable choice: you could use, for instance, `mean_squared_error`. But `crossentropy` is usually the best choice when you are dealing with models that output probabilities. `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.

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.

```
[51]: model.compile(optimizer='rmsprop',
                  loss='binary_crossentropy',
                  metrics=['accuracy'])
```

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:

```
[52]: from keras import optimizers

model.compile(optimizer=optimizers.RMSprop(lr=0.001),
              loss='binary_crossentropy',
              metrics=['accuracy'])
```

```
[53]: from keras import losses
      from keras import metrics

      model.compile(optimizer=optimizers.RMSprop(lr=0.001),
                    loss=losses.binary_crossentropy,
                    metrics=[metrics.binary_accuracy])
```

### 0.0.6 Validating our approach

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:

```
[54]: x_val = x_train[:10000]
      partial_x_train = x_train[10000:]

      y_val = y_train[:10000]
      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:

```
[55]: history = model.fit(partial_x_train,
                          partial_y_train,
                          epochs=20,
                          batch_size=512,
                          validation_data=(x_val, y_val))
```

```
Epoch 1/20
30/30 [=====] - 1s 30ms/step - loss: 0.5057 -
binary_accuracy: 0.7887 - val_loss: 0.3782 - val_binary_accuracy: 0.8666
Epoch 2/20
30/30 [=====] - 1s 26ms/step - loss: 0.2929 -
binary_accuracy: 0.9023 - val_loss: 0.3097 - val_binary_accuracy: 0.8772
Epoch 3/20
30/30 [=====] - 1s 27ms/step - loss: 0.2145 -
binary_accuracy: 0.9287 - val_loss: 0.2788 - val_binary_accuracy: 0.8906
Epoch 4/20
30/30 [=====] - 1s 27ms/step - loss: 0.1681 -
binary_accuracy: 0.9467 - val_loss: 0.2880 - val_binary_accuracy: 0.8824
Epoch 5/20
30/30 [=====] - 1s 25ms/step - loss: 0.1399 -
binary_accuracy: 0.9556 - val_loss: 0.3045 - val_binary_accuracy: 0.8794
Epoch 6/20
30/30 [=====] - 1s 26ms/step - loss: 0.1155 -
binary_accuracy: 0.9648 - val_loss: 0.3011 - val_binary_accuracy: 0.8855
Epoch 7/20
```

```

30/30 [=====] - 1s 27ms/step - loss: 0.0979 -
binary_accuracy: 0.9707 - val_loss: 0.3184 - val_binary_accuracy: 0.8846
Epoch 8/20
30/30 [=====] - 1s 27ms/step - loss: 0.0768 -
binary_accuracy: 0.9791 - val_loss: 0.3402 - val_binary_accuracy: 0.8828
Epoch 9/20
30/30 [=====] - 1s 26ms/step - loss: 0.0658 -
binary_accuracy: 0.9824 - val_loss: 0.3657 - val_binary_accuracy: 0.8726
Epoch 10/20
30/30 [=====] - 1s 24ms/step - loss: 0.0551 -
binary_accuracy: 0.9869 - val_loss: 0.3851 - val_binary_accuracy: 0.8762
Epoch 11/20
30/30 [=====] - 1s 29ms/step - loss: 0.0444 -
binary_accuracy: 0.9896 - val_loss: 0.4179 - val_binary_accuracy: 0.8719
Epoch 12/20
30/30 [=====] - 1s 25ms/step - loss: 0.0344 -
binary_accuracy: 0.9932 - val_loss: 0.4761 - val_binary_accuracy: 0.8727
Epoch 13/20
30/30 [=====] - 1s 28ms/step - loss: 0.0297 -
binary_accuracy: 0.9939 - val_loss: 0.4887 - val_binary_accuracy: 0.8681
Epoch 14/20
30/30 [=====] - 1s 25ms/step - loss: 0.0245 -
binary_accuracy: 0.9948 - val_loss: 0.5079 - val_binary_accuracy: 0.8709
Epoch 15/20
30/30 [=====] - 1s 27ms/step - loss: 0.0215 -
binary_accuracy: 0.9955 - val_loss: 0.5438 - val_binary_accuracy: 0.8681
Epoch 16/20
30/30 [=====] - 1s 28ms/step - loss: 0.0119 -
binary_accuracy: 0.9989 - val_loss: 0.6631 - val_binary_accuracy: 0.8499
Epoch 17/20
30/30 [=====] - 1s 27ms/step - loss: 0.0142 -
binary_accuracy: 0.9977 - val_loss: 0.6145 - val_binary_accuracy: 0.8688
Epoch 18/20
30/30 [=====] - 1s 27ms/step - loss: 0.0072 -
binary_accuracy: 0.9997 - val_loss: 0.6815 - val_binary_accuracy: 0.8570
Epoch 19/20
30/30 [=====] - 1s 26ms/step - loss: 0.0097 -
binary_accuracy: 0.9983 - val_loss: 0.6824 - val_binary_accuracy: 0.8652
Epoch 20/20
30/30 [=====] - 1s 29ms/step - loss: 0.0042 -
binary_accuracy: 0.9999 - val_loss: 0.7129 - val_binary_accuracy: 0.8640

```

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:

```
[56]: history_dict = history.history
history_dict.keys()
#dict_keys(['loss', 'val_loss', 'binary_accuracy', 'val_binary_accuracy'])

#history_dict = history.history
print(history_dict)
```

```
{'loss': [0.505662739276886, 0.2929249703884125, 0.21446317434310913,
0.1681060940027237, 0.13985630869865417, 0.11549120396375656,
0.09791756421327591, 0.07680363953113556, 0.06578858196735382,
0.05514388531446457, 0.04441479220986366, 0.03441550210118294,
0.02969759702682495, 0.02454623207449913, 0.021527240052819252,
0.011946368031203747, 0.014245566911995411, 0.007232708856463432,
0.00965186394751072, 0.004207812715321779], 'binary_accuracy':
[0.7886666655540466, 0.9023333191871643, 0.9286666512489319, 0.9466666579246521,
0.9556000232696533, 0.9648000001907349, 0.9706666469573975, 0.9791333079338074,
0.9824000000953674, 0.9868666529655457, 0.9896000027656555, 0.9932000041007996,
0.9938666820526123, 0.9947999715805054, 0.9954666495323181, 0.9989333152770996,
0.9976666569709778, 0.9997333288192749, 0.9983333349227905, 0.9998666644096375],
'val_loss': [0.37815892696380615, 0.3096763491630554, 0.27876749634742737,
0.2880069315433502, 0.3044685423374176, 0.30110716819763184, 0.3184034526348114,
0.3401671051979065, 0.3656865656375885, 0.3850671648979187, 0.41786128282546997,
0.47608858346939087, 0.4887121617794037, 0.507949948310852, 0.5438354015350342,
0.6630561947822571, 0.6145161986351013, 0.6815347671508789, 0.682390034198761,
0.7129063010215759], 'val_binary_accuracy': [0.866599977016449,
0.8772000074386597, 0.8906000256538391, 0.8823999762535095, 0.8794000148773193,
0.8855000138282776, 0.8845999836921692, 0.8827999830245972, 0.8726000189781189,
0.8762000203132629, 0.8719000220298767, 0.8726999759674072, 0.8680999875068665,
0.8708999752998352, 0.8680999875068665, 0.8499000072479248, 0.8687999844551086,
0.8569999933242798, 0.8651999831199646, 0.8640000224113464]}
```

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:

```
[57]: import matplotlib.pyplot as plt

acc = history.history['binary_accuracy']
val_acc = history.history['val_binary_accuracy']

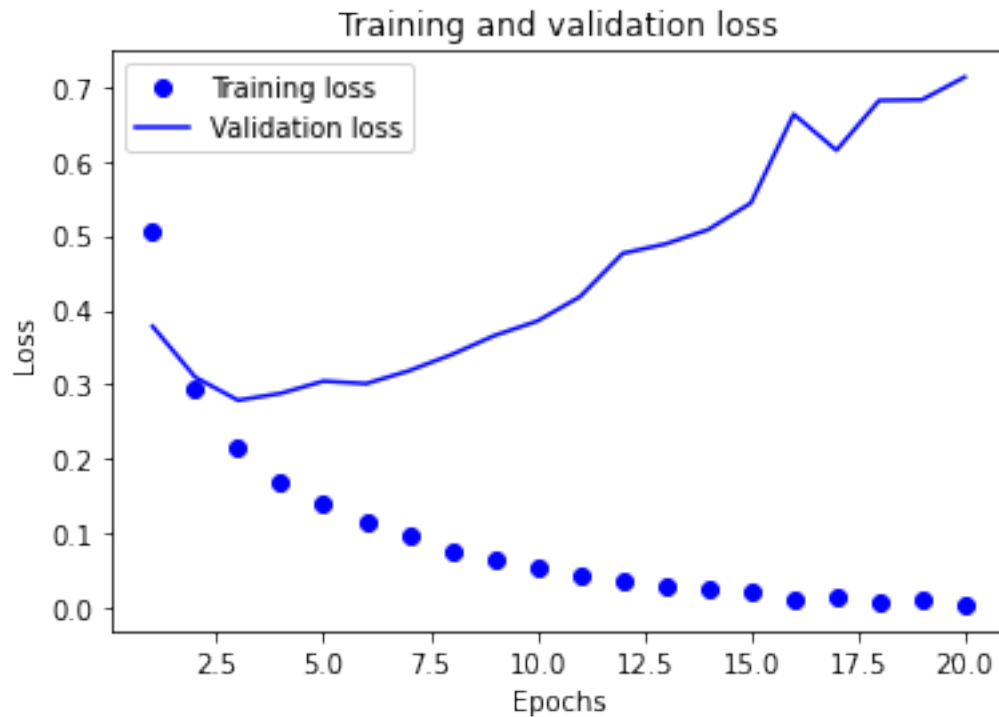
#acc = history.history['acc']
#val_acc = history.history['val_acc']
loss = history.history['loss']
val_loss = history.history['val_loss']

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()

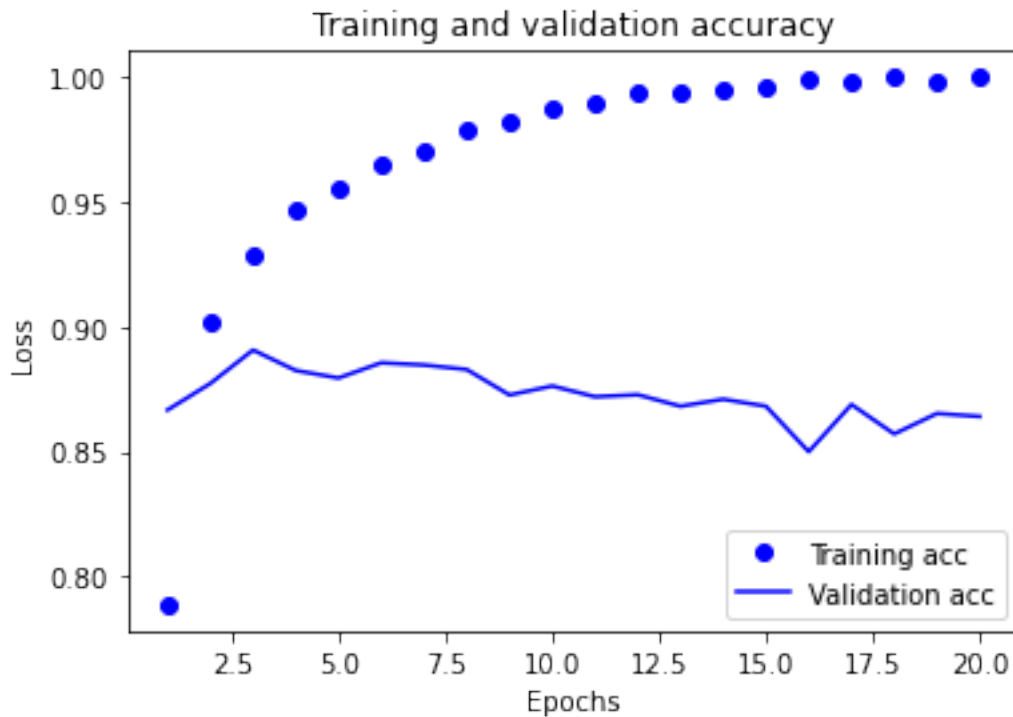
plt.show()
```



```
[58]: plt.clf() # clear figure
#acc_values = history_dict['acc']
#val_acc_values = history_dict['val_acc']
acc = history.history['binary_accuracy']
val_acc = history.history['val_binary_accuracy']

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('Loss')
plt.legend()
```

```
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:

```
[59]: model = models.Sequential()
      model.add(layers.Dense(16, activation='relu', input_shape=(10000,)))
      model.add(layers.Dense(16, activation='relu'))
      model.add(layers.Dense(1, activation='sigmoid'))
```

```

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

model.fit(x_train, y_train, epochs=4, batch_size=512)
results = model.evaluate(x_test, y_test)

```

```

Epoch 1/4
49/49 [=====] - 0s 8ms/step - loss: 0.4550 - accuracy:
0.8254
Epoch 2/4
49/49 [=====] - 0s 8ms/step - loss: 0.2605 - accuracy:
0.9088
Epoch 3/4
49/49 [=====] - 0s 8ms/step - loss: 0.2021 - accuracy:
0.9273
Epoch 4/4
49/49 [=====] - 0s 7ms/step - loss: 0.1690 - accuracy:
0.9407
782/782 [=====] - 1s 2ms/step - loss: 0.3034 -
accuracy: 0.8822

```

```
[60]: results
```

```
[60]: [0.30343276262283325, 0.8822399973869324]
```

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

### 0.0.7 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:

```
[61]: model.predict(x_test)
```

```
[61]: array([[0.14747846],
             [0.99994195],
             [0.64993894],
             ...,
             [0.09259456],
             [0.05922219],
             [0.57290345]], 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).

### 0.0.8 Further experiments

We were using 2 hidden layers. Try to use 1 or 3 hidden layers and see how it affects validation

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)

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

### 0.0.9 Conclusions

Here's what you should take away from this example:

There's usually quite a bit of preprocessing you need to do on your raw data in order to be able to

Stacks of Dense layers with relu activations can solve a wide range of problems (including sentiment

In a binary classification problem (two output classes), your network should end with a Dense layer

With such a scalar sigmoid output, on a binary classification problem, the loss function you should use

The rmsprop optimizer is generally a good enough choice of optimizer, whatever your problem. The

As they get better on their training data, neural networks eventually start overfitting and ending

END

<https://github.com/souwadeGit/deep-learning-with-python-notebooks/blob/master/3.6-classifying-newswires.ipynb>

<https://github.com/fchollet/deep-learning-with-python-notebooks/issues/16>