

Assignment 5.2 Implement the news classifier - classifying-newswires

January 18, 2021

Assignment 5.2 Implement the news classifier found in section 3.5 of Deep Learning with Python as a Luigi workflow. Example code and results can be found in [dsc650/assignments/assignment05/](https://dsc650.org/assignments/assignment05/).

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

```
[2]: '2.4.3'
```

Classifying newswires: a multi-class classification example

The Reuters dataset We will be working with the *Reuters dataset*, a set of short newswires and their topics, published by Reuters in 1986. It's a very simple, widely used toy dataset for text classification. There are 46 different topics; some topics are more represented than others, but each topic has at least 10 examples in the training set.

Like IMDB and MNIST, the Reuters dataset comes packaged as part of Keras. Let's take a look right away:

```
[3]: from keras.datasets import reuters

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

Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets/reuters.npz>

2113536/2110848 [=====] - 0s 0us/step

Like with the IMDB dataset, the argument `num_words=10000` restricts the data to the 10,000 most frequently occurring words found in the data.

We have 8,982 training examples and 2,246 test examples:

```
[4]: len(train_data)
```

```
[4]: 8982
```

```
[5]: len(test_data)
```

```
[5]: 2246
```

As with the IMDB reviews, each example is a list of integers (word indices):

```
[6]: train_data[10]
```

```
[6]: [1,
      245,
      273,
      207,
      156,
      53,
      74,
      160,
      26,
      14,
      46,
      296,
      26,
      39,
      74,
      2979,
      3554,
      14,
      46,
      4689,
      4329,
      86,
      61,
      3499,
      4795,
      14,
      61,
      451,
      4329,
      17,
      12]
```

Here's how you can decode it back to words, in case you are curious:

```
[7]: word_index = reuters.get_word_index()
reverse_word_index = dict([(value, key) for (key, value) in word_index.items()])
# Note that our indices were offset by 3
# because 0, 1 and 2 are reserved indices for "padding", "start of sequence",
→ and "unknown".
decoded_newswire = ' '.join([reverse_word_index.get(i - 3, '?') for i in
→ train_data[0]])
```

Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/reuters_word_index.json

557056/550378 [=====] - 0s 0us/step

```
[8]: decoded_newswire
```

```
[8]: '? ? ? said as a result of its december acquisition of space co it expects
earnings per share in 1987 of 1 15 to 1 30 dlrs per share up from 70 cts in 1986
the company said pretax net should rise to nine to 10 mln dlrs from six mln dlrs
in 1986 and rental operation revenues to 19 to 22 mln dlrs from 12 5 mln dlrs it
said cash flow per share this year should be 2 50 to three dlrs reuter 3'
```

The label associated with an example is an integer between 0 and 45: a topic index.

```
[9]: train_labels[10]
```

```
[9]: 3
```

Preparing the data We can vectorize the data with the exact same code as in our previous example:

```
[10]: import numpy as np

def vectorize_sequences(sequences, dimension=10000):
    results = np.zeros((len(sequences), dimension))
    for i, sequence in enumerate(sequences):
        results[i, sequence] = 1.
    return results

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

To vectorize the labels, there are two possibilities: we could just cast the label list as an integer tensor, or we could use a “one-hot” encoding. One-hot encoding is a widely used format for categorical data, also called “categorical encoding”. For a more detailed explanation of one-hot encoding, you can refer to Chapter 6, Section 1. In our case, one-hot encoding of our labels consists in embedding each label as an all-zero vector with a 1 in the place of the label index, e.g.:

```
[11]: def to_one_hot(labels, dimension=46):
        results = np.zeros((len(labels), dimension))
        for i, label in enumerate(labels):
            results[i, label] = 1.
        return results

# Our vectorized training labels
one_hot_train_labels = to_one_hot(train_labels)
# Our vectorized test labels
one_hot_test_labels = to_one_hot(test_labels)
```

Note that there is a built-in way to do this in Keras, which you have already seen in action in our MNIST example:

```
[12]: from keras.utils.np_utils import to_categorical

one_hot_train_labels = to_categorical(train_labels)
one_hot_test_labels = to_categorical(test_labels)
```

Building our network This topic classification problem looks very similar to our previous movie review classification problem: in both cases, we are trying to classify short snippets of text. There is however a new constraint here: the number of output classes has gone from 2 to 46, i.e. the dimensionality of the output space is much larger.

In a stack of `Dense` layers like what we were using, each layer can only access information present in the output of the previous layer. If one layer drops some information relevant to the classification problem, this information can never be recovered by later layers: each layer can potentially become an “information bottleneck”. In our previous example, we were using 16-dimensional intermediate layers, but a 16-dimensional space may be too limited to learn to separate 46 different classes: such small layers may act as information bottlenecks, permanently dropping relevant information.

For this reason we will use larger layers. Let’s go with 64 units:

```
[13]: from keras import models
      from keras import layers

model = models.Sequential()
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(46, activation='softmax'))
```

There are two other things you should note about this architecture:

- We are ending the network with a `Dense` layer of size 46. This means that for each input sample, our network will output a 46-dimensional vector. Each entry in this vector (each dimension) will encode a different output class.
- The last layer uses a `softmax` activation. You have already seen this pattern in the MNIST example. It means that the network will output a *probability distribution* over the 46 different output classes, i.e. for every input sample, the network will produce a 46-dimensional output vector where `output[i]` is the probability that the sample belongs to class `i`. The 46 scores will sum to 1.

The best loss function to use in this case is `categorical_crossentropy`. It measures the distance between two probability distributions: in our case, between the probability distribution output by our network, and the true distribution of the labels. By minimizing the distance between these two distributions, we train our network to output something as close as possible to the true labels.

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

Validating our approach Let's set apart 1,000 samples in our training data to use as a validation set:

```
[15]: x_val = x_train[:1000]
      partial_x_train = x_train[1000:]

      y_val = one_hot_train_labels[:1000]
      partial_y_train = one_hot_train_labels[1000:]
```

Now let's train our network for 20 epochs:

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

```
Epoch 1/20
16/16 [=====] - 1s 38ms/step - loss: 2.6559 - accuracy:
0.5318 - val_loss: 1.7841 - val_accuracy: 0.6590
Epoch 2/20
16/16 [=====] - 0s 20ms/step - loss: 1.4204 - accuracy:
0.7146 - val_loss: 1.3010 - val_accuracy: 0.7280
Epoch 3/20
16/16 [=====] - 0s 18ms/step - loss: 1.0321 - accuracy:
0.7871 - val_loss: 1.1327 - val_accuracy: 0.7560
Epoch 4/20
16/16 [=====] - 0s 17ms/step - loss: 0.8092 - accuracy:
0.8299 - val_loss: 1.0334 - val_accuracy: 0.7760
Epoch 5/20
16/16 [=====] - 0s 16ms/step - loss: 0.6396 - accuracy:
0.8692 - val_loss: 0.9611 - val_accuracy: 0.8020
Epoch 6/20
16/16 [=====] - 0s 18ms/step - loss: 0.5131 - accuracy:
0.8946 - val_loss: 0.9130 - val_accuracy: 0.8120
Epoch 7/20
16/16 [=====] - 0s 18ms/step - loss: 0.4142 - accuracy:
0.9147 - val_loss: 0.9025 - val_accuracy: 0.8080
Epoch 8/20
16/16 [=====] - 0s 19ms/step - loss: 0.3401 - accuracy:
0.9312 - val_loss: 0.8976 - val_accuracy: 0.8140
Epoch 9/20
16/16 [=====] - 0s 19ms/step - loss: 0.2818 - accuracy:
0.9396 - val_loss: 0.8819 - val_accuracy: 0.8220
Epoch 10/20
16/16 [=====] - 0s 17ms/step - loss: 0.2373 - accuracy:
0.9454 - val_loss: 0.9086 - val_accuracy: 0.8120
Epoch 11/20
```

```

16/16 [=====] - 0s 16ms/step - loss: 0.2070 - accuracy:
0.9496 - val_loss: 0.9782 - val_accuracy: 0.7930
Epoch 12/20
16/16 [=====] - 0s 17ms/step - loss: 0.1856 - accuracy:
0.9516 - val_loss: 0.9478 - val_accuracy: 0.8050
Epoch 13/20
16/16 [=====] - 0s 17ms/step - loss: 0.1645 - accuracy:
0.9536 - val_loss: 0.9664 - val_accuracy: 0.8080
Epoch 14/20
16/16 [=====] - 0s 18ms/step - loss: 0.1475 - accuracy:
0.9559 - val_loss: 0.9776 - val_accuracy: 0.8070
Epoch 15/20
16/16 [=====] - 0s 18ms/step - loss: 0.1407 - accuracy:
0.9560 - val_loss: 1.0031 - val_accuracy: 0.8070
Epoch 16/20
16/16 [=====] - 0s 18ms/step - loss: 0.1375 - accuracy:
0.9555 - val_loss: 1.0161 - val_accuracy: 0.8040
Epoch 17/20
16/16 [=====] - 0s 14ms/step - loss: 0.1226 - accuracy:
0.9575 - val_loss: 1.0834 - val_accuracy: 0.7890
Epoch 18/20
16/16 [=====] - 0s 16ms/step - loss: 0.1216 - accuracy:
0.9595 - val_loss: 1.0357 - val_accuracy: 0.8090
Epoch 19/20
16/16 [=====] - 0s 15ms/step - loss: 0.1176 - accuracy:
0.9579 - val_loss: 1.0458 - val_accuracy: 0.8170
Epoch 20/20
16/16 [=====] - 0s 15ms/step - loss: 0.1089 - accuracy:
0.9577 - val_loss: 1.0637 - val_accuracy: 0.8100

```

Let's display its loss and accuracy curves:

```

[17]: import matplotlib.pyplot as plt

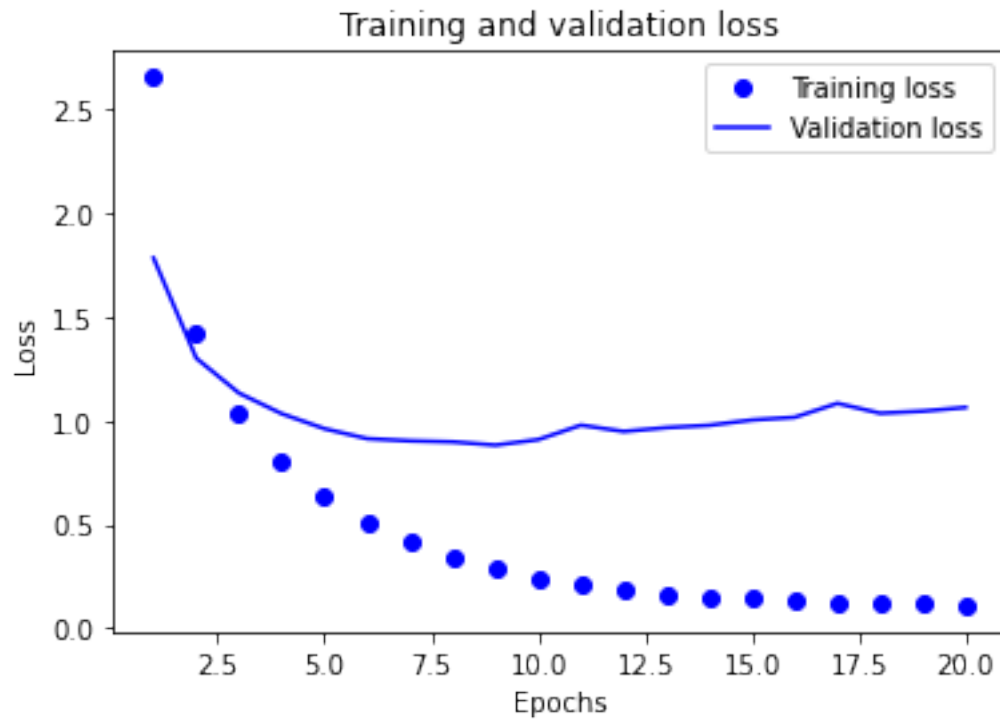
loss = history.history['loss']
val_loss = history.history['val_loss']

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

plt.plot(epochs, loss, 'bo', label='Training loss')
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()

```

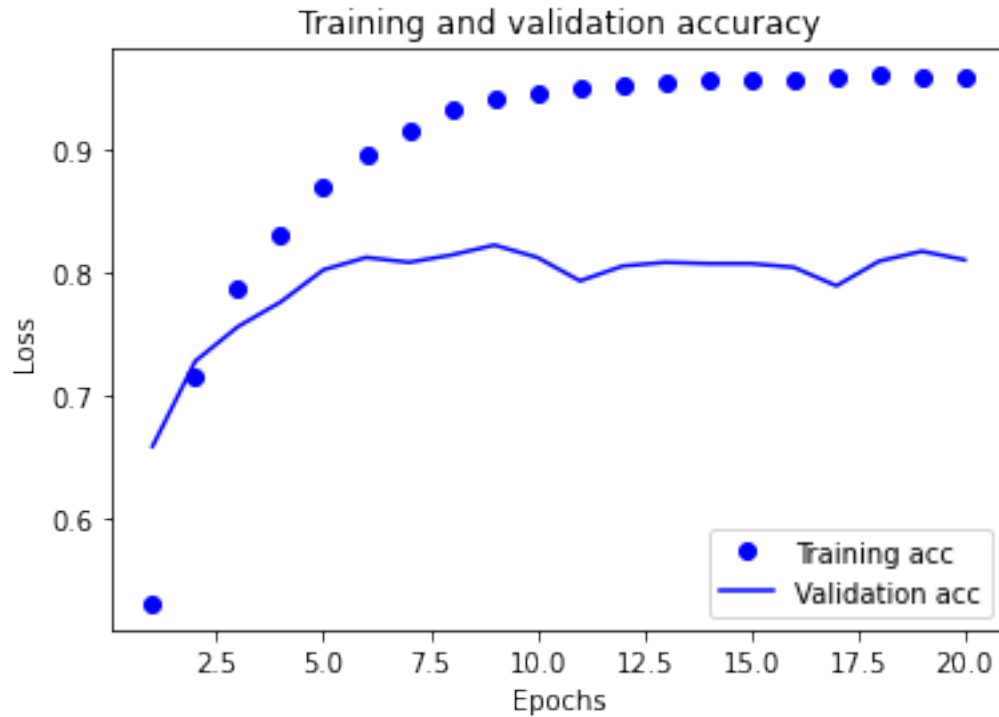


```
[19]: plt.clf()      # clear figure

acc = history.history['accuracy']
val_acc = history.history['val_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()
```



It seems that the network starts overfitting after 8 epochs. Let's train a new network from scratch for 8 epochs, then let's evaluate it on the test set:

```
[20]: model = models.Sequential()
model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
model.add(layers.Dense(64, activation='relu'))
model.add(layers.Dense(46, activation='softmax'))

model.compile(optimizer='rmsprop',
              loss='categorical_crossentropy',
              metrics=['accuracy'])
model.fit(partial_x_train,
          partial_y_train,
          epochs=8,
          batch_size=512,
          validation_data=(x_val, y_val))
results = model.evaluate(x_test, one_hot_test_labels)
```

Epoch 1/8

16/16 [=====] - 0s 23ms/step - loss: 2.6092 - accuracy: 0.5436 - val_loss: 1.7309 - val_accuracy: 0.6570

Epoch 2/8

16/16 [=====] - 0s 16ms/step - loss: 1.4066 - accuracy: 0.7124 - val_loss: 1.2981 - val_accuracy: 0.7310


```

Epoch 3/8
16/16 [=====] - 0s 21ms/step - loss: 1.0364 - accuracy:
0.7813 - val_loss: 1.1487 - val_accuracy: 0.7520
Epoch 4/8
16/16 [=====] - 0s 19ms/step - loss: 0.8248 - accuracy:
0.8235 - val_loss: 1.0525 - val_accuracy: 0.7760
Epoch 5/8
16/16 [=====] - 0s 17ms/step - loss: 0.6606 - accuracy:
0.8592 - val_loss: 0.9863 - val_accuracy: 0.7940
Epoch 6/8
16/16 [=====] - 0s 19ms/step - loss: 0.5364 - accuracy:
0.8849 - val_loss: 0.9433 - val_accuracy: 0.8090
Epoch 7/8
16/16 [=====] - 0s 15ms/step - loss: 0.4348 - accuracy:
0.9077 - val_loss: 0.9208 - val_accuracy: 0.8210
Epoch 8/8
16/16 [=====] - 0s 14ms/step - loss: 0.3525 - accuracy:
0.9250 - val_loss: 0.9402 - val_accuracy: 0.8020
71/71 [=====] - 0s 2ms/step - loss: 0.9906 - accuracy:
0.7841

```

```
[21]: results
```

```
[21]: [0.9906436204910278, 0.784060537815094]
```

Our approach reaches an accuracy of ~78%. With a balanced binary classification problem, the accuracy reached by a purely random classifier would be 50%, but in our case it is closer to 19%, so our results seem pretty good, at least when compared to a random baseline:

```

[22]: import copy

test_labels_copy = copy.copy(test_labels)
np.random.shuffle(test_labels_copy)
float(np.sum(np.array(test_labels) == np.array(test_labels_copy))) /
    ↳ len(test_labels)

```

```
[22]: 0.18432769367764915
```

Generating predictions on new data We can verify that the `predict` method of our model instance returns a probability distribution over all 46 topics. Let's generate topic predictions for all of the test data:

```
[23]: predictions = model.predict(x_test)
```

Each entry in `predictions` is a vector of length 46:

```
[24]: predictions[0].shape
```

```
[24]: (46,)
```

The coefficients in this vector sum to 1:

```
[25]: np.sum(predictions[0])
```

```
[25]: 0.9999999
```

The largest entry is the predicted class, i.e. the class with the highest probability:

```
[26]: np.argmax(predictions[0])
```

```
[26]: 3
```

A different way to handle the labels and the loss We mentioned earlier that another way to encode the labels would be to cast them as an integer tensor, like such:

```
[27]: y_train = np.array(train_labels)
      y_test = np.array(test_labels)
```

The only thing it would change is the choice of the loss function. Our previous loss, `categorical_crossentropy`, expects the labels to follow a categorical encoding. With integer labels, we should use `sparse_categorical_crossentropy`:

```
[28]: model.compile(optimizer='rmsprop', loss='sparse_categorical_crossentropy',
      ↪ metrics=['acc'])
```

This new loss function is still mathematically the same as `categorical_crossentropy`; it just has a different interface.

On the importance of having sufficiently large intermediate layers We mentioned earlier that since our final outputs were 46-dimensional, we should avoid intermediate layers with much less than 46 hidden units. Now let's try to see what happens when we introduce an information bottleneck by having intermediate layers significantly less than 46-dimensional, e.g. 4-dimensional.

```
[29]: model = models.Sequential()
      model.add(layers.Dense(64, activation='relu', input_shape=(10000,)))
      model.add(layers.Dense(4, activation='relu'))
      model.add(layers.Dense(46, activation='softmax'))

      model.compile(optimizer='rmsprop',
                    loss='categorical_crossentropy',
                    metrics=['accuracy'])
      model.fit(partial_x_train,
                partial_y_train,
                epochs=20,
                batch_size=128,
                validation_data=(x_val, y_val))
```

Epoch 1/20
63/63 [=====] - 1s 11ms/step - loss: 3.5444 - accuracy:
0.0724 - val_loss: 3.2885 - val_accuracy: 0.1080
Epoch 2/20
63/63 [=====] - 0s 8ms/step - loss: 2.9024 - accuracy:
0.1468 - val_loss: 2.5931 - val_accuracy: 0.2570
Epoch 3/20
63/63 [=====] - 0s 7ms/step - loss: 2.0825 - accuracy:
0.3929 - val_loss: 1.8874 - val_accuracy: 0.3880
Epoch 4/20
63/63 [=====] - 1s 8ms/step - loss: 1.4172 - accuracy:
0.6209 - val_loss: 1.4323 - val_accuracy: 0.7010
Epoch 5/20
63/63 [=====] - 0s 7ms/step - loss: 1.0772 - accuracy:
0.7422 - val_loss: 1.3447 - val_accuracy: 0.7050
Epoch 6/20
63/63 [=====] - 0s 7ms/step - loss: 0.9485 - accuracy:
0.7651 - val_loss: 1.3232 - val_accuracy: 0.7120
Epoch 7/20
63/63 [=====] - 0s 7ms/step - loss: 0.8626 - accuracy:
0.7880 - val_loss: 1.3370 - val_accuracy: 0.7200
Epoch 8/20
63/63 [=====] - 0s 7ms/step - loss: 0.7927 - accuracy:
0.8063 - val_loss: 1.3376 - val_accuracy: 0.7120
Epoch 9/20
63/63 [=====] - 0s 7ms/step - loss: 0.7339 - accuracy:
0.8162 - val_loss: 1.3856 - val_accuracy: 0.7190
Epoch 10/20
63/63 [=====] - 0s 7ms/step - loss: 0.6844 - accuracy:
0.8294 - val_loss: 1.4644 - val_accuracy: 0.7090
Epoch 11/20
63/63 [=====] - 0s 7ms/step - loss: 0.6411 - accuracy:
0.8390 - val_loss: 1.4823 - val_accuracy: 0.7180
Epoch 12/20
63/63 [=====] - 0s 7ms/step - loss: 0.6023 - accuracy:
0.8460 - val_loss: 1.5135 - val_accuracy: 0.7170
Epoch 13/20
63/63 [=====] - 0s 7ms/step - loss: 0.5700 - accuracy:
0.8557 - val_loss: 1.5496 - val_accuracy: 0.7190
Epoch 14/20
63/63 [=====] - 0s 7ms/step - loss: 0.5402 - accuracy:
0.8606 - val_loss: 1.6238 - val_accuracy: 0.7240
Epoch 15/20
63/63 [=====] - 0s 7ms/step - loss: 0.5114 - accuracy:
0.8682 - val_loss: 1.6981 - val_accuracy: 0.7040
Epoch 16/20
63/63 [=====] - 0s 7ms/step - loss: 0.4894 - accuracy:
0.8750 - val_loss: 1.7280 - val_accuracy: 0.7190

```
Epoch 17/20
63/63 [=====] - 0s 7ms/step - loss: 0.4690 - accuracy:
0.8777 - val_loss: 1.7849 - val_accuracy: 0.7080
Epoch 18/20
63/63 [=====] - 0s 7ms/step - loss: 0.4501 - accuracy:
0.8815 - val_loss: 1.8519 - val_accuracy: 0.7130
Epoch 19/20
63/63 [=====] - 1s 10ms/step - loss: 0.4342 - accuracy:
0.8839 - val_loss: 1.8981 - val_accuracy: 0.7160
Epoch 20/20
63/63 [=====] - 0s 7ms/step - loss: 0.4184 - accuracy:
0.8889 - val_loss: 1.9624 - val_accuracy: 0.7110
```

[29]: <tensorflow.python.keras.callbacks.History at 0x7f5b2b3105b0>

Our network now seems to peak at ~71% test accuracy, a 8% absolute drop. This drop is mostly due to the fact that we are now trying to compress a lot of information (enough information to recover the separation hyperplanes of 46 classes) into an intermediate space that is too low-dimensional. The network is able to cram *most* of the necessary information into these 8-dimensional representations, but not all of it.

Further experiments

- Try using larger or smaller layers: 32 units, 128 units...
- We were using two hidden layers. Now try to use a single hidden layer, or three hidden layers.

Wrapping up Here’s what you should take away from this example:

- If you are trying to classify data points between N classes, your network should end with a **Dense** layer of size N.
- In a single-label, multi-class classification problem, your network should end with a **softmax** activation, so that it will output a probability distribution over the N output classes.
- *Categorical crossentropy* is almost always the loss function you should use for such problems. It minimizes the distance between the probability distributions output by the network, and the true distribution of the targets.
- There are two ways to handle labels in multi-class classification: **** Encoding the labels via “categorical encoding” (also known as “one-hot encoding”) and using `categorical_crossentropy` as your loss function. ** Encoding the labels as integers and using the `sparse_categorical_crossentropy` loss function.**
- If you need to classify data into a large number of categories, then you should avoid creating information bottlenecks in your network by having intermediate layers that are too small.

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

[]: