

< Deep Learning - PART3 TF2 RNNs >

Ch 6. RNNs Workshop 2 - NLP - IMDB : Word Embeddings

2021/10/01

[Reference] : François Chollet, **Deep Learning with Python**, Chapter 6, Section 1, Manning, 2018.

[Code] : <https://github.com/fchollet/deep-learning-with-python-notebooks>
(<https://github.com/fchollet/deep-learning-with-python-notebooks>)

Another popular and powerful way to associate a vector with a word is the use of dense "word vectors", also called "word embeddings". While the vectors obtained through one-hot encoding are binary, sparse (mostly made of zeros) and very high-dimensional (same dimensionality as the number of words in the vocabulary), "word embeddings" are low-dimensional floating point vectors (i.e. "dense" vectors, as opposed to sparse vectors). Unlike word vectors obtained via one-hot encoding, word embeddings are learned from data. It is common to see word embeddings that are 256-dimensional, 512-dimensional, or 1024-dimensional when dealing with very large vocabularies. On the other hand, one-hot encoding words generally leads to vectors that are 20,000-dimensional or higher (capturing a vocabulary of 20,000 token in this case). So, word embeddings pack more information into far fewer dimensions.

[1. Learning word embeddings with the Embedding layer](#)

[2. Using pre-trained word embeddings](#)

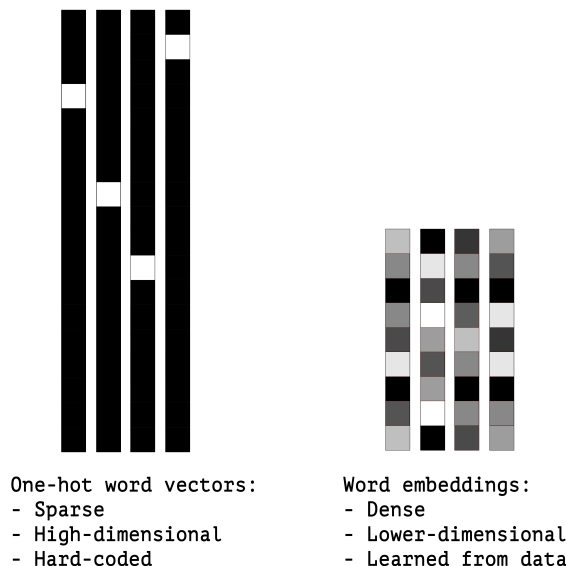
In [1]:



```
1 import tensorflow as tf
2 tf.__version__
```

Out[1]:

'2.4.1'



There are two ways to obtain word embeddings:

- Learn word embeddings jointly with the main task you care about (e.g. document classification or sentiment prediction). In this setup, you would start with random word vectors, then learn your word vectors in the same way that you learn the weights of a neural network.
- Load into your model word embeddings that were pre-computed using a different machine learning task than the one you are trying to solve. These are called "pre-trained word embeddings".

Let's take a look at both.

1. Learning word embeddings with the Embedding layer

The simplest way to associate a dense vector to a word would be to pick the vector at random. The problem with this approach is that the resulting embedding space would have no structure: for instance, the words "accurate" and "exact" may end up with completely different embeddings, even though they are interchangeable in most sentences. It would be very difficult for a deep neural network to make sense of such a noisy, unstructured embedding space.

To get a bit more abstract: the geometric relationships between word vectors should reflect the semantic relationships between these words. Word embeddings are meant to map human language into a geometric space. For instance, in a reasonable embedding space, we would expect synonyms to be embedded into similar word vectors, and in general we would expect the geometric distance (e.g. L2 distance) between any two word vectors to relate to the semantic distance of the associated words (words meaning very different things would be embedded to points far away from each other, while related words would be closer). Even beyond mere distance, we may want specific **directions** in the embedding space to be meaningful.

In real-world word embedding spaces, common examples of meaningful geometric transformations are "gender vectors" and "plural vector". For instance, by adding a "female vector" to the vector "king", one obtain the vector "queen". By adding a "plural vector", one obtain "kings". Word embedding spaces typically feature thousands of such interpretable and potentially useful vectors.

Is there some "ideal" word embedding space that would perfectly map human language and could be used for any natural language processing task? Possibly, but in any case, we have yet to compute anything of the sort. Also, there isn't such a thing as "human language", there are many different languages and they are not isomorphic, as a language is the reflection of a specific culture and a specific context. But more pragmatically, what makes a good word embedding space depends heavily on your task: the perfect word embedding space for an English-language movie review sentiment analysis model may look very different from the perfect embedding space for an English-language legal document classification model, because the importance of certain semantic relationships varies from task to task.

It is thus reasonable to **learn** a new embedding space with every new task. Thankfully, backpropagation makes this really easy, and Keras makes it even easier. It's just about learning the weights of a layer: the `Embedding` layer.

In [4]:



```
1 from tensorflow.keras.layers import Embedding
2
3 # The Embedding layer takes at least two arguments:
4 # the number of possible tokens, here 1000 (1 + maximum word index),
5 # and the dimensionality of the embeddings, here 64.
6 embedding_layer = Embedding(1000, 64)
```

The `Embedding` layer is best understood as a dictionary mapping integer indices (which stand for specific words) to dense vectors. It takes as input integers, it looks up these integers into an internal dictionary, and it returns the associated vectors. It's effectively a dictionary lookup.

The `Embedding` layer takes as input a 2D tensor of integers, of shape `(samples, sequence_length)`, where each entry is a sequence of integers. It can embed sequences of variable lengths, so for instance we could feed into our embedding layer above batches that could have shapes `(32, 10)` (batch of 32 sequences of length 10) or `(64, 15)` (batch of 64 sequences of length 15). All sequences in a batch must have the same length, though (since we need to pack them into a single tensor), so sequences that are shorter than others should be padded with zeros, and sequences that are longer should be truncated.

This layer returns a 3D floating point tensor, of shape `(samples, sequence_length, embedding_dimensionality)`. Such a 3D tensor can then be processed by a RNN layer or a 1D convolution layer (both will be introduced in the next sections).

When you instantiate an `Embedding` layer, its weights (its internal dictionary of token vectors) are initially random, just like with any other layer. During training, these word vectors will be gradually adjusted via backpropagation, structuring the space into something that the

downstream model can exploit. Once fully trained, your embedding space will show a lot of structure -- a kind of structure specialized for the specific problem you were training your model for.

Let's apply this idea to the IMDB movie review sentiment prediction task that you are already familiar with. Let's quickly prepare the data. We will restrict the movie reviews to the top 10,000 most common words (like we did the first time we worked with this dataset), and cut the reviews after only 20 words. Our network will simply learn 8-dimensional embeddings for each of the 10,000 words, turn the input integer sequences (2D integer tensor) into embedded sequences (3D float tensor), flatten the tensor to 2D, and train a single Dense layer on top for classification.

In [5]:



```
1 from tensorflow.keras.datasets import imdb
2 from tensorflow.keras import preprocessing
3
4 # Number of words to consider as features
5 max_features = 10000
6 # Cut texts after this number of words
7 # (among top max_features most common words)
8 maxlen = 20
9
10 # Load the data as lists of integers.
11 (x_train, y_train), (x_test, y_test) = imdb.load_data(num_words=max_features)
12
13 # This turns our lists of integers
14 # into a 2D integer tensor of shape `(samples, maxlen)`
15 x_train = preprocessing.sequence.pad_sequences(x_train, maxlen=maxlen)
16 x_test = preprocessing.sequence.pad_sequences(x_test, maxlen=maxlen)
```

In [6]:



```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Flatten, Dense
3
4 model = Sequential()
5 # We specify the maximum input length to our Embedding layer
6 # so we can later flatten the embedded inputs
7 model.add(Embedding(10000, 8, input_length=maxlen))
8 # After the Embedding layer,
9 # our activations have shape `(samples, maxlen, 8)`.
10
11 # We flatten the 3D tensor of embeddings
12 # into a 2D tensor of shape `(samples, maxlen * 8)`
13 model.add(Flatten())
14
15 # We add the classifier on top
16 model.add(Dense(1, activation='sigmoid'))
17 model.compile(optimizer='rmsprop', loss='binary_crossentropy', metrics=['ac
18 model.summary()
19
20 history = model.fit(x_train, y_train,
21                     epochs=10,
22                     batch_size=32,
23                     validation_split=0.2)
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
embedding_1 (Embedding)	(None, 20, 8)	80000
flatten (Flatten)	(None, 160)	0
dense (Dense)	(None, 1)	161

=====
Total params: 80,161
Trainable params: 80,161
Non-trainable params: 0

Epoch 1/10
625/625 [=====] - 1s 1ms/step - loss: 0.6826 - acc: 0.5812 - val_loss: 0.6011 - val_acc: 0.7092
Epoch 2/10
625/625 [=====] - 0s 765us/step - loss: 0.5522 - acc: 0.7508 - val_loss: 0.5158 - val_acc: 0.7378
Epoch 3/10
625/625 [=====] - 0s 790us/step - loss: 0.4573 - acc: 0.7910 - val_loss: 0.4946 - val_acc: 0.7482
Epoch 4/10

```
625/625 [=====] - 0s 792us/step - loss: 0.4203 - acc: 0.8070 - val_loss: 0.4904 - val_acc: 0.7548
Epoch 5/10
625/625 [=====] - 0s 768us/step - loss: 0.4008 - acc: 0.8174 - val_loss: 0.4911 - val_acc: 0.7560
Epoch 6/10
625/625 [=====] - 0s 769us/step - loss: 0.3763 - acc: 0.8346 - val_loss: 0.4926 - val_acc: 0.7588
Epoch 7/10
625/625 [=====] - 1s 821us/step - loss: 0.3638 - acc: 0.8410 - val_loss: 0.4976 - val_acc: 0.7578
Epoch 8/10
625/625 [=====] - 0s 759us/step - loss: 0.3407 - acc: 0.8505 - val_loss: 0.5020 - val_acc: 0.7570
Epoch 9/10
625/625 [=====] - 0s 723us/step - loss: 0.3216 - acc: 0.8626 - val_loss: 0.5089 - val_acc: 0.7544
Epoch 10/10
625/625 [=====] - 1s 893us/step - loss: 0.3051 - acc: 0.8714 - val_loss: 0.5127 - val_acc: 0.7554
```

We get to a validation accuracy of ~76%, which is pretty good considering that we only look at the first 20 words in every review. But note that merely flattening the embedded sequences and training a single Dense layer on top leads to a model that treats each word in the input sequence separately, without considering inter-word relationships and structure sentence (e.g. it would likely treat both *"this movie is shit"* and *"this movie is the shit"* as being negative "reviews"). It would be much better to add recurrent layers or 1D convolutional layers on top of the embedded sequences to learn features that take into account each sequence as a whole. That's what we will focus on in the next few sections.

2. Using pre-trained word embeddings

Sometimes, you have so little training data available that could never use your data alone to learn an appropriate task-specific embedding of your vocabulary. What to do then?

Instead of learning word embeddings jointly with the problem you want to solve, you could be loading embedding vectors from a pre-computed embedding space known to be highly structured and to exhibit useful properties -- that captures generic aspects of language structure. The rationale behind using pre-trained word embeddings in natural language processing is very much the same as for using pre-trained convnets in image classification: we don't have enough data available to learn truly powerful features on our own, but we expect the features that we need to be fairly generic, i.e. common visual features or semantic features. In this case it makes sense to reuse features learned on a different problem.

Such word embeddings are generally computed using word occurrence statistics (observations about what words co-occur in sentences or documents), using a variety of techniques, some involving neural networks, others not. The idea of a dense, low-dimensional embedding space for words, computed in an unsupervised way, was initially explored by Bengio et al. in the early 2000s, but it only started really taking off in research and industry applications after the release of

one of the most famous and successful word embedding scheme: the Word2Vec algorithm, developed by Mikolov at Google in 2013. Word2Vec dimensions capture specific semantic properties, e.g. gender.

There are various pre-computed databases of word embeddings that can download and start using in a Keras `Embedding` layer. Word2Vec is one of them. Another popular one is called "GloVe", developed by Stanford researchers in 2014. It stands for "Global Vectors for Word Representation", and it is an embedding technique based on factorizing a matrix of word co-occurrence statistics. Its developers have made available pre-computed embeddings for millions of English tokens, obtained from Wikipedia data or from Common Crawl data.

Let's take a look at how you can get started using GloVe embeddings in a Keras model. The same method will of course be valid for Word2Vec embeddings or any other word embedding database that you can download. We will also use this example to refresh the text tokenization techniques we introduced a few paragraphs ago: we will start from raw text, and work our way up.

Putting it all together: from raw text to word embeddings

We will be using a model similar to the one we just went over -- embedding sentences in sequences of vectors, flattening them and training a `Dense` layer on top. But we will do it using pre-trained word embeddings, and instead of using the pre-tokenized IMDB data packaged in Keras, we will start from scratch, by downloading the original text data.

Download the IMDB data as raw text

First, head to <http://ai.stanford.edu/~amaas/data/sentiment/> and download the raw IMDB dataset (if the URL isn't working anymore, just Google "IMDB dataset"). Uncompress it.

Now let's collect the individual training reviews into a list of strings, one string per review, and let's also collect the review labels (positive / negative) into a `labels` list:

In [5]:



```
1 import os
2
3 imdb_dir = './aclImdb' # directory for file folder : aclImdb
4 train_dir = os.path.join(imdb_dir, 'train')
5
6 labels = []
7 texts = []
8
9 for label_type in ['neg', 'pos']:
10     dir_name = os.path.join(train_dir, label_type)
11     for fname in os.listdir(dir_name):
12         if fname[-4:] == '.txt':
13             f = open(os.path.join(dir_name, fname), encoding='utf8')
14             texts.append(f.read())
15             f.close()
16             if label_type == 'neg':
17                 labels.append(0)
18             else:
19                 labels.append(1)
```

Tokenize the data

Let's vectorize the texts we collected, and prepare a training and validation split. We will merely be using the concepts we introduced earlier in this section.

Because pre-trained word embeddings are meant to be particularly useful on problems where little training data is available (otherwise, task-specific embeddings are likely to outperform them), we will add the following twist: we restrict the training data to its first 200 samples. So we will be learning to classify movie reviews after looking at just 200 examples...

In [6]:



```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
3 import numpy as np
4
5 maxlen = 100 # We will cut reviews after 100 words
6 training_samples = 200 # We will be training on 200 samples
7 validation_samples = 10000 # We will be validating on 10000 samples
8 max_words = 10000 # We will only consider the top 10,000 words in the data
9
10 tokenizer = Tokenizer(num_words=max_words)
11 tokenizer.fit_on_texts(texts)
12 sequences = tokenizer.texts_to_sequences(texts)
13
14 word_index = tokenizer.word_index
15 print('Found %s unique tokens.' % len(word_index))
16
17 data = pad_sequences(sequences, maxlen=maxlen)
18
19 labels = np.asarray(labels)
20 print('Shape of data tensor:', data.shape)
21 print('Shape of label tensor:', labels.shape)
22
23 # Split the data into a training set and a validation set
24 # But first, shuffle the data, since we started from data
25 # where sample are ordered (all negative first, then all positive).
26 indices = np.arange(data.shape[0])
27 np.random.shuffle(indices)
28 data = data[indices]
29 labels = labels[indices]
30
31 x_train = data[:training_samples]
32 y_train = labels[:training_samples]
33 x_val = data[training_samples: training_samples + validation_samples]
34 y_val = labels[training_samples: training_samples + validation_samples]
```

Found 88582 unique tokens.
Shape of data tensor: (25000, 100)
Shape of label tensor: (25000,)

Download the GloVe word embeddings

Head to <https://nlp.stanford.edu/projects/glove/> (where you can learn more about the GloVe algorithm), and download the pre-computed embeddings from 2014 English Wikipedia. It's a 822MB zip file named `glove.6B.zip`, containing 100-dimensional embedding vectors for 400,000 words (or non-word tokens). Un-zip it.

Pre-process the embeddings

Let's parse the un-zipped file (it's a `txt` file) to build an index mapping words (as strings) to their vector representation (as number vectors).

In [7]:



```
1 glove_dir = './glove.6B'
2
3 embeddings_index = {}
4 f = open(os.path.join(glove_dir, 'glove.6B.100d.txt'), encoding='utf-8')
5 for line in f:
6     values = line.split()
7     word = values[0]
8     coefs = np.asarray(values[1:], dtype='float32')
9     embeddings_index[word] = coefs
10 f.close()
11
12 print('Found %s word vectors.' % len(embeddings_index))
```

Found 400000 word vectors.

Now let's build an embedding matrix that we will be able to load into an `Embedding` layer. It must be a matrix of shape `(max_words, embedding_dim)`, where each entry `i` contains the `embedding_dim`-dimensional vector for the word of index `i` in our reference word index (built during tokenization). Note that the index `0` is not supposed to stand for any word or token -- it's a placeholder.

In [8]:



```
1 embedding_dim = 100
2
3 embedding_matrix = np.zeros((max_words, embedding_dim))
4 for word, i in word_index.items():
5     embedding_vector = embeddings_index.get(word)
6     if i < max_words:
7         if embedding_vector is not None:
8             # Words not found in embedding index will be all-zeros.
9             embedding_matrix[i] = embedding_vector
```

Define a model

We will be using the same model architecture as before:

In [9]:



```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Embedding, Flatten, Dense
3
4 model = Sequential()
5 model.add(Embedding(max_words, embedding_dim, input_length=maxlen))
6 model.add(Flatten())
7 model.add(Dense(32, activation='relu'))
8 model.add(Dense(1, activation='sigmoid'))
9 model.summary()
```

Model: "sequential_1"

Layer (type)	Output Shape	Param #
=====		
embedding_2 (Embedding)	(None, 100, 100)	1000000

flatten_1 (Flatten)	(None, 10000)	0

dense_1 (Dense)	(None, 32)	320032

dense_2 (Dense)	(None, 1)	33
=====		
Total params: 1,320,065		
Trainable params: 1,320,065		
Non-trainable params: 0		

Load the GloVe embeddings in the model

The `Embedding` layer has a single weight matrix: a 2D float matrix where each entry `i` is the word vector meant to be associated with index `i`. Simple enough. Let's just load the GloVe matrix we prepared into our `Embedding` layer, the first layer in our model:

In [10]:



```
1 model.layers[0].set_weights([embedding_matrix])
2 model.layers[0].trainable = False
```

Additionally, we freeze the embedding layer (we set its `trainable` attribute to `False`), following the same rationale as what you are already familiar with in the context of pre-trained convnet features: when parts of a model are pre-trained (like our `Embedding` layer), and parts are randomly initialized (like our classifier), the pre-trained parts should not be updated during training to avoid forgetting what they already know. The large gradient update triggered by the randomly initialized layers would be very disruptive to the already learned features.

Train and evaluate

Let's compile our model and train it:

In [11]:



```
1 model.compile(optimizer='rmsprop',
2               loss='binary_crossentropy',
3               metrics=['acc'])
4 history = model.fit(x_train, y_train,
5                     epochs=10,
6                     batch_size=32,
7                     validation_data=(x_val, y_val))
8 model.save_weights('pre_trained_glove_model.h5')
```

Train on 200 samples, validate on 10000 samples

Epoch 1/10

200/200 [=====] - 1s 6ms/sample - loss: 1.5738 - acc: 0.4800 - val_loss: 0.7025 - val_acc: 0.4944

Epoch 2/10

200/200 [=====] - 1s 4ms/sample - loss: 0.5967 - acc: 0.7400 - val_loss: 1.0009 - val_acc: 0.4960

Epoch 3/10

200/200 [=====] - 1s 4ms/sample - loss: 0.5493 - acc: 0.7000 - val_loss: 0.7034 - val_acc: 0.5347

Epoch 4/10

200/200 [=====] - 1s 4ms/sample - loss: 0.3423 - acc: 0.8550 - val_loss: 0.8030 - val_acc: 0.5150

Epoch 5/10

200/200 [=====] - 1s 4ms/sample - loss: 0.3813 - acc: 0.7700 - val_loss: 0.7004 - val_acc: 0.5548

Epoch 6/10

200/200 [=====] - 1s 3ms/sample - loss: 0.1321 - acc: 0.9950 - val_loss: 0.7264 - val_acc: 0.5550

Epoch 7/10

200/200 [=====] - 1s 3ms/sample - loss: 0.0850 - acc: 1.0000 - val_loss: 1.5290 - val_acc: 0.5067

Epoch 8/10

200/200 [=====] - 1s 3ms/sample - loss: 0.3065 - acc: 0.8400 - val_loss: 0.7852 - val_acc: 0.5451

Epoch 9/10

200/200 [=====] - 1s 3ms/sample - loss: 0.0427 - acc: 1.0000 - val_loss: 0.7541 - val_acc: 0.5568

Epoch 10/10

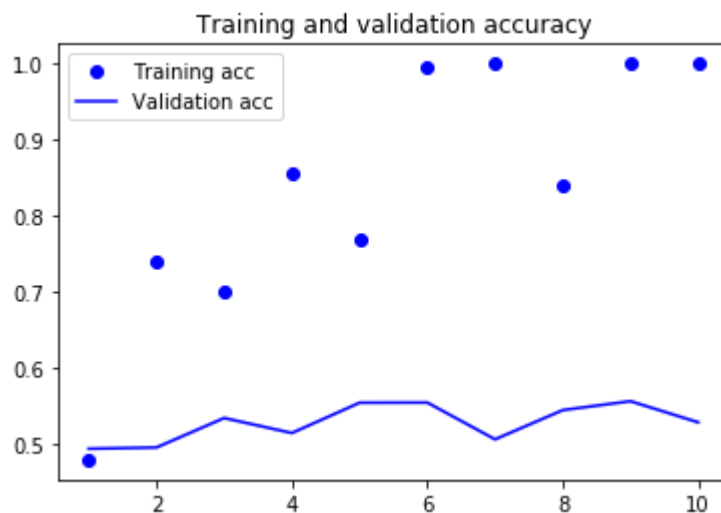
200/200 [=====] - 1s 3ms/sample - loss: 0.0293 - acc: 1.0000 - val_loss: 0.8507 - val_acc: 0.5290

Let's plot its performance over time:

In [12]:



```
1 import matplotlib.pyplot as plt
2 %matplotlib inline
3
4 acc = history.history['acc']
5 val_acc = history.history['val_acc']
6 loss = history.history['loss']
7 val_loss = history.history['val_loss']
8
9 epochs = range(1, len(acc) + 1)
10
11 plt.plot(epochs, acc, 'bo', label='Training acc')
12 plt.plot(epochs, val_acc, 'b', label='Validation acc')
13 plt.title('Training and validation accuracy')
14 plt.legend()
15
16 plt.figure()
17
18 plt.plot(epochs, loss, 'bo', label='Training loss')
19 plt.plot(epochs, val_loss, 'b', label='Validation loss')
20 plt.title('Training and validation loss')
21 plt.legend()
22
23 plt.show()
```





The model quickly starts overfitting, unsurprisingly given the small number of training samples. Validation accuracy has high variance for the same reason, but seems to reach high 50s.

Note that your mileage may vary: since we have so few training samples, performance is heavily dependent on which exact 200 samples we picked, and we picked them at random. If it worked really poorly for you, try picking a different random set of 200 samples, just for the sake of the exercise (in real life you don't get to pick your training data).

We can also try to train the same model without loading the pre-trained word embeddings and without freezing the embedding layer. In that case, we would be learning a task-specific embedding of our input tokens, which is generally more powerful than pre-trained word embeddings when lots of data is available. However, in our case, we have only 200 training samples. Let's try it:

In [13]:



```
1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import Embedding, Flatten, Dense
3
4 model = Sequential()
5 model.add(Embedding(max_words, embedding_dim, input_length=maxlen))
6 model.add(Flatten())
7 model.add(Dense(32, activation='relu'))
8 model.add(Dense(1, activation='sigmoid'))
9 model.summary()
10
11 model.compile(optimizer='rmsprop',
12               loss='binary_crossentropy',
13               metrics=['acc'])
14 history = model.fit(x_train, y_train,
15                     epochs=10,
16                     batch_size=32,
17                     validation_data=(x_val, y_val))
```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
embedding_3 (Embedding)	(None, 100, 100)	1000000
flatten_2 (Flatten)	(None, 10000)	0
dense_3 (Dense)	(None, 32)	320032
dense_4 (Dense)	(None, 1)	33

=====
Total params: 1,320,065
Trainable params: 1,320,065
Non-trainable params: 0

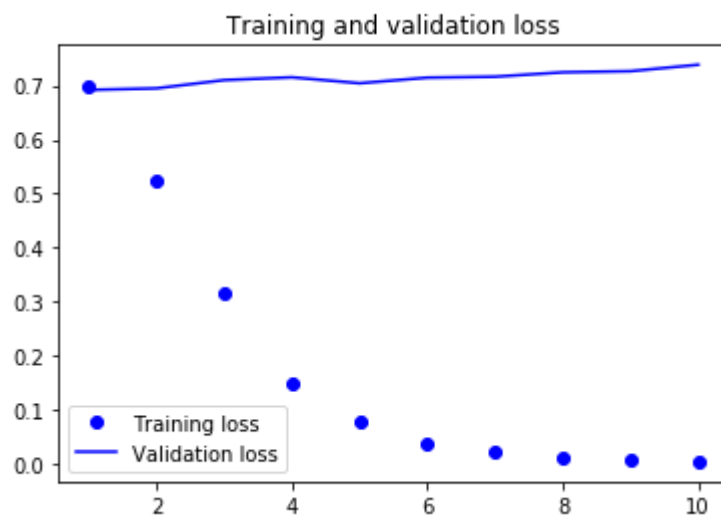
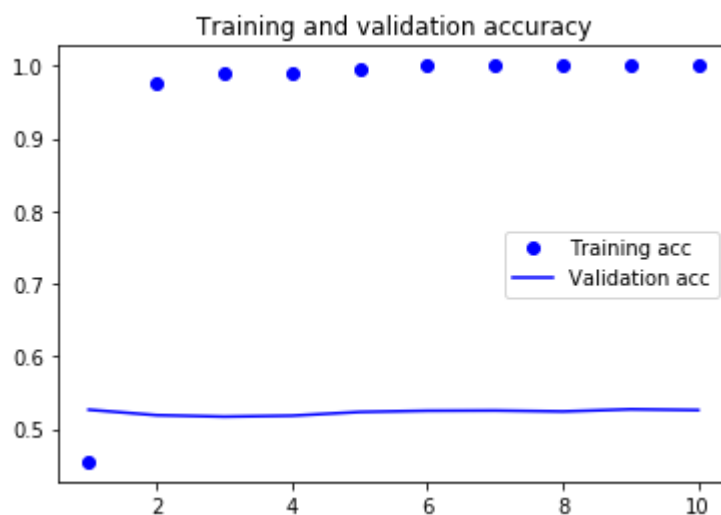
Train on 200 samples, validate on 10000 samples
Epoch 1/10
200/200 [=====] - 1s 7ms/sample - loss: 0.6968 - acc: 0.4550 - val_loss: 0.6920 - val_acc: 0.5265
Epoch 2/10
200/200 [=====] - 1s 4ms/sample - loss: 0.5245 - acc: 0.9750 - val_loss: 0.6949 - val_acc: 0.5190
Epoch 3/10
200/200 [=====] - 1s 4ms/sample - loss: 0.3176 - acc: 0.9900 - val_loss: 0.7101 - val_acc: 0.5170
Epoch 4/10
200/200 [=====] - 1s 4ms/sample - loss: 0.1494 - acc: 0.9900 - val_loss: 0.7153 - val_acc: 0.5183

Epoch 5/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0770 - acc: 0.9950 - val_loss: 0.7043 - val_acc: 0.5234
Epoch 6/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0376 - acc: 1.0000 - val_loss: 0.7147 - val_acc: 0.5251
Epoch 7/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0211 - acc: 1.0000 - val_loss: 0.7164 - val_acc: 0.5254
Epoch 8/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0121 - acc: 1.0000 - val_loss: 0.7245 - val_acc: 0.5241
Epoch 9/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0071 - acc: 1.0000 - val_loss: 0.7266 - val_acc: 0.5271
Epoch 10/10
200/200 [=====] - 1s 4ms/sample - loss: 0.0043 - acc: 1.0000 - val_loss: 0.7385 - val_acc: 0.5260

In [14]:



```
1 acc = history.history['acc']
2 val_acc = history.history['val_acc']
3 loss = history.history['loss']
4 val_loss = history.history['val_loss']
5
6 epochs = range(1, len(acc) + 1)
7
8 plt.plot(epochs, acc, 'bo', label='Training acc')
9 plt.plot(epochs, val_acc, 'b', label='Validation acc')
10 plt.title('Training and validation accuracy')
11 plt.legend()
12
13 plt.figure()
14
15 plt.plot(epochs, loss, 'bo', label='Training loss')
16 plt.plot(epochs, val_loss, 'b', label='Validation loss')
17 plt.title('Training and validation loss')
18 plt.legend()
19
20 plt.show()
```



Validation accuracy stalls in the low 50s. So in our case, pre-trained word embeddings does outperform jointly learned embeddings. If you increase the number of training samples, this will quickly stop being the case -- try it as an exercise.

Finally, let's evaluate the model on the test data. First, we will need to tokenize the test data:

In [15]:

```
1 test_dir = os.path.join(imdb_dir, 'test')
2
3 labels = []
4 texts = []
5
6 for label_type in ['neg', 'pos']:
7     dir_name = os.path.join(test_dir, label_type)
8     for fname in sorted(os.listdir(dir_name)):
9         if fname[-4:] == '.txt':
10             f = open(os.path.join(dir_name, fname), encoding='utf-8')
11             texts.append(f.read())
12             f.close()
13             if label_type == 'neg':
14                 labels.append(0)
15             else:
16                 labels.append(1)
17
18 sequences = tokenizer.texts_to_sequences(texts)
19 x_test = pad_sequences(sequences, maxlen=maxlen)
20 y_test = np.asarray(labels)
```

And let's load and evaluate the first model:

In [16]:

```
1 model.load_weights('pre_trained_glove_model.h5')
2 model.evaluate(x_test, y_test, verbose=2)
```

25000/1 - 1s - loss: 0.5036 - acc: 0.5312

Out[16]:

[0.8453588010597229, 0.5312]

We get an appalling test accuracy of 54%. Working with just a handful of training samples is hard!