# How to Develop a Word Embedding Model for Predicting Movie Review Sentiment

by **Jason Brownlee** in **Deep Learning for Natural Language Processing**

**Develop a Deep Learning Model to Automatically Classify Movie Reviews as Positive or Negative in Python with Keras, Step-by-Step.**

Word embeddings are a technique for representing text where different words with similar meaning have a similar real-valued vector representation.

They are a key breakthrough that has led to great performance of neural network models on a suite of challenging natural language processing problems.

In this tutorial, you will discover how to develop word embedding models for neural networks to classify movie reviews.

After completing this tutorial, you will know:

- How to prepare movie review text data for classification with deep learning methods.
- How to learn a word embedding as part of fitting a deep learning model.
- How to learn a standalone word embedding and how to use a pre-trained embedding in a neural network model.

Let's get started.

**Note**: This is an excerpt from: "Deep Learning for Natural Language Processing".

Take a look, if you want more step-by-step tutorials on getting the most out of deep learning methods when working with text data.

How to Develop a Word Embedding Model for Predicting Movie Review Sentiment

Photo by Katrina Br*?#*!@nd, some rights reserved.

# Tutorial Overview

This tutorial is divided into 5 parts; they are:

1. Movie Review Dataset
2. Data Preparation
3. Train Embedding Layer
4. Train word2vec Embedding
5. Use Pre-trained Embedding

## Python Environment

This tutorial assumes you have a Python SciPy environment installed, ideally with Python 3.

You must have Keras (2.2 or higher) installed with either the TensorFlow or Theano backend.

The tutorial also assumes you have scikit-learn, Pandas, NumPy, and Matplotlib installed.

If you need help with your environment, see this tutorial:

- [How to Setup a Python Environment for Machine Learning and Deep Learning with Anaconda](#)
A GPU is not required for this tutorial, nevertheless, you can access GPUs cheaply on Amazon Web Services. Learn how in this tutorial:

- [How to Setup Amazon AWS EC2 GPUs to Train Keras Deep Learning Models (step-by-step)](#)
Let's dive in.

# 1. Movie Review Dataset

The Movie Review Data is a collection of movie reviews retrieved from the imdb.com website in the early 2000s by Bo Pang and Lillian Lee. The reviews were collected and made available as part of their research on natural language processing.

The reviews were originally released in 2002, but an updated and cleaned up version were released in 2004, referred to as "v2.0".

The dataset is comprised of 1,000 positive and 1,000 negative movie reviews drawn from an archive of the rec.arts.movies.reviews newsgroup hosted at [imdb.com](#). The authors refer to this dataset as the "polarity dataset."
*Our data contains 1000 positive and 1000 negative reviews all written before 2002, with a cap of 20 reviews per author (312 authors total) per category. We refer to this corpus as the polarity dataset.*

— [A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts](#), 2004.
The data has been cleaned up somewhat, for example:

- The dataset is comprised of only English reviews.
- All text has been converted to lowercase.
- There is white space around punctuation like periods, commas, and brackets.
- Text has been split into one sentence per line.

The data has been used for a few related natural language processing tasks. For classification, the performance of machine learning models (such as Support Vector Machines) on the data is in the range of high 70% to low 80% (e.g. 78%-82%).

More sophisticated data preparation may see results as high as 86% with 10-fold cross validation. This gives us a ballpark of low-to-mid 80s if we were looking to use this dataset in experiments of modern methods.

*… depending on choice of downstream polarity classifier, we can achieve highly statistically significant improvement (from 82.8% to 86.4%)*

— A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts, 2004.
You can download the dataset from here:

- Movie Review Polarity Dataset (review_polarity.tar.gz, 3MB)
  After unzipping the file, you will have a directory called "*txt_sentoken*" with two sub-directories containing the text "*neg*" and "*pos*" for negative and positive reviews. Reviews are stored one per file with a naming convention cv000 to cv999 for each neg and pos. Next, let's look at loading and preparing the text data.

# 2. Data Preparation

In this section, we will look at 3 things:

1. Separation of data into training and test sets.
2. Loading and cleaning the data to remove punctuation and numbers.
3. Defining a vocabulary of preferred words.

## Split into Train and Test Sets

We are pretending that we are developing a system that can predict the sentiment of a textual movie review as either positive or negative.

This means that after the model is developed, we will need to make predictions on new textual reviews. This will require all of the same data preparation to be performed on those new reviews as is performed on the training data for the model.

We will ensure that this constraint is built into the evaluation of our models by splitting the training and test datasets prior to any data preparation. This means that any knowledge in

the data in the test set that could help us better prepare the data (e.g. the words used) are unavailable in the preparation of data used for training the model.

That being said, we will use the last 100 positive reviews and the last 100 negative reviews as a test set (100 reviews) and the remaining 1,800 reviews as the training dataset.

This is a 90% train, 10% split of the data.

The split can be imposed easily by using the filenames of the reviews where reviews named 000 to 899 are for training data and reviews named 900 onwards are for test.

## Loading and Cleaning Reviews

The text data is already pretty clean; not much preparation is required.

If you are new to cleaning text data, see this post:

- [How to Clean Text for Machine Learning with Python](#)

Without getting bogged down too much in the details, we will prepare the data using the following way:

- Split tokens on white space.
- Remove all punctuation from words.
- Remove all words that are not purely comprised of alphabetical characters.
- Remove all words that are known stop words.
- Remove all words that have a length <= 1 character.

We can put all of these steps into a function called *clean_doc()* that takes as an argument the raw text loaded from a file and returns a list of cleaned tokens. We can also define a function *load_doc()* that loads a document from file ready for use with the *clean_doc()* function.
An example of cleaning the first positive review is listed below.

```
 1  from nltk.corpus import stopwords
 2  import string
 3
 4  # load doc into memory
 5  def load_doc(filename):
 6          # open the file as read only
 7          file = open(filename, 'r')
 8          # read all text
 9          text = file.read()
10          # close the file
11          file.close()
12          return text
13
14  # turn a doc into clean tokens
```

```
15 def clean_doc(doc):
16          # split into tokens by white space
17          tokens = doc.split()
18          # remove punctuation from each token
19          table = str.maketrans('', '', string.punctuation)
20          tokens = [w.translate(table) for w in tokens]
21          # remove remaining tokens that are not alphabetic
22          tokens = [word for word in tokens if word.isalpha()]
23          # filter out stop words
24          stop_words = set(stopwords.words('english'))
25          tokens = [w for w in tokens if not w in stop_words]
26          # filter out short tokens
27          tokens = [word for word in tokens if len(word) > 1]
28          return tokens
29
30 # load the document
31 filename = 'txt_sentoken/pos/cv000_29590.txt'
32 text = load_doc(filename)
33 tokens = clean_doc(text)
34 print(tokens)
```

Running the example prints a long list of clean tokens.

There are many more cleaning steps we may want to explore and I leave them as further exercises.

I'd love to see what you can come up with.
Post your approaches and findings in the comments at the end.

```
    …
1   'creepy', 'place', 'even', 'acting', 'hell', 'solid', 'dreamy', 'depp', 'turning', 'typically', 'strong', 'performance', 'deftly',
2   'handling', 'british', 'accent', 'ians', 'holm', 'joe', 'goulds', 'secret', 'richardson', 'dalmatians', 'log', 'great', 'supporting',
    'roles', 'big', 'surprise', 'graham', 'cringed', 'first', 'time', 'opened', 'mouth', 'imagining', 'attempt', 'irish', 'accent',
    'actually', 'wasnt', 'half', 'bad', 'film', 'however', 'good', 'strong', 'violencegore', 'sexuality', 'language', 'drug', 'content']
```

# Define a Vocabulary

It is important to define a vocabulary of known words when using a bag-of-words or embedding model.

The more words, the larger the representation of documents, therefore it is important to constrain the words to only those believed to be predictive. This is difficult to know beforehand and often it is important to test different hypotheses about how to construct a useful vocabulary.

We have already seen how we can remove punctuation and numbers from the vocabulary in the previous section. We can repeat this for all documents and build a set of all known words.

We can develop a vocabulary as a Counter, which is a dictionary mapping of words and their counts that allow us to easily update and query.

Each document can be added to the counter (a new function called *add_doc_to_vocab()*) and we can step over all of the reviews in the negative directory and then the positive directory (a new function called *process_docs()*).
The complete example is listed below.

```
1  from string import punctuation
2  from os import listdir
3  from collections import Counter
4  from nltk.corpus import stopwords
5
6  # load doc into memory
7  def load_doc(filename):
8          # open the file as read only
9          file = open(filename, 'r')
10         # read all text
11         text = file.read()
12         # close the file
13         file.close()
14         return text
15
16 # turn a doc into clean tokens
17 def clean_doc(doc):
18         # split into tokens by white space
19         tokens = doc.split()
20         # remove punctuation from each token
21         table = str.maketrans('', '', punctuation)
22         tokens = [w.translate(table) for w in tokens]
23         # remove remaining tokens that are not alphabetic
24         tokens = [word for word in tokens if word.isalpha()]
25         # filter out stop words
26         stop_words = set(stopwords.words('english'))
27         tokens = [w for w in tokens if not w in stop_words]
28         # filter out short tokens
29         tokens = [word for word in tokens if len(word) > 1]
30         return tokens
31
32 # load doc and add to vocab
33 def add_doc_to_vocab(filename, vocab):
34         # load doc
35         doc = load_doc(filename)
36         # clean doc
37         tokens = clean_doc(doc)
38         # update counts
39         vocab.update(tokens)
40
41 # load all docs in a directory
42 def process_docs(directory, vocab, is_trian):
43         # walk through all files in the folder
44         for filename in listdir(directory):
45                 # skip any reviews in the test set
46                 if is_trian and filename.startswith('cv9'):
47                         continue
48                 if not is_trian and not filename.startswith('cv9'):
49                         continue
```

```
50                      # create the full path of the file to open
51                      path = directory + '/' + filename
52                      # add doc to vocab
53                      add_doc_to_vocab(path, vocab)
54
55 # define vocab
56 vocab = Counter()
57 # add all docs to vocab
58 process_docs('txt_sentoken/neg', vocab, True)
59 process_docs('txt_sentoken/pos', vocab, True)
60 # print the size of the vocab
61 print(len(vocab))
62 # print the top words in the vocab
63 print(vocab.most_common(50))
```

Running the example shows that we have a vocabulary of 44,276 words.

We also can see a sample of the top 50 most used words in the movie reviews.

Note, that this vocabulary was constructed based on only those reviews in the training dataset.

```
  44276
  [('film', 7983), ('one', 4946), ('movie', 4826), ('like', 3201), ('even', 2262), ('good', 2080), ('time', 2041), ('story', 1907),
  ('films', 1873), ('would', 1844), ('much', 1824), ('also', 1757), ('characters', 1735), ('get', 1724), ('character', 1703),
1 ('two', 1643), ('first', 1588), ('see', 1557), ('way', 1515), ('well', 1511), ('make', 1418), ('really', 1407), ('little', 1351),
2 ('life', 1334), ('plot', 1288), ('people', 1269), ('could', 1248), ('bad', 1248), ('scene', 1241), ('movies', 1238), ('never',
  1201), ('best', 1179), ('new', 1140), ('scenes', 1135), ('man', 1131), ('many', 1130), ('doesnt', 1118), ('know', 1092),
  ('dont', 1086), ('hes', 1024), ('great', 1014), ('another', 992), ('action', 985), ('love', 977), ('us', 967), ('go', 952),
  ('director', 948), ('end', 946), ('something', 945), ('still', 936)]
```

We can step through the vocabulary and remove all words that have a low occurrence, such as only being used once or twice in all reviews.

For example, the following snippet will retrieve only the tokens that of appears 2 or more times in all reviews.

```
1 # keep tokens with a min occurrence
2 min_occurane = 2
3 tokens = [k for k,c in vocab.items() if c >= min_occurane]
4 print(len(tokens))
```

Running the above example with this addition shows that the vocabulary size drops by a little more than half its size from 44,276 to 25,767 words.

```
1 25767
```

Finally, the vocabulary can be saved to a new file called *vocab.txt* that we can later load and use to filter movie reviews prior to encoding them for modeling. We define a new function called *save_list()* that saves the vocabulary to file, with one word per file.
For example:

```
1 # save list to file
2 def save_list(lines, filename):
3           # convert lines to a single blob of text
4           data = '\n'.join(lines)
```

```
 5          # open file
 6          file = open(filename, 'w')
 7          # write text
 8          file.write(data)
 9          # close file
10          file.close()
11
12 # save tokens to a vocabulary file
13 save_list(tokens, 'vocab.txt')
```

Running the min occurrence filter on the vocabulary and saving it to file, you should now have a new file called *vocab.txt* with only the words we are interested in.
The order of words in your file will differ, but should look something like the following:

```
 1 aberdeen
 2 dupe
 3 burt
 4 libido
 5 hamlet
 6 arlene
 7 available
 8 corners
 9 web
10 columbia
11 …
```

We are now ready to look at learning features from the reviews.

# 3. Train Embedding Layer

In this section, we will learn a word embedding while training a neural network on the classification problem.

A word embedding is a way of representing text where each word in the vocabulary is represented by a real valued vector in a high-dimensional space. The vectors are learned in such a way that words that have similar meanings will have similar representation in the vector space (close in the vector space). This is a more expressive representation for text than more classical methods like bag-of-words, where relationships between words or tokens are ignored, or forced in bigram and trigram approaches.

The real valued vector representation for words can be learned while training the neural network. We can do this in the Keras deep learning library using the [Embedding layer](#).
If you are new to word embeddings, see the post:

- [What Are Word Embeddings for Text?](#)
  If you are new to word embedding layers in Keras, see the post:

- [How to Use Word Embedding Layers for Deep Learning with Keras](#)

The first step is to load the vocabulary. We will use it to filter out words from movie reviews that we are not interested in.

If you have worked through the previous section, you should have a local file called '*vocab.txt*' with one word per line. We can load that file and build a vocabulary as a set for checking the validity of tokens.

```
1  # load doc into memory
2  def load_doc(filename):
3          # open the file as read only
4          file = open(filename, 'r')
5          # read all text
6          text = file.read()
7          # close the file
8          file.close()
9          return text
10
11 # load the vocabulary
12 vocab_filename = 'vocab.txt'
13 vocab = load_doc(vocab_filename)
14 vocab = vocab.split()
15 vocab = set(vocab)
```

Next, we need to load all of the training data movie reviews. For that we can adapt the *process_docs()* from the previous section to load the documents, clean them, and return them as a list of strings, with one document per string. We want each document to be a string for easy encoding as a sequence of integers later.

Cleaning the document involves splitting each review based on white space, removing punctuation, and then filtering out all tokens not in the vocabulary.

The updated *clean_doc()* function is listed below.

```
1  # turn a doc into clean tokens
2  def clean_doc(doc, vocab):
3          # split into tokens by white space
4          tokens = doc.split()
5          # remove punctuation from each token
6          table = str.maketrans('', '', punctuation)
7          tokens = [w.translate(table) for w in tokens]
8          # filter out tokens not in vocab
9          tokens = [w for w in tokens if w in vocab]
10         tokens = ' '.join(tokens)
11         return tokens
```

The updated *process_docs()* can then call the *clean_doc()* for each document on the '*pos*' and '*neg*' directories that are in our training dataset.

```
1  # load all docs in a directory
2  def process_docs(directory, vocab, is_trian):
3          documents = list()
4          # walk through all files in the folder
5          for filename in listdir(directory):
6                  # skip any reviews in the test set
7                  if is_trian and filename.startswith('cv9'):
```

```
 8                             continue
 9                 if not is_trian and not filename.startswith('cv9'):
10                             continue
11                 # create the full path of the file to open
12                 path = directory + '/' + filename
13                 # load the doc
14                 doc = load_doc(path)
15                 # clean doc
16                 tokens = clean_doc(doc, vocab)
17                 # add to list
18                 documents.append(tokens)
19         return documents
20
21 # load all training reviews
22 positive_docs = process_docs('txt_sentoken/pos', vocab, True)
23 negative_docs = process_docs('txt_sentoken/neg', vocab, True)
24 train_docs = negative_docs + positive_docs
```

The next step is to encode each document as a sequence of integers.

The Keras Embedding layer requires integer inputs where each integer maps to a single token that has a specific real-valued vector representation within the embedding. These vectors are random at the beginning of training, but during training become meaningful to the network.

We can encode the training documents as sequences of integers using the Tokenizer class in the Keras API.
First, we must construct an instance of the class then train it on all documents in the training dataset. In this case, it develops a vocabulary of all tokens in the training dataset and develops a consistent mapping from words in the vocabulary to unique integers. We could just as easily develop this mapping ourselves using our vocabulary file.

```
1 # create the tokenizer
2 tokenizer = Tokenizer()
3 # fit the tokenizer on the documents
4 tokenizer.fit_on_texts(train_docs)
```

Now that the mapping of words to integers has been prepared, we can use it to encode the reviews in the training dataset. We can do that by calling the *texts_to_sequences()* function on the Tokenizer.

```
1 # sequence encode
2 encoded_docs = tokenizer.texts_to_sequences(train_docs)
```

We also need to ensure that all documents have the same length.

This is a requirement of Keras for efficient computation. We could truncate reviews to the smallest size or zero-pad (pad with the value '0') reviews to the maximum length, or some hybrid. In this case, we will pad all reviews to the length of the longest review in the training dataset.

First, we can find the longest review using the *max()* function on the training dataset and take its length. We can then call the Keras function *pad_sequences()* to pad the sequences to the maximum length by adding 0 values on the end.

```
1 # pad sequences
2 max_length = max([len(s.split()) for s in train_docs])
3 Xtrain = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
```

Finally, we can define the class labels for the training dataset, needed to fit the supervised neural network model to predict the sentiment of reviews.

```
1 # define training labels
2 ytrain = array([0 for _ in range(900)] + [1 for _ in range(900)])
```

We can then encode and pad the test dataset, needed later to evaluate the model after we train it.

```
 1  # load all test reviews
 2  positive_docs = process_docs('txt_sentoken/pos', vocab, False)
 3  negative_docs = process_docs('txt_sentoken/neg', vocab, False)
 4  test_docs = negative_docs + positive_docs
 5  # sequence encode
 6  encoded_docs = tokenizer.texts_to_sequences(test_docs)
 7  # pad sequences
 8  Xtest = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
 9  # define test labels
10  ytest = array([0 for _ in range(100)] + [1 for _ in range(100)])
```

We are now ready to define our neural network model.

The model will use an Embedding layer as the first hidden layer. The Embedding requires the specification of the vocabulary size, the size of the real-valued vector space, and the maximum length of input documents.

The vocabulary size is the total number of words in our vocabulary, plus one for unknown words. This could be the vocab set length or the size of the vocab within the tokenizer used to integer encode the documents, for example:

```
1 # define vocabulary size (largest integer value)
2 vocab_size = len(tokenizer.word_index) + 1
```

We will use a 100-dimensional vector space, but you could try other values, such as 50 or 150. Finally, the maximum document length was calculated above in the *max_length* variable used during padding.

The complete model definition is listed below including the Embedding layer.

We use a Convolutional Neural Network (CNN) as they have proven to be successful at document classification problems. A conservative CNN configuration is used with 32 filters (parallel fields for processing words) and a kernel size of 8 with a rectified linear ('relu') activation function. This is followed by a pooling layer that reduces the output of the convolutional layer by half.

Next, the 2D output from the CNN part of the model is flattened to one long 2D vector to represent the 'features' extracted by the CNN. The back-end of the model is a standard Multilayer Perceptron layers to interpret the CNN features. The output layer uses a sigmoid activation function to output a value between 0 and 1 for the negative and positive sentiment in the review.

For more advice on effective deep learning model configuration for text classification, see the post:

[Best Practices for Document Classification with Deep Learning](#)

```
1 # define model
2 model = Sequential()
3 model.add(Embedding(vocab_size, 100, input_length=max_length))
4 model.add(Conv1D(filters=32, kernel_size=8, activation='relu'))
5 model.add(MaxPooling1D(pool_size=2))
6 model.add(Flatten())
7 model.add(Dense(10, activation='relu'))
8 model.add(Dense(1, activation='sigmoid'))
9 print(model.summary())
```

Running just this piece provides a summary of the defined network.

We can see that the Embedding layer expects documents with a length of 442 words as input and encodes each word in the document as a 100 element vector.

```
 1 _____
 2 Layer (type)              Output Shape          Param #
 3 =====================================================================
 4 embedding_1 (Embedding)     (None, 442, 100)        2576800
 5 _____
 6 conv1d_1 (Conv1D)           (None, 435, 32)          25632
 7 _____
 8 max_pooling1d_1 (MaxPooling1 (None, 217, 32)            0
 9 _____
10 flatten_1 (Flatten)         (None, 6944)               0
11 _____
12 dense_1 (Dense)             (None, 10)              69450
13 _____
14 dense_2 (Dense)             (None, 1)               11
15 =====================================================================
16 Total params: 2,671,893
17 Trainable params: 2,671,893
18 Non-trainable params: 0
19 _____
```

Next, we fit the network on the training data.

We use a binary cross entropy loss function because the problem we are learning is a binary classification problem. The efficient Adam implementation of stochastic gradient

descent is used and we keep track of accuracy in addition to loss during training. The model is trained for 10 epochs, or 10 passes through the training data.

The network configuration and training schedule were found with a little trial and error, but are by no means optimal for this problem. If you can get better results with a different configuration, let me know.

```
1 # compile network
2 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
3 # fit network
4 model.fit(Xtrain, ytrain, epochs=10, verbose=2)
```

After the model is fit, it is evaluated on the test dataset. This dataset contains words that we have not seen before and reviews not seen during training.

```
1 # evaluate
2 loss, acc = model.evaluate(Xtest, ytest, verbose=0)
3 print('Test Accuracy: %f' % (acc*100))
```

We can tie all of this together.

The complete code listing is provided below.

```
 1   from string import punctuation
 2   from os import listdir
 3   from numpy import array
 4   from keras.preprocessing.text import Tokenizer
 5   from keras.preprocessing.sequence import pad_sequences
 6   from keras.models import Sequential
 7   from keras.layers import Dense
 8   from keras.layers import Flatten
 9   from keras.layers import Embedding
10   from keras.layers.convolutional import Conv1D
11   from keras.layers.convolutional import MaxPooling1D
12
13   # load doc into memory
14   def load_doc(filename):
15           # open the file as read only
16           file = open(filename, 'r')
17           # read all text
18           text = file.read()
19           # close the file
20           file.close()
21           return text
22
23   # turn a doc into clean tokens
24   def clean_doc(doc, vocab):
25           # split into tokens by white space
26           tokens = doc.split()
27           # remove punctuation from each token
28           table = str.maketrans('', '', punctuation)
29           tokens = [w.translate(table) for w in tokens]
30           # filter out tokens not in vocab
31           tokens = [w for w in tokens if w in vocab]
32           tokens = ' '.join(tokens)
33           return tokens
34
35   # load all docs in a directory
```

```python
36  def process_docs(directory, vocab, is_trian):
37          documents = list()
38          # walk through all files in the folder
39          for filename in listdir(directory):
40                  # skip any reviews in the test set
41                  if is_trian and filename.startswith('cv9'):
42                          continue
43                  if not is_trian and not filename.startswith('cv9'):
44                          continue
45                  # create the full path of the file to open
46                  path = directory + '/' + filename
47                  # load the doc
48                  doc = load_doc(path)
49                  # clean doc
50                  tokens = clean_doc(doc, vocab)
51                  # add to list
52                  documents.append(tokens)
53          return documents
54
55  # load the vocabulary
56  vocab_filename = 'vocab.txt'
57  vocab = load_doc(vocab_filename)
58  vocab = vocab.split()
59  vocab = set(vocab)
60
61  # load all training reviews
62  positive_docs = process_docs('txt_sentoken/pos', vocab, True)
63  negative_docs = process_docs('txt_sentoken/neg', vocab, True)
64  train_docs = negative_docs + positive_docs
65
66  # create the tokenizer
67  tokenizer = Tokenizer()
68  # fit the tokenizer on the documents
69  tokenizer.fit_on_texts(train_docs)
70
71  # sequence encode
72  encoded_docs = tokenizer.texts_to_sequences(train_docs)
73  # pad sequences
74  max_length = max([len(s.split()) for s in train_docs])
75  Xtrain = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
76  # define training labels
77  ytrain = array([0 for _ in range(900)] + [1 for _ in range(900)])
78
79  # load all test reviews
80  positive_docs = process_docs('txt_sentoken/pos', vocab, False)
81  negative_docs = process_docs('txt_sentoken/neg', vocab, False)
82  test_docs = negative_docs + positive_docs
83  # sequence encode
84  encoded_docs = tokenizer.texts_to_sequences(test_docs)
85  # pad sequences
86  Xtest = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
87  # define test labels
88  ytest = array([0 for _ in range(100)] + [1 for _ in range(100)])
89
90  # define vocabulary size (largest integer value)
91  vocab_size = len(tokenizer.word_index) + 1
92
93  # define model
94  model = Sequential()
95  model.add(Embedding(vocab_size, 100, input_length=max_length))
```

```
 96  model.add(Conv1D(filters=32, kernel_size=8, activation='relu'))
 97  model.add(MaxPooling1D(pool_size=2))
 98  model.add(Flatten())
 99  model.add(Dense(10, activation='relu'))
100  model.add(Dense(1, activation='sigmoid'))
101  print(model.summary())
102  # compile network
103  model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
104  # fit network
105  model.fit(Xtrain, ytrain, epochs=10, verbose=2)
106  # evaluate
107  loss, acc = model.evaluate(Xtest, ytest, verbose=0)
108  print('Test Accuracy: %f' % (acc*100))
```

Running the example prints the loss and accuracy at the end of each training epoch. We can see that the model very quickly achieves 100% accuracy on the training dataset.

At the end of the run, the model achieves an accuracy of 84.5% on the test dataset, which is a great score.

Given the stochastic nature of neural networks, your specific results will vary. Consider running the example a few times and taking the average score as the skill of the model.

```
 1  …
 2  Epoch 6/10
 3  2s - loss: 0.0013 - acc: 1.0000
 4  Epoch 7/10
 5  2s - loss: 8.4573e-04 - acc: 1.0000
 6  Epoch 8/10
 7  2s - loss: 5.8323e-04 - acc: 1.0000
 8  Epoch 9/10
 9  2s - loss: 4.3155e-04 - acc: 1.0000
10  Epoch 10/10
11  2s - loss: 3.3083e-04 - acc: 1.0000
12  Test Accuracy: 84.500000
```

We have just seen an example of how we can learn a word embedding as part of fitting a neural network model.

Next, let's look at how we can efficiently learn a standalone embedding that we could later use in our neural network.

# 4. Train word2vec Embedding

In this section, we will discover how to learn a standalone word embedding using an efficient algorithm called word2vec.

A downside of learning a word embedding as part of the network is that it can be very slow, especially for very large text datasets.

The word2vec algorithm is an approach to learning a word embedding from a text corpus in a standalone way. The benefit of the method is that it can produce high-quality word embeddings very efficiently, in terms of space and time complexity.

The first step is to prepare the documents ready for learning the embedding.

This involves the same data cleaning steps from the previous section, namely splitting documents by their white space, removing punctuation, and filtering out tokens not in the vocabulary.

The word2vec algorithm processes documents sentence by sentence. This means we will preserve the sentence-based structure during cleaning.

We start by loading the vocabulary, as before.

```
1  # load doc into memory
2  def load_doc(filename):
3          # open the file as read only
4          file = open(filename, 'r')
5          # read all text
6          text = file.read()
7          # close the file
8          file.close()
9          return text
10
11 # load the vocabulary
12 vocab_filename = 'vocab.txt'
13 vocab = load_doc(vocab_filename)
14 vocab = vocab.split()
15 vocab = set(vocab)
```

Next, we define a function named *doc_to_clean_lines()* to clean a loaded document line by line and return a list of the cleaned lines.

```
1  # turn a doc into clean tokens
2  def doc_to_clean_lines(doc, vocab):
3          clean_lines = list()
4          lines = doc.splitlines()
5          for line in lines:
6                  # split into tokens by white space
7                  tokens = line.split()
8                  # remove punctuation from each token
9                  table = str.maketrans('', '', punctuation)
10                 tokens = [w.translate(table) for w in tokens]
11                 # filter out tokens not in vocab
12                 tokens = [w for w in tokens if w in vocab]
13                 clean_lines.append(tokens)
14         return clean_lines
```

Next, we adapt the process_docs() function to load and clean all of the documents in a folder and return a list of all document lines.

The results from this function will be the training data for the word2vec model.

```
1  # load all docs in a directory
2  def process_docs(directory, vocab, is_trian):
3          lines = list()
4          # walk through all files in the folder
5          for filename in listdir(directory):
6                  # skip any reviews in the test set
7                  if is_trian and filename.startswith('cv9'):
8                          continue
9                  if not is_trian and not filename.startswith('cv9'):
10                         continue
11                 # create the full path of the file to open
12                 path = directory + '/' + filename
13                 # load and clean the doc
14                 doc = load_doc(path)
15                 doc_lines = doc_to_clean_lines(doc, vocab)
16                 # add lines to list
17                 lines += doc_lines
18         return lines
```

We can then load all of the training data and convert it into a long list of 'sentences' (lists of tokens) ready for fitting the word2vec model.

```
1 # load training data
2 positive_lines = process_docs('txt_sentoken/pos', vocab, True)
3 negative_lines = process_docs('txt_sentoken/neg', vocab, True)
4 sentences = negative_docs + positive_docs
5 print('Total training sentences: %d' % len(sentences))
```

We will use the word2vec implementation provided in the Gensim Python library. Specifically the [Word2Vec class](#).

For more on training a standalone word embedding with Gensim, see the post:

- [How to Develop Word Embeddings in Python with Gensim](#)

The model is fit when constructing the class. We pass in the list of clean sentences from the training data, then specify the size of the embedding vector space (we use 100 again), the number of neighboring words to look at when learning how to embed each word in the training sentences (we use 5 neighbors), the number of threads to use when fitting the model (we use 8, but change this if you have more or less CPU cores), and the minimum occurrence count for words to consider in the vocabulary (we set this to 1 as we have already prepared the vocabulary).

After the model is fit, we print the size of the learned vocabulary, which should match the size of our vocabulary in vocab.txt of 25,767 tokens.

```
1 # train word2vec model
2 model = Word2Vec(sentences, size=100, window=5, workers=8, min_count=1)
3 # summarize vocabulary size in model
4 words = list(model.wv.vocab)
5 print('Vocabulary size: %d' % len(words))
```

Finally, we save the learned embedding vectors to file using
the save_word2vec_format() on the model's '*wv*' (word vector) attribute. The embedding is
saved in ASCII format with one word and vector per line.
The complete example is listed below.

```
1  from string import punctuation
2  from os import listdir
3  from gensim.models import Word2Vec
4
5  # load doc into memory
6  def load_doc(filename):
7          # open the file as read only
8          file = open(filename, 'r')
9          # read all text
10         text = file.read()
11         # close the file
12         file.close()
13         return text
14
15 # turn a doc into clean tokens
16 def doc_to_clean_lines(doc, vocab):
17         clean_lines = list()
18         lines = doc.splitlines()
19         for line in lines:
20                 # split into tokens by white space
21                 tokens = line.split()
22                 # remove punctuation from each token
23                 table = str.maketrans('', '', punctuation)
24                 tokens = [w.translate(table) for w in tokens]
25                 # filter out tokens not in vocab
26                 tokens = [w for w in tokens if w in vocab]
27                 clean_lines.append(tokens)
28         return clean_lines
29
30 # load all docs in a directory
31 def process_docs(directory, vocab, is_trian):
32         lines = list()
33         # walk through all files in the folder
34         for filename in listdir(directory):
35                 # skip any reviews in the test set
36                 if is_trian and filename.startswith('cv9'):
37                         continue
38                 if not is_trian and not filename.startswith('cv9'):
39                         continue
40                 # create the full path of the file to open
41                 path = directory + '/' + filename
42                 # load and clean the doc
43                 doc = load_doc(path)
44                 doc_lines = doc_to_clean_lines(doc, vocab)
45                 # add lines to list
46                 lines += doc_lines
47         return lines
48
49 # load the vocabulary
50 vocab_filename = 'vocab.txt'
51 vocab = load_doc(vocab_filename)
52 vocab = vocab.split()
53 vocab = set(vocab)
```

```
54
55  # load training data
56  positive_lines = process_docs('txt_sentoken/pos', vocab, True)
57  negative_lines = process_docs('txt_sentoken/neg', vocab, True)
58  sentences = negative_docs + positive_docs
59  print('Total training sentences: %d' % len(sentences))
60
61  # train word2vec model
62  model = Word2Vec(sentences, size=100, window=5, workers=8, min_count=1)
63  # summarize vocabulary size in model
64  words = list(model.wv.vocab)
65  print('Vocabulary size: %d' % len(words))
66
67  # save model in ASCII (word2vec) format
68  filename = 'embedding_word2vec.txt'
69  model.wv.save_word2vec_format(filename, binary=False)
```

Running the example loads 58,109 sentences from the training data and creates an embedding for a vocabulary of 25,767 words.

You should now have a file 'embedding_word2vec.txt' with the learned vectors in your current working directory.

```
1  Total training sentences: 58109
2  Vocabulary size: 25767
```

Next, let's look at using these learned vectors in our model.

# 5. Use Pre-trained Embedding

In this section, we will use a pre-trained word embedding prepared on a very large text corpus.

We can use the pre-trained word embedding developed in the previous section and the CNN model developed in the section before that.

The first step is to load the word embedding as a directory of words to vectors. The word embedding was saved in so-called '*word2vec*' format that contains a header line. We will skip this header line when loading the embedding.
The function below named *load_embedding()* loads the embedding and returns a directory of words mapped to the vectors in NumPy format.

```
1   # load embedding as a dict
2   def load_embedding(filename):
3           # load embedding into memory, skip first line
4           file = open(filename,'r')
5           lines = file.readlines()[1:]
6           file.close()
7           # create a map of words to vectors
8           embedding = dict()
9           for line in lines:
10                  parts = line.split()
```

```
11                    # key is string word, value is numpy array for vector
12                    embedding[parts[0]] = asarray(parts[1:], dtype='float32')
13          return embedding
```

Now that we have all of the vectors in memory, we can order them in such a way as to match the integer encoding prepared by the Keras Tokenizer.

Recall that we integer encode the review documents prior to passing them to the Embedding layer. The integer maps to the index of a specific vector in the embedding layer. Therefore, it is important that we lay the vectors out in the Embedding layer such that the encoded words map to the correct vector.

Below defines a function *get_weight_matrix()* that takes the loaded embedding and the tokenizer.word_index vocabulary as arguments and returns a matrix with the word vectors in the correct locations.

```
 1  # create a weight matrix for the Embedding layer from a loaded embedding
 2  def get_weight_matrix(embedding, vocab):
 3          # total vocabulary size plus 0 for unknown words
 4          vocab_size = len(vocab) + 1
 5          # define weight matrix dimensions with all 0
 6          weight_matrix = zeros((vocab_size, 100))
 7          # step vocab, store vectors using the Tokenizer's integer mapping
 8          for word, i in vocab.items():
 9                  weight_matrix[i] = embedding.get(word)
10          return weight_matrix
```

Now we can use these functions to create our new Embedding layer for our model.

```
 1  ...
 2  # load embedding from file
 3  raw_embedding = load_embedding('embedding_word2vec.txt')
 4  # get vectors in the right order
 5  embedding_vectors = get_weight_matrix(raw_embedding, tokenizer.word_index)
 6  # create the embedding layer
 7  embedding_layer = Embedding(vocab_size, 100, weights=[embedding_vectors], input_length=max_length,
    trainable=False)
```

Note that the prepared weight matrix *embedding_vectors* is passed to the new Embedding layer as an argument and that we set the '*trainable*' argument to '*False*' to ensure that the network does not try to adapt the pre-learned vectors as part of training the network.
We can now add this layer to our model. We also have a slightly different model configuration with a lot more filters (128) in the CNN model and a kernel that matches the 5 words used as neighbors when developing the word2vec embedding. Finally, the back-end of the model was simplified.

```
1 # define model
2 model = Sequential()
3 model.add(embedding_layer)
4 model.add(Conv1D(filters=128, kernel_size=5, activation='relu'))
5 model.add(MaxPooling1D(pool_size=2))
6 model.add(Flatten())
7 model.add(Dense(1, activation='sigmoid'))
8 print(model.summary())
```

These changes were found with a little trial and error.

The complete code listing is provided below.

```
 1   from string import punctuation
 2   from os import listdir
 3   from numpy import array
 4   from numpy import asarray
 5   from numpy import zeros
 6   from keras.preprocessing.text import Tokenizer
 7   from keras.preprocessing.sequence import pad_sequences
 8   from keras.models import Sequential
 9   from keras.layers import Dense
10   from keras.layers import Flatten
11   from keras.layers import Embedding
12   from keras.layers.convolutional import Conv1D
13   from keras.layers.convolutional import MaxPooling1D
14
15   # load doc into memory
16   def load_doc(filename):
17           # open the file as read only
18           file = open(filename, 'r')
19           # read all text
20           text = file.read()
21           # close the file
22           file.close()
23           return text
24
25   # turn a doc into clean tokens
26   def clean_doc(doc, vocab):
27           # split into tokens by white space
28           tokens = doc.split()
29           # remove punctuation from each token
30           table = str.maketrans('', '', punctuation)
31           tokens = [w.translate(table) for w in tokens]
32           # filter out tokens not in vocab
33           tokens = [w for w in tokens if w in vocab]
34           tokens = ' '.join(tokens)
35           return tokens
36
37   # load all docs in a directory
38   def process_docs(directory, vocab, is_trian):
39           documents = list()
40           # walk through all files in the folder
41           for filename in listdir(directory):
42                   # skip any reviews in the test set
43                   if is_trian and filename.startswith('cv9'):
44                           continue
45                   if not is_trian and not filename.startswith('cv9'):
46                           continue
47                   # create the full path of the file to open
48                   path = directory + '/' + filename
49                   # load the doc
50                   doc = load_doc(path)
51                   # clean doc
52                   tokens = clean_doc(doc, vocab)
53                   # add to list
54                   documents.append(tokens)
55           return documents
```

```python
56
57  # load embedding as a dict
58  def load_embedding(filename):
59          # load embedding into memory, skip first line
60          file = open(filename,'r')
61          lines = file.readlines()[1:]
62          file.close()
63          # create a map of words to vectors
64          embedding = dict()
65          for line in lines:
66                  parts = line.split()
67                  # key is string word, value is numpy array for vector
68                  embedding[parts[0]] = asarray(parts[1:], dtype='float32')
69          return embedding
70
71  # create a weight matrix for the Embedding layer from a loaded embedding
72  def get_weight_matrix(embedding, vocab):
73          # total vocabulary size plus 0 for unknown words
74          vocab_size = len(vocab) + 1
75          # define weight matrix dimensions with all 0
76          weight_matrix = zeros((vocab_size, 100))
77          # step vocab, store vectors using the Tokenizer's integer mapping
78          for word, i in vocab.items():
79                  weight_matrix[i] = embedding.get(word)
80          return weight_matrix
81
82  # load the vocabulary
83  vocab_filename = 'vocab.txt'
84  vocab = load_doc(vocab_filename)
85  vocab = vocab.split()
86  vocab = set(vocab)
87
88  # load all training reviews
89  positive_docs = process_docs('txt_sentoken/pos', vocab, True)
90  negative_docs = process_docs('txt_sentoken/neg', vocab, True)
91  train_docs = negative_docs + positive_docs
92
93  # create the tokenizer
94  tokenizer = Tokenizer()
95  # fit the tokenizer on the documents
96  tokenizer.fit_on_texts(train_docs)
97
98  # sequence encode
99  encoded_docs = tokenizer.texts_to_sequences(train_docs)
100 # pad sequences
101 max_length = max([len(s.split()) for s in train_docs])
102 Xtrain = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
103 # define training labels
104 ytrain = array([0 for _ in range(900)] + [1 for _ in range(900)])
105
106 # load all test reviews
107 positive_docs = process_docs('txt_sentoken/pos', vocab, False)
108 negative_docs = process_docs('txt_sentoken/neg', vocab, False)
109 test_docs = negative_docs + positive_docs
110 # sequence encode
111 encoded_docs = tokenizer.texts_to_sequences(test_docs)
112 # pad sequences
113 Xtest = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
114 # define test labels
115 ytest = array([0 for _ in range(100)] + [1 for _ in range(100)])
```

```
116
117 # define vocabulary size (largest integer value)
118 vocab_size = len(tokenizer.word_index) + 1
119
120 # load embedding from file
121 raw_embedding = load_embedding('embedding_word2vec.txt')
122 # get vectors in the right order
123 embedding_vectors = get_weight_matrix(raw_embedding, tokenizer.word_index)
124 # create the embedding layer
125 embedding_layer = Embedding(vocab_size, 100, weights=[embedding_vectors], input_length=max_length,
126 trainable=False)
127
128 # define model
129 model = Sequential()
130 model.add(embedding_layer)
131 model.add(Conv1D(filters=128, kernel_size=5, activation='relu'))
132 model.add(MaxPooling1D(pool_size=2))
133 model.add(Flatten())
134 model.add(Dense(1, activation='sigmoid'))
135 print(model.summary())
136 # compile network
137 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
138 # fit network
139 model.fit(Xtrain, ytrain, epochs=10, verbose=2)
140 # evaluate
141 loss, acc = model.evaluate(Xtest, ytest, verbose=0)
    print('Test Accuracy: %f' % (acc*100))
```

Running the example shows that performance was not improved.

In fact, performance was a lot worse. The results show that the training dataset was learned successfully, but evaluation on the test dataset was very poor, at just above 50% accuracy.

The cause of the poor test performance may be because of the chosen word2vec configuration or the chosen neural network configuration.

```
 1  …
 2  Epoch 6/10
 3  2s - loss: 0.3306 - acc: 0.8778
 4  Epoch 7/10
 5  2s - loss: 0.2888 - acc: 0.8917
 6  Epoch 8/10
 7  2s - loss: 0.1878 - acc: 0.9439
 8  Epoch 9/10
 9  2s - loss: 0.1255 - acc: 0.9750
10  Epoch 10/10
11  2s - loss: 0.0812 - acc: 0.9928
12  Test Accuracy: 53.000000
```

The weights in the embedding layer can be used as a starting point for the network, and adapted during the training of the network. We can do this by setting '*trainable=True*' (the default) in the creation of the embedding layer.

Repeating the experiment with this change shows slightly better results, but still poor.

I would encourage you to explore alternate configurations of the embedding and network to see if you can do better. Let me know how you do.

```
 1  …
 2  Epoch 6/10
 3  4s - loss: 0.0950 - acc: 0.9917
 4  Epoch 7/10
 5  4s - loss: 0.0355 - acc: 0.9983
 6  Epoch 8/10
 7  4s - loss: 0.0158 - acc: 1.0000
 8  Epoch 9/10
 9  4s - loss: 0.0080 - acc: 1.0000
10  Epoch 10/10
11  4s - loss: 0.0050 - acc: 1.0000
12  Test Accuracy: 57.500000
```

It is possible to use pre-trained word vectors prepared on very large corpora of text data.

For example, both Google and Stanford provide pre-trained word vectors that you can download, trained with the efficient word2vec and GloVe methods respectively.

Let's try to use pre-trained vectors in our model.

You can download pre-trained GloVe vectors from the Stanford webpage. Specifically, vectors trained on Wikipedia data:

- glove.6B.zip (822 Megabyte download)

Unzipping the file, you will find pre-trained embeddings for various different dimensions. We will load the 100 dimension version in the file '*glove.6B.100d.txt*'
The Glove file does not contain a header file, so we do not need to skip the first line when loading the embedding into memory. The updated *load_embedding()* function is listed below.

```
 1  # load embedding as a dict
 2  def load_embedding(filename):
 3          # load embedding into memory, skip first line
 4          file = open(filename,'r')
 5          lines = file.readlines()
 6          file.close()
 7          # create a map of words to vectors
 8          embedding = dict()
 9          for line in lines:
10                  parts = line.split()
11                  # key is string word, value is numpy array for vector
12                  embedding[parts[0]] = asarray(parts[1:], dtype='float32')
13          return embedding
```

It is possible that the loaded embedding does not contain all of the words in our chosen vocabulary. As such, when creating the Embedding weight matrix, we need to skip words that do not have a corresponding vector in the loaded GloVe data. Below is the updated, more defensive version of the *get_weight_matrix()* function.

```
 1  # create a weight matrix for the Embedding layer from a loaded embedding
 2  def get_weight_matrix(embedding, vocab):
```

```
3          # total vocabulary size plus 0 for unknown words
4          vocab_size = len(vocab) + 1
5          # define weight matrix dimensions with all 0
6          weight_matrix = zeros((vocab_size, 100))
7          # step vocab, store vectors using the Tokenizer's integer mapping
8          for word, i in vocab.items():
9                  vector = embedding.get(word)
10                 if vector is not None:
11                         weight_matrix[i] = vector
12         return weight_matrix
```

We can now load the GloVe embedding and create the Embedding layer as before.

```
1  # load embedding from file
2  raw_embedding = load_embedding('glove.6B.100d.txt')
3  # get vectors in the right order
4  embedding_vectors = get_weight_matrix(raw_embedding, tokenizer.word_index)
5  # create the embedding layer
6  embedding_layer = Embedding(vocab_size, 100, weights=[embedding_vectors], input_length=max_length,
   trainable=False)
```

We will use the same model as before.

The complete example is listed below.

```
1   from string import punctuation
2   from os import listdir
3   from numpy import array
4   from numpy import asarray
5   from numpy import zeros
6   from keras.preprocessing.text import Tokenizer
7   from keras.preprocessing.sequence import pad_sequences
8   from keras.models import Sequential
9   from keras.layers import Dense
10  from keras.layers import Flatten
11  from keras.layers import Embedding
12  from keras.layers.convolutional import Conv1D
13  from keras.layers.convolutional import MaxPooling1D
14
15  # load doc into memory
16  def load_doc(filename):
17          # open the file as read only
18          file = open(filename, 'r')
19          # read all text
20          text = file.read()
21          # close the file
22          file.close()
23          return text
24
25  # turn a doc into clean tokens
26  def clean_doc(doc, vocab):
27          # split into tokens by white space
28          tokens = doc.split()
29          # remove punctuation from each token
30          table = str.maketrans('', '', punctuation)
31          tokens = [w.translate(table) for w in tokens]
32          # filter out tokens not in vocab
33          tokens = [w for w in tokens if w in vocab]
34          tokens = ' '.join(tokens)
35          return tokens
36
```

```python
37  # load all docs in a directory
38  def process_docs(directory, vocab, is_trian):
39          documents = list()
40          # walk through all files in the folder
41          for filename in listdir(directory):
42                  # skip any reviews in the test set
43                  if is_trian and filename.startswith('cv9'):
44                          continue
45                  if not is_trian and not filename.startswith('cv9'):
46                          continue
47                  # create the full path of the file to open
48                  path = directory + '/' + filename
49                  # load the doc
50                  doc = load_doc(path)
51                  # clean doc
52                  tokens = clean_doc(doc, vocab)
53                  # add to list
54                  documents.append(tokens)
55          return documents
56
57  # load embedding as a dict
58  def load_embedding(filename):
59          # load embedding into memory, skip first line
60          file = open(filename,'r')
61          lines = file.readlines()
62          file.close()
63          # create a map of words to vectors
64          embedding = dict()
65          for line in lines:
66                  parts = line.split()
67                  # key is string word, value is numpy array for vector
68                  embedding[parts[0]] = asarray(parts[1:], dtype='float32')
69          return embedding
70
71  # create a weight matrix for the Embedding layer from a loaded embedding
72  def get_weight_matrix(embedding, vocab):
73          # total vocabulary size plus 0 for unknown words
74          vocab_size = len(vocab) + 1
75          # define weight matrix dimensions with all 0
76          weight_matrix = zeros((vocab_size, 100))
77          # step vocab, store vectors using the Tokenizer's integer mapping
78          for word, i in vocab.items():
79                  vector = embedding.get(word)
80                  if vector is not None:
81                          weight_matrix[i] = vector
82          return weight_matrix
83
84  # load the vocabulary
85  vocab_filename = 'vocab.txt'
86  vocab = load_doc(vocab_filename)
87  vocab = vocab.split()
88  vocab = set(vocab)
89
90  # load all training reviews
91  positive_docs = process_docs('txt_sentoken/pos', vocab, True)
92  negative_docs = process_docs('txt_sentoken/neg', vocab, True)
93  train_docs = negative_docs + positive_docs
94
95  # create the tokenizer
96  tokenizer = Tokenizer()
```

```
 97  # fit the tokenizer on the documents
 98  tokenizer.fit_on_texts(train_docs)
 99
100  # sequence encode
101  encoded_docs = tokenizer.texts_to_sequences(train_docs)
102  # pad sequences
103  max_length = max([len(s.split()) for s in train_docs])
104  Xtrain = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
105  # define training labels
106  ytrain = array([0 for _ in range(900)] + [1 for _ in range(900)])
107
108  # load all test reviews
109  positive_docs = process_docs('txt_sentoken/pos', vocab, False)
110  negative_docs = process_docs('txt_sentoken/neg', vocab, False)
111  test_docs = negative_docs + positive_docs
112  # sequence encode
113  encoded_docs = tokenizer.texts_to_sequences(test_docs)
114  # pad sequences
115  Xtest = pad_sequences(encoded_docs, maxlen=max_length, padding='post')
116  # define test labels
117  ytest = array([0 for _ in range(100)] + [1 for _ in range(100)])
118
119  # define vocabulary size (largest integer value)
120  vocab_size = len(tokenizer.word_index) + 1
121
122  # load embedding from file
123  raw_embedding = load_embedding('glove.6B.100d.txt')
124  # get vectors in the right order
125  embedding_vectors = get_weight_matrix(raw_embedding, tokenizer.word_index)
126  # create the embedding layer
127  embedding_layer = Embedding(vocab_size, 100, weights=[embedding_vectors], input_length=max_length,
128  trainable=False)
129
130  # define model
131  model = Sequential()
132  model.add(embedding_layer)
133  model.add(Conv1D(filters=128, kernel_size=5, activation='relu'))
134  model.add(MaxPooling1D(pool_size=2))
135  model.add(Flatten())
136  model.add(Dense(1, activation='sigmoid'))
137  print(model.summary())
138  # compile network
139  model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
140  # fit network
141  model.fit(Xtrain, ytrain, epochs=10, verbose=2)
142  # evaluate
143  loss, acc = model.evaluate(Xtest, ytest, verbose=0)
     print('Test Accuracy: %f' % (acc*100))
```

Running the example shows better performance.

Again, the training dataset is easily learned and the model achieves 76% accuracy on the test dataset. This is good, but not as good as using a learned Embedding layer.

This may be cause of the higher quality vectors trained on more data and/or using a slightly different training process.

Your specific results may vary given the stochastic nature of neural networks. Try running the example a few times.

```
 1  …
 2  Epoch 6/10
 3  2s - loss: 0.0278 - acc: 1.0000
 4  Epoch 7/10
 5  2s - loss: 0.0174 - acc: 1.0000
 6  Epoch 8/10
 7  2s - loss: 0.0117 - acc: 1.0000
 8  Epoch 9/10
 9  2s - loss: 0.0086 - acc: 1.0000
10  Epoch 10/10
11  2s - loss: 0.0068 - acc: 1.0000
12  Test Accuracy: 76.000000
```

In this case, it seems that learning the embedding as part of the learning task may be a better direction than using a specifically trained embedding or a more general pre-trained embedding.

# Further Reading

This section provides more resources on the topic if you are looking go deeper.

## Dataset

- Movie Review Data
- A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts, 2004.
- Movie Review Polarity Dataset (.tgz)
- Dataset Readme v2.0 and v1.1.

## APIs

- collections API – Container datatypes
- Tokenizer Keras API
- Embedding Keras API
- Gensim Word2Vec API
- Gensim WordVector API

## Embedding Methods

- word2vec on Google Code
- GloVe on Stanford

## Related Posts

- Using pre-trained word embeddings in a Keras model, 2016.
- Implementing a CNN for Text Classification in TensorFlow, 2015.

# Summary

In this tutorial, you discovered how to develop word embeddings for the classification of movie reviews.

Specifically, you learned:

- How to prepare movie review text data for classification with deep learning methods.
- How to learn a word embedding as part of fitting a deep learning model.
- How to learn a standalone word embedding and how to use a pre-trained embedding in a neural network model.

Do you have any questions?

Ask your questions in the comments below and I will do my best to answer.

**Note**: This post is an excerpt chapter from: "[Deep Learning for Natural Language Processing](#)". Take a look, if you want more step-by-step tutorials on getting the most out of deep learning methods when working with text data.