

[Click to Take the FREE NLP Crash-Course](#)

How to Develop a Multichannel CNN Model for Text Classification

by **Jason Brownlee** on January 12, 2018 in **Deep Learning for Natural Language Processing**

[Tweet](#)[Share](#)[Share](#)

Last Updated on December 4, 2019

A standard deep learning model for [text classification](#) and [sentiment analysis](#) uses a word embedding layer and one-dimensional convolutional neural network.

The model can be expanded by using multiple parallel convolutional neural networks that read the source document using different kernel sizes. This, in effect, creates a multichannel convolutional neural network for text that reads text with different n-gram sizes (groups of words).

In this tutorial, you will discover how to develop a multichannel convolutional neural network for sentiment prediction on text movie review data.

After completing this tutorial, you will know:

- How to prepare movie review text data for modeling.
- How to develop a multichannel convolutional neural network for text in Keras.
- How to evaluate a fit model on unseen movie review data.

Discover how to develop deep learning models for text classification, translation, photo captioning and more [in my new book](#), with 30 step-by-step tutorials and full source code.

Let's get started.

- **Update Feb/2018:** Small code change to reflect changes in Keras 2.1.3 API.



How to Develop an N-gram Multichannel Convolutional Neural Network for Sentiment Analysis

Photo by [Ed Dunens](#), some rights reserved.

Tutorial Overview

This tutorial is divided into 4 parts; they are:

1. Movie Review Dataset
2. Data Preparation
3. Develop Multichannel Model
4. Evaluate Model

Python Environment

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

You must have Keras (2.0 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 post:

- [How to Setup a Python Environment for Machine Learning and Deep Learning with Anaconda](#)

Need help with Deep Learning for Text Data?

Take my free 7-day email crash course now (with code).

Click to sign-up and also get a free PDF Ebook version of the course.

[Start Your FREE Crash-Course Now](#)

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 was 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 containing the text “neg” and “pos” for negative and positive reviews, respectively. Each file contains a review with a naming convention cv000 to cv999 for each neg and pos.

Next, let's look at loading and preparing the text data.

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. Prepare all reviews and save to file.

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) is 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.

Without getting bogged down too much by the details, we will prepare the data in 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 loads and cleans one movie review.

The tokens from the clean review are printed for review.

```

1 ...
2 'creepy', 'place', 'even', 'acting', 'hell', 'solid', 'dreamy', 'depp', 'turning', 'typically

```

Clean All Reviews and Save

We can now use the function to clean reviews and apply it to all reviews.

To do this, we will develop a new function named *process_docs()* below that will walk through all reviews in a directory, clean them, and return them as a list.

We will also add an argument to the function to indicate whether the function is processing train or test reviews, that way the filenames can be filtered (as described above) and only those train or test reviews requested will be cleaned and returned.

The full function is listed below.

```

1 # load all docs in a directory
2 def process_docs(directory, 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)
17        # add to list
18        documents.append(tokens)
19    return documents

```

We can call this function with negative training reviews as follows:

```

1 negative_docs = process_docs('txt_sentoken/neg', True)

```

Next, we need labels for the train and test documents. We know that we have 900 training documents and 100 test documents. We can use a Python list comprehension to create the labels for the negative (0) and positive (1) reviews for both train and test sets.

```

1 trainy = [0 for _ in range(900)] + [1 for _ in range(900)]
2 testY = [0 for _ in range(100)] + [1 for _ in range(100)]

```

Finally, we want to save the prepared train and test sets to file so that we can load them later for modeling and model evaluation.

The function below-named `save_dataset()` will save a given prepared dataset (X and y elements) to a file using the pickle API.

```

1 # save a dataset to file
2 def save_dataset(dataset, filename):
3     dump(dataset, open(filename, 'wb'))
4     print('Saved: %s' % filename)

```

Complete Example

We can tie all of these data preparation steps together.

The complete example is listed below.

```

1 from string import punctuation
2 from os import listdir
3 from nltk.corpus import stopwords
4 from pickle import dump
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        tokens = ' '.join(tokens)
31        return tokens
32
33    # load all docs in a directory
34    def process_docs(directory, is_trian):
35        documents = list()
36        # walk through all files in the folder
37        for filename in listdir(directory):
38            # skip any reviews in the test set
39            if is_trian and filename.startswith('cv9'):
40                continue
41            if not is_trian and not filename.startswith('cv9'):
42                continue
43            # create the full path of the file to open
44            path = directory + '/' + filename
45            # load the doc
46            doc = load_doc(path)
47            # clean doc
48            tokens = clean_doc(doc)
49            # add to list
50            documents.append(tokens)
51        return documents
52
53    # save a dataset to file
54    def save_dataset(dataset, filename):
55        dump(dataset, open(filename, 'wb'))
56        print('Saved: %s' % filename)
57
58    # load all training reviews
59    negative_docs = process_docs('txt_sentoken/neg', True)
60    positive_docs = process_docs('txt_sentoken/pos', True)
61    trainX = negative_docs + positive_docs
62    trainy = [0 for _ in range(900)] + [1 for _ in range(900)]
63    save_dataset([trainX, trainy], 'train.pkl')
64
65    # load all test reviews

```



```

66 negative_docs = process_docs('txt_sentoken/neg', False)
67 positive_docs = process_docs('txt_sentoken/pos', False)
68 testX = negative_docs + positive_docs
69 testY = [0 for _ in range(100)] + [1 for _ in range(100)]
70 save_dataset([testX, testY], 'test.pkl')

```

Running the example cleans the text movie review documents, creates labels, and saves the prepared data for both train and test datasets in *train.pkl* and *test.pkl* respectively.

Now we are ready to develop our model.

Develop Multichannel Model

In this section, we will develop a multichannel convolutional neural network for the sentiment analysis prediction problem.

This section is divided into 3 parts:

1. Encode Data
2. Define Model.
3. Complete Example.

Encode Data

The first step is to load the cleaned training dataset.

The function below-named *load_dataset()* can be called to load the pickled training dataset.

```

1 # load a clean dataset
2 def load_dataset(filename):
3     return load(open(filename, 'rb'))
4
5 trainLines, trainLabels = load_dataset('train.pkl')

```

Next, we must fit a Keras Tokenizer on the training dataset. We will use this tokenizer to both define the vocabulary for the *Embedding layer* and encode the review documents as integers.

The function *create_tokenizer()* below will create a Tokenizer given a list of documents.

```

1 # fit a tokenizer
2 def create_tokenizer(lines):
3     tokenizer = Tokenizer()
4     tokenizer.fit_on_texts(lines)
5     return tokenizer

```

We also need to know the maximum length of input sequences as input for the model and to pad all

sequences to the fixed length.

The function `max_length()` below will calculate the maximum length (number of words) for all reviews in the training dataset.

```
1 # calculate the maximum document length
2 def max_length(lines):
3     return max([len(s.split()) for s in lines])
```

We also need to know the size of the vocabulary for the Embedding layer.

This can be calculated from the prepared Tokenizer, as follows:

```
1 # calculate vocabulary size
2 vocab_size = len(tokenizer.word_index) + 1
```

Finally, we can integer encode and pad the clean movie review text.

The function below named `encode_text()` will both encode and pad text data to the maximum review length.

```
1 # encode a list of lines
2 def encode_text(tokenizer, lines, length):
3     # integer encode
4     encoded = tokenizer.texts_to_sequences(lines)
5     # pad encoded sequences
6     padded = pad_sequences(encoded, maxlen=length, padding='post')
7     return padded
```

Define Model

A standard model for document classification is to use an Embedding layer as input, followed by a one-dimensional convolutional neural network, pooling layer, and then a prediction output layer.

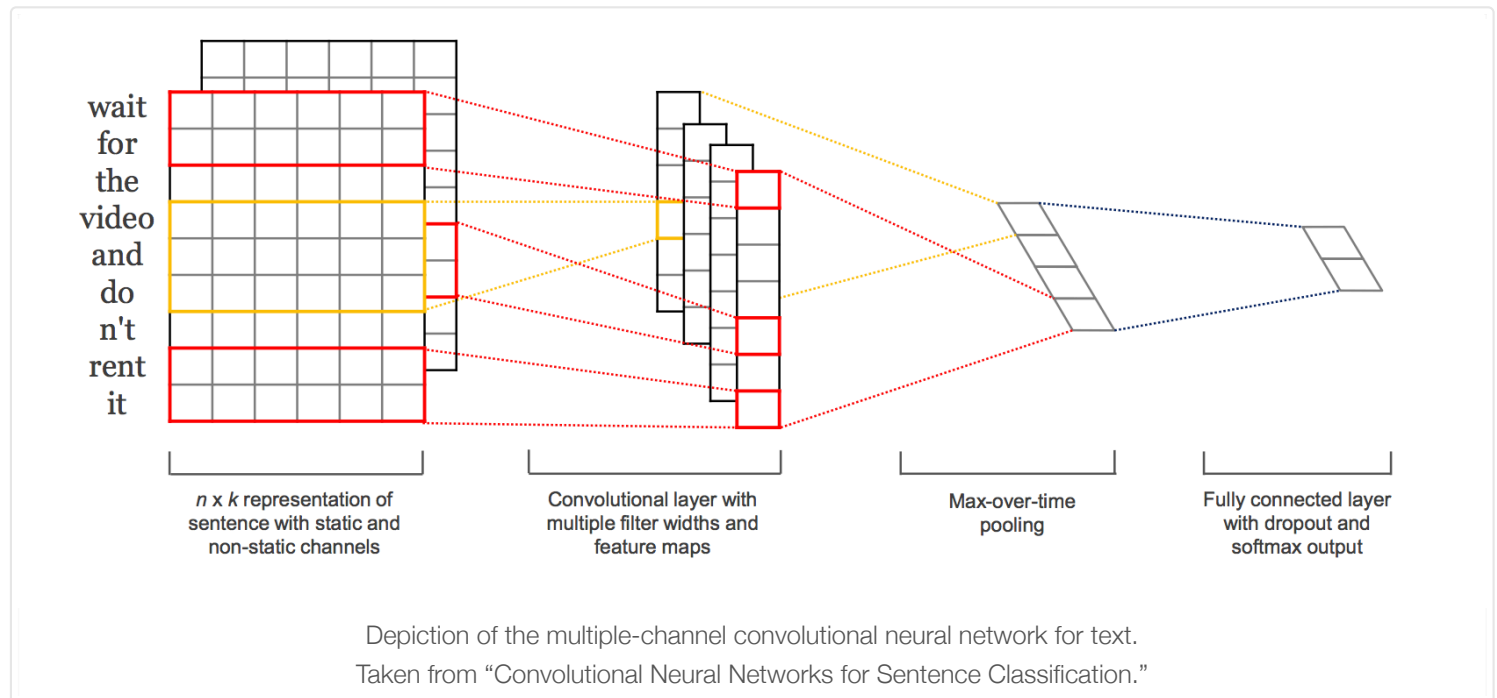
The kernel size in the [convolutional layer](#) defines the number of words to consider as the convolution is passed across the input text document, providing a grouping parameter.

A multi-channel convolutional neural network for document classification involves using multiple versions of the standard model with different sized kernels. This allows the document to be processed at different resolutions or different n-grams (groups of words) at a time, whilst the model learns how to best integrate these interpretations.

This approach was first described by Yoon Kim in his 2014 paper titled “[Convolutional Neural Networks for Sentence Classification](#).”

In the paper, Kim experimented with static and dynamic (updated) embedding layers, we can simplify the approach and instead focus only on the use of different kernel sizes.

This approach is best understood with a diagram taken from Kim's paper:



In Keras, a multiple-input model can be defined using the [functional API](#).

We will define a model with three input channels for processing 4-grams, 6-grams, and 8-grams of movie review text.

Each channel is comprised of the following elements:

- Input layer that defines the length of input sequences.
- Embedding layer set to the size of the vocabulary and 100-dimensional real-valued representations.
- One-dimensional convolutional layer with 32 filters and a kernel size set to the number of words to read at once.
- Max Pooling layer to consolidate the output from the convolutional layer.
- Flatten layer to reduce the three-dimensional output to two dimensional for concatenation.

The output from the three channels are concatenated into a single vector and process by a Dense layer and an output layer.

The function below defines and returns the model. As part of defining the model, a summary of the

defined model is printed and a plot of the model graph is created and saved to file.

```

1  # define the model
2  def define_model(length, vocab_size):
3      # channel 1
4      inputs1 = Input(shape=(length,))
5      embedding1 = Embedding(vocab_size, 100)(inputs1)
6      conv1 = Conv1D(filters=32, kernel_size=4, activation='relu')(embedding1)
7      drop1 = Dropout(0.5)(conv1)
8      pool1 = MaxPooling1D(pool_size=2)(drop1)
9      flat1 = Flatten()(pool1)
10     # channel 2
11     inputs2 = Input(shape=(length,))
12     embedding2 = Embedding(vocab_size, 100)(inputs2)
13     conv2 = Conv1D(filters=32, kernel_size=6, activation='relu')(embedding2)
14     drop2 = Dropout(0.5)(conv2)
15     pool2 = MaxPooling1D(pool_size=2)(drop2)
16     flat2 = Flatten()(pool2)
17     # channel 3
18     inputs3 = Input(shape=(length,))
19     embedding3 = Embedding(vocab_size, 100)(inputs3)
20     conv3 = Conv1D(filters=32, kernel_size=8, activation='relu')(embedding3)
21     drop3 = Dropout(0.5)(conv3)
22     pool3 = MaxPooling1D(pool_size=2)(drop3)
23     flat3 = Flatten()(pool3)
24     # merge
25     merged = concatenate([flat1, flat2, flat3])
26     # interpretation
27     dense1 = Dense(10, activation='relu')(merged)
28     outputs = Dense(1, activation='sigmoid')(dense1)
29     model = Model(inputs=[inputs1, inputs2, inputs3], outputs=outputs)
30     # compile
31     model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
32     # summarize
33     print(model.summary())
34     plot_model(model, show_shapes=True, to_file='multichannel.png')
35     return model

```

Complete Example

Pulling all of this together, the complete example is listed below.

```

1  from pickle import load
2  from numpy import array
3  from keras.preprocessing.text import Tokenizer
4  from keras.preprocessing.sequence import pad_sequences
5  from keras.utils.vis_utils import plot_model
6  from keras.models import Model
7  from keras.layers import Input
8  from keras.layers import Dense
9  from keras.layers import Flatten
10 from keras.layers import Dropout
11 from keras.layers import Embedding
12 from keras.layers.convolutional import Conv1D

```

```

13 from keras.layers.convolutional import MaxPooling1D
14 from keras.layers.merge import concatenate
15
16 # load a clean dataset
17 def load_dataset(filename):
18     return load(open(filename, 'rb'))
19
20 # fit a tokenizer
21 def create_tokenizer(lines):
22     tokenizer = Tokenizer()
23     tokenizer.fit_on_texts(lines)
24     return tokenizer
25
26 # calculate the maximum document length
27 def max_length(lines):
28     return max([len(s.split()) for s in lines])
29
30 # encode a list of lines
31 def encode_text(tokenizer, lines, length):
32     # integer encode
33     encoded = tokenizer.texts_to_sequences(lines)
34     # pad encoded sequences
35     padded = pad_sequences(encoded, maxlen=length, padding='post')
36     return padded
37
38 # define the model
39 def define_model(length, vocab_size):
40     # channel 1
41     inputs1 = Input(shape=(length,))
42     embedding1 = Embedding(vocab_size, 100)(inputs1)
43     conv1 = Conv1D(filters=32, kernel_size=4, activation='relu')(embedding1)
44     drop1 = Dropout(0.5)(conv1)
45     pool1 = MaxPooling1D(pool_size=2)(drop1)
46     flat1 = Flatten()(pool1)
47     # channel 2
48     inputs2 = Input(shape=(length,))
49     embedding2 = Embedding(vocab_size, 100)(inputs2)
50     conv2 = Conv1D(filters=32, kernel_size=6, activation='relu')(embedding2)
51     drop2 = Dropout(0.5)(conv2)
52     pool2 = MaxPooling1D(pool_size=2)(drop2)
53     flat2 = Flatten()(pool2)
54     # channel 3
55     inputs3 = Input(shape=(length,))
56     embedding3 = Embedding(vocab_size, 100)(inputs3)
57     conv3 = Conv1D(filters=32, kernel_size=8, activation='relu')(embedding3)
58     drop3 = Dropout(0.5)(conv3)
59     pool3 = MaxPooling1D(pool_size=2)(drop3)
60     flat3 = Flatten()(pool3)
61     # merge
62     merged = concatenate([flat1, flat2, flat3])
63     # interpretation
64     dense1 = Dense(10, activation='relu')(merged)
65     outputs = Dense(1, activation='sigmoid')(dense1)
66     model = Model(inputs=[inputs1, inputs2, inputs3], outputs=outputs)
67     # compile
68     model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
69     # summarize

```

```

70     print(model.summary())
71     plot_model(model, show_shapes=True, to_file='multichannel.png')
72     return model
73
74 # load training dataset
75 trainLines, trainLabels = load_dataset('train.pkl')
76 # create tokenizer
77 tokenizer = create_tokenizer(trainLines)
78 # calculate max document length
79 length = max_length(trainLines)
80 # calculate vocabulary size
81 vocab_size = len(tokenizer.word_index) + 1
82 print('Max document length: %d' % length)
83 print('Vocabulary size: %d' % vocab_size)
84 # encode data
85 trainX = encode_text(tokenizer, trainLines, length)
86 print(trainX.shape)
87
88 # define model
89 model = define_model(length, vocab_size)
90 # fit model
91 model.fit([trainX, trainX, trainX], array(trainLabels), epochs=10, batch_size=16)
92 # save the model
93 model.save('model.h5')

```

Running the example first prints a summary of the prepared training dataset.

```

1 Max document length: 1380
2 Vocabulary size: 44277
3 (1800, 1380)

```

Next, a summary of the defined model is printed.

| Layer (type) | Output Shape | Param # | Connected to |
|-------------------------|-------------------|---------|-------------------|
| input_1 (InputLayer) | (None, 1380) | 0 | |
| input_2 (InputLayer) | (None, 1380) | 0 | |
| input_3 (InputLayer) | (None, 1380) | 0 | |
| embedding_1 (Embedding) | (None, 1380, 100) | 4427700 | input_1[0][0] |
| embedding_2 (Embedding) | (None, 1380, 100) | 4427700 | input_2[0][0] |
| embedding_3 (Embedding) | (None, 1380, 100) | 4427700 | input_3[0][0] |
| conv1d_1 (Conv1D) | (None, 1377, 32) | 12832 | embedding_1[0][0] |
| conv1d_2 (Conv1D) | (None, 1375, 32) | 19232 | embedding_2[0][0] |
| conv1d_3 (Conv1D) | (None, 1373, 32) | 25632 | embedding_3[0][0] |
| dropout_1 (Dropout) | (None, 1377, 32) | 0 | conv1d_1[0][0] |

| | | | | |
|----|--------------------------------|------------------|--------|-----------------------|
| 23 | | | | |
| 24 | dropout_2 (Dropout) | (None, 1375, 32) | 0 | conv1d_2[0][0] |
| 25 | | | | |
| 26 | dropout_3 (Dropout) | (None, 1373, 32) | 0 | conv1d_3[0][0] |
| 27 | | | | |
| 28 | max_pooling1d_1 (MaxPooling1D) | (None, 688, 32) | 0 | dropout_1[0][0] |
| 29 | | | | |
| 30 | max_pooling1d_2 (MaxPooling1D) | (None, 687, 32) | 0 | dropout_2[0][0] |
| 31 | | | | |
| 32 | max_pooling1d_3 (MaxPooling1D) | (None, 686, 32) | 0 | dropout_3[0][0] |
| 33 | | | | |
| 34 | flatten_1 (Flatten) | (None, 22016) | 0 | max_pooling1d_1[0][0] |
| 35 | | | | |
| 36 | flatten_2 (Flatten) | (None, 21984) | 0 | max_pooling1d_2[0][0] |
| 37 | | | | |
| 38 | flatten_3 (Flatten) | (None, 21952) | 0 | max_pooling1d_3[0][0] |
| 39 | | | | |
| 40 | concatenate_1 (Concatenate) | (None, 65952) | 0 | flatten_1[0][0] |
| 41 | | | | flatten_2[0][0] |
| 42 | | | | flatten_3[0][0] |
| 43 | | | | |
| 44 | dense_1 (Dense) | (None, 10) | 659530 | concatenate_1[0][0] |
| 45 | | | | |
| 46 | dense_2 (Dense) | (None, 1) | 11 | dense_1[0][0] |
| 47 | ===== | | | |
| 48 | Total params: 14,000,337 | | | |
| 49 | Trainable params: 14,000,337 | | | |
| 50 | Non-trainable params: 0 | | | |
| 51 | ===== | | | |

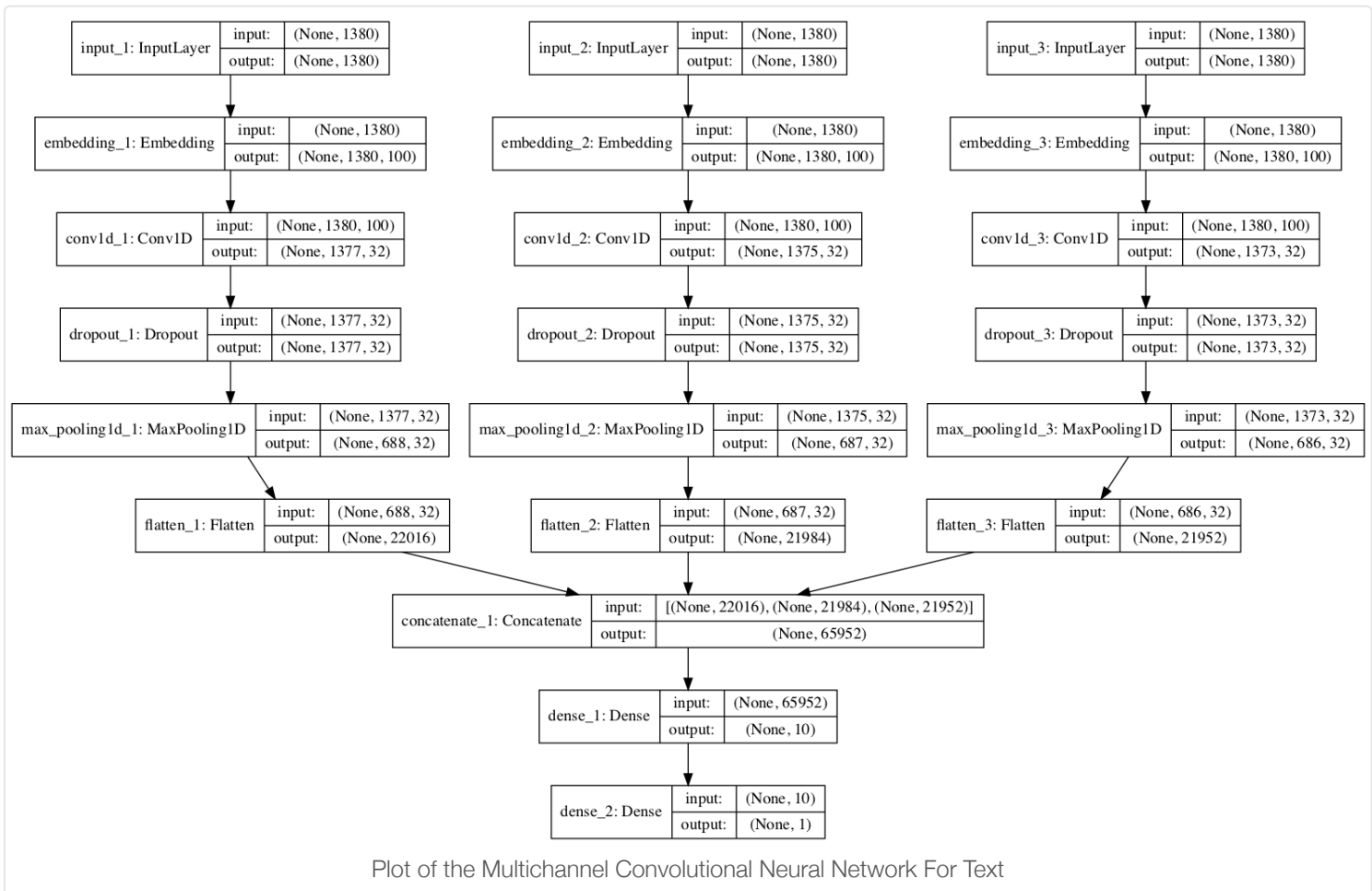
The model is fit relatively quickly and appears to show good skill on the training dataset.

```

1 ...
2 Epoch 6/10
3 1800/1800 [=====] - 30s - loss: 9.9093e-04 - acc: 1.0000
4 Epoch 7/10
5 1800/1800 [=====] - 29s - loss: 5.1899e-04 - acc: 1.0000
6 Epoch 8/10
7 1800/1800 [=====] - 28s - loss: 3.7958e-04 - acc: 1.0000
8 Epoch 9/10
9 1800/1800 [=====] - 29s - loss: 3.0534e-04 - acc: 1.0000
10 Epoch 10/10
11 1800/1800 [=====] - 29s - loss: 2.6234e-04 - acc: 1.0000

```

A plot of the defined model is saved to file, clearly showing the three input channels for the model.



The model is fit for a number of epochs and saved to the file *model.h5* for later evaluation.

Evaluate Model

In this section, we can evaluate the fit model by predicting the sentiment on all reviews in the unseen test dataset.

Using the data loading functions developed in the previous section, we can load and encode both the training and test datasets.

```

1 # load datasets
2 trainLines, trainLabels = load_dataset('train.pkl')
3 testLines, testLabels = load_dataset('test.pkl')
4
5 # create tokenizer
6 tokenizer = create_tokenizer(trainLines)
7 # calculate max document length
8 length = max_length(trainLines)
9 # calculate vocabulary size
10 vocab_size = len(tokenizer.word_index) + 1
11 print('Max document length: %d' % length)

```

```

12 print('Vocabulary size: %d' % vocab_size)
13 # encode data
14 trainX = encode_text(tokenizer, trainLines, length)
15 testX = encode_text(tokenizer, testLines, length)
16 print(trainX.shape, testX.shape)

```

We can load the saved model and evaluate it on both the training and test datasets.

The complete example is listed below.

```

1 from pickle import load
2 from numpy import array
3 from keras.preprocessing.text import Tokenizer
4 from keras.preprocessing.sequence import pad_sequences
5 from keras.models import load_model
6
7 # load a clean dataset
8 def load_dataset(filename):
9     return load(open(filename, 'rb'))
10
11 # fit a tokenizer
12 def create_tokenizer(lines):
13     tokenizer = Tokenizer()
14     tokenizer.fit_on_texts(lines)
15     return tokenizer
16
17 # calculate the maximum document length
18 def max_length(lines):
19     return max([len(s.split()) for s in lines])
20
21 # encode a list of lines
22 def encode_text(tokenizer, lines, length):
23     # integer encode
24     encoded = tokenizer.texts_to_sequences(lines)
25     # pad encoded sequences
26     padded = pad_sequences(encoded, maxlen=length, padding='post')
27     return padded
28
29 # load datasets
30 trainLines, trainLabels = load_dataset('train.pkl')
31 testLines, testLabels = load_dataset('test.pkl')
32
33 # create tokenizer
34 tokenizer = create_tokenizer(trainLines)
35 # calculate max document length
36 length = max_length(trainLines)
37 # calculate vocabulary size
38 vocab_size = len(tokenizer.word_index) + 1
39 print('Max document length: %d' % length)
40 print('Vocabulary size: %d' % vocab_size)
41 # encode data
42 trainX = encode_text(tokenizer, trainLines, length)
43 testX = encode_text(tokenizer, testLines, length)
44 print(trainX.shape, testX.shape)
45

```

```
46 # load the model
47 model = load_model('model.h5')
48
49 # evaluate model on training dataset
50 loss, acc = model.evaluate([trainX,trainX,trainX], array(trainLabels), verbose=0)
51 print('Train Accuracy: %f' % (acc*100))
52
53 # evaluate model on test dataset dataset
54 loss, acc = model.evaluate([testX,testX,testX],array(testLabels), verbose=0)
55 print('Test Accuracy: %f' % (acc*100))
```

Running the example prints the skill of the model on both the training and test datasets.

```
1 Max document length: 1380
2 Vocabulary size: 44277
3 (1800, 1380) (200, 1380)
4
5 Train Accuracy: 100.000000
6 Test Accuracy: 87.500000
```

We can see that, as expected, the skill on the training dataset is excellent, here at 100% accuracy.

We can also see that the skill of the model on the unseen test dataset is also very impressive, achieving 87.5%, which is above the skill of the model reported in the 2014 paper (although not a direct apples-to-apples comparison).

Extensions

This section lists some ideas for extending the tutorial that you may wish to explore.

- **Different n-grams.** Explore the model by changing the kernel size (number of n-grams) used by the channels in the model to see how it impacts model skill.
- **More or Fewer Channels.** Explore using more or fewer channels in the model and see how it impacts model skill.
- **Deeper Network.** Convolutional neural networks perform better in computer vision when they are deeper. Explore using deeper models here and see how it impacts model skill.

Further Reading

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

- [Convolutional Neural Networks for Sentence Classification, 2014.](#)
- [Convolutional Neural Networks for Sentence Classification \(code\).](#)
- [Keras Functional API](#)

Summary

In this tutorial, you discovered how to develop a multichannel convolutional neural network for sentiment prediction on text movie review data.

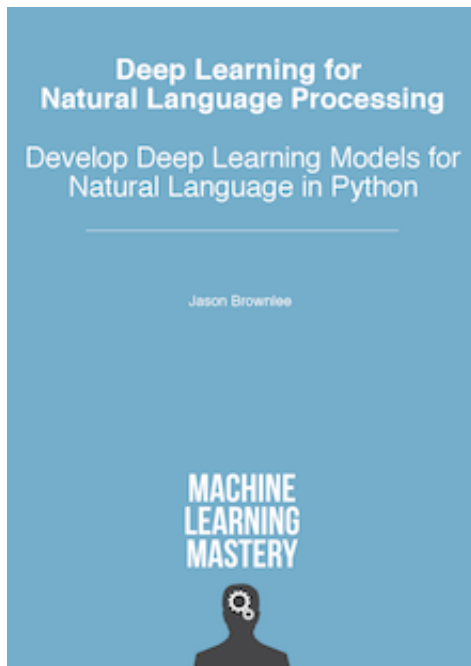
Specifically, you learned:

- How to prepare movie review text data for modeling.
- How to develop a multichannel convolutional neural network for text in Keras.
- How to evaluate a fit model on unseen movie review data.

Do you have any questions?

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

Develop Deep Learning models for Text Data Today!



Develop Your Own Text models in Minutes

...with just a few lines of python code

Discover how in my new Ebook:

[Deep Learning for Natural Language Processing](#)

It provides **self-study tutorials** on topics like:

Bag-of-Words, Word Embedding, Language Models, Caption Generation, Text Translation and much more...

Finally Bring Deep Learning to your Natural Language Processing Projects

Skip the Academics. Just Results.

[SEE WHAT'S INSIDE](#)

Tweet

Share

Share

About Jason Brownlee