# Using Keras and LSTM to Implement Sentence Comparison

B02201041 何庭昀  B02901136 慕家志

## Abstract

In Kaggle data analysis competition, there is a topic we are interested in, Quora Question pairs ( https://www.kaggle.com/c/quora-question-pairs ). The main purpose of this problem is to predict whether two sentences have the same meaning, that is, to implement sentence comparison. We learned a lot about RNN and LSTM in the class, so we want to use this technic to solve our problem. And then, we want to compare effect of different combination of LSTM and DNN.

## Introduction

Quora is a place to gain and share knowledge about anything. It's a platform to ask questions and connect with people who contribute unique insights and quality answers. Over 100 million people visit Quora every month, so it's no surprise that many people ask similarly worded questions. Multiple questions with the same intent can cause seekers to spend more time finding the best answer to their question, and make writers feel they need to answer multiple versions of the same question.

In this problem, we are challenged to tackle this natural language processing problem by applying RNN and LSTM to classify whether question pairs are duplicates or not. Moreover, we want to compare effect of different number of CNN layers and to find out if bidirection LSTM is better than single LSTM layer. Doing so will make it easier to find high quality answers to questions resulting in an improved experience for Quora writers, seekers, and readers.

# Data description

The goal of this problem is to predict which of the provided pairs of questions contain two questions with the same meaning. The training data similarity labels have been supplied by human, which are inherently subjective, as the true meaning of sentences can never be known with certainty. Human labeling is also a 'noisy' process, and reasonable people will disagree. As a result, the similarity labels on this dataset should be taken to be 'informed' but not 100% accurate, and may include incorrect labeling. Data given in this problem is (1.) a training data with 404290 pairs of sentences and (2.) a testing data with 2345796 pairs of sentences. ( https://www. kaggle.com/c/quora-question-pairs/data )

For example, in training data, we have

1. "What is the step by step guide to invest in share market in India?"

   "What is the step by step guide to invest in share market?"

   "0"

2. "Astrology: I am a Capricorn Sun Cap moon and cap rising... what does that say about me?"

   "I'm a triple Capricorn (Sun, Moon and ascendant in Capricorn) What does this say about me?"

   "1"

The labels "0" means previous two sentences are not duplicates, and "1" means previous two sentences are duplicates.

In testing data, we have

1. "How does the Surface Pro himself 4 compare with iPad Pro?"

   "Why did Microsoft choose core m3 and not core i3 home Surface Pro 4?"

2. "Should I have a hair transplant at age 24? How much would it cost?"

   "How much cost does hair transplant require?"

Our final purpose is to predict whether question pairs are duplicates or not.

## Method

We decide to use google library, word2vec, and word2vec model, "GoogleNews-vectors-negative300.bin" (https://github.com/mmihaltz/word2vec-GoogleNews-vectors ), to transform every word in sentences to a distributed represent vector, which means we can evaluate correlation of two words by their vectors value. For example, "telephone" and "phone" will be very close in word2vec form, and "king" + "women" - "men" will be very close to "queen". If a word in sentence cannot be found in google word2vec model, it will be corrected by Peter Norvig's word corrector. And then, we train our model by these sentence. After finishing the training, we predict test data by our model, and evaluate similarity of each sentence pair.

## Code and model structure

Our code is modified by listdo's code on ( https://www.kaggle.com/lystdo/lstm -with-word2vec-embeddings ). In the first part we import some necessary library. Noted that we also import some Keras model such as Dense and LSTM.

```
1.   ########################################
2.   ## import packages
3.   ########################################
4.   import os
5.   import re
```

```
6.   import csv
7.   import codecs
8.   import numpy as np
9.   import pandas as pd
10.
11.  from nltk.corpus import stopwords
12.  from nltk.stem import SnowballStemmer
13.  from string import punctuation
14.  from gensim.models import KeyedVectors, Word2Vec
15.  from keras.preprocessing.text import Tokenizer
16.  from keras.preprocessing.sequence import pad_sequences
17.  from keras.layers import Dense, Input, LSTM, Embedding, Dropout, Activation
18.  from keras.layers.merge import concatenate
19.  from keras.models import Model
20.  from keras.layers.normalization import BatchNormalization
21.  from keras.callbacks import EarlyStopping, ModelCheckpoint
22.  import re, collections
23.  import sys
24.
```

And then we define Peter Norvig's word corrector. The way that the corrector correct the word is to split and combine the word in some way (function edits1) by one or two times (function known_edits2), delete those word which are not exist in dictionary data "big.txt" (function known)( https://github.com/dscape/spell/blob/master /test/resources/big.txt ), and then calculate which word in these candidates has the most frequency in dictionary data "big.txt" (function correct). If you want to learn more, please check the website of Peter Norvig's word corrector ( http://norvig.com/spell-correct.html ).

```
25.  ######################################
26.  ## word corrector define
27.  ######################################
28.  def words(text):
29.      return re.findall('[a-z]+', text.lower())
30.
```

```
31.  def train(features):
32.      model = collections.defaultdict(lambda: 1)
33.      for f in features:
34.          model[f] += 1
35.      return model
36.
37.  NWORDS = train(words(open('big.txt').read()))
38.  alphabet = 'abcdefghijklmnopqrstuvwxyz'
39.
40.  def edits1(word):
41.      s = [(word[:i], word[i:]) for i in range(len(word) + 1)]
42.      deletes    = [a + b[1:] for a, b in s if b]
43.      transposes = [a + b[1] + b[0] + b[2:] for a, b in s if len(b)>1]
44.      replaces   = [a + c + b[1:] for a, b in s for c in alphabet if b]
45.      inserts    = [a + c + b     for a, b in s for c in alphabet]
46.      return set(deletes + transposes + replaces + inserts)
47.
48.  def known_edits2(word):
49.      return set(e2 for e1 in edits1(word) for e2 in edits1(e1) if e2 in NWORDS)
50.
51.  def known(words):
52.      return set(w for w in words if w in NWORDS)
53.
54.  def correct(word):
55.      candidates = known([word]) or known(edits1(word)) or known_edits2(word) or [word]
56.      return max(candidates, key=NWORDS.get)
57.
```

we set directories of google word2vec model "GoogleNews-vectors-negative300.bin", training data, and test data. And then we set some parameters. In this part, we set activation function as "relu".

```
58.  ########################################
59.  ## set directories and parameters
60.  ########################################
61.  BASE_DIR = './'
62.  EMBEDDING_FILE = BASE_DIR + 'GoogleNews-vectors-negative300.bin'
```

```
63.  TRAIN_DATA_FILE = BASE_DIR + 'train.csv'

64.  TEST_DATA_FILE = BASE_DIR + 'test.csv'

65.  MAX_SEQUENCE_LENGTH = 30

66.  MAX_NB_WORDS = 200000

67.  EMBEDDING_DIM = 200

68.  VALIDATION_SPLIT = 0.1

69.

70.  num_lstm = np.random.randint(200, 250) #175~275

71.  num_dense = np.random.randint(100, 150) #100~150

72.  rate_drop_lstm = 0.15 + np.random.rand() * 0.25

73.  rate_drop_dense = 0.15 + np.random.rand() * 0.25

74.

75.  act = 'relu'

76.  re_weight = True # whether to re-weight classes to fit the 17.5% share in test set

77.

78.  STAMP = 'lstm_%d_%d_%.2f_%.2f'%(num_lstm, num_dense, rate_drop_lstm,

79.        rate_drop_dense)

80.
```

In this part, we load word2vec model "GoogleNews-vectors- negative300.bin".

```
81.  ########################################

82.  ## index word vectors

83.  ########################################

84.  print('Indexing word vectors')

85.

86.  word2vec = KeyedVectors.load_word2vec_format(EMBEDDING_FILE, binary=True)

87.  print('Found %s word vectors of word2vec' % len(word2vec.vocab))

88.
```

We use text_to_wordlist function ( https://www.kaggle.com/currie32/quora-question-pairs/the-importance-of-cleaning-text ), to clean the text. That is, delete some unnecessary characters. And then we make the sentences to the form that can be sent to Keras model.

```
89.  ########################################

90.  ## text_to_wordlist
```

```python
91.   #########################################
92.   print('Processing text dataset')
93.
94.   # The function "text_to_wordlist" is from
95.   # https://www.kaggle.com/currie32/quora-question-pairs/the-importance-of-cleaning-text

96.   def text_to_wordlist(text, remove_stopwords=False, stem_words=False):
97.       # Clean the text, with the option to remove stopwords and to stem words.
98.       # Convert words to lower case and split them
99.       text = text.lower().split()
100.
101.      # Optionally, remove stop words
102.      if remove_stopwords:
103.          stops = set(stopwords.words("english"))
104.          text = [w for w in text if not w in stops]
105.
106.      text = " ".join(text)
107.
108.      # Clean the text
109.      text = re.sub(r"[^A-Za-z0-9^,!.\/'+-=]", " ", text)
110.      text = re.sub(r"what's", "what is ", text)
111.      text = re.sub(r"\'s", " ", text)
112.      text = re.sub(r"\'ve", " have ", text)
113.      text = re.sub(r"can't", "cannot ", text)
114.      text = re.sub(r"n't", " not ", text)
115.      text = re.sub(r"i'm", "i am ", text)
116.      text = re.sub(r"\'re", " are ", text)
117.      text = re.sub(r"\'d", " would ", text)
118.      text = re.sub(r"\'ll", " will ", text)
119.      text = re.sub(r",", " ", text)
120.      text = re.sub(r"\.", " ", text)
121.      text = re.sub(r"!", " ! ", text)
122.      text = re.sub(r"\/", " ", text)
123.      text = re.sub(r"\^", " ^ ", text)
124.      text = re.sub(r"\+", " + ", text)
125.      text = re.sub(r"\-", " - ", text)
126.      text = re.sub(r"\=", " = ", text)
127.      text = re.sub(r"'", " ", text)
```

```python
128.        text = re.sub(r"(\d+)(k)", r"\g<1>000", text)
129.        text = re.sub(r":", " : ", text)
130.        text = re.sub(r" e g ", " eg ", text)
131.        text = re.sub(r" b g ", " bg ", text)
132.        text = re.sub(r" u s ", " american ", text)
133.        text = re.sub(r"\0s", "0", text)
134.        text = re.sub(r" 9 11 ", "911", text)
135.        text = re.sub(r"e - mail", "email", text)
136.        text = re.sub(r"j k", "jk", text)
137.        text = re.sub(r"\s{2,}", " ", text)
138.
139.        # Optionally, shorten words to their stems
140.        if stem_words:
141.            text = text.split()
142.            stemmer = SnowballStemmer('english')
143.            stemmed_words = [stemmer.stem(word) for word in text]
144.            text = " ".join(stemmed_words)
145.
146.        # Return a list of words
147.        return(text)
148.
149. ########################################
150. ## text preprocessing
151. ########################################
152. texts_1 = []
153. texts_2 = []
154. labels = []
155. with codecs.open(TRAIN_DATA_FILE, encoding='utf-8') as f:
156.     reader = csv.reader(f, delimiter=',')
157.     header = next(reader)
158.     for values in reader:
159.         texts_1.append(text_to_wordlist(values[3]))
160.         texts_2.append(text_to_wordlist(values[4]))
161.         labels.append(int(values[5]))
162. print('Found %s texts in train.csv' % len(texts_1))
163.
164. test_texts_1 = []
165. test_texts_2 = []
```

```python
166. test_ids = []
167. with codecs.open(TEST_DATA_FILE, encoding='utf-8') as f:
168.     reader = csv.reader(f, delimiter=',')
169.     header = next(reader)
170.     for values in reader:
171.         test_texts_1.append(text_to_wordlist(values[1]))
172.         test_texts_2.append(text_to_wordlist(values[2]))
173.         test_ids.append(values[0])
174. print('Found %s texts in test.csv' % len(test_texts_1))
175.
176. tokenizer = Tokenizer(num_words=MAX_NB_WORDS)
177. tokenizer.fit_on_texts(texts_1 + texts_2 + test_texts_1 + test_texts_2)
178.
179. sequences_1 = tokenizer.texts_to_sequences(texts_1)
180. sequences_2 = tokenizer.texts_to_sequences(texts_2)
181. test_sequences_1 = tokenizer.texts_to_sequences(test_texts_1)
182. test_sequences_2 = tokenizer.texts_to_sequences(test_texts_2)
183.
184. word_index = tokenizer.word_index
185. print('Found %s unique tokens' % len(word_index))
186.
187. data_1 = pad_sequences(sequences_1, maxlen=MAX_SEQUENCE_LENGTH)
188. data_2 = pad_sequences(sequences_2, maxlen=MAX_SEQUENCE_LENGTH)
189. labels = np.array(labels)
190. print('Shape of data tensor:', data_1.shape)
191. print('Shape of label tensor:', labels.shape)
192.
193. test_data_1 = pad_sequences(test_sequences_1, maxlen=MAX_SEQUENCE_LENGTH)
194. test_data_2 = pad_sequences(test_sequences_2, maxlen=MAX_SEQUENCE_LENGTH)
195. test_ids = np.array(test_ids)
196.
```

In this part, we transform the sentences to word2vec vectors. If words cannot be found in word2vec model, we will correct it by Peter Norvig's word corrector, and then check the new words are in word2vec model again. And then split some valuation data from training data.

```python
197. ########################################
198. ## prepare embeddings
199. ########################################
200. print('Preparing embedding matrix')
201.
202. nb_words = min(MAX_NB_WORDS, len(word_index))+1
203.
204. embedding_matrix = np.zeros((nb_words, EMBEDDING_DIM))
205. success_count=0
206. for word, i in word_index.items():
207.     if word in word2vec.vocab:
208.         embedding_matrix[i] = word2vec.word_vec(word)
209.     else:
210.         tempword=correct(word)
211.         if tempword in word2vec.vocab:
212.             success_count+=1
213.             print ('success:'+word+' to '+ tempword)
214.             embedding_matrix[i] = word2vec.word_vec(tempword)
215.
216. print('Null word embeddings: %d' % np.sum(np.sum(embedding_matrix, axis=1) == 0))
217. print('successfully correct ' + str(success_count) + ' token')
218.
219. ########################################
220. ## sample train/validation data
221. ########################################
222. perm = np.random.permutation(len(data_1))
223. idx_train = perm[:int(len(data_1)*(1-VALIDATION_SPLIT))]
224. idx_val = perm[int(len(data_1)*(1-VALIDATION_SPLIT)):]
225.
226. data_1_train = np.vstack((data_1[idx_train], data_2[idx_train]))
227. data_2_train = np.vstack((data_2[idx_train], data_1[idx_train]))
228. labels_train = np.concatenate((labels[idx_train], labels[idx_train]))
229.
230. data_1_val = np.vstack((data_1[idx_val], data_2[idx_val]))
231. data_2_val = np.vstack((data_2[idx_val], data_1[idx_val]))
232. labels_val = np.concatenate((labels[idx_val], labels[idx_val]))
233.
234. weight_val = np.ones(len(labels_val))
```

```
235. if re_weight:
236.     weight_val *= 0.472001959
237.     weight_val[labels_val==0] = 1.309028344
238.
```
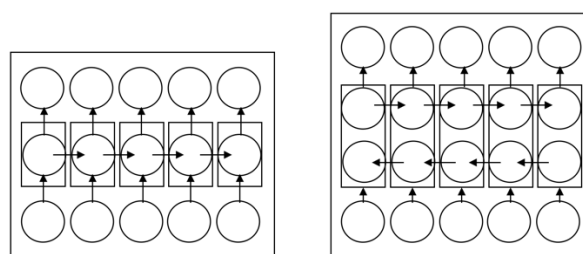
In this part, we construct our model. First, we take two sentences as two input layers, and connect to two embedding layers, and then connect to two LSTM layers, and take two LSTM output to DNN, and then a one-node output layer.

There are two feature we want to compare

1. LSTM vs. Bidirectional LSTM
2. 2 vs. 3 vs. 4 hidden layers of DNN

So we want to compare effect of 6 models

1. LSTM+2 hidden layers of DNN (50+20)
2. LSTM+3 hidden layers of DNN (100+50+20)
3. LSTM+4 hidden layers of DNN (150+100+50+20)
4. Bidirectional LSTM+2 hidden layers of DNN (50+20)
5. Bidirectional LSTM+3 hidden layers of DNN (100+50+20)
6. Bidirectional LSTM+4 hidden layers of DNN (150+100+50+20)



(a)                          (b)

Structure overview
(a) unidirectional RNN
(b) bidirectional RNN

After defining models, we train it, and then make prediction. And then submit the result to Kaggle to get accuracy.

```python
239. #########################################
240. ## define the model structure
241. ## 1. LSTM vs. Bidirectional LSTM
242. ## 2. 2 vs. 3 vs. 4 hidden layers of DNN
243. #########################################
244. embedding_layer = Embedding(nb_words,
245.         EMBEDDING_DIM,
246.         weights=[embedding_matrix],
247.         input_length=MAX_SEQUENCE_LENGTH,
248.         trainable=False)
249.
250. lstm_layer = LSTM(num_lstm, dropout=rate_drop_lstm, recurrent_dropout=rate_drop_lstm)
251. '''''
252. lstm_layer = Bidirectional(LSTM(num_lstm, dropout=rate_drop_lstm, recurrent_dropout=rat
     e_drop_lstm))
253. '''
254.
255. sequence_1_input = Input(shape=(MAX_SEQUENCE_LENGTH,), dtype='int32')
256. embedded_sequences_1 = embedding_layer(sequence_1_input)
257. x1 = lstm_layer(embedded_sequences_1)
258.
259. sequence_2_input = Input(shape=(MAX_SEQUENCE_LENGTH,), dtype='int32')
260. embedded_sequences_2 = embedding_layer(sequence_2_input)
261. y1 = lstm_layer(embedded_sequences_2)
262.
263. merged = concatenate([x1, y1])
264. merged = Dropout(rate_drop_dense)(merged)
265. merged = BatchNormalization()(merged)
266.
267. '''''
268. merged = Dense(150, activation=act)(merged)
269. merged = Dropout(rate_drop_dense)(merged)
```

```python
270. merged = BatchNormalization()(merged)

271. merged = Dense(100, activation=act)(merged)

272. merged = Dropout(rate_drop_dense)(merged)

273. merged = BatchNormalization()(merged)

274. '''

275. merged = Dense(50, activation=act)(merged)

276. merged = Dropout(rate_drop_dense)(merged)

277. merged = BatchNormalization()(merged)

278. merged = Dense(20, activation=act)(merged)

279. merged = Dropout(rate_drop_dense)(merged)

280. merged = BatchNormalization()(merged)

281.

282. preds = Dense(1, activation='sigmoid')(merged)

283. ########################################

284. ## add class weight

285. ########################################

286. if re_weight:

287.     class_weight = {0: 1.309028344, 1: 0.472001959}

288. else:

289.     class_weight = None

290.

291. ########################################

292. ## train the model

293. ########################################

294. model = Model(inputs=[sequence_1_input, sequence_2_input],

295.         outputs=preds)

296. model.compile(loss='binary_crossentropy',

297.         optimizer='nadam',

298.         metrics=['acc'])

299. model.summary()

300. print(STAMP)

301.

302. early_stopping =EarlyStopping(monitor='val_loss', patience=3)

303. bst_model_path = STAMP + '.h5'

304. model_checkpoint = ModelCheckpoint(bst_model_path, save_best_only=True, save_weights_on
    ly=True)

305.

306. hist = model.fit([data_1_train, data_2_train], labels_train,
```

```
307.          validation_data=([data_1_val, data_2_val], labels_val, weight_val),
308.          epochs=50, batch_size=1024, shuffle=True,
309.          class_weight=class_weight, callbacks=[early_stopping, model_checkpoint])
310.
311. model.load_weights(bst_model_path)
312. bst_val_score = min(hist.history['val_loss'])
313.
314. ########################################
315. ## make the submission
316. ########################################
317. print('Start making the submission before fine-tuning')
318.
319. preds = model.predict([test_data_1, test_data_2], batch_size=1024, verbose=1)
320. preds += model.predict([test_data_2, test_data_1], batch_size=1024, verbose=1)
321. preds /= 2
322.
323. submission = pd.DataFrame({'test_id':test_ids, 'is_duplicate':preds.ravel()})
324. submission.to_csv('%.4f_'%(bst_val_score)+STAMP+'.csv', index=False)
325.
```

## Comparison and result

We construct 6 models, and here are the results.

| model | Log_loss |
| --- | --- |
| LSTM+2 hidden layers of DNN | 0.31136 |
| LSTM+3 hidden layers of DNN | 0.30114 |
| LSTM+4 hidden layers of DNN | 0.30531 |
| Bidirectional LSTM+2 hidden layers of DNN | 0.31074 |
| Bidirectional LSTM+3 hidden layers of DNN | 0.30629 |
| Bidirectional LSTM+4 hidden layers of DNN | 0.31404 |

As you can see, in both LSTM and bidirectional LSTM models, 3 hidden layers of DNN has the best accuracy. Maybe it is the result of deeper DNN is harder to train, or even overfitting.

On the other hand, the accuracy of single layer LSTM model is higher than bidirectional LSTM model. It is different from the results of most of the theses. We think this is because we haven't found the right hyperparameters of bidirectional LSTM. We believe that if we have more time, we can find better hyperparameters of all models, and maybe our result will match results of theses.

Most important, all our models have loss around 0.3, and in Kaggle data analysis competition, we have rank of 963 out of 3307. There are some reasons why our accuracy is so low.

First, we haven't found the right hyperparameters. We didn't have enough time to try others hyperparameters because training deep learning models cost a lot of times (6 to 8 hours per time).

Second, using deep learning is not the best idea. In the e-book "neural network and deep learning", we learned that deep learning can do everything, however, it is always hard to train. After reading a lot of articles in Kaggle discuss board, we found that most people have higher rank are not using deep learning, they use "Xgboost", a model that combine many easier model (and therefore easier to train) to get better prediction, instead. Furthermore, I have heard that in Kaggle, half of the winner use DNN, and the others use Xgboost. I believe that both of them are strong ML model, and maybe some problem suit DNN, and some suit Xgboost.

Third, training data is not big enough. We only have 404290 sentence pairs in training data, but 2345797 testing data (over 5 times), and it seems like we should do more to handle such situation.

## Conclusion

In this project, we used deep learning to implement sentences comparison, and we compare results of different structure of models. We know that deep learning is a good way to solve use, since it is not so hard to write DNN code in Keras, and can easily get a not-so-bad result. However, if we need to be more accurate, we have to spend a lot of time to tune our model, do more preprocessing, or even try different model to choose the best.

## Reference

[1] Jonas Mueller and Aditya Thyagarajan, "Siamese Recurrent Architectures for Learning Sentence Similarity," Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence (AAAI-16)