



SHANGHAI JIAO TONG UNIVERSITY

X033525

MACHINE LEARNING

---

## Quora Question Pairs @ Kaggle

---

*Author:*

Yuliang Xiu  
Xiaoting Jiang  
Weiyu Cheng  
Bowen Zhang

*Student Number:*

116033910045  
116033910035  
016033910007  
016033910035

June 7, 2017

Contents

<b>1</b>	<b>Problem Description</b>	<b>2</b>
1.1	Background . . . . .	2
1.2	Definition . . . . .	2
<b>2</b>	<b>Related Work</b>	<b>3</b>
<b>3</b>	<b>Team Division</b>	<b>4</b>
<b>4</b>	<b>Dataset Preprocessing</b>	<b>4</b>
4.1	Dataset Attribute . . . . .	4
4.2	Data Cleaning . . . . .	5
<b>5</b>	<b>Manual Feature Extraction</b>	<b>5</b>
<b>6</b>	<b>Sentence Embedding</b>	<b>6</b>
<b>7</b>	<b>Classification Model</b>	<b>7</b>
<b>8</b>	<b>Final Result</b>	<b>8</b>

# 1 Problem Description

## 1.1 Background

Where else but Quora can a physicist help a chef with a math problem and get cooking tips in return? 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. This empowers people to learn from each other and to better understand the world.

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. Quora values canonical questions because they provide a better experience to active seekers and writers, and offer more value to both of these groups in the long term.

Currently, Quora uses a Random Forest model to identify duplicate questions. In this competition, Kagglers are challenged to tackle this natural language processing problem by applying advanced techniques to classify whether question pairs are duplicates or not. Doing so will make it easier to find high quality answers to questions resulting in an improved experience for Quora writers, seekers, and readers.

## 1.2 Definition

More formally, the duplicate detection problem can be defined as follows: given a pair of questions  $q_1$  and  $q_2$ , train a model that learns the function, where 1 represents that  $q_1$  and  $q_2$  have the same intent and 0 otherwise:

$$f(q_1, q_2) \rightarrow 0 \quad \text{or} \quad 1$$

Submission is evaluated by logloss between predicted value and ground truth:

$$\text{logloss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(P_{ij})$$

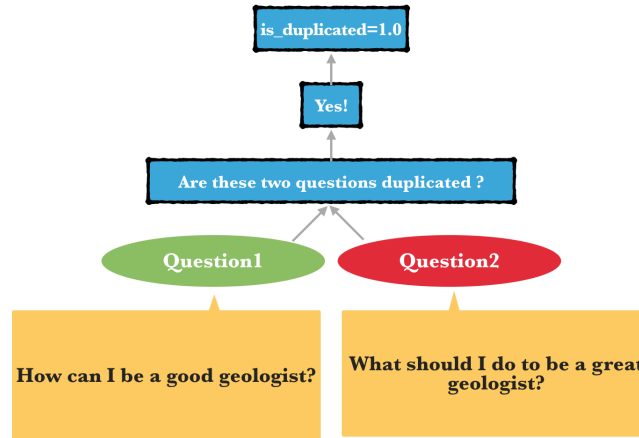


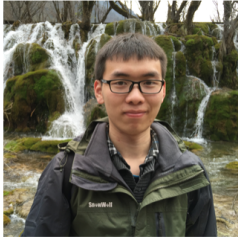
Figure 1: Problem definition diagram

## 2 Related Work

Some paraphrase identification related papers: [6][4][10][4][2][9][1][3]

Model	Source of Word Embeddings	Accuracy
"BiMPM model" [5]	GloVe Common Crawl (840B tokens, 300D)	<b>0.88</b>
"LSTM with concatenation" [6]	"Quora's text corpus"	0.87
"LSTM with distance and angle" [6]	"Quora's text corpus"	0.87
"Decomposable attention" [6]	"Quora's text corpus"	0.86
"L.D.C." [5]	GloVe Common Crawl (840B tokens, 300D)	0.86
Max bag-of-embeddings ( <i>this work</i> )	GloVe Common Crawl (840B tokens, 300D)	0.83
"Multi-Perspective-LSTM" [5]	GloVe Common Crawl (840B tokens, 300D)	0.83
"Siamese-LSTM" [5]	GloVe Common Crawl (840B tokens, 300D)	0.83
"Neural bag-of-words" (max) [7]	GloVe Common Crawl pruned to 1M vocab. (spaCy default)	0.83
"Neural bag-of-words" (max & mean) [7]	GloVe Common Crawl pruned to 1M vocab. (spaCy default)	0.83
"Max-out Window Encoding" with depth 2 [7]	GloVe Common Crawl pruned to 1M vocab. (spaCy default)	0.83
"Neural bag-of-words" (mean) [7]	GloVe Common Crawl pruned to 1M vocab. (spaCy default)	0.81
"Multi-Perspective-CNN" [5]	GloVe Common Crawl (840B tokens, 300D)	0.81
"Siamese-CNN" [5]	GloVe Common Crawl (840B tokens, 300D)	0.80
"Spacy + TD-IDF + Siamese" [8]	GloVe (6B tokens, 300D)	0.79

### 3 Team Division



**Bowen Zhang**  
DDST Ph.D.



**Weiyu Cheng**  
DDST Ph.D.



**Xiaoting Jiang**  
DDST M.Phil.



**Yuliang Xiu\***  
MVIG M.Phil.

\* Team leader

**Original  
Dataset  
Cleaning**

**Manual  
Feature  
Engineering**

**Sentence to  
Vector  
Embedding**

**Classifier  
Design and  
Boosting**

### 4 Dataset Preprocessing

#### 4.1 Dataset Attribute

- Training dataset size : 404290
- Testing dataset size: 2345896
- Dataset format:

	id	qid1	qid2	question1	question2	is_duplicate
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0
3	3	7	8	Why am I mentally very lonely? How can I solve...	Find the remainder when $23^{24}$ i...	0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

## 4.2 Data Cleaning

- **Correct wrong labels:** There are many wrong labeled records in original dataset, which are caused by redundant blank. Before training phrase, we must correct these records to make sure all records are labeled correctly.

question1: How can I be a good geologist?

question2: How can I be a good geologist ?

i\_duplicated: 0 → 1

- **Remove special characters:** ℵ, ♣, © ...
- **Change abbreviation:** what's → what is can't → cannot
- **Use standard format:** e-mail → email
- **Replace reference words:** Due to quora is an american QA website, so when users use "us", it always refers to american, so we change all us → american.

## 5 Manual Feature Extraction

We extract 28 kinds of manual features between ques1 and ques2 as follows:

Manual Features	Description	Manual Features	Description
word_match	number of common words	fuzz_qratio	<b>Fuzz Ratio</b>
tf_idf_word_match	$\text{sum}(\text{tf\_idf}[q1 \& q2]) / \text{sum}(\text{tf\_idf}[q1]   [q2])$	fuzz_wratio	
sentiment	sentiment ( positive or negative )	fuzz_partial_ratio	
diff_len	<b>Length Difference</b>	fuzz_partial_token_set_ratio	
diff_len_char		fuzz_partial_token_sort_ratio	
diff_len_word		fuzz_token_set_ratio	
q1_freq	frequency of question1 in all data	fuzz_token_sort_ratio	
q2_freq	frequency of question2 in all data	cosine_distance	<b>Distance</b>
q1_q2_intersect	intersection of q1&q2 related sents	cityblock_distance	
q1_q2_wm_ratio	$\text{sum}(\text{inter}[q1 \& q2]) / \text{sum}(\text{inter}[q1]   [q2])$	jaccard_distance	
wmd	Word Mover's Distance	canberra_distance	
norm_wmd	norm of Word Mover's Distance	euclidean_distance	
skew_vec	skewness of sentence	minkowski_distance	
kur_vec	kurtosis of sentence	braycurtis_distance	

Figure 2: Manual feature table

## 6 Sentence Embedding

During sentence embedding phrase, we firstly remove stop words and punctuations from question sentence, and then do stemming on the rest sentence. Secondly, we use keras *text\_to\_sequence* function to convert sentence to number sequence, and then do padding operation (30 dimensions) on sequence. Thirdly, use pre-trained embedding matrix such as **Glove**[7] or **FastText**[5] to do embedding operation, which return a  $30 \times 300$  matrix. Finally, put this matrix into Bidirection-LSTM[8] to extract sequential information from this sequence.

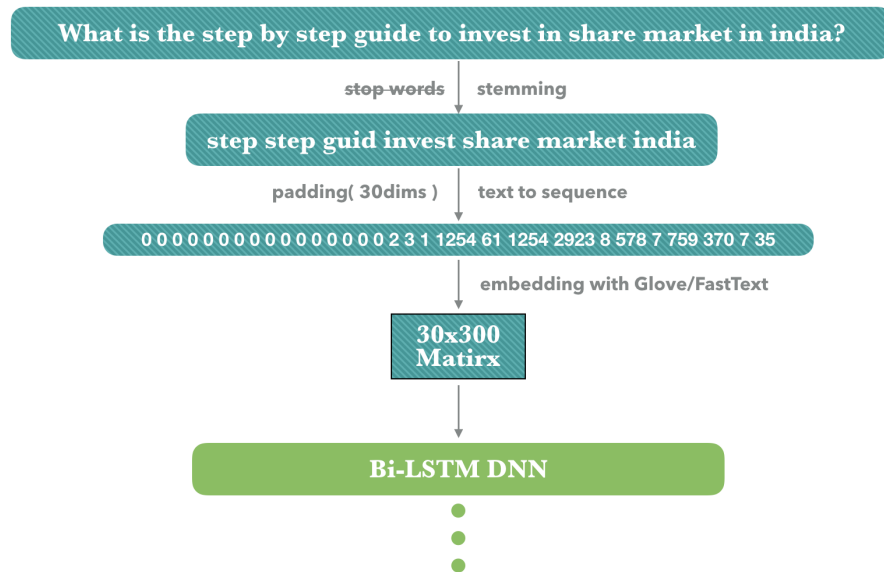


Figure 3: Sentence embedding procedure

---

```

embedding_layer = Embedding(nb_words,
    EMBEDDING_DIM,
    weights=[embedding_matrix],
    input_length=MAX_SEQUENCE_LENGTH,
    trainable=False)
lstm_layer = Bidirectional(LSTM(num_lstm, dropout=rate_drop_lstm,
    ↪ recurrent_dropout=rate_drop_lstm))
  
```

---

## 7 Classification Model

Input of network is one question pair. To extract sequential information from questions, we use two separately Bidirection-LSTM on question1 and question2 embedding matrix. Meanwhile, we calculate manual features or traditional features between question1 and question2, and use fully connected layer to align dimension to 200. Finally we concatenate two sequential features and manual features together and put forward through one fully connected layer to get final predicted value.

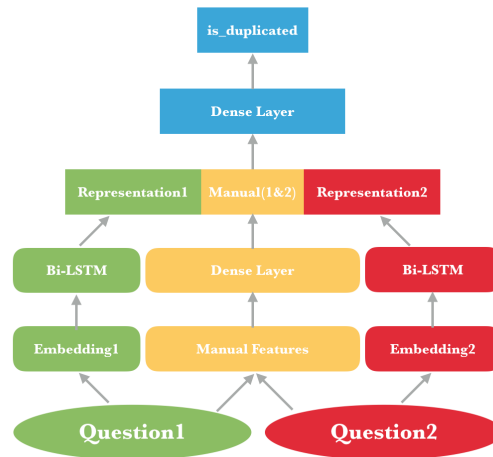


Figure 4: Classifier network framework

---

```

sequence_1_input = Input(shape=(MAX_SEQUENCE_LENGTH,), dtype='int32')
embedded_sequences_1 = embedding_layer(sequence_1_input)
x1 = lstm_layer(embedded_sequences_1)
sequence_2_input = Input(shape=(MAX_SEQUENCE_LENGTH,), dtype='int32')
embedded_sequences_2 = embedding_layer(sequence_2_input)
y1 = lstm_layer(embedded_sequences_2)
z1 = Input(shape=(x_train.shape[1],), dtype='float32')
z1_dense = Dense(num_dense/2, activation=act)(z1)
merged = concatenate([x1, y1, z1_dense])
merged = BatchNormalization()(merged)
merged = Dropout(rate_drop_dense)(merged)
merged = Dense(num_dense, activation=act)(merged)
merged = BatchNormalization()(merged)
merged = Dropout(rate_drop_dense)(merged)
preds = Dense(1, activation='sigmoid')(merged)

```

---



# 8 Final Result

- Best Rank: 91/3379 (top 3%)
- Prize: Silver Medal (within top 5%)
- Loss value: 0.14128
- Accuracy: 90.13%
- Competition page: [Leaderboard of quora question pair](#)
- Github code: [kaggle-quora@github](#)

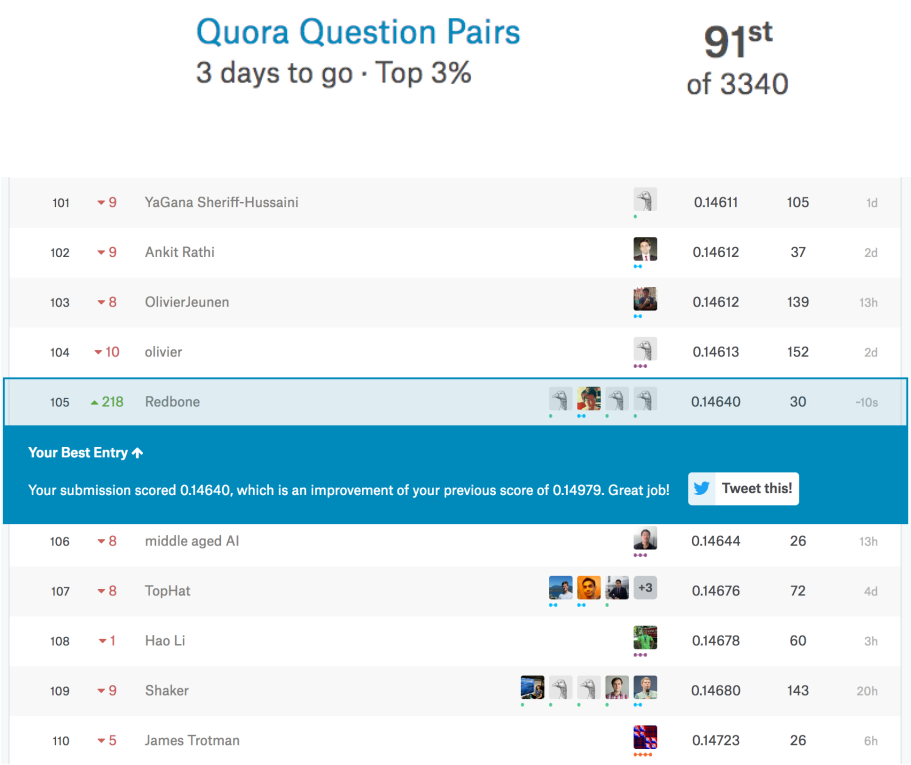


Figure 5: Final rank

## References

- [1] *Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks*, 2015.
- [2] *A Decomposable Attention Model for Natural Language Inference*, 2016.
- [3] William Blacoe and Mirella Lapata. A Comparison of Vector-based Representations for Semantic Composition. *EMNLP-CoNLL*, 2012.
- [4] Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. *arXiv.org*, page arXiv:1508.05326, August 2015.
- [5] Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, H erve J egou, and Tomas Mikolov. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016.
- [6] Tomas Mikolov, Kai Chen 0010, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. *CoRR*, cs.CL:arXiv:1301.3781, 2013.
- [7] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove - Global Vectors for Word Representation. *EMNLP*, 2014.
- [8] M Schuster and K K Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45:2673–2681, November 1997.
- [9] Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved Semantic Representations From Tree-Structured Long Short-Term Memory Networks. *CoRR*, 2015.
- [10] Zhiguo Wang, Wael Hamza, and Radu Florian. Bilateral Multi-Perspective Matching for Natural Language Sentences. *CoRR*, cs.AI:arXiv:1702.03814, 2017.