

机器学习 Quora Question Pairs 实验报告

数据科学与计算机学院 计算机科学与技术 2016 级

王凯祺 16337233

2019 年 6 月 12 日

1 使用的系统 + 编码语言 + 环境说明

操作系统: macOS Mojave 10.14.5

代码语言: Python 3

环境说明:

```
1 pip3 freeze
```

```
1 boto==2.49.0
2 boto3==1.9.162
3 botocore==1.12.162
4 certifi==2019.3.9
5 chardet==3.0.4
6 cycler==0.10.0
7 docutils==0.14
8 gensim==3.7.3
9 idna==2.8
10 jmespath==0.9.4
11 joblib==0.13.2
12 kiwisolver==1.1.0
13 matplotlib==3.1.0
14 nltk==3.4.1
15 numpy==1.16.4
16 pandas==0.24.2
17 pyparsing==2.4.0
18 python-dateutil==2.8.0
19 python-Levenshtein==0.12.0
20 pytz==2019.1
21 requests==2.22.0
22 s3transfer==0.2.1
23 scikit-learn==0.21.2
24 scipy==1.3.0
25 seaborn==0.9.0
26 six==1.12.0
27 sklearn==0.0
28 smart-open==1.8.4
29 urllib3==1.25.3
```

2 简单的流程图

2.1 数据处理方法

2.1.1 观察数据特征

训练集共 404,290 行、6 列，每行包含 id、问题 1 id、问题 2 id、问题 1、问题 2、是否重复 (标签)。

测试集共 2,345,796 行、3 列，每行包含 id、问题 1、问题 2。

大致看了几个问题，就觉得处理起来挺棘手的。比如“What is the step by step guide to invest in share market in india?”和“What is the step by step guide to invest in share market?”是不重复的两个问题，但它们之间只差了一个“in India”。这直接导致一个问题，如果我们只根据两个句子的相似度来判定是否重复，这样的问题组合就会被判为重复。再比如，“How can I be a good geologist?”和“What should I do to be a great geologist?”是重复的，我们要让机器自己学到“How”和“What should I do”是同一个意思。这必须要有相当数量的同类样本才能学到的。

在如此大量的测试集面前，却只有这么点训练集。我觉得很难从训练集中提取到太多有用的东西，然后应用到测试集中。

2.1.2 数据清洗

我参考了一个名为 The Importance of Cleaning Text 的 Kernel (<https://www.kaggle.com/currie32/the-importance-of-cleaning-text>)，对输入数据进行清洗。

清洗主要做的任务是：

- 将非 ascii 字符全部删去，规范化单词，如将“re”替换为“are”。
- 删除标点符号。
- 删除停词，如 is、the、a。
- 将单词处理成原型。

主要代码：

```
1 def text_to_wordlist(text, remove_stop_words=True, stem_words=False):
2     # Clean the text, with the option to remove stop_words and to stem words.
3     global stops
4     text = re.sub(r"[^A-Za-z0-9]", "_", text)
5     text = re.sub(r"what's", "", text)
6     text = re.sub(r"What's", "", text)
7     text = re.sub(r"\'s", "_", text)
8     text = re.sub(r"\'ve", "_have_", text)
9     text = re.sub(r"can't", "cannot_", text)
10    text = re.sub(r"n't", "_not_", text)
11    text = re.sub(r"I'm", "I_am", text)
12    text = re.sub(r"_m_", "_am_", text)
```

```

13 text = re.sub(r"\'re", "_are_", text)
14 text = re.sub(r"\'d", "_would_", text)
15 text = re.sub(r"\'ll", "_will_", text)
16 text = re.sub(r"60k", "_60000_", text)
17 text = re.sub(r"_e_g_", "_eg_", text)
18 text = re.sub(r"_b_g_", "_bg_", text)
19 text = re.sub(r"\0s", "0", text)
20 text = re.sub(r"_9_ll_", "911", text)
21 text = re.sub(r"e-mail", "email", text)
22 text = re.sub(r"\s{2,}", "_", text)
23 text = re.sub(r"quikly", "quickly", text)
24 text = re.sub(r"_usa_", "_America_", text)
25 text = re.sub(r"_USA_", "_America_", text)
26 text = re.sub(r"_u_s_", "_America_", text)
27 text = re.sub(r"_uk_", "_England_", text)
28 text = re.sub(r"_UK_", "_England_", text)
29 text = re.sub(r"india", "India", text)
30 text = re.sub(r"switzerland", "Switzerland", text)
31 text = re.sub(r"china", "China", text)
32 text = re.sub(r"chinese", "Chinese", text)
33 text = re.sub(r"imrovement", "improvement", text)
34 text = re.sub(r"intially", "initially", text)
35 text = re.sub(r"quora", "Quora", text)
36 text = re.sub(r"_dms_", "direct_messages_", text)
37 text = re.sub(r"demonitization", "demonetization", text)
38 text = re.sub(r"actived", "active", text)
39 text = re.sub(r"kms", "_kilometers_", text)
40 text = re.sub(r"KMs", "_kilometers_", text)
41 text = re.sub(r"_cs_", "_computer_science_", text)
42 text = re.sub(r"_upvotes_", "_up_votes_", text)
43 text = re.sub(r"_iPhone_", "_phone_", text)
44 text = re.sub(r"\0rs_", "_rs_", text)
45 text = re.sub(r"calender", "calendar", text)
46 text = re.sub(r"ios", "operating_system", text)
47 text = re.sub(r"gps", "GPS", text)
48 text = re.sub(r"gst", "GST", text)
49 text = re.sub(r"programing", "programming", text)
50 text = re.sub(r"bestfriend", "best_friend", text)
51 text = re.sub(r"dna", "DNA", text)
52 text = re.sub(r"III", "3", text)
53 text = re.sub(r"the_US", "America", text)
54 text = re.sub(r"Astrology", "astrology", text)
55 text = re.sub(r"Method", "method", text)
56 text = re.sub(r"Find", "find", text)
57 text = re.sub(r"banglore", "Banglore", text)
58 text = re.sub(r"_J_K_", "_JK_", text)
59
60 # Remove punctuation from text
61 text = ''.join([c for c in text if c not in punctuation])
62

```

```

63     # Optionally, remove stop words
64     if remove_stop_words:
65         text = text.split()
66         text = [w for w in text if not w in stops]
67         text = "_".join(text)
68
69     # Optionally, shorten words to their stems
70     if stem_words:
71         text = text.split()
72         stemmer = SnowballStemmer('english')
73         stemmed_words = [stemmer.stem(word) for word in text]
74         text = "_".join(stemmed_words)
75
76     # Return a list of words
77     return text
78
79
80 def clean_data_frame(row):
81     row['question1'] = text_to_wordlist(str(row['question1']))
82     row['question2'] = text_to_wordlist(str(row['question2']))
83     return row
84
85
86 df_train = df_train.apply(clean_data_frame, axis=1, raw=True)
87 df_test = df_test.apply(clean_data_frame, axis=1, raw=True)

```

2.1.3 数据处理

我参考了一个名为 Data Analysis & XGBoost Starter (0.35460 LB) 的 Kernel，它能在二维坐标轴上反映训练集的词语匹配程度以及 TF-IDF。我在它的基础上加了 word2vec 相似度匹配程度、word2vec 词语间平均距离、编辑距离、编辑距离比例这些评估参数。

词语匹配程度 设问题 1 的单词总数为 n ，问题 2 的单词总数为 m ，问题 1 出现的单词中也在问题 2 出现的个数为 a ，问题 2 出现的单词中也在问题 1 出现的个数为 b ，词语匹配程度是这样定义的：

$$\frac{a+b}{n+m}$$

实现代码如下：

```

1 def word_match_share(row):
2     global stops
3     q1words = {}
4     q2words = {}
5     for word in str(row['question1']).lower().split():
6         if word not in stops:
7             q1words[word] = 1
8     for word in str(row['question2']).lower().split():
9         if word not in stops:
10            q2words[word] = 1

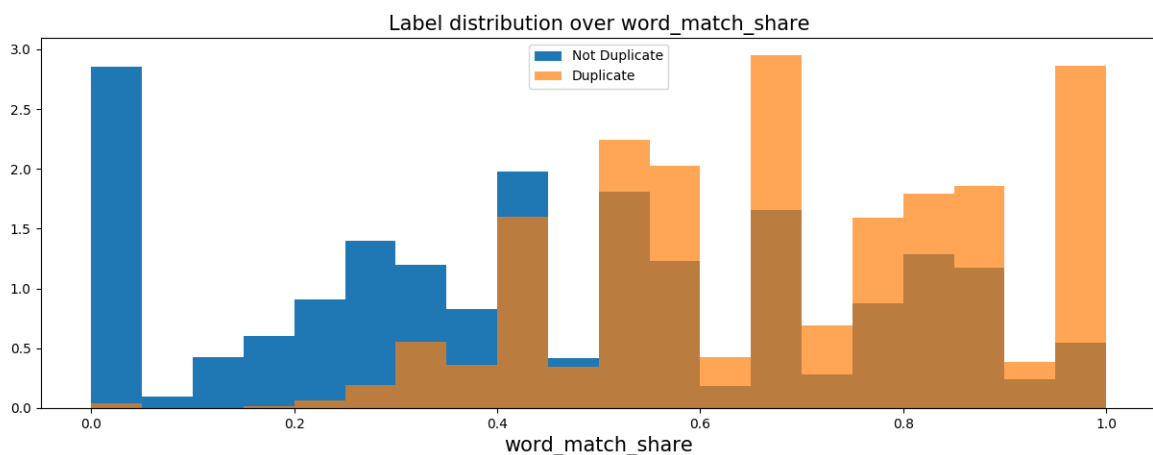
```

```

11     if len(q1words) == 0 or len(q2words) == 0:
12         # The computer-generated chaff includes a few questions that are nothing but
           stopwords
13         return 0
14     shared_words_in_q1 = [w for w in q1words.keys() if w in q2words]
15     shared_words_in_q2 = [w for w in q2words.keys() if w in q1words]
16     R = (len(shared_words_in_q1) + len(shared_words_in_q2)) / (len(q1words) + len(
           q2words))
17     return R
18
19
20 def print_plot(data, df_train, title, filename=None):
21     plt.figure(figsize=(15, 5))
22     plt.hist(data[df_train['is_duplicate'] == 0], bins=20, density=True, label='Not_
           Duplicate')
23     plt.hist(data[df_train['is_duplicate'] == 1], bins=20, density=True, alpha=0.7,
           label='Duplicate')
24     plt.legend()
25     plt.title(title, fontsize=15)
26     plt.xlabel('word_match_share', fontsize=15)
27     if filename is None or filename == '':
28         plt.show()
29     else:
30         plt.savefig(filename)
31
32
33 train_word_match = df_train.apply(word_match_share, axis=1, raw=True)
34 print_plot(train_word_match, df_train, 'Label_distribution_over_word_match_share', '
           1.png')

```

可视化效果如下：



我们可以看到，左侧几乎都是非重复的，当两个句子的单词差别较大时，容易判断它是非重复的；但是右侧重复和不重复的有很多重叠部分，分类效果并不好。

TF-IDF TF-IDF 是一种统计方法，用以评估一字词对于一个文件集或一个语料库中的其中一份文件的重要程度。我们记 $n_{i,j}$ 为第 i 个词在文件 d_j 中的出现次数， $|D|$ 表示语料库的文件总数。有公式：

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

$$tfidf_{i,j} = tf_{i,j} \times idf_i$$

我们把“在问题 2 中出现过的”问题 1 的单词的 TF-IDF 值、以及“在问题 1 中出现过的”问题 2 的单词的 TF-IDF 值求和，然后除以问题 1、问题 2 所有单词的 TF-IDF 值总和，以此作为这个问题组合的特征值。

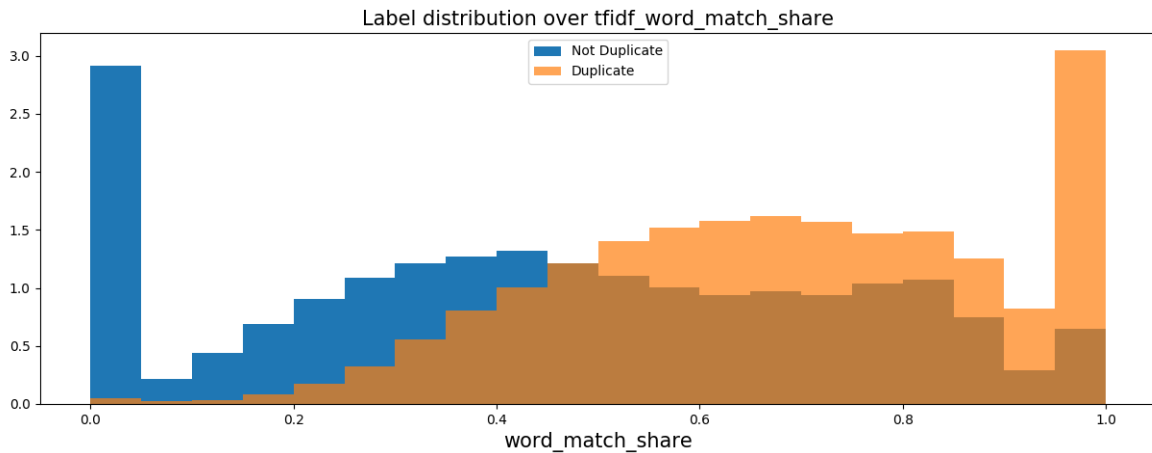
实现代码如下：

```

1 def tfidf_word_match_share(row):
2     global stops, weights, counts, words
3     q1words = {}
4     q2words = {}
5     for word in str(row['question1']).lower().split():
6         if word not in stops:
7             q1words[word] = 1
8     for word in str(row['question2']).lower().split():
9         if word not in stops:
10             q2words[word] = 1
11     if len(q1words) == 0 or len(q2words) == 0:
12         # The computer-generated chaff includes a few questions that are nothing but
13         # stopwords
14         return 0
15     shared_weights = [weights.get(w, 0) for w in q1words.keys() if w in q2words] + [
16         weights.get(w, 0) for w in q2words.keys() if w in q1words]
17     total_weights = [weights.get(w, 0) for w in q1words] + [weights.get(w, 0) for w
18         in q2words]
19     R = np.sum(shared_weights) / np.sum(total_weights)
20     return R
21
22 train_qs = pd.Series(df_train['question1'].tolist() + df_train['question2'].tolist()
23     ).astype(str)
24 test_qs = pd.Series(df_test['question1'].tolist() + df_test['question2'].tolist()).
25     astype(str)
26 words = (" ".join(train_qs)).lower().split()
27 counts = Counter(words)
28 weights = {word: get_weight(count) for word, count in counts.items()}
29 tfidf_train_word_match = df_train.apply(tfidf_word_match_share, axis=1, raw=True)
30 print_plot(tfidf_train_word_match, df_train, 'Label_distribution_over_
31     tfidf_word_match_share', '2.png')

```

我们看看这个特征值对样本的分类效果如何。可视化效果如下：



word2vec 相似度匹配 我们可以将问题里的句子导入 word2vec 模型，每个单词可以计算出一个向量。我们可以对于两个问题的每个单词组，计算一个平均相似度。

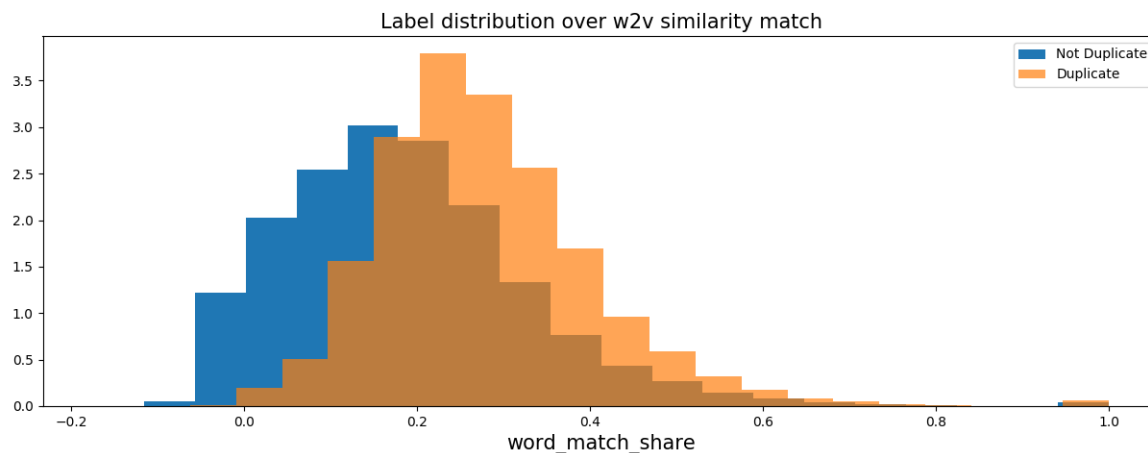
代码如下：

```

1 def build_corpus(data):
2     corpus = []
3     for col in ['question1', 'question2']:
4         for sentence in data[col].iteritems():
5             word_list = str(sentence[1]).lower().split("_")
6             corpus.append(word_list)
7     return corpus
8
9 corpus_train = build_corpus(df_train)
10 corpus_test = build_corpus(df_test)
11 corpus = corpus_train + corpus_test
12 w2vmodel = word2vec.Word2Vec(corpus, size=100, window=20, min_count=500, workers=4)
13
14 def word2vec_similarity_match(row):
15     global stops, w2vmodel
16     qcnt = 0
17     scnt = 0
18     for word1 in str(row['question1']).lower().split():
19         if word1 not in stops:
20             for word2 in str(row['question2']).lower().split():
21                 if word2 not in stops:
22                     qcnt += 1
23                     if word1 in w2vmodel and word2 in w2vmodel:
24                         scnt += w2vmodel.wv.similarity(word1, word2)
25     if qcnt == 0:
26         return 0
27     return scnt / qcnt
28
29 train_w2v_similarity_match = df_train.apply(word2vec_similarity_match, axis=1, raw=
    True)
30 print_plot(train_w2v_similarity_match, df_train, 'Label_distribution_over_w2v_
    similarity_match', '3.png')

```

可视化效果如下：



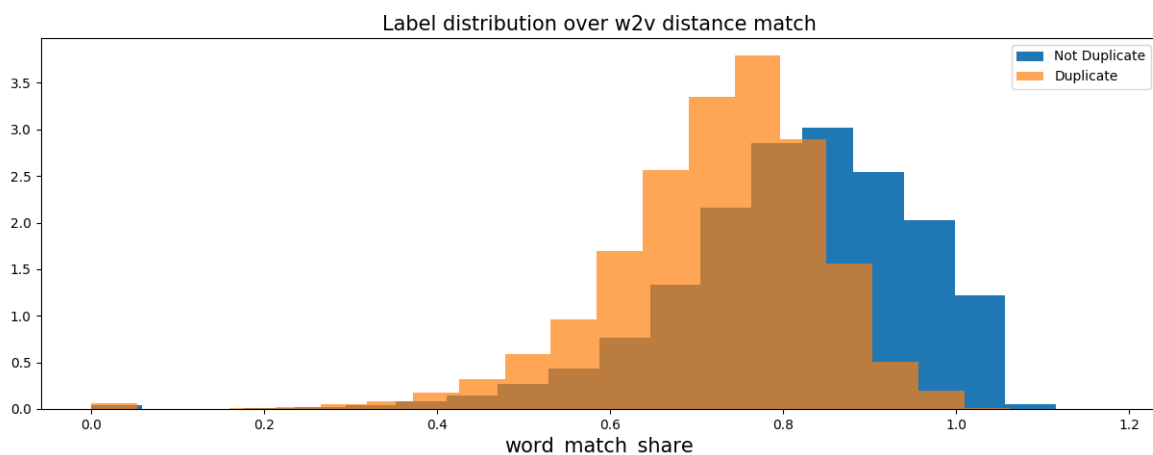
重叠部分还是相当多的，利用这个也挺难区分的。我们要想出尽可能多的分类特征，以提高分类准确度。

word2vec 平均距离 对于两个问题的每个单词组，计算一个平均余弦距离。

代码如下：

```
1 def word2vec_distance_match(row):
2     global stops, w2vmodel
3     qcnt = 0
4     scnt = 0
5     for word1 in str(row['question1']).lower().split():
6         if word1 not in stops:
7             for word2 in str(row['question2']).lower().split():
8                 if word2 not in stops:
9                     qcnt += 1
10                    if word1 in w2vmodel and word2 in w2vmodel:
11                        scnt += w2vmodel.wv.distance(word1, word2)
12                    else:
13                        scnt += 1.0
14
15     if qcnt == 0:
16         return 1
17     return scnt / qcnt
18
19 train_w2v_distance_match = df_train.apply(word2vec_distance_match, axis=1, raw=True)
20 print_plot(train_w2v_distance_match, df_train, 'Label_distribution_over_w2v_distance_match', '4.png')
```

可视化效果如下：



跟上一个图好像正好反过来，好像并没有什么太大的变化。

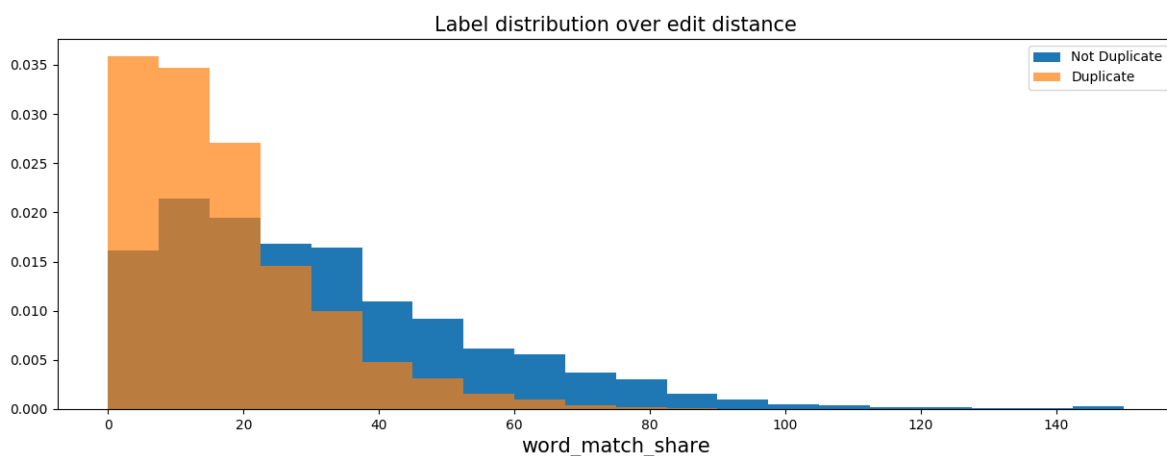
编辑距离 我看了一下网上的评论，发现编辑距离也可以作为一项指标，就也拿来试试。

两个字符串的编辑距离是这样定义的：由一个字符串开始进行插入、修改、删除操作，变成另一个字符串的最少修改次数。这是一个比较简单的动态规划问题。考虑到两问题串的长度都较短， $O(n^2)$ 的时间复杂度我们可以接受，就尝试这种方法。

代码如下：

```
1 def edit_distance(row):
2     ans = Levenshtein.distance(str(row['question1']), str(row['question2']))
3     if ans > 150:
4         ans = 150
5     return ans
6
7 train_edit_distance = df_train.apply(edit_distance, axis=1, raw=True)
8 print_plot(train_edit_distance, df_train, 'Label_distribution_over_edit_distance', '
    5.png')
```

可视化效果如下：

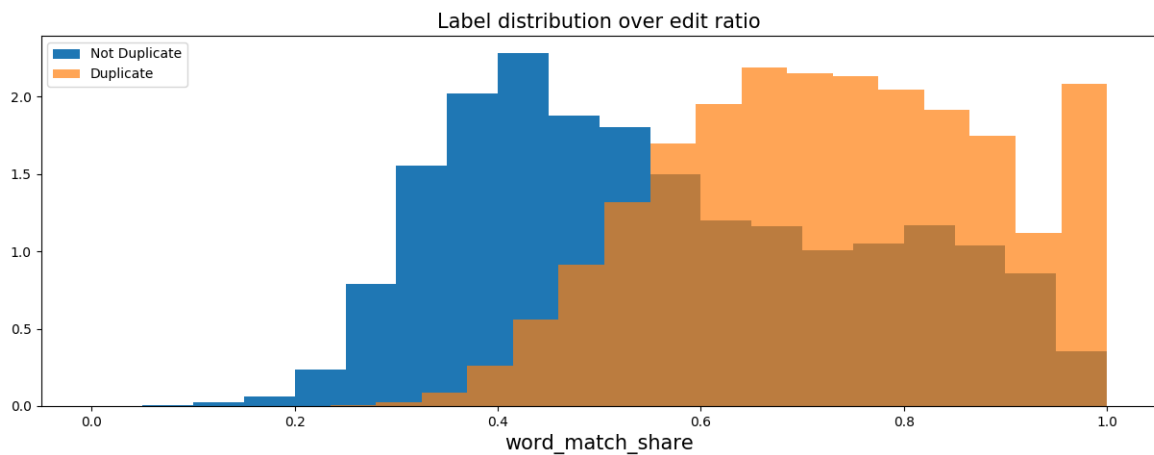


编辑距离比例 单纯的编辑距离可能在数值上过于绝对，我们将编辑距离除以两字符串长度的最大值，得到编辑距离比例。

代码如下：

```
1 def edit_ratio(row):
2     ans = Levenshtein.ratio(str(row['question1']), str(row['question2']))
3     return ans
4
5 train_edit_ratio = df_train.apply(edit_ratio, axis=1, raw=True)
6 print_plot(train_edit_ratio, df_train, 'Label_distribution_over_edit_ratio', '6.png'
7            )
```

可视化效果如下：



2.2 XGBoost 模型

2.2.1 实现

对于每条数据，只利用上述的 6 个特征，我们写一段简单的 XGBoost 代码，跑一下训练集，并预测测试集。

```
1 x_train = pd.DataFrame()
2 x_test = pd.DataFrame()
3 x_train['word_match'] = train_word_match
4 x_train['tfidf_word_match'] = tfidf_train_word_match
5 x_train['w2v_similarity_match'] = train_w2v_similarity_match
6 x_train['w2v_distance_match'] = train_w2v_distance_match
7 x_train['edit_distance'] = train_edit_distance
8 x_train['edit_ratio'] = train_edit_ratio
9 x_test['word_match'] = df_test.apply(word_match_share, axis=1, raw=True)
10 x_test['tfidf_word_match'] = df_test.apply(tfidf_word_match_share, axis=1, raw=True)
11 x_test['w2v_similarity_match'] = df_test.apply(word2vec_similarity_match, axis=1,
12        raw=True)
12 x_test['w2v_distance_match'] = df_test.apply(word2vec_distance_match, axis=1, raw=
13        True)
13 x_test['edit_distance'] = df_test.apply(edit_distance, axis=1, raw=True)
```

```

14 x_test['edit_ratio'] = df_test.apply(edit_ratio, axis=1, raw=True)
15 y_train = df_train['is_duplicate'].values
16
17 # Set our parameters for xgboost
18 params = {}
19 params['objective'] = 'binary:logistic'
20 params['eval_metric'] = 'logloss'
21 params['eta'] = 0.02
22 params['max_depth'] = 4
23
24 d_train = xgb.DMatrix(x_train, label=y_train)
25 d_valid = xgb.DMatrix(x_valid, label=y_valid)
26
27 watchlist = [(d_train, 'train'), (d_valid, 'valid')]
28
29 bst = xgb.train(params, d_train, 10000, watchlist, early_stopping_rounds=50,
30                 verbose_eval=10)
31
32 d_test = xgb.DMatrix(x_test)
33 p_test = bst.predict(d_test)
34
35 sub = pd.DataFrame()
36 sub['test_id'] = df_test['test_id']
37 sub['is_duplicate'] = p_test
38 sub.to_csv('xgb_6.csv', index=False)

```

2.2.2 参数的调整

为了调整参数，我们采用控制变量法，固定其它所有变量，调整某一变量。我们假设各个变量之间是相互独立的。

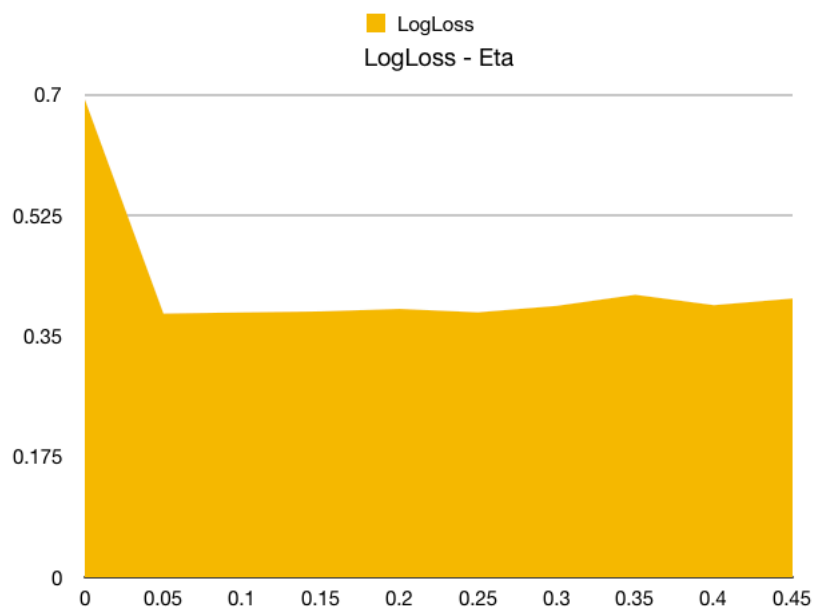
eta (收缩步长) 我们首先从 eta (收缩步长) 入手。eta 的取值范围是从 [0, 1]，默认值是 0.3。我们对 0 - 0.45 每隔 0.05 的数值（共计 10 个数值）都跑一遍，然后交到 kaggle 上测试 log loss 值。

跑一次要花 5 个小时……心好累……

幸运的是，我还是跑出来了。下表为 eta 和 LogLoss 的关系：

Eta	LogLoss
0	0.69314
0.05	0.38264
0.1	0.38432
0.15	0.38579
0.2	0.38946
0.25	0.38444
0.3	0.39383
0.35	0.40996
0.4	0.39487
0.45	0.40454

可视化的结果：

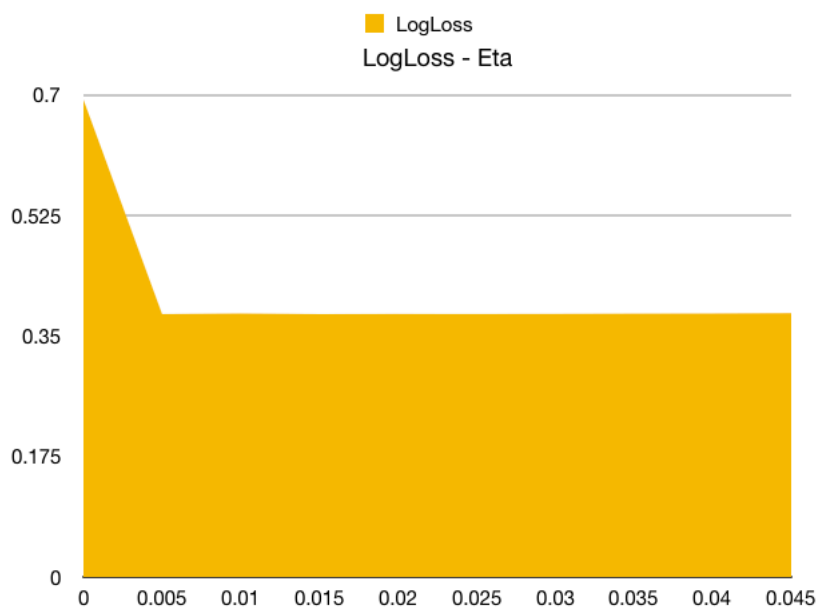


我们可以判断出，当 eta 大于 0.05 时，LogLoss 的值几乎没有太大变化；但在 $[0, 0.05)$ 无法得出结论。因此我们还需要在 $[0, 0.05)$ 区间再做一次实验。

下表为 eta 和 LogLoss 在 $[0, 0.05)$ 区间的关系：

Eta	LogLoss
0	0.69314
0.005	0.38198
0.01	0.38283
0.015	0.38194
0.02	0.38217
0.025	0.38206
0.03	0.38224
0.035	0.38261
0.04	0.38283
0.045	0.38335

可视化的结果：



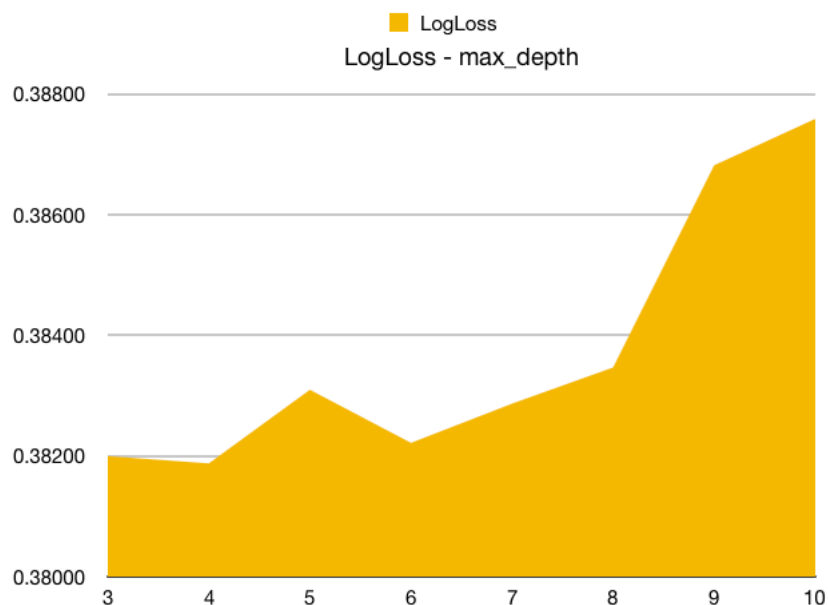
我们可以看到，eta 只要不取 0，对 LogLoss 几乎没有太大影响。但 eta 取较大值（大于 0.05）时，LogLoss 呈略微上升趋势。我们取这些 eta 中 LogLoss 最小的，即 $\eta = 0.015$ 。

max_depth（最大深度） 我们将 eta 取上一步骤中的最佳 $\eta = 0.015$ 。最大深度深度的取值范围是从 $[1, \infty)$ ，默认值为 6。我们对 3, 4, 5, 6, 7, 8, 9, 10 这 8 个整数都试一次，然后交到 kaggle 上测试 log loss 值。

下表为 max_depth 和 LogLoss 在 $[3, 10]$ 区间的关系：

max_depth	LogLoss
3	0.38200
4	0.38188
5	0.38310
6	0.38222
7	0.38287
8	0.38347
9	0.38682
10	0.38759

可视化的结果：



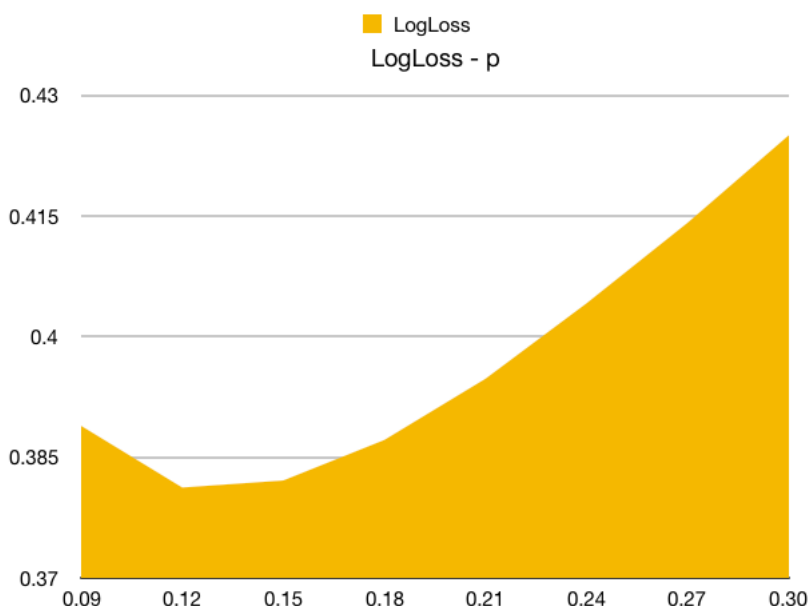
可以看到，树并不是越深越好。max_depth 值取 4 时，LogLoss 达到最小。

训练集中正负样本的比例 训练集中有 36.92% 的正样本（重复），以及 63.08% 的负样本（不重复）。我们可以适当调整正负样本的比例，使得训练出来的 log loss 值最优。我们设参数 $p : 0.6308$ 为正负样本的比例，并使 p 取值 0.09, 0.12, ..., 0.30 共 8 个取值。固定 $\text{eta} = 0.015$, $\text{max_depth} = 4$ ，然后测试 log loss 值。

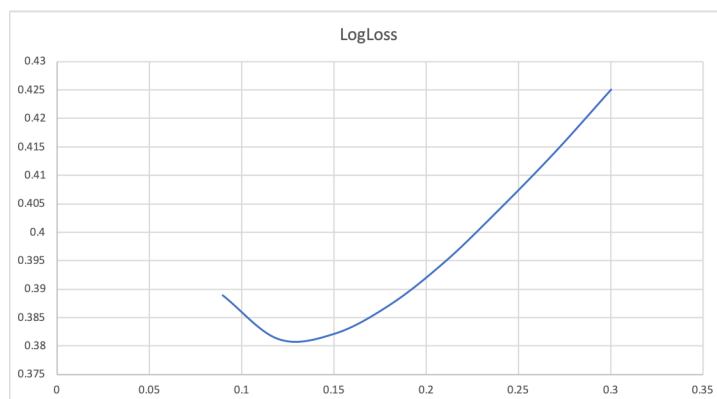
下表为 p 和 LogLoss 在 $[0, 0.05)$ 区间的关系：

p	LogLoss
0.09	0.38893
0.12	0.38126
0.15	0.38214
0.18	0.38718
0.21	0.39479
0.24	0.40417
0.27	0.41419
0.30	0.42506

可视化的结果：



我们可以看到，最小值点应在 $[0.12, 0.15]$ 区间内。我们把折线图改成曲线图，如下图所示。



我们将图片拿到 Photoshop 中放大来看，并用标尺工具测出最小值点约为 $p = 0.1293$ 。

2.2.3 运行结果

使用最优的参数 $\eta = 0.015$ 、 $\text{max_depth} = 4$ 、 $p = 0.1293$ 来训练，在 kaggle 的最终得分为 0.38142，这个得分明显太差。

2.3 LSTM + XGBoost 模型

前一个模型有一个非常严重的问题，即我们只根据词语出现的次数、词语的匹配程度来判断是否重复，而在一定程度上忽略了词语的顺序，也忽略了同义词之间的关系。因此，我们使用 LSTM 来提取特征，与原来的 6 个特征加起来，一起交给 XGBoost 来学习。我们期望 log loss 能进一步地减少。

2.3.1 数据处理

本部分参考了 Kaggle 上的 Kernel (Question Similarity using LSTM Embedding)。

```
1 import numpy as np
2 import pandas as pd
3
4 from sklearn.feature_extraction.text import CountVectorizer
5 import itertools
6
7 import re
8 from keras.preprocessing.sequence import pad_sequences
9 from sklearn.model_selection import train_test_split
10
11 import keras.layers as lyr
12 from keras.models import Model
```

输入数据 .

```
1 df_train = pd.read_csv('../input/train.csv', encoding='utf-8')
2 df_train['id'] = df_train['id'].apply(str)
3
4 df_test = pd.read_csv('../input/test.csv', encoding='utf-8')
5 df_test['test_id'] = df_test['test_id'].apply(str)
6
7 df_all = pd.concat((df_train, df_test))
8 df_all['question1'].fillna('', inplace=True)
9 df_all['question2'].fillna('', inplace=True)
```

创建词库 .

```
1 counts_vectorizer = CountVectorizer(max_features=10000-1).fit(
2     itertools.chain(df_all['question1'], df_all['question2']))
3 other_index = len(counts_vectorizer.vocabulary_)
4 words_tokenizer = re.compile(counts_vectorizer.token_pattern)
```

准备训练数据 .

```
1 def create_padded_seqs(texts, max_len=10):
2     seqs = texts.apply(lambda s:
3         [counts_vectorizer.vocabulary_[w] if w in counts_vectorizer.vocabulary_ else
4         other_index
5         for w in words_tokenizer.findall(s.lower())])
6     return pad_sequences(seqs, maxlen=max_len)
7
8 X1_train, X1_val, X2_train, X2_val, y_train, y_val = \
9     train_test_split(create_padded_seqs(df_all[df_all['id'].notnull()][['question1']])
10 ,
```



```

10         create_padded_seqs(df_all[df_all['id'].notnull()][['question2']])
11         ,
12         df_all[df_all['id'].notnull()][['is_duplicate']].values,
13         stratify=df_all[df_all['id'].notnull()][['is_duplicate']].values,
14         test_size=0.1, random_state=4242)

```

创建模型、训练

```

1 input1_tensor = lyr.Input(X1_train.shape[1:])
2 input2_tensor = lyr.Input(X2_train.shape[1:])
3 words_embedding_layer = lyr.Embedding(X1_train.max() + 1, 100)
4 seq_embedding_layer = lyr.LSTM(64, activation='tanh')
5 seq_embedding = lambda tensor: seq_embedding_layer(words_embedding_layer(tensor))
6 merge_layer = lyr.multiply([seq_embedding(input1_tensor), seq_embedding(
    input2_tensor)])
7 dense1_layer = lyr.Dense(16, activation='sigmoid')(merge_layer)
8 ouput_layer = lyr.Dense(1, activation='sigmoid')(dense1_layer)
9 model = Model([input1_tensor, input2_tensor], ouput_layer)
10 model.compile(loss='binary_crossentropy', optimizer='adam')
11 model.summary()
12
13 model.fit([X1_train, X2_train], y_train,
14         validation_data=([X1_val, X2_val], y_val),
15         batch_size=128, epochs=6, verbose=2)
16
17 features_model = Model([input1_tensor, input2_tensor], merge_layer)
18 features_model.compile(loss='mse', optimizer='adam')

```

用模型预测参数，输出

```

1 X1_train = create_padded_seqs(df_all[df_all['id'].notnull()][['question1']])
2 X2_train = create_padded_seqs(df_all[df_all['id'].notnull()][['question2']])
3 X1_test = create_padded_seqs(df_all[df_all['test_id'].notnull()][['question1']])
4 X2_test = create_padded_seqs(df_all[df_all['test_id'].notnull()][['question2']])
5
6 F_train = features_model.predict([X1_train, X2_train], batch_size=128)
7 F_test = features_model.predict([X1_test, X2_test], batch_size=128)
8
9 name = []
10 for i in range(0, 64):
11     name.append('lstm' + str(i))
12 F_train_df = pd.DataFrame(F_train, columns=name)
13 F_train_df.index.name = 'id'
14 F_train_df.to_csv("../input/lstm_train.csv")
15
16 F_test_df = pd.DataFrame(F_test, columns=name)
17 F_test_df.index.name = 'test_id'
18 F_test_df.to_csv("../input/lstm_test.csv")

```

这样，我们就对每个问题组合求出一个 64 维向量，将它们与上一步骤中的 6 个特征结合起来，一起放进 XGBoost 中训练，我们期望能对结果有一定的提升。

2.3.2 参数的选择

我们使用上一模型中的参数 $\eta = 0.015$ 、 $\max_depth = 4$ 、 $p = 0.1293$ 来训练，在我迭代了 10000 次，看到验证集的 Log Loss 下降到 0.1473 左右，非常惊喜。交上去一看，居然还比之前差，得分 0.39674，手动微笑，手动再见。

```

Run: main x
[9880] train-logloss:0.132415 valid-logloss:0.147405
[9890] train-logloss:0.132396 valid-logloss:0.147397
[9900] train-logloss:0.132373 valid-logloss:0.147386
[9910] train-logloss:0.132352 valid-logloss:0.147378
[9920] train-logloss:0.132328 valid-logloss:0.147367
[9930] train-logloss:0.13231 valid-logloss:0.147357
[9940] train-logloss:0.132284 valid-logloss:0.147346
[9950] train-logloss:0.132258 valid-logloss:0.147336
[9960] train-logloss:0.13224 valid-logloss:0.147328
[9970] train-logloss:0.132219 valid-logloss:0.147319
[9980] train-logloss:0.132201 valid-logloss:0.147312
[9990] train-logloss:0.132171 valid-logloss:0.147297
[9999] train-logloss:0.132147 valid-logloss:0.147286
  
```

Name	Submitted	Wait time	Execution time	Score
simple_xgb_6_lstm.csv	a few seconds ago	1 seconds	11 seconds	0.39296
Complete				

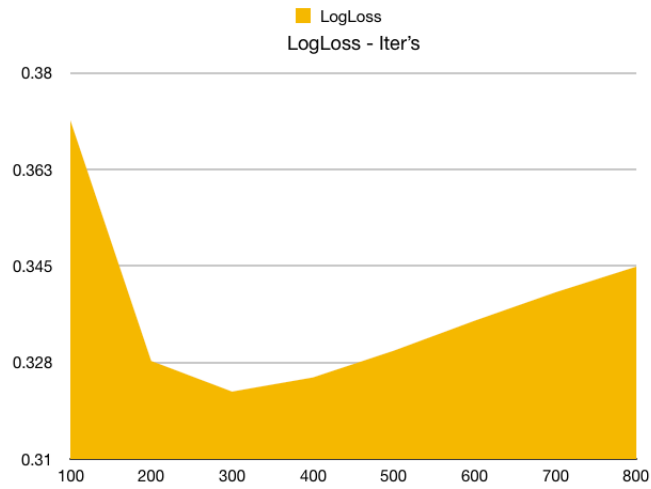
变更模型是需要重新选择参数的，即重新确定新的 η 、 \max_depth 、 p 。我认为 η 、 \max_depth 对结果的影响不大，但正负样本的比例 p 的影响可能比较大。另外就是验证集的 Log Loss 明显偏低，不排除有过拟合的可能。

我们先测试是否为过拟合的问题。先迭代 200 次，验证集的 Log Loss 为 0.257135，交到 kaggle 上的结果是 0.32778。原来真的是过拟合了！

我将 100 次到 800 次每隔 100 次迭代的结果提交到 kaggle 上，Log Loss 关系表如下：

Iter's	LogLoss
100	0.37142
200	0.32781
300	0.32225
400	0.32481
500	0.32966
600	0.33509
700	0.34024
800	0.34490

可视化结果：



可见，最低点在 200 到 400 之间。我们再对这个区间细分，将 200 次到 400 次每隔 10 次的迭代结果提交到 kaggle 上，Log Loss 关系表如下：

Iter's	LogLoss
200	0.32781
210	0.32643
220	0.32542
230	0.32453
240	0.32383
250	0.32327
260	0.32283
270	0.32255
280	0.32233
290	0.32228
300	0.32225
310	0.32227
320	0.32240
330	0.32250
340	0.32267
350	0.32298
360	0.32336
370	0.32365
380	0.32401
390	0.32438
400	0.32481

