

一个呆萌的Query纠错算法

原创 少华 少华幽居 2020-07-16

预想

我所理解的Query纠错是在拥有一个在海量文本上训练的单向语言模型，或者双向的，例如Bert，并在对应任务场景的海量数据集上预训练之后，对于给定的query，以beam search的方式做query生成。但实际情况是，对应任务场景并没有海量的数据集。一种简单直接的方式是把query转成pinyin，并与候选的query集合计算编辑距离，取改写最小的，根据候选的query的热度做排序选最优的。这样做效果不易于保证，其一，query改写往往是多段的，当该query与候选query有很多相似段的时候，这个候选集合就很难以有效的手段评估；其二，这种方式没法做到查漏补缺；其三，这种方式计算量随着热门query的增长而激增。

基于此，有一个新的想法，既可以保证query能够多段改写并且效果更优，又可以保证能够查漏补缺，还可以保证计算量不会随着热门query的增长而激增。最重要的意义是，该方法可以在对应任务场景没有海量数据集的情况下，做冷启动，为后续上模型积累自动标注的数据。

该想法便是用ngram和masked ngram做检索，然后依据ngram拼音的编辑距离来做小段的ngram纠错，通过beam search的方式保存最优的K个结果，这里的K为beam size。

算法简概

step 1

从对应任务场景的线上日志里聚合热门query，同时需要统计出每条query的频次，可以用于后面的决策过程。

step 2

基于热门query构建ngram检索词表和masked ngram检索词表，在本例中， $N=3$ 。例如：

你知道五道口地铁站怎么走吗

加入开始标记和结束标记，本例以 _ 为标记，改写query为：

你知道五道口地铁站怎么走吗

创建ngram检索词表

```
1    _你知
2    你知道
3    知道五
4    道五道
5    五道口
6    道口地
7    口地铁
8    地铁站
9    铁站怎
10   怎么走
11   么走吗
12   走吗_
```

创建masked ngram检索词表

```
1    _你知[mask]    =>    _你知
2    你知[mask]     =>    你知道
3    知道[mask]     =>    知道五
4    道五[mask]     =>    道五道
5    五道[mask]     =>    五道口
6    道口[mask]     =>    道口地
7    口地[mask]     =>    口地铁
8    地铁[mask]     =>    地铁站
9    铁站[mask]     =>    铁站怎
10   怎么[mask]     =>    怎么走
11   么走[mask]     =>    么走吗
12   走吗[mask]     =>    走吗_
```

step 3给定query，查询query的ngram是否存在于ngram检索词表，如果存在，则不用纠错，反之，mask掉该ngram最后一个字，查询该ngram是否存在于masked ngram检索词表里，如果不存在，则不做纠错，反之，则找到最多beam size个候选。伪代码如下：

```
1  # ngram纠错，返回候选ngram集合
2  rewrite_ngram(ngrams):
3      candidates = []
4      for ngram in ngrams:
5          if ngram in ngram检索词表:
6              ngram加入候选
7          else:
```

```

8         new = mask掉ngram最后一个字
9         if new in masked_ngram检索词表:
10             计算候选改写的ngram和ngram的编辑距离
11             依据编辑距离排序候选集合
12             编辑距离最小的都加入到候选集合里
13         else:
14             ngram加入候选
15         shuffle(candidates)
16         return candidates[:beam_size]
17
18 # 这里K表示beam size
19 query_correct(query):
20     q = '_' + query + '_'
21     query_cands.add(q[:K-1])
22     # 已完成纠错的query候选
23     complete_cands = set()
24     while query_cands:
25         ngrams = 根据query_cands和query构建ngrams
26         new_ngrams = rewrite_ngram(ngrams)
27         依据new_ngrams和query_cands再构建query_cands
28         并将已完成纠错的query加入complete_cands
29
30     从complete_cands获取最优候选
31     根据该候选的长度, 原query的长度, 以及编辑距离决定是否做改写
32     返回改写后的query

```

例子以及代码

本例随机抽取了QA训练集的100条query

1	吃什么长睫毛	吃什么长睫毛快	1
2	自己做葡萄酒可以用电饭锅发酵吗	做葡萄酒的比例	0
3	再见美人鱼电影歌曲	再见美人鱼电影插曲歌名	1
4	宾馆的wifi怎么连接	酒店的wifi怎么连接??	1
5	...		
6	...		
7	...		
8	湛江去海南多少公里	湛江到海南有多远	1
9	虚拟号段是什么意思	虚拟号码是什么意思?	1

10	墙头马上的孩子叫 笔记本电脑开机怎样进入安全模式	0
11	小米起了虫子还能吃吗 总有一天我要撕去这虚假的星空日语怎么读	0

生成ngram检索词表如下:

1	_吃什	2
2	吃什么	2
3	什么长	2
4	么长睫	2
5	长睫毛	2
6	睫毛_	1
7	睫毛快	1
8	毛快_	1
9	...	
10	...	
11	...	
12	虚假的	1
13	假的星	1
14	的星空	1
15	星空日	1
16	空日语	1
17	日语怎	1
18	语怎么	1
19	么读_	1

生成masked ngram检索词表如下(其中 ^A 和 ^B 是单个字符):

1	_吃^A	_吃了^B2^A_吃的^B1^A_吃什^B2
2	吃什^A	吃什么^B2
3	什么^A	什么颜^B1^A什么穿^B1^A什么关^B1^A什么号^B1^A什么意^B3^A什么品^B1^A什么
4	_ ^B4^A什么?^B4^A什么灯^B1^A什么长^B2^A什么马^B1^A什么是^B1^A什么证^B1^A什么原	
5	^B2^A什么叫^B1^A什么我^B1^A什么牌^B1	
6	么长^A	么长睫^B2
7	长睫^A	长睫毛^B2
8	睫毛^A	睫毛快^B1^A睫毛_^B1
9	毛快^A	毛快_^B1
10	...	
11	...	
12	...	

```
13  撕去^A   撕去这^B1
14  去这^A   去这虚^B1
15  这虚^A   这虚假^B1
16  虚假^A   虚假的^B1
17  假的^A   假的星^B1
18  的星^A   的星空^B1
19  星空^A   星空日^B1
20  空日^A   空日语^B1
    日语^A   日语怎^B1
    语怎^A   语怎么^B1
```

尝试纠错Query 你知道吗：

```
1  query: 你知道吗
2  ngrams: ['_你知']
3  ngram: _你知
4  replace new ngram [_你] failed
5  candidates: ['_你知']
6  ngrams: ['你知道']
7  ngram: 你知道
8  replace new ngram [你知] failed
9  candidates: ['你知道']
10 ngrams: ['知道吗']
11 ngram: 知道吗
12 rewrite [知道吗] to [知道手]
13 candidates: ['知道手']
14 ngrams: ['道手_']
15 ngram: 道手_
16 rewrite [道手_] to [道手机]
17 candidates: ['道手机']
18 ngrams: ['手机_']
19 ngram: 手机_
20 rewrite [手机_] to [手机实]
21 candidates: ['手机实']
22 ngrams: ['机实_']
23 ngram: 机实_
24 rewrite [机实_] to [机实名]
25 candidates: ['机实名']
26 ngrams: ['实名_']
```

```
27 ngram: 实名_  
28 rewrite [实名_] to [实名制]  
29 candidates: ['实名制']  
30 ngrams: ['名制_']  
31 ngram: 名制_  
32 has ngram: 名制_  
33 candidates: ['名制_']  
34 new query: 你知道手机实名制
```

算法最终推导出的改写Query是 [你知道手机实名制](#)，与原Query不等长，并且与原query的编辑距离很大，会在最终的决策中过滤掉，所以最终对该Query不会做Query改写。

尝试纠错Query [宾馆的Wafi怎么链接](#)：

```
1 query: 宾馆的Wafi怎么链接  
2 ngrams: ['_宾馆']  
3 ngram: _宾馆  
4 has ngram: _宾馆  
5 candidates: ['_宾馆']  
6 ngrams: ['宾馆的']  
7 ngram: 宾馆的  
8 has ngram: 宾馆的  
9 candidates: ['宾馆的']  
10 ngrams: ['馆的W']  
11 ngram: 馆的W  
12 rewrite [馆的W] to [馆的w]  
13 candidates: ['馆的w']  
14 ngrams: ['的wa']  
15 ngram: 的wa  
16 rewrite [的wa] to [的wi]  
17 candidates: ['的wi']  
18 ngrams: ['wif']  
19 ngram: wif  
20 has ngram: wif  
21 candidates: ['wif']  
22 ngrams: ['ifi']  
23 ngram: ifi  
24 has ngram: ifi  
25 candidates: ['ifi']  
26 ngrams: ['fi怎']
```

```
27 ngram: fi怎
28 has ngram: fi怎
29 candidates: ['fi怎']
30 ngrams: ['i怎么']
31 ngram: i怎么
32 has ngram: i怎么
33 candidates: ['i怎么']
34 ngrams: ['怎么链']
35 ngram: 怎么链
36 rewrite [怎么链] to [怎么连]
37 candidates: ['怎么连']
38 ngrams: ['么连接']
39 ngram: 么连接
40 has ngram: 么连接
41 candidates: ['么连接']
42 ngrams: ['连接_']
43 ngram: 连接_
44 has ngram: 连接_
45 candidates: ['连接_']
46 new query: 宾馆的wifi怎么连接
```

该纠错结果非常**NICE**!!!!!!!

算法代码见https://github.com/thelostpeace/query_correction

喜欢此内容的人还喜欢

太全面了！汽轮机结构及运行控制原理，建议收藏
机电要参

外面吃它人均至少300+，在家做，50元任吃！

菜菜美食日记