

6 【中文分词系列】7. 深度学习分词？只需一个词典！

Mar By 苏剑林 | 2017-03-06 | 47661位读者 引用

这个系列慢慢写到第7篇，基本上也把分词的各种模型理清楚了，除了一些细微的调整（比如最后的分类器换成CRF）外，剩下的就看怎么玩了。基本上来说，要速度，就用基于词典的分词，要较好地解决组合歧义何和新词识别，则用复杂模型，比如之前介绍的LSTM、FCN都可以。但问题是，用深度学习训练分词器，需要标注语料，这费时费力，仅有的公开的几个标注语料，又不可能赶得上时效，比如，几乎没有哪几个公开的分词系统能够正确切分出“扫描二维码，关注微信号”来。

本文就是做了这样的一个实验，**仅用一个词典，就完成了深度学习分词器的训练，居然效果还不错！**这种方案可以称得上是半监督的，甚至是无监督的。

随机组合就可以

做法很简单，既然深度学习需要语料，那我就自己生成一批语料。怎么生成？我把词典中的词随机组合就行了。对不对，随机组合生成的不是自然语言呀？我开始也怀疑这一点，但实验之后发现，这样做出来的效果特别好，甚至有胜于标注语料的结果的现象。

事不宜迟，我们来动手。首先得准备一个带词频的词表，**词频是必须的，如果没有词频，则效果会大打折扣。**然后写一个函数，以正比于词频的概率，随机挑选词表中的词语，组合成“句子”。

```
1 import numpy as np
2 import pandas as pd
3
4 class Random_Choice:
5     def __init__(self, elements, weights):
6         d = pd.DataFrame(zip(elements, weights))
7         self.elements, self.weights = [], []
8         for i,j in d.groupby(1):
9             self.weights.append(len(j)*i)
10            self.elements.append(tuple(j[0]))
11            self.weights = np.cumsum(self.weights).astype(np.float64)/sum(self.weights)
12     def choice(self):
13         r = np.random.random()
14         w = self.elements[np.where(self.weights >= r)[0][0]]
15         return w[np.random.randint(0, len(w))]
```

注意，这里按了权重（weights）来分组，从而实现了随机抽。**随机抽的速度取决于分组的数目**，所以，词典的词频最好可以做一些预处理，使得它有“扎堆”的现象，即可以把出现10001,10002,10003次的词语都可以取整为10000，把出现12001,12002,12003,12004次的词语都取整为12000，类似这样的预处理。这步是相当重要的，因为如果有GPU的话，后面模型的**训练速度瓶颈就是在这里了**。

接着，统计字表，然后写生成器，这些都是很常规的，用的方法还是4tag字标注法，然后添加x标签代表填充标注，不清楚的读者，可以回头阅读LSTM分词的文章：

```
1 import pickle
2 words = pd.read_csv('dict.txt', delimiter='\t', header=None, encoding='utf-8')
3 words[0] = words[0].apply(unicode)
```

```

4 words = words.set_index(0)[1]
5
6 try:
7     char2id = pickle.load(open('char2id.dic'))
8 except:
9     from collections import defaultdict
10    print u'fail to load old char2id.'
11    char2id = pd.Series(list(''.join(words.index))).value_counts()
12    char2id[:] = range(1, len(char2id)+1)
13    char2id = defaultdict(int, char2id.to_dict())
14    pickle.dump(char2id, open('char2id.dic', 'w'))
15
16 word_size = 128
17 maxlen = 48
18 batch_size = 1024
19
20 def word2tag(s):
21     if len(s) == 1:
22         return 's'
23     elif len(s) >= 2:
24         return 'b'+ 'm'*(len(s)-2)+'e'
25
26 tag2id = {'s':[1,0,0,0,0], 'b':[0,1,0,0,0], 'm':[0,0,1,0,0], 'e':[0,0,0,1,0]}
27
28 def data_generator():
29     wc = Random_Choice(words.index, words)
30     x, y = [], []
31     while True:
32         n = np.random.randint(1, 17)
33         seq = [wc.choice() for i in range(n)]
34         tag = ''.join([word2tag(i) for i in seq])
35         seq = [char2id[i] for i in ''.join(seq)]
36         if len(seq) > maxlen:
37             continue
38         else:
39             seq = seq + [0]*(maxlen-len(seq))
40             tag = [tag2id[i] for i in tag]
41             tag = tag + [[0,0,0,0,1]]*(maxlen-len(tag))
42             x.append(seq)
43             y.append(tag)
44         if len(x) == batch_size:
45             yield np.array(x), np.array(y)
46             x, y = [], []

```

还是以前的模型

模型的话，用回以前写过的LSTM或者CNN都行，我这里用了LSTM，结果表明LSTM具有很强的记忆能力。

```

1 # Keras 2.0 + Tensorflow 1.0 运行通过
2
3 from keras.layers import Dense, Embedding, LSTM, TimeDistributed, Input, Bidirectional
4 from keras.models import Model
5
6 sequence = Input(shape=(maxlen,), dtype='int32')

```

```

7 embedded = Embedding(len(char2id)+1, word_size, input_length=maxlen, mask_zero=True)(sequence)
8 blstm = Bidirectional(LSTM(64, return_sequences=True))(embedded)
9 output = TimeDistributed(Dense(5, activation='softmax'))(blstm)
10 model = Model(inputs=sequence, outputs=output)
11 model.compile(loss='categorical_crossentropy', optimizer='adam')
12
13 try:
14     model.load_weights('model.weights')
15 except:
16     print u'fail to load old weights.'
17
18 for i in range(100):
19     print i
20     model.fit_generator(data_generator(), steps_per_epoch=100, epochs=10)
21     model.save_weights('model.weights')

```

用我的GTX1060，结合我的词典（**50万不同词语**），每轮大概要花70s，这里每10轮保存一次模型，而range(100)这个100是随便写的，因为每隔10轮就保存一次，所以读者可以随时看效果中断程序。

关于准确率的问题，**要注意的是因为使用了mask_zero=True，所以x标签被忽略了，训练过程但最后显示的训练准确率，又把x标签算进去了，所以最后现实的准确率只能不到0.3，大概就是0.28左右。**这不重要的，我们训练完成后，再测试即可。

最后是一个经验：**字向量维度越大，对长词的识别效果一般越好。**

结合动态规划输出

到这里，就结合一下viterbi算法，通过动态规划来输出最后的结果。动态规划能够保证输出最优的结果，但是会降低效率。而直接输出分类器预测的最大结果，也能够得到类似的结果（但理论上有可能出现bbbb这样的标注结果）。这个看情况而定吧，这里毕竟是实验，所以就使用了viterbi，如果**真的在生产环境中，为了追求速度，应该不用为好（放弃一点精度，大幅提高速度）。**

```

1 zy = {'be':0.5,
2       'bm':0.5,
3       'eb':0.5,
4       'es':0.5,
5       'me':0.5,
6       'mm':0.5,
7       'sb':0.5,
8       'ss':0.5
9       }
10
11 zy = {i:np.log(zy[i]) for i in zy.keys()}
12
13 def viterbi(nodes):
14     paths = {'b':nodes[0]['b'], 's':nodes[0]['s']}
15     for l in range(1,len(nodes)):
16         paths_ = paths.copy()
17         paths = {}
18         for i in nodes[l].keys():
19             nows = {}

```

```

20         for j in paths_.keys():
21             if j[-1]+i in zy.keys():
22                 nows[j+i]= paths_[j]+nodes[1][i]+zy[j[-1]+i]
23             k = np.argmax(nows.values())
24             paths[nows.keys()[k]] = nows.values()[k]
25     return paths.keys()[np.argmax(paths.values())]
26
27 def simple_cut(s):
28     if s:
29         s = s[:maxlen]
30         r = model.predict(np.array([[char2id[i] for i in s]+[0]*(maxlen-len(s))]), vert
31         r = np.log(r)
32         nodes = [dict(zip(['s','b','m','e'], i[:4])) for i in r]
33         t = viterbi(nodes)
34         words = []
35         for i in range(len(s)):
36             if t[i] in ['s', 'b']:
37                 words.append(s[i])
38             else:
39                 words[-1] += s[i]
40         return words
41     else:
42         return []
43
44 import re
45 not_cuts = re.compile(u'([\da-zA-Z ]+)|[。、？！\.\?;!""']')
46 def cut_word(s):
47     result = []
48     j = 0
49     for i in not_cuts.finditer(s):
50         result.extend(simple_cut(s[j:i.start()]))
51         result.append(s[i.start():i.end()])
52         j = i.end()
53     result.extend(simple_cut(s[j:]))
54     return result

```

这里的代码跟之前一样，基本没怎么改。效率并不是很高，不过都说了，这里是实验，如果有兴趣用到生产环境的朋友，自己想办法优化吧。

来测试一下

结合我自己整理的词典，最终的模型在backoff2005的评测集上达到85%左右的准确率（backoff2005提供的score脚本算出的准确率），这个准确率取决于你的词典。

看上去很糟糕？这不重要。首先，我们并没有用它的训练集，纯粹使用词典无监督的训练，这个准确率，已经很让人满意了。其次，准确率看上去低，但实际情况更好，因为这不过是语料的分词标准问题而已，比如

- 1、评测集的标准答案，将“古老的中华文化”分为“古老/的/中华/文化”，而模型将它分为“古老/的/中华文化”；
- 2、评测集的标准答案，将“在录入上走向中西文求同的道路”分为“在/录入/上/走/向/中/西文/求/同/的/道路”，而模型将它分为“在/录入/上/走向/中西文/求同/的/道路”；

3、评测集的标准答案，将“更是教育学家关心的问题”分为“更/是/教育学/家/关心/的/问题”，而模型将它分为“更是/教育学家/关心/的/问题”；

这些例子随便扫一下，还能举出很多，这说明backoff2005的标注本身也不是特别标准，我们不用太在意这个准确率，而更应该留意到，我们得到的模型，对新词、长词，有着更好的识别效果，而达到这个效果，只需要一个词典，这比标注数据容易多了，这对于一些特定领域（比如医学、旅游等等）定制分词系统，是非常有效的。

下面再给一些测试例子，这些例子基本比之前有监督训练的结果还要好：

罗斯福是第二次世界大战期间同盟国阵营的重要领导人之一。1941年珍珠港事件发生后，罗斯福力主对日本宣战，并引进了价格管制和配给。罗斯福以租借法案使美国转变为“民主国家的兵工厂”，使美国成为同盟国主要的军火供应商和融资者，也使得美国国内产业大幅扩张，实现充分就业。二战后期同盟国逐渐扭转形势后，罗斯福对塑造战后世界秩序发挥了关键作用，其影响力在雅尔塔会议及联合国的成立中尤其明显。后来，在美国协助下，盟军击败德国、意大利和日本。

苏剑林是科学空间的博主

结婚的和尚未结婚的

大肠杆菌是人和许多动物肠道中最主要且数量最多的一种细菌

九寨沟国家级自然保护区位于四川省阿坝藏族羌族自治州南坪县境内，距离成都市400多公里，是一条纵深40余公里的山沟谷地

现代内地人，因莫高窟而知敦煌。敦煌因莫高窟也在近代蜚声海外。但莫高窟始凿于四世纪，到1900年才为世人瞩目，而敦煌早从汉武帝时，即公元前一百多年起，就已是西北名城了

从前对巴特农神庙怎么干，现在对圆明园也怎么干，只是更彻底，更漂亮，以至于荡然无存。我们所有大教堂的财宝加在一起，也许还抵不上东方这座了不起的富丽堂皇的博物馆。那儿不仅仅有艺术珍品，还有大堆的金银制品。丰功伟绩！收获巨大！两个胜利者，一个塞满了腰包，这是看得见的，另一个装满了箱子。

别忘了，这仅仅是用一个词典做出来的，分出来的不少词，比如尤其是人名，都是词典中没有的。这效果，足可让人满意了吧？

想想为什么

回到我们一开始的疑惑，为什么随机组合的文本也能够训练出一个很棒的分词器出来？原因就在于，我们一开始，基于词典的分词，就是做了个假设：句子是由词随机组合起来的。这样，我们分词，就要对字符串进行切分，使得如下概率：

$$p(w_1)p(w_2)\dots p(w_n)$$

最大。最后求解的过程就用到了动态规划。

而这里，我们不外乎也沿用了这种假设——文本是随机组合的——这个假设严格上来说不成立，但一般来说够用了，效果也显示在这里了。可能让人意外的是，这样出来的分词器，居然也能把“结婚的和尚未结婚的”的组合歧义句分出来。其实不难理解，我们随机组合的时候，是按照词频来挑选词语的，这就导致了高频出现更多，低频出现更少，这样经过大量重复操作后，事实上我们是**通过LSTM来学习到了动态规划这个过程**！

这是很惊人的，这表示，我们可以**用LSTM，来学习传统的一些优化算法**！更甚地，通过改进RNN用来解决一些传统cs问题，比如凸包，三角剖分，甚至是TSP，最神奇的地方在于这玩意效果竟然还不错，甚至比一些近似算法效果好。（<https://www.zhihu.com/question/47563637>）注意，诸如TSP之类的问题，是一个NP问题，原则上没有多项式时间的解法，但利用LSTM，甚至我们有可能得到一个线性的有效解法，这对于传统算法领域是多么大的冲击！用神经网络来设计优化算法，最终用优化算法来优化神经网络，实现一个自己优化自己，那才叫智能！

呃～扯远了，总之，效果是唯一标注吧。

已经训练好的模型

最后分享一个已经训练好的模型，有Keras的读者可以下载测试：

~~seg.zip（该权重最新版本的tensorflow+Keras已经不可用，请大家根据上述代码自行训练～）~~

转载到请包括本文地址：<https://kexue.fm/archives/4245>

更详细的转载事宜请参考：《科学空间FAQ》

如果您需要引用本文，请参考：

苏剑林. (2017, Mar 06). 《【中文分词系列】7. 深度学习分词？只需一个词典！》[Blog post]. Retrieved from <https://kexue.fm/archives/4245>