

6 获取并处理中文维基百科语料

Jan By 苏剑林 | 2017-01-06 | 12665位读者 引用

中文语料库中，质量高而又容易获取的语料库，应该就是维基百科的中文语料了，而且维基百科相当厚道，每个月把所有条目都打包一次（下载地址在这里：<https://dumps.wikimedia.org/zhwiki/>），供全世界使用，这才是真正的“取之于民，回馈于民”呀。遗憾的是，由于天朝的无理封锁，中文维基百科的条目到目前只有91万多条，而百度百科、互动百科都有千万条了（英文维基百科也有上千万了）。尽管如此，这并没有阻挡中文维基百科成为几乎是最高质量的中文语料库。（百度百科、互动百科它们只能自己用爬虫爬取，而且不少记录质量相当差，几乎都是互相复制甚至抄袭。）

门槛

尽量下载很容易，但是使用维基百科语料还是有一定门槛的。直接下载下来的维基百科语料是一个带有诸多html和markdown标记的文本压缩包，基本不能直接使用。幸好，已经有热心的高手为我们写好了处理工具，主要有两个：1、[Wikipedia Extractor](#)；2、gensim的wikicorpus库。它们都是基于python的。

然而，这两个主流的处理方法都不能让我满意。首先，Wikipedia Extractor提取出来的结果，会去掉{}标记的内容，这样会导致下面的情形

西方语言中“数学”（；）一词源自于古希腊语的（）

这是因为括号里的词带有{}标记，被清空了；而按照网上的教程，直接用gensim.corpora.wikicorpus.WikiCorpus处理，问题更严重，因为它连所有标点都去掉了。对于追求一份高质量语料库的、具有强迫症的笔者来说，这都是不能接受的。因此，自己动手结合gensim，写了一个处理脚本。

代码

```
1 from gensim.corpora.wikicorpus import extract_pages,filter_wiki
2 import bz2file
3 import re
4 import opencc
5 from tqdm import tqdm
6 import codecs
7
8 wiki = extract_pages(bz2file.open('zhwiki-latest-pages-articles.xml.bz2'))
9
10 def wiki_replace(d):
11     s = d[1]
12     s = re.sub('.*{|\[s\S]*?|}', '', s)
13     s = re.sub('<gallery>[s\S]*?</gallery>', '', s)
14     s = re.sub('(\.){([\^{}\\n]*?|[\^{}\\n]*?)}', '\\1[[\\2]]', s)
15     s = filter_wiki(s)
16     s = re.sub('\\* *\\n\\{2,}', '', s)
17     s = re.sub('\\n+', '\\n', s)
18     s = re.sub('\\n[:;]|\\n +', '\\n', s)
19     s = re.sub('\\n==', '\\n\\n==', s)
20     s = u' [' + d[0] + u'] \\n' + s
```

```
21     return opencv.convert(s).strip()
22
23 i = 0
24 f = codecs.open('wiki.txt', 'w', encoding='utf-8')
25 w = tqdm(wiki, desc=u'已获取0篇文章')
26 for d in w:
27     if not re.findall('[a-zA-Z]+:', d[0]) and d[0] and not re.findall(u'^#', d[
28         s = wiki_replace(d)
29         f.write(s+'\n\n\n')
30         i += 1
31         if i % 100 == 0:
32             w.set_description(u'已获取%s篇文章'%i)
33
34 f.close()
```

注释

可见，代码的主要部分是正则表达式。首先通过bz2file直接不解压来读取下载下来的语料，然后用gensim的extract_pages来提取每个页面，提取后，先处理页面的一些特殊的、非文本的标记，然后将部分有用的{{}}标记替换为[]，因为[]标记不会被完全清空（具体原理读者得自己测试了），然后用gensim的filter_wiki函数直接清理，接着再处理一下换行的问题，最后通过opencv将繁体转化为简体。

后面的循环中，re.findall('[a-zA-Z]+:', d[0])这个条件是去掉那些帮助页面，re.findall(u'^#', d[1])这个条件是去掉重定向的页面，最后得到大概就是91.9万个页面。tqdm是用来显示进度的，这个必须有。程序在我的机器上运行了大概40分钟，得到了1.5G左右的纯文本语料。运行时间不重要，因为预处理是一次性的。

值得注意的是，opencv不能用sudo apt-get install opencv来安装，这个默认版本太低，要用源码编译安装，然后pip install opencv安装python接口，这时候在python中调用opencv可能会报“段错误”，这时候要运行

```
1 | cp /usr/lib/libopencv.so.1.0.0 /usr/lib/x86_64-linux-gnu/
```

副产品

上面提到了重定向，重定向意味着两个词具有同样的意思，这样我把中文维基里边所有的重定向都提取出来了，做了个匹配表。也就是说，词表中每一行的两个词都是相同含义的。这算是一个副产品了。

基于中文维基重定向的同义词表：[wiki_cn_mapping.7z](#)

转载到请包括本文地址：<https://spaces.ac.cn/archives/4176>

如果您需要引用本文，请参考：

苏剑林. (2017, Jan 06). 《获取并处理中文维基百科语料》[Blog post]. Retrieved from <https://spaces.ac.cn/archives/4176>