

词性标注

班级：2016211303 姓名：黄若鹏 学号：2016212901

1.题目要求

Part-of-speech tagging: 30 points

-This data set contains one month of Chinese daily which are segment and POS tagged under Peking Univ. standard.

-Project idea:

- * Design a sequence learning method to predicate a POS tags for each word in sentence.

- * Use 80% data for model training and other 20% for testing (or 5-fold cross validation to test learner's performance. So it could be interesting to separate dataset.)

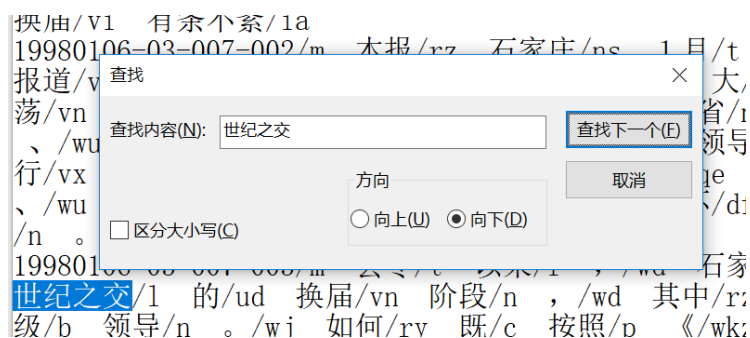
首先我想说这个数据集有太多坑了，导致训练时大部分都在处理数据，实在是耗费了大量时间

总是报各种由于数据集格式不规范导致的错误。数据量有这么大，实在是心累。

为了 debug 采用了 python 中的异常处理结构，这样一步一步 de 由于数据格式不规范的 bug。

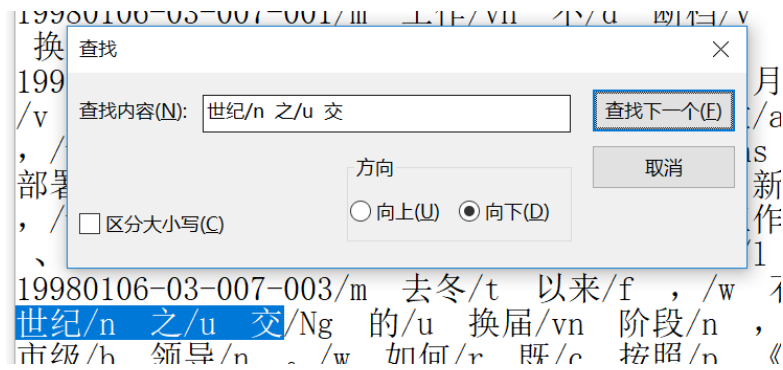
2.处理数据

2.1.这个语料好像还有不同版本，有些细节，很不同



上图中世纪之交是一个词

而请看下图

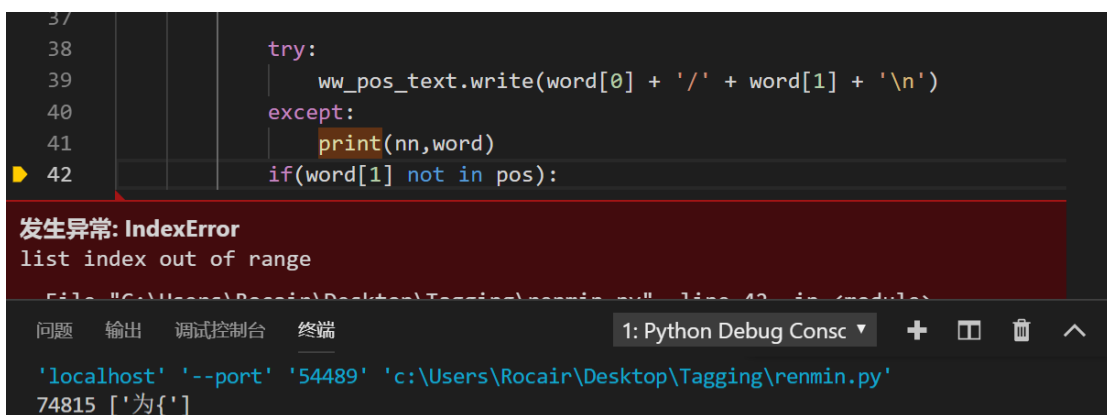


世纪之交却拆成了三个词。

这个问题如果统一语料，不算问题，但是在实验中，我交叉使用了不同的数据，使得，训练过程中，大量词无法正确匹配！因此只能自己编写处理数据的程序。

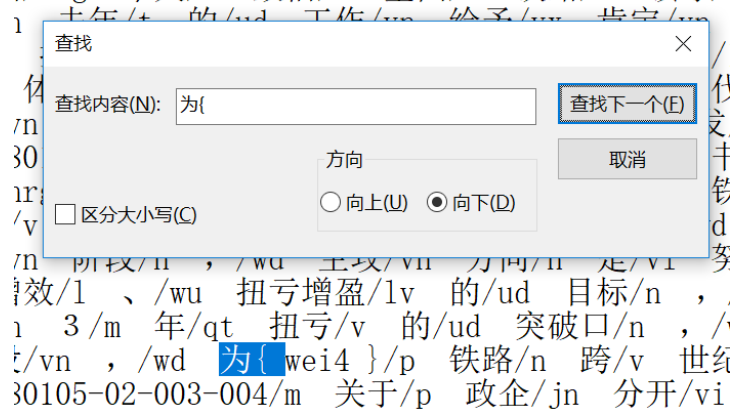
2.2.数据格式不规范

我本以为处理同一数据集就不会遇到啥奇葩问题了，看来我还是太单纯了。



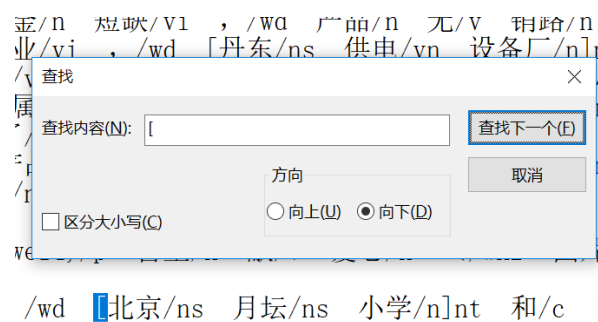
上图是我进行处理数据，然后按相应的格式输出到另一个文件中。方法是先对数据按空格 split 然后在按"/"符号 split 将词和词性分割开，但是却遇到上述问题，一开始，是花费了很长时间检查代码，是不是哪里错了，可是代码逻辑很简单，应该不会出错，面对这种 bug 只能采用异常处理，try except 方式了（C++中是 try catch），打印结果。

后来去检查数据集，奇葩问题出现，在同一个文档中，对于一个词的描述既存在正确形式，又存在错误格式，请看下图



虽然差异不大，左图是正确表达格式，而右图是错误表达格式，“{”后面不应该出现空格，这样会导致 split 函数操作后结果不正确。虽说对于人类来说很好分别，但是对于代码来判断，却要考虑大量逻辑！

2.3 复合词问题

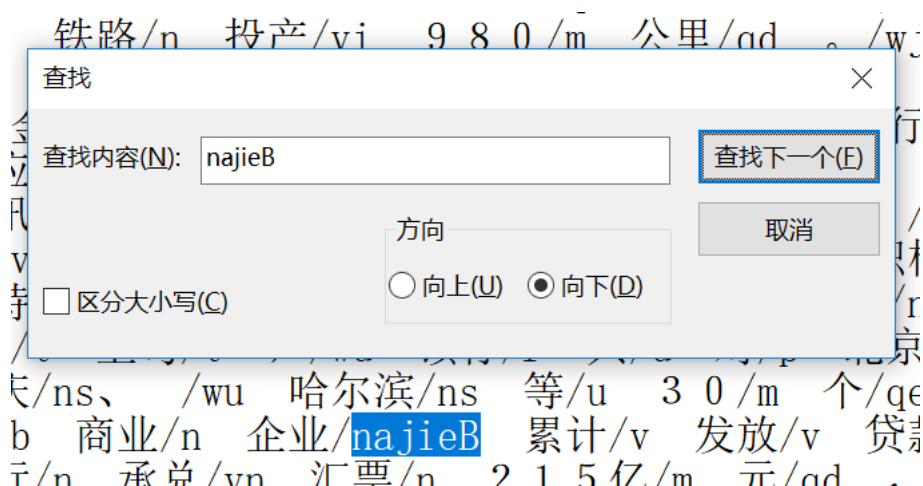


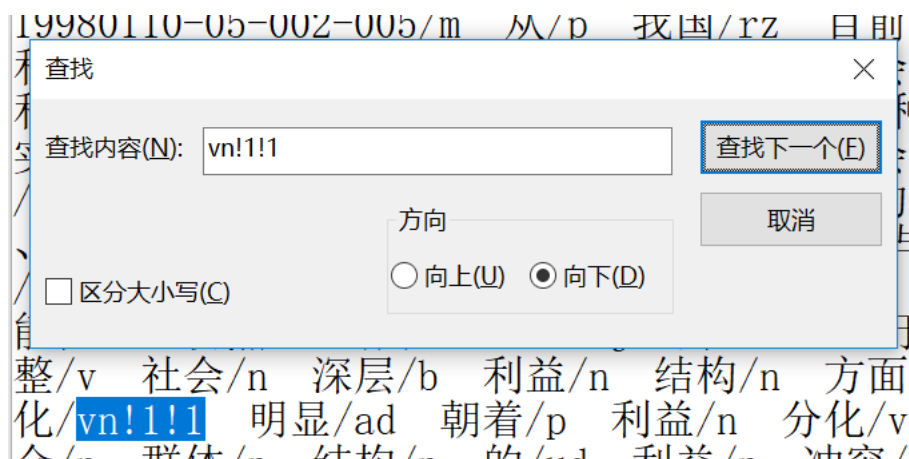
就是还存在这种复合词的情况，这里确实很麻烦。

这里最合理的处理方法应该是，根据复合词，建立一个新词，但是这里，为了简便，不生成新的词汇，只处理这个复合词成为不同的小词，并分别处理

```
word = tmp[i].split('/')
index = word[1].find('"', 0) #处理复合词
if(index!=-1):
    word[1]=word[1][:index]
word[0]=word[0].replace('[', ' ')#处理复合词
```

2.4.其他各种奇葩问题





不知道这些词性是如何创造出来的

总之，做其他题时，可能面对这些奇葩问题，无法感知，因此也就略过了，而做这道题如果就用提供的数据集，要解决一堆这样的问题，恰逢期末考试周，本作业难以做的完美，望见谅。

3.运行环境

系统：windows10

工具：VsCode

4.输入输出

输入：只有一个 1998-01-105-带音.txt 文件，其余文件均有，DataProcessing.py

程序生成

输出：词性标注正确率

5.概要

词性标注 (part-of-speech tagging) ,又称为词类标注或者简称标注, 是指为分词结果中的每个单词标注一个正确的词性的程序, 也即确定每个词是名词、动词、形容词或者其他词性的过程。

我通过上网查询, 发现主要有以下几种统计方法:

- (1) 基于最大熵的词性标注
- (2) 基于统计最大概率输出词性
- (3) 基于 HMM 的词性标注

我是通过 HMM 隐马尔科夫模型来进行词性标注。观测序列即为分词后的语句, 隐藏序列即为经过标注后的词性标注序列。起始概率 发射概率和转移概率和分词中的含义大同小异, 可以通过大规模语料统计得到。观测序列到隐藏序列的计算可以通过 viterbi 算法, 利用统计得到的起始概率 发射概率和转移概率来得到。得到隐藏序列后, 就完成了词性标注过程。

数据处理

写程序时, 我首先对提供的数据进行了清洗和整理, 把时间信息给去掉, 得到总训练语料之后, 我又建立了一个所有词的集合 (所有出现的词, 但不能重复), 和一个所有词性的集合, 这里都是为了 HMM 模型建立那个大表做准备。

建立模型

使用 HMM 隐马尔可夫模型, 进行词性标注。

训练集是这个数据的 80%, 这里我也查了 5-fold cross validation 的概念, 这里由于程序训练比较耗时, 因此在程序中是将整个文件读进去, 为了方便直接取前 80%为训练集, 后 20%为测试集, 也可以自定义训练集和测试集

输入输出

这次实验使用的就是老师提供的新闻报纸语料

输入训练集是标注好词性的文件的前 80%

测试集是没有标注词性的文件后 20%

输出为, 标注好词性的测试集。

检测性能时 可以通过模型输出与元数据集做比较, 统计标注正确性。

HMM 的模型表示

HMM 由隐含状态 S 、可观测状态 O 、初始状态概率矩阵 π 、隐含状态概率转移矩阵 A 、可观测值转移矩阵 B (混淆矩阵) 组成。

π 和 A 决定了状态序列, B 决定了观测序列, 因此, HMM 可以由三元符号表示:

$$\lambda = (A, B, \pi)$$

HMM 的两个性质：

1. 齐次假设：

$$P(i_t | i_{t-1}, o_{t-1}, i_{t-2}, o_{t-2} \cdots i_1, o_1) = P(i_t | i_{t-1})$$

2. 观测独立性假设：

$$P(o_t | i_T, o_T, i_{T-1}, o_{T-1} \cdots i_1, o_1) = P(o_t | i_t)$$

viterbi 用于词性标注

词性标注问题映射到隐马模型可以表述为：模型中状态(词性)的数目为词性符号的个数 N ；从每个状态可能输出的不同符号(单词)的数目为词汇的个数 M 。假设在统计意义上每个词性的概率分布只与上一个词的词性有关(即词性的二元语法)，而每个单词的概率分布只与其词性相关。那么，我们就可以通过对已分词并做了词性标注的训练语料进行统计，统计出 HMM 的参数 λ ，当然这就是上述学习问题。

然后可以根据已知的词语，通过 viterbi 算法，求出每个词语对应的词性，即完成词性标注。

6.部分数据

6.1.词性表

v
n
ud
a
wp
t
wkz
m
qe
wky
nt
wu
nrf
nrg
wd
wj
Vg
k|
wm
p

vi
f
rr
dc
ui
vn
ns
c
rz
s
wt
vl
df
qt
d
ad
wyz
jn
wyy
lv
ul

qv
mq
j
r
vx
an
vu
u
b
l
jb
i
ia
in
vq
iv
wf
q
vd
ue

nr
Tg
uz
qd
y
nx
ry
qz
ld
lb
qc
ww
uo
jv
h
Ag
ib
qb
o
e
ws
Dg
qr
rzw

6.2.所有词汇的“集合”，所有词不能重复

迈向 充满 希望 的 新 世 纪 —— 一九九八年 新年 讲话 （ 附 图 片 1 张 ） 中 共 中 央 总 书 记 、 国 家 主 席 江 泽 民	一九九七年 十二月 三十一日 12月 31日 ， 发表 1998年 《 》 。新华社 记者 兰 红 光 摄 同 胞 们 朋 友 女 士 先 生	中央 民 播 台 国 际 电 视 台 向 全 国 各 族 各 地 香 港 特 别 行 政 区 澳 门 台 湾 海 外 胞 界 各 国 以 挚 候	1997年 是 发 展 历 史 上 非 常 重 要 很 不 平 凡 一 年 决 心 继 承 邓 小 平 同 志 遗 志 继 续 把 建 设 有 特 色 社 会 主 义 事 业
---	---	---	--

oooooooo

6.3.训练集

迈向/v 充满/v 希望/n 的/ud 新/a 世纪/n ——/wp 一九九八年/t 新年/t 讲话/n (/wkz 附/v 图片/n 1/m 张/qe)/wky 中共中央/nt 总书记/n、/wu 国家/n 主席/n 江/nrf 泽民/nrg (/wkz 一九九七年/t 十二月/t 三十一日/t)/wky 12月/t 31日/t，/wd 中共中央/nt 总书记/n、/wu 国家/n 主席/n 江/nrf 泽民/nrg 发表/v 1998年/t 新年/t 讲话/n 《/wkz 迈向/v 充满/v 希望/n 的/ud 新/a 世纪/n 》/wky。/wj (/wkz 新华社/nt 记者/n 兰/nrf 红光/nrg 摄/Vg)/wky 同胞/n 们/k、/wu 朋友/n 们/k、/wu 女士/n 们/k、/wu 先生/n 们/k：/wm 在/p 1998年/t 来临/vi 之际/f，/wd 我/rr 十分/dc 高兴/a 地/ui 通过/p 中央/n 人民/n 广播/vn 电台/n、/wu 中国/ns 国际/n 广播/vn 电台/n 和/c 中央/n 电视台/n，/wd 向/p 全国/n 各族/rz 人民/n，/wd 向/p 香港/ns 特别/a 行政区/n 同胞/n、/wu 澳门/ns 和/c 台湾/ns 同胞/n、/wu 海外/s 侨胞/n，/wd 向/p 世界/n 各国/rz 的/ud 朋友/n 们/k，/wd 致以/v 诚挚/a 的/ud 问候/vn 和/c 良好/a 的/ud 祝愿/vn！/wt 1997年/t，/wd 是/vl 中国/ns 发展/vn 历史/n 上/f 非常/dc 重要/a 的/ud 很/dc 不/df 平凡/a 的/ud 一/m 年/qt。/wj 中国/ns 人民/n 决心/d 继承/v 邓/nrf 小平/nrg 同志/n 的/ud 遗志/n，/wd 继续/v 把/p 建设/v 有/v 中国/ns 特色/n 社会主义/n 事业/n 推向/v 前进/vi。/wj 中国/ns 政府/n 顺利/ad 恢复/v 对/p 香港/ns 行使/v 主权/n，/wd 并/c 按照/p “/wyz 一国两制/jn”/wyy、/wu “/wyz 港人治港/lv”/wyy、/wu 高度/d 自治/vi 的/ud 方针/n 保持/v 香港/ns 的/ud 繁荣/vn 稳定/vn。/wj 中国/ns 共产党/n 成功/a 地/ui 召开/v 了/ul 第十五/m 次/qv 全国/n 代表大会/n，/wd 高举/v 邓小平理论/n 伟大/a 旗帜/n，/wd 总结/v 百年/mq 历史/n，/wd 展望/v 新/a 的/ud 世纪/n，/wd 制定/v 了/ul 中国/ns 跨/v 世纪/n 发展/v 的

这里，是整理数据过后的结果，但是绝对这个整理过后的数据依然不完美，肯定存在一堆标记有问题的地方。

6.4.测试集

交通银行坚持改革，加强管理，努力防范金融风险，1997年继续以良好的经营业绩令同行们刮目相看。去年该行全年完成利润47.84亿元，剔除政策性因素，按同口径计算，利润为63.59亿元，比上年增长18.57%，名列银行同业盈利前茅。1997年交通银行认真贯彻国家经济金融政策，自觉服从金融监管，全面推进内部管理制度，明显加快结构调整步伐，各项业务发展迅速。1997年各项人民币存款余额2700亿元，比1996年增长20%；各项人民币贷款余额1800多亿元，比上年增长18%。全行完成80亿元凭证式国债承销计划，超额完成420亿元现金回笼计划。本报讯著名家用电器生产商瑞典伊莱克斯公司于一月十四日在海口向来自全国的六百多名代理商宣布：推出九八冰箱新系列『新静界』系列。其最突出的特点是低噪音设计。（马艳）本报讯深圳先科电子有限公司努力提高产品质量，该公司生产的VCD视盘机最近通过电子工业部优等品检测。为(wei4)促进销售，先科公司最近推出一元行动，凡购买先科产品的消费者，只要再加一元钱，即可得到麦克风一支或卡拉OK碟一套。（唐纯）本报讯辽河油田在稳定石油、天然气生产的同时，以国家利益为重，克服资金紧张等困难因素，及时申报缴纳税款，去年共实现增值税、消费税十三点七亿元，比上年增长一点六八亿元。（夏才源 李兴华）本报讯一九九七年，青岛市国税局深入推行税收征管改革，完善构建了税收收入全过程监控体系，确保了税收足额、均衡入库。全年实现工商各税七十六点九六亿元，比去年同期增长百分之二十三点四一，增收十五点六亿元。（张梦谦）本报讯茂名石化三十万吨乙烯装置投料试车第二年成功地生产三十万吨乙烯，实现了一九九七年度产量达标，而且投产后装置

测试集是不含任何标注，词语和符号文件。

7. 结果分析

通过多次测试，在抽取的测试集上，平均准确率为 0.80622

```
训练结束！
正确率为：0.810321
PS C:\Users\Rocair\Desktop\NLProc>
```

```
训练结束！
正确率为：0.803652
PS C:\Users\Rocair\Desktop\NLProc>
```

8. 部分代码

```
#计算概率矩阵
line = 0

while(True):
    text = fin.readline()
    if(text == '\n'):
        continue
    if(text == ""): #结束
        break
    tmp = text.split() #按空格分割
    len_temp = len(tmp) #语料长度
    line += 1
    for i in range(0, len_temp):
        word = tmp[i].split('/') #按/分割
        pre = tmp[i-1].split('/') #前一个词按/分割]
        #print(word[1])
        fre[word[1]] += 1 #统计该词出现频率
        if(i == 1):
            pi[word[1]] += 1 #先验概率
        elif(i > 0):
            A[pre[1]][word[1]] += 1 #状态转移矩阵概率
            B[word[1]][word[0]] += 1 #观察概率矩阵
            B1[word[1]][word[0]] += 1
    fin.close()
```

```

for i in pos:

    for j in pos:#避免概率 0
        A[i][j] += 1
        # print(A[i][j])
    for j in ww:
        B[i][j] += 1
        # print(B[i][j])

for i in pos:
    pi[i] = pi[i]*1.0/line
    for j in pos:
        A[i][j] = A[i][j]*1.0/(fre[i])
    for j in ww:
        B[i][j] = B[i][j]*1.0/(fre[i])
        # print(B[i][j])
print("训练结束！ ")

```

#维特比算法

```

max_end = ""
for i in range(0, num-1):
    max = 0

    for each in res[text[i]]:
        for each1 in res[text[i+1]]:
            if(A[each][each1] > max):
                max = A[each][each1]
                max_res[i] = each
            if(i == num-2):
                max_end = each1

max_res[num-1] = max_end

```

#结果分析

```

tag_result = open('tag_result.txt','r')#标注结果
reference = open('ww_pos.txt', 'r')#训练集

a=tag_result.read().split()[start:end]
b=reference.read().split()[start:end]
correctnum=0,num=0
for t1,t2 in zip(a,b):
    num+=1
    if t1.split('/')[1]==t2.split('/')[1]:

```

```
correctnum+=1
```

```
print("正确率为:",correctnum/num)
```