

Homework2 实验报告

姓名:	程倩楠	学号:	201834858
实验课程:	Data Mining	实验名称:	NBC
指导老师:	尹建华	完成时间:	2018.11.16

一、实验目的

理解朴素贝叶斯的基本原理，能够编程实现朴素贝叶斯分类器。

二、实验环境

python3.6

三、实验内容

实现朴素贝叶斯分类器，测试在 20 Newsgroups 数据集上的效果。

四、实验过程

（一）数据预处理

在 Homework1 中借助 Lucene 工具用 java 程序实现新闻文档的预处理，这里借助 nltk 模块用 python 程序实现预处理，处理步骤如下：

1. 词条化（利用非字母的字符对文本进行分割）
2. 去除停用词
3. 小写化
4. 词干化（Porter Stemmer）
5. 统计出每个单词的集合频率（Collection Frequency, CF），即单词在所有新闻文档中出现的次数。去除文档中 $CF < 5$ 的所有低频单词。

（二）划分数据

划分数据依靠函数 `dataSeg(fold, rightCateFile, trainDataProp=0.8)` 来实现，函数的主要功能为划分出一定比例的测试数据和训练数据，并将标注的测试数据记录在特定文件中。

参数	解释
fold	将第几份数据作为测试数据 (用于 n-Fold Cross Validation)
rightCateFile	记录标注测试数据的文件

	<newsName_newsClass newsClass>
trainDataProp	划分为训练数据的比例
<pre> 15 def dataSeg(fold, rightCateFile, trainDataProp=0.8): 16 fw=open(rightCateFile, 'w') 17 srcDir='used_news' 18 srcClassList=os.listdir(srcDir) 19 for i in range(len(srcClassList)): 20 srcClassDir=srcDir+'/'+srcClassList[i] 21 srcNewsList=os.listdir(srcClassDir) 22 m=len(srcNewsList) #新闻类中所包含的新闻文档的数目 23 testBeginIndex=fold*(m*(1-trainDataProp)) #测试数据的起始索引 24 testEndIndex=(fold+1)*(m*(1-trainDataProp)) #测试数据的结束索引 25 for j in range(m): #遍历新闻类下的所有新闻文档 26 if (j>=testBeginIndex) and (j<testEndIndex): 27 #记录标注: 文档id(文档名_所属类) 所属类 28 fw.write('%s %s\n' % (srcNewsList[j]+'_'+srcClassList[i], \ 29 srcClassList[i])) 30 targetClassDir='TestData'+str(fold)+'/'+srcClassList[i] 31 else: 32 targetClassDir='TrainData'+str(fold)+'/'+srcClassList[i] 33 if os.path.exists(targetClassDir)==False: 34 os.makedirs(targetClassDir) 35 nf=open(targetClassDir+'/'+srcNewsList[j], 'w') 36 lineList=open(srcClassDir+'/'+srcNewsList[j], 'rb').readlines() 37 for line in lineList: 38 line=line.decode('utf-8', 'ignore') 39 nf.write('%s\n' % line.strip('\n')) 40 nf.close() 41 fw.close() </pre>	
<p>由于20news-18828中含有重名文档,这里利用newsName_newsClass作为文档标识符</p>	

(三) 朴素贝叶斯分类器

朴素贝叶斯的实现主要包括以下几个函数:

1. 函数 trainNB
<p>训练朴素贝叶斯分类器: 对训练样本数据进行统计, 得到类 cate 下单词 word 出现的次数、每个新闻类中单词总数、训练数据中不重复的单词总数。</p>
<pre> 15 def trainNB(trainDir): 16 cateWordCount={} 17 cateWordNum={} 18 vocab=set() 19 classList=os.listdir(trainDir) 20 for i in range(len(classList)): 21 count=0 #记录每个新闻类中单词总数 22 classDir=trainDir+'/'+classList[i] #类目录 23 newList=os.listdir(classDir) 24 for j in range(len(newList)): 25 newsDir=classDir+'/'+newList[j] 26 lines=open(newsDir, 'rb').readlines() 27 for line in lines: </pre>

```

28         line=line.decode('utf-8','ignore')
29         count+=1
30         word=line.strip('\n')
31         vocab.add(word) #记录训练数据中不重复词汇
32         key=classList[i]+'_'+word
33         cateWordCount[key]=cateWordCount.get(key,0)+1 #记录每个类中每个单词
34         cateWordNum[classList[i]]=count
35     vocabNum=len(vocab)
36     return cateWordCount,cateWordNum,vocabNum

```

2. 函数 calCateProb

计算测试文档属于某个类的概率：

$$p(\text{cate}|\text{doc})=p(w_1,w_2,\dots|\text{cate})*p(\text{cate})$$

多项式模型+平滑技术+取对数（防止下溢）：

$$p(\text{word}|\text{cate})=(\text{类 cate 下单词 word 出现的次数}+1)/(\text{类 cate 下单词总数}+\text{训练数据中不重复的单词总数})$$

$$p(\text{cate})=\text{类 cate 下单词总数}/\text{训练数据中单词总数}$$

$$p(\text{cate}|\text{doc})=\sum_{i=1}^n \log(p(w_i|\text{cate}))+\log(p(\text{cate}))$$

```

42 def calCateProb(k,testNewsWords,cateWordCount,cateWordNum,totalNum,vocabNum):
43     prob=0
44     wordNumInCate=cateWordNum[k] #新闻类k中单词总数
45     for i in range(len(testNewsWords)):
46         key=k+'_'+testNewsWords[i]
47         if key in cateWordCount:
48             wordCountInCate=cateWordCount[key]
49         else:
50             wordCountInCate=0.0
51         xcProb=np.log((wordCountInCate+1)/(wordNumInCate+vocabNum))
52         prob=prob+xcProb
53     res=prob+np.log(wordNumInCate)-np.log(totalNum)
54     return res

```

3. 函数 classifyNB

朴素贝叶斯分类器通过比较 $p(\text{cate}_i|\text{doc})$, $i=1,2,\dots,20$ 对每个测试文档进行分类，将分类结果记录在特定文件中。

```

56 #朴素贝叶斯对测试文档进行分类
57 def classifyNB(trainDir,testDir,resultCateFile):
58     fw=open(resultCateFile,'w')
59     #训练分类器
60     cateWordCount,cateWordNum,vocabNum=trainNB(trainDir)
61     #得到训练数据单词总数
62     totalNum=sum(cateWordNum.values())
63     #对测试文档做分类
64     testClassList=os.listdir(testDir)
65     for i in range(len(testClassList)):
66         testClassDir=testDir+'/'+testClassList[i]
67         testNewsList=os.listdir(testClassDir)

```

```

68     for j in range(len(testNewsList)):
69         testNewsWords=[] #测试文档的单词列表
70         testNewsDir=testClassDir+'/'+testNewsList[j]
71         lines=open(testNewsDir,'rb').readlines()
72         for line in lines:
73             line=line.decode('utf-8','ignore')
74             word=line.strip('\n')
75             testNewsWords.append(word)
76         maxP=0.0
77         trainClassList=os.listdir(trainDir)
78         for k in range(len(trainClassList)):
79             p=calCateProb(trainClassList[k],testNewsWords,cateWordCount, \
80                           cateWordNum,totalNum,vocabNum)
81             if k==0:
82                 maxP=p
83                 bestCate=trainClassList[k]
84                 continue
85             if p>maxP:
86                 maxP=p
87                 bestCate=trainClassList[k]
88             fw.write('%s %s\n' % (testNewsList[j]+'_'+testClassList[i], \
89                                   bestCate))
90     fw.close()

```

4. 函数 errorRate

依据标注测试样本文件与朴素贝叶斯分类结果文件，计算机朴素贝叶斯分类器的错误率。

```

93 def errorRate(rightCateFile,resultCateFile):
94     rightCateDict={}
95     resultCateDict={}
96     errorCount=0.0
97     for line in open(rightCateFile,'rb').readlines():
98         line=line.decode('utf-8','ignore')
99         (newsID,cate)=line.strip('\n').split()
100        rightCateDict[newsID]=cate
101    for line in open(resultCateFile,'rb').readlines():
102        line=line.decode('utf-8','ignore')
103        (newsID,cate)=line.strip('\n').split()
104        resultCateDict[newsID]=cate
105    for key in rightCateDict.keys():
106        #输出分类结果
107        print('新闻ID: '+key)
108        print('朴素贝叶斯分类: '+resultCateDict[key])
109        print('新闻真实所属类: '+rightCateDict[key])
110        if rightCateDict[key]!=resultCateDict[key]:
111            errorCount+=1.0
112    errorRate=errorCount/len(rightCateDict)
113    print('error rate: %f' % (errorRate))
114    return errorRate

```

(四) 5 折交叉验证

将数据分为五折，选择其中一折作为测试数据，其余作为训练数据计算错误率，取五次错误率的平均值，作为最终朴素贝叶斯分类器的错误率。包括以下几个过程：

1. 生成五次实验的测试数据、训练数据、标注文件；

- 朴素贝斯分类器对五次实验的测试数据进行分类；
- 计算并记录五次实验的错误率；
- 绘制条形图：可视化每次实验的错误率；
- 计算五次实验的平均错误率；

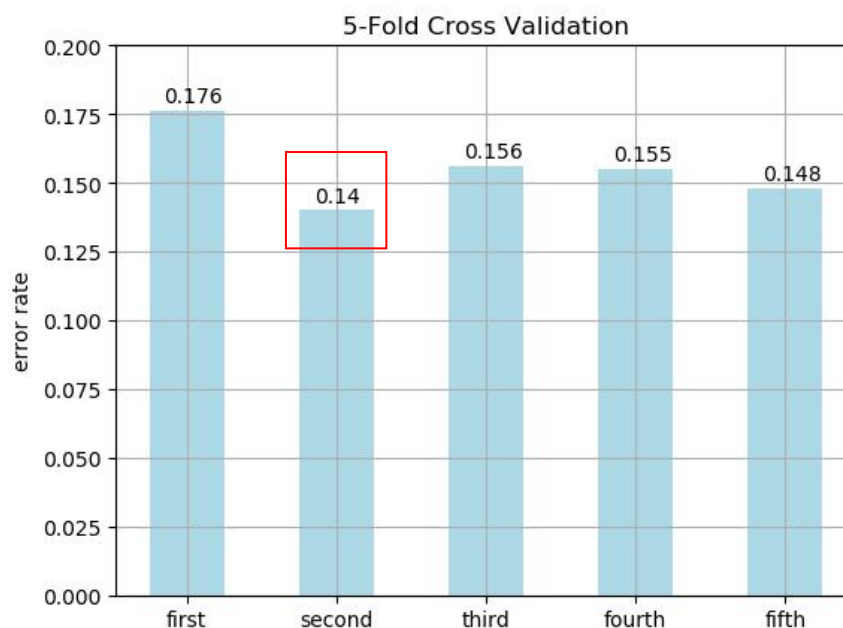
五、实验结果

（一）随机划分数据（20%测试数据、80%训练数据）

```
新闻ID: 83528_talk.religion.misc  
朴素贝叶斯分类: talk.religion.misc  
新闻真实所属类: talk.religion.misc  
新闻ID: 83529_talk.religion.misc  
朴素贝叶斯分类: soc.religion.christian  
新闻真实所属类: talk.religion.misc  
新闻ID: 83535_talk.religion.misc  
朴素贝叶斯分类: talk.religion.misc  
新闻真实所属类: talk.religion.misc  
新闻ID: 83544_talk.religion.misc  
朴素贝叶斯分类: talk.religion.misc  
新闻真实所属类: talk.religion.misc  
新闻ID: 83547_talk.religion.misc  
朴素贝叶斯分类: soc.religion.christian  
新闻真实所属类: talk.religion.misc  
新闻ID: 83558_talk.religion.misc  
朴素贝叶斯分类: talk.politics.mideast  
新闻真实所属类: talk.religion.misc  
error rate: 0.175504  
Naive Bayes Finished!
```

（二）五折交叉验证

- 条形图可视化五次实验的错误率



2. 输出五次实验的错误率和平均错误率

```
errorRate0 = 0.175504  
errorRate1 = 0.140239  
errorRate2 = 0.156358  
errorRate3 = 0.155113  
errorRate4 = 0.147912  
averageErrorRate = 0.155025  
5-Fold Cross Validation Finished!
```