

# 文本检索大作业

2000011714 王骥一 化学与分子工程学院

1.分词。读入数据集，对数据集中的文章进行去标点、分词、去除停用词及低频词，英文数据集考虑词性还原。

**1. 数据处理 (20')。需完成：**读入数据集，对数据集中的文章进行去标点、分词、去除停用词及低频词，英文数据集可以考虑提取词根等。构建词典（下述检索词及相似词都只需在词典范围内完成），保存为vocab.txt。（该部分可直接调包）

```
1 English_f=pd.read_csv('all_news.csv')
2
3
4 stoplist = ['very', 'ourselves', 'am', 'doesn', 'through', 'me', 'against', 'up', 'just', 'her', 'ours',
5             'couldn', 'because', 'is', 'isn', 'it', 'only', 'in', 'such', 'too', 'mustn', 'under', 'their',
6             'if', 'to', 'my', 'himself', 'after', 'why', 'while', 'can', 'each', 'itself', 'his', 'all', 'once',
7             'herself', 'more', 'our', 'they', 'hasn', 'on', 'ma', 'them', 'its', 'where', 'did', 'll', 'you',
8             'didn', 'nor', 'as', 'now', 'before', 'those', 'yours', 'from', 'who', 'was', 'm', 'been', 'will',
9             'into', 'same', 'how', 'some', 'of', 'out', 'with', 's', 'being', 't', 'mightn', 'she', 'again', 'be',
10            'by', 'shan', 'have', 'yourselves', 'needn', 'and', 'are', 'o', 'these', 'further', 'most', 'yourself',
11            'having', 'aren', 'here', 'he', 'were', 'but', 'this', 'myself', 'own', 'we', 'so', 'i', 'does', 'both',
12            'when', 'between', 'd', 'had', 'the', 'y', 'has', 'down', 'off', 'than', 'haven', 'whom', 'wouldn',
13            'should', 've', 'over', 'themselves', 'few', 'then', 'hadn', 'what', 'until', 'won', 'no', 'about',
14            'any', 'that', 'for', 'shouldn', 'don', 'do', 'there', 'doing', 'an', 'or', 'ain', 'hers', 'wasn',
15            'weren', 'above', 'a', 'at', 'your', 'theirs', 'below', 'other', 'not', 're', 'him', 'during', 'which']
16
```

```
#分词
lemmatizer=nltk.stem.WordNetLemmatizer()
def wordcount(news):
    freq=[]
    news=re.sub(r'\W+|\d+', ' ', news)
    for word in news.split():
        if word not in stoplist:
            freq.append(lemmatizer.lemmatize(word))
    return ' '.join(freq)
English_f['word']=English_f['body'].apply(wordcount)
words=English_f['word'].str.split(' ', expand=True).stack().rename('words').reset_index()
words['level_0']=words['level_0']+1
words.columns=['id', 'level_1', 'words']
print(words)
```

这一部分首先利用正则表达式去除了除字母之外的所有字符（包括数字和标点符号）然后利用网上的停词表（后来了解到可以直接调用 nltk 包）去除停用词，使用 lemmatizer 进行词性的还原，分词后整理、重命名列得到每篇文章的分词结果

## 2.构建 TF-IDF 矩阵

```
32 #计算TF
33 def counttf(x):
34     l=len(x)
35     p=x['words'].value_counts()
36     d={'words':p.index,
37       'TF':p.values/l}
38     p_2d=pd.DataFrame(data=d)
39
40     r=pd.merge(x, p_2d, on='words')
41     return r
42
43 TF_data=words.groupby('id').apply(counttf)
44 #统计完频率之后 去除重复词
45 TF_data.drop_duplicates(subset=['id', 'words'], inplace=True)
46
47 #计算IDF
48 doc_num=len(English_f)
49 def countidf(x):
50     df_t=len(x['id'].value_counts())
51     return np.log(doc_num/df_t)
52
53 IDF_data=words.groupby('words').apply(countidf).reset_index()
54 IDF_data.columns=['words', 'IDF']
55 tf_idf=pd.merge(TF_data, IDF_data)
56 tf_idf['TF-IDF']=tf_idf['TF']*tf_idf['IDF']
57 tf_idf.to_csv('tfidf.csv')
```

```
1 #转化为词向量
2 print(tf_idf['words'])
3 a=list(set(tf_idf['words']))
4 k=len(a)
5 word_list=dict(zip(a, list(range(k))))
6 wordlst=pd.Series(word_list)
7
8 def wordtovec(x):
9     vec=np.zeros(k,)
10    for i in range(x.shape[0]):
11        idx=word_list[x.iloc[i]['words']]
12        vec[idx]=x.iloc[i]['TF-IDF']
13    return vec
14 newdata=tf_idf.groupby('id').apply(wordtovec)
15 print(newdata)
16
17 #降维
18 new_data=np.vstack(newdata.values)
19
20 model=PCA(n_components=300)
21 dev=model.fit_transform(new_data)
22 print(dev)
```

对每篇文章计算其 TF-IDF，并生成一个词典（这里生成了一个 key-value 对，以便之后进行查询），之后利用这个生成了 29060 维的词向量，并降到 300 维。

### 3.文章相似度计算以扩展搜索结果，词语相似度计算以扩展搜索词

```
1 #计算各个文章之间的相似度
2 for i in range(dev.shape[0]):
3     dev[i]/=np.linalg.norm(dev[i])
4 cos_sim=dev@dev.T
5
6 print(cos_sim.shape)
7 cos_sim=cos_sim-np.eye(2225)
8 print(cos_sim)
```

```
1 #利用tf-idf矩阵，计算词之间的相似度
2 new_data=np.load('tf_idf.npy')
3 matrix=new_data.T@new_data
4
5 x=matrix.shape[0]
6 idxar=np.argsort(matrix,axis=1)[:,-4:-1]#取相似度排序前三个为模糊词
7 synonym=word_list
8 for word in word_list:
9     idxlst=idxar[word_list[word]]
10    lst=[]
11    for k in idxlst:
12        lst.append(a[k]) #a是前面提到的词典，它和word_list互为补充，使查询时间缩短到O(1)
13    synonym[word]=lst
```

### 4.降维效果评价

```
1 import numpy as np
2 from sklearn.cluster import KMeans
3 import pandas as pd
4
5 English_f=pd.read_csv('all_news.csv')
6 dev=np.load('dev.npy')
7 dict1=dict(zip(set(English_f['topic']),list(range(5))))
8 def replace(x):
9     return dict1[x]
10 y=English_f['topic'].apply(replace)
11
```

```
1 oper=KMeans(n_clusters=5)
2 oper.fit(dev)
3 newX=oper.predict(dev)
4 df=pd.DataFrame({'origin_class':y,'new_class':newX})
5 def maxcount(x):
6     return x.value_counts().iloc[0]
7 ml=df.groupby('new_class').apply(maxcount)
8 purity=sum(ml)/len(y)
9 print(purity)
```

0.8179775280898877

这里重新开了一个 notebook 进行评价，purity 为 0.82，还是比较高的，说明降到 300 维没有损失太多信息

## 5.服务器搭建

```
1 class LocalServer(object):
2     def __init__(self, host, port):
3         self.address = (host, port)
4
48    def run(self):
49        server = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
50        server.setsockopt(socket.SOL_SOCKET, socket.SO_REUSEADDR, 1)
51        server.bind(self.address)
52        i=0
53        server.listen(5)
54        while True:
55            conn, add=server.accept()
56            print(f'Connected to {add}')
57
58            words=json.loads(conn.recv(1024).decode('utf-8')) #接收传参
59            if 'break_000' in words: #此处认为break_000为结束符
60                break
61
62            print(words)
63            p=Thread(target=self.lookupfor, args=(words, conn)) #多线程实现并发处理
64            p.start()
65            p.join()
66            server.close()
```

使用 socket 套接字搭建服务器端, 接收到客户端的 words 列表后, 将其和连接名称 conn 一同传入 lookupfor 函数进行搜索, lookupfor 函数负责向客户端返回搜索结果。这里使用 thread 进行多线程的并发处理 (使用 Process 多进程处理失败了, 不知道为什么, 可能是拷贝连接 conn 的时候发生了占线问题)

在提交的作业中, 多个 gui 是为了测试并发处理 (但似乎并不能很好地测试)

当输入的 words 中包含 break000 时, 断开连接, 终止服务器运行

## 6.服务器端检索文本并返回结果

```
def lookupfor(self, words, conn):
    lemmatizer=nlk.stem.WordNetLemmatizer()
    for i in range(len(words)):
        words[i]=lemmatizer.lemmatize(words[i]) #使用lemmatizer进行词性还原
        words+=synonym[words[i]] #将模糊词添加到列表中
    global topics
    global titles
    k=len(words)
    l=len(topics)
    returnlist=[]
    #首先检索主题和题目，这里认为二者重要性相同，如果又出现了一次，认为其重要性更高，排在前面
    for j in range(k):
        for i in range(l):
            if topics[i]==words[j]:
                if i in returnlist:
                    returnlist[returnlist.index[i]],returnlist[0]=returnlist[0],returnlist[returnlist.index[i]]
                else:
                    returnlist.append(i)
            if titles[i]==words[j]:
                if i in returnlist:
                    returnlist[returnlist.index[i]],returnlist[0]=returnlist[0],returnlist[returnlist.index[i]]
                else:
                    returnlist.append(i)
```

lookupfor 函数首先对传入的词进行词形还原和模糊词扩充，接着开始检索，将检索到的文章 id 放入 returnlist

首先检索主题和题目，将含有关键词的主题和题目放入 returnlist，如果 returnlist 中已经出现该文章，则将该文章放到最前面（一个粗糙的排序方法）

```
returnlist.append(i)
#其次检索与已有文章相似度最高的一篇文章，依次放入returnlist
for num in returnlist:
    idx=np.argmax(cos_sim[num])
    if idx in returnlist:
        returnlist[returnlist.index[idx]],returnlist[0]=returnlist[0],returnlist[returnlist.index[idx]]
    else:
        returnlist.append(idx)
#最后按关键字进行查询，按TF-IDF由高到低依次输出
for j in range(k):
    if words[j] in word_list:
        df=tf_idf[tf_idf['words']==words[j]]
        for i in df['id']:
            if i not in returnlist:
                returnlist.append(i-1)
articles=[]
for k in returnlist:
    articles.append(tuple(English_f.iloc[k,0:2]))

conn.send(json.dumps(articles).encode('utf-8')) #将结果转化为json字符串后返回客户端
return articles
```

接着将和已有文章相似度最高的文章放入 returnlist，与前相同，如果已经出现，则将该文章放到最前面

最后按关键字进行查询，按 TF-IDF 由高到低依次输出，不过这里不再改变原有文章的次序（最后一个重要性比较低）

接着将文章 id 转换为 article，打包发送给客户端