

# 现代程序设计第2周作业

谢奕飞 20377077

利用python数据结构（list, dict, set等）完成简单的文本分析任务。弹幕是现下视频网站，尤其是短视频网站提供的关键功能之一。以B站为例，其有着特殊的弹幕文化，且在视频的不同部分往往会有不同话题的弹幕：比如在视频开头会出现“来啦”“x小时前”“第一!”;在up主暗示一键三连之后常常会出现“下次一定”或者“你币有了”;和up主建立默契之后，观众可以判断视频是否有恰饭，往往在广告之前会出现“要素察觉”“恰饭”“快跑”等等。因此，弹幕经常被作为测度用户（viewer）与视频作者（up主）之间交互行为的关键数据。本次作业提供的数据来自B站某知名up主，已上传至课程资料的数据目录下，数据格式说明如下。

a. 弹幕文件：danmuku.csv，为2799000 rows × 3 columns，本次作业仅使用第一列，即弹幕的文本内容。

b. 停用词表示例，stopwords\_list.txt

请大家尝试完成以下数据分析任务：

1. 使用danmuku.csv，其中一个弹幕可以视为一个文档（document），读入文档并分词（可以使用jieba或pyltp）。
2. 过滤停用词（可用stopwords\_list.txt，或自己进一步扩充）并统计词频，输出特定数目的高频词和低频词进行观察。建议将停用词提前加入到jieba等分词工具的自定义词典中，避免停用词未被正确分词。
3. 根据词频进行特征词筛选，如只保留高频词，删除低频词（出现次数少于5之类），并得到特征词组成的特征集。
4. 利用特征集为每一条弹幕生成向量表示，可以是0，1表示（one-hot，即该特征词在弹幕中是否出现）也可以是出现次数的表示（该特征词在弹幕中出现了多少次）。注意，可能出现一些过短的弹幕，建议直接过滤掉。
5. 利用该向量表示，随机找几条弹幕，计算不同弹幕间的语义相似度，可尝试多种方式，如欧几里得距离或者余弦相似度等，并观察距离小的样本对和距离大的样本对是否在语义上确实存在明显的差别。请思考，这种方法有无可能帮助我们找到最有代表性的弹幕？

## 代码

### CaculateWordsFrequency函数

```
1 def CaculateWordsFrequency(file_name, stopfile_name):
2     f = open(file_name)
3     df = pd.read_csv(f, encoding='utf-8', usecols=[0])#读取文件
4     comments = list(df['content'])
5     stop_f = open(stopfile_name, 'r', encoding='utf-8')#导入分词表
6     jieba.load_userdict(stopfile_name)
7     stop_words = set()
8     for line in stop_f:
9         line = line.rstrip('\n')
10        if len(line):
11            stop_words.add(line)
12    stop_f.close()
13    counts = {}#分词+记录频率
14    seg_list = []
15    documents = []
```

```

16     for content in comments:
17         cur_seg_list = jieba.lcut(content)
18         seg_list.append(cur_seg_list)
19         documents.append(' '.join(cur_seg_list))
20         for seg in cur_seg_list:
21             if seg not in stop_words:
22                 counts[seg] = counts.get(seg, 0) + 1
23     counts = dict(sorted(counts.items(), key=lambda dc: dc[1], reverse=True))
24     print('TOP10: {}'.format(list(counts.keys())[:10]))
25     dic_word = {'word':counts.keys(), 'frequency':counts.values()}
26     pd.DataFrame(dic_word).to_csv('word_counts.csv')
27     f.close()
28     stop_f.close()

```

## OutputWordCloud函数

```

1 def OutputWordCloud(counts, n):#绘制词云图
2     wordpng = wordcloud.WordCloud(background_color='white', height=700, width=1000)
3     wordpng.generate(" ".join(list(counts.keys())[:n]))
4     wordpng.to_file("result.png")

```

## OutputMainWord函数

```

1 def OutputMainWord(counts):#打印关键词
2     word_list = [key for key, value in counts.items() if value > 3]
3     word_list = list(set(word_list))
4     word_id = {}
5     for i in range(len(word_list)):
6         word_id[word_list[i]] = i
7     return set(word_list), word_id
8

```

## CaculateEuclid\_distance函数

```

1 def CaculateEuclid_distance(v1, v2):#计算欧氏距离
2     dis=distance.euclidean(np.array(v1), np.array(v2))
3     return dis
4

```

## CaculateCosine\_distance函数

```

1 def CaculateCosine_distance(v1, v2):#计算余弦相似度
2     if np.dot(np.array(v1), np.array(v2)) > 0:
3         return distance.cosine(np.array(v1), np.array(v2))
4     return MAXDIS
5

```

## RandomComment函数

```

1 def RandomComment(seg_list):#随机选取弹幕
2     id = rd.randint(0, len(seg_list) - 1)

```

```

3     random_vector = [0] * (len(main_word) + 5)
4     for word in seg_list[id]:
5         if word in main_word:
6             random_vector[word_id[word]] = 1
7     return comments[id], random_vector
8

```

## 主函数

```

1  dic_words, counts, documents, seg_list, stop_words, comments = CaculateWords
2  OutputWordCloud(counts, 50)
3  main_word, word_id = OutputMainWord(counts)
4  print('词频特征词提取:', len(main_word), end='')
5  for i in range(50):
6      print(rd.choice(list(main_word)), end=',')
7  print()
8  random_list = []
9  for i in range(5):
10     content, vector = RandomComment(seg_list)
11     random_list.append(vector)
12     print('random comment: ', i+1, ': ', content, sep='')
13  euclid_dis = [[CaculateEuclid_distance(c1, c2) for c2 in random_list] for c1 in
14  cos_dis = [[CaculateEuclid_distance(c1, c2) for c2 in random_list] for c1 in
15  print('Euclid_distance:\n', euclid_dis)
16  print('Cosine_distance:\n', cos_dis)

```

## 运行结果

```

d:\Project\Python - VS Code 控制台
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\NO_THA\1\AppData\Local\Temp\jieba.cache
Loading model cost 0.577 seconds.
Prefix dict has been built successfully.
TOP10: ['哈哈哈哈哈', '武汉', '吃', '加油', '藕', '蒜', '好吃', '真的', '萝卜', '热情']
main word: 7676快用,汤头,醋点,闻不戴,复兴,致死量,盘算着,表现,别演,念叨,MMP,书,贵得,手足无措,恰饭,社牛,给你个,非说,穿着,
里加,团结,少说,小纸,和谐,发明,热热闹闹,看母,一圈,咱俩,bbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbbb,带货,点熟,猪脚,蒜汁,很浓,
长居,变猪,江湖义气,苏州,六小时,行,贾玲,互吹,慈祥,茴香,脏脏,宝,衡阳,十年,疼,
random comment: 1: 都尼玛滴滴推翻鸟
random comment: 2: 干与干寻笑吐
random comment: 3: 这才90?
random comment: 4: 我半夜就
random comment: 5: 看起来就好好吃
Euclid distance:
[[0.0, 2.6457513110645907, 2.23606797749979, 2.23606797749979, 2.23606797749979], [2.6457513110645907, 0.0, 2.0, 2.0, 2
.0], [2.23606797749979, 2.0, 0.0, 1.4142135623730951, 1.4142135623730951], [2.23606797749979, 2.0, 1.4142135623730951, 0
.0, 1.4142135623730951], [2.23606797749979, 2.0, 1.4142135623730951, 1.4142135623730951, 0.0]]
Cosine distance:
[[0.0, 1.0, 1.0, 1.0, 1.0], [1.0, 0.0, 1.0, 1.0, 1.0], [1.0, 1.0, 0.0, 1.0, 1.0], [1.0, 1.0, 1.0, 0.0, 1.0], [1.0, 1.0, 1.0, 1.0, 0.0]]
请按任意键继续. . .

```

result.png

word\_counts.csv

A1	B	C
word	frequency	
0	哈哈哈哈哈	533854
1	武汉	293070
2	吃	147886
3	加油	115976
4	藕	94067
5	蒜	91532
6	好吃	90721
7	真的	64410
8	萝卜	63036
9	热情	49344
10	想	41101
11	帮	37316
12	大哥	33496
13	啊啊啊	33019
14	猴头	31158
15	喜欢	30474
16	牛杂	29741
17	考古	29655
18	真	29193
19	笑	29182
20	死	26230
21	户部	26137
22	爱	24257
23	感受	22660
24	天	22312
25	挂	21658
26	街	21805
27	树梢	21004
28	大派	20065
29	可爱	19814