

作业 2: 词云

姓名: 白植权 学号: 221050020

1. 作业简介

本作业读取作业 1 中网络爬虫程序爬取并存储在 CSV 文件中的数据, 对这些数据进行分词处理并生成词云图。其中词云图 (Word Cloud) 是一种常见的可视化工具, 用于显示文本数据。在词云图中, 每个单词的大小表示其在文本中出现的频率。出现频率较高的词语以较大的字体显示, 反之则以较小的字体显示, 这使得用户可以一眼看出文本中的主要主题和关键词。

2. Python 代码

```
import pandas as pd
import jieba
from wordcloud import WordCloud
from stopwords import stop_words
import os

def chinese_jieba(text):
    wordlist_jieba = jieba.cut(text)
    space_wordlist = " ".join(wordlist_jieba)
    return space_wordlist

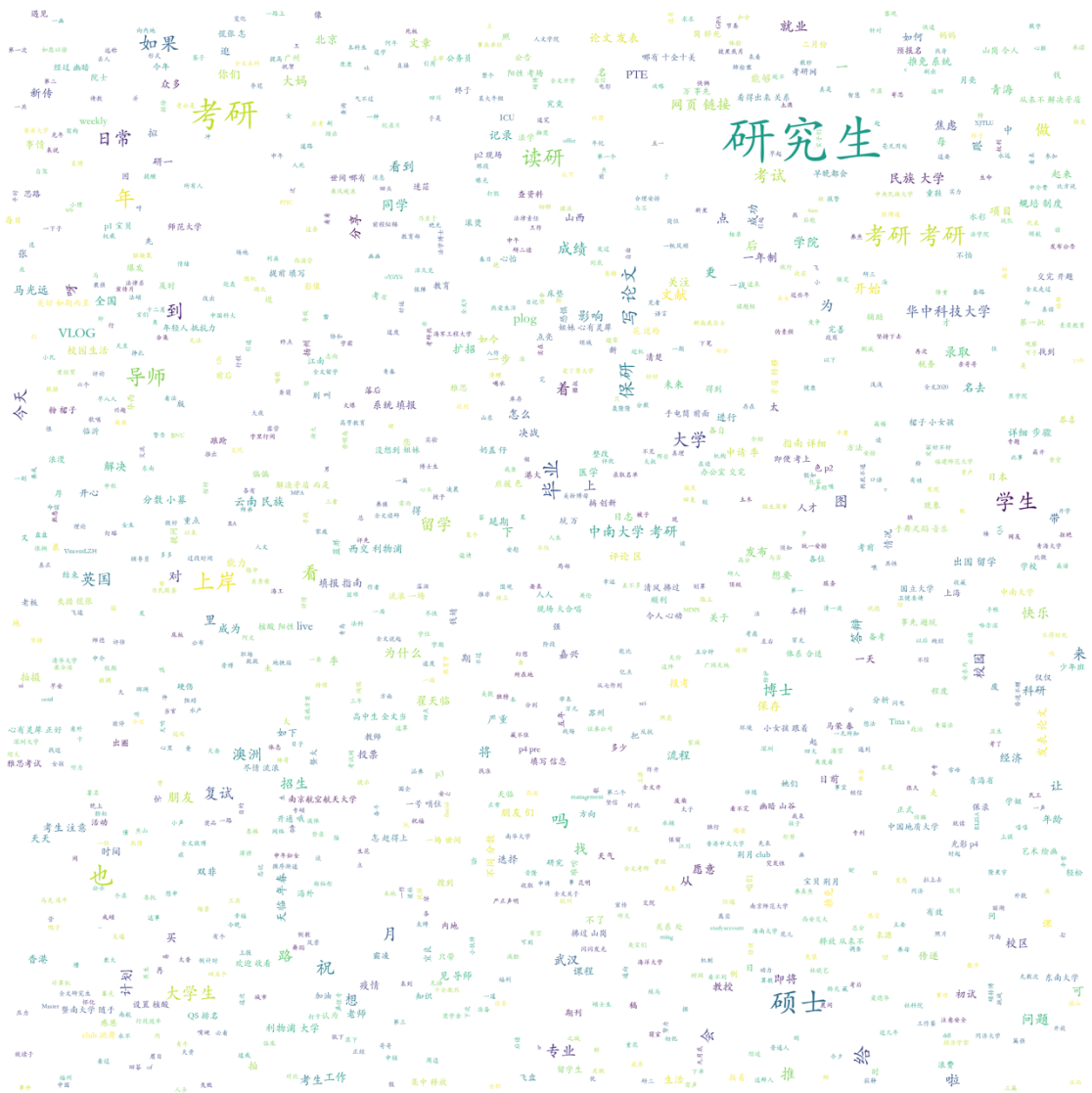
def ciyun():
    topic_list = [f for f in os.listdir('../work1/') if f.endswith('.csv')]
    for topic in topic_list:
        df = pd.read_csv('../work1/' + topic)
        comment_list = df['comment'].values.tolist()
        text = ""
        for jj in range(len(comment_list)):
            text = text + chinese_jieba(comment_list[jj])
        print(text)
        # 调用包 PIL 中的 open 方法, 读取图片文件, 通过 numpy 中的 array 方法生成数组
        wordcloud = WordCloud(width=3000, height=3000,
                               font_path="font.ttf", # 字体文件
                               background_color="white", # 设置背景颜色
                               max_font_size=150, # 设置字体最大值
                               max_words=2000, # 设置最大显示的字数
                               stopwords=stop_words, # 设置停用词, 停用词则不再词云图中表示
                               ).generate(text)
        image = wordcloud.to_image()
        wordcloud.to_file(topic.split('.')[0] + '-ciyun.png') # 导出文件
        image.show()
if __name__ == "__main__":
    ciyun()
```

```

1 # -*- coding: utf-8 -*-
2 # @author : zhang
3 # @time   : 2023/4/3 下午4:14
4 # @function:
5
6 stop_words = {'微博', '视频', '接', '真的', '什么', '的', '在', '是', '我', '你', '他', '她', '它', '和', '与', '去', '把', '请', '都', '啊', '知道', '用'
7             , '呢', '这', '那', '个', '些', '很', '多', '挺', '赞', '出', '自己', '结果', '最近', '不', '一个', '能', '谁', '呢', '吧', '怎么', '有',
8             , '关于', '人', '全文', '就是', '说', '我们', '大家', '被', '所以', '就', '一些', '一定', '一点', '一直', '一样',
9             , '一次', '一起', '一些', '一份', '一下', '一些', '两个', '三个', '几个', '这些', '那些', '这个', '那个', '所有',
10            , '所有的', '任何', '任何人', '任何事情', '任何时候', '任何方式', '任何情况', '别人', '另外', '其他', '其他人', '其事情', '其他时候',
11            , '其他方式', '其他情况', '自己的', '有些', '有时候', '有着', '只有', '只是', '有', '就是说', '好的', '很多', '许多',
12            , '大部分', '绝大部分', '少数', '少部分', '一切', '每个', '每一个', '每件事', '每一时刻', '任何地方', '任何原因',
13            , '每个人', '每一件东西', '每种', '每种情况', '怎么样', '那样', '这样', '那样时', '这样时', '那个', '这个', '那种', '这种',
14            , '哪', '哪一个', '哪种', '哪一种', '这一些', '那一些', '什么样的', '这样的', '一样的', '全部的', '每个人都',
15            , '每件事情', '任人的', '任何事情的', '任何情况的', '有关', '当讲', '之前', '之后', '之内', '之外', '金条', '可以', '可能',
16            , '不能', '实现', '实际上', '现在', '同时', '但是', '然而', '虽然', '尽管', '无论如何', '之间', '而且', '此外',
17            , '即使', '因为', '所以', '由于', '要求', '需要', '推荐', '请求', '建议', '希望', '尽可能', '最好', '通常', '一般',
18            , '有时', '很少', '偶尔', '有时候', '可能会', '有一些', '有点', '好像', '似乎', '不错', '不行', '不要', '不必', '不能', '不想',
19            , '不用', '不必要', '没有', '不是', '不真', '不知', '不会', '可以', '可能', '应该', '应当', '一定', '必须', '需要', '要求', '要是', '要不然', '比如',
20            , '以及', '尤其', '并且', '但是', '然而', '之后', '接着', '而且', '除了', '以外', '依然', '依旧', '仍然', '仍旧', '尽管', '虽然', '但', '还', '这样',
21            , '这时', '这里', '那里', '那么', '那时', '哪里', '只有', '就要', '才能', '才会', '越来越', '已经', '正在', '没有', '就靠', '时候', '自然', '应该', '为了',
22            , '对于', '同时', '一样', '经常', '常常', '大概', '不常', '不时', '不曾', '偶尔', '突然', '马上', '立即', '立刻', '顿时', '每次', '总是',
23            , '往往', '通常', '一般', '常常', '总的来说', '总的说来', '总的来讲', '绝对', '完全', '彻底', '全部', '从不', '从来', '决不', '绝不', '一点', '有点', '稍微', '有些',
24            , '比较', '更加', '非常', '极其', '异常', '相当', '特别', '十分', '细心', '特意', '其实', '事实上', '确实', '实际上', '当然', '无论', '不管', '不论', '任何',
25            , '只要', '就算', '即使', '这个', '它们', '他们', '自己', '的话', '什么的', '等等', '之', '之类', '之类的', '谁的', '什么的', '别的', '或者', '或是'}

```

[illegible]



(3) “杭州”主题词云图：

