



西安交通大学

XI'AN JIAOTONG UNIVERSITY

新闻与新媒体学院

python 编程语言课程报告

共青团中央微博数据分析

姓名：杨忠强

学号：2215310366

2023 年 7 月 23 日

摘要

微博作为快速便捷的社交媒体平台，在网民舆论思潮中发挥着重要的信息引导作用。本次作业分析微博平台共青团中央账号的发微情况。以 2019 年至今其所发布的所有微博作为主要研究对象，统计分析其发微规律与时间构成情况，同时利用 jieba 分词提取所发微博中的高频词汇，分词后通过 LDA 主题模型分析其所发微博的几大主题构成，最后通过 pyecharts 制作可视化展现相应的结果提高交互性。

关键词: 微博, 共青团中央, 数据分析, 可视化



目录

摘要	I
第 1 章 目标简介	1
1.1 数据爬取	1
1.2 数据分析可视化	1
1.3 运行环境	2
第 2 章 数据爬取	3
2.1 爬取对象	3
2.2 爬虫实现	3
2.3 数据采集结果	5
第 3 章 数据预处理	6
3.1 数据去重与格式转换	6
3.2 文本内容清洗	7
第 4 章 数据分析与可视化	8
4.1 逐年发微热力图	8
4.2 发微时间统计	10
4.2.1 逐年统计	10
4.2.2 每年逐月统计	11
4.2.3 每月逐天统计	12
4.2.4 每周逐天统计	13
4.2.5 每天逐时统计	13
4.3 周时热力统计	14
4.4 高频词统计	15
4.5 LDA 主题分析	16
第 5 章 总结展望	19

第 1 章 目标简介

本报告目的是抓取社交媒体平台数据，并进行计算分析，而后对所爬取数据进行分析 and 可视化展示。

1.1 数据爬取

网络爬虫的主要目的是将互联网上的网页下载到本地形成一个或联网内容的镜像备份，通过对内容的解析之后格式化处理进行存储。本报告使用爬虫项目所用框架为 scrapy。

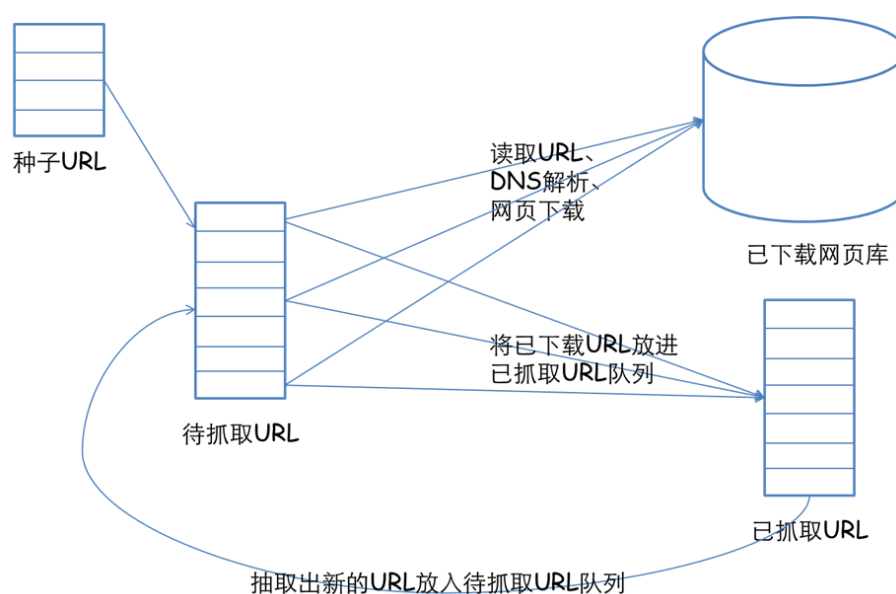


图 1.1: 爬虫示意图

1.2 数据分析可视化

本报告使用 pandas、jieba 等库进行数据处理，具体方案将在后文阐述。可视化工具使用为 pyecharts1.2，可视化数据生成 html 文件后能够较好实现交互。

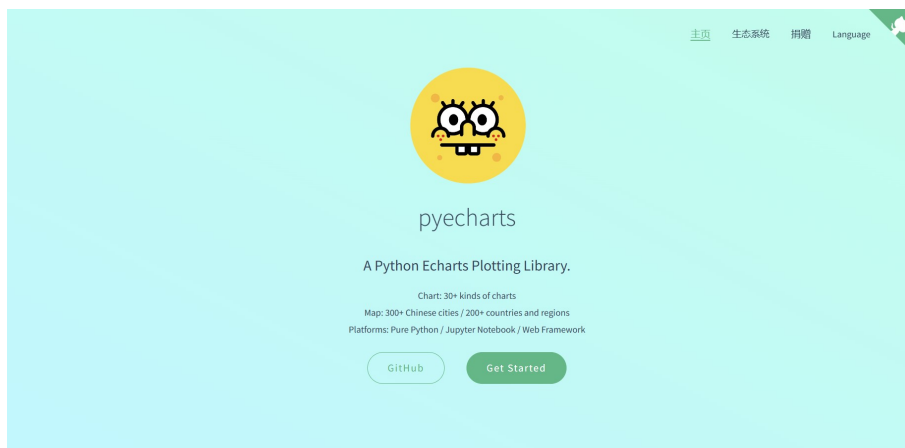


图 1.2: pyecharts

1.3 运行环境

本报告在 pycharm（Community Edition 2022.2.2）的 windows 版本下运行，环境配置为 python3.10，界面如图 1.3 所示。

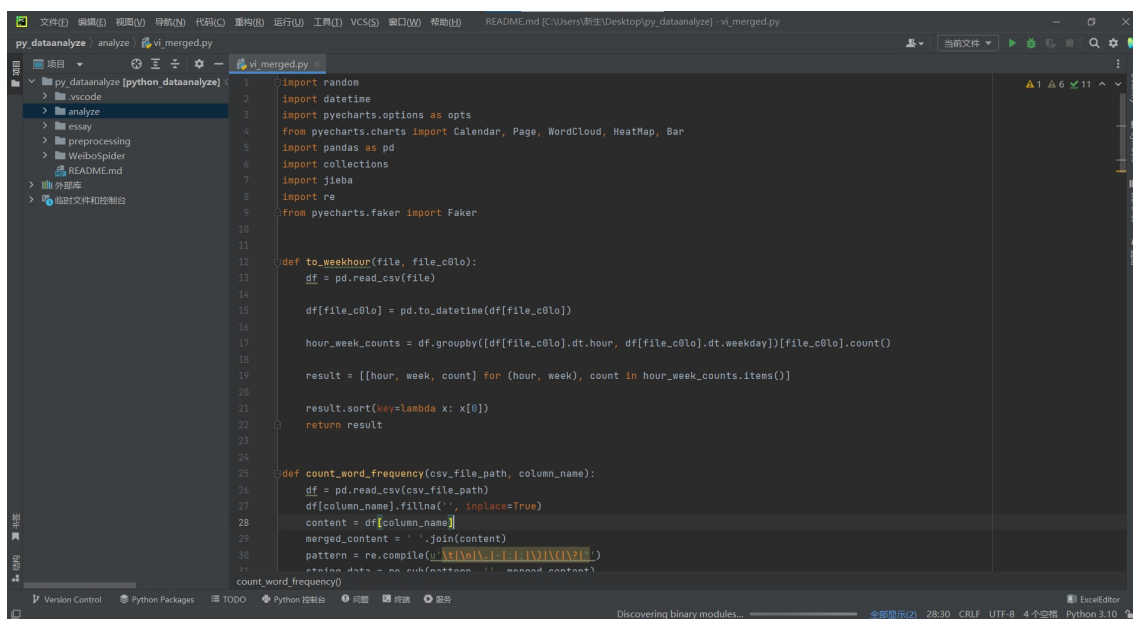


图 1.3: pycharm 运行环境示意图

第 2 章 数据爬取

2.1 爬取对象

本报告选取共青团中央的官方微博如图 2.1，爬取时间范围为 2019 年 1 月 1 日到 2023 年 7 月 16 日。

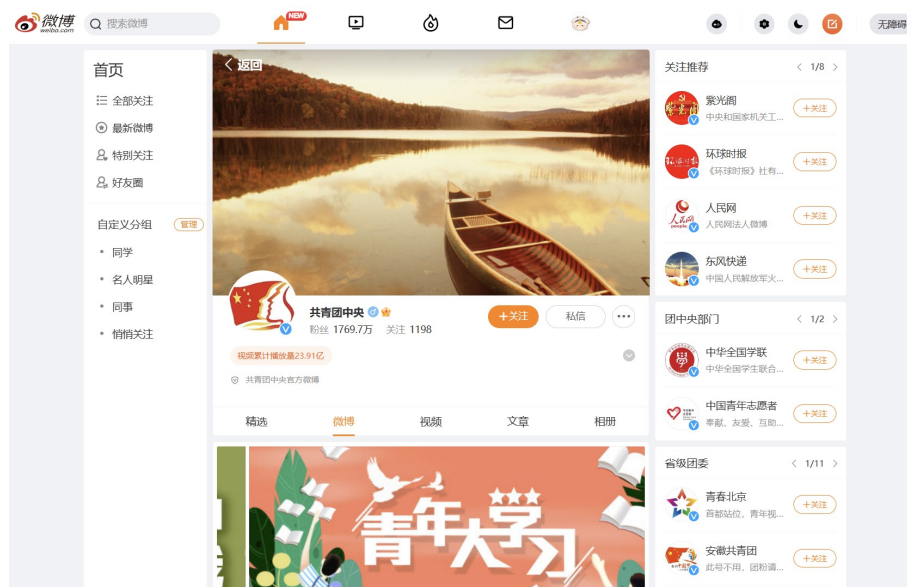


图 2.1: 共青团中央官方微博

2.2 爬虫实现

本报告数据爬取使用 Github 基于 weibo.com 的新版 API 构建的爬虫项目，项目地址为：<https://github.com/ngnhuyong/WeiboSpider>，项目主页如图2.2所示。主要使用该项目下的 tweet.py 基于用户 id 进行爬虫，项目具体使用可参照其 README 文件。



图 2.2: 项目主页示例

该项目将爬取数据解析后存储为 json 格式，其单条数据示例如下所示。

```

1  {
2      "_id": "4762810834227120",
3      "mblogid": "LqlZNhJFm",
4      "created_at": "2022-04-27 10:20:54",
5      "geo": null,
6      "ip_location": null,
7      "reposts_count": 1890,
8      "comments_count": 1924,
9      "attitudes_count": 12167,
10     "source": "三星Galaxy S22 Ultra",
11     "content": "生于乱世纵横四海，义之所在不计生死，孤勇者陈恭一生当如是。#风起陇西
              今日开播# #风起陇西# 今晚，恭候你！"
12     "pic_urls": [],
13     "pic_num": 0,
14     "isLongText": false,
15     "user": {
16         "_id": "1087770692",
17         "avatar_hd": "https://tvax1.sinaimg.cn/crop.0.0.1080.1080.1024/
              40d610441y8gbhxwgy419j20u00u0goc.
              jpg?KID=imgbed,tva&Expires=
              1682768013&ssig=r1QurGoc2L",
18         "nick_name": "陈坤",
19         "verified": true,
20         "mbrank": 7,
21         "mbtype": 12,
22         "verified_type": 0
23     },
24     "video": "http://f.video.weibocdn.com/o0/CmQEWK1ylx07VAm0nrxe01041200YDIc0E010.
              mp4?label=mp4_720p&template=1280x720."

```

```

25         "url": "https://weibo.com/1087770692/LqlZNhJFm",
26         "crawl_time": 1682757213
27     }

```

2.3 数据采集结果

本报告共采集到共青团中央从 2019 年 1 月 1 日至 2023 年 7 月 16 日发微数据 25244 条，采集结果如图 2.3 所示。

The image shows a large JSON array of Weibo posts. Each post is an object with the following fields: 'ip_location', 'reposts_count', 'comments_count', 'attitudes_count', 'source', and 'content'. The 'content' field contains the text of the post, which includes various topics such as national holidays, social issues, and official statements. The posts are numbered from 25289 down to 25244.

图 2.3: 数据采集结果

第 3 章 数据预处理

与一般的数据处理全流程类似，本报告的实现过程主要包括数据的爬取、数据预处理、计算结果分析和数据的可视化展示等几个主要步骤。

3.1 数据去重与格式转换

首先，将 json 文件转换为 csv 格式，先写入表头，后循环读取每一行的数据的 key 值所对应的 value 后输出，对应函数实现如下。

```

1 def json_to_csv(input_file, output_file):
2     with open(input_file, 'r', encoding='utf-8') as json_file, open(output_file, 'w',
3                                     encoding='utf-8',
4                                     newline='') as
5
6         writer = csv.writer(csv_file)
7         writer.writerow(['user_id', 'mblogid', 'created_at', 'reposts_count', '
8                             comments_count', 'attitudes_count
9                             ',
10                             'content', 'pic_urls', 'pic_num', 'isLongText', 'url', '
11                             crawl_time'])
12
13     for line in json_file:
14         data = json.loads(line)
15         user_id = data['_id']
16         mblogid = data['mblogid']
17         created_at = data['created_at']
18         reposts_count = data['reposts_count']
19         comments_count = data['comments_count']
20         attitudes_count = data['attitudes_count']
21         content = data['content'].replace('\n', ' ') # 将换行符替换为空格，避
22                                                         免输出时自动换行
23         pic_urls = data['pic_urls']
24         pic_num = data['pic_num']
25         islongtext = data['isLongText']
26         url = data['url']
27         crawl_time = data['crawl_time']
28
29         writer.writerow([user_id, mblogid, created_at, reposts_count,
30                             comments_count,
31                             attitudes_count,
32                             content, pic_urls, pic_num, islongtext, url,
33                             crawl_time])

```

此外，使用 pandas 的内置函数 `drop_duplicates` 进行去重处理并使用 `sort_value` 实现按照 'created_at' 列的时间降序排列，核心代码如下。

```
1 df = pd.read_csv('gqt_output.csv')
2 df.drop_duplicates(inplace=True)
3 df = df.sort_values('created_at', ascending=False)
4 df.to_csv('gqt_output.csv', index=False)
```

3.2 文本内容清洗

定义正则表达式匹配去除 content 中的无效文本。

```
1 def clean_text(text):
2     text = re.sub(r'【[^】】+】', '', text)
3     text = re.sub(r'\[[^\]\[]+\]', '', text)
4     text = re.sub(r"(回复)?(//)?\s*@\S*?\s*(?:|:| |!$)", " ", text) # 去除正文中的@
                                     和回复/转发中的用户名
5     text = re.sub(r'\([^()]+\)', '', text)
6     text = re.sub(r'\([^\(]+\)', '', text)
7     text = re.sub(r"#([^\#]+)#", "", text)
8     text = re.sub(r" ", "", text)
9     text = text.replace("\n", " ")
10    text = re.sub(r"↓↓", "", text)
11    text = re.sub(r"(\s)+", r"\1", text)
12    link_regex = r'http://t\.cn/[\w]+(,)?(?:=[^\w]|,|$)'
13    text = re.sub(link_regex, '', text)
14    stop_terms = ['展开', '全文', '展开全部', '一个', '网页', '链接', '?【', 'ue627',
                  'c【', '10', '一下', '一直',
15                  'u3000', '24', '12',
16                  '30', '?我', '15', '11', '17', '?\\', '显示地图', '原图']
17    for x in stop_terms:
18        text = text.replace(x, "")
19    return text.strip()
```

第 4 章 数据分析与可视化

本报告使用可视化库为 pyecharts，能够在 web 端实现数据交互查看，最终 html 文件位于 analyze 目录下，预览动态交互可运行 vi_merged.html 可查看，以下部分仅展示静态图片。

4.1 逐年发微热力图

统计分析每年中月日发微博频率，主要代码如下。

```

1 df = pd.read_csv('cleaned_text.csv', parse_dates=['created_at'])
2 date_counts = df['created_at'].dt.date.value_counts()
3 cal_result = [(date, count) for date, count in date_counts.items()]
4 page = Page(page_title='共青团中央微博分析', layout=Page.DraggablePageLayout)
5 for year in range(2019, 2024):
6     calendar = (
7         Calendar()
8         .add("",
9             cal_result,
10            calendar_opts=opts.CalendarOpts(range_=str(year)),
11            daylabel_opts=opts.CalendarDayLabelOpts(name_map="cn"),
12            monthlabel_opts=opts.CalendarMonthLabelOpts(name_map="cn"),
13        )
14        .set_global_opts(
15            title_opts=opts.TitleOpts(title=f"{year} 年每天发微频率"),
16            visualmap_opts=opts.VisualMapOpts(
17                max_=55, min_=5, orient="horizontal", is_pieewise=True,
18                pos_top="230px",
19                pos_left="100px",
20            )
21        )

```

逐年可视化热力图如下。

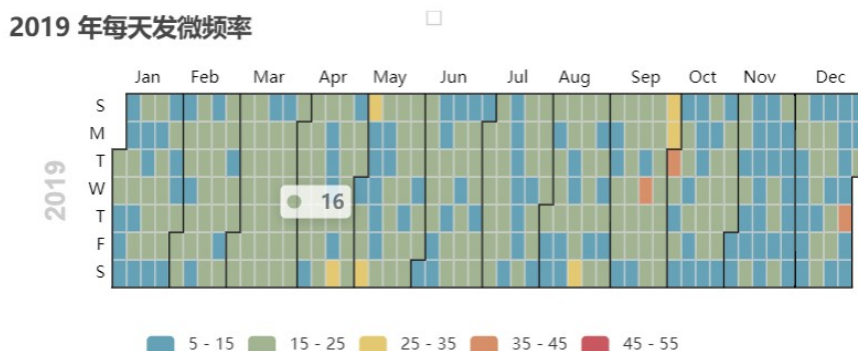


图 4.1: 2019 热力图

2020 年每天发微频率

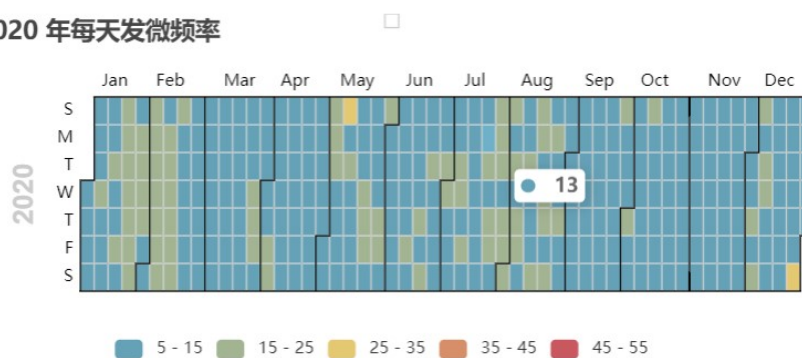


图 4.2: 2020 热力图

2021 年每天发微频率

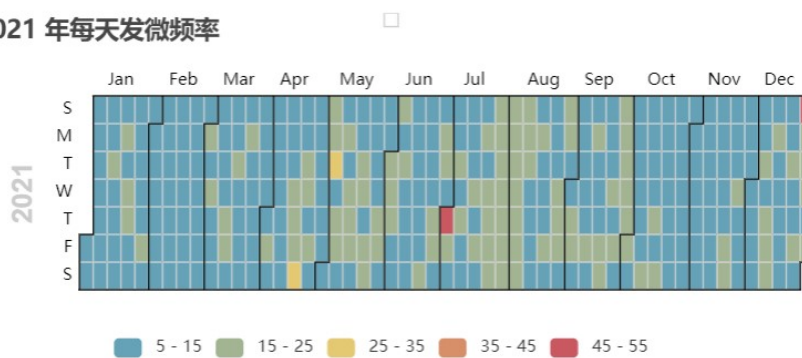


图 4.3: 2021 热力图

2022 年每天发微频率

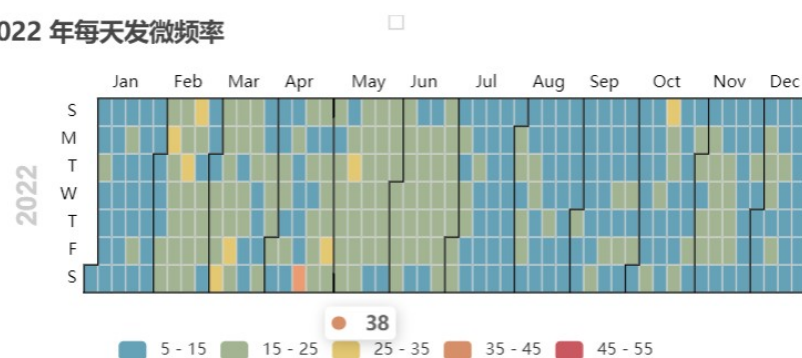


图 4.4: 2022 热力图

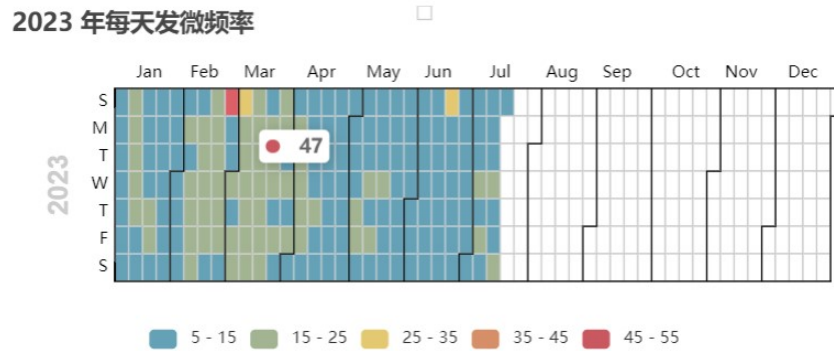


图 4.5: 2023 热力图

4.2 发微时间统计

使用 pandas 内置函数按照年、月、周、日、时分别统计共青团中央发微时间。

```
1 df = pd.read_csv('cleaned_text.csv', parse_dates=['created_at'])
2 dateyear_counts = df['created_at'].dt.year.value_counts()
3 datemonth_counts = df['created_at'].dt.month.value_counts()
4 dateweek_counts = df['created_at'].dt.weekday.value_counts()
5 dateday_counts = df['created_at'].dt.day.value_counts()
6 datehour_counts = df['created_at'].dt.hour.value_counts()
```

4.2.1 逐年统计

利用 pyecharts 制作逐年统计可视化，结果如图4.6所示。

```
1 c_year = (
2     Bar()
3     .add_xaxis([i[0] for i in resultyear])
4     .add_yaxis("共青团中央", [i[1] for i in resultyear])
5     .set_global_opts(
6         title_opts=opts.TitleOpts(title="每年发微频率"),
7         yaxis_opts=opts.AxisOpts(name="次数"),
8         xaxis_opts=opts.AxisOpts(name="时间"),
9     )
10
11 )
```



图 4.6: 逐年统计

4.2.2 每年逐月统计

逐月统计可视化。

```
1  c_month = (  
2      Bar()  
3      .add_xaxis(Faker.months)  
4      .add_yaxis("共青团中央", [i[1] for i in resultmonth])  
5      .set_global_opts(  
6          title_opts=opts.TitleOpts(title="每月发微频率"),  
7          yaxis_opts=opts.AxisOpts(name="次数"),  
8          xaxis_opts=opts.AxisOpts(name="时间"),  
9      )  
10  
11 )
```

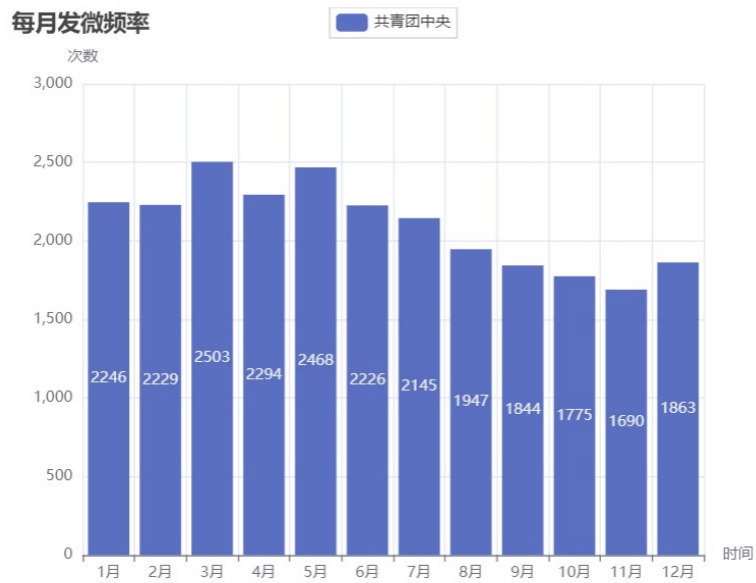


图 4.7: 逐月统计

4.2.3 每月逐天统计

本小节后续部分与上述实现过程相似，不再赘述，仅展示最终可视化结果。

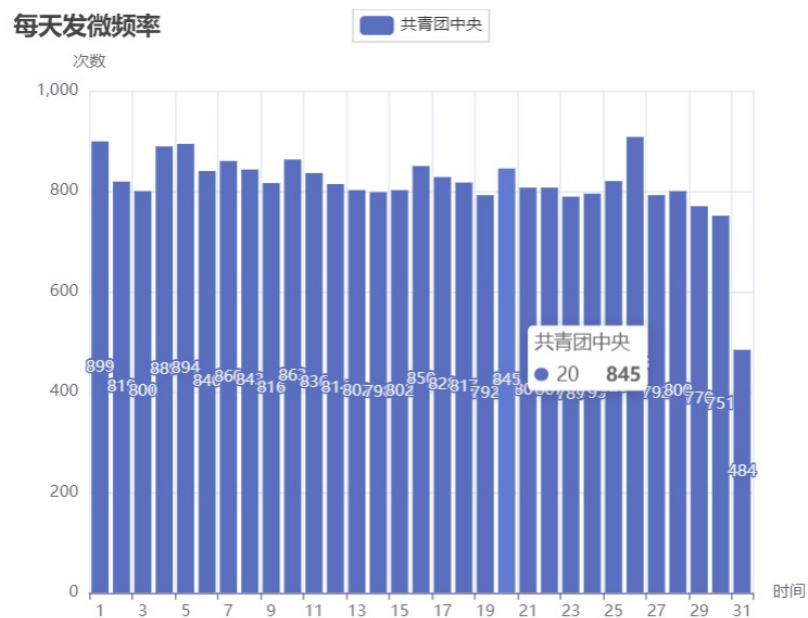


图 4.8: 逐天统计

4.2.4 每周逐天统计

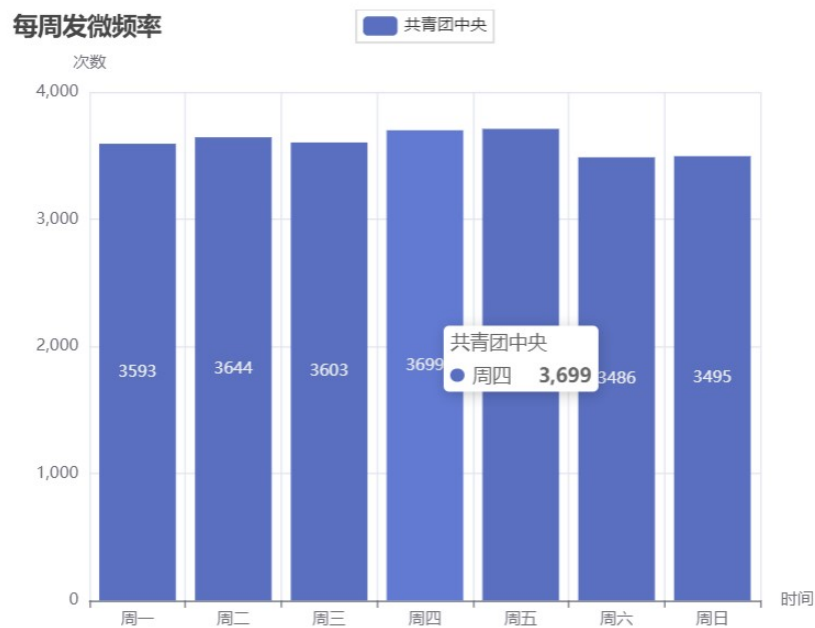


图 4.9: 逐周统计

4.2.5 每天逐时统计

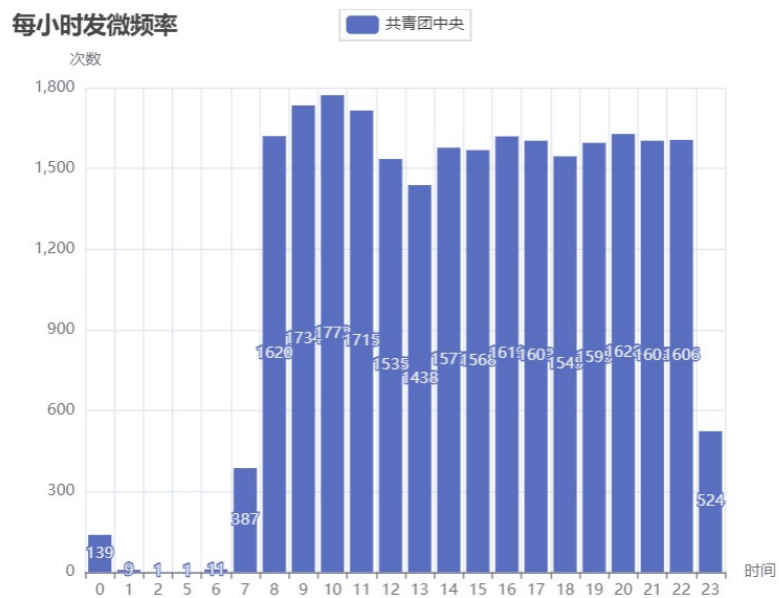


图 4.10: 逐时统计

4.3 周时热力统计

使用 pandas 的 groupby 多列分组先按照 hour 来分组。每组内，再按 weekday 来分组，并统计每个周时数量返回。

```

1 def to_weekhour(file, file_colm):
2     df = pd.read_csv(file)
3     df[file_colm] = pd.to_datetime(df[file_colm])
4     hour_week_counts = df.groupby([df[file_colm].dt.hour, df[file_colm].dt.weekday]
5                                   )[file_colm].count()
6     result = [[hour, week, count] for (hour, week), count in hour_week_counts.items()
7               ]
8     result.sort(key=lambda x: x[0]) #按照第一个排序
9     return result

```

调用 to_weekhour 函数将分类统计结果传回 pyecharts 进行可视化绘图，核心代码如下。

```

1 c_heatmap = (
2     HeatMap()
3     .add_xaxis(Faker.clock)
4     .add_yaxis(
5         "发微次数", Faker.week, to_weekhour('cleaned_text.csv', 'created_at'),
6         label_opts=opts.LabelOpts(position="middle")
7     )
8     .set_global_opts(
9         title_opts=opts.TitleOpts(title="共青团中央发微周时热力图"),
10        visualmap_opts=opts.VisualMapOpts(),
11    )
12 )

```

最终可视化效果如图4.11

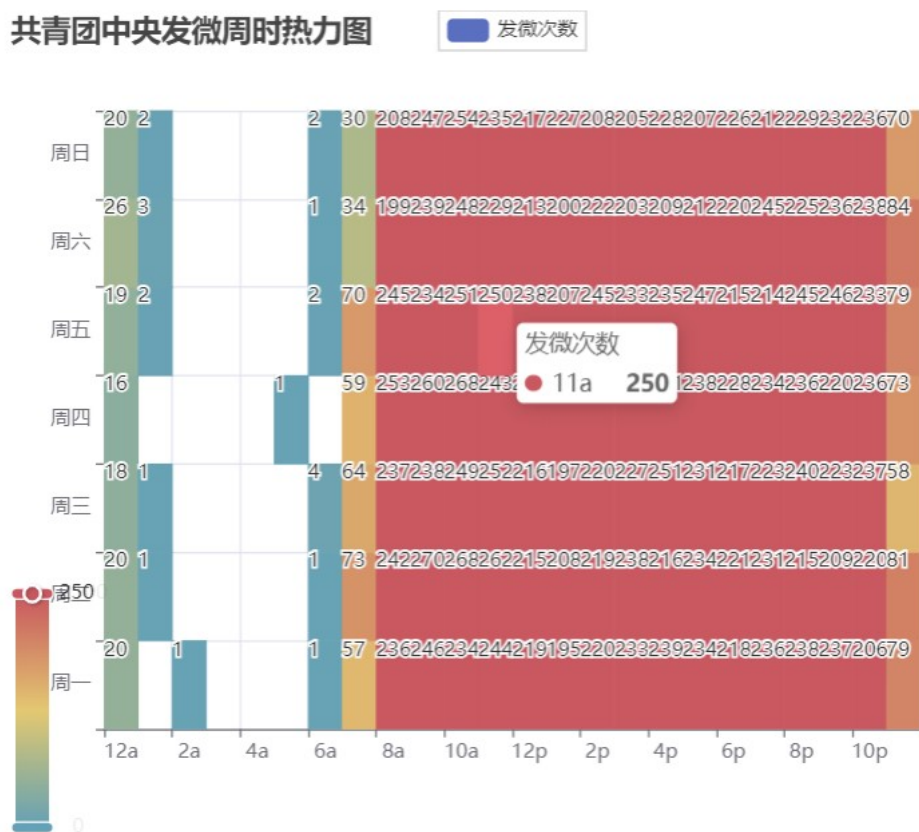


图 4.11: 周时热力图

4.4 高频词统计

定义函数统计词频，添加自定义词典和停用词词典，使用 jieba 库进行分词，jieba 分词工具有三种主要分词模式，本报告选取 HMM 模式。

```

1 def count_word_frequency(csv_file_path, column_name):
2     df = pd.read_csv(csv_file_path)
3     df[column_name].fillna('', inplace=True)
4     content = df[column_name]
5     merged_content = ' '.join(content)
6     pattern = re.compile(u'\t|\n|\.|-|:|;|\)|\(|\?|!|"')
7     string_data = re.sub(pattern, '', merged_content)
8     jieba.suggest_freq('共青团', True)
9     jieba.load_userdict('user_words.txt')
10    seg_list_exact = jieba.cut(string_data, cut_all=False, HMM=True) # 分词模式
11    object_list = []
12    # 去除停用词
13    with open('stop_words.txt', 'r', encoding='UTF-8') as meaninglessFile:
14        stopwords = set(meaninglessFile.read().split('\n'))
15    stopwords.add(' ')
16    for word in seg_list_exact:

```

```

17     if word not in stopwords:
18         object_list.append(word)
19 word_counts = collections.Counter(object_list)
20 for word in list(word_counts.keys()):
21     if len(word) == 1:
22         del word_counts[word]
23 word_counts_top = word_counts.most_common(300)
24 return word_counts_top

```

使用 pyecharts 制作词云可视化，此外本报告也采用了自定义功能更强大的 word-cloud 库结合 matplotlib 等库制作了共青团中央微博内容词云可视化，相关结果为摘要页微博图标背景的词云可视化图片。

```

1 c_wordcloud = (
2     WordCloud()
3     .add(
4         "",
5         count_word_frequency('cleaned_text.csv', 'content'),
6         word_size_range=[20, 100],
7         textstyle_opts=opts.TextStyleOpts(font_family="cursive"),
8     )
9     .set_global_opts(title_opts=opts.TitleOpts(title="top200 词云"))
10 )

```

使用 pyecharts 实现的可视化效果如图所示，能够实现在网页端点击悬停等交互。



图 4.12: 词频统计词云

4.5 LDA 主题分析

首先对文本内容进行分词，分词程序位于 analyze/LDA/cutword.py，具体实现与词

云处一致调用 jieba 库。分词后，基于已分词文档使用 gensim 库进行 LDA 主题分析，并使用 pyLDAvis 进行可视化展示。思路阐述参考下列代码注释

```
1 import pyLDAvis.gensim_models
2 from gensim.corpora.dictionary import Dictionary
3 from gensim.models.ldamodel import LdaModel
4 from gensim import corpora, models
5
6 if __name__ == '__main__':
7     words_ls = []
8     with open('outdemowei.txt', 'r', encoding='UTF-8') as f:
9         data = f.readlines()
10        for line in data:
11            words_ls.append(line.split(' '))
12    f.close()
13    # 构造词典，将每个唯一的词语映射到一个唯一的整数ID
14    dictionary = corpora.Dictionary(words_ls)
15    # 基于词典，使【词】→【稀疏向量】，并将向量放入列表，形成【稀疏向量集】（即语料库）
16    corpus = [dictionary.doc2bow(words) for words in words_ls]
17    # tf-idf编码
18    tfidf = models.TfidfModel(corpus)
19    # 将语料库转换为tf-idf编码形式
20    corpus_tf = tfidf[corpus]
21    # lda模型，num_topics设置主题的个数
22    lda = LdaModel(corpus=corpus_tf, id2word=dictionary, num_topics=4, random_state=100, iterations=50)
23
24    for topic in lda.print_topics(num_words=10):
25        print(topic)
26    # 用pyLDAvis可视化
27    plot = pyLDAvis.gensim_models.prepare(lda, corpus_tf, dictionary)
28    pyLDAvis.save_html(plot, 'ldaweibo.html')
```

通过 LDAModel 抽取 4 个主题后，最终可视化效果如图4.13所示，html 可视化源文件位于 analyze/LDA/ldaweibo.html。

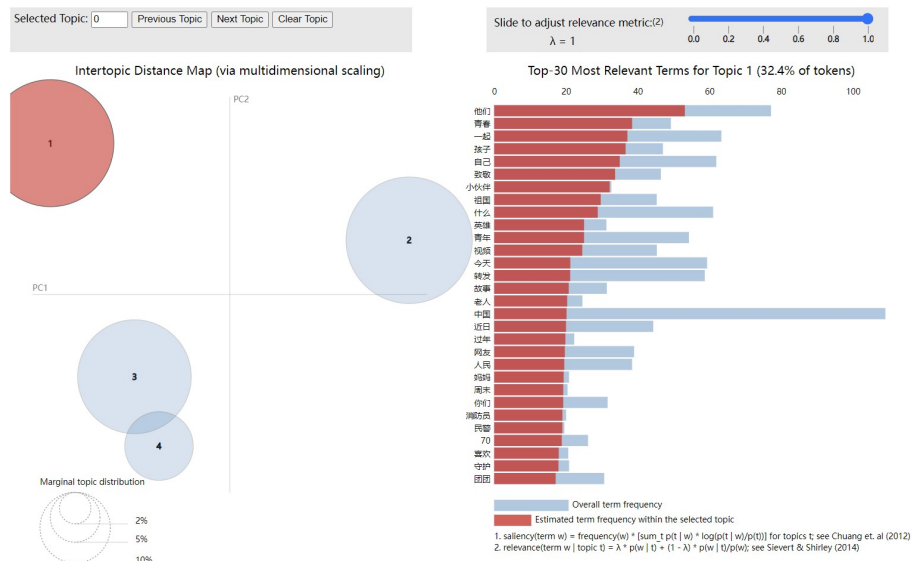


图 4.13: LDA 分析结果

第 5 章 总结展望

本文在 python3.10 环境下，基于对 Github 开源爬虫项目的改动实现数据爬取，调用了 pandas,pyecharts,jieba,gensim 等软件包，对所获得数据进行了统计分析和可视化交互展示。国内新闻传播顶刊有计算传播论文中写到利用 python 进行情感分析，使用 snownlp 进行分析，作为一个封装好的项目，再去调用已经没有难度，基础还是在于循环，判断等语句的应用。本报告并不试图去调用一些可能是用家具评论数据进行有监督机器学习的工具通过一行代码就能实现情感分析来用到自己所爬取的可能和模型预训练预料毫无关系的数据，试问这样的实现难度和简单写个函数让每一条数据调用分析后输出结果可视化有什么不同呢？背后更深的应该是整个模型实现的过程，有什么条件需要满足，能不能用到我的项目上。本报告也存在遗憾，在 LDA 分析的时候 gensim 库是做的无监督，但是确定主题数的时候没能实现一些暴力搜索之类的方法来确定主题数，还是人为尝试了很多数字之后确定的主题数，后续有时间一定解决这部分的遗憾实现编程的严谨。

写完报告主体想到读过的很多新闻传播论文，真诚希望很多文章在使用某个方法非要去赋予概念的时候能够多一点技术理性，首先理解技术逻辑、判断自己的项目是否适合使用该数据分析技术、最终分析社会问题，否则如果仅仅是冠以计算和智能的 tag 毫无判断也无法读懂方法核心为了让自己的研究显得更“高端”、方法更新，在用一个错误的手段分析社会解读社会，这个世界真的会因为所谓学术更好嘛？

老师对待计算机，对待编程，对待技术、对待社会的严谨，从中收获很多。希望国内计算传播越来越好！