



# 南昌大学实验报告

学生姓名：丁俊 学号：8003119100 专业班级：信息安全 193 班  
实验类型：☐ 验证 ☒ 综合 ☐ 设计 ☐ 创新 实验日期：2021.11.21 实验成绩：

## 一、实验项目名称

数据爬取与分析

## 二、实验目的

1. 掌握 request、BeautifulSoup 和 re 等爬虫实现工具
2. 掌握 Numpy、Pandas 和 Matplotlib 基本使用
3. 掌握表示、清洗、统计和展示数据的能力

## 三、实验任务

- 1、参照《青春有你 2》、《乘风破浪的姐姐》、《安家》数据爬取与分析，自选一个主题完成数据爬取与分析。
- 2、参照“数据可视化 - 中国的天气数据”掌握 Basemap、seaborn、pyecharts 等可视化库的使用。（可选）

## 四、主要仪器设备及耗材

软件：Anaconda 或者 pycharm 等

## 五、实验步骤与结果

### 1、爬虫目标

- 猫眼电影榜单 Top100 榜，网址为 <https://www.maoyan.com/board/4?offset=0>  
从网页中提取出 top100 的排名、电影名称、评分、上映时间、详情链接、上映地区、演员、剧照图片等。
- 根据爬虫结果，进行简单的可视化分析。

### 2、爬取步骤

- 网址 url 分析

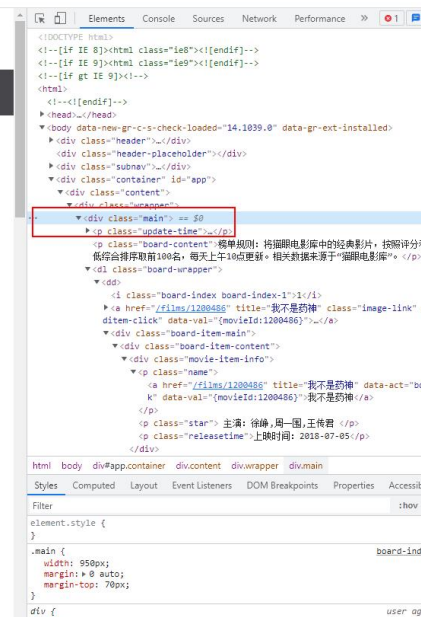
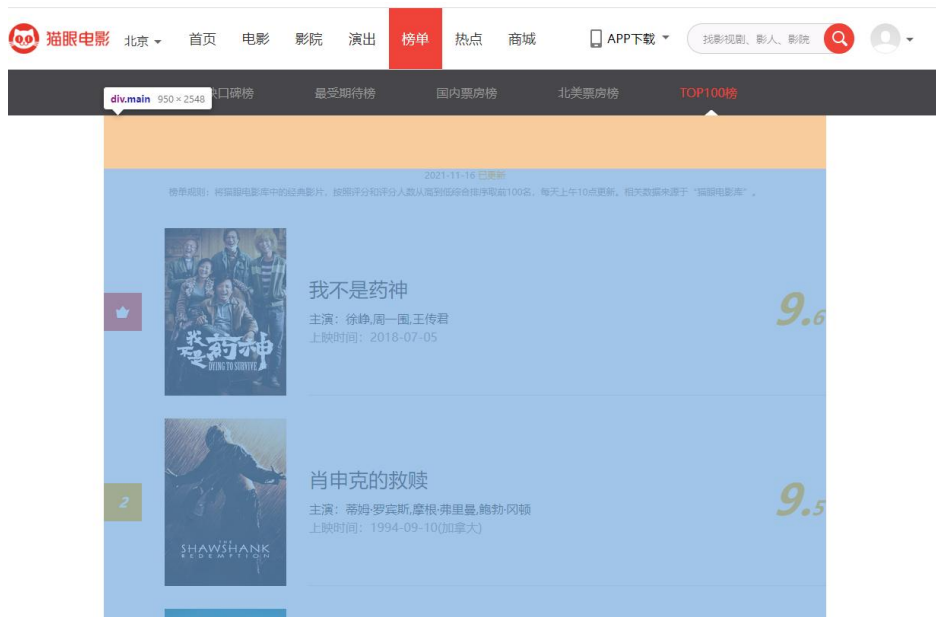
首先，打开猫眼 top100 榜单电影的地址 <https://www.maoyan.com/board/4?offset=0>  
页面非常简单，页面底部有分页显示栏目，我们要做得就是进行分页爬取，总共需要爬取 10 页数据，每一页的 html 样式和结构都类似，这里使用一个偏移量 offset 对构造 10 页地址，分别进行访问。

- Requests 进行访问

```
headers = {  
    'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/94.0.4606.81 Safari/537.36 ',  
}  
  
proxy = {  
    'http': '114.106.171.160:13456',  
    'https': '114.104.121.169:13456',  
    'http': '113.124.92.234:13456'  
}  
  
# 爬取网页并返回网页内容div块中的内容  
def crawl_data(url):  
    try:  
        response = requests.get(url, headers = headers)  
        if response.status_code == 200:  
            soup = BeautifulSoup(response.text, 'lxml')  
            div_pos = soup.find_all('div', {'class': 'main'})  
            return div_pos  
        else:  
            return None  
    except:  
        return '错误'
```

定义一个函数 `crawl_data()`，传入 `url` 参数，在 `main` 函数中设置 `url` 和偏移量。则合理使用了 `proxy` 代理，因为猫眼电影网有反爬机制，避免同一 `ip` 过多地访问而锁 `ip`，即使如此，也很难一次性爬取这么多的数据，所以我分了两次进行爬取。

```
def main(offset):
    url = 'http://maoyan.com/board/4?offset=' + str(offset)
    html = crawl_data(url) # 访问一次
    # print(html)
    # 下载数据
    # dum_data(html)
    # url = get_img_urls(info_url('https://www.maoyan.com/films/1200486'))
    # get_other_info(info_url('https://www.maoyan.com/films/1297'))
    show_score()
    show_actors()
    show_year()
```



从网页结构可以看出，电影内容是包含在一个'class=main'的div块中的，直接使用find\_all函数找到这个div返回。

## ● 提取电影字段信息



可以看到每个电影信息都被包含在一个 dd 内部

```
bs = BeautifulSoup(str(div_html), 'lxml')
movies = [] # 存储电影信息
for item in range(10):
    movie = {}
    movie['rank'] = bs.select('dd i.board-index')[item].string
    movie['name'] = bs.select('.name a')[item].string
    movie['score'] = bs.select('.integer')[item].string + bs.select('.fraction')[item].string
    movie['time'] = get_release_time(bs.select('.releasetime')[item].string.strip()[5:])
    movie['link'] = 'https://www.maoyan.com' + bs.select('a.image-link')[item].get('href')
    link = movie['link']
    name = movie['name']
    next_bs = info_url(link) # 访问电影详情页面
    movie['area'] = get_other_info(next_bs)
    # 访问电影详情页面
    movie_pics = get_img_urls(next_bs) # 获得图片链接
    down_pics(name, movie_pics) # 下载图片保存文件夹
    movie['actors'] = bs.find_all(name='p', attrs={'class': 'star'})[item].string.strip()[3:]
    # 相关演员信息
```

这是提取网页中电影的相关字段的代码，使用 select CSS 选择器进行选取，当然还有一些正则表达式用于提取时间和地区字段，如下，`r'(.?)(\($))'` 匹配“上映时间:2018-07-05”，`re` 将其分成两组，前者表示任意字符且非贪婪匹配，然后进行分片选取第五位开始的数据，即年份。

```
# 提取上映时间函数
def get_release_time(data):
    pattern = re.compile(r'(.?)(\($))')
    item = re.search(pattern, data)
    if item is None:
        return '未知'
    return item.group(1)

# 提取中文字段
def get_area(data):
    item = re.search('[\u4e00-\u9fa5]*', data)
    if item is None:
        return '未知'
    # print(item)
    return item.group()
```

```
tmp_area = bs.select('.movie-brief-container ul li')[1].string.strip()
# print(tmp_area)
area = get_area(tmp_area)
```

电影的上映地区需要点进电影详情链接才能爬取到，然后传到 `get_area` 中获取其中的中文字段即可。



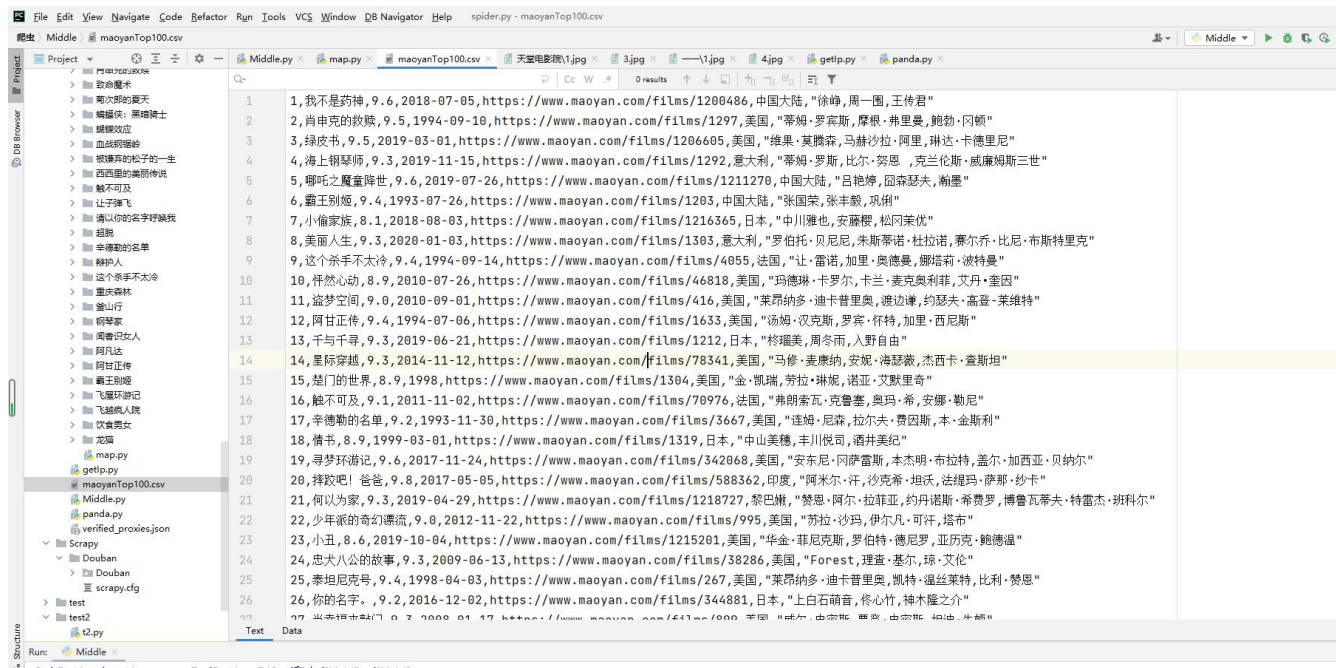


## ● 保存数据到 csv 文件

# 把字典形式的数据添加到csv文件中

```
def write_to_file(item):  
    with open('maoyanTop100.csv', 'a', encoding='utf_8_sig', newline='') as f:  
        # a为追加模式  
        fieldnames = ['rank', 'name', 'score', 'time', 'link', 'area', 'actors']  
        w = csv.DictWriter(f, fieldnames=fieldnames)  
        w.writerow(item)
```

这里采用将字典形式的内容保存到 csv 文件中，如图

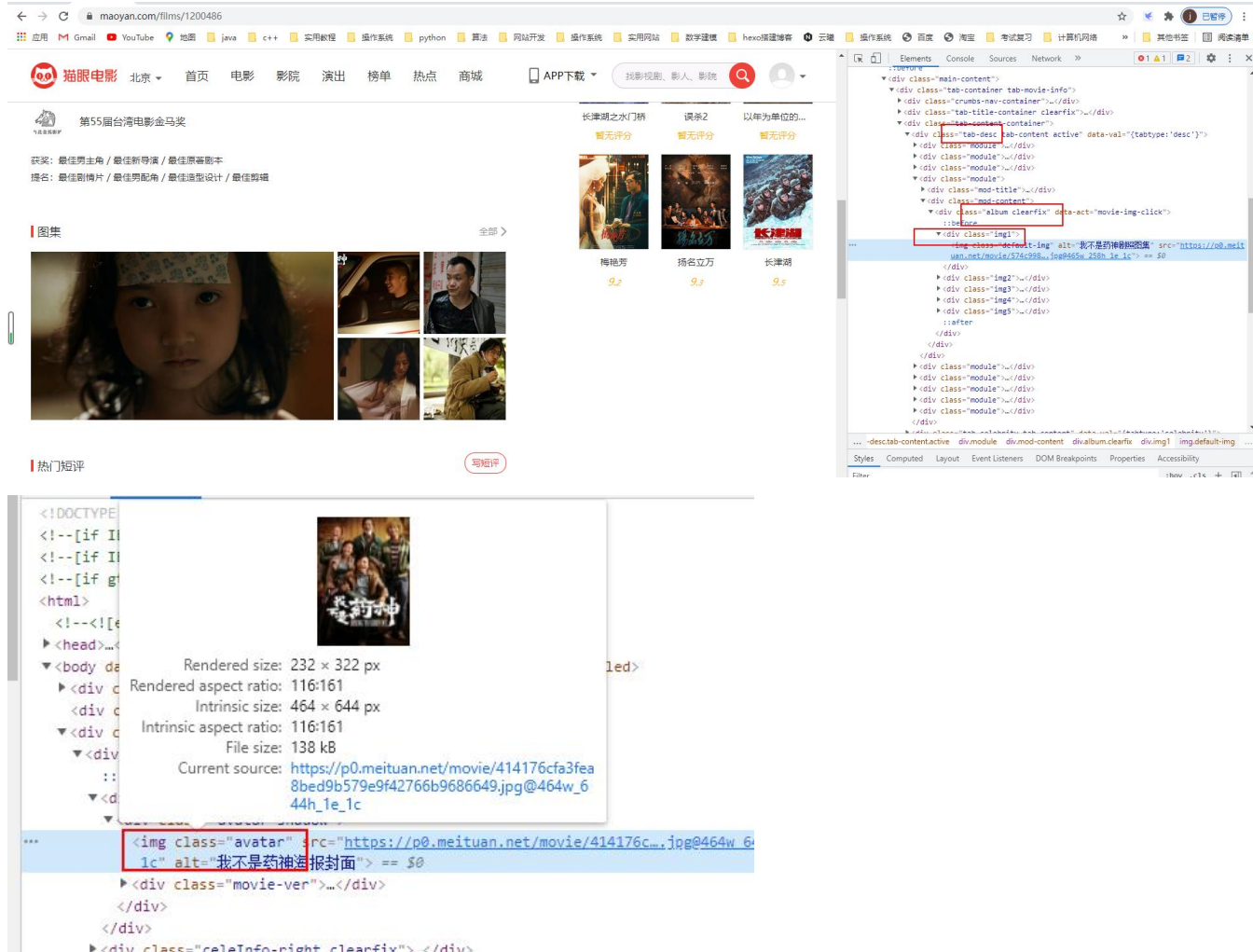


## ● 获取电影的图片链接

# 访问电影剧照，并保存图片，先建立图片的urls集合

```
def get_img_urls(bs):  
    # 访问网页  
    # response_detail = requests.get(movies_links,headers = headers)  
    # bs = BeautifulSoup(response_detail.text,'lxml')  
    pic_urls = [] # 存储图片的链接集合  
    main_img = bs.find(name='img',attrs={'class':'avatar'})  
    main_url = main_img.get('src')  
    pic_urls.append(main_url)  
  
    divs_html = bs.select('.tab-desc .album div')  
    for div in divs_html:  
        img_url = div.img.get('data-src')  
        pic_urls.append(img_url)  
    print(pic_urls) # 打印出电影的剧照集  
    return pic_urls # 返回链接集合
```

首先找到电影的封面链接，获取“src”字段，并保存到图片列表中，然后在电影图集中，找到相关字段“data-src”依次保存在列表中并返回。



访问电影详情链接函数，为了避免代码冗余，直接将这多条语句封装成一个函数，返回 BeautifulSoup 对象。

# 访问电影详细函数、避免多次访问, 返回Beautiful对象

```
def info_url(movies_link):
```

```
    """
```

```
    封装请求，避免多次访问，需要人工验证
```

```
    :param movies_link:
```

```
    :return:
```

```
    """
```

```
    response_detail = requests.get(movies_link, proxies=_proxy, headers=_headers)
```

```
    bs = BeautifulSoup(response_detail.text, 'lxml')
```

```
    return bs
```

## ● 下载图片并保存到文件夹

```
# 下载图片保存到每个电影的文件夹中
def down_pics(name, pic_urls):
    print(pic_urls)
    pic_path = '电影图片/' + name + '/'
    # 在图片->电影名->图片
    if not os.path.exists(pic_path):
        os.makedirs(pic_path) # 递归创建目录
    num = 0
    for i, pic_url in enumerate(pic_urls):
        print('-----') # 分割线
        try:
            pic = requests.get(pic_url, proxies=proxy, timeout=20)
            pic_str = str(i + 1) + '.jpg'
            with open(pic_path + pic_str, 'wb') as f:
                f.write(pic.content)
            print('成功下载第' + str(i+1) + '张图片' + str(pic_url))
            num = num + 1
        except Exception as e:
            print('下载' + str(pic_url) + '张图片失败')
            print(e)
            continue
```

传入一个图片链接列表，首先创建电影的图片文件夹，依次访问列表中的图片链接写入文件夹。

## 3、数据分析

### ● 电影评分分布图

```
1. # 各电影的评分饼状图
2. def show_score():
3.     matplotlib.rcParams['font.sans-serif'] = ['KaiTi']
4.     plt.style.use('ggplot')
5.     columns = ['rank', 'name', 'score', 'time', 'link', 'area', 'actors']
6.     df = pd.read_csv('maoyanTop100.csv', encoding='utf-8', header = None, names = columns, index_col='rank')
7.     # 8.0~8.5 8.5~9.0 9.0~9.5 >9.5
8.     size1 = 0
9.     size2 = 0
10.    size3 = 0
11.    size4 = 0
12.    plt.style.use('ggplot')
13.    plt.figure(figsize=(14,14))
14.    plt.tight_layout()
15.    for item in df.score:
16.        if 8.0 <= item < 8.5:
17.            size1 += 1
18.        elif 8.5 <= item < 9.0:
19.            size2 += 1
```

```

20.     elif 9.0 <= item < 9.5:
21.         size3 += 1
22.     elif 9.5 <= item:
23.         size4 +=1
24.     labels = ['8.0~8.5', '8.5~9.0', '9.0~9.5', '>=9.5']
25.     sizes = [size1,size2,size3,size4]
26.     explodes = [0,0,0.05,0]
27.     patches, l_text, p_text = plt.pie(sizes, labels=labels,explode = explodes, autopct='%1.1f%%', shadow=
        False, startangle=90, labeldistance=1)
28.     # 返回值得到 p_text 是饼图内部文字、l_texts 是图外的文本
29.     for t in l_text:
30.         t.set_size(40)
31.         t.set_rotation(30)
32.     l_text[1].set_rotation(90)
33.     for s in p_text:
34.         s.set_size(40)
35.
36.     plt.title('top100 电影评分分布',fontsize = 45)
37.     plt.tight_layout()
38.     plt.axis('equal')
39.     plt.savefig('top100 电影评分分布')
40.     plt.show()

```

分别定义四个级别 8.0~8.5(含 8.0)、8.5~9.0(含 8.5)、9.0~9.5(含 9.0) >=9.5

## ● 演员电影数量 top10

```

1. # 电影数量最多的演员
2. def show_actors():
3.     matplotlib.rcParams['font.sans-serif'] = ['KaiTi']
4.     plt.style.use('ggplot')
5.     columns = ['rank', 'name', 'score', 'time', 'link', 'area', 'actors']
6.     df = pd.read_csv('maoyanTop100.csv', encoding='utf-
7.     8', header=None, names=columns, index_col='rank')
7.     actor_total = df.actors
8.     # print(actor_total)
9.     actor_list = []
10.    for i in actor_total.str.replace(',').str.split(','):
11.        actor_list.extend(i)
12.    # print(actor_list)
13.    # print(len(actor_list))
14.    # print(actor_list.count('刘德华'))
15.    # 去除重复的演员名

```



```

16. actor = set(actor_list)
17. print(actor)
18. actor_all = {}
19. # actor_all 字典 {演员:数量}
20. for i in actor: # 遍历演员 (i)
21.     if actor_list.count(i) > 1:
22.         # 选出电影数量超过 1 部的演员，并记录数量
23.         actor_all[i] = actor_list.count(i)
24. actor_all = sorted(actor_all.items(), key = lambda actor_list:actor_list[1], reverse = True)
25. # items()方法返回可遍历的(键，值)元组数组,lambda 表示按字典的第二个值排序(即电影数)
26. # 元组变为字典
27. actor_all = dict(actor_all[:10])
28. print(actor_all)
29. # 绘图
30. x = list(actor_all.keys())
31. y = list(actor_all.values())
32. plt.bar(range(10), y, tick_label = x)
33. plt.xticks(rotation = 270)
34. for x,y in enumerate(y):
35.     plt.text(x, y + 0.1, '%s' %round(y, 1), ha = 'center', color = '#6D6D6D')
36. plt.title('演员电影作品数量排名 top10', color = '#6D6D6D')
37. plt.xlabel('演员')
38. plt.ylabel('数量(部)')
39. plt.tight_layout() # 填充整个图像区域
40. plt.savefig('演员电影作品数量排名 top10.jpg')
41. plt.show()

```

这个函数画的是一个柱形图。首先 `pd.read_csv` 读取 csv 文件的内容转化为 `DataFrame` 类型，因为在 csv 中同一部电影的演员不止一个，需要使用 `split` 函数以“，”分割开每一个演员，继续使用 `set` 集合进行去重，使用 `sorted` 函数进行排序，然后 `plt.text()` 在每一个柱上添加注释数据。

## ● 电影年份与数量关系

```

1. # 电影作品数量集中的年份
2. def show_year():
3.     plt.rcParams['font.sans-serif'] = ['SimHei']
4.     columns = ['rank', 'name', 'score', 'time', 'link', 'area', 'actors']
5.     df = pd.read_csv('maoyanTop100.csv', encoding='utf-8', header=None, names=columns, index_col='rank')
6.     # print(df['time'])
7.     df['year'] = df['time'].map(lambda x:x.split('-')[0]) # 取时期格式中的第一个数字年份作为处理结果 year 列值
8.     # print(df['year'])
9.     # 统计各年份电影数量

```

```

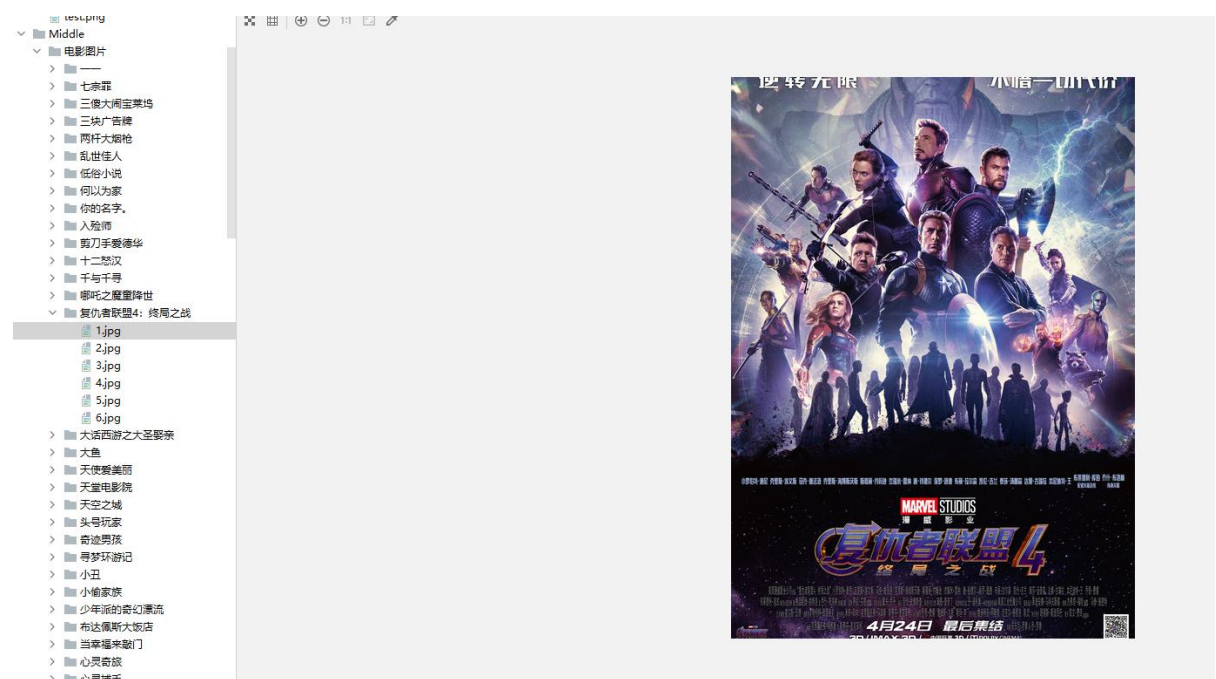
10. group_year = df.groupby('year')
11. group_year_mount = group_year.year.count() # Series 类型
12. for x,y in enumerate(list(group_year_mount.values)):
13.     plt.text(x, y+0.2, '%s' %round(y,1), ha = 'center', color = "#6D6D6D", fontsize = 13)
14. # print(group_year_mount)
15. plt.title('年份电影数量分布',fontsize = 15)
16. plt.xticks(rotation = 270, fontsize = 10)
17. plt.xlabel('年份(年)')
18. plt.ylabel('数量(部)')
19. plt.plot(group_year_mount, 'o-g',ms = 5)
20. plt.tight_layout()
21. plt.savefig('年份电影数量分布.jpg')
22. plt.show()

```

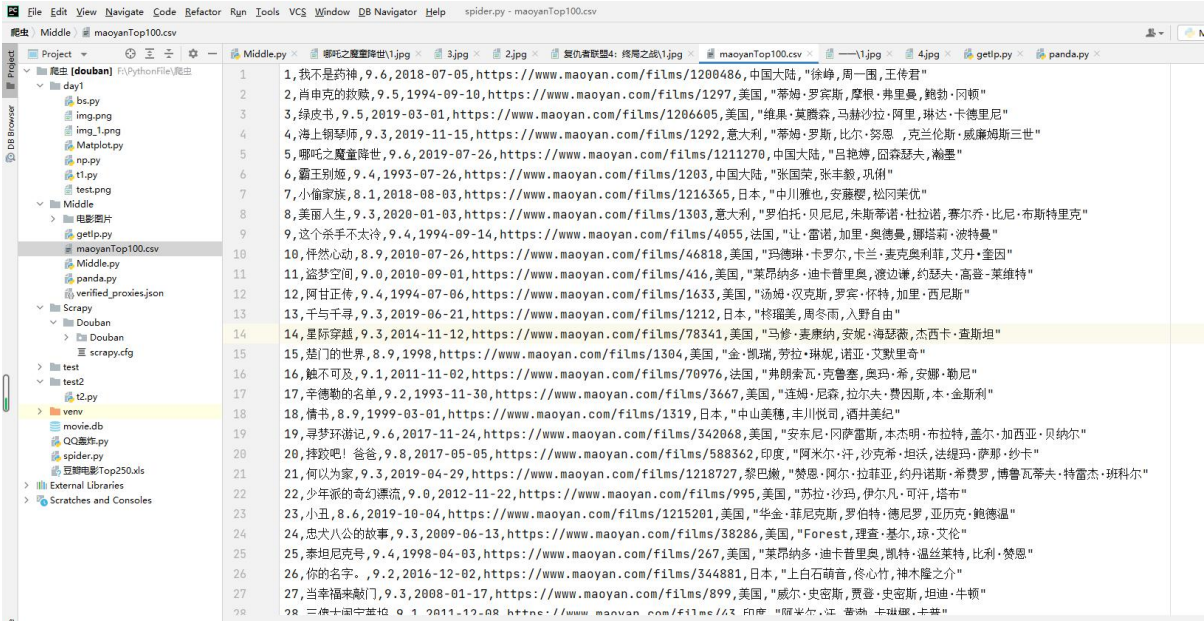
这个函数画的是折线图，这里采用了 map 函数将“time”这一列映射为“year”新的一列，df.groupby(‘year’) 对年份进行分组，每一组对应其中这一年的电影数量(year.count())，生成的 group\_year\_mount 是一个 Series 类型，直接传入 plt 自动生成年份—数量图像。

## 4、实验结果

### 电影图片

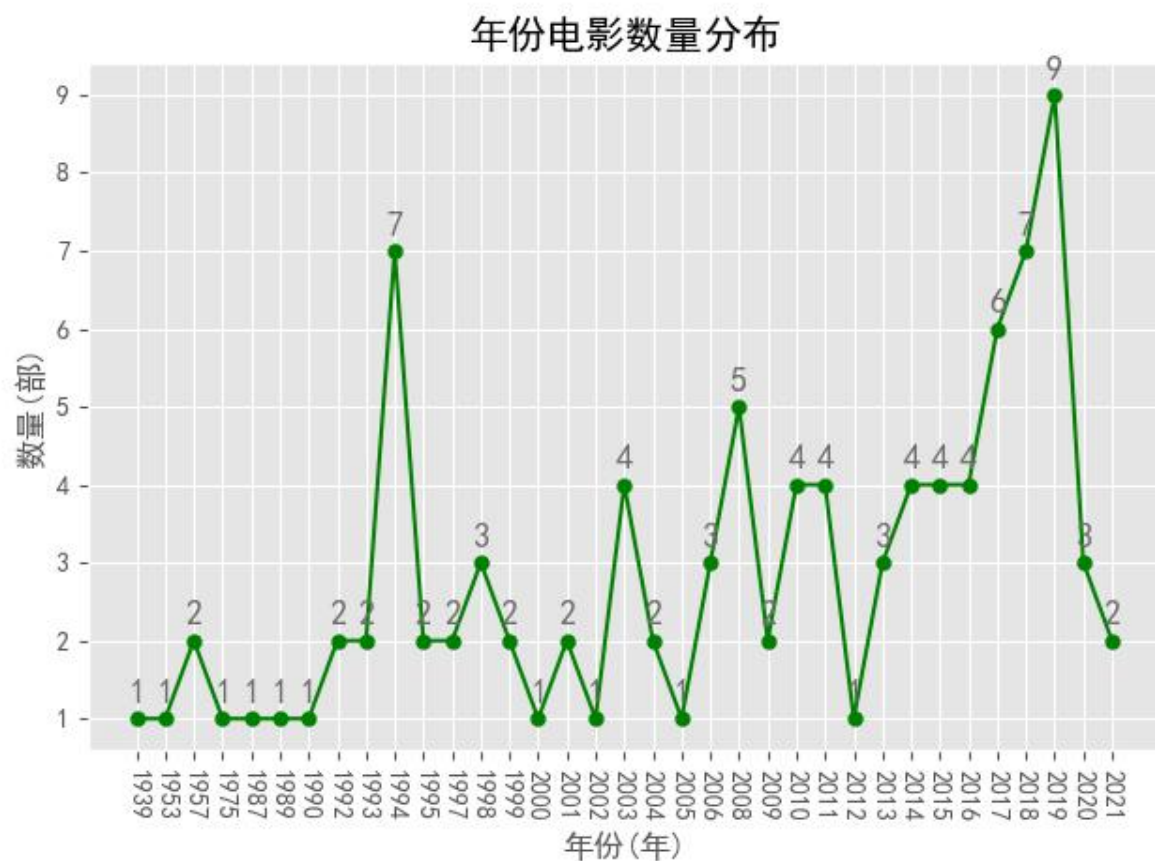
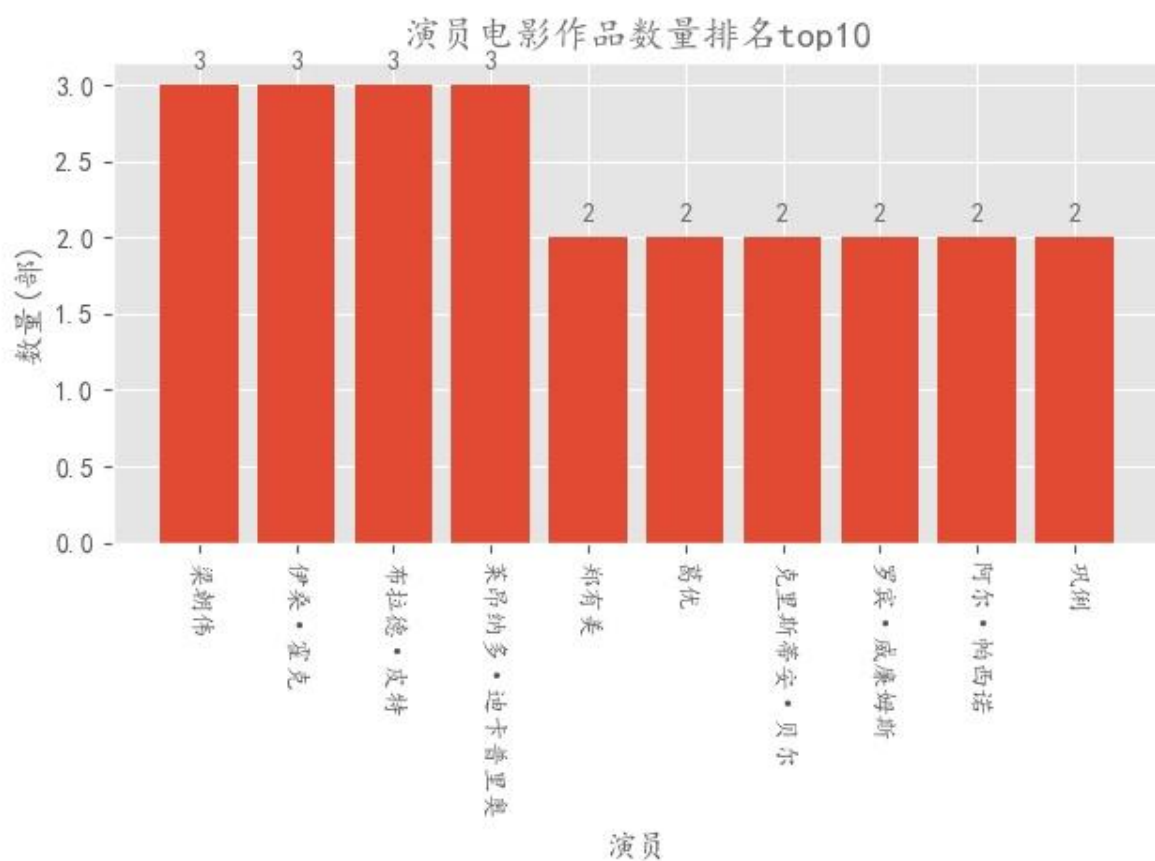


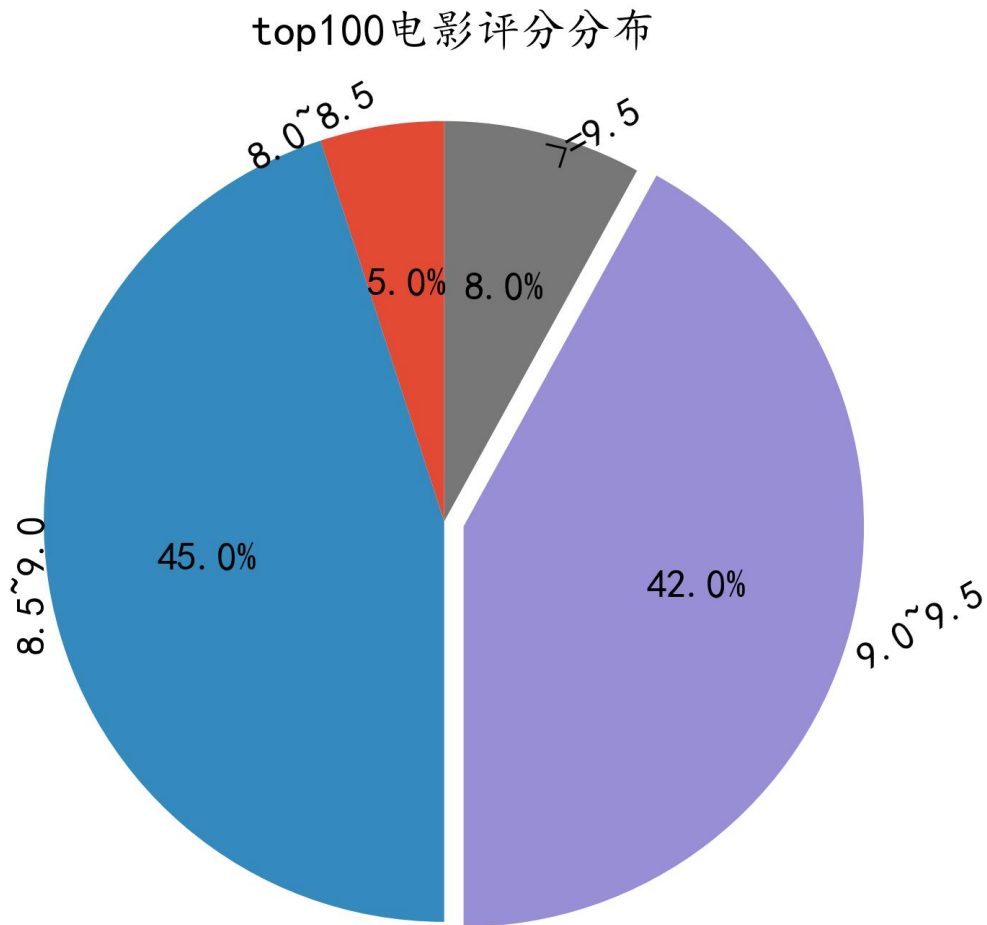
数据 csv 文件



分析图表







## 六、参考资料

- [1] 朝乐门 著.Python 编程从数据分析到数据科学.电子工业出版社.2019.1
- [2] [美] 阿曼多·凡丹戈 (Armando Fandango) 著, 韩波 译.Python 数据分析(第 2 版). 人民邮电出版社.2018.6
- [3] 嵩天, 礼欣, 黄天羽 著.Python 语言程序设计基础 (第 2 版). 高等教育出版社.2017.2