

用 Python 来了解一下《安家》

原创 野客 Python技术 1周前

如果要选一部近期最火的电视剧，一定非《安家》莫属，你可能没有具体看过，但如果你看微博的话一定听过这个名字，这部电视剧多次登上微博热搜榜，好像还有几次冲上了热搜榜首，该剧主要讲述的是关于房产中介卖房的故事，电视剧原名也是叫卖房子的人。

使用 Python 分析这部电视剧，主要包括两个步骤：获取数据和分析数据，数据来源我们选取《安家》的豆瓣评论区数据。

获取数据

豆瓣中《安家》的地址是：<https://movie.douban.com/subject/30482003/>，我们打开看一下，如下图所示：

从图中我们可以直观的看出截止目前有 9 万多人进行了打分，从评分上来看，打三星和四星的人数居多，总体评分 6.2 属于及格分，算是中规中矩吧。

我们把页面向下拉到评论区位置，如下图所示：

我们可以看到目前有 3 万多条评论数据，豆瓣对查看评论数据的限制是：未登录时最多可以查看 200 条，登录用户最多可以查看 500 条，也就是说我们最多可以抓取 500 条评论相关的信息，我们要抓取的数据项包括：用户昵称、星级、评论时间、评论内容，将这些信息抓取后我们再将它们存到 csv 文件，代码实现如下：

```
1 import requests, time, random, pandas as pd
2 from lxml import etree
3
4 def spider():
5     url = 'https://accounts.douban.com/j/mobile/login/basic'
6     headers = {"User-Agent": 'Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.0; Trident/4.0)'}
7     # 安家评论网址，为了动态翻页，start 后加了格式化数字，短评页面有 20 条数据，每页增加 20 条
8     url_comment = 'https://movie.douban.com/subject/30482003/comments?start=%d&limit=20&sort=new_score'
9     data = {
10         'ck': '',
11         'name': '自己的用户',
12         'password': '自己的密码',
13         'remember': 'false',
14         'ticket': ''
15     }
16     session = requests.session()
17     session.post(url=url, headers=headers, data=data)
18     # 初始化 4 个 list 分别存用户名、评星、时间、评论文字
19     users = []
20     stars = []
21     times = []
22     content = []
23     # 抓取 500 条，每页 20 条，这也是豆瓣给的上限
24     for i in range(0, 500, 20):
```

```

24     # 获取 HTML
25     data = session.get(url_comment % i, headers=headers)
26     # 状态码 200 表示成功
27     print('第', i, '页', '状态码: ', data.status_code)
28     # 暂停 0-1 秒时间, 防止 IP 被封
29     time.sleep(random.random())
30     # 解析 HTML
31     selector = etree.HTML(data.text)
32     # 用 xpath 获取单页所有评论
33     comments = selector.xpath('//div[@class="comment"]')
34     # 遍历所有评论, 获取详细信息
35     for comment in comments:
36         # 获取用户名
37         user = comment.xpath('./h3/span[2]/a/text()')[0]
38         # 获取评星
39         star = comment.xpath('./h3/span[2]/span[2]/@class')[0][7:8]
40         # 获取时间
41         date_time = comment.xpath('./h3/span[2]/span[3]/@title')
42         # 有的时间为空, 需要判断下
43         if len(date_time) != 0:
44             date_time = date_time[0]
45         else:
46             date_time = None
47         # 获取评论文字
48         comment_text = comment.xpath('./p/span/text()')[0].strip()
49         # 添加所有信息到列表
50         users.append(user)
51         stars.append(star)
52         times.append(date_time)
53         content.append(comment_text)
54     # 用字典包装
55     comment_dic = {'user': users, 'star': stars, 'time': times, 'comments': content}
56     # 转换成 DataFrame 格式
57     comment_df = pd.DataFrame(comment_dic)
58     # 保存数据
59     comment_df.to_csv('data.csv')
60     # 将评论单独再保存下来
61     comment_df['comments'].to_csv('comment.csv', index=False)
62

```

分析数据

现在数据取到了, 我们使用 Python 来对这些数据进行分析一下。

评论数量

首先, 我们来统计一下这 500 条数据每天的评论数量, 然后利用折线图进行数据展示, 代码实现如下:

```

1 import pandas as pd, matplotlib.pyplot as plt

```

```

2
3 csv_data = pd.read_csv('data.csv')
4 df = pd.DataFrame(csv_data)
5 df_gp = df.groupby(['time']).size()
6 values = df_gp.values.tolist()
7 index = df_gp.index.tolist()
8 # 设置画布大小
9 plt.figure(figsize=(10, 6))
10 # 数据
11 plt.plot(index, values, label='评论数')
12 # 设置数字标签
13 for a, b in zip(index, values):
14     plt.text(a, b, b, ha='center', va='bottom', fontsize=13, color='black')
15 plt.title('评论数随时间变化折线图')
16 plt.xticks(rotation=330)
17 plt.tick_params(labelsize=10)
18 plt.ylim(0, 200)
19 plt.legend(loc='upper right')
20 plt.show()

```

看一下效果图：

从图中我们可以看出 2 月 21、22 这两天评论数最多，其中 2 月 21 号为开播日，评论数较多很正常，2 月 22 号评论数多于开播日，我们大致可以推测是开播后网络等渠道进一步扩散的因素，之后随着时间的推移热度有所下降，评论数量呈下降至相对平稳的趋势。

角色分析

我们接着统计评论区中几个主要角色被提及的次数，然后再利用柱状图进行数据展示，代码实现如下所示：

```

1 import pandas as pd, jieba, matplotlib.pyplot as plt
2
3 csv_data = pd.read_csv('data.csv')
4 roles = {'姑姑':0, '房似锦':0, '王子':0, '闪闪':0, '老油条':0, '楼山关':0, '鱼化龙':0}
5 names = list(roles.keys())
6 for name in names:
7     jieba.add_word(name)
8 for row in csv_data['comments']:
9     row = str(row)
10    for name in names:
11        count = row.count(name)
12        roles[name] += count
13 plt.figure(figsize=(8, 5))
14 # 数据
15 plt.bar(list(roles.keys()), list(roles.values()), width=0.5, label='提及次数', color=['g', 'r', 'dodg
16 # 设置数字标签
17 for a, b in zip(list(roles.keys()), list(roles.values())):
18     plt.text(a, b, b, ha='center', va='bottom', fontsize=13, color='black')
19 plt.title('角色被提及次数柱状图')
20 plt.xticks(rotation=270)

```

```

plt.xticks(rotation=270)
20 plt.tick_params(labelsize=10)
21 plt.ylim(0, 30)
22 plt.legend(loc='upper right')
23 plt.show()
24

```

看一下效果图：

我们从角色被提及的次数可以大致推测出角色的受欢迎程度。

星级变化

我们接着根据获取数据来看一下这几天星级变化的大致趋势，一天中如果有多条星评我们取其平均值即可，代码实现如下所示：

```

1 import pandas as pd, numpy as np, matplotlib.pyplot as plt
2
3 csv_data = pd.read_csv('data.csv')
4 df_time = csv_data.groupby(['time']).size()
5 df_star = csv_data.groupby(['star']).size()
6 index = df_time.index.tolist()
7 value = [0] * len(index)
8 # 生成字典
9 dic = dict(zip(index, value))
10 for k, v in dic.items():
11     stars = csv_data.loc[csv_data['time'] == str(k), 'star']
12     # 平均值
13     avg = np.mean(list(map(int, stars.values.tolist())))
14     dic[k] = round(avg, 2)
15 # 设置画布大小
16 plt.figure(figsize=(9, 6))
17 # 数据
18 plt.plot(list(dic.keys()), list(dic.values()), label='星级')
19 plt.title('星级随时间变化折线图')
20 plt.xticks(rotation=330)
21 plt.tick_params(labelsize=10)
22 plt.ylim(0, 5)
23 plt.legend(loc='upper right')
24 plt.show()

```

看一下效果图：

从现有数据来看，《安家》的星级整体维持在 2 星左右，我们可以发现尽管该剧比较热，但观众对该剧的满意并不是很高。

词云展示

最后，我们对所有评论进行词云效果展示，这样可以让我们更加直观的看出评论区哪些词汇出现的频率较高，实现代码

如下所示:

```
1 from wordcloud import WordCloud
2 import numpy as np, jieba
3 from PIL import Image
4
5 def jieba_():
6     # 打开评论数据文件
7     content = open('comment.csv', 'rb').read()
8     # jieba 分词
9     word_list = jieba.cut(content)
10    words = []
11    # 过滤掉的词
12    remove_words = ['以及', '不会', '一些', '那个', '只有',
13                    '不过', '东西', '这个', '所有', '这么',
14                    '但是', '全片', '一点', '一部', '一个',
15                    '什么', '虽然', '一切', '样子', '一样',
16                    '只能', '不是', '一种', '这个', '为了']
17    for word in word_list:
18        if word not in remove_words:
19            words.append(word)
20    global word_cloud
21    # 用逗号隔开词语
22    word_cloud = ', '.join(words)
23
24 def cloud():
25     # 打开词云背景图
26     cloud_mask = np.array(Image.open('bg.jpg'))
27     # 定义词云的一些属性
28     wc = WordCloud(
29         # 背景图分割颜色为白色
30         background_color='white',
31         # 背景图样
32         mask=cloud_mask,
33         # 显示最大词数
34         max_words=100,
35         # 显示中文
36         font_path='./fonts/simhei.ttf',
37         # 最大尺寸
38         max_font_size=80
39     )
40    global word_cloud
41    # 词云函数
42    x = wc.generate(word_cloud)
43    # 生成词云图片
44    image = x.to_image()
45    # 展示词云图片
46    image.show()
47    # 保存词云图片
48    wc.to_file('anjia.png')
```

```
49
50 jieba_()
51 cloud()
```

看一下效果图：

总结

本文通过爬取豆瓣中《安家》的评论区数据并对其进行可视化，我们可以大致了解观众对《安家》这部电视大致的评价情况，当然因为我们所获取的样本数量有限，可能或多或少还会与用户实际的评价情况有一点偏差。

【代码获取方式】

识别文末二维码，回复：666

PS： 公号内回复「Python」即可进入 Python 新手学习交流群，一起 **100天计划！**

-END-

Python 技术
关于 Python 都在这里

文章已于2020-03-25修改