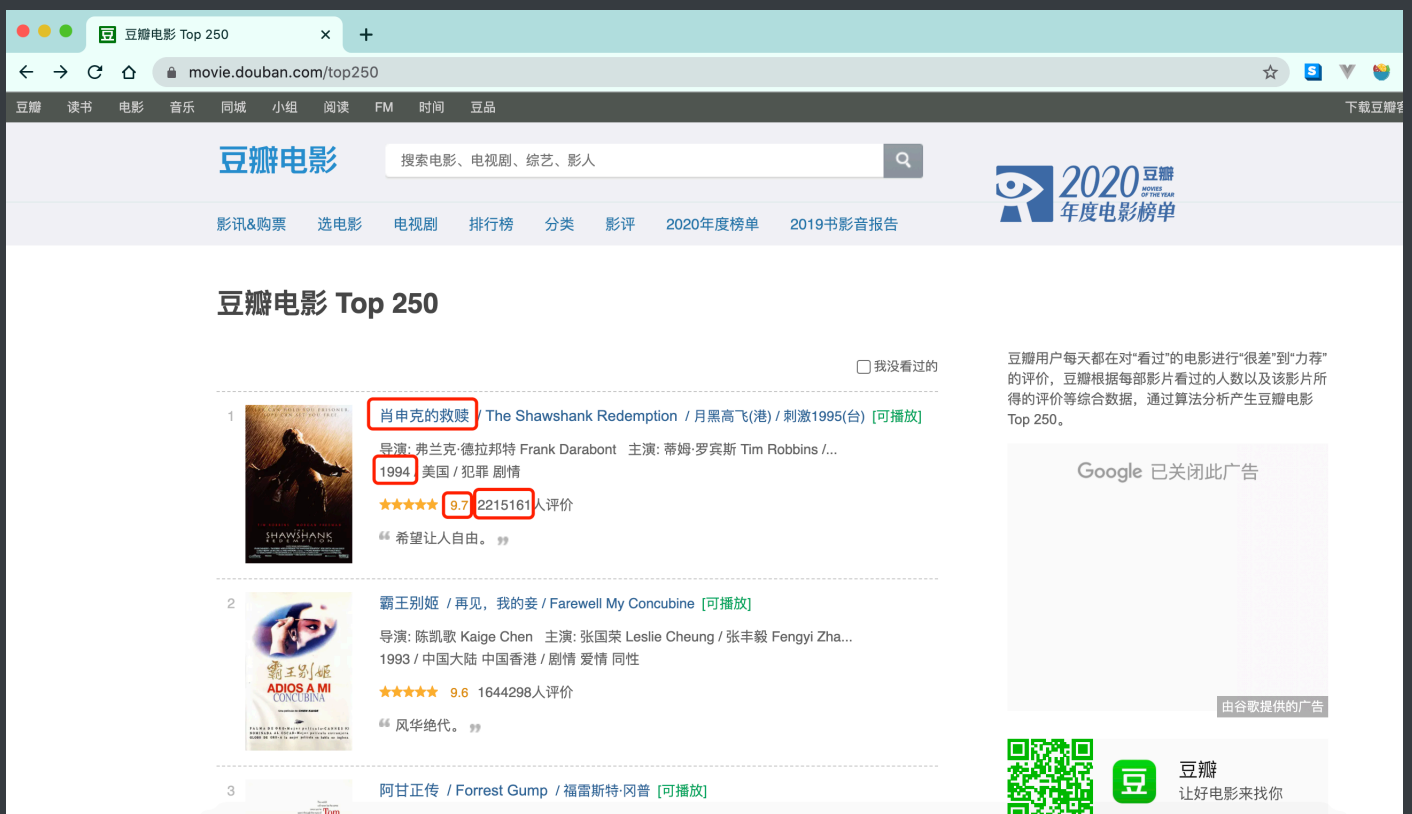


手刃豆瓣TOP250电影信息

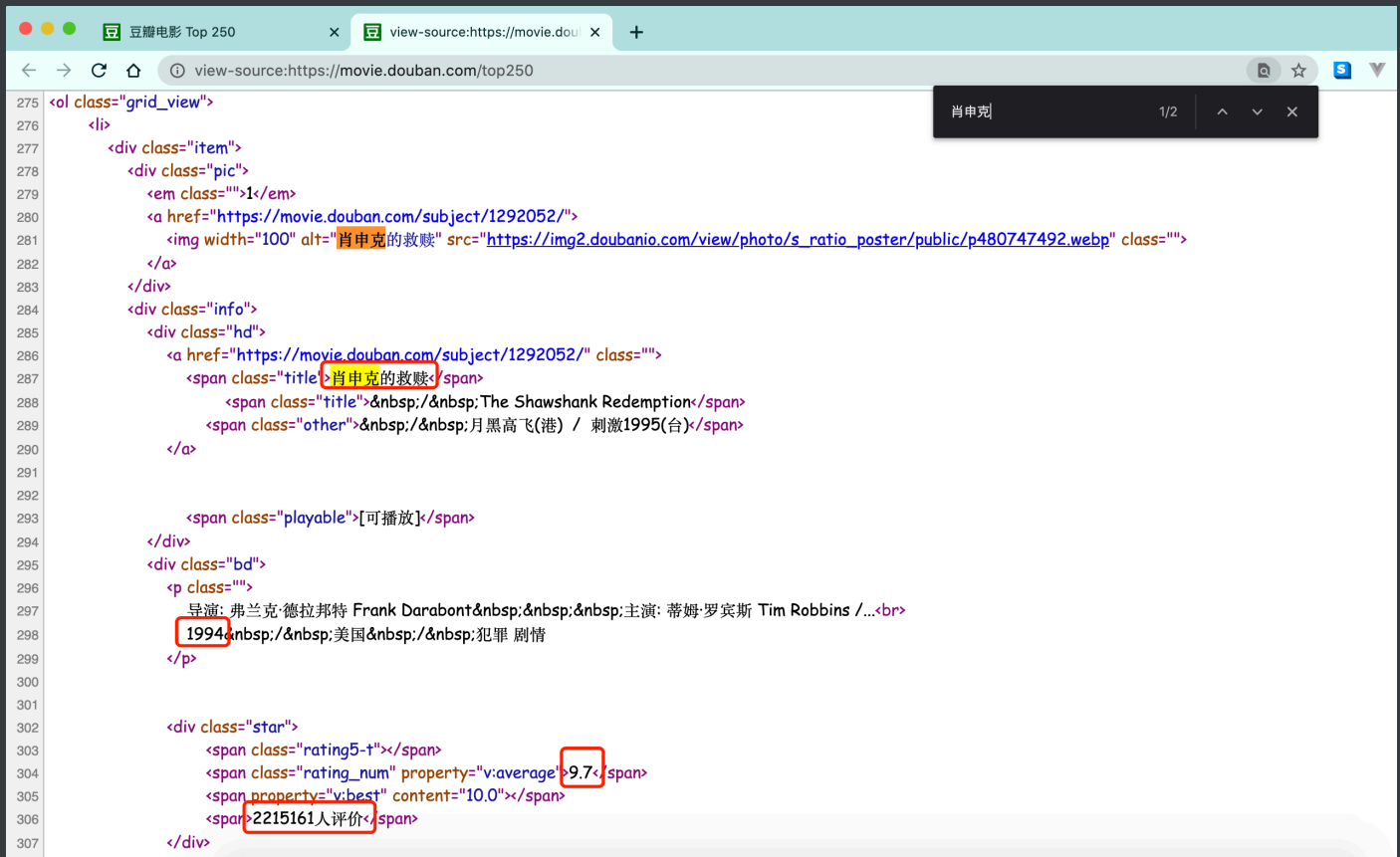
终于可以放开手脚干一番事业了. 今天我们的目标是豆瓣电影TOP250排行榜. 没别的意思, 练手而已

先看需求:



目标: 抓取"电影名称","上映年份","评分","评分人数"四项内容.

怎么做呢? 首先, 先看一下页面源代码. 数据是否是直接怼在源代码上的?



很明显, 我们想要的`数据`全部都在`页面源代码`中体现了. 所以, 我们不需要考虑`js动态加载数据`的情况了. 那么接下来就是编写爬虫代码的第一步了. 拿到`页面源代码`:

```
1 import requests
2
3 headers = {
4     "user-agent": "Mozilla/5.0 (Macintosh; Intel
    Mac OS X 10_15_4) AppleWebKit/537.36 (KHTML, like
    Gecko) Chrome/87.0.4280.88 Safari/537.36"
5 }
6
7 url = "https://movie.douban.com/top250?
    start=0&filter="
8 resp = requests.get(url, headers=headers)
9 print(resp.text)
10
```

然后呢. 从页面源代码中提取我们需要的内容. 这时候我们就可以去写正则了.

```

1 obj = re.compile(r'<li>.*?<div class="item">.*?
  <div class="pic">.*?<em class="">(P<num>\d+)
  </em>'
2
  r'.*?<span class="title">(P
  <name>.*?)</span>'
3
  r'.*?<p class="">.*?<br>\n(P
  <year>.*?)&nbsp;'
4
  r'.*?property="v:average">(P
  <average>.*?)</span>'
5
  r'.*?<span>(P<people>\d+)人评价
  </span>', re.S)
6

```

开始匹配, 将最终完整的数据按照自己喜欢(需要)的方式写入文件.

```

1 it = obj.finditer(resp.text)
2 with open("movie.csv", mode="w", encoding="utf-8")
  as f:
3     csvwriter = csv.writer(f) # 创建csv文件写入工具,
  也可以直接f.write()
4     for item in it:
5         dic = item.groupdict()
6         dic['year'] = dic['year'].strip()
7         csvwriter.writerow(dic.values()) # 写入数
  据

```

代码还有优化空间, 各位可以思考一下如何进一步对代码优化(时间复杂度, 空间复杂度), 各位可以自行想办法将豆瓣TOP250条数据全部抓取到. 我这里就偷工减料了. 嘿嘿~