

# 专栏：009：高评分电影抓取

## 用理工科思维看待这个世界

### 系列爬虫专栏

崇尚的学习思维是：输入，输出平衡，且平衡点不断攀升。

曾经有大神告诫说：没事别瞎写文章；所以，很认真的写的是能力范围内的，看客要是看不懂，不是你的问题，问题在我，得持续输入，再输出。

今天的主题是：实战爬取电影，并存储至MySQL数据库

## 1：框架

序号	目标	说明
01	抓取目标分析	-目标是什么
02	分解任务	—
03	MySQL建表操作	本地建表
04	实战抓取	—
05	参考及总结	—

## 2：目标

任务是：抓取网站数据，存放至MySQL数据库中。

效果显示：抓取的目标存放至本地MySQL数据库

1 信息2 表数据3 信息

限制行 第一行: 0 行数: 1000

Film	Director	R.	Number	Url	Describe
肖申克的救赎	导演: 弗兰克·德拉邦特Frank Darabont主演: 蒂姆·罗宾斯Tim Robbins / 摩根·弗里曼Morgan Freeman	9.6	691626人评价	https://movie.douban.com/subject/1292052/	希望让人自由。
美丽人生	导演: 罗伯托·贝尼尼Roberto Benigni主演: 罗伯托·贝尼尼Roberto Benigni / 尼可莱塔·布拉斯奇Nicola Pizzani	9.5	327128人评价	https://movie.douban.com/subject/1292063/	最美的谎言。
控方证人	导演: 比利·怀尔德Billy Wilder主演: 查尔斯·劳顿Charles Laughton / 玛琳·奥哈拉Maureen O'Hara	9.5	42603人评价	https://movie.douban.com/subject/1296141/	比利·怀尔德满分作品。
这个杀手不太冷	导演: 吕克·贝松Luc Besson主演: 让·雷诺Jean Reno / 娜塔丽·波特曼Natalie Portman	9.4	661379人评价	https://movie.douban.com/subject/1295644/	猛男泰和小萝莉不得不让
阿甘正传	导演: Robert Zemeckis主演: Tom Hanks / Rob	9.4	579883人评价	https://movie.douban.com/subject/1292720/	一部美国近现代史。
霸王别姬	导演: 陈凯歌Kaige Chen主演: 张国荣Leslie Cheun / 巩俐Li Gongli	9.4	477559人评价	https://movie.douban.com/subject/1291546/	风华绝代。
辛德勒的名单	导演: 史蒂文·斯皮尔伯格Steven Spielberg主演: 连姆·尼森Liam Neeson / 伊迪娜·梅泽Idina Menzel	9.4	306330人评价	https://movie.douban.com/subject/1295124/	拯救一个人，就是拯救整
机器人总动员	导演: 安德鲁·斯坦顿Andrew Stanton主演: 本·贝尔特Ben Burtt / 艾迪·墨菲Eddie Murphy	9.3	421133人评价	https://movie.douban.com/subject/2131459/	小瓦力，大人生。
十二怒汉	导演: Sidney Lumet主演: 亨利·方达Henry Fonda / 李·科布Lee Remick	9.3	134597人评价	https://movie.douban.com/subject/1293182/	1957年的理想主义。
海豚湾	导演: Louie Psihoyos主演: John Chisholm /	9.3	159096人评价	https://movie.douban.com/subject/3442220/	海豚的微笑，是世界上第
千与千寻	导演: 宫崎骏Hayao Miyazaki主演: 柊瑠美Rumi Hir	9.2	524533人评价	https://movie.douban.com/subject/1291561/	最好的宫崎骏，最好的少
海上钢琴师	导演: 朱塞佩·托纳多雷Giuseppe Tornatore主演: 蒂姆·罗斯Tim Roth / 比尔·默瑞Bill Murray	9.2	500680人评价	https://movie.douban.com/subject/1292001/	每个人都要走一条自己选
盗梦空间	导演: 克里斯托弗·诺兰Christopher Nolan主演: 莱昂纳多·迪卡普里奥Leonardo DiCaprio / 约瑟夫·高登-莱维特Joseph Gordon-Levitt	9.2	641152人评价	https://movie.douban.com/subject/3541415/	诺兰给了我们一场无法逃
放牛班的春天	导演: 克里斯托夫·巴莱特Christophe Barratier主	9.2	369983人评价	https://movie.douban.com/subject/1291549/	天籁一般的童声，是最挂
忠犬八公的故事	导演: 莱塞·霍尔斯道姆Lasse Hallström主演: 理查·基尔Richard Gere / 琼·艾伦Joan Allen	9.2	351890人评价	https://movie.douban.com/subject/3011091/	永远都不能忘记你所爱的
教父	导演: 弗朗西斯·福特·科波拉Francis Ford Coppola	9.2	280241人评价	https://movie.douban.com/subject/1291841/	千万不要记恨你的对手。
乱世佳人	导演: Victor Fleming / George Cukor主演: 克	9.2	225730人评价	https://movie.douban.com/subject/1300267/	Tomorrow is another d
大闹天宫	导演: 万籁天Laiming Wan / 唐澄Cheng Tang主演: 邱	9.2	74724人评价	https://movie.douban.com/subject/1418019/	经典之作，历久弥新。
城市之光	导演: Charles Chaplin主演: 查理·卓别林Charle	9.2	31022人评价	https://movie.douban.com/subject/1293908/	永远的小人物，伟大的精
泰坦尼克号	导演: James Cameron主演: Leonardo DiCapr	9.1	534587人评价	https://movie.douban.com/subject/1292722/	失去的才是永恒的。
三傻大闹宝莱坞	导演: 拉库马·希拉尼Rajkumar Hirani主演: 阿米尔·汗Amir Khan / 安努舒卡·谢蒂Anushka Sharma	9.1	549230人评价	https://movie.douban.com/subject/3793023/	英俊版憨豆，高情商版语
龙猫	导演: 宫崎骏Hayao Miyazaki主演: 日高法子Moriko	9.1	343707人评价	https://movie.douban.com/subject/1291560/	人人心中都有个龙猫，重

初始URL

url = https://movie.douban.com/top250

字段：

Film：电影名称

Director: 电影导演

Rates：评分数

Numbers：评分人数

Url: 电影链接

Describe: 电影介绍 (网站的一句话，经典台词之类的)

### 3：任务分解

具体点击网页审查元素：[链接](#)

- 字段的正则表达式

电影名称：

Film\_pattern = r'<span class="title">(.\*?)</span>'

电影导演：

Director\_pattern = r'<p class="">(.\*?)</p>'

评分数：

Rates\_pattern = r'<span class="rating\_num" property="v:average">(.\*?)</span>'

评分人数：先抓大，再在大的里面匹配所需的文本信息

Number\_pattern\_large = r'<div class="star">(.\*?)</div>'

Number\_pattern\_small = r'<span>(.\*?)</span>'

电影链接：先抓大，再在大的里面匹配所需的文本信息

Urlfilm\_pattern\_large = r'<div class="hd">(.\*?)</div>'

Urlfilm\_pattern\_small = r'<a href="(.\*?)>'

电影介绍：

Describe\_pattern = r'<span class="inq">(.\*?)</span>'

- URL的分析：

由翻页，自己匹配出网址：(也可以网址抓取)。10页。

```
urls = ["https://movie.douban.com/top250?start={}&filter=".format(i) for i
in range(0,250,25)]
```

网址分析完成，正则分析完成。任务完成了大半。

---

## 4：数据库建表操作

在本地：数据库名为：exercise 下创建一个表名为：douban\_film

建表的SQL语法：参照[w3school](#)

```
CREATE TABLE `douban_film` (
  `Film` CHAR(32) DEFAULT NULL COMMENT '电影名称',
  `Director` CHAR(32) DEFAULT NULL COMMENT '电影导演',
  `Rates` FLOAT DEFAULT NULL COMMENT '评分数',
  `Number` CHAR(16) DEFAULT NULL COMMENT '评分人数',
  `Url` CHAR(128) DEFAULT NULL COMMENT '评分人数',
  `Describe` CHAR(128) DEFAULT NULL COMMENT '电影介绍'
) ENGINE=INNODB DEFAULT CHARSET=utf8 COMMENT='豆瓣电影250介绍'
`douban_film`
# 执行sql语句就可以创建一个表，各字段及其属性如上示
```

---

## 5：实战抓取

单独使用正则，会出现很多难以匹配(可能没有尝试其他匹配规则)。需要对网页进行不断的分析。

抓取核心代码：(大神轻拍代码...)

# 网页抓取字段示例

```
def content_json(self, content):
    Film_all = re.findall(self.Film_pattern, content, re.S)
    Film = []
    for one_film in Film_all:
        if "&nbsp;" not in one_film:
            Film.append(one_film)
    Director_all = re.findall(self.Director_pattern, content, re.S)
    Director = []
    for one_Director in Director_all:
        one = self.str_replace(one_Director)
        Director.append(one)
    Rates = re.findall(self.Rates_pattern, content, re.S)
    Number_large = re.findall(self.Number_pattern_large, content, re.S)
    Number = []
    for one_number in Number_large:
        Number_one = re.findall(self.Number_pattern_small, one_number, re.
S)[0]
        Number.append(Number_one)
    Describe = re.findall(self.Describe_pattern, content, re.S)
    Url_large = re.findall(self.Urlfilm_pattern_large, content, re.S)
    Url = []
    for one_url in Url_large:
        Url_one = re.findall(self.Urlfilm_pattern_small, one_url, re.S)[0]
        Url.append(Url_one)
    Film_collection = []
    for Film, Director, Rates, Number, Describe, Url in zip(Film, Directo
r, Rates, Number, Describe, Url):
        data = {
            "Film": Film,
            "Director": Director,
            "Rates": Rates,
            "Number": Number,
            "Describe": Describe,
            "Url": Url
        }
        Film_collection.append(data)
    return Film_collection
```

# 文本匹配会有很多不需要的字段，如下函数实现数据清洗

```
def str_replace(self, str_one):
    str_one_1 = str_one.replace("\n", '')
    str_one_2 = str_one_1.replace("<br>", '')
    str_one_3 = str_one_2.replace("&nbsp;", '')
    str_one_4 = str_one_3.replace("\t", '')
    str_one_5 = str_one_4.replace(" ", '').strip()
    return str_one_5
```

# 单独写sql语句比较繁琐，如下函数实现JSON格式数据转换成SQL语句

```
import copy
def json_to_mysql(json_obj, table, sql_type="insert"):
    local_copy = copy.deepcopy(json_obj)
    if sql_type == "insert":
        sql_part1 = "insert into " + table
        keys = local_copy.keys()

        sql_part2 = "("
        for key in keys:
            sql_part2 += "`%s`"%(key)
            sql_part2 += ","
        sql_part2 = sql_part2.rstrip(",")
        sql_part2 += ")"

        sql_part3 = "("
        for key in keys:
            sql_part3 += "'" + (local_copy[key]) + "'"
            sql_part3 += ","
        sql_part3 = sql_part3.rstrip(",")
        sql_part3 += ")"

        sql = sql_part1 + " " + sql_part2 + " values " + sql_part3
    return sql
```

核心代码已经完成：

整体思路如下：

- 分析首页文本信息的正则表达式
- 抓取首页的字段
- 对字段进行数据的清洗，去掉不需要的信息
- 将数据结构化
- 循环操作
- 获取的全部信息执行sql语句，存入已经建表的MySQL数据库中

完整版代码：[完整版代码](#)

另一款数据库可视化工具显示效果：

