

第一次作业

2019年9月26日 10:37

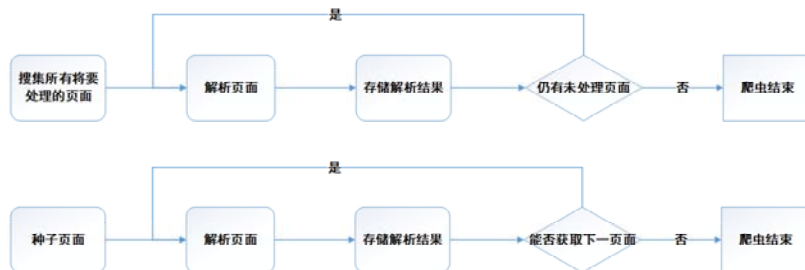
必做1.

- 对于时间精力有限的小白来说，想要在最短的时间内搜寻到各类开源爬虫库及软件，最好的办法就是将前人（多指技术大牛）已经写好的总结性文章找出来，然后加以自己的想法，得到最终的结果。
- 根据经验，这些总结性的文章一定存在于**各类优质技术内容生产及分享平台**，比如**github**，**知乎**，**CSDN**，**cnblogs**，其余网站水平层次不齐，且有较高概率是从上述平台转载（抄）过来的。因此在搜索过程中，为了提高搜索的效率及结果的质量，我只收集上述网站的内容，其余网站不看。
- 为了找到这些优质内容所在的具体页面，将关键字输入搜索引擎。这里的关键字使用“**开源 爬虫 库 框架 软件**”，同时为了提高搜索的多样性，同时使用百度和谷歌。

大体的搜索过程如下图：



阅读这些优质内容，我大致了解了爬虫**较为常见**的流程：



但无论是何种流程，无论如何绕不开的有2个阶段：**获取页面**，**解析页面**。

- 获取页面，首先要获得该页面的URL，然后通过互联网得到该页面的代码，这其中涉及到网络中的各类协议。
- 解析页面，网络获取的页面一般是XML类型的文件，解析页面就是分析其中的字符串，而各类字符串的分析逻辑较为复杂，无论何种高级程序语言处理起来都很棘手。

如果一个人不借助于第三方库，直接从头开始写，那不论是对小白还是技术大牛都是极为繁琐的。因此封装了各类处理方法的库就在这个过程中诞生了，有了这些库，处理上述问题的过程就被大大简化了，而且也无需了解底层的处理逻辑和实现细节，只需要知道输入输出是什么即可。

用于获取页面的库：

urllib、**requests**、**urllib2**、**urllib3**、**grab**、**pycurl**、**httplib2**、**RoboBrowser**

用于解析页面的库：

beautifulsoup、**lxml**、**re**、**cssselect**

为了实现整个爬虫，可以自由组合上面的库，例如

urllib+beautifulsoup、urllib+lxml、urllib+re、requests+re

当然也有提供一整条龙服务的框架，

scrapy、**pyspider**、**nutch**

甚至有专门的软件：

八爪鱼采集器、后羿采集器、火车头

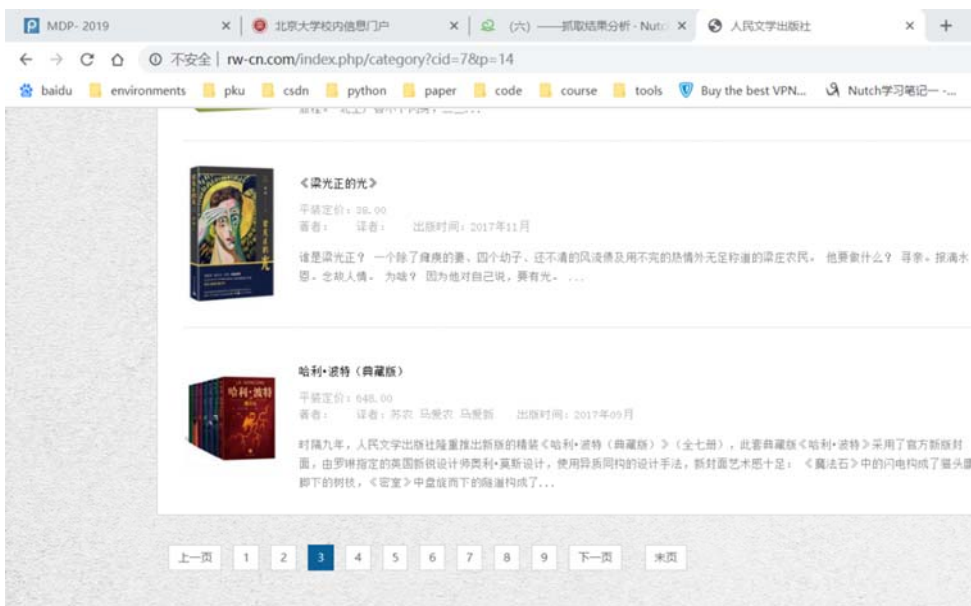
ps：由于爬虫软件对提高编程能力几乎没有帮助，后面不再提及

至于什么是最好用的，任何一种定义式的结论都不具备普遍意义。因此，我觉得要对不同的用户来看。

- 对于企业：企业的目的是赚钱，所以那么采用何种方法，必定要结合企业自身的情况来对爬虫中的各项指标进行综合分析，如**稳定性**、**多线程**、**分布式**、**速度**、**开发成本**等等进行选择。因此对于企业，“最好用”——**企业利益最大化**
- 对于个人，如果是技术大牛，那“最好用”——**功能齐全**；如果是小白，“最好用”——**上手容易**。
- 对于我自己，好不好用得试过了才知道。

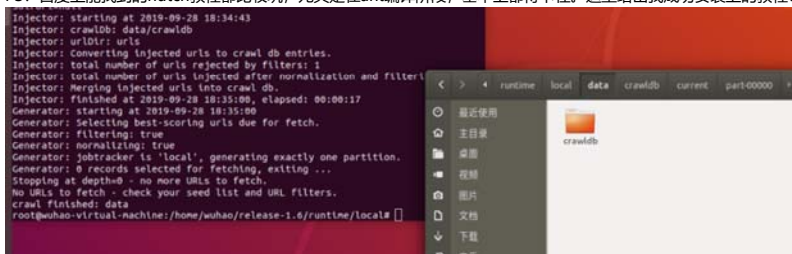
因此尝试使用**requests+beautifulsoup**，**scrapy**，**nutch**分别实现选做中的**人民文学出版社的书籍：要求获取书名、ISBN号码、销售网站及当前价格**

- 分析人民文学出版社的网页，首先要爬取图书列表上的所有图书的详细信息URL，虽然图书列表被分为多页，但是每一页都是静态页面，很容易就能推算出每一页的URL，进而爬取所有的图书详细地址；随后就是循环爬取每一本书的内容。



- nutch的环境安装是真的繁琐，在虚拟机上折腾了一天半，占据了做本次作业一半以上的时间好不容易安装上了，结果按着教程却怎么也出不来结果。结论：nutch虽然功能强大，但是对于我这样仅仅是了解入门爬取网页的人来说，学习成本太高。

PS：百度上能找到的nutch教程都比较坑，尤其是在ant编译阶段，基本上都得卡住。这里给我成功安装上的教程：<https://www.cnblogs.com/huligong1234/p/3464371.html>



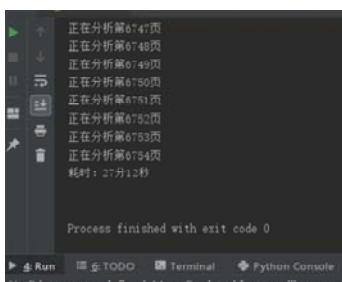
- requests+beautifulsoup:

```

1 #!/usr/bin/env python
2 # coding:utf-8
3 # Author: WuHao
4
5 from bs4 import BeautifulSoup
6 import requests
7 import time
8
9 start_time = time.time()
10 urlmodel = "http://www.rw-cn.com/index.php/category?cid=7&p="
11 urllist = []
12
13 print('正在获取所有图书URL')
14 for i in range(0, 6749, 7):
15     urlcur = urlmodel + str(i)
16     f = requests.get(urlcur) # Get该网页从而获取该html内容
17     soup = BeautifulSoup(f.content, "lxml")
18     for j in soup.find_all('a', class_='a_7 fl'):
19         urllist.append(j['href'])
20
21 fd = open('book.txt', 'w', encoding='utf-8')
22 for i, url in enumerate(urllist):
23     print('正在分析第%d页' % i)
24     f = requests.get(url) # Get该网页从而获取该html内容
25     soup = BeautifulSoup(f.content, "lxml") # 用lxml解析器解析该网页的内容, 好像f.text也是返回的html
26
27     for j in soup.find_all('h2', class_='h_10'):
28         # print(j.text)
29         name=j.text.split()[0]
30         price=j.text.split()[1]
31         for k in soup.find_all('div', class_='div_47 fix'): # , 找到div并且class为p10的标签
32             a = k.find_all('span') # 在每个对应div标签下找span标签, 会发现, 一个a里面有四组span
33             isbn=a[1].text.replace('ISBN:', '')
34             fd.write('书名: ' + name + '\n')
35             fd.write('ISBN号码: ' + isbn + '\n')
36             fd.write('当前价格: ' + price + '\n')
37             fd.write('销售网站: ' + url + '\n')
38             fd.write('-----\n')
39
40 fd.close()
41 end_time=time.time()
42 print("耗时: " + str(int((end_time-start_time)/60)) + '分' + str(int((end_time-start_time)%60)) + '秒\n')

```

for i in range(0, 6749, 7)



book.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

书名: 应物兄
 ISBN号码: 9787020147465
 当前价格: 75.00元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7370>

书名: 牵风记
 ISBN号码: 9787020148240
 当前价格: 43.00元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7367>

书名: 《世间生活》
 ISBN号码: 9787020151967
 当前价格: 元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7388>

书名: 《仰郁生花》
 ISBN号码: 9787020151820
 当前价格: 42元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7387>

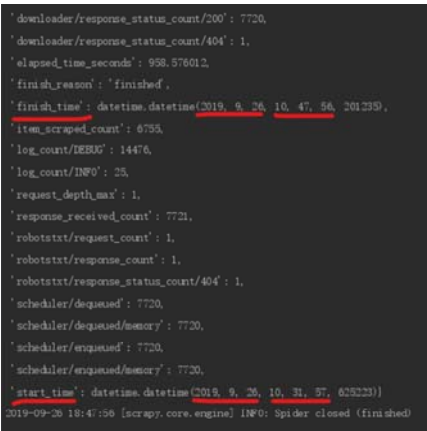
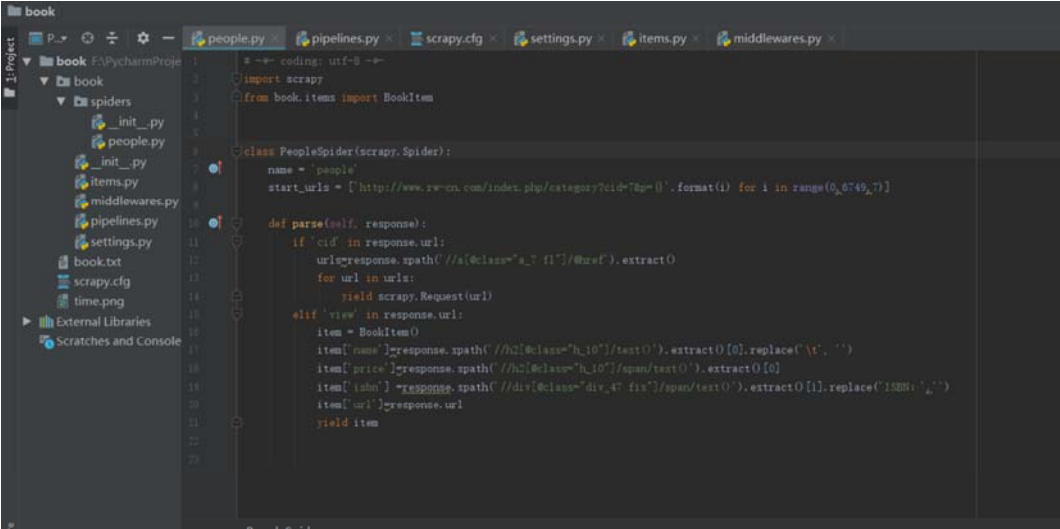
书名: 《辛格自选集》
 ISBN号码: 9787020141210
 当前价格: 89.0元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7386>

书名: 《穗子的动物园》
 ISBN号码: 978-7-02-015136-3
 当前价格: 58元
 销售网站: <http://www.rw-cn.com/index.php/category/view?id=7384>

书名: 《耶稣的学生时代》
 ISBN号码: 9787020149421

结论：灵活，流程逻辑清晰，上手难度低。PS：代码结果：<https://github.com/wuhao9714/crawl/tree/master/requests%2Bbf>

scrapy



耗时：15m59s

book.txt - 记事本

文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

书名：《耶稣的学生时代》

ISBN号码：9787020149421

当前价格：58元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7383>

书名：《燕子最后飞去了哪里》

ISBN号码：9787020122257

当前价格：43.00元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7297>

书名：《彩色面纱》

ISBN号码：

当前价格：¥ 36.00元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7253>

书名：哈利·波特（典藏版）

ISBN号码：9787020127993

当前价格：648.00元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7331>

书名：《吃鲷鱼让我打瞌》

ISBN号码：9787020120963

当前价格：49.00元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7290>

书名：《湮没的时尚·花想容》

ISBN号码：9787020114849

当前价格：元

销售网站：<http://www.rw-cn.com/index.php/category/view?id=7312>

书名：《周涛散文》

ISBN号码：978-7-02-010815-2

结论：上手难度中等，缺点是，url页面的获取顺序是不确定的，每次运行程序得到的爬虫结果均不一样。PS:代码+结果：<https://github.com/wuhao9714/crawl/tree/master/scrapy>

综上，我觉得最好用的是requests+beautifulsoup（或xml）

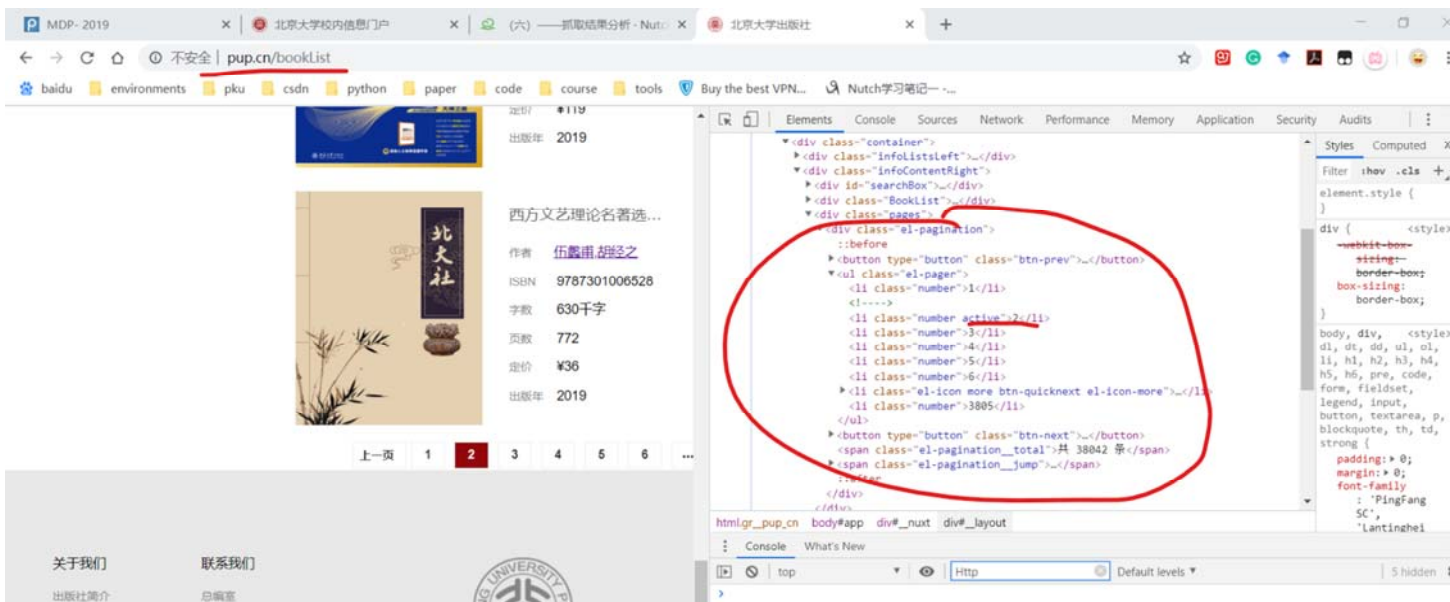
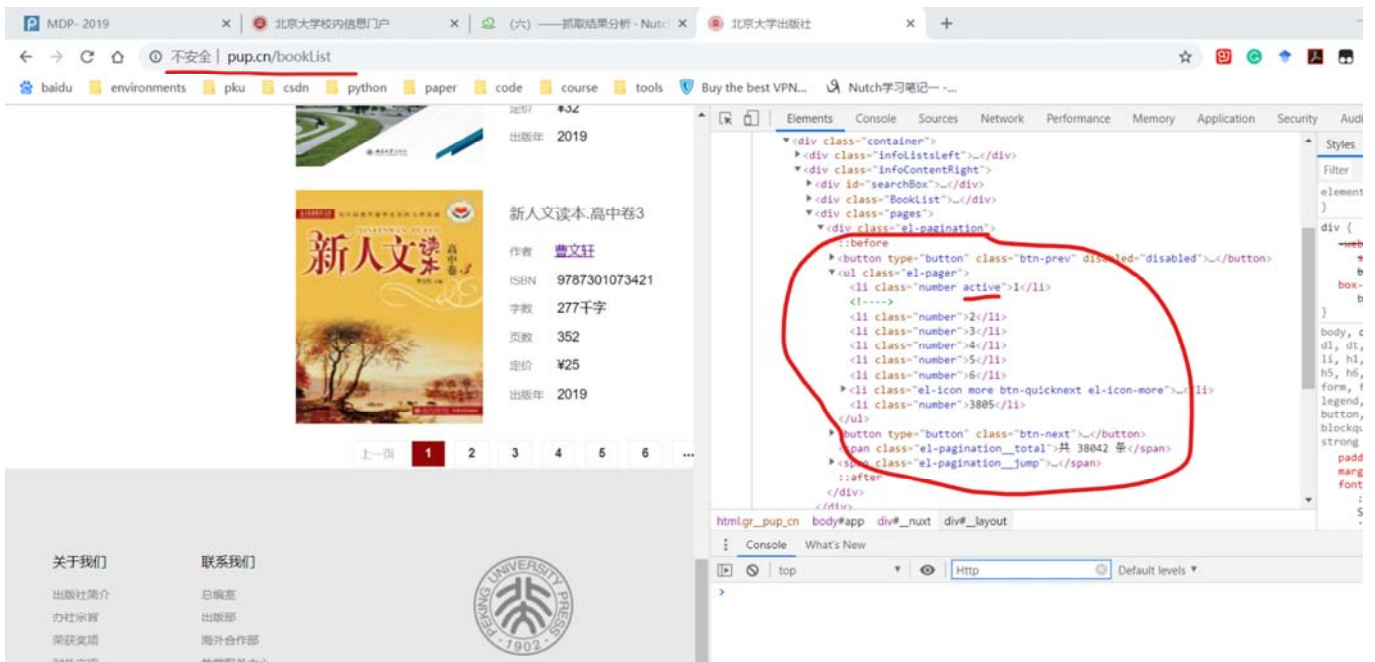
必做二：已阅读

选做一：精力有限，国庆期间完成其余课程作业后如果还有时间再尝试做做

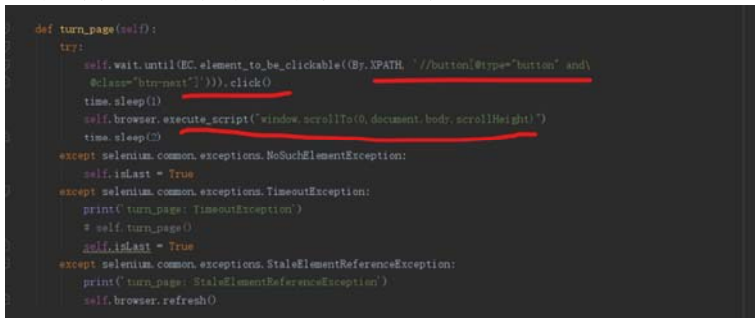
选作二：

1.利用selenium爬取动态网页

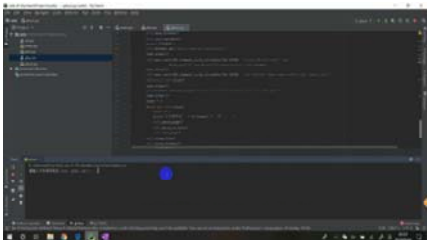
爬取人民文学出版社是比较容易的，毕竟可以通过地址推算出所有的图书列表网址，进而获取每一本图书的详细地址信息。然而北京大学出版社的图书列表是动态的，每一页的地址都是相同的，那么此时继续使用人民文学出版社时的方法就失效了。



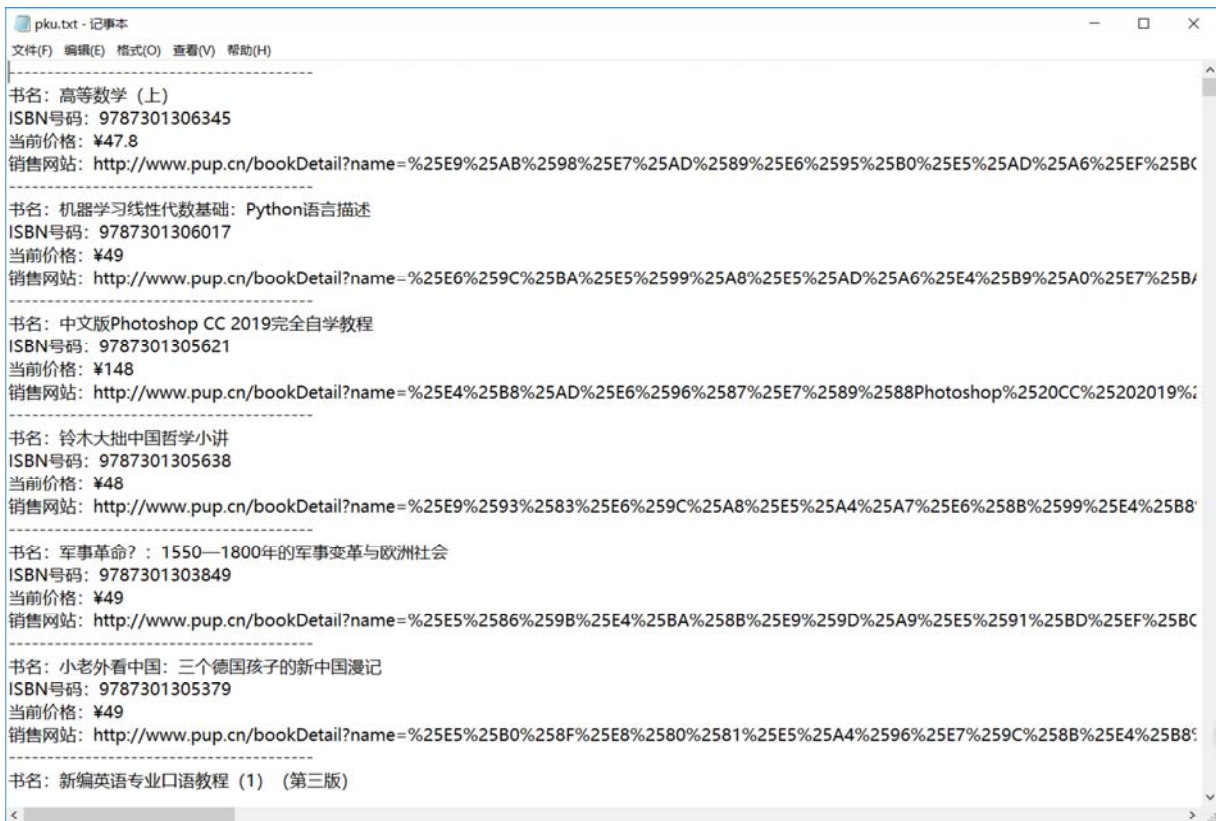
针对动态网页，selenium可以模拟浏览器动作，自动点击下一页按钮。



此处有程序执行时的屏幕录制，里面有模拟浏览器的动作。PS:将录屏视频上传到了B站: <https://www.bilibili.com/video/av69082404>



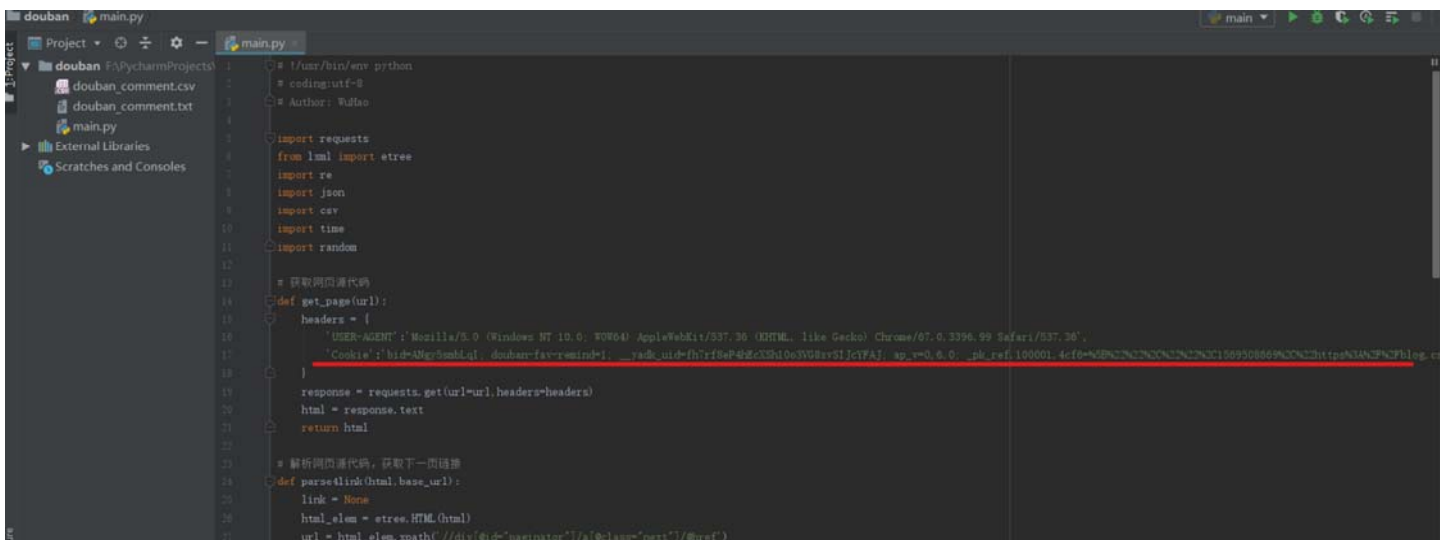
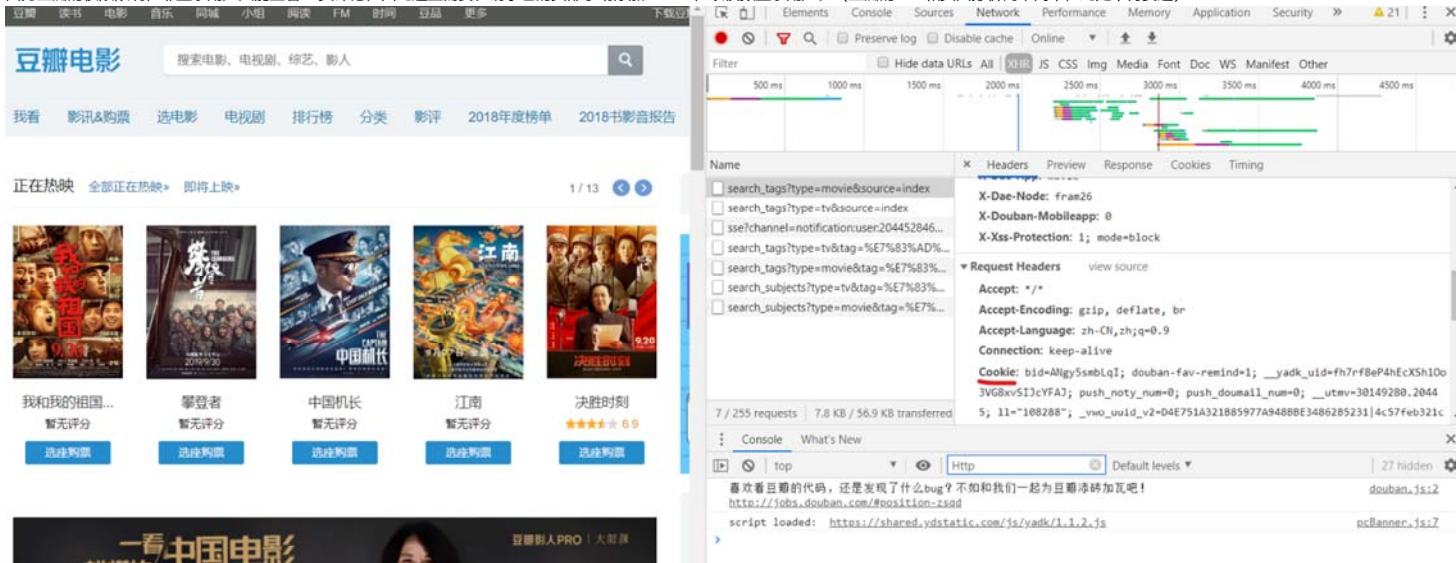
PS:由于图书总量太多，此处只爬取北京大学出版社2019年出版的图书



PS:代码+结果: <https://github.com/wuhao9714/crawl/tree/master/selenium>

2.使用requests+ lxml+ cookie模拟登录来爬豆瓣上《上海堡垒》的所有差评。

因为豆瓣的权限限制,非登录用户只能查看10页评论,因此这里需要在请求包的头部手动添加cookie,以模拟登录用户。(豆瓣的URL和页面分析向来简单,此处不再赘述)



```
douban_comment.txt - 记事本
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)

-----
agree: 19357
author: 郊外的耶稣
star: 较差
content: 头发这么长, 有没有一点军人的样子!

-----
agree: 15167
author: 西楼尘
star: 很差
content: 如果故事不是发生在陆家嘴而是在象牙山, 江洋暗恋林澜就像王长贵喜欢谢大脚, 这背景除了土一点并不违和。与其搞什么德尔塔母舰, 还不如拍个5

-----
agree: 15903
author: 咕叽咕叽
star: 很差
content: 毫无逻辑的剧情, 生硬尴尬的表演, 这是一部很标准的烂片

-----
agree: 14322
author: 淋
star: 很差
content: 科幻片来说, 它基本上是把中国科幻能犯的错误都集齐了; 要说爱情片, 它倒是有连七夕档都不敢上的自知之明

-----
agree: 12743
author: 表姐电影
star: 很差
content: 我们以为《流浪地球》是中国科幻元年的起点, 但这年刚过一半, 《上海堡垒》就给科幻元年提前划上终点了。

-----
agree: 12292
author: 良良
star: 很差
content: 披着言情皮的科幻片, 剧情毫无章法, 鹿晗演技好差。

-----
agree: 10721

< >
```

PS:代码+结果: <https://github.com/wuhao9714/crawl/tree/master/requests%2Blxml%2Bcookie>