

# SPIDER-DAY03

## 1. lxml解析库

### 1.1 安装使用流程

```
1  【1】安装
2      sudo pip3 install lxml
3
4  【2】使用流程
5      2.1》导模块：          from lxml import etree
6      2.2》创建解析对象：    parse_html = etree.HTML(html)
7      2.3》解析对象调用xpath：r_list = parse_html.xpath('xpath表达式')
```

### 1.2 lxml+xpath使用

```
1  【1】基准xpath：匹配所有电影信息的节点对象列表
2      //dl[@class="board-wrapper"]/dd
3      [<element dd at xxx>,<element dd at xxx>,...]
4
5  【2】遍历对象列表，依次获取每个电影信息
6      item = {}
7      for dd in dd_list:
8          item['name'] = dd.xpath('..//p[@class="name"]/a/text()').strip()
9          item['star'] = dd.xpath('..//p[@class="star"]/text()').strip()
10         item['time'] = dd.xpath('..//p[@class="releasetime"]/text()').strip()
```

## 2. 豆瓣图书爬虫

### 2.1 需求分析

```
1 【1】 抓取目标 - 豆瓣图书top250的图书信息
2     https://book.douban.com/top250?start=0
3     https://book.douban.com/top250?start=25
4     https://book.douban.com/top250?start=50
5     ... ...
6
7 【2】 抓取数据
8     2.1) 书籍名称 : 红楼梦
9     2.2) 书籍描述 : [清] 曹雪芹 著 / 人民文学出版社 / 1996-12 / 59.70元
10    2.3) 书籍评分 : 9.6
11    2.4) 评价人数 : 286382人评价
12    2.5) 书籍类型 : 都云作者痴, 谁解其中味?
```

## 2.2 实现流程

```
1 【1】 确认数据来源 - 响应内容存在
2 【2】 分析URL地址规律 - start为0 25 50 75 ...
3 【3】 xpath表达式
4     3.1) 基准xpath, 匹配每本书籍的节点对象列表
5         //div[@class="indent"]/table
6
7     3.2) 依次遍历每本书籍的节点对象, 提取具体书籍数据
8         书籍名称 : .//div[@class="p12"]/a/@title
9         书籍描述 : .//p[@class="p1"]/text()
10        书籍评分 : .//span[@class="rating_nums"]/text()
11        评价人数 : .//span[@class="p1"]/text()
12        书籍类型 : .//span[@class="inq"]/text()
```

## 2.3 代码实现

```
1 import requests
2 from lxml import etree
3 import time
4 import random
5 from fake_useragent import UserAgent
6
7 class DoubanBookSpider:
8     def __init__(self):
9         self.url = 'https://book.douban.com/top250?start={}'
10
11    def get_html(self, url):
12        headers = { 'User-Agent':UserAgent().random }
13        html = requests.get(url=url, headers=headers).content.decode('utf-8', 'ignore')
14        # 直接调用解析函数
15        self.parse_html(html)
16
17    def parse_html(self, html):
18        p = etree.HTML(html)
```

```

19     # 基准xpath,匹配每本书的节点对象列表
20     table_list = p.xpath('//div[@class="indent"]/table')
21     for table in table_list:
22         item = {}
23         # 书名
24         name_list = table.xpath('.//div[@class="pl2"]/a/@title')
25         item['book_name'] = name_list[0].strip() if name_list else None
26         # 信息
27         info_list = table.xpath('.//p[@class="pl"]/text()')
28         item['book_info'] = info_list[0].strip() if info_list else None
29         # 评分
30         score_list = table.xpath('.//span[@class="rating_nums"]/text()')
31         item['book_score'] = score_list[0].strip() if score_list else None
32         # 人数
33         number_list = table.xpath('.//span[@class="pl"]/text()')
34         item['book_number'] = number_list[0].strip()[1:-1].strip() if number_list else
None
35         # 描述
36         comment_list = table.xpath('.//span[@class="inq"]/text()')
37         item['book_comment'] = comment_list[0].strip() if comment_list else None
38
39         print(item)
40
41     def run(self):
42         for i in range(10):
43             start = i * 25
44             page_url = self.url.format(start)
45             self.get_html(url=page_url)
46             # 控制数据抓取的频率,uniform生成指定范围内浮点数
47             time.sleep(random.uniform(0, 3))
48
49
50 if __name__ == '__main__':
51     spider = DoubanBookSpider()
52     spider.run()

```

## 3. 百度贴吧小视频爬虫

### 3.1 需求分析

- 1 【1】官网地址：进入某个百度贴吧，寻找有视频的帖子，比如如下帖子链接：
- 2 <https://tieba.baidu.com/p/7185877941>
- 3
- 4 【2】目标
- 5 2.1> 在此帖子中提取中具体视频的链接(src)
- 6 2.2> 将视频抓取保存到本地文件(向src发请求获取bytes数据类型,以wb方式保存到本地)

### 3.2 实现流程

```
1 【1】 确认数据来源：静态!!!
2 【2】 帖子中视频的xpath表达式
3
4 ### 重要：页面中xpath不能全信，一切以响应内容为主
5 ### 重要：页面中xpath不能全信，一切以响应内容为主
6 ### 重要：页面中xpath不能全信，一切以响应内容为主
7 ### 重要：页面中xpath不能全信，一切以响应内容为主
8 ### 重要：页面中xpath不能全信，一切以响应内容为主
9 ### 重要：页面中xpath不能全信，一切以响应内容为主
10 ### 重要：页面中xpath不能全信，一切以响应内容为主
11 ### 重要：页面中xpath不能全信，一切以响应内容为主
12 ### 重要：页面中xpath不能全信，一切以响应内容为主
13 ### 重要：页面中xpath不能全信，一切以响应内容为主
14 ### 重要：页面中xpath不能全信，一切以响应内容为主
```

### 3.3 代码实现

```
1 import requests
2 from lxml import etree
3
4 # 向帖子链接发请求
5 url = 'https://tieba.baidu.com/p/7185877941'
6 headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
7 like Gecko) Chrome/87.0.4280.88 Safari/537.36'}
8 html = requests.get(url=url, headers=headers).text
9
10 # 视频链接的xpath表达式 - 一切以响应内容为准
11 x = '//div[@class="video_src_wrapper"]/embed/@data-video'
12 eobj = etree.HTML(html)
13 video_url_list = eobj.xpath(x)
14
15 # 将视频抓取保存到本地文件
16 if video_url_list:
17     video_url = video_url_list[0]
18     video_html = requests.get(url=video_url, headers=headers).content
19     with open('girl.mp4', 'wb') as f:
20         f.write(video_html)
21 else:
22     print('提取视频链接失败')
```

## 4. 代理参数

### 4.1 代理IP概述

```
1 【1】 定义
2     代替你原来的IP地址去对接网络的IP地址
3
4 【2】 作用
```

```
5     隐藏自身真实IP,避免被封
6
7     【3】获取代理IP网站
8         快代理、全网代理、代理精灵、... ...
9
10    【4】参数类型
11        proxies
12        proxies = { '协议':'协议://IP:端口号' }
13        proxies = { '协议':'协议://用户名:密码@IP:端口号' }
14
```

## 4.2 代理分类

### 4.2.1 普通代理

```
1     【1】代理格式
2         proxies = { '协议':'协议://IP:端口号' }
3
4     【2】使用免费普通代理IP访问测试网站: http://httpbin.org/get
5
6     import requests
7     url = 'http://httpbin.org/get'
8     headers = {'User-Agent':'Mozilla/5.0'}
9     # 定义代理,在代理IP网站中查找免费代理IP
10    proxies = {
11        'http': 'http://112.85.164.220:9999',
12        'https': 'https://112.85.164.220:9999'
13    }
14    html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
15    print(html)
16
```

### 4.2.2 私密代理和独享代理

```
1     【1】代理格式
2         proxies = { '协议':'协议://用户名:密码@IP:端口号' }
3
4     【2】使用私密代理或独享代理IP访问测试网站: http://httpbin.org/get
5
6     import requests
7     url = 'http://httpbin.org/get'
8     proxies = {
9        'http': 'http://309435365:szayclhp@106.75.71.140:16816',
10        'https': 'https://309435365:szayclhp@106.75.71.140:16816',
11    }
12    headers = {
13        'User-Agent' : 'Mozilla/5.0',
14    }
15
16    html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
17    print(html)
```

## 4.3 建立代理IP池

```

1  """
2  建立开放代理的代理ip池
3  """
4  import requests
5
6  class ProxyPool:
7      def __init__(self):
8          self.api_url = 'http://dev.kdlapi.com/api/getproxy/?
9          orderid=999955248138592&num=20&protocol=2&method=2&an_ha=1&sep=1'
10         self.test_url = 'http://httpbin.org/get'
11         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
12         AppleWebKit/537.36 (KHTML, like Gecko) Chrome/85.0.4183.83 Safari/537.36'}
13
14     def get_proxy(self):
15         html = requests.get(url=self.api_url, headers=self.headers).text
16         # proxy_list: ['1.1.1.1:8888', '2.2.2.2:9999,...]
17         proxy_list = html.split('\r\n')
18         for proxy in proxy_list:
19             # 测试proxy是否可用
20             self.test_proxy(proxy)
21
22     def test_proxy(self, proxy):
23         """测试1个代理ip是否可用"""
24         proxies = {
25             'http' : 'http://{0}'.format(proxy),
26             'https': 'https://{0}'.format(proxy),
27         }
28         try:
29             resp = requests.get(url=self.test_url, proxies=proxies, headers=self.headers,
30             timeout=3)
31             if resp.status_code == 200:
32                 print(proxy, '\033[31m可用\033[0m')
33             else:
34                 print(proxy, '不可用')
35         except Exception as e:
36             print(proxy, '不可用')
37
38     def run(self):
39         self.get_proxy()
40
41 if __name__ == '__main__':
42     spider = ProxyPool()
43     spider.run()

```

## 5. requests.post()

### 5.1 POST请求

```
1  【1】适用场景：Post类型请求的网站
2
3  【2】参数：data={}
4      2.1) Form表单数据：字典
5      2.2) res = requests.post(url=url,data=data,headers=headers)
6
7  【3】POST请求特点：Form表单提交数据
8
```

### 5.2 控制台抓包

#### ■ 打开方式及常用选项

```
1  【1】打开浏览器，F12打开控制台，找到Network选项卡
2
3  【2】控制台常用选项
4      2.1) Network：抓取网络数据包
5          a> ALL：抓取所有的网络数据包
6          b> XHR：抓取异步加载的网络数据包
7          c> JS：抓取所有的JS文件
8      2.2) Sources：格式化输出并打断点调试JavaScript代码，助于分析爬虫中一些参数
9      2.3) Console：交互模式，可对JavaScript中的代码进行测试
10
11 【3】抓取具体网络数据包后
12     3.1) 单击左侧网络数据包地址，进入数据包详情，查看右侧
13     3.2) 右侧：
14         a> Headers：整个请求信息
15             General、Response Headers、Request Headers、Query String、Form Data
16         b> Preview：对响应内容进行预览
17         c> Response：响应内容
18
```

## 6. 今日作业

```
1  【1】链家二手房爬虫
2      # 注意：一切以响应内容为准
3      1.1> 官网地址：进入链家官网，点击二手房：https://bj.lianjia.com/ershoufang/
4      1.2> 目标：抓取100页的二手房源信息，包含房源的：
5          名称
6          地址
7          户型、面积、方位、是否精装、楼层、年代、类型
8          总价
9          单价
```

```
10 1.3> 数据处理
11 将数据分别存入：MySQL、MongoDB、csv文件中
12
13 【2】抓取快代理网站免费高匿代理，并测试是否可用来建立自己的代理IP池
14 https://www.kuaidaili.com/free/
```

# SPIDER-DAY04

## 1. 有道翻译爬虫

### 1.1 项目需求

```
1 破解有道翻译接口，抓取翻译结果
2
3 # 结果展示
4 请输入要翻译的词语：elephant
5 翻译结果：大象
6 *****
7 请输入要翻译的词语：喵喵叫
8 翻译结果：mews
```

### 1.2 项目分析流程

```
1 【1】准备抓包：F12开启控制台，刷新页面
2 【2】寻找地址
3 2.1) 页面中输入翻译单词，控制台中抓取到网络数据包，查找并分析返回翻译数据的地址
4 F12-Network-XHR-Headers-General-Request URL
5 【3】发现规律
6 3.1) 找到返回具体数据的地址，在页面中多输入几个单词，找到对应URL地址
7 3.2) 分析对比 Network - All(或者XHR) - Form Data，发现对应的规律
8 【4】寻找JS加密文件
9 控制台右上角 ...->Search->搜索关键字->单击->跳转到Sources，左下角格式化符号{}
10 【5】查看JS代码
11 搜索关键字，找到相关加密方法，用python实现加密算法
12 【6】断点调试
13 JS代码中部分参数不清楚可通过断点调试来分析查看
14 【7】Python实现JS加密算法
```

### 1.3 项目步骤

1、开启F12抓包，找到Form表单数据如下：



```
1 i: 喵喵叫
2 from: AUTO
3 to: AUTO
4 smartresult: dict
5 client: fanyideskweb
6 salt: 15614112641250
7 sign: 94008208919faa19bd531acde36aac5d
8 ts: 1561411264125
9 bv: f4d62a2579ebb44874d7ef93ba47e822
10 doctype: json
11 version: 2.1
12 keyfrom: fanyi.web
13 action: FY_BY_REALTIME
```

## 2、在页面中多翻译几个单词，观察Form表单数据变化

```
1 salt: 15614112641250
2 sign: 94008208919faa19bd531acde36aac5d
3 ts: 1561411264125
4 bv: f4d62a2579ebb44874d7ef93ba47e822
5 # 但是bv的值不变
```

## 3、一般为本地js文件加密，刷新页面，找到js文件并分析JS代码

```
1 控制台右上角 - Search - 搜索salt - 查看文件 - 格式化输出
2
3 【结果】：最终找到相关JS文件：fanyi.min.js
```

## 4、打开JS文件，分析加密算法，用Python实现

```
1 【ts】经过分析为13位的时间戳，字符串类型
2   js代码实现) "" + (new Date).getTime()
3   python实现) str(int(time.time()*1000))
4
5 【salt】
6   js代码实现) ts + parseInt(10 * Math.random(), 10);
7   python实现) ts + str(random.randint(0, 9))
8
9 【sign】('设置断点调试，来查看 e 的值，发现 e 为要翻译的单词')
10  js代码实现) n.md5("fanyideskweb" + e + salt + "]BjuETDhU)zqSxf-=B#7m")
11
12   from hashlib import md5
13   s = md5()
14   s.update('').encode())
15   sign = s.hexdigest()
```

## 5、pycharm中正则处理headers和formdata

```
1 【1】pycharm进入方法：Ctrl + r，选中 Regex
2 【2】处理headers和formdata
3     (.*)：(.*)
4     "$1": "$2",
5 【3】点击 Replace All
```

## 1.4 代码实现

```
1 import requests
2 import time
3 import random
4 from hashlib import md5
5
6 class YdSpider(object):
7     def __init__(self):
8         # url一定为F12抓到的 headers -> General -> Request URL
9         self.url = 'http://fanyi.youdao.com/translate_o?smartresult=dict&smartresult=rule'
10        self.headers = {
11            # 检查频率最高 - 3个
12            "Cookie": "OUTFOX_SEARCH_USER_ID=970246104@10.169.0.83;
OUTFOX_SEARCH_USER_ID_NCOO=570559528.1224236;
_ntes_nnid=96bc13a2f5ce64962adfd6a278467214,1551873108952; JSESSIONID=aaae9i7plXP1KaJH_gkYw;
td_cookie=18446744072941336803; SESSION_FROM_COOKIE=unknown;
__rl__test__cookies=1565689460872",
13            "Referer": "http://fanyi.youdao.com/",
14            "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
Gecko) Chrome/76.0.3809.100 Safari/537.36",
15        }
16
17        # 获取salt,sign,ts
18        def get_salt_sign_ts(self, word):
19            # ts
20            ts = str(int(time.time()*1000))
21            # salt
22            salt = ts + str(random.randint(0,9))
23            # sign
24            string = "fanyideskweb" + word + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
25            s = md5()
26            s.update(string.encode())
27            sign = s.hexdigest()
28
29            return salt, sign, ts
30
31        # 主函数
32        def attack_yd(self, word):
33            # 1. 先拿到salt,sign,ts
34            salt, sign, ts = self.get_salt_sign_ts(word)
35            # 2. 定义form表单数据为字典: data={}
36            # 检查了salt sign
37            data = {
38                "i": word,
39                "from": "AUTO",
```

```

40         "to": "AUTO",
41         "smartresult": "dict",
42         "client": "fanyideskweb",
43         "salt": salt,
44         "sign": sign,
45         "ts": ts,
46         "bv": "7e3150ecbdf9de52dc355751b074cf60",
47         "doctype": "json",
48         "version": "2.1",
49         "keyfrom": "fanyi.web",
50         "action": "FY_BY_REALTIME",
51     }
52     # 3. 直接发请求: requests.post(url, data=data, headers=xxx)
53     html = requests.post(
54         url=self.url,
55         data=data,
56         headers=self.headers
57     ).json()
58     # res.json() 将json格式的字符串转为python数据类型
59     result = html['translateResult'][0][0]['tgt']
60
61     print(result)
62
63     # 主函数
64     def run(self):
65         # 输入翻译单词
66         word = input('请输入要翻译的单词: ')
67         self.attack_yd(word)
68
69     if __name__ == '__main__':
70         spider = YdSpider()
71         spider.run()

```

## 2. 百度翻译JS逆向爬虫

### 2.1 JS逆向详解

```

1  【1】应用场景
2      当JS加密的代码过于复杂,没有办法破解时,考虑使用JS逆向思想
3
4  【2】模块
5      2.1》模块名: execjs
6      2.2》安装:  sudo pip3 install pyexecjs
7      2.3》使用流程
8          import execjs
9          with open('xxx.js', 'r') as f:
10              js_code = f.read()
11
12              js_obj = execjs.compile(js_code)
13              js_obj.eval('函数名("参数")')

```

## 2.2 JS代码调试

- 抓到JS加密文件，存放到 translate.js 文件中

```
1 // e(r, gtk) 增加了gtk参数
2 // i = window[1] 改为了 i = gtk
3 function a(r) {
4     if (Array.isArray(r)) {
5         for (var o = 0, t = Array(r.length); o < r.length; o++)
6             t[o] = r[o];
7         return t
8     }
9     return Array.from(r)
10 }
11 function n(r, o) {
12     for (var t = 0; t < o.length - 2; t += 3) {
13         var a = o.charAt(t + 2);
14         a = a >= "a" ? a.charCodeAt(0) - 87 : Number(a),
15         a = "+" === o.charAt(t + 1) ? r >>> a : r << a,
16         r = "+" === o.charAt(t) ? r + a & 4294967295 : r ^ a
17     }
18     return r
19 }
20 function e(r,gtk) {
21     var o = r.match(/[\uD800-\uDBFF][\uDC00-\uDFFF]/g);
22     if (null === o) {
23         var t = r.length;
24         t > 30 && (r = "" + r.substr(0, 10) + r.substr(Math.floor(t / 2) - 5, 10) +
25 r.substr(-10, 10))
26     } else {
27         for (var e = r.split(/[\uD800-\uDBFF][\uDC00-\uDFFF]/), C = 0, h = e.length, f
28 = []; h > C; C++)
29             "" !== e[C] && f.push.apply(f, a(e[C].split(""))),
30             C !== h - 1 && f.push(o[C]);
31         var g = f.length;
32         g > 30 && (r = f.slice(0, 10).join("") + f.slice(Math.floor(g / 2) - 5,
33 Math.floor(g / 2) + 5).join("") + f.slice(-10).join(""))
34     }
35     var u = void 0
36     , l = "" + String.fromCharCode(103) + String.fromCharCode(116) +
37 String.fromCharCode(107);
38     u = null !== i ? i : (i = gtk || "") || "";
39     for (var d = u.split("."), m = Number(d[0]) || 0, s = Number(d[1]) || 0, S = [], c
40 = 0, v = 0; v < r.length; v++) {
41         var A = r.charCodeAt(v);
42         128 > A ? S[c++] = A : (2048 > A ? S[c++] = A >> 6 | 192 : (55296 === (64512 &
43 A) && v + 1 < r.length && 56320 === (64512 & r.charCodeAt(v + 1)) ? (A = 65536 + ((1023
44 & A) << 10) + (1023 & r.charCodeAt(++v)),
45         S[c++] = A >> 18 | 240,
46         S[c++] = A >> 12 & 63 | 128) : S[c++] = A >> 12 | 224,
47         S[c++] = A >> 6 & 63 | 128),
48         S[c++] = 63 & A | 128)
49     }
```

```

43     for (var p = m, F = "" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(97) + ("" + String.fromCharCode(94) + String.fromCharCode(43) +
String.fromCharCode(54)), D = "" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(51) + ("" + String.fromCharCode(94) + String.fromCharCode(43) +
String.fromCharCode(98)) + ("" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(102)), b = 0; b < S.length; b++)
44         p += S[b],
45         p = n(p, F);
46     return p = n(p, D),
47     p ^= s,
48     0 > p && (p = (2147483647 & p) + 2147483648),
49     p %= 1e6,
50     p.toString() + "." + (p ^ m)
51 }
52 var i = null;

```

#### ▪ test\_translate.py调试JS文件

```

1  import execjs
2
3  with open('translate.js', 'r', encoding='utf-8') as f:
4      jscode = f.read()
5
6  jsobj = execjs.compile(jscode)
7  sign = jsobj.eval('e("hello","320305.131321201")')
8  print(sign)

```

## 2.3 百度翻译代码实现

```

1  import requests
2  import execjs
3  import re
4
5  class BaiduTranslateSpider:
6      def __init__(self):
7          self.url = 'https://fanyi.baidu.com/v2transapi?from=en&to=zh'
8          self.index_url = 'https://fanyi.baidu.com/'
9          self.post_headers = {
10              "accept": "*/*",
11              "accept-encoding": "gzip, deflate, br",
12              "accept-language": "zh-CN,zh;q=0.9",
13              "cache-control": "no-cache",
14              "content-length": "135",
15              "content-type": "application/x-www-form-urlencoded; charset=UTF-8",

```



```

45         gtk = re.findall("window.gtk = '(.*?)'", html, re.S)[0]
46         token = re.findall("token: '(.*?)'", html, re.S)[0]
47
48         return gtk, token
49
50     def get_sign(self, word):
51         """功能函数:生成sign"""
52         # 先获取到gtk和token
53         gtk, token = self.get_gtk_token()
54         with open('translate.js', 'r', encoding='utf-8') as f:
55             js_code = f.read()
56
57         js_obj = execjs.compile(js_code)
58         sign = js_obj.eval('e("{}","{}").format(word, gtk)')
59
60         return sign
61
62     def attack_bd(self, word):
63         """爬虫逻辑函数"""
64         gtk, token = self.get_gtk_token()
65         sign = self.get_sign(word)
66         data = {
67             "from": "en",
68             "to": "zh",
69             "query": word,
70             "transtype": "realtime",
71             "simple_means_flag": "3",
72             "sign": sign,
73             "token": token,
74             "domain": "common",
75         }
76         # json():把json格式的字符串转为python数据类型
77         html = requests.post(url=self.url,
78                             data=data,
79                             headers=self.post_headers).json()
80         result = html['trans_result'][0]['dst']
81
82         return result
83
84     def run(self):
85         word = input('请输入要翻译的单词:')
86         print(self.attack_bd(word))
87
88 if __name__ == '__main__':
89     spider = BaiduTranslateSpider()
90     spider.run()

```

## 3. 动态加载数据抓取

### 3.1 AJAX动态加载

#### ■ 数据特点

- 1 【1】 右键 -> 查看网页源码中没有具体数据
- 2 【2】 滚动鼠标滑轮或其他动作时加载,或者页面局部刷新

#### ■ 分析流程

- 1 【1】 F12打开控制台, 页面动作抓取网络数据包
- 2 【2】 抓取json文件URL地址
- 3 2.1) 控制台中 XHR : 异步加载的数据包
- 4 2.2) XHR -> QueryStringParameters(查询参数)

## 3.2 豆瓣电影爬虫

### 3.2.1 项目需求

- 1 【1】 地址: 豆瓣电影 - 排行榜 - 剧情
- 2 【2】 目标: 电影名称、电影评分
- 3
- 4 `<span><a href=".*?type_name=(.*?)&type=(.*?)&`
- 5 `https://movie.douban.com/chart`

### 3.2.2 抓包分析

- 1 【1】 Request URL(基准URL地址) : `https://movie.douban.com/j/chart/top_list?`
- 2 【2】 Query String(查询参数)
- 3 # 抓取的查询参数如下:
- 4 `type: 13 # 电影类型`
- 5 `interval_id: 100:90`
- 6 `action: ''`
- 7 `start: 0 # 每次加载电影的起始索引值 0 20 40 60`
- 8 `limit: 20 # 每次加载的电影数量`

### 3.2.3 代码实现

```
1 """
2 抓取豆瓣电影数据 - 全站抓取
3 """
4 import requests
5 import json
6 import time
7 import random
8 from fake_useragent import UserAgent
9 import re
10
11 class DoubanSpider:
12     def __init__(self):
13         self.url = 'https://movie.douban.com/j/chart/top_list?type=
14         {}&interval_id=100%3A90&action=&start={}&limit=20'
```



```

15 def get_html(self, url):
16     """功能函数1: 获取html"""
17     headers = {'User-Agent': UserAgent().random}
18     html = requests.get(url=url, headers=headers).text
19
20     return html
21
22 def parse_html(self, url):
23     # 提取数据函数
24     # html: [{},{},{},...]
25     html = json.loads(self.get_html(url=url))
26     for one_film_dict in html:
27         item = {}
28         item['rank'] = one_film_dict['rank']
29         item['name'] = one_film_dict['title']
30         item['time'] = one_film_dict['release_date']
31         item['score'] = one_film_dict['score']
32
33         print(item)
34
35 def get_total(self, typ):
36     """获取电影总数"""
37     total_url = 'https://movie.douban.com/j/chart/top_list_count?type=
38 {}&interval_id=100%3A90'.format(typ)
39     total_html = json.loads(self.get_html(url=total_url))
40
41     return total_html['total']
42
43 def get_all_film_dict(self):
44     """获取所有电影类别及对应的type值的字典"""
45     all_type_url = 'https://movie.douban.com/chart'
46     all_type_html = self.get_html(url=all_type_url)
47     regex = '<span><a href=.*?type_name=(.*?)&type=(.*?)&interval_id=100:90&action=">'
48     pattern = re.compile(regex, re.S)
49     # all_list: [('剧情', '11'), ('喜剧', '5'), ...]
50     all_list = pattern.findall(all_type_html)
51     all_film_dict = {}
52     for one in all_list:
53         all_film_dict[one[0]] = one[1]
54
55     return all_film_dict
56
57 def run(self):
58     # {'剧情': '5', '喜剧': '23', '爱情': '13', ... ...}
59     all_film_dict = self.get_all_film_dict()
60     # 生成提示菜单
61     menu = ''
62     for key in all_film_dict:
63         menu = menu + key + '|'
64     print(menu)
65     # 接收用户输入,并获取对应的type的值
66     film_type = input('请输入电影类别:')
67     typ = all_film_dict[film_type]
68     # 获取此类别下的电影总数
69     total = self.get_total(typ)
70     for start in range(0, total, 20):
71         page_url = self.url.format(typ, start)

```

```
71         self.parse_html(url=page_url)
72         # 控制频率
73         time.sleep(random.randint(1, 2))
74
75
76 if __name__ == '__main__':
77     spider = DoubanSpider()
78     spider.run()
```

## 4. 今日作业

- 1 **【1】肯德基餐厅门店信息抓取（POST请求练习）**
- 2 1.1) URL地址: <http://www.kfc.com.cn/kfccda/storelist/index.aspx>
- 3 1.2) 所抓数据: 餐厅编号、餐厅名称、餐厅地址、城市
- 4 1.3) 数据存储: 保存到数据库
- 5 1.4) 程序运行效果:
- 6 请输入城市名称: 北京
- 7 把北京的所有肯德基门店的信息保存到数据库中