

Day05回顾

■ Ajax动态加载数据抓取流程

```
1  【1】 F12打开控制台，执行页面动作抓取网络数据包
2
3  【2】 抓取json文件URL地址
4      2.1) 控制台中 XHR ： 异步加载的数据包
5      2.2) XHR -> Query String Parameters(查询参数)
```

■ json模块

```
1  【1】 抓取的json数据转为python数据类型
2      1.1) html = json.loads('[{},{},{}]')
3      1.2) html = requests.get(url=url,headers=headers).json()
4      1.3) html = requests.post(url=url,headers=headers).json()
5
6  【2】 抓取数据保存到json文件
7      import json
8      with open('xxx.json','w') as f:
9          json.dump([{}},{},{}],f,ensure_ascii=False)
```

■ 数据抓取最终梳理

```
1  【1】 响应内容中存在
2      1.1) 确认抓取数据在响应内容中是否存在
3
4      1.2) 分析页面结构，观察URL地址规律
5          a) 大体查看响应内容结构，查看是否有更改 --（百度视频案例）
6          b) 查看页面跳转时URL地址变化，查看是否新跳转 --（民政部案例）
7
8      1.3) 开始写代码进行数据抓取
9
10 【2】 响应内容中不存在
11     2.1) 确认抓取数据在响应内容中是否存在
12
13     2.2) F12抓包，开始刷新页面或执行某些行为，主要查看XHR异步加载数据包
14         a) GET请求：Request Headers、Query String Parameters
15         b) POST请求：Request Headers、FormData
16
17     2.3) 观察查询参数或者Form表单数据规律，如果需要进行进一步抓包分析处理
18         a) 比如有道翻译的 salt+sign，抓取并分析JS做进一步处理
19         b) 此处注意请求头中的cookie和referer以及User-Agent
20
21     2.4) 使用res.json()获取数据，利用列表或者字典的方法获取所需数据
```

■ 多线程爬虫思路梳理

```
1  【1】所用到的模块
2      1.1) from threading import Thread
3      1.2) from threading import Lock
4      1.3) from queue import Queue
5
6  【2】整体思路
7      2.1) 创建URL队列: q = Queue()
8      2.2) 产生URL地址,放入队列: q.put(url)
9      2.3) 线程事件函数: 从队列中获取地址,开始抓取: url = q.get()
10     2.4) 创建多线程,并运行
11
12  【3】代码结构
13     def __init__(self):
14         """创建URL队列"""
15         self.q = Queue()
16
17     def url_in(self):
18         """生成待爬取的URL地址,入队列"""
19         pass
20
21     def parse_html(self):
22         """线程事件函数,获取地址,进行数据抓取"""
23         while True:
24             if not self.q.empty():
25                 url = self.q.get()
26             else:
27                 break
28
29     def run(self):
30         self.url_in()
31         t_list = []
32         for i in range(3):
33             t = Thread(target=self.parse_html)
34             t_list.append(t)
35             t.start()
36
37         for th in t_list:
38             th.join()
```

Day06笔记

多线程爬虫

■ 多线程回顾

```
1  【1】应用场景
2      1.1) 多进程 : CPU密集程序
```

```

3      1.2) 多线程：爬虫(网络I/O)、本地磁盘I/O
4
5      【2】队列要点：q.get()防止阻塞方式
6      2.1) 方法1: q.get(block=False)
7      2.2) 方法2: q.get(block=True, timeout=3)
8      2.3) 方法3:
9          if not q.empty():
10             q.get()
11
12      【3】线程锁的使用:
13      from threading import Lock
14      lock = Lock()
15      lock.acquire()
16      lock.release()
17      【注意】：上锁后再次上锁会阻塞,多线程操作共享资源时需要加锁和释放锁

```

小米应用商店抓取(多线程)

目标

```

1      【1】网址：百度搜 - 小米应用商店, 进入官网 http://app.mi.com/
2
3      【2】目标：抓取聊天社交分类下的
4      2.1) 应用名称
5      2.2) 应用链接

```

实现步骤

```

1      【1】确认是否为动态加载
2      1.1) 页面局部刷新
3      1.2) 右键查看网页源代码, 搜索关键字未搜到, 为动态加载, 需要抓取网络数据包分析
4
5      【2】F12抓取网络数据包
6      2.1) 抓取返回json数据的URL地址 (Headers中的Request URL)
7          http://app.mi.com/categoryAllListApi?page={}&categoryId=2&pageSize=30
8
9      2.2) 查看并分析查询参数 (headers中的Query String Parameters)
10         page: 1          只有page在变, 0 1 2 3 ... ..
11         categoryId: 2
12         pageSize: 30
13
14      【3】将抓取数据保存到csv文件 - 注意线程锁问题
15      from threading import Lock
16      lock = Lock()
17      # 加锁 + 释放锁
18      lock.acquire()
19      lock.release()

```

整体思路

```

1      【1】在 __init__(self) 中创建文件对象, 多线程操作此对象进行文件写入
2      self.f = open('xiaomi.csv', 'a')

```

```

3     self.writer = csv.writer(self.f)
4     self.lock = Lock()
5
6     【2】 每个线程抓取1页数据后将数据进行文件写入，写入文件时需要加锁
7     def parse_html(self):
8         app_list = []
9         for xxx in xxx:
10             app_list.append([name,link,typ])
11             self.lock.acquire()
12             self.wirter.writerow(app_list)
13             self.lock.release()
14
15     【3】 所有数据抓取完成关闭文件
16     def run(self):
17         self.f.close()

```

■ 代码实现

```

1 class XiaomiSpider:
2     def __init__(self):
3         self.url = 'http://app.mi.com/categotyAllListApi?page={}&categoryId=2&pageSize=30'
4         self.headers = { 'User-Agent':UserAgent().random }
5         self.q = Queue()
6         self.lock = Lock()
7         # 计数
8         self.i = 0
9         # 存入csv文件
10        self.f = open('xiaomi.csv', 'w', encoding='utf-8')
11        self.writer = csv.writer(self.f)
12
13        # 获取html
14        def get_html(self,url):
15            html = requests.get(url=url,headers=self.headers).text
16
17            return html
18
19        # URL入队列
20        def url_in(self):
21            for page in range(67):
22                url = self.url.format(page)
23                self.q.put(url)
24
25        # 线程事件函数
26        def parse_html(self):
27            while True:
28                # 加锁 - 防止出现死锁(self.q中剩余1个地址,但是被多个线程判断的情况)
29                self.lock.acquire()
30                if not self.q.empty():
31                    url = self.q.get()
32                    print(url)
33                    # 获取地址成功后马上释放锁,给其他线程机会,安全前提下提升效率
34                    self.lock.release()
35                    app_list = []
36                    # 请求 + 解析  html: {'count':2000,'data':[{},{},{}]}
37                    html = self.get_html(url=url)
38                    print(html)

```

```

39         html = json.loads(html)
40         item = {}
41         for one_app in html['data']:
42             item['app_name'] = one_app['displayName']
43             item['app_type'] = one_app['level1CategoryName']
44             item['app_link'] = one_app['packageName']
45             print(item)
46             # 一页的app信息
47             app_list.append( (item['app_name'],item['app_type'],item['app_link']) )
48             # 加锁+释放锁
49             self.lock.acquire()
50             self.i += 1
51             self.lock.release()
52             # 把1页的数据写入到csv文件
53             self.lock.acquire()
54             self.writer.writerow(app_list)
55             self.lock.release()
56             # 简单控制一下数据抓取频率,因为我们没有代理IP,容易被封掉IP
57             time.sleep(random.uniform(2,4))
58         else:
59             # 如果队列为空了,上面已经上锁,所以此处释放锁
60             self.lock.release()
61             break
62
63     # 入口函数
64     def run(self):
65         # 1.先让URL地址入队列
66         self.url_in()
67         # 2.多线程,开始执行
68         t_list = []
69         for i in range(2):
70             t = Thread(target=self.parse_html)
71             t_list.append(t)
72             t.start()
73
74         for j in t_list:
75             j.join()
76
77         print('数量:',self.i)
78         # 最终关闭文件
79         self.f.close()
80
81     if __name__ == '__main__':
82         start_time = time.time()
83         spider = XiaomiSpider()
84         spider.run()
85         end_time = time.time()
86         print('执行时间:%.2f' % (end_time-start_time))

```

腾讯招聘数据抓取(多线程)

- 确定URL地址及目标

```
1 【1】URL：百度搜索腾讯招聘 - 查看工作岗位
2 【2】目标:抓取职位的如下信息
3     a> 职位名称
4     b> 职位地址
5     c> 职位类别（技术类、销售类...）
6     d> 发布时间
7     e> 工作职责
8     f> 工作要求
```

■ 要求与分析

```
1 【1】通过查看网页源码,得知所需数据均为动态加载
2 【2】通过F12抓取网络数据包,进行分析
3 【3】一级页面抓取数据: postid
4 【4】二级页面抓取数据: 名称+地址+类别+时间+职责+要求
```

■ 一级页面json地址

```
1 ""index在变,timestamp未检查""
2 https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1563912271089&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=&attrId=&keyword={}&pageIndex={}&pageSize=10&language=zh-cn&area=cn
```

■ 二级页面地址

```
1 ""postId在变,在一级页面中可拿到""
2 https://careers.tencent.com/tencentcareer/api/post/ByPostId?timestamp=1563912374645&postId=
{}&language=zh-cn
```

■ 多线程编写思路提示

```
1 【思考】 两级页面是否需要指定两个队列分别存放?
2 提示1: 建立2个队列,分别存放不同级的URL地址
3 提示2: 从对列中get地址,最好使用timeout参数
```

■ 代码实现

```
1 import requests
2 import json
3 import time
4 from fake_useragent import UserAgent
5 from queue import Queue
6 from threading import Thread, Lock
7 from urllib import parse
8 import csv
9
10 class TencentSpider(object):
11     def __init__(self):
12         self.one_url = 'https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1563912271089&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=
&attrId=&keyword={}&pageIndex={}&pageSize=10&language=zh-cn&area=cn'
```

```

13         self.two_url = 'https://careers.tencent.com/tencentcareer/api/post/ByPostId?
timestamp=1563912374645&postId={}&language=zh-cn'
14         self.one_q = Queue()
15         self.two_q = Queue()
16         self.lock1 = Lock()
17         self.lock2 = Lock()
18         self.i = 0
19         self.f = open('tencent.csv', 'w', encoding='utf-8')
20         self.writer = csv.writer(self.f)
21         # 存放所有数据的大列表,用于writerows()方法
22         self.item_list = []
23
24
25     def get_html(self, url):
26         headers = { 'User-Agent':UserAgent().random }
27         html = requests.get(url=url, headers=headers).text
28         return html
29
30     def url_in(self):
31         keyword = input('请输入职位类别:')
32         keyword = parse.quote(keyword)
33         total = self.get_total(keyword)
34         for page in range(1, total+1):
35             one_url = self.one_url.format(keyword, page)
36             print(one_url)
37             self.one_q.put(one_url)
38
39     # 获取总页数
40     def get_total(self, keyword):
41         url = self.one_url.format(keyword, 1)
42         html = requests.get(url=url, headers={'User-Agent':UserAgent().random}).json()
43         n = int(html['Data']['Count'])
44         total = n//10 if n%10==0 else n//10+1
45
46         return total
47
48     # 线程1事件函数
49     def parse_one_page(self):
50         while True:
51             self.lock1.acquire()
52             if not self.one_q.empty():
53                 one_url = self.one_q.get()
54                 self.lock1.release()
55                 html = json.loads(self.get_html(one_url))
56                 for job in html['Data']['Posts']:
57                     post_id = job['PostId']
58                     two_url = self.two_url.format(post_id)
59                     self.lock1.acquire()
60                     self.two_q.put(two_url)
61                     self.lock1.release()
62             else:
63                 self.lock1.release()
64                 break
65
66     # 线程2事件函数
67     def parse_two_page(self):
68         while True:

```

```

69         try:
70             self.lock2.acquire()
71             two_url = self.two_q.get(block=True, timeout=3)
72             self.lock2.release()
73             html = json.loads(self.get_html(two_url))
74             # 名称+地址+类别+时间+职责+要求
75             item = {}
76             item['name'] = html['Data']['RecruitPostName']
77             item['address'] = html['Data']['LocationName']
78             item['type'] = html['Data']['CategoryName']
79             item['time'] = html['Data']['LastUpdateTime']
80             item['duty'] = html['Data']['Responsibility']
81             item['require'] = html['Data']['Requirement']
82             # 写入csv文件数据类型: [(), (), (), ..., ()]
83             item_tuple = (item['name'], item['duty'], item['require'])
84             self.item_list.append(item_tuple)
85
86             print(item)
87             self.lock2.acquire()
88             self.i += 1
89             self.lock2.release()
90         except Exception as e:
91             self.lock2.release()
92             print(e, end="")
93             break
94
95     def run(self):
96         self.url_in()
97         t1_list = []
98         t2_list = []
99         for i in range(5):
100             t = Thread(target=self.parse_one_page)
101             t1_list.append(t)
102             t.start()
103
104         for i in range(5):
105             t = Thread(target=self.parse_two_page)
106             t2_list.append(t)
107             t.start()
108
109         for t in t1_list:
110             t.join()
111
112         for t in t2_list:
113             t.join()
114
115         print('数量:', self.i)
116         # 将所有数据一次性写入文件
117         self.writer.writerows(self.item_list)
118         self.f.close()
119
120     if __name__ == '__main__':
121         start_time = time.time()
122         spider = TencentSpider()
123         spider.run()
124         end_time = time.time()
125         print('执行时间: %.2f' % (end_time - start_time))

```


cookie模拟登录

1 | 【1】适用网站及场景：抓取需要登录才能访问的页面

豆瓣网登录案例

■ 方法一 - 登录网站手动抓取Cookie

```
1  【1】先登录成功1次,获取到携带登录信息的Cookie
2      登录成功 - 我的豆瓣 - F12抓包 - 刷新主页 - 找到主页的包(一般为第1个网络数据包)
3
4  【2】headers中携带着Cookie发请求
5      headers = {
6          'Cookie': '',
7          'User-Agent': ''
8      }
```

```
1  """方法一代码实现"""
2
3  # 1、将url改为 个人主页的URL地址
4  # 2、将Cookie的值改为 登录成功的Cookie值
5  import requests
6
7  def login():
8      url = '个人主页的URL地址'
9      headers = {
10         'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like
11         Gecko) Chrome/80.0.3987.116 Safari/537.36',
12         'Cookie': '自己抓到的Cookie值',
13     }
14     html = requests.get(url=url,headers=headers).text
15     # 查看html中是否包含个人主页的信息 - 比如搜索 "个人主页"
16     print(html)
17     login()
```

■ 方法二

```

1  【1】原理
2      1.1) 把抓取到的cookie处理为字典
3      1.2) 使用requests.get()中的参数:cookies - 格式为字典
4
5  【2】处理cookie为字典
6      cookies = {}
7      cookies_str = 'xxxx'
8      for kv in cookies_str.split('; '):
9          key = kv.split('=')[0]
10         value = kv.split('=')[1]
11         cookies[key] = value

```

```

1  """方法二代码实现"""
2
3  import requests
4
5  def login():
6      url = '自己账号的个人主页'
7      headers = {
8          'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like
          Gecko) Chrome/80.0.3987.116 Safari/537.36',
9      }
10     # 处理cookie为字典
11     cookies_str = '自己抓到的Cookie'
12     cookies = {}
13     for kv in cookies_str.split('; '):
14         key = kv.split('=')[0]
15         value = kv.split('=')[1]
16         cookies[key] = value
17
18     # 确认html
19     html = requests.get(url=url,headers=headers,cookies=cookies).text
20     print(html)
21
22 login()

```

■ 方法三 - requests模块处理Cookie

```

1  【1】 思路 : requests模块提供了session类,来实现客户端和服务端的会话保持
2
3  【2】原理
4      2.1) 实例化session对象 : s = requests.session()
5      2.2) 让session对象发送get或者post请求
6          res = session.post(url=url,data=data,headers=headers)
7          res = session.get(url=url,headers=headers)
8
9  【3】思路梳理
10     3.1) 浏览器原理: 访问需要登录的页面会带着之前登录过的cookie
11     3.2) 程序原理: 同样带着之前登录的cookie去访问 - 由session对象完成
12     3.3) 具体步骤
13         a> 实例化session对象
14         b> 登录网站: 由session对象发送请求,登录对应网站
15         c> 访问页面: 由session对象请求需要登录才能访问的页面
16

```

```
17 【4】如何把用户名和密码信息提交给服务器
18 4.1) 输入用户名和错误密码,登录1次进行抓包
19 4.2) 在网络数据包中找到具体提交用户名和密码信息的地址,一般为POST请求
20 4.3) 将正确的用户名和密码信息POST到网络数据包的URL地址 - Request URL
```

```
1 """方法三代码实现"""
2 import requests
3
4 session = requests.session()
5
6 def login():
7     post_url = 'https://accounts.douban.com/j/mobile/login/basic'
8     post_data = {
9         'ck': '',
10        'name': '自己的用户名',
11        'password': '自己的密码',
12        'remember': 'false',
13        'ticket': '',
14    }
15     headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/80.0.3987.116 Safari/537.36'}
16     session.post(url=post_url,data=post_data,headers=headers)
17     url = '自己个人主页的URL地址'
18     html = session.get(url=url,headers=headers).text
19     print(html)
20
21 login()
```

selenium+phantomjs/Chrome/Firefox

■ selenium

```
1 【1】定义
2 1.1) Web自动化测试工具,可运行在浏览器,根据指令操作浏览器
3 1.2) 只是工具,必须与第三方浏览器结合使用
4
5 【2】安装
6 2.1) Linux: sudo pip3 install selenium
7 2.2) Windows: python -m pip install selenium
```

■ phantomjs浏览器

```
1 【1】定义
2 phantomjs为无界面浏览器(又称无头浏览器),在内存中进行页面加载,高效
3
4 【2】下载地址
5 2.1) chromedriver : 下载对应版本
6      http://chromedriver.storage.googleapis.com/index.html
7 2.2) geckodriver
8      https://github.com/mozilla/geckodriver/releases
9 2.3) phantomjs
```

```

10         https://phantomjs.org/download.html
11
12     【3】Ubuntu安装
13     3.1) 下载后解压 : tar -zxvf geckodriver.tar.gz
14
15     3.2) 拷贝解压后文件到 /usr/bin/ (添加环境变量)
16         sudo cp geckodriver /usr/bin/
17
18     3.3) 添加可执行权限
19         sudo chmod 777 /usr/bin/geckodriver
20
21     【4】Windows安装
22     4.1) 下载对应版本的phantomjs、chromedriver、geckodriver
23     4.2) 把chromedriver.exe拷贝到python安装目录的Scripts目录下(添加到系统环境变量)
24         # 查看python安装路径: where python
25     4.3) 验证
26         cmd命令行: chromedriver

```

■ 示例代码

```

1  """示例代码一: 使用 selenium+浏览器 打开百度"""
2
3  # 导入selenium的webdriver接口
4  from selenium import webdriver
5  import time
6
7  # 创建浏览器对象
8  browser = webdriver.Chrome()
9  browser.get('http://www.baidu.com/')
10 # 5秒钟后关闭浏览器
11 time.sleep(5)
12 browser.quit()

```

```

1  """示例代码二: 打开百度, 搜索赵丽颖, 点击搜索, 查看"""
2
3  from selenium import webdriver
4  import time
5
6  # 1.创建浏览器对象 - 已经打开了浏览器
7  browser = webdriver.Chrome()
8  # 2.输入: http://www.baidu.com/
9  browser.get('http://www.baidu.com/')
10 # 3.找到搜索框, 向这个节点发送文字: 赵丽颖
11 browser.find_element_by_xpath('//*[@id="kw"]').send_keys('赵丽颖')
12 # 4.找到 百度一下 按钮, 点击一下
13 browser.find_element_by_xpath('//*[@id="su"]').click()

```

■ 浏览器对象(browser)方法

```

1  【1】 browser = webdriver.Chrome(executable_path='path')
2  【2】 browser.get(url)
3  【3】 browser.page_source # HTML结构源码
4  【4】 browser.page_source.find('字符串')
5      # 从html源码中搜索指定字符串,没有找到返回: -1
6  【5】 browser.close() # 关闭当前页
7  【6】 browser.quit() # 关闭浏览器

```

■ 定位节点

```

1  【1】 单元素查找('结果为1个节点对象')
2      1.1) browser.find_element_by_id('')
3      1.2) browser.find_element_by_name('')
4      1.3) browser.find_element_by_class_name('')
5      1.4) browser.find_element_by_xpath('')
6      1.5) browser.find_element_by_link_text('')
7      ... ...
8
9  【2】 多元素查找('结果为[节点对象列表]')
10     2.1) browser.find_elements_by_id('')
11     2.2) browser.find_elements_by_name('')
12     2.3) browser.find_elements_by_class_name('')
13     2.4) browser.find_elements_by_xpath('')
14     ... ...

```

■ 猫眼电影示例

```

1  from selenium import webdriver
2  import time
3
4  url = 'https://maoyan.com/board/4'
5  browser = webdriver.Chrome()
6  browser.get(url)
7
8  def get_data():
9      # 基准xpath: [<selenium xxx li at xxx>,<selenium xxx li at>]
10     li_list = browser.find_elements_by_xpath('//*[@id="app"]/div/div/div[1]/dl/dd')
11     for li in li_list:
12         item = {}
13         # info_list: ['1', '霸王别姬', '主演: 张国荣', '上映时间: 1993-01-01', '9.5']
14         info_list = li.text.split('\n')
15         item['number'] = info_list[0]
16         item['name'] = info_list[1]
17         item['star'] = info_list[2]
18         item['time'] = info_list[3]
19         item['score'] = info_list[4]
20
21         print(item)
22
23  while True:
24     get_data()
25     try:
26         browser.find_element_by_link_text('下一页').click()
27         time.sleep(2)
28     except Exception as e:

```

```
29     print('恭喜你!抓取结束')
30     browser.quit()
31     break
```

■ 节点对象操作

```
1  【1】ele.send_keys('') # 搜索框发送内容
2  【2】ele.click()
3  【3】ele.text # 获取文本内容, 包含子节点和后代节点的文本内容
4  【4】ele.get_attribute('src') # 获取属性值
```

作业1 - 京东爬虫

■ 目标

```
1  【1】目标网址 : https://www.jd.com/
2  【2】抓取目标 : 商品名称、商品价格、评价数量、商品商家
```

■ 思路提醒

```
1  【1】打开京东, 到商品搜索页
2  【2】匹配所有商品节点对象列表
3  【3】把节点对象的文本内容取出来, 查看规律, 是否有更好的处理方法?
4  【4】提取完1页后, 判断如果不是最后1页, 则点击下一页
5      # 如何判断是否为最后1页? ??
```

■ 实现步骤-参考与提示

```
1  # 1. 找节点
2  1、首页搜索框 : /*[@id="key"]
3  2、首页搜索按钮 : /*[@id="search"]/div/div[2]/button
4  3、商品页的 商品信息节点对象列表 : /*[@id="J_goodsList"]/ul/li
5  4、for循环遍历后
6     名称: ./div[@class="p-name"]/a/em
7     价格: ./div[@class="p-price"]
8     评论: ./div[@class="p-commit"]/strong
9     商家: ./div[@class="p-shopnum"]
10
11 # 2. 执行JS脚本, 获取动态加载数据
12 browser.execute_script(
13     'window.scrollTo(0,document.body.scrollHeight)'
14 )
```

■ 代码实现

```
1  browser.excute_script(
2  'window.scrollTo(0,document.body.scrollHeight)'
3  )
4  time.sleep(2)
5
```

```
6  【1】 搜索内容：爬虫书
7      li_list = [<li1>,<li2>,...<lin>]
8      for li in li_list:
9          方法1: print(li.text)
10         方法2: item['name']=li.find_element_by_xpath('')
11
12  【2】一定要注意给页面元素加载预留时间
13
14  【3】执行JS脚本
```

作业2

- 1 【1】多线程改写 - 链家二手房案例
- 2 【2】多线程改写 - 汽车之家案例