

Day05回顾

增量爬取思路

- 1 1、将爬取过的地址存放到数据库中
- 2 2、程序爬取时先到数据库中查询比对，如果已经爬过则不会继续爬取

动态加载网站数据抓取

- 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 # 控制台中 XHR : 异步加载的数据包
- 4 # XHR -> Query String(查询参数)

day06笔记

作业1 - 小米应用商店

将抓取数据保存到csv文件

- 1 注意多线程写入的线程锁问题
- 2 `from threading import Lock`
- 3 `lock = Lock()`
- 4 `lock.acquire()`
- 5 `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 2、每个线程抓取数据后将数据进行文件写入，写入文件时需要加锁
- 3 3、所有数据抓取完成关闭文件

代码实现

```
1 import requests
2 from threading import Thread
3 from queue import Queue
4 import time
5 from fake_useragent import UserAgent
6 from lxml import etree
7 from threading import Lock
8 import csv
9 import random
10
11 class XiaomiSpider(object):
12     def __init__(self):
13         self.url = 'http://app.mi.com/categoryAllListApi?page={}&categoryId={}&pageSize=30'
14         # 存放所有URL地址的队列
15         self.q = Queue()
16         self.ua = UserAgent()
17         self.i = 0
18         # 存放所有类型id的空列表
19         self.id_list = []
20         self.lock = Lock()
21         self.f = open('小米.csv', 'a', encoding='gb18030', newline='')
22         self.writer = csv.writer(self.f)
23
24     def get_cateid(self):
25         # 请求
26         url = 'http://app.mi.com/'
27         headers = { 'User-Agent':self.ua.random }
28         html = requests.get(url=url,headers=headers).text
29         # 解析
30         parse_html = etree.HTML(html)
31         xpath_bds = '//ul[@class="category-list"]/li'
32         li_list = parse_html.xpath(xpath_bds)
33         for li in li_list:
34             typ_name = li.xpath('./a/text()')[0]
35             typ_id = li.xpath('./a/@href')[0].split('/')[1]
36             # 计算每个类型的页数
37             pages = self.get_pages(typ_id)
38             self.id_list.append( (typ_id,pages) )
39
40         # 入队列
41         self.url_in()
42
43     # 获取counts的值并计算页数
44     def get_pages(self,typ_id):
45         # 每页返回的json数据中,都有count这个key
46         url = self.url.format(0,typ_id)
47         html = requests.get(url=url,headers={'User-Agent':self.ua.random}).json()
48         count = html['count']
49         pages = int(count) // 30 + 1
50
51         return pages
52
53     # url入队列
54     def url_in(self):
```

```

55     for id in self.id_list:
56         # id为元组,('2',pages)
57         for page in range(id[1]):
58             url = self.url.format(page,id[0])
59             print(url)
60             # 把URL地址入队列
61             self.q.put(url)
62
63     # 线程事件函数: get() - 请求 - 解析 - 处理数据
64     def get_data(self):
65         while True:
66             if not self.q.empty():
67                 url = self.q.get()
68                 headers = {'User-Agent':self.ua.random}
69                 html = requests.get(url=url,headers=headers).json()
70                 self.parse_html(html)
71                 # 每爬取一页随机休眠0.5秒钟
72                 time.sleep(random.uniform(0,1))
73             else:
74                 break
75
76     # 解析函数
77     def parse_html(self,html):
78         app_info = []
79         for app in html['data']:
80             # 应用名称
81             name = app['displayName']
82             link = 'http://app.mi.com/details?id=' + app['packageName']
83             # 应用类别
84             type_name = app['level1CategoryName']
85             app_info.append([name,type_name,link])
86             print(name,type_name,link)
87             self.i += 1
88
89     # 在文件中写入数据中进行加锁处理
90     self.lock.acquire()
91     self.writer.writerow(app_info)
92     self.lock.release()
93
94
95
96     # 主函数
97     def main(self):
98         # URL入队列
99         self.get_cateid()
100         t_list = []
101         # 创建多个线程
102         for i in range(1):
103             t = Thread(target=self.get_data)
104             t_list.append(t)
105             t.start()
106
107         # 回收线程
108         for t in t_list:
109             t.join()
110
111         print('数量:',self.i)

```

```

112
113     # 所有数据写入完成后释放锁
114     self.f.close()
115
116 if __name__ == '__main__':
117     start = time.time()
118     spider = XiaomiSpider()
119     spider.main()
120     end = time.time()
121     print('执行时间:%.2f' % (end-start))

```

作业2 - 腾讯招聘数据抓取

确定URL地址及目标

- 1、URL：百度搜索腾讯招聘 - 查看工作岗位
- 2、目标：职位名称、工作职责、岗位要求

要求与分析

- 1、通过查看网页源码,得知所需数据均为 Ajax 动态加载
- 2、通过F12抓取网络数据包,进行分析
- 3、一级页面抓取数据: 职位名称
- 4、二级页面抓取数据: 工作职责、岗位要求

一级页面json地址(index在变,timestamp未检查)

- 1 `https://careers.tencent.com/tencentcareer/api/post/Query?timestamp=1563912271089&countryId=&cityId=&bgIds=&productId=&categoryId=&parentId=&attrId=&keyword=&pageIndex={}&pageSize=10&language=zh-cn&area=cn`

二级页面地址(postId在变,在一级页面中可拿到)

- 1 `https://careers.tencent.com/tencentcareer/api/post/ByPostId?timestamp=1563912374645&postId={}&language=zh-cn`

代码实现

```

1 import requests
2 import json
3 import time
4 import random
5 from fake_useragent import UserAgent
6
7 class TencentSpider(object):
8     def __init__(self):

```

```

9         self.one_url = 'https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1563912271089&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=&att
rId=&keyword=&pageIndex={}&pageSize=10&language=zh-cn&area=cn'
10         self.two_url = 'https://careers.tencent.com/tencentcareer/api/post/ByPostId?
timestamp=1563912374645&postId={}&language=zh-cn'
11
12     # 获取User-Agent
13     def get_headers(self):
14         ua = UserAgent()
15         headers = { 'User-Agent': ua.random }
16         return headers
17
18     # 获取响应内容函数
19     def get_page(self,url):
20         html = requests.get(url=url,headers=self.get_headers()).content.decode('utf-8','ignore')
21         # json.loads()把json格式的字符串转为python数据类型
22         html = json.loads(html)
23         return html
24
25     # 主线函数: 获取所有数据
26     def parse_page(self,one_url):
27         html = self.get_page(one_url)
28         item = {}
29         for job in html['Data']['Posts']:
30             item['name'] = job['RecruitPostName']
31             item['address'] = job['LocationName']
32             # 拿postId为了拼接二级页面地址
33             post_id = job['PostId']
34             # 职责和要求(二级页面)
35             two_url = self.two_url.format(post_id)
36             item['duty'],item['requirement'] = self.parse_two_page(two_url)
37             print(item)
38
39     def parse_two_page(self,two_url):
40         html = self.get_page(two_url)
41         # 职责 + 要求
42         duty = html['Data']['Responsibility']
43         requirement = html['Data']['Requirement']
44
45         return duty,requirement
46
47     # 获取总页数
48     def get_pages(self):
49         url = self.one_url.format(1)
50         html = self.get_page(url)
51         pages = int(html['Data']['Count']) // 10 + 1
52
53         return pages
54
55     def main(self):
56         # 总页数
57         pages = self.get_pages()
58         for index in range(1,pages):
59             one_url = self.one_url.format(index)
60             self.parse_page(one_url)
61             time.sleep(random.uniform(0.5,1.5))
62

```



```

40         item['address'] = job['LocationName']
41         # 拿postid为了拼接二级页面地址
42         post_id = job['PostId']
43         # 职责和要求(二级页面)
44         two_url = self.two_url.format(post_id)
45         item['duty'], item['requirement'] = self.parse_two_page(two_url)
46         print(item)
47         self.i += 1
48         # 每爬取按完成1页随机休眠
49         time.sleep(random.uniform(0,1))
50     else:
51         break
52
53 def parse_two_page(self, two_url):
54     html = self.get_page(two_url)
55     # 职责 + 要求
56     duty = html['Data']['Responsibility']
57     requirement = html['Data']['Requirement']
58
59     return duty, requirement
60
61 # 获取总页数
62 def get_pages(self):
63     url = self.one_url.format(1)
64     html = self.get_page(url)
65     pages = int(html['Data']['Count']) // 10 + 1
66
67     return pages
68
69 def main(self):
70     # one_url入队列
71     pages = self.get_pages()
72     for index in range(1, pages):
73         one_url = self.one_url.format(index)
74         self.q.put(one_url)
75
76     t_list = []
77     for i in range(5):
78         t = Thread(target=self.parse_page)
79         t_list.append(t)
80         t.start()
81
82     for t in t_list:
83         t.join()
84
85     print('数量:', self.i)
86
87 if __name__ == '__main__':
88     start = time.time()
89     spider = TencentSpider()
90     spider.main()
91     end = time.time()
92     print('执行时间: %.2f' % (end-start))

```

cookie模拟登录

适用网站及场景

- 1 抓取需要登录才能访问的页面

cookie和session机制

- 1 # http协议为无连接协议
- 2 cookie: 存放在客户端浏览器
- 3 session: 存放在Web服务器

人人网登录案例

■ 方法一 - 登录网站手动抓取Cookie

- 1 1、先登录成功1次,获取到携带登陆信息的Cookie
- 2 登录成功 - 个人主页 - F12抓包 - 刷新个人主页 - 找到主页的包(profile)
- 3 2、携带着cookie发请求
- 4 ** Cookie
- 5 ** User-Agent

```
1 import requests
2
3 class RenRenLogin(object):
4     def __init__(self):
5         # url为需要登录才能正常访问的地址
6         self.url = 'http://www.renren.com/967469305/profile'
7         # headers中的cookie为登录成功后抓取到的cookie
8         self.headers = {
9             # 此处注意cookie, 要自己抓取
10            "Cookie": "xxx",
11            "User-Agent": "Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/74.0.3729.169 Safari/537.36",
12        }
13
14        # 获取个人主页响应
15        def get_html(self):
16            html = requests.get(url=self.url,headers=self.headers,verify=False).text
17            print(html)
18            self.parse_html(html)
19
20        # 可获取并解析整个人人网内需要登录才能访问的地址
21        def parse_html(self,html):
22            pass
23
24 if __name__ == '__main__':
25     spider = RenRenLogin()
26     spider.get_html()
```


▪ 方法二 - requests模块处理Cookie

原理思路及实现

```
1 # 1. 思路
2 requests模块提供了session类,来实现客户端和服务端的会话保持
3
4 # 2. 原理
5 1、实例化session对象
6     session = requests.session()
7 2、让session对象发送get或者post请求
8     res = session.post(url=url,data=data,headers=headers)
9     res = session.get(url=url,headers=headers)
10
11 # 3. 思路梳理
12 浏览器原理: 访问需要登录的页面会带着之前登录过的cookie
13 程序原理: 同样带着之前登录的cookie去访问 - 由session对象完成
14 1、实例化session对象
15 2、登录网站: session对象发送请求,登录对应网站,把cookie保存在session对象中
16 3、访问页面: session对象请求需要登录才能访问的页面,session能够自动携带之前的这个cookie,进行请求
```

具体步骤

```
1 1、寻找Form表单提交地址 - 寻找登录时POST的地址
2     查看网页源码,查看form表单,找action对应的地址: http://www.renren.com/PLogin.do
3
4 2、发送用户名和密码信息到POST的地址
5     * 用户名和密码信息以什么方式发送? -- 字典
6     键 : <input>标签中name的值(email,password)
7     值 : 真实的用户名和密码
8     post_data = {'email':'','password':''}
```

程序实现

```
1 整体思路
2 1、先POST: 把用户名和密码信息POST到某个地址中
3 2、再GET: 正常请求去获取页面信息
```

```
1 import requests
2
3 class RenrenLogin(object):
4     def __init__(self):
5         # 定义常用变量
6         self.post_url = 'http://www.renren.com/PLogin.do'
7         self.get_url = 'http://www.renren.com/967469305/profile'
8         self.post_data = {
9             'email' : 'xxx',
10            'password' : 'xxx'
11        }
12        self.headers = {
13            'User-Agent' : 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
14            like Gecko) Chrome/74.0.3729.169 Safari/537.36',
15            'Referer' : 'http://www.renren.com/SysHome.do'
16        }
```

```

16
17     # 实例化session会话保持对象
18     self.session = requests.session()
19
20     # 先post 再get
21     def post_get_data(self):
22         # 先POST,把用户名和密码信息POST到一个地址
23         self.session.post(url=self.post_url,data=self.post_data,headers=self.headers)
24
25         # 再get个人主页
26         html = self.session.get(url=self.get_url,headers=self.headers).text
27         print(html)
28
29     if __name__ == '__main__':
30         spider = RenrenLogin()
31         spider.post_get_data()

```

▪ 方法三

原理

- 1、把抓取到的cookie处理为字典
- 2、使用requests.get()中的参数:cookies

代码实现

```

1  import requests
2
3  class RenRenLogin(object):
4      def __init__(self):
5          # url为需要登录才能正常访问的地址
6          self.url = 'http://www.renren.com/967469305/profile'
7          self.headers = {
8              "User-Agent": "Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML,
              like Gecko) Chrome/74.0.3729.169 Safari/537.36",
9          }
10
11     # 将字符串cookie转为字典格式
12     def cookie_to_dict(self,cookie_str):
13         cookie_dic = {}
14         for i in cookie_str.split('; '):
15             cookie_dic[i.split('=')[0]] = i.split('=')[1]
16         return cookie_dic
17
18     # 获取个人主页响应
19     def get_html(self):
20         cookie_str = 'xxx'
21         cookie_dict = self.cookie_to_dict(cookie_str)
22         html =
requests.get(url=self.url,headers=self.headers,verify=False,cookies=cookie_dict).text
23         print(html)
24         self.parse_html(html)
25
26     # 可获取并解析整个人人网内需要登录才能访问的地址
27     def parse_html(self,html):

```

```
28     pass
29
30 if __name__ == '__main__':
31     spider = RenRenLogin()
32     spider.get_html()
```

json解析模块

json.loads(json)

■ 作用

```
1 把json格式的字符串转为Python数据类型
```

■ 示例

```
1 html_json = json.loads(res.text)
```

json.dumps(python)

■ 作用

```
1 把 python 类型 转为 json 类型
```

■ 示例

```
1 import json
2
3 # json.dumps()之前
4 item = {'name':'QQ','app_id':1}
5 print('before dumps',type(item))
6 # json.dumps之后
7 item = json.dumps(item)
8 print('after dumps',type(item))
```

json.load(f)

作用

```
1 将json文件读取,并转为python类型
```

示例

```
1 import json
2
3 with open('D:\\spider_test\\xiaomi.json','r') as f:
4     data = json.load(f)
5
6 print(data)
```

json.dump(python,f,ensure_ascii=False)

■ 作用

```
1 把python数据类型 转为 json格式的字符串
2 # 一般让你把抓取的数据保存为json文件时使用
```

■ 参数说明

```
1 第1个参数: python类型的数据(字典, 列表等)
2 第2个参数: 文件对象
3 第3个参数: ensure_ascii=False # 序列化时编码
```

■ 示例1

```
1 import json
2
3 item = {'name':'QQ','app_id':1}
4 with open('小米.json','a') as f:
5     json.dump(item,f,ensure_ascii=False)
```

■ 示例2

```
1 import json
2
3 item_list = []
4 for i in range(3):
5     item = {'name':'QQ','id':i}
6     item_list.append(item)
7
8 with open('xiaomi.json','a') as f:
9     json.dump(item_list,f,ensure_ascii=False)
```

练习: 将腾讯招聘数据存入到json文件

```

1  # 1. __init__()
2      self.f = open('tencent.json', 'a')
3      self.item_list = []
4  # 2. parse_page()
5      self.item_list.append(item)
6  # 3. main()
7      json.dump(self.item_list, self.f, ensure_ascii=False)
8      self.f.close()

```

json模块总结

```

1  # 爬虫最常用
2  1、数据抓取 - json.loads(html)
3      将响应内容由: json 转为 python
4  2、数据保存 - json.dump(item_list, f, ensure_ascii=False)
5      将抓取的数据保存到本地 json文件
6
7  # 抓取数据一般处理方式
8  1、txt文件
9  2、csv文件
10 3、json文件
11 4、MySQL数据库
12 5、MongoDB数据库
13 6、Redis数据库

```

selenium+phantomjs/Chrome/Firefox

selenium

▪ 定义

- 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
- 2 只是工具，必须与第三方浏览器结合使用

▪ 安装

- 1 Linux: `sudo pip3 install selenium`
- 2 Windows: `python -m pip install selenium`

*phantomjs*浏览器

▪ 定义

- 1 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

■ 安装(phantomjs、chromedriver、geckodriver)

Windows

```
1 1、下载对应版本的phantomjs、chromedriver、geckodriver
2 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)
3   # 查看python安装路径: where python
4 3、验证
5   cmd命令行: chromedriver
6
7   # 下载地址
8 1、chromedriver : 下载对应版本
9   http://chromedriver.storage.googleapis.com/index.html
10 2、geckodriver
11   https://github.com/mozilla/geckodriver/releases
12 3、phantomjs
13   https://phantomjs.org/download.html
```

Linux

```
1 1、下载后解压
2   tar -zxvf geckodriver.tar.gz
3 2、拷贝解压后文件到 /usr/bin/ (添加环境变量)
4   sudo cp geckodriver /usr/bin/
5 3、更改权限
6   sudo -i
7   cd /usr/bin/
8   chmod 777 geckodriver
```

■ 使用

示例代码一：使用 selenium+浏览器 打开百度

```
1 from selenium import webdriver
2 import time
3
4 browser = webdriver.Chrome()
5 browser.get('http://www.baidu.com/')
6 browser.save_screenshot('baidu.png')
7 browser.quit()
```

示例代码二：打开百度，搜索赵丽颖，查看

```
1 from selenium import webdriver
2 import time
3
4 browser = webdriver.Chrome()
5 browser.get('http://www.baidu.com/')
6
7 # 向搜索框(id kw)输入 赵丽颖
8 ele = browser.find_element_by_xpath('//*[@id="kw"]')
9 ele.send_keys('赵丽颖')
10
11 time.sleep(1)
```

```

12 # 点击 百度一下 按钮(id su)
13 su = browser.find_element_by_xpath('//*[@id="su"]')
14 su.click()
15
16 # 截图
17 browser.save_screenshot('赵丽颖.png')
18 # 关闭浏览器
19 browser.quit()

```

■ 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 # 查看响应内容
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.quit() # 关闭浏览器

```

■ 定位节点

单元素查找(1个节点对象)

```

1 1、 browser.find_element_by_id('')
2 2、 browser.find_element_by_name('')
3 3、 browser.find_element_by_class_name('')
4 4、 browser.find_element_by_xpath('')
5 ... ..

```

多元素查找([节点对象列表])

```

1 1、 browser.find_elements_by_id('')
2 2、 browser.find_elements_by_name('')
3 3、 browser.find_elements_by_class_name('')
4 4、 browser.find_elements_by_xpath('')
5 ... ..

```

■ 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
4 4、 ele.get_attribute('src') # 获取属性值

```

京东爬虫案例

■ 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，是否有更好的处理办法？
- 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

1. 找节点

- 1 1、首页搜索框：//*[@id="key"]
- 2 2、首页搜索按钮：//*[@id="search"]/div/div[2]/button
- 3 3、商品页的商品信息节点对象列表：//*[@id="J_goodsList"]/ul/li

2. 执行JS脚本，获取动态加载数据

```
1 browser.execute_script(  
2     'window.scrollTo(0,document.body.scrollHeight)'  
3 )
```

3. 代码实现

```
1 from selenium import webdriver  
2 import time  
3  
4 class JdSpider(object):  
5     def __init__(self):  
6         self.browser = webdriver.Chrome()  
7         self.url = 'https://www.jd.com/'  
8         self.i = 0  
9  
10        # 获取商品页面  
11        def get_page(self):  
12            self.browser.get(self.url)  
13            # 找2个节点  
14            self.browser.find_element_by_xpath('//*[@id="key"]').send_keys('爬虫书籍')  
15            self.browser.find_element_by_xpath('//*[@id="search"]/div/div[2]/button').click()  
16            time.sleep(2)  
17  
18        # 解析页面  
19        def parse_page(self):  
20            # 把下拉菜单拉到底部,执行JS脚本  
21            self.browser.execute_script(  
22                'window.scrollTo(0,document.body.scrollHeight)'  
23            )  
24            time.sleep(2)  
25            # 匹配所有商品节点对象列表  
26            li_list = self.browser.find_elements_by_xpath('//*[@id="J_goodsList"]/ul/li')
```



```

27         for li in li_list:
28             li_info = li.text.split('\n')
29             if li_info[0][0:2] == '每满':
30                 price = li_info[1]
31                 name = li_info[2]
32                 commit = li_info[3]
33                 market = li_info[4]
34             else:
35                 price = li_info[0]
36                 name = li_info[1]
37                 commit = li_info[2]
38                 market = li_info[3]
39             print('\033[31m*****\033[0m')
40             print(price)
41             print(commit)
42             print(market)
43             print(name)
44             self.i += 1
45
46     # 主函数
47     def main(self):
48         self.get_page()
49         while True:
50             self.parse_page()
51             # 判断是否该点击下一页,没有找到说明不是最后一页
52             if self.browser.page_source.find('pn-next disabled') == -1:
53                 self.browser.find_element_by_class_name('pn-next').click()
54                 time.sleep(2)
55             else:
56                 break
57         print(self.i)
58
59 if __name__ == '__main__':
60     spider = JdSpider()
61     spider.main()

```

chromedriver设置无界面模式

```

1  from selenium import webdriver
2
3  options = webdriver.ChromeOptions()
4  # 添加无界面参数
5  options.add_argument('--headless')
6  browser = webdriver.Chrome(options=options)
7  browser.get('http://www.baidu.com/')
8  browser.save_screenshot('baidu.png')

```

作业

- 1、使用selenium+浏览器抓取 民政部 数据
- 2、尝试去破解一下百度翻译