

Day05回顾

增量爬取思路

- 1、将爬取过的地址存放到数据库中
- 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、F12打开控制台，页面动作抓取网络数据包
- 2、抓取json文件URL地址
- # 控制台中 XHR : 异步加载的数据包
- # XHR -> Query String(查询参数)

多线程爬虫

■ 使用流程

```
1 # 1、URL队列
2 q.put(url)
3 # 2、线程事件函数
4 while True:
5     if not url_queue.empty():
6         ...get()、请求、解析
7     else:
8         break
9 # 创建并启动线程
10 t_list = []
11 for i in range(5):
12     t = Thread(target=parse_page)
13     t_list.append(t)
14     t.start()
15 # 阻塞等待回收线程
16 for i in t_list:
17     i.join()
```

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、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、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、所有数据抓取完成关闭文件

代码实现

```
1 |
```

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=&parentCategoryId=&attr
Id=&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 |
```

多线程有什么思路?

```
1 把所有一级页面链接提交到队列,进行多线程数据抓取
```

代码实现

```
1 |
```

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 |
```

■ 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

```

方法三

原理

```

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

```

代码实现

```

1

```

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 |
```

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 |
```

示例代码二: 打开百度, 搜索赵丽颖, 查看

```
1 |
```

■ 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

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 |
```

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