

Day04回顾

▪ requests.get()参数

```
1 【1】 url
2 【2】 proxies -> {}
3     proxies = {
4         'http': 'http://1.1.1.1:8888',
5         'https': 'https://1.1.1.1:8888'
6     }
7 【3】 timeout
8 【4】 headers
```

▪ requests.post()

```
1 data : 字典, Form表单数据
```

▪ 常见的反爬机制及处理方式

```
1 【1】 Headers反爬虫
2     1.1) 检查: Cookie、Referer、User-Agent
3     1.2) 解决方案: 通过F12获取headers,传给requests.get()方法
4
5 【2】 IP限制
6     2.1) 网站根据IP地址访问频率进行反爬,短时间内限制IP访问
7     2.2) 解决方案:
8         a) 构造自己IP代理池,每次访问随机选择代理,经常更新代理池
9         b) 购买开放代理或私密代理IP
10        c) 降低爬取的速度
11
12 【3】 User-Agent限制
13     3.1) 类似于IP限制, 检测频率
14     3.2) 解决方案: 构造自己的User-Agent池,每次访问随机选择
15         a> fake_useragent模块
16         b> 新建py文件,存放大量User-Agent
17         c> 程序中定义列表,存放大量的User-Agent
18
19 【4】 对响应内容做处理
20     4.1) 页面结构和响应内容不同
21     4.2) 解决方案: 打印并查看响应内容,用xpath或正则做处理
22
23 【5】 JS加密 - 有道翻译案例
24     5.1) 抓取到对应的JS文件,寻找加密算法
25     5.2) 用Python实现加密算法,生成指定的参数
26
```

```
27 【6】JS逆向 - 百度翻译案例
28     6.1) 抓取到对应的JS文件, 寻找加密算法
29     6.2) 利用pyexecjs模块执行JS代码
30
31 【7】动态加载 - 豆瓣电影、小米应用商店、腾讯招聘
32     7.1) F12抓取网络数据包再进行分析, 找返回实际json数据的地址
```

■ 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 : 找到异步加载的数据包
5     2.2) GET请求: Network -> XHR -> URL 和 Query String Parameters(查询参数)
6     2.3) POST请求: Network -> XHR -> URL 和 Form Data
```

■ 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存在
3
4     1.2) 分析页面结构, 观察URL地址规律
5         大体查看响应内容结构, 查看是否有更改, 以响应内容为准写正则或xpath表达式
6
7     1.3) 开始写代码进行数据抓取
8
9 【2】响应内容中不存在
10    2.1) 确认抓取数据在响应内容中不存在
11
12    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
13        a) GET请求: Request URL、Request Headers、Query String Parameters
14        b) POST请求: Request URL、Request Headers、FormData
15
16    2.3) 观察查询参数或者Form表单数据规律, 如果需要进一步抓包分析处理
17        a) 比如有道翻译的 salt+sign, 抓取并分析JS做进一步处理
18        b) 此处注意请求头中的Cookie和Referer以及User-Agent
19
20    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         self.lock = Lock()
17
18     def url_in(self):
19         """生成待爬取的URL地址,入队列"""
20         pass
21
22     def parse_html(self):
23         """线程事件函数,获取地址,进行数据抓取"""
24         while True:
25             self.lock.acquire()
26             if not self.q.empty():
27                 url = self.q.get()
28                 self.lock.release()
29             else:
30                 self.lock.release()
31                 break
32
33     def run(self):
34         self.url_in()
35         t_list = []
36         for i in range(3):
37             t = Thread(target=self.parse_html)
38             t_list.append(t)
39             t.start()
40
41         for th in t_list:
42             th.join()
43
44     【4】队列要点: q.get()防止阻塞方式
45         4.1) 方法1: q.get(block=False)
46         4.2) 方法2: q.get(block=True, timeout=3)
47         4.3) 方法3:
48             if not q.empty():
49                 q.get()

```

Day05笔记

多线程抓取百度图片

目标

- 1 【1】多线程实现抓取百度图片指定关键字的所有图片
- 2 【2】URL地址: <https://image.baidu.com/>
- 3 【3】实现效果
- 4 请输入关键字: 赵丽颖
- 5 则: 在当前目录下创建 ./images/赵丽颖/ 目录,并将所有赵丽颖图片保存到此目录下

■ 实现步骤

```
1  【1】 输入关键字,进入图片页面,滚动鼠标滑轮发现所有图片链接均为动态加载
2  【2】 F12抓取到返回实际数据的json地址如下:
3      https://image.baidu.com/search/acjson?
tn=resultjson_com&ipn=rj&ct=201326592&is=&fp=result&queryWord=
{}&cl=2&lm=-1&hd=&latest=&copyright=&ie=utf-8&oe=utf-8&adpicid=&st=-1&z=&ic=0&word=
{}&s=&se=&tab=&width=&height=&face=0&istype=2&qc=&nc=1&fr=&pn=
{}&rn=30&gsm=1e&1599728538999=
4
5  【3】 分析有效的查询参数如下:
6      queryWord: 赵丽颖
7      word: 赵丽颖
8      pn: 30    pn的值为 0 30 60 90 120 ... ...
9
10 【4】 图片总数如何获取?
11     在json数据中能够查到图片总数
12     displayNum: 1037274
```

■ 代码实现

```
1  import requests
2  from threading import Thread, Lock
3  from queue import Queue
4  from fake_useragent import UserAgent
5  import os
6  from urllib import parse
7  import time
8
9  class BaiduImageSpider:
10     def __init__(self):
11         self.keyword = input('请输入关键字:')
12         self.directory = './images/{}/'.format(self.keyword)
13         if not os.path.exists(self.directory):
14             os.makedirs(self.directory)
15         # 编码
16         self.params = parse.quote(self.keyword)
17         self.url = 'https://image.baidu.com/search/acjson?
tn=resultjson_com&ipn=rj&ct=201326592&is=&fp=result&queryWord=
{}&cl=2&lm=-1&hd=&latest=&copyright=&ie=utf-8&oe=utf-8&adpicid=&st=-1&z=&ic=0&word=
{}&s=&se=&tab=&width=&height=&face=0&istype=2&qc=&nc=1&fr=&pn=
{}&rn=30&gsm=5a&1599724879800='
18         # 创建url队列、线程锁
19         self.url_queue = Queue()
20         self.lock = Lock()
21
22     def get_html(self, url):
23         """请求功能函数"""
24         headers = {'User-Agent': UserAgent().random}
25         html = requests.get(url=url, headers=headers).json()
26
27         return html
28
29     def get_total_image(self):
30         """获取图片总数量"""
31         total_page_url = self.url.format(self.params, self.params, 0)
```

```

32     total_page_html = self.get_html(url=total_page_url)
33     total_page_nums = total_page_html['displayNum']
34
35     return total_page_nums
36
37 def url_in(self):
38     total_image_nums = self.get_total_image()
39     for pn in range(0, total_image_nums, 30):
40         page_url = self.url.format(self.params, self.params, pn)
41         self.url_queue.put(page_url)
42
43 def parse_html(self):
44     """线程事件函数"""
45     while True:
46         # 加锁
47         self.lock.acquire()
48         if not self.url_queue.empty():
49             page_url = self.url_queue.get()
50             # 释放锁
51             self.lock.release()
52             html = self.get_html(url=page_url)
53             for one_image_dict in html['data']:
54                 try:
55                     image_url = one_image_dict['hoverURL']
56                     self.save_image(image_url)
57                 except Exception as e:
58                     continue
59             else:
60                 self.lock.release()
61                 break
62
63 def save_image(self, image_url):
64     """保存一张图片到本地"""
65     headers = {'User-Agent': UserAgent().random}
66     image_html = requests.get(url=image_url, headers=headers).content
67     # 加锁、释放锁
68     self.lock.acquire()
69     filename = self.directory + image_url[-24:]
70     self.lock.release()
71     with open(filename, 'wb') as f:
72         f.write(image_html)
73     print(filename, '下载成功')
74
75 def run(self):
76     # 让URL地址入队列
77     self.url_in()
78     # 创建多线程
79     t_list = []
80     for i in range(1):
81         t = Thread(target=self.parse_html)
82         t_list.append(t)
83         t.start()
84
85     for t in t_list:
86         t.join()
87
88 if __name__ == '__main__':

```

```
89     start_time = time.time()
90     spider = BaiduImageSpider()
91     spider.run()
92     end_time = time.time()
93     print('time:%.2f' % (end_time - start_time))
```

selenium+PhantomJS/Chrome/Firefox

■ selenium

```
1  【1】定义
2      1.1) 开源的Web自动化测试工具
3
4  【2】用途
5      2.1) 对Web系统进行功能性测试, 版本迭代时避免重复劳动
6      2.2) 兼容性测试(测试web程序在不同操作系统和不同浏览器中是否运行正常)
7      2.3) 对web系统进行大数量测试
8
9  【3】特点
10     3.1) 可根据指令操控浏览器
11     3.2) 只是工具, 必须与第三方浏览器结合使用
12
13  【4】安装
14     4.1) Linux: sudo pip3 install selenium
15     4.2) Windows: python -m pip install selenium
```

■ PhantomJS浏览器

```
1  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
```

■ 环境安装

```
1  【1】下载驱动
2      2.1) chromedriver : 下载对应版本
3          http://npm.taobao.org/mirrors/chromedriver/
4      2.2) geckodriver
5          https://github.com/mozilla/geckodriver/releases
6      2.3) phantomjs
7          https://phantomjs.org/download.html
8
9  【2】添加到系统环境变量
10     2.1) windows: 拷贝到Python安装目录的Scripts目录中
11             windows查看python安装目录(cmd命令行): where python
12     2.2) Linux : 拷贝到/usr/bin目录中 : sudo cp chromedriver /usr/bin/
13
14  【3】Linux中需要修改权限
15     sudo chmod 777 /usr/bin/chromedriver
16
17  【4】验证
18     4.1) Ubuntu | Windows
19         from selenium import webdriver
20         webdriver.Chrome()
```

```

21         webdriver.Firefox()
22
23     4.2) Mac
24         from selenium import webdriver
25         webdriver.Chrome(executable_path='/Users/xxx/chromedriver')
26         webdriver.Firefox(executable_path='/User/xxx/geckodriver')

```

■ 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  driver = webdriver.Chrome()
9  driver.get('http://www.baidu.com/')
10 # 5秒钟后关闭浏览器
11 time.sleep(5)
12 driver.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  driver = webdriver.Chrome()
8  # 2.输入：http://www.baidu.com/
9  driver.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 driver.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 driver.find_element_by_xpath('//*[@id="su"]').click()

```

■ 浏览器对象(browser)方法

```

1  【1】 driver.get(url=url)    - 地址栏输入url地址并确认
2  【2】 driver.quit()          - 关闭浏览器
3  【3】 driver.close()         - 关闭当前页
4  【4】 driver.page_source     - HTML结构源码
5  【5】 driver.page_source.find('字符串')
6      从html源码中搜索指定字符串,没有找到返回: -1,经常用于判断是否为最后一页
7  【6】 driver.maximize_window() - 浏览器窗口最大化

```

■ 定位节点八种方法

```

1  【1】 单元素查找('结果为1个节点对象')
2      1.1) 【最常用】 driver.find_element_by_id('id属性值')
3      1.2) 【最常用】 driver.find_element_by_name('name属性值')
4      1.3) 【最常用】 driver.find_element_by_class_name('class属性值')
5      1.4) 【最万能】 driver.find_element_by_xpath('xpath表达式')

```

```

6     1.5) 【匹配a节点时常用】driver.find_element_by_link_text('链接文本')
7     1.6) 【匹配a节点时常用】driver.find_element_by_partical_link_text('部分链接文本')
8     1.7) 【最没用】driver.find_element_by_tag_name('标记名称')
9     1.8) 【较常用】driver.find_element_by_css_selector('css表达式')
10
11 【2】多元素查找('结果为[节点对象列表]')
12     2.1) driver.find_elements_by_id('id属性值')
13     2.2) driver.find_elements_by_name('name属性值')
14     2.3) driver.find_elements_by_class_name('class属性值')
15     2.4) driver.find_elements_by_xpath('xpath表达式')
16     2.5) driver.find_elements_by_link_text('链接文本')
17     2.6) driver.find_elements_by_partical_link_text('部分链接文本')
18     2.7) driver.find_elements_by_tag_name('标记名称')
19     2.8) driver.find_elements_by_css_selector('css表达式')
20
21 【3】注意
22     当属性值中存在 空格 时,我们要使用 . 去代替空格
23     页面中class属性值为: btn btn-account
24     driver.find_element_by_class_name('btn.btn-account').click()

```

■ 猫眼电影示例

```

1  from selenium import webdriver
2  import time
3
4  url = 'https://maoyan.com/board/4'
5  driver = webdriver.Chrome()
6  driver.get(url)
7
8  def get_data():
9      # 基准xpath: [<selenium xxx li at xxx>,<selenium xxx li at>]
10     li_list = driver.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         driver.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         driver.quit()
31         break

```

■ 节点对象操作


```
1 【1】 node.send_keys('') - 向文本框发送内容
2 【2】 node.click() - 点击
3 【3】 node.get_attribute('属性名') - 获取节点的属性值
```

*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 driver = webdriver.Chrome(options=options)
```

selenium - 鼠标操作

```
1 from selenium import webdriver
2 # 导入鼠标事件类
3 from selenium.webdriver import ActionChains
4
5 driver = webdriver.Chrome()
6 driver.get('http://www.baidu.com/')
7
8 # 移动到 设置, perform()是真正执行操作, 必须有
9 element = driver.find_element_by_xpath('//*[@id="u1"]/a[8]')
10 ActionChains(driver).move_to_element(element).perform()
11
12 # 单击, 弹出的Ajax元素, 根据链接节点的文本内容查找
13 driver.find_element_by_link_text('高级搜索').click()
```

selenium - 切换页面

■ 适用网站+应对方案

```
1 【1】 适用网站类型
2 页面中点开链接出现新的窗口, 但是浏览器对象driver还是之前页面的对象, 需要切换到不同的窗口进行
  操作
3
4 【2】 应对方案 - driver.switch_to.window()
5
6 # 获取当前所有句柄 (窗口) - [handle1,handle2]
7 all_handles = driver.window_handles
8 # 切换browser到新的窗口, 获取新窗口的对象
9 driver.switch_to.window(all_handles[1])
```

■ 民政部网站案例-selenium

```
1  """
2  适用selenium+Chrome抓取民政部行政区划代码
3  http://www.mca.gov.cn/article/sj/xzqh/2019/
4  """
5  from selenium import webdriver
6
7  class GovSpider(object):
8      def __init__(self):
9          # 设置无界面
10         options = webdriver.ChromeOptions()
11         options.add_argument('--headless')
12         # 添加参数
13         self.driver = webdriver.Chrome(options=options)
14         self.one_url = 'http://www.mca.gov.cn/article/sj/xzqh/2019/'
15
16     def get_incr_url(self):
17         self.driver.get(self.one_url)
18         # 提取最新链接节点对象并点击
19
20     self.driver.find_element_by_xpath('//td[@class="arlisttd"]/a[contains(@title,"代
21     码")]').click()
22     # 切换句柄
23     all_handlers = self.driver.window_handles
24     self.driver.switch_to.window(all_handlers[1])
25     self.get_data()
26
27     def get_data(self):
28         tr_list = self.driver.find_elements_by_xpath('//tr[@height="19"]')
29         for tr in tr_list:
30             code = tr.find_element_by_xpath('./td[2]').text.strip()
31             name = tr.find_element_by_xpath('./td[3]').text.strip()
32             print(name,code)
33
34     def run(self):
35         self.get_incr_url()
36         self.driver.quit()
37
38 if __name__ == '__main__':
39     spider = GovSpider()
40     spider.run()
```

selenium - iframe

■ 特点+方法

- 1 【1】 特点
- 2 网页中嵌套了网页，先切换到iframe，然后再执行其他操作
- 3
- 4 【2】 处理步骤
- 5 2.1) 切换到要处理的Frame

```

6      2.2) 在Frame中定位页面元素并进行操作
7      2.3) 返回当前处理的Frame的上一级页面或主页面
8
9      【3】常用方法
10     3.1) 切换到frame - driver.switch_to.frame(frame节点对象)
11     3.2) 返回上一级 - driver.switch_to.parent_frame()
12     3.3) 返回主页面 - driver.switch_to.default_content()
13
14     【4】使用说明
15     4.1) 方法一：默认支持id和name属性值：switch_to.frame(id属性值|name属性值)
16     4.2) 方法二：
17         a> 先找到frame节点：frame_node = driver.find_element_by_xpath('xxx')
18         b> 在切换到frame：driver.switch_to.frame(frame_node)

```

■ 示例1 - 登录豆瓣网

```

1  """
2  登录豆瓣网
3  """
4  from selenium import webdriver
5  import time
6
7  # 打开豆瓣官网
8  driver = webdriver.Chrome()
9  driver.get('https://www.douban.com/')
10
11 # 切换到iframe子页面
12 login_frame = driver.find_element_by_xpath('//*[id="anony-reg-
13 new"]/div/div[1]/iframe')
14 driver.switch_to.frame(login_frame)
15
16 # 密码登录 + 用户名 + 密码 + 登录豆瓣
17 driver.find_element_by_xpath('/html/body/div[1]/div[1]/ul[1]/li[2]').click()
18 driver.find_element_by_xpath('//*[id="username"]').send_keys('自己的用户名')
19 driver.find_element_by_xpath('//*[id="password"]').send_keys('自己的密码')
20 driver.find_element_by_xpath('/html/body/div[1]/div[2]/div[1]/div[5]/a').click()
21 time.sleep(3)
22
23 # 点击我的豆瓣
24 driver.find_element_by_xpath('//*[id="db-nav-sns"]/div/div/div[3]/ul/li[2]/a').click()

```

■ selenium+phantomjs|chrome|firefox小总结

```

1  【1】设置无界面模式
2      options = webdriver.ChromeOptions()
3      options.add_argument('--headless')
4      driver =
5  webdriver.Chrome(executable_path='/home/tarena/chromedriver',options=options)
6
7  【2】browser执行JS脚本
8      driver.execute_script('window.scrollTo(0,document.body.scrollHeight)')
9
10 【3】鼠标操作
11     from selenium.webdriver import ActionChains
12     ActionChains(driver).move_to_element('node').perform()

```

```
12
13 【4】 切换句柄 - switch_to.frame(handle)
14     all_handles = driver.window_handles
15     driver.switch_to.window(all_handles[1])
16
17 【5】 iframe子页面
18     driver.switch_to.frame(frame_node)
```

▪ lxml中的xpath 和 selenium中的xpath的区别

```
1  【1】 lxml中的xpath用法 - 推荐自己手写
2      div_list = p.xpath('//div[@class="abc"]/div')
3      item = {}
4      for div in div_list:
5          item['name'] = div.xpath('./a/@href')[0]
6          item['likes'] = div.xpath('./a/text()')[0]
7
8  【2】 selenium中的xpath用法 - 推荐copy - copy xpath
9      div_list = browser.find_elements_by_xpath('//div[@class="abc"]/div')
10     item = {}
11     for div in div_list:
12         item['name'] = div.find_element_by_xpath('./a').get_attribute('href')
13         item['likes'] = div.find_element_by_xpath('./a').text
```

今日作业

- 1 【1】 使用selenium+浏览器 获取有道翻译结果
- 2 【2】 使用selenium+浏览器 登录网易qq邮箱 : <https://mail.qq.com/>
- 3 【3】 使用selenium+浏览器 登录网易163邮箱 : <https://mail.163.com/>