

Day04回顾

requests.get()参数

```
1 1、url
2 2、params -> {} : 查询参数 Query String
3 3、proxies -> {}
4     proxies = {
5         'http': 'http://username:password@1.1.1.1:8888',
6         'https': 'https://username:password@1.1.1.1:8888'
7     }
8 4、auth -> ('tarenacode', 'code_2013')
9 5、verify -> True/False
10 6、timeout
```

requests.post()

```
1 data -> {} Form表单数据 : Form Data
```

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 2、控制台常用选项
3     1、Network：抓取网络数据包
4         1、ALL：抓取所有的网络数据包
5         2、XHR：抓取异步加载的网络数据包
6         3、JS：抓取所有的JS文件
7     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
8     3、Console：交互模式，可对JavaScript中的代码进行测试
9 3、抓取具体网络数据包后
10     1、单击左侧网络数据包地址，进入数据包详情，查看右侧
11     2、右侧：
12         1、Headers：整个请求信息
13             General、Response Headers、Request Headers、Query String、Form Data
14         2、Preview：对响应内容进行预览
15         3、Response：响应内容
```

■ 有道翻译过程梳理

1.
 - 1 1. 打开首页
 - 2 2. 准备抓包: F12开启控制台
 - 3 3. 寻找地址
 - 4 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
 - 5 4. 发现规律
 - 6 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, 分析对比 Network - All(或者XHR) - Form Data, 发现对应的规律
 - 7 5. 寻找JS文件
 - 8 右上角 ... -> Search -> 搜索关键字 -> 单击 -> 跳转到Sources, 左下角格式化符号{}
 - 9 6. 查看JS代码
 - 10 搜索关键字, 找到相关加密方法
 - 11 7. 断点调试
 - 12 8. 完善程序

常见的反爬机制及处理方式

- 1 1. Headers反爬虫 : Cookie、Referer、User-Agent
- 2 解决方案: 通过F12获取headers, 传给requests.get()方法
- 3
- 4 2. IP限制 : 网站根据IP地址访问频率进行反爬, 短时间内进制IP访问
- 5 解决方案:
- 6 1、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
- 7 2、购买开放代理或私密代理IP
- 8 3、降低爬取的速度
- 9
- 10 3. User-Agent限制 : 类似于IP限制
- 11 解决方案: 构造自己的User-Agent池, 每次访问随机选择
- 12
- 13 5. 对查询参数或Form表单数据认证(salt、sign)
- 14 解决方案: 找到JS文件, 分析JS处理方法, 用Python按同样方式处理
- 15
- 16 6. 对响应内容做处理
- 17 解决方案: 打印并查看响应内容, 用xpath或正则做处理

python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

Day05笔记

有道翻译代码实现

有道翻译验证了什么？ - headers

- 1 1、Cookie
- 2 2、Referer
- 3 3、User-Agent

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            "Cookie": "OUTFOX_SEARCH_USER_ID=970246104@10.169.0.83;
OUTFOX_SEARCH_USER_ID_NCOO=570559528.1224236;
_ntes_nnid=96bc13a2f5ce64962adfd6a278467214,1551873108952; JSESSIONID=aaaae9i7p1XP1KaJH_gkYw;
td_cookie=18446744072941336803; SESSION_FROM_COOKIE=unknown;
__rl_test_cookies=1565689460872",
12            "Referer": "http://fanyi.youdao.com/",
13            "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",
14        }
15
16        # 获取salt,sign,ts
17        def get_salt_sign_ts(self,word):
18            # salt
19            salt = str(int(time.time()*1000)) + str(random.randint(0,9))
20            # sign
21            string = "fanyideskweb" + word + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
22            s = md5()
23            s.update(string.encode())
24            sign = s.hexdigest()
25            # ts
26            ts = str(int(time.time()*1000))
27            return salt,sign,ts
28
29        # 主函数
30        def attack_yd(self,word):
31            # 1. 先拿到salt,sign,ts
32            salt, sign, ts = self.get_salt_sign_ts(word)
33            # 2. 定义form表单数据为字典: data={}
34            data = {
35                'i': word,
```

```

36         'from': 'AUTO',
37         'to': 'AUTO',
38         'smartresult': 'dict',
39         'client': 'fanyideskweb',
40         'salt': salt,
41         'sign': sign,
42         'ts': ts,
43         'bv': 'cf156b581152bd0b259b90070b1120e6',
44         'doctype': 'json',
45         'version': '2.1',
46         'keyfrom': 'fanyi.web',
47         'action': 'FY_BY_REALTIME'
48     }
49     # 3. 直接发请求:requests.post(url,data=data,headers=xxx)
50     json_html = requests.post(self.url, data=data, headers=self.headers).json()
51     result = json_html['translateResult'][0][0]['tgt']
52     return result
53
54     # 主函数
55     def main(self):
56         # 输入翻译单词
57         word = input('请输入要翻译的单词: ')
58         result = self.attack_yd(word)
59         print('翻译结果:', result)
60
61     if __name__ == '__main__':
62         spider = YdSpider()
63         spider.main()

```

民政部网站数据抓取

目标

- 1、URL: <http://www.mca.gov.cn/> - 民政数据 - 行政区划代码
即: <http://www.mca.gov.cn/article/sj/xzqh/2019/>
- 2、目标: 抓取最新中华人民共和国县以上行政区划代码

实现步骤

■ 1、从民政数据网站中提取最新行政区划代码链接

```

1 # 特点
2 1、最新的在上面
3 2、命名格式: 2019年x月中华人民共和国县以上行政区划代码
4 # 代码实现
5 import requests
6 from lxml import etree
7 import re
8
9 url = 'http://www.mca.gov.cn/article/sj/xzqh/2019/'
10 headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/74.0.3729.169 Safari/537.36'}

```

```

11 html = requests.get(url, headers=headers).text
12 parse_html = etree.HTML(html)
13 article_list = parse_html.xpath('//a[@class="artitlelist"]')
14
15 for article in article_list:
16     title = article.xpath('./@title')[0]
17     # 正则匹配title中包含这个字符串的链接
18     if title.endswith('代码'):
19         # 获取到第1个就停止即可, 第1个永远是最新的链接
20         two_link = 'http://www.mca.gov.cn' + article.xpath('./@href')[0]
21         print(two_link)
22         break

```

■ 2、从二级页面链接中提取真实链接（反爬-响应内容中嵌入JS，指向新的链接）

```

1 1、向二级页面链接发请求得到响应内容，并查看嵌入的JS代码
2 2、正则提取真实的二级页面链接
3 # 相关思路代码
4 two_html = requests.get(two_link, headers=headers).text
5 # 从二级页面的响应中提取真实的链接（此处为JS动态加载跳转的地址）
6 new_two_link = re.findall(r'window.location.href="(.*?)"', html2, re.S)[0]

```

■ 3、在数据库表中查询此条链接是否已经爬取，建立增量爬虫

```

1 1、数据库中建立version表，存储爬取的链接
2 2、每次执行程序时和version表中记录核对，查看是否已经爬取过
3 # 思路代码
4 cursor.execute('select * from version')
5 result = self.cursor.fetchall()
6 if result:
7     if result[-1][0] == two_link:
8         print('已是最新')
9     else:
10         # 有更新，开始抓取
11         # 将链接再重新插入version表记录

```

■ 4、代码实现

```

1 '''民政部网站数据抓取（增量爬虫）'''
2 import requests
3 from lxml import etree
4 import re
5 import pymysql
6
7 class Govement(object):
8     def __init__(self):
9         self.one_url = 'http://www.mca.gov.cn/article/sj/xzqh/2019/'
10        self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/74.0.3729.169 Safari/537.36'}
11        self.db = pymysql.connect('192.168.153.138', 'tiger', '123456', 'govdb')
12        self.cursor = self.db.cursor()
13
14        # 获取二级页面真实链接，并和数据库中比对
15        def get_flase_link(self):

```

```

16     html = requests.get(self.one_url, headers=self.headers).text
17     # 此处隐藏了真实的二级页面的url链接, 通过js脚本生成, 保存本地文件查看
18     parse_html = etree.HTML(html)
19     a_list = parse_html.xpath('//a[@class="artitlelist"]')
20     for a in a_list:
21         title = a.get('title')
22         # 正则匹配title中包含这个字符串的链接
23         if title.endswith('代码'):
24             # 获取到第1个就停止即可, 第1个永远是最新的链接
25             false_link = 'http://www.mca.gov.cn' + a.get('href')
26             break
27
28     # 提取二级页面真实链接
29     self.get_real_link(false_link)
30
31 def get_real_link(self, false_link):
32     # 从已提取的two_link中提取二级页面的真实链接
33     two_html = requests.get(false_link, headers=self.headers).text
34     # 从二级页面的响应中提取真实的链接 (此处为JS动态加载跳转的地址)
35     real_two_link = re.findall(r'window.location.href="(.*?)"', two_html, re.S)[0]
36
37     self.incre_spider(real_two_link)
38
39 def incre_spider(self, real_two_link):
40     # 实现增量爬取
41     self.cursor.execute('select * from version')
42     result = self.cursor.fetchall()
43     if result:
44         if result[0][0] == real_two_link:
45             print('已是最新, 无须爬取')
46     else:
47         self.get_data(real_two_link)
48         self.cursor.execute('delete from version')
49         self.cursor.execute('insert into version values(%)', [real_two_link])
50         self.db.commit()
51
52     # 用xpath直接提取数据
53     def get_data(self, real_two_link):
54         real_two_html = requests.get(real_two_link, headers=self.headers).text
55         parse_html = etree.HTML(real_two_html)
56         # 基准xpath, 提取每个信息的节点列表对象
57         tr_list = parse_html.xpath('//tr[@height=19]')
58         city_info = {}
59         for tr in tr_list:
60             city_info['code'] = tr.xpath('./td[2]/text()')[0]
61             city_info['name'] = tr.xpath('./td[3]/text()')[0]
62             print(city_info)
63
64
65
66 if __name__ == '__main__':
67     spider = Govement()
68     spider.get_flase_link()

```

动态加载数据抓取-Ajax

■ 特点

- 1、右键 -> 查看网页源码中没有具体数据
- 2、滚动鼠标滑轮或其他动作时加载

■ 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
- 3、# 控制台中 XHR : 异步加载的数据包
- 4、# XHR -> QueryStringParameters(查询参数)

豆瓣电影数据抓取案例

■ 目标

- 1、地址：豆瓣电影 - 排行榜 - 剧情
- 2、目标：电影名称、电影评分

■ F12抓包 (XHR)

- 1、Request URL(基准URL地址) : `https://movie.douban.com/j/chart/top_list?`
- 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
- 8、`limit: 20` # 每次加载的电影数量

■ 代码实现

```
1 import requests
2 import json
3
4 class DoubanSpider(object):
5     def __init__(self):
6         self.url = 'https://movie.douban.com/j/chart/top_list?'
7         self.headers = {'User-Agent' : 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/72.0.3626.119 Safari/537.36'}
8         self.i = 0
9
10    # 获取页面
11    def get_page(self,params):
12        res = requests.get(url=self.url,params=params,headers=self.headers,verify=True)
13        res.encoding = 'utf-8'
14        # 返回 python 数据类型
```

```

15         html = res.json()
16         self.parse_page(html)
17
18     # 解析并保存数据
19     def parse_page(self,html):
20         item = {}
21         # html为大列表 [{电影1信息},{},{}]
22         for one in html:
23             # 名称
24             item['name'] = one['title'].strip()
25             # 评分
26             item['score'] = float(one['score'].strip())
27             # 打印测试
28             print(item)
29             self.i += 1
30
31     # 主函数
32     def main(self):
33         for start in range(0,41,20):
34             params = {
35                 'type' : '24',
36                 'interval_id' : '100:90',
37                 'action' : '',
38                 'start' : str(start),
39                 'limit' : '20'
40             }
41             # 调用函数,传递params参数
42             self.get_page(params)
43             print('电影数量:',self.i)
44
45 if __name__ == '__main__':
46     spider = DoubanSpider()
47     spider.main()

```

练习: 能否抓取指定类型的所有电影信息? - 无须指定数量

```

1 import requests
2 import json
3
4 class DoubanSpider(object):
5     def __init__(self):
6         self.url = 'https://movie.douban.com/j/chart/top_list?'
7         self.headers = {'User-Agent' : 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/72.0.3626.119 Safari/537.36'}
8         self.i = 0
9
10    # 获取页面
11    def get_page(self,params):
12        res = requests.get(url=self.url,params=params,headers=self.headers,verify=True)
13        res.encoding = 'utf-8'
14        # 返回 python 数据类型
15        html = res.json()
16        self.parse_page(html)
17
18    # 解析并保存数据
19    def parse_page(self,html):

```



```

20         item = {}
21         # html为大列表 [{电影1信息},{},{}]
22         for one in html:
23             # 名称
24             item['name'] = one['title'].strip()
25             # 评分
26             item['score'] = float(one['score'].strip())
27             # 打印测试
28             print(item)
29             self.i += 1
30
31     # 获取电影总数
32     def total_number(self):
33         # F12抓包抓到的地址
34         url = 'https://movie.douban.com/j/chart/top_list_count?type=24&interval_id=100%3A90'
35         html = requests.get(url=url, headers=self.headers, verify=False).json()
36         total = int(html['total'])
37
38         return total
39
40     # 主函数
41     def main(self):
42         # 获取type的值
43         name = input('请输入电影类型(剧情|喜剧|爱情):')
44         typ_dic = {'剧情': '11', '喜剧': '24', '爱情': '13'}
45         typ = str(typ_dic[name])
46         # 获取电影总数
47         total = self.total_number()
48         for start in range(0, (total+1), 20):
49             params = {
50                 'type': typ,
51                 'interval_id': '100:90',
52                 'action': '',
53                 'start': str(start),
54                 'limit': '20'
55             }
56             # 调用函数,传递params参数
57             self.get_page(params)
58             print('电影数量:', self.i)
59
60 if __name__ == '__main__':
61     spider = DoubanSpider()
62     spider.main()

```

多线程爬虫

应用场景

- 1、多进程：CPU密集程序
- 2、多线程：爬虫(网络I/O)、本地磁盘I/O

知识点回顾

■ 队列

```
1 # 导入模块
2 from queue import Queue
3 # 使用
4 q = Queue()
5 q.put(url)
6 q.get() # 当队列为空时, 阻塞
7 q.empty() # 判断队列是否为空, True/False
```

■ 线程模块

```
1 # 导入模块
2 from threading import Thread
3
4 # 使用流程
5 t = Thread(target=函数名) # 创建线程对象
6 t.start() # 创建并启动线程
7 t.join() # 阻塞等待回收线程
8
9 # 如何创建多线程, 如下方法你觉得怎么样????
10 for i in range(5):
11     t = Thread(target=函数名)
12     t.start()
13     t.join()
```

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
2 2、目标 : 应用分类 - 聊天社交
3     应用名称
4     应用链接
```

■ 实现步骤

1. 确认是否为动态加载

```
1 1、页面局部刷新
2 2、右键查看网页源代码, 搜索关键字未搜到
3 # 此网站为动态加载网站, 需要抓取网络数据包分析
```

2. F12抓取网络数据包

```
1 1、抓取返回json数据的URL地址 (Headers中的Request URL)
2 http://app.mi.com/categotyAlllistApi?page={}&categoryId=2&pageSize=30
3
4 2、查看并分析查询参数 (headers中的Query String Parameters)
5 page: 1
6 categoryId: 2
7 pageSize: 30
8 # 只有page再变, 0 1 2 3 ... , 这样我们就可以通过控制page的直拼接多个返回json数据的URL地址
```

■ 代码实现

```
1 import requests
2 from threading import Thread
3 from queue import Queue
4 import json
5 import time
6
7 class XiaomiSpider(object):
8     def __init__(self):
9         self.headers = {'User-Agent': 'Mozilla/5.0'}
10        self.url = 'http://app.mi.com/categotyAlllistApi?page={}&categoryId=2&pageSize=30'
11        # 定义队列, 用来存放URL地址
12        self.url_queue = Queue()
13
14        # URL入队列
15        def url_in(self):
16            # 拼接多个URL地址, 然后put()到队列中
17            for i in range(67):
18                self.url.format((str(i)))
19                self.url_queue.put(self.url)
20
21        # 线程事件函数(请求, 解析提取数据)
22        def get_page(self):
23            # 先get()URL地址, 发请求
24            # json模块做解析
25            while True:
26                # 当队列不为空时, 获取url地址
27                if not self.url_queue.empty():
28                    url = self.url_queue.get()
29                    html = requests.get(url, headers=self.headers).text
30                    self.parse_page(html)
31                else:
32                    break
33
34        # 解析函数
35        def parse_page(self, html):
36            app_json = json.loads(html)
37            item = {}
38            for app in app_json['data']:
39                # 应用名称
40                item['name'] = app['displayName']
41                # 应用链接
42                item['link'] = 'http://app.mi.com/details?id={}'.format(app['packageName'])
43
44            print(item)
```

```

45     # 主函数
46     def main(self):
47         self.url_in()
48         # 存放所有线程的列表
49         t_list = []
50
51         for i in range(10):
52             t = Thread(target=self.get_page)
53             t.start()
54             t_list.append(t)
55
56         # 统一回收线程
57         for p in t_list:
58             p.join()
59
60 if __name__ == '__main__':
61     start = time.time()
62     spider = XiaomiSpider()
63     spider.main()
64     end = time.time()
65     print('执行时间:%.2f' % (end-start))

```

今日作业

- 1 1、有道翻译案例复写一遍
- 2 2、抓取腾讯招聘数据(两级页面 - 职位名称、岗位职责、工作要求)
- 3 3、把腾讯招聘案例改写为多线程
- 4 4、把链家二手房案例改写为多线程
- 5 5、民政部数据抓取案例完善
- 6 # 1、将抓取的数据存入数据库，最好分表按照层级关系去存
- 7 # 2、增量爬取时表中数据也要更新