

# DAY05

## 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://1.1.1.1:8888',
6         'https': 'https://1.1.1.1:8888'
7     }
8 4、auth -> ('tarenacode', 'code_2013')
9 5、verify -> True/False
10 6、timeout
```

### *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,短时间内限制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
```

# Day05笔记

## 代理参数-proxies

### ■ 定义

- 1、定义：代替你原来的IP地址去对接网络的IP地址。
- 2、作用：隐藏自身真实IP,避免被封。

### 普通代理

#### ■ 获取代理IP网站

- 1 西刺代理、快代理、全网代理、代理精灵、... ..

#### ■ 参数类型

```
1 1、语法结构
2     proxies = {
3         '协议': '协议://IP:端口号'
4     }
5 2、示例
6     proxies = {
7         'http': 'http://IP:端口号',
8         'https': 'https://IP:端口号'
9     }
```

#### ■ 示例

使用免费普通代理IP访问测试网站: <http://httpbin.org/get>

```
1 import requests
2
3 url = 'http://httpbin.org/get'
4 headers = {
5     'User-Agent': 'Mozilla/5.0'
6 }
7 # 定义代理,在代理IP网站中查找免费代理IP
8 proxies = {
9     'http': 'http://112.85.164.220:9999',
10    'https': 'https://112.85.164.220:9999'
11 }
12 html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
13 print(html)
```

**思考: 建立一个自己的代理IP池, 随时更新用来抓取网站数据**

- 1、从西刺代理IP网站上, 抓取免费代理IP
- 2、测试抓取的IP, 可用的保存在文件中

### 思考 - 代码实现

```

1 import requests
2 from lxml import etree
3 import time
4 import random
5 from fake_useragent import UserAgent
6
7 class GetProxyIP(object):
8     def __init__(self):
9         self.url = 'https://www.xicidaili.com/nn/{}'.format(
10
11         # 随机生成1个User-Agent
12         def get_headers(self):
13             ua = UserAgent()
14             headers = { 'User-Agent':ua.random }
15
16             return headers
17
18         # 获取可用代理IP文件
19         def get_ip_file(self,url):
20             html = requests.get(url=url,headers=self.get_headers(),timeout=5).text
21
22             parse_html = etree.HTML(html)
23             tr_list = parse_html.xpath('//tr')
24             for tr in tr_list[1:]:
25                 ip = tr.xpath('./td[2]/text()')[0]
26                 port = tr.xpath('./td[3]/text()')[0]
27                 # 测试ip:port是否可用
28                 self.test_ip(ip,port)
29
30         def test_ip(self,ip,port):
31             proxies = {
32                 'http': 'http://{ip}:{port}'.format(ip,port),
33                 'https': 'https://{ip}:{port}'.format(ip, port),
34             }
35             test_url = 'http://www.baidu.com/'
36             try:
37                 res = requests.get(url = test_url,proxies = proxies,timeout = 8)
38                 if res.status_code == 200:
39                     print(ip,port,'Success')
40                     with open('proxies.txt','a') as f:
41                         f.write(ip + ':' + port + '\n')
42             except Exception as e:
43                 print(ip,port,'Failed')
44
45         # 主函数
46         def main(self):
47             for i in range(1,1001):
48                 url = self.url.format(i)
49                 self.get_ip_file(url)
50                 time.sleep(random.randint(5,10))
51
52 if __name__ == '__main__':
53     spider = GetProxyIP()
54     spider.main()

```

写一个获取收费开放代理的接口

```

1  # 获取开放代理的接口
2  import requests
3
4  def test_ip(ip):
5      url = 'http://www.baidu.com/'
6      proxies = {
7          'http': 'http://{ip}'.format(ip),
8          'https': 'https://{ip}'.format(ip),
9      }
10
11     try:
12         res = requests.get(url=url, proxies=proxies, timeout=8 )
13         if res.status_code == 200:
14             return True
15     except Exception as e:
16         return False
17
18 # 提取代理IP
19 def get_ip_list():
20     api_url = 'http://dev.kdlapi.com/api/getproxy/?
21     orderid=946562662041898&num=100&protocol=1&method=2&an_an=1&an_ha=1&sep=2'
22     html = requests.get(api_url).content.decode('utf-8', 'ignore')
23     # ip_port_list: ['IP:PORT', 'IP:PORT', '']
24     ip_port_list = html.split('\r\n')
25
26     # 依次遍历代理IP, 并进行测试
27     with open('proxy_ip.txt', 'a') as f:
28         for ip in ip_port_list:
29             if test_ip(ip):
30                 f.write(ip + '\n')
31
32 if __name__ == '__main__':
33     get_ip_list()

```

## 私密代理

### ■ 语法格式

```

1  1、语法结构
2  proxies = {
3      '协议': '协议://用户名:密码@IP:端口号'
4  }
5
6  2、示例
7  proxies = {
8      'http': 'http://用户名:密码@IP:端口号',
9      'https': 'https://用户名:密码@IP:端口号'
10 }

```

## 示例代码

```

1 import requests
2 url = 'http://httpbin.org/get'
3 proxies = {
4     'http': 'http://309435365:szayclhp@1.195.160.232:17509',
5     'https': 'https://309435365:szayclhp@1.195.160.232:17509'
6 }
7 headers = {
8     'User-Agent' : 'Mozilla/5.0',
9 }
10
11 html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
12 print(html)

```

## 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

```

## requests.post() 参数

### ■ 适用场景

```

1 Post类型请求的网站

```

### ■ 参数-data

```

1 response = requests.post(url,data=data,headers=headers)
2 # data : post数据 (Form表单数据-字典格式)

```

### ■ 请求方式的特点

```
1 # 一般
2 GET请求 : 参数在URL地址中有显示
3 POST请求: Form表单提交数据
```

## 有道翻译破解案例(post)

### ■ 目标

```
1 破解有道翻译接口, 抓取翻译结果
2 # 结果展示
3 请输入要翻译的词语: elephant
4 翻译结果: 大象
5 *****
6 请输入要翻译的词语: 喵喵叫
7 翻译结果: mews
```

### ■ 实现步骤

```
1 1、浏览器F12开启网络抓包, Network-All, 页面翻译单词后找Form表单数据
2 2、在页面中多翻译几个单词, 观察Form表单数据变化 (有数据是加密字符串)
3 3、刷新有道翻译页面, 抓取并分析JS代码 (本地JS加密)
4 4、找到JS加密算法, 用Python按同样方式加密生成加密数据
5 5、将Form表单数据处理为字典, 通过requests.post()的data参数发送
```

### ■ 具体实现

#### ■ 1、开启F12抓包, 找到Form表单数据如下:

```
1 i: 喵喵叫
2 from: AUTO
3 to: AUTO
4 smartresult: dict
5 client: fanyideskweb
6 salt: 15614112641250
7 sign: 94008208919faa19bd531acde36aac5d
8 ts: 1561411264125
9 bv: f4d62a2579ebb44874d7ef93ba47e822
10 doctype: json
11 version: 2.1
12 keyfrom: fanyi.web
13 action: FY_BY_REALTIME
```

#### ■ 2、在页面中多翻译几个单词, 观察Form表单数据变化

```
1 salt: 15614112641250
2 sign: 94008208919faa19bd531acde36aac5d
3 ts: 1561411264125
4 bv: f4d62a2579ebb44874d7ef93ba47e822
5 # 但是bv的值不变
```

### ■ 3、一般为本地js文件加密，刷新页面，找到js文件并分析JS代码

```
1 # 方法1
2 Network - JS选项 - 搜索关键词salt
3 # 方法2
4 控制台右上角 - Search - 搜索salt - 查看文件 - 格式化输出
5 # 最终找到相关JS文件 : fanyi.min.js
```

### ■ 4、打开JS文件，分析加密算法，用Python实现

```
1 # ts : 经过分析为13位的时间戳，字符串类型
2 js代码实现: "" + (new Date).getTime()
3 python实现: str(int(time.time()*1000))
4
5 # salt
6 js代码实现: ts+parseInt(10 * Math.random(), 10);
7 python实现: ts+str(random.randint(0,9))
8
9 # sign (设置断点调试, 来查看 e 的值, 发现 e 为要翻译的单词)
10 js代码实现: n.md5("fanyideskweb" + e + salt + "%A-rKaT5fb[Gy?;N5@Tj")
11 python实现:
12 from hashlib import md5
13 string = "fanyideskweb" + word + salt + "%A-rKaT5fb[Gy?;N5@Tj"
14 s = md5()
15 s.update(string.encode())
16 sign = s.hexdigest()
```

### ■ 5、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            # 检查频率最高 - 3个
12            "Cookie": "OUTFOX_SEARCH_USER_ID=970246104@10.169.0.83;
13            OUTFOX_SEARCH_USER_ID_NCOO=570559528.1224236;
14            _ntes_nnid=96bc13a2f5ce64962adfd6a278467214,1551873108952; JSESSIONID=aaae9i7p1XP1KaJH_gkYw;
15            td_cookie=18446744072941336803; SESSION_FROM_COOKIE=unknown;
16            __rl__test__cookies=1565689460872",
17            "Referer": "http://fanyi.youdao.com/",
18            "User-Agent": "Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
19            Gecko) Chrome/76.0.3809.100 Safari/537.36",
20        }
21
22        # 获取salt,sign,ts
23        def get_salt_sign_ts(self, word):
24            # ts
25            ts = str(int(time.time()*1000))
26            # salt
```

```

22     salt = ts + str(random.randint(0,9))
23     # sign
24     string = "fanyideskweb" + word + salt + "%A-rKaT5fb[Gy?;N5@Tj"
25     s = md5()
26     s.update(string.encode())
27     sign = s.hexdigest()
28
29     return salt,sign,ts
30
31 # 主函数
32 def attack_yd(self,word):
33     # 1. 先拿到salt,sign,ts
34     salt,sign,ts = self.get_salt_sign_ts(word)
35     # 2. 定义form表单数据为字典: data={}
36     # 检查了salt sign
37     data = {
38         "i": word,
39         "from": "AUTO",
40         "to": "AUTO",
41         "smartresult": "dict",
42         "client": "fanyideskweb",
43         "salt": salt,
44         "sign": sign,
45         "ts": ts,
46         "bv": "7e3150ecbdf9de52dc355751b074cf60",
47         "doctype": "json",
48         "version": "2.1",
49         "keyfrom": "fanyi.web",
50         "action": "FY_BY_REALTIME",
51     }
52     # 3. 直接发请求:requests.post(url,data=data,headers=xxx)
53     res = requests.post(
54         url=self.url,
55         data=data,
56         headers=self.headers
57     )
58     # res.json() 将json格式的字符串转为python数据类型
59     html = res.json()
60     # html: {'translateResult': [[{'tgt': '你好', 'src': 'hello'}]], 'errorCode': 0, 'type':
    'en2zh-CHS', 'smartResult': {'entries': ['', 'n. 表示问候, 惊奇或唤起注意时的用语\r\n', 'int.
    喂; 哈罗\r\n', 'n. (Hello)人名; (法)埃洛\r\n'], 'type': 1}}
61     result = html['translateResult'][0][0]['tgt']
62
63     print(result)
64
65 # 主函数
66 def main(self):
67     # 输入翻译单词
68     word = input('请输入要翻译的单词:')
69     self.attack_yd(word)
70
71 if __name__ == '__main__':
72     spider = YdSpider()
73     spider.main()

```



## 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
```

## 民政部网站数据抓取

### 目标

```
1 1、URL: http://www.mca.gov.cn/ - 民政数据 - 行政区划代码
2  即: http://www.mca.gov.cn/article/sj/xzqh/2019/
3 2、目标: 抓取最新中华人民共和国县级以上行政区划代码
```

### 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级以上行政区划代码
```

#### 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、真实链接中提取所需数据

#### 4、代码实现

```
1 import requests
2 from lxml import etree
3 import re
4
5 class GovementSpider(object):
6     def __init__(self):
7         self.url = 'http://www.mca.gov.cn/article/sj/xzqh/2019/'
8         self.headers = {'User-Agent': 'Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.1; WOW64; Trident/4.0; SLCC2; .NET CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center PC 6.0; .NET4.0C; InfoPath.3)'}
9
10
11     # 获取假链接
12     def get_false_link(self):
```

```

13     html = requests.get(url = self.url,headers = self.headers).text
14     # 解析
15     parse_html = etree.HTML(html)
16     # a_list: [<element a at xx>,<element a at xxx>]
17     a_list = parse_html.xpath('//a[@class="artitlelist"]')
18     for a in a_list:
19         # get()方法:获取某个属性的值
20         title = a.get('title')
21         if title.endswith('代码'):
22             false_link = 'http://www.mca.gov.cn' + a.get('href')
23             self.get_true_link(false_link)
24             break
25
26
27     # 获取真链接
28     def get_true_link(self,false_link):
29         # 先获取假链接的响应,然后根据响应获取真链接
30         html = requests.get(url = false_link,headers = self.headers).text
31         # 利用正则提取真实链接
32         re_bds = r'window.location.href="(.*?)"'
33         pattern = re.compile(re_bds,re.S)
34         true_link = pattern.findall(html)[0]
35
36         self.parse_html(true_link)
37
38     # 数据提取
39     def parse_html(self,true_link):
40         html = requests.get(url = true_link,headers = self.headers).text
41
42         # xpath提取数据
43         parse_html = etree.HTML(html)
44         tr_list = parse_html.xpath('//tr[@height="19"]')
45         for tr in tr_list:
46             code = tr.xpath('./td[2]/text()')[0].strip()
47             name = tr.xpath('./td[3]/text()')[0].strip()
48
49             print(name,code)
50
51
52     # 主函数
53     def main(self):
54         self.get_false_link()
55
56 if __name__ == '__main__':
57     spider = GovementSpider()
58     spider.main()

```

## 扩展

1. 建立增量爬虫 - 网站有更新时抓取, 否则不抓
2. 所抓数据存到数据库, 按照层级关系分表存储 - 省、市、县表

## 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,或者页面局部刷新

### ■ 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 0 20 40 60
- 8、`limit: 20` # 每次加载的电影数量

### ■ 代码实现 - 全站抓取 - 完美接口 - 指定类型所有电影信息

```
1 import requests
2 import time
3 import random
4 import re
5 from useragents import ua_list
6
7 class DoubanSpider(object):
8     def __init__(self):
9         self.url = 'https://movie.douban.com/j/chart/top_list?'
10        self.i = 0
11
12        # 获取随机headers
13        def get_headers(self):
14            headers = {'User-Agent':random.choice(ua_list)}
15
```

```

16         return headers
17
18     # 获取页面
19     def get_page(self, params):
20         # 返回 python 数据类型
21         html = requests.get(url=self.url, params=params, headers=self.get_headers()).json()
22         self.parse_page(html)
23
24     # 解析并保存数据
25     def parse_page(self, html):
26         item = {}
27         # html为大列表 [{电影1信息},{},{}]
28         for one in html:
29             # 名称 + 评分
30             item['name'] = one['title'].strip()
31             item['score'] = float(one['score'].strip())
32             # 打印测试
33             print(item)
34             self.i += 1
35
36     # 获取电影总数
37     def total_number(self, type_number):
38         # F12抓包抓到的地址
39         url = 'https://movie.douban.com/j/chart/top_list_count?type=
40 {}&interval_id=100%3A90'.format(type_number)
41         headers = self.get_headers()
42         html = requests.get(url=url, headers=headers).json()
43         total = int(html['total'])
44
45         return total
46
47     # 获取所有电影的名字和对应type值
48     def get_all_type_films(self):
49         # 获取 类型和类型码
50         url = 'https://movie.douban.com/chart'
51         headers = self.get_headers()
52         html = requests.get(url=url, headers=headers).text
53         re_bds = r'<a href=.*?type_name=(.*?)&type=(.*?)&.*?>'
54         pattern = re.compile(re_bds, re.S)
55         r_list = pattern.findall(html)
56         # 存放所有类型和对应类型码大字典
57         type_dict = {}
58         menu = ''
59         for r in r_list:
60             type_dict[r[0].strip()] = r[1].strip()
61             # 获取input的菜单, 显示所有电影类型
62             menu += r[0].strip() + '|'
63
64         return type_dict, menu
65
66     # 主函数
67     def main(self):
68         # 获取type的值
69         type_dict, menu = self.get_all_type_films()
70         menu = menu + '\n请做出你的选择:'
71         name = input(menu)

```

```

72         type_number = type_dict[name]
73         # 获取电影总数
74         total = self.total_number(type_number)
75         for start in range(0, (total+1), 20):
76             params = {
77                 'type' : type_number,
78                 'interval_id' : '100:90',
79                 'action' : '',
80                 'start' : str(start),
81                 'limit' : '20'
82             }
83             # 调用函数,传递params参数
84             self.get_page(params)
85             # 随机休眠1-3秒
86             time.sleep(random.randint(1,3))
87             print('电影数量:', self.i)
88
89     if __name__ == '__main__':
90         spider = DoubanSpider()
91         spider.main()

```

## 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、抓取百度图片 - 所有,而不是30张
- 4 4、民政部数据抓取案例完善
- 5 # 1、将抓取的数据存入数据库,最好分表按照层级关系去存
- 6 # 2、增量爬取时表中数据也要更新