

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
```

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 4. Ajax动态加载 : 从url加载网页的源代码后, 会在浏览器执行JavaScript程序, 这些程序会加载更多内容
- 14 解决方案: F12或抓包工具抓包处理
- 15
- 16 5. 对查询参数加密
- 17 解决方案: 找到JS文件, 分析加密算法, 用Python实现加密执行JS文件中的代码, 返回加密数据
- 18
- 19 6. 对响应内容做处理
- 20 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

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`
- 5、`interval_id: 100:90`
- 6、`action: ''`
- 7、`start: 0`
- 8、`limit`: 用户输入的电影数量

■ json模块的使用

- 1、`json.loads(json格式的字符串)`: 把json格式的字符串转为python数据类型
- 2、# 示例
- 3、`html = json.loads(res.text)`
- 4、`print(type(html))`

■ 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
9     # 获取页面
10    def get_page(self,params):
11        res = requests.get(
12            url=self.url,
13            params=params,
14            headers=self.headers,
15            verify=True
16        )
17        res.encoding = 'utf-8'
18        # json.loads() json格式->Python格式
19        html = res.json()
20        self.parse_page(html)
21
22    # 解析并保存数据
23    def parse_page(self,html):
24        # html为大列表 [{电影1信息},{},{}]
25        for h in html:
26            # 名称
27            name = h['title'].strip()
28            # 评分
29            score = float(h['score'].strip())
30            # 打印测试
31            print([name,score])
32
33    # 主函数
34    def main(self):
35        limit = input('请输入电影数量:')
36        params = {
37            'type' : '24',
38            'interval_id' : '100:90',
39            'action' : '',
40            'start' : '0',
41            'limit' : limit
42        }
43        # 调用函数,传递params参数
44        self.get_page(params)
45
46    if __name__ == '__main__':
47        spider = DoubanSpider()
48        spider.main()

```

思考: 实现用户在终端输入电影类型和电影数量, 将对应电影信息抓取到数据库

```

1    # 主函数
2    def main(self):
3        # 定义字典,来存放类型和对应的数字
4        ty_dict = {'剧情': '11', '喜剧': '24', '爱情': '13'}
5        ty = input('请输入电影类型(剧情|喜剧|爱情):')
6        if ty in ['剧情', '喜剧', '爱情']:
7            limit = input('请输入电影数量:')
8            params = {
9                'type' : ty_dict[ty],
10               'interval_id' : '100:90',

```

```

11         'action' : '',
12         'start' : '0',
13         'limit' : limit
14     }
15     # 调用函数,传递params参数
16     self.get_page(params)
17     else:
18         print('输入有误')

```

练习: 腾讯招聘案例抓包看看?

■ URL地址及目标

1. 确定URL地址及目标

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

2. F12抓包

3. 一级页面json地址(index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

```

4. 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

```

■ 具体代码实现

```

1 import requests
2 import json
3 import time
4 import random
5
6 class TencentSpider(object):
7     def __init__(self):
8         self.headers = {'User-Agent': 'Mozilla/5.0'}
9         self.one_url = 'https://careers.tencent.com/tencentcareer/api/post/Query?
timestamp=1563912271089&countryId=&cityId=&bgIds=&productId=&categoryId=&parentCategoryId=&attr
rId=&keyword=&pageIndex={}&pageSize=10&language=zh-cn&area=cn'
10
11     def get_page(self, url):
12         res = requests.get(url, headers=self.headers)
13         res.encoding = 'utf-8'
14         # json.loads()把json格式的字符串转为python数据类型
15         html = json.loads(res.text)
16         return html
17
18     def parse_one_page(self, html):
19         job_info = {}

```

```

20     for job in html['Data']['Posts']:
21         job_info['job_name'] = job['RecruitPostName']
22         job_info['job_address'] = job['LocationName']
23         # 拿postid为了拼接二级页面地址
24         post_id = job['PostId']
25         # 职责和要求(二级页面)
26         # 得到二级页面链接
27         two_url = 'https://careers.tencent.com/tencentcareer/api/post/ByPostId?
timestamp=1563912374645&postId={}&language=zh-cn'.format(post_id)
28         # 发请求解析
29         job_info['job_duty'], job_info['job_requirement'] = self.parse_two_page(two_url)
30         print(job_info)
31
32     def parse_two_page(self, two_url):
33         html = self.get_page(two_url)
34         # 职责
35         job_duty = html['Data']['Responsibility']
36         # 要求
37         job_requirement = html['Data']['Requirement']
38
39         return job_duty, job_requirement
40
41     def main(self):
42         for index in range(1, 11):
43             url = self.one_url.format(index)
44             one_html = self.get_page(url)
45             self.parse_one_page(one_html)
46             time.sleep(random.uniform(0.5, 1.5))
47
48 if __name__ == '__main__':
49     spider = TencentSpider()
50     spider.main()

```

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

■ 方法一

```

1 1、先登录成功1次,获取到携带登陆信息的Cookie
2   F12打开控制台,在页面输入用户名、密码,登录成功,找到/home(一般在抓到地址的上面)
3 2、携带着cookie发请求
4   ** Cookie
5   ** Referer(源,代表你从哪里转过来的)
6   ** User-Agent

```

```

1 import requests
2 from lxml import etree
3

```

```

4 # url为需要登录才能正常访问的地址
5 url = 'http://www.renren.com/969255813/profile'
6 # headers中的cookie为登录成功后抓取到的cookie
7 headers = {
8     "Accept":
9     "text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3",
10    "Accept-Encoding": "gzip, deflate",
11    "Accept-Language": "zh-CN,zh;q=0.9",
12    "Connection": "keep-alive",
13    # 此处注意cookie, 要自己抓取
14    "Cookie": "",
15    "Host": "www.renren.com",
16    "Referer": "http://www.renren.com/SysHome.do",
17    "Upgrade-Insecure-Requests": "1",
18    "User-Agent": "Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/74.0.3729.169 Safari/537.36",
19 }
20 html = requests.get(url,headers=headers).text
21 # 解析
22 parse_html = etree.HTML(html)
23 result = parse_html.xpath('//*[@id="operate_area"]/div[1]/ul/li[1]/span/text()')[0].strip()
24 # result:就读于中央戏剧学院
25 print(result)

```

▪ 方法二

1. 知识点

1 利用requests模块中的session会话保持功能

2. session会话使用流程

```

1 1、实例化session对象
2 session = requests.session()
3 2、让session对象发送get或者post请求
4 res = session.get(url,headers=headers)

```

3. 具体步骤

```

1 1、寻找登录时POST的地址
2 查看网页源码,查看form,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':''}

```

4. 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 from lxml import etree
3
4 # 定义常用变量
5 post_url = 'http://www.renren.com/PLogin.do'
6 post_data = {
7     'email' : 'xxxxxx',
8     'password' : 'xxxxxx'
9 }
10 headers = {
11     'User-Agent' : 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like
12     Gecko) Chrome/74.0.3729.169 Safari/537.36',
13     'Referer' : 'http://www.renren.com/SysHome.do'
14 }
15 # 实例化session会话保持对象
16 session = requests.session()
17 # 先POST,把用户名和密码信息POST到一个地址
18 session.post(post_url,data=post_data,headers=headers)
19
20 # 再get个人主页
21 url = 'http://www.renren.com/970294164/profile'
22 html = session.get(url,headers=headers).text
23
24 parse_html = etree.HTML(html)
25 result = parse_html.xpath('//*[@id="operate_area"]/div[1]/ul/li[1]/span/text()')[0].strip()
26 print(result)
```

requests.post()

■ 适用场景a

- 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)

1. 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
```

2. 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
6 # 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: r+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 + "n%A-rKaT5fb[Gy?;N5@Tj")
11 python实现:
12 from hashlib import md5
13 s = md5()
14 s.update("fanyideskweb" + e + salt + "n%A-rKaT5fb[Gy?;N5@Tj".encode())
15 sign = s.hexdigest()

```

■ 5、代码实现

```

1 import requests
2 import time
3 from hashlib import md5
4 import random
5
6 # 获取相关加密算法的结果
7 def get_salt_sign_ts(word):
8     # salt
9     salt = str(int(time.time()*1000)) + str(random.randint(0,9))
10    # sign
11    string = "fanyideskweb" + word + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
12    s = md5()
13    s.update(string.encode())
14    sign = s.hexdigest()
15    # ts
16    ts = str(int(time.time()*1000))
17    return salt,sign,ts
18
19 # 攻克有道
20 def attack_yd(word):
21     salt,sign,ts = get_salt_sign_ts(word)
22     # url为抓包抓到的地址 F12 -> translate_o -> post
23     url = 'http://fanyi.youdao.com/translate_o?smartresult=dict&smartresult=rule'
24     headers = {
25         "Accept": "application/json, text/javascript, */*; q=0.01",
26         "Accept-Encoding": "gzip, deflate",
27         "Accept-Language": "zh-CN,zh;q=0.9",

```

```

28         "Connection": "keep-alive",
29         "Content-Length": "238",
30         "Content-Type": "application/x-www-form-urlencoded; charset=UTF-8",
31         "Cookie": "OUTFOX_SEARCH_USER_ID=-1449945727@10.169.0.82;
OUTFOX_SEARCH_USER_ID_NCOO=1492587933.976261; JSESSIONID=aaa5_Lj5jzfQZ_IPPuaSw;
__rl__test__cookies=1559193524685",
32         "Host": "fanyi.youdao.com",
33         "Origin": "http://fanyi.youdao.com",
34         "Referer": "http://fanyi.youdao.com/",
35         "User-Agent": "Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like
Gecko) Chrome/74.0.3729.169 Safari/537.36",
36         "X-Requested-With": "XMLHttpRequest",
37     }
38     # Form表单数据
39     data = {
40         'i': word,
41         'from': 'AUTO',
42         'to': 'AUTO',
43         'smartresult': 'dict',
44         'client': 'fanyideskweb',
45         'salt': salt,
46         'sign': sign,
47         'ts': ts,
48         'bv': 'cf156b581152bd0b259b90070b1120e6',
49         'doctype': 'json',
50         'version': '2.1',
51         'keyfrom': 'fanyi.web',
52         'action': 'FY_BY_REALTIME'
53     }
54
55     json_html = requests.post(url,data=data,headers=headers).json()
56     result = json_html['translateResult'][0][0]['tgt']
57     return result
58
59 if __name__ == '__main__':
60     word = input('请输入要翻译的单词: ')
61     result = attack_yd(word)
62     print(result)

```

代理IP使用

购买开放代理使用

■ 开放代理接口

```

1 # getip.py
2 # 获取开放代理的接口
3 import requests
4
5 # 提取代理IP

```

```

6 def get_ip_list():
7     api_url = 'http://dev.kdlapi.com/api/getproxy/?
    orderid=996140620552954&num=100&protocol=2&method=2&an_an=1&an_ha=1&sep=1'
8     res = requests.get(api_url)
9     ip_port_list = res.text.split('\r\n')
10
11     return ip_port_list
12
13 if __name__ == '__main__':
14     proxy_ip_list = get_ip_list()
15     print(proxy_ip_list)

```

测试开放代理接口: <http://httpbin.org/get>

- 1 1、从代理网站上获取购买的普通代理的api链接
- 2 2、从api链接中提取出IP
- 3 3、随机选择代理IP访问网站进行数据抓取

```

1 from getip import *
2 import time
3 import random
4
5 url = 'http://httpbin.org/get'
6 headers = {'User-Agent' : 'Mozilla/5.0'}
7 proxy_ip_list = get_ip_list()
8
9 while True:
10     # 判断是否还有可用代理
11     if not proxy_ip_list:
12         proxy_ip_list = get_ip_list()
13
14     proxy_ip = random.choice(proxy_ip_list)
15     proxies = {
16         'http' : 'http://{}'.format(proxy_ip),
17         'https' : 'https://{}'.format(proxy_ip)
18     }
19     print(proxies)
20
21     try:
22         html =
requests.get(url=url,proxies=proxies,headers=headers,timeout=5,verify=False).text
23         print(html)
24         break
25     except:
26         print('正在更换代理IP, 请稍后... ..')
27         # 及时把不可用的代理IP移除
28         proxy_ip_list.remove(proxy_ip)
29         continue

```

购买私密代理使用

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@122.114.67.136:16819',
5     'https': 'https://309435365:szayclhp@122.114.67.136:16819',
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)
```

今日作业

作业

- 1、仔细复习并总结有道翻译案例，抓包流程，代码实现
- 2、通过百度翻译，来再次熟练抓包流程，分析，断点调试等操作