

Day05回顾

动态加载网站数据抓取

- 1 1、F12打开控制台，页面动作抓取网络数据包
- 2 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 -> Query String(查询参数)

有道翻译流程梳理

- 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)
- 7 - Form Data, 发现对应的规律
- 8 5. 寻找JS文件
- 9 右上角 ... -> Search -> 搜索关键字 -> 单击 -> 跳转到Sources, 左下角格式化符号{}
- 10 6. 查看JS代码
- 11 搜索关键字，找到相关加密方法，分析并用python实现
- 12 7. 断点调试
8. 完善程序

cookie模拟登陆

- 1 1、适用网站类型：爬取网站页面时需要登录后才能访问，否则获取不到页面的实际响应数据
- 2 2、方法1 (利用cookie)
- 3 1、先登录成功1次, 获取到携带登陆信息的Cookie (处理headers)
- 4 2、利用处理的headers向URL地址发请求
- 5 3、方法2 (利用session会话保持)
- 6 1、登陆，找到POST地址: form -> action对应地址
- 7 2、定义字典，创建session实例发送请求
- 8 # 字典key : <input>标签中name的值(email,password)
- 9 # post_data = {'email':'','password':''}

Day06笔记

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)
```

百度翻译破解案例

目标

- 1 破解百度翻译接口，抓取翻译结果数据

实现步骤

■ 1、F12抓包,找到json的地址,观察查询参数

```
1 1、POST地址: https://fanyi.baidu.com/v2transapi
2 2、Form表单数据 (多次抓取在变的字段)
3   from: zh
4   to: en
5   sign: 54706.276099 #这个是如何生成的?
6   token: a927248ae7146c842bb4a94457ca35ee # 基本固定,但也想办法获取
```

■ 2、抓取相关JS文件

```
1 右上角 - 搜索 - sign: - 找到具体JS文件(index_c8a141d.js) - 格式化输出
```

3、在JS中寻找sign的生成代码

```
1 1、在格式化输出的JS代码中搜索: sign: 找到如下JS代码: sign: m(a),
2 2、通过设置断点, 找到m(a)函数的位置, 即生成sign的具体函数
3   # 1. a 为要翻译的单词
4   # 2. 鼠标移动到 m(a) 位置处, 点击可进入具体m(a)函数代码块
```

4、生成sign的m(a)函数具体代码如下(在一个大的define中)

```
1 function a(r) {
2     if (Array.isArray(r)) {
3         for (var o = 0, t = Array(r.length); o < r.length; o++)
4             t[o] = r[o];
5         return t
6     }
7     return Array.from(r)
8 }
9 function n(r, o) {
10     for (var t = 0; t < o.length - 2; t += 3) {
11         var a = o.charAt(t + 2);
12         a = a >= "a" ? a.charCodeAt(0) - 87 : Number(a),
13         a = "+" === o.charAt(t + 1) ? r >>> a : r << a,
14         r = "+" === o.charAt(t) ? r + a & 4294967295 : r ^ a
15     }
16     return r
17 }
18 function e(r) {
19     var o = r.match(/[\uD800-\uDBFF][\uDC00-\uDFFF]/g);
```

```

20     if (null === o) {
21         var t = r.length;
22         t > 30 && (r = "" + r.substr(0, 10) + r.substr(Math.floor(t / 2) - 5, 10) +
r.substr(-10, 10))
23     } else {
24         for (var e = r.split(/[\uD800-\uDBFF][\uDC00-\uDFFF]/), C = 0, h = e.length, f = []; h
> C; C++)
25             "" !== e[C] && f.push.apply(f, a(e[C].split(""))),
26             C !== h - 1 && f.push(o[C]);
27         var g = f.length;
28         g > 30 && (r = f.slice(0, 10).join("") + f.slice(Math.floor(g / 2) - 5, Math.floor(g /
2) + 5).join("") + f.slice(-10).join(""))
29     }
30     // var u = void 0
31     // , l = "" + String.fromCharCode(103) + String.fromCharCode(116) +
String.fromCharCode(107);
32     // u = null !== i ? i : (i = window[l] || "") || "";
33     // 断点调试,然后从网页源码中找到 window.gtk的值
34     var u = '320305.131321201'
35
36     for (var d = u.split("."), m = Number(d[0]) || 0, s = Number(d[1]) || 0, S = [], c = 0, v
= 0; v < r.length; v++) {
37         var A = r.charCodeAt(v);
38         128 > A ? S[c++] = A : (2048 > A ? S[c++] = A >> 6 | 192 : (55296 === (64512 & A) && v
+ 1 < r.length && 56320 === (64512 & r.charCodeAt(v + 1)) ? (A = 65536 + ((1023 & A) << 10) +
(1023 & r.charCodeAt(++v)),
39             S[c++] = A >> 18 | 240,
40             S[c++] = A >> 12 & 63 | 128) : S[c++] = A >> 12 | 224,
41             S[c++] = A >> 6 & 63 |
128),
42             S[c++] = 63 & A | 128)
43     }
44     for (var p = m, F = "" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(97) + ("" + String.fromCharCode(94) + String.fromCharCode(43) +
String.fromCharCode(54)), D = "" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(51) + ("" + String.fromCharCode(94) + String.fromCharCode(43) +
String.fromCharCode(98)) + ("" + String.fromCharCode(43) + String.fromCharCode(45) +
String.fromCharCode(102)), b = 0; b < S.length; b++)
45         p += S[b],
46         p = n(p, F);
47     return p = n(p, D),
48     p ^= s,
49     0 > p && (p = (2147483647 & p) + 2147483648),
50     p %= 1e6,
51     p.toString() + "." + (p ^ m)
52 }
53 var i = null;
54 //此行报错,直接注释掉即可
55 //t.exports = e

```

- 5、直接将代码写入本地js文件,利用pyexecjs模块执行js代码进行调试

```
1 import execjs
2
3 with open('node.js','r') as f:
4     js_data = f.read()
5     # 创建对象
6     exec_object = execjs.compile(js_data)
7     sign = exec_object.eval('e("hello")')
8     print(sign)
```

■ 获取token

```
1 # 在js中
2 token: window.common.token
3 # 在响应中想办法获取此值
4 token_url = 'https://fanyi.baidu.com/?aldtype=16047'
5 regex: "token: '(.*?)'"
```

■ 具体代码实现

```
1 import requests
2 import re
3 import execjs
4
5 class BaiduTranslateSpider(object):
6     def __init__(self):
7         self.token_url = 'https://fanyi.baidu.com/?aldtype=16047'
8         self.post_url = 'https://fanyi.baidu.com/v2transapi'
9         self.headers = {
10             'accept':
11             'text/html,application/xhtml+xml,application/xml;q=0.9,image/webp,image/apng,*/*;q=0.8,application/signed-exchange;v=b3',
12             # 'accept-encoding': 'gzip, deflate, br',
13             'accept-language': 'zh-CN,zh;q=0.9',
14             'cache-control': 'no-cache',
```

```

14         'cookie': 'BAIDUID=52920E829C1F64EE98183B703F4E37A9:FG=1;
    BIDUPSID=52920E829C1F64EE98183B703F4E37A9; PSTM=1562657403;
    to_lang_often=%5B%7B%22value%22%3A%22en%22%2C%22text%22%3A%22u82F1%u8BED%22%7D%2C%7B%22va
    ue%22%3A%22zh%22%2C%22text%22%3A%22u4E2D%u6587%22%7D%5D; REALTIME_TRANS_SWITCH=1;
    FANYI_WORD_SWITCH=1; HISTORY_SWITCH=1; SOUND_SPD_SWITCH=1; SOUND_PREFER_SWITCH=1; delPer=0;
    BDORZ=B490B5EBF6F3CD402E515D22BCDA1598; BCLID=6890774803653935935;
    BDSFRCVID=4XAsJeCCxG3DLCbwbJrKDGwjNA0UN_I3KhXZ3J;
    H_BDCLCKID_SF=tRk8oIDaJCvSe6r1MtQ_M4F_qxby26nUQ5neaJ5n0-
    nnhnL4W46bqJKFLtozKMoI3C7fotJJ5nololIRy6CKjjb-jaDqJ5n3bTnjstcS2RREHJrg-
    trSMDCSHGRGWl09WDTm_D_KfxnkOnc6qJj0-jjXqqo8K5Ljaa5n-
    pPKKRAaQD04bPbZL4DdMa7HLtA03mkjbnczfn020P5P51J_e-4syPRG2xRnWivrKfA-
    b4ncjRcTehom3xI8LNj4050Tt2LEoDPMJKIbMI_rMbbfhKC3hqJfaI62aKDs_RCMbhqcEIL4eJOIb6_w5gcq0T_Httj
    tXR0atn7ZSMB5j4Qo5pK95p38bxnDK2rQLb5zah5nhMJ53j7JDPMP0-4rJhxy523i5J6vQpnJ8hQ3DRoWXPiQbN7P-
    p5Z5mAqKl0MLl0kbC_6j5DWDtvLeU7J-n8XbI60XRj85-
    ohHjrFMtQ_q4tehHRMBUo9WDTm_DoTttt5fUj6qJj855jXqqo8KMtHJaFf-pPKKRAashnzWjrKqQ0Q5pj-
    WnQr3mkjbn5yfn020pjPX6joht4syPRG2xRnWivrKfA-
    b4ncjRcTehom3xI8LNj4050Tt2LEoC0XtIDhMDvPMCTSMt_HMxrKetJyaR0JhpjbWJ5TEPnjDUOdLPDW-
    46HBM3xbKQw5CJGBf7zhpvdWhC5y6ISKx-_J68Dtf5; ZD_ENTRY=baidu; PSINO=2;
    H_PS_PSSID=26525_1444_21095_29578_29521_28518_29098_29568_28830_29221_26350_29459;
    locale=zh; Hm_lvt_64ecd82404c51e03dc91cb9e8c025574=1563426293,1563996067;
    from_lang_often=%5B%7B%22value%22%3A%22zh%22%2C%22text%22%3A%22u4E2D%u6587%22%7D%2C%7B%22v
    alue%22%3A%22en%22%2C%22text%22%3A%22u82F1%u8BED%22%7D%5D;
    Hm_lpvT_64ecd82404c51e03dc91cb9e8c025574=1563999768;
    yjs_js_security_passport=2706b5b03983b8fa12fe756b8e4a08b98fb43022_1563999769_js',
15         'pragma': 'no-cache',
16         'upgrade-insecure-requests': '1',
17         'user-agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML,
    like Gecko) Chrome/75.0.3770.142 Safari/537.36',
18     }
19
20     # 获取token
21     def get_token(self):
22         token_url = 'https://fanyi.baidu.com/?aldtype=16047'
23         # 定义请求头
24         r = requests.get(self.token_url, headers=self.headers)
25         token = re.findall(r"token: '(.*?)'", r.text)
26         window_gtk = re.findall(r"window.*?gtk = '(.*?)';</script>", r.text)
27         if token:
28             return token[0], window_gtk[0]
29
30     # 获取sign
31     def get_sign(self, word):
32         with open('百度翻译.js', 'r') as f:
33             js_data = f.read()
34
35             exec_object = execjs.compile(js_data)
36             sign = exec_object.eval('e("{}")'.format(word))
37
38             return sign
39
40     # 主函数
41     def main(self, word, fro, to):
42         token, gtk = self.get_token()
43         sign = self.get_sign(word)
44         # 找到form表单数据如下, sign和token需要想办法获取
45         form_data = {
46             'from': fro,

```

```

47         'to': to,
48         'query': word,
49         'transtype': 'realtime',
50         'simple_means_flag': '3',
51         'sign': sign,
52         'token': token
53     }
54     r = requests.post(self.post_url, data=form_data, headers=self.headers)
55     print(r.json()['trans_result']['data'][0]['dst'])
56
57 if __name__ == '__main__':
58     spider = BaiduTranslateSpider()
59     menu = '1. 英译汉 2. 汉译英'
60     choice = input('1. 英译汉 2. 汉译英 : ')
61     word = input('请输入要翻译的单词:')
62     if choice == '1':
63         fro = 'en'
64         to = 'zh'
65     elif choice == '2':
66         fro = 'zh'
67         to = 'en'
68
69     spider.main(word, fro, to)

```

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

```

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 re.findall(r'.*?中华人民共和国县级以上行政区划代码.*?', title, re.S):
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('127.0.0.1', 'root', '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_two_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.xpath('./@title')[0]
22         # 正则匹配title中包含这个字符串的链接
23         if re.findall(r'.*?中华人民共和国县级以上行政区划代码.*?', title, re.S):
24             # 获取到第1个就停止即可, 第1个永远是最新的链接
25             two_link = 'http://www.mca.gov.cn' + a.xpath('./@href')[0]
26             break
27
28     # 从已提取的two_link中提取二级页面的真实链接
29     two_html = requests.get(two_link, headers=self.headers).text
30     # 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)
31     real_two_link = re.findall(r'window.location.href="(.*?)"', two_html, re.S)[0]
32     # 实现增量爬取
33     self.cursor.execute('select * from version')
34     result = self.cursor.fetchall()
35     if result:
36         if result[-1][0] == real_two_link:
37             print('已是最新')
38     else:
39         self.get_data(real_two_link)
40         self.cursor.execute('insert into version values(%)', [real_two_link])
41         self.db.commit()
42
43     # 用xpath直接提取数据
44     def get_data(self, real_two_link):
45         real_two_html = requests.get(real_two_link, headers=self.headers).text
46         parse_html = etree.HTML(real_two_html)
47         # 基准xpath,提取每个信息的节点列表对象
48         tr_list = parse_html.xpath('//tr[@height=19]')
49         city_info = {}
50         for tr in tr_list:
51             city_info['code'] = tr.xpath('./td[2]/text()')
52             city_info['name'] = tr.xpath('./td[3]/text()')
53             print(city_info)
54
55
56
57 if __name__ == '__main__':
58     spider = Govement()
59     spider.get_two_link()

```

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            for app in app_json['data']:
38                # 应用名称
39                name = app['displayName']
40                # 应用链接
41                link = 'http://app.mi.com/details?id={}'.format(app['packageName'])
42                d = { '名称' : name, '链接' : link }
43                print(d)
44
45        # 主函数
```

```

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

```

json解析模块

json.loads(json 格式字符串)

■ 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.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 import json
2
3 app_dict = {
4     '应用名称' : 'QQ',
5     '应用链接' : 'http://qq.com'
6 }
7 with open('小米.json','a') as f:
8     json.dump(app_dict,f,ensure_ascii=False)
```

今日作业

- 1 1、有道翻译案例复写一遍
- 2 2、百度翻译案例复写一遍
- 3 3、民政部数据抓取案例完善
- 4 # 1、将抓取的数据存入数据库，最好分表按照层级关系去存
- 5 # 2、增量爬取时表中数据也要更新
- 6 4、把链家二手房案例改写为多线程