

*Day02回顾**

爬取网站思路

- 1、先确定是否为动态加载网站
- 2、找URL规律
- 3、正则表达式
- 4、定义程序框架，补全并测试代码

存入csv文件

```
1 import csv
2 with open('xxx.csv','w') as f:
3     writer = csv.writer(f)
4     writer.writerow([])
5     writer.writerows([( ),( ),( )])
```

持久化存储之MySQL

```
1 db = pymysql.connect('IP',...)
2 cursor = db.cursor()
3 # cursor.execute('SQL',[ ])
4 # cursor.executemany('SQL',[( ),( ),( )])
5 db.commit()
6 cursor.close()
7 db.close()
```

requests模块

■ get()

- 1、发请求并获取响应对象
- 2、res = requests.get(url,headers=headers)

■ 响应对象res属性

```
1 res.text : 字符串
2 res.content : bytes
3 res.encoding: 字符编码 res.encoding='utf-8'
4 res.status_code : HTTP响应码
5 res.url : 实际数据URL地址
```

■ 非结构化数据保存

```
1 with open('xxx.jpg','wb') as f:
2     f.write(res.content)
```

多级页面数据抓取

```
1 1、先爬去一级页面,提取链接,继续跟进
2 2、爬取二级页面,提取数据
3 3、... ...
```

Chrome浏览器安装插件

■ 安装方法

```
1 1、把下载的相关插件（对应操作系统浏览器）后缀改为 .zip
2 2、打开Chrome浏览器 -> 右上角设置 -> 更多工具 -> 扩展程序 -> 点开开发者模式
3 3、把相关插件 拖拽 到浏览器中，释放鼠标即可安装
4 4、重启浏览器，使插件生效
```

■ 需要安装插件

```
1 1、Xpath Helper：轻松获取HTML元素的XPath路径
2 2、Proxy SwitchyOmega：Chrome浏览器中的代理管理扩展程序
3 3、JsonView：格式化输出json格式数据
```

Day03笔记

xpath解析

■ 定义

1 XPath即为XML路径语言，它是一种用来确定XML文档中某部分位置的语言，同样适用于HTML文档的检索

■ 示例HTML代码

```
1 <ul class="book_list">
2   <li>
3     <title class="book_001">Harry Potter</title>
4     <author>J K. Rowling</author>
5     <year>2005</year>
6     <price>69.99</price>
7   </li>
8
9   <li>
10    <title class="book_002">Spider</title>
11    <author>Forever</author>
12    <year>2019</year>
13    <price>49.99</price>
14  </li>
15 </ul>
```

■ 匹配演示

```
1 1、查找所有的li节点
2 //li
3 2、查找li节点下的title子节点中,class属性值为'book_001'的节点
4 //li/title[@class="book_001"]
5 3、查找li节点下所有title节点的,class属性的值
6 //li//title/@class
7
8 # 只要涉及到条件,加 []
9 # 只要获取属性值,加 @
```

■ 选取节点

```
1 1、// : 从所有节点中查找 (包括子节点和后代节点)
2 2、@ : 获取属性值
3 # 使用场景1 (属性值作为条件)
4 //div[@class="movie"]
5 # 使用场景2 (直接获取属性值)
6 //div/a/@src
```

■ 匹配多路径 (或)

1 xpath表达式1 | xpath表达式2 | xpath表达式3

■ 常用函数

```

1 1、contains()：匹配属性值中包含某些字符串节点
2   # 查找class属性值中包含"book_"的title节点
3   //title[contains(@class,"book_")]
4 2、text()：获取节点的文本内容
5   # 查找所有书籍的名称
6   //ul[@class="book_list"]/li/title/text()

```

lxml解析库

■ 安装

```
1 sudo pip3 install lxml
```

■ 使用流程

```

1 1、导模块
2   from lxml import etree
3 2、创建解析对象
4   parse_html = etree.HTML(html)
5 3、解析对象调用xpath
6   r_list = parse_html.xpath('xpath表达式')

```

■ 练习

```

1 from lxml import etree
2
3 html = '''<div class="wrapper">
4     <i class="iconfont icon-back" id="back"></i>
5     <a href="/" id="channel">新浪社会</a>
6     <ul id="nav">
7         <li><a href="http://domestic.firefox.sina.com/" title="国内">国内</a></li>
8         <li><a href="http://world.firefox.sina.com/" title="国际">国际</a></li>
9         <li><a href="http://mil.firefox.sina.com/" title="军事">军事</a></li>
10        <li><a href="http://photo.firefox.sina.com/" title="图片">图片</a></li>
11        <li><a href="http://society.firefox.sina.com/" title="社会">社会</a></li>
12        <li><a href="http://ent.firefox.sina.com/" title="娱乐">娱乐</a></li>
13        <li><a href="http://tech.firefox.sina.com/" title="科技">科技</a></li>
14        <li><a href="http://sports.firefox.sina.com/" title="体育">体育</a></li>
15        <li><a href="http://finance.firefox.sina.com/" title="财经">财经</a></li>
16        <li><a href="http://auto.firefox.sina.com/" title="汽车">汽车</a></li>
17    </ul>
18    <i class="iconfont icon-liebiao" id="menu"></i>
19 </div>'''
20 # 创建解析对象
21 parseHtml = etree.HTML(html)
22 # 调用xpath返回结果,text()为文本内容
23 rList = parseHtml.xpath('//a/text()')
24 #print(rList)
25
26 # 提取所有的href的属性值
27 r2 = parseHtml.xpath('//a/@href')

```

```

28 #print(r2)
29
30 # 提取所有href的值,不包括 /
31 r3 = parseHtml.xpath('//ul[@id="nav"]/li/a/@href')
32 #print(r3)
33
34 # 获取 图片、军事、...,不包括新浪社会
35 r4 = parseHtml.xpath('//ul[@id="nav"]/li/a/text()')
36 for r in r4:
37     print(r)

```

猫眼电影 (xpath)

■ 目标

- 1、地址: 猫眼电影 - 榜单 - top100榜
- 2、目标: 电影名称、主演、上映时间

■ 步骤

- 1、确定是否为静态页面 (右键-查看网页源代码, 搜索关键字确认)
- 2、写xpath表达式
- 3、写程序框架

■ xpath表达式

- 1、基准xpath: 匹配所有电影信息的节点对象列表
`//dl[@class="board-wrapper"]/dd`
- 2、遍历对象列表, 依次获取每个电影信息

```

for dd in dd_list:
    电影名称 : dd.xpath('./a/@title')[0].strip()
    电影主演 : dd.xpath('./p[@class="star"]/text()')[0].strip()
    上映时间 : dd.xpath('./p[@class="releasetime"]/text()')[0].strip()

```

■ 代码实现 (修改之前urllib库代码)

- 1、将urllib库改为requests模块实现
- 2、改写parse_page()方法

```

1 import requests
2 from lxml import etree
3 import time
4 import random
5
6 class MaoyanSpider(object):
7     def __init__(self):
8         self.url = 'https://maoyan.com/board/4?offset={}'
9         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'}

```

```

10     # 添加计数(页数)
11     self.page = 1
12
13     # 获取页面
14     def get_page(self,url):
15         # random.choice一定要写在这里,每次请求都会随机选择
16         res = requests.get(url,headers=self.headers)
17         res.encoding = 'utf-8'
18         html = res.text
19         self.parse_page(html)
20
21     # 解析页面
22     def parse_page(self,html):
23         # 创建解析对象
24         parse_html = etree.HTML(html)
25         # 基准xpath节点对象列表
26         dd_list = parse_html.xpath('//dl[@class="board-wrapper"]/dd')
27         movie_dict = {}
28         # 依次遍历每个节点对象,提取数据
29         for dd in dd_list:
30             movie_dict['name'] = dd.xpath('./p/a/@title')[0].strip()
31             movie_dict['star'] = dd.xpath('./p[@class="star"]/text()')[0].strip()
32             movie_dict['time'] = dd.xpath('./p[@class="releasetime"]/text()')[0].strip()
33
34             print(movie_dict)
35
36     # 主函数
37     def main(self):
38         for offset in range(0,31,10):
39             url = self.url.format(str(offset))
40             self.get_page(url)
41             print('第%d页完成' % self.page)
42             time.sleep(random.randint(1,3))
43             self.page += 1
44
45 if __name__ == '__main__':
46     spider = MaoyanSpider()
47     spider.main()

```

链家二手房案例 (xpath)

■ 实现步骤

1. 确定是否为静态

1 | 打开二手房页面 -> 查看网页源码 -> 搜索关键字

2. xpath表达式

```
1 1、修改方法：右键 -> copy xpath -> 测试修改
2 2、基准xpath表达式(匹配每个房源信息节点列表)
3 //ul[@class="sellListContent"]/li[@class="clear LOGCLICKDATA"] |
  //ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]
4 3、依次遍历后每个房源信息xpath表达式
5 * 名称: .//a[@data-el="region"]/text()
6 * 总价: .//div[@class="totalPrice"]/span/text()
7 * 单价: .//div[@class="unitPrice"]/span/text()
```

3. 代码实现

```
1 import requests
2 from lxml import etree
3 import time
4
5 class LianjiaSpider(object):
6     def __init__(self):
7         self.url = 'https://bj.lianjia.com/ershoufang/pg{}/'
8         self.headers = {'User-Agent': 'Mozilla/5.0'}
9
10    def get_page(self,url):
11        res = requests.get(url,headers=self.headers,timeout=10)
12        res.encoding = 'utf-8'
13        html = res.text
14        self.parse_page(html)
15
16    def parse_page(self,html):
17        parse_html = etree.HTML(html)
18        # 基准xpath
19        li_list = parse_html.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGCLICKDATA"]
20    | //ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]')
21        print(len(li_list))
22        house_dict = {}
23        # 遍历依次匹配每个房源信息
24        for li in li_list:
25            house_dict['house_name'] = li.xpath('.//a[@data-el="region"]/text()')[0].strip()
26            house_dict['total_price'] = li.xpath('.//div[@class="totalPrice"]/span/text()')
27            [0].strip()
28            house_dict['unit_price'] = li.xpath('.//div[@class="unitPrice"]/span/text()')[0].strip()
29
30            print(house_dict)
31
32    def main(self):
33        for pg in range(1,4):
34            url = self.url.format(str(pg))
35            self.get_page(url)
36            print('第%d页爬取成功' % pg)
37            time.sleep(0.5)
38
39 if __name__ == '__main__':
40     spider = LianjiaSpider()
41     spider.main()
```

百度贴吧图片抓取

■ 目标

```
1  抓取指定贴吧所有图片
```

■ 思路

```
1  1、获取贴吧主页URL,下一页,找到不同页的URL规律
2  2、获取1页中所有帖子URL地址: [帖子链接1,帖子链接2,...]
3  3、对每个帖子链接发请求,获取图片URL
4  4、向图片的URL发请求,以wb方式写入本地文件
```

■ 实现步骤

1. 贴吧URL规律

```
1  http://tieba.baidu.com/f?kw=??&pn=50
```

2. xpath表达式

```
1  1、帖子链接xpath
2      //*[@id="thread_list"]/li[@class=" j_thread_list clearfix"]/div/div[2]/div[1]/div[1]/a/@href
3
4  2、图片链接xpath
5      //div[@class="d_post_content j_d_post_content  clearfix"]/img[@class="BDE_Image"]/@src
6
7  3、视频链接xpath
8      //div[@class="video_src_wrapper"]/embed/@data-video
9      # 注意: 此处视频链接前端对响应内容做了处理,需要查看网页源代码来查看, 复制HTML代码在线格式化
```

3. 代码实现

```
1  import requests
2  from urllib import parse
3  from lxml import etree
4
5  class BaiduImgSpider(object):
6      def __init__(self):
7          self.url = 'http://tieba.baidu.com/f?{'
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     # 获取帖子链接
11     def get_tlink(self,url):
12         html = requests.get(url,headers=self.headers).text
13         # 提取帖子链接
14         parse_html = etree.HTML(html)
15         tlink_list = parse_html.xpath('//*[id="thread_list"]/li[@class=" j_thread_list clearfix"]/div/div[2]/div[1]/div[1]/a/@href')
16         # tlink_list: ['/p/23234','/p/9032323']
```



```

17     for tlink in tlink_list:
18         t_url = 'http://tieba.baidu.com' + tlink
19         # 提取图片链接并保存
20         self.get_imglink(t_url)
21
22     # 获取图片链接
23     def get_imglink(self, t_url):
24         res = requests.get(t_url, headers=self.headers)
25         res.encoding = 'utf-8'
26         html = res.text
27         # 提取图片链接
28         parse_html = etree.HTML(html)
29         # 匹配图片和视频的xpath表达式,中间加或 |
30         imglink_list = parse_html.xpath('//*[@class="d_post_content j_d_post_content clearfix"]/img/@src | //div[@class="video_src_wrapper"]/embed/@data-video')
31
32         for imglink in imglink_list:
33             self.write_img(imglink)
34
35     # 保存图片
36     def write_img(self, imglink):
37         res = requests.get(imglink, headers=self.headers)
38         # 截取链接的后10位作为文件名
39         filename = imglink[-10:]
40         with open(filename, 'wb') as f:
41             f.write(res.content)
42             print('%s下载成功' % filename)
43
44     # 指定贴吧名称,起始页和终止页,爬取图片
45     def main(self):
46         name = input('请输入贴吧名:')
47         begin = int(input('请输入起始页:'))
48         end = int(input('请输入终止页:'))
49         for page in range(begin, end+1):
50             # 查询参数编码
51             params = {
52                 'kw' : name,
53                 'pn' : str( (page-1)*50 )
54             }
55             params = parse.urlencode(params)
56             url = self.url.format(params)
57             # 开始获取图片
58             self.get_tlink(url)
59
60 if __name__ == '__main__':
61     spider = BaiduImgSpider()
62     spider.main()

```

requests.get()参数

查询参数-params

■ 参数类型

```
1 字典,字典中键值对作为查询参数
```

■ 使用方法

```
1 1、res = requests.get(url,params=params,headers=headers)
2 2、特点:
3 * url为基准的url地址, 不包含查询参数
4 * 该方法会自动对params字典编码,然后和url拼接
```

■ 示例

```
1 import requests
2
3 baseurl = 'http://tieba.baidu.com/f?'
4 params = {
5     'kw' : '赵丽颖吧',
6     'pn' : '50'
7 }
8 headers = {'User-Agent' : 'Mozilla/4.0 (compatible; MSIE 8.0; Windows NT 6.1; WOW64;
9     Trident/4.0; SLCC2; .NET CLR 2.0.50727; .NET CLR 3.5.30729; .NET CLR 3.0.30729; Media Center
10 PC 6.0; .NET4.0C; InfoPath.3)'}
11 # 自动对params进行编码,然后自动和url进行拼接,去发请求
12 res = requests.get(baseurl,params=params,headers=headers)
13 res.encoding = 'utf-8'
14 print(res.text)
```

Web 客户端验证 参数-auth

■ 作用及类型

```
1 1、针对于需要web客户端用户名密码认证的网站
2 2、auth = ('username','password')
```

■ 达内code课程方向案例

```
1 import requests
2 import re
3
4 class NoteSpider(object):
5     def __init__(self):
6         self.url = 'http://code.tarena.com.cn/'
7         self.headers = {'User-Agent': 'Mozilla/5.0'}
8         self.auth = ('tarenacode', 'code_2013')
9
10     # 获取+解析
11     def get_parse_page(self):
12         res = requests.get(
13             url=self.url,
```

```

14         auth=self.auth,
15         headers=self.headers
16     )
17     res.encoding = 'utf-8'
18     html = res.text
19     # 解析
20     p = re.compile('<a href=.*?>(.*?)</a>',re.S)
21     r_list = p.findall(html)
22     # r_list : ['..', 'AIDCode', 'ACCCode']
23     for r in r_list:
24         if r != '..':
25             print({ '课程方向' : r })
26
27 if __name__ == '__main__':
28     spider = NoteSpider()
29     spider.get_parse_page()

```

思考：爬取具体的笔记文件？

SSL证书认证参数-verify

■ 适用网站及场景

- 1、适用网站：https类型网站但是没有经过 证书认证机构 认证的网站
- 2、适用场景：抛出 `SSLError` 异常则考虑使用此参数

■ 参数类型

```

1 1、verify=True(默认)    : 检查证书认证
2 2、verify=False (常用) : 忽略证书认证
3 # 示例
4 response = requests.get(
5     url=url,
6     params=params,
7     headers=headers,
8     verify=False
9 )

```

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

1. 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://115.171.85.221:9000',
10    'https': 'https://115.171.85.221:9000'
11 }
12 html = requests.get(url,proxies=proxies,headers=headers,timeout=5).text
13 print(html)
```

2. 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 import requests
2 import random
3 from lxml import etree
4 from fake_useragent import UserAgent
5 import time
6
7 # 生成随机的User-Agent
8 def get_random_ua():
9     # 创建User-Agent对象
10    ua = UserAgent()
11    # 随机生成1个User-Agent
12    return ua.random
13
14 # 请求头
15 headers = {
16     'User-Agent': get_random_ua()
```

```

17 }
18 url = 'http://httpbin.org/get'
19
20 # 从西刺代理网站上获取随机的代理IP
21 def get_ip_list():
22     # 访问西刺代理网站国内高匿代理，找到所有的tr节点对象
23     res = requests.get('https://www.xicidaili.com/nn/', headers=headers)
24     parse_html = etree.HTML(res.text)
25     # 基准xpath，匹配每个代理IP的节点对象列表
26     ipobj_list = parse_html.xpath('//tr')
27     # 定义空列表，获取网页中所有代理IP地址及端口号
28     ip_list = []
29     # 从列表中第2个元素开始遍历，因为第1个为：字段名（国家、IP、... ...）
30     for ip in ipobj_list[1:]:
31         ip_info = ip.xpath('./td[2]/text()')[0]
32         port_info = ip.xpath('./td[3]/text()')[0]
33         ip_list.append(
34             {
35                 'http' : 'http://' + ip_info + ':' + port_info,
36                 'https' : 'https://' + ip_info + ':' + port_info
37             }
38         )
39     # 随机选择一个代理
40     proxies = random.choice(ip_list)
41     # 返回代理IP及代理池（列表ip_list）
42     return ip_list
43
44 # 主程序
45 def main_print():
46     # 我的IP代理池
47     ip_list = get_ip_list()
48     while True:
49         try:
50             # 设置超时时间，如果代理不能使用则切换下一个
51             proxies = random.choice(ip_list)
52             res = requests.get(url=url, headers=headers, proxies=proxies, timeout=5)
53             res.encoding = 'utf-8'
54             print(res.text)
55
56         except Exception as e:
57             # 此代理IP不能使用，从代理池中移除
58             ip_list.remove(proxies)
59             print('%s不能用，已经移除' % proxies)
60             # 继续循环获取最后1个代理IP
61             continue
62
63
64 if __name__ == '__main__':
65     main_print()

```

今日作业

糗事百科 (xpath)

- | | |
|---|--|
| 1 | 1、URL地址: https://www.qiushibaike.com/text/ |
| 2 | 2、目标 : 用户昵称、段子内容、好笑数量、评论数量 |

电影天堂 (xpath)