

*Day02回顾**

爬取网站思路

- 1 1、先确定是否为动态加载网站
- 2 2、找URL规律
- 3 3、正则表达式
- 4 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 1、发请求并获取响应对象
- 2 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) - 如下xpath表达式含义?
4 //div[@class="movie"]
5 # 使用场景2 (直接获取属性值) - 如下xpath表达式含义?
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 - 自己来写一写
6
```

lxml解析库

■ 安装

```
1 sudo pip3 install lxml
```

■ 使用流程

```
1 1、导模块
2
3 2、创建解析对象
4
5 3、解析对象调用xpath
6
```

■ 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
21 # 问题1：获取所有 a 节点的文本内容
22
23 # 问题2：获取所有 a 节点的 href 的属性值
24
25 # 问题3： 获取所有 a 节点的href的属性值，但是不包括 /
26
27 # 问题4： 获取 图片、军事、...,不包括新浪社会
```

猫眼电影 (xpath)

目标

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

步骤

- 1 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 2、写xpath表达式
- 3 3、写程序框架

xpath表达式

```

1 1、基准xpath：匹配所有电影信息的节点对象列表
2   //dl[@class="board-wrapper"]/dd
3
4 2、遍历对象列表，依次获取每个电影信息
5   for dd in dd_list:
6       电影名称 : dd.xpath('./a/@title')[0].strip()
7       电影主演 : dd.xpath('./p[@class="star"]/text()')[0].strip()
8       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 1、将urllib库改为requests模块实现
- 2 2、改写parse_page()方法

```

1

```

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、基准xpath表达式(匹配每个房源信息节点列表)
2 //ul[@class="sellListContent"]/li[@class="clear LOGCLICKDATA"] |
  //ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]
3 2、依次遍历后每个房源信息xpath表达式
4 * 名称: .//a[@data-el="region"]/text()
5 * 总价: .//div[@class="totalPrice"]/span/text()
6 * 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         pass
8
9     def get_page(self,url):
10        pass
11
12    def parse_page(self,html):
13        pass
14
15    def main(self):
16        pass
17
18 if __name__ == '__main__':
19     start = time.time()
20     spider = LianjiaSpider()
21     spider.main()
22     end = time.time()
23     print('执行时间:%.2f' % (end-start))
```

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
```

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; 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 # 自动对params进行编码,然后自动和url进行拼接,去发请求
10 res = requests.get(baseurl,params=params,headers=headers)
11 res.encoding = 'utf-8'
12 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 |
```

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 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 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 |

今日作业

糗事百科 (xpath)

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

电影天堂 (xpath)