

Day03回顾

数据抓取

■ 思路步骤

- 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 | xpath表达式
- 4 【4】定义程序框架，补全并测试代码

■ 增量爬虫实现思路

- 1 【1】原理
- 2 利用Redis集合特性，可将抓取过的指纹添加到redis集合中，根据返回值来判定是否需要抓取
- 3 【2】实现(根据sadd的返回值)
- 4 返回值为1：代表之前未抓取过，需要进行抓取
- 5 返回值为0：代表已经抓取过，无须再次抓取

■ 目前反爬处理

- 1 【1】基于User-Agent反爬
- 2 1.1) 发送请求携带请求头: headers={'User-Agent': 'Mozilla/5.0 xxxxxx'}
- 3 1.2) 多个请求时随机切换User-Agent
- 4 a) 定义py文件存放大量User-Agent，导入后使用random.choice()每次随机选择
- 5 b) 使用fake_useragent模块每次访问随机生成User-Agent
- 6 from fake_useragent import UserAgent
- 7 agent = UserAgent().random

数据持久化

■ csv

```
1 import csv
2 with open('xxx.csv', 'w', encoding='utf-8', newline='') as f:
3     writer = csv.writer(f)
4     writer.writerow([])
```

■ MySQL

```

1  import pymysql
2
3  # __init__(self):
4      self.db = pymysql.connect('IP',... ...)
5      self.cursor = self.db.cursor()
6
7  # save_html(self,r_list):
8      self.cursor.execute('sql',[data1])
9      self.db.commit()
10
11 # run(self):
12     self.cursor.close()
13     self.db.close()

```

▪ MongoDB

```

1  import pymongo
2
3  # __init__(self):
4      self.conn = pymongo.MongoClient('localhost', 27017)
5      self.db = self.conn['库名']
6      self.myset = self.db['集合名']
7
8  # save_html(self,r_list):
9      self.myset.insert_one({})

```

Day04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

▪ 匹配演示 - 猫眼电影top100

```

1  【1】 查找所有的dd节点
2      //dd
3  【2】 获取所有电影的名称的a节点：所有class属性值为name的a节点
4      //p[@class="name"]/a
5  【3】 获取dl节点下第2个dd节点的电影节点
6      //dl[@class="board-wrapper"]/dd[2]
7  【4】 获取所有电影详情页链接：获取每个电影的a节点的href的属性值
8      //p[@class="name"]/a/@href
9
10 【注意】
11     1> 只要涉及到条件,加 [] : //dl[@class="xxx"]    //dl/dd[2]
12     2> 只要获取属性值,加 @ : //dl[@class="xxx"]    //p/a/@href

```

■ 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      2.1> 使用场景1 (属性值作为条件)
4          //div[@class="movie-item-info"]
5      2.2> 使用场景2 (直接获取属性值)
6          //div[@class="movie-item-info"]/a/img/@src
7
8  【3】 练习 - 猫眼电影top100
9      3.1> 匹配电影名称
10         //div[@class="movie-item-info"]/p[1]/a/@title
11      3.2> 匹配电影主演
12         //div[@class="movie-item-info"]/p[2]/text()
13      3.3> 匹配上映时间
14         //div[@class="movie-item-info"]/p[3]/text()
15      3.4> 匹配电影链接
16         //div[@class="movie-item-info"]/p[1]/a/@href

```

■ 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】 text() : 获取节点的文本内容
2      xpath表达式末尾不加 /text() : 则得到的结果为节点对象
3      xpath表达式末尾加 /text() 或者 /@href : 则得到结果为字符串
4
5  【2】 contains() : 匹配属性值中包含某些字符串节点
6      匹配class属性值中包含 'movie-item' 这个字符串的 div 节点
7      //div[contains(@class,"movie-item")]

```

■ 终极总结

```
1 【1】 字符串: xpath表达式的末尾为: /text() 、 /@href 得到的列表中为 '字符串'
2
3 【2】 节点对象: 其他剩余所有情况得到的列表中均为 '节点对象'
4     [<element dd at xxxa>,<element dd at xxxb>,<element dd at xxxc>]
5     [<element div at xxxa>,<element div at xxxb>]
6     [<element p at xxxa>,<element p at xxxb>,<element p at xxxc>]
```

■ 课堂练习

```
1 【1】 匹配汽车之家-二手车,所有汽车的链接 :
2     //li[@class="cards-li list-photo-li"]/a[1]/@href
3     //a[@class="carinfo"]/@href
4 【2】 匹配汽车之家-汽车详情页中,汽车的
5     2.1)名称: //div[@class="car-box"]/h3/text()
6     2.2)里程: //ul/li[1]/h4/text()
7     2.3)时间: //ul/li[2]/h4/text()
8     2.4)挡位+排量: //ul/li[3]/h4/text()
9     2.5)所在地: //ul/li[4]/h4/text()
10    2.6)价格: //div[@class="brand-price-item"]/span[@class="price"]/text()
```

lxml解析库

■ 安装

```
1 【1】 Ubuntu: sudo pip3 install lxml
2 【2】 Windows: python -m pip 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')
```

■ xpath最常用

```
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     [<element dd at xxx>,<element dd at xxx>,...]
4
5 【2】 遍历对象列表, 依次获取每个电影信息
6     item = {}
7     for dd in dd_list:
8         item['name'] = dd.xpath('.//p[@class="name"]/a/text()').strip()
9         item['star'] = dd.xpath('.//p[@class="star"]/text()').strip()
10        item['time'] = dd.xpath('.//p[@class="releasetime"]/text()').strip()
```

豆瓣图书信息抓取 - xpath

需求分析

```
1  【1】 抓取目标 - 豆瓣图书top250的图书信息
2      https://book.douban.com/top250?start=0
3      https://book.douban.com/top250?start=25
4      https://book.douban.com/top250?start=50
5      ... ..
6
7  【2】 抓取数据
8      2.1) 书籍名称 : 红楼梦
9      2.2) 书籍描述 : [清] 曹雪芹 著 / 人民文学出版社 / 1996-12 / 59.70元
10     2.3) 书籍评分 : 9.6
11     2.4) 评价人数 : 286382人评价
12     2.5) 书籍类型 : 都云作者痴, 谁解其中味?
```

步骤分析

```
1  【1】 确认数据来源 - 响应内容存在
2  【2】 分析URL地址规律 - start为0 25 50 75 ...
3  【3】 xpath表达式
4      3.1) 基准xpath, 匹配每本书籍的节点对象列表
5           //div[@class="indent"]/table
6
7      3.2) 依次遍历每本书籍的节点对象, 提取具体书籍数据
8           书籍名称 : .//div[@class="pl2"]/a/@title
9           书籍描述 : .//p[@class="pl1"]/text()
10          书籍评分 : .//span[@class="rating_nums"]/text()
11          评价人数 : .//span[@class="pl1"]/text()
12          书籍类型 : .//span[@class="inq"]/text()
```

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 DoubanBookSpider:
8      def __init__(self):
9          self.url = 'https://book.douban.com/top250?start={}'
10
11     def get_html(self, url):
12         headers = { 'User-Agent':UserAgent().random }
13         html = requests.get(url=url, headers=headers).content.decode('utf-8','ignore')
14         # 直接调用解析函数
15         self.parse_html(html)
16
17     def parse_html(self, html):
18         p = etree.HTML(html)
19         # 基准xpath, 匹配每本书的节点对象列表
20         table_list = p.xpath('//div[@class="indent"]/table')
```

```

21         for table in table_list:
22             item = {}
23             # 书名
24             name_list = table.xpath('..//div[@class="p12"]/a/@title')
25             item['book_name'] = name_list[0].strip() if name_list else None
26             # 信息
27             info_list = table.xpath('..//p[@class="p1"]/text()')
28             item['book_info'] = info_list[0].strip() if info_list else None
29             # 评分
30             score_list = table.xpath('..//span[@class="rating_nums"]/text()')
31             item['book_score'] = score_list[0].strip() if score_list else None
32             # 人数
33             number_list = table.xpath('..//span[@class="p1"]/text()')
34             item['book_number'] = number_list[0].strip()[1:-1].strip() if number_list
else None
35             # 描述
36             comment_list = table.xpath('..//span[@class="inq"]/text()')
37             item['book_comment'] = comment_list[0].strip() if comment_list else None
38
39             print(item)
40
41         def run(self):
42             for i in range(10):
43                 start = i * 25
44                 page_url = self.url.format(start)
45                 self.get_html(url=page_url)
46                 # 控制数据抓取的频率,uniform生成指定范围内浮点数
47                 time.sleep(random.uniform(0, 3))
48
49
50     if __name__ == '__main__':
51         spider = DoubanBookSpider()
52         spider.run()

```

链家二手房案例 (xpath)

■ 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

■ 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      '此处滚动鼠标滑轮时,li节点的class属性值会发生变化,通过查看网页源码确定xpath表达式'
3      //ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]
4
5  【2】依次遍历后每个房源信息xpath表达式
6      2.1)名称: ../div[@class="positionInfo"]/a[1]/text()
7      2.2)地址: ../div[@class="positionInfo"]/a[2]/text()
8      2.3)户型+面积+方位+是否精装+楼层+年代+类型
9          info_list: '../div[@class="houseInfo"]/text()' -> [0].strip().split('|')
10         a)户型: info_list[0]

```

```

11     b)面积: info_list[1]
12     c)方位: info_list[2]
13     d)精装: info_list[3]
14     e)楼层: info_list[4]
15     f)年代: info_list[5]
16     g)类型: info_list[6]
17
18     2.4)总价+单价
19     a)总价: .//div[@class="totalPrice"]/span/text()
20     b)单价: .//div[@class="unitPrice"]/span/text()
21
22     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
23     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
24     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
25     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
26     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
27     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
28     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
29     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
30     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
31     ### 重要: 页面中xpath不能全信, 一切以响应内容为主
32     ### 重要: 页面中xpath不能全信, 一切以响应内容为主

```

■ 示意代码

```

1  import requests
2  from lxml import etree
3  from fake_useragent import UserAgent
4
5  # 1.定义变量
6  url = 'https://bj.lianjia.com/ershoufang/pg1/'
7  headers = {'User-Agent':UserAgent().random}
8  # 2.获取响应内容
9  html = requests.get(url=url,headers=headers).text
10 # 3.解析提取数据
11 parse_obj = etree.HTML(html)
12 # 3.1 基准xpath,得到每个房源信息的li节点对象列表, 如果此处匹配出来空, 则一定要查看响应内容
13 li_list = parse_obj.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]')
14 for li in li_list:
15     item = {}
16     # 名称
17     name_list = li.xpath('.//div[@class="positionInfo"]/a[1]/text()')
18     item['name'] = name_list[0].strip() if name_list else None
19     # 地址
20     add_list = li.xpath('.//div[@class="positionInfo"]/a[2]/text()')
21     item['add'] = add_list[0].strip() if add_list else None
22     # 户型 + 面积 + 方位 + 是否精装 + 楼层 + 年代 + 类型
23     house_info_list = li.xpath('.//div[@class="houseInfo"]/text()')
24     item['content'] = house_info_list[0].strip() if house_info_list else None
25     # 总价
26     total_list = li.xpath('.//div[@class="totalPrice"]/span/text()')
27     item['total'] = total_list[0].strip() if total_list else None
28     # 单价
29     unit_list = li.xpath('.//div[@class="unitPrice"]/span/text()')
30     item['unit'] = unit_list[0].strip() if unit_list else None

```

```
31
32     print(item)
```

■ 完整代码实现 - 自己实现

```
1  import requests
2  from lxml import etree
3  import time
4  import random
5  from fake_useragent import UserAgent
6
7  class LianjiaSpider(object):
8      def __init__(self):
9          self.url = 'https://bj.lianjia.com/ershoufang/pg{}/'
10
11      def parse_html(self, url):
12          headers = {'User-Agent': UserAgent().random}
13          html = requests.get(url=url, headers=headers).content.decode('utf-8', 'ignore')
14          self.get_data(html)
15
16
17      def get_data(self, html):
18          p = etree.HTML(html)
19          # 基准xpath: [<element li at xxx>,<element li>]
20          li_list = p.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]')
21          # for遍历, 依次提取每个房源信息, 放到字典item中
22          item = {}
23          for li in li_list:
24              # 名称+区域
25              name_list = li.xpath('.//div[@class="positionInfo"]/a[1]/text()')
26              item['name'] = name_list[0].strip() if name_list else None
27              address_list = li.xpath('.//div[@class="positionInfo"]/a[2]/text()')
28              item['address'] = address_list[0].strip() if address_list else None
29              # 户型+面积+方位+是否精装+楼层+年代+类型
30              # h_list: ['']
31              h_list = li.xpath('.//div[@class="houseInfo"]/text()')
32              if h_list:
33                  info_list = h_list[0].split('|')
34                  if len(info_list) == 7:
35                      item['model'] = info_list[0].strip()
36                      item['area'] = info_list[1].strip()
37                      item['direct'] = info_list[2].strip()
38                      item['perfect'] = info_list[3].strip()
39                      item['floor'] = info_list[4].strip()
40                      item['year'] = info_list[5].strip()[:-2]
41                      item['type'] = info_list[6].strip()
42                  else:
43                      item['model'] = item['area'] = item['direct'] = item['perfect'] =
item['floor'] = item['year'] = item['type'] = None
44                  else:
45                      item['model'] = item['area'] = item['direct'] = item['perfect'] =
item['floor'] = item['year'] = item['type'] = None
46
47              # 总价+单价
48              total_list = li.xpath('.//div[@class="totalPrice"]/span/text()')
```



```

49         item['total'] = total_list[0].strip() if total_list else None
50         unit_list = li.xpath('..//div[@class="unitPrice"]/span/text()')
51         item['unit'] = unit_list[0].strip() if unit_list else None
52
53         print(item)
54
55     def run(self):
56         for pg in range(1,101):
57             url = self.url.format(pg)
58             self.parse_html(url)
59             time.sleep(random.randint(1,2))
60
61 if __name__ == '__main__':
62     spider = LianjiaSpider()
63     spider.run()

```

■ 持久化到数据库中 - 自己实现

- 1 【1】 将数据存入MongoDB数据库
- 2 【2】 将数据存入MySQL数据库

代理参数-proxies

■ 定义及分类

- 1 【1】 定义 ： 代替你原来的IP地址去对接网络的IP地址
- 2
- 3 【2】 作用 ： 隐藏自身真实IP,避免被封

■ 普通代理

- 1 【1】 获取代理IP网站
- 2 快代理、全网代理、代理精灵、... ..
- 3
- 4 【2】 参数类型
- 5 proxies = { '协议': '协议://IP:端口号' }
- 6 proxies = {
- 7 'http': 'http://IP:端口号',
- 8 'https': 'https://IP:端口号',
- 9 }

■ 普通代理 - 示例

```

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

```

■ 代理IP池建立

```

1  """
2  建立开放代理的代理ip池
3  """
4  import requests
5
6  class ProxyPool:
7      def __init__(self):
8          self.api_url = 'http://dev.kdlapi.com/api/getproxy/?
orderid=999955248138592&num=20&protocol=2&method=2&an_ha=1&sep=1'
9          self.test_url = 'http://httpbin.org/get'
10         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/85.0.4183.83 Safari/537.36'}
11
12         def get_proxy(self):
13             html = requests.get(url=self.api_url, headers=self.headers).text
14             # proxy_list: ['1.1.1.1:8888', '2.2.2.2:9999,...]
15             proxy_list = html.split('\r\n')
16             for proxy in proxy_list:
17                 # 测试proxy是否可用
18                 self.test_proxy(proxy)
19
20         def test_proxy(self, proxy):
21             """测试1个代理ip是否可用"""
22             proxies = {
23                 'http' : 'http://{0}'.format(proxy),
24                 'https': 'https://{0}'.format(proxy),
25             }
26             try:
27                 resp = requests.get(url=self.test_url, proxies=proxies,
headers=self.headers, timeout=3)
28                 if resp.status_code == 200:
29                     print(proxy, '\033[31m可用\033[0m')
30                 else:
31                     print(proxy, '不可用')
32             except Exception as e:
33                 print(proxy, '不可用')
34
35         def run(self):
36             self.get_proxy()
37

```

```
38 if __name__ == '__main__':
39     spider = ProxyPool()
40     spider.run()
```

■ 私密代理+独享代理

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

■ 私密代理+独享代理 - 示例代码

```
1  import requests
2  url = 'http://httpbin.org/get'
3  proxies = {
4      'http': 'http://309435365:szayclhp@106.75.71.140:16816',
5      'https': 'https://309435365:szayclhp@106.75.71.140:16816',
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)
```

requests.post()

■ 适用场景

```
1  【1】适用场景 : Post类型请求的网站
2
3  【2】参数 : data={}
4      2.1) Form表单数据: 字典
5      2.2) res = requests.post(url=url,data=data,headers=headers)
6
7  【3】POST请求特点 : Form表单提交数据
```

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
```

```

3  【2】控制台常用选项
4    2.1) Network: 抓取网络数据包
5      a> ALL: 抓取所有的网络数据包
6      b> XHR: 抓取异步加载的网络数据包
7      c> JS : 抓取所有的JS文件
8    2.2) 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
9    2.3) Console: 交互模式, 可对JavaScript中的代码进行测试
10
11  【3】抓取具体网络数据包后
12    3.1) 单击左侧网络数据包地址, 进入数据包详情, 查看右侧
13    3.2) 右侧:
14      a> Headers: 整个请求信息
15            General、Response Headers、Request Headers、Query String、Form Data
16      b> Preview: 对响应内容进行预览
17      c> Response: 响应内容

```

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】准备抓包: F12开启控制台, 刷新页面
2  【2】寻找地址
3    2.1) 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
4          F12-Network-XHR-Headers-General-Request URL
5  【3】发现规律
6    3.1) 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
7    3.2) 分析对比 Network - All(或者XHR) - Form Data, 发现对应的规律
8  【4】寻找JS加密文件
9    控制台右上角 ...->Search->搜索关键字->单击->跳转到Sources, 左下角格式化符号{}
10 【5】查看JS代码
11    搜索关键字, 找到相关加密方法, 用python实现加密算法
12 【6】断点调试
13    JS代码中部分参数不清楚可通过断点调试来分析查看
14 【7】Python实现JS加密算法

```

今日作业

- 1 【1】将汽车之家案例使用 lxml + xpath 实现
- 2 【2】完善链家二手房案例，使用 lxml + xpath
- 3 【3】抓取快代理网站免费高匿代理，并测试是否可用来建立自己的代理IP池
- 4 <https://www.kuaidaili.com/free/>
- 5 【4】仔细熟悉有道翻译案例抓包及流程分析（至少操作5遍）