

概述

- 1 **【1】定义**
- 2 1.1) 网络蜘蛛、网络机器人, 抓取网络数据的程序
- 3 1.2) 其实就是用Python程序模仿人点击浏览器并访问网站, 而且模仿的越逼真越好
- 4
- 5 **【2】爬取数据的目的**
- 6 2.1) 公司项目的测试数据, 公司业务所需数据
- 7 2.2) 获取大量数据, 用来做数据分析
- 8
- 9 **【3】企业获取数据方式**
- 10 3.1) 公司自有数据
- 11 3.2) 第三方数据平台购买 (数据堂、贵阳大数据交易所)
- 12 3.3) 爬虫爬取数据
- 13
- 14 **【4】Python做爬虫优势**
- 15 4.1) Python : 请求模块、解析模块丰富成熟, 强大的Scrapy网络爬虫框架
- 16 4.2) PHP : 对多线程、异步支持不太好
- 17 4.3) JAVA: 代码笨重, 代码量大
- 18 4.4) C/C++: 虽然效率高, 但是代码成型慢
- 19
- 20 **【5】爬虫分类**
- 21 5.1) 通用网络爬虫 (搜索引擎使用, 遵守robots协议)
- 22 robots协议: 网站通过robots协议告诉搜索引擎哪些页面可以抓取, 哪些页面不能抓取, 通用网络爬虫
- 23 需要遵守robots协议 (君子协议)
- 24 示例: <https://www.baidu.com/robots.txt>
- 25 5.2) 聚焦网络爬虫 : 自己写的爬虫程序
- 26
- 27 **【6】爬取数据步骤**
- 28 6.1) 确定需要爬取的URL地址
- 29 6.2) 由请求模块向URL地址发出请求, 并得到网站的响应
- 30 6.3) 从响应内容中提取所需数据
- 31 a> 所需数据, 保存
- b> 页面中有其他需要继续跟进的URL地址, 继续第2步去发请求, 如此循环

- 重大问题思考

网站如何来判定是人类正常访问还是爬虫程序访问? --检查请求头!!!

```

1  # 请求头 (headers) 中的 User-Agent
2  # 测试案例：向测试网站http://httpbin.org/get发请求，查看请求头(User-Agent)
3  import requests
4
5  url = 'http://httpbin.org/get'
6  res = requests.get(url=url)
7  html = res.text
8  print(html)
9  # 请求头中:User-Agent为-> python-requests/2.22.0 那第一个被网站干掉的是谁??? 我们是不是需要
    发送请求时重构一下User-Agent??? 添加 headers 参数!!!

```

• 重大问题解决

```

1  """
2  包装好请求头后,向测试网站发请求,并验证
3  养成好习惯,发送请求携带请求头,重构User-Agent      User-Agent参数详解
4  """
5  import requests
6
7  url = 'http://httpbin.org/get'
8  headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/535.1
    (KHTML, like Gecko) Chrome/14.0.835.163 Safari/535.1'}
9  html = requests.get(url=url, headers=headers).content.decode('utf-8', 'ignore') #
    'ignore' 忽略无法转码的字符串 防止网页中带有无法识别字符串而报错
10 # Uni codeDecodeError: utf-8 xxx cannot decode char \xxx in.
11     ignore 可解决
12 # UnicodeEncodeError: gbk code cannot encode char \xxx in,
13     windows 写入文件时常报错误
14     # with open('xxx.txt', 'w', encoding='gb18030') as f:
15 print(html)

```

• 小总结

```

1  【1】 什么是robots协议,爬虫分为通用网络爬虫和聚焦网络爬虫,只有通用爬虫需要遵守协议
2  【2】 requests模块使用
3      res = requests.get(url=url, headers={'User-Agent': 'xxx'})
4      响应对象res属性:
5          a> res.text      # 字符串文本
6          b> res.content   # 二进制文本
7          c> res.status_code # 响应码
8          d> res.url       # 真实url
9  【3】网站乱码解析
10  方法1:
11      res = requests.get(url=url, headers=headers)
12      res.encoding = 'utf-8'
13      file.write(res.text)
14  方法2:  # 推荐使用方式
15      获取bytes数据,手动转码

```

```
16 requests.get(url=url, headers=headers).content.decode('utf-8')
```

- 抓取步骤

- 1 【1】确定所抓取数据在响应中是否存在（右键 - 查看网页源码 - 搜索关键字）
- 2 【2】数据存在：查看URL地址规律
- 3 【3】写正则表达式，来匹配数据
- 4 【4】程序结构
- 5 a>每爬取1个页面后随机休眠一段时间
- 6
- 7

```
1 # 程序结构
2 class xxxSpider(object):
3     def __init__(self):
4         # 定义常用变量,url,headers及计数等
5
6     def get_html(self):
7         # 获取响应内容函数,使用随机User-Agent
8
9     def parse_html(self):
10        # 使用正则表达式来解析页面,提取数据
11
12    def save_html(self):
13        # 将提取的数据按要求保存, csv、MySQL数据库等
14
15    def run(self):
16        # 程序入口函数,用来控制整体逻辑
17
18    if __name__ == '__main__':
19        # 程序开始运行时间戳
20        start = time.time()
21        spider = xxxSpider()
22        spider.run()
23        # 程序运行结束时间戳
24        end = time.time()
25        print('执行时间:%.2f' % (end-start))
```

方法

- 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

- 多级页面数据抓取思路

```
1  【1】整体思路
2      1.1> 爬取一级页面,提取 所需数据+链接,继续跟进
3      1.2> 爬取二级页面,提取 所需数据+链接,继续跟进
4      1.3> ... ...
5
6  【2】代码实现思路
7      2.1> 避免重复代码 - 请求、解析需定义函数
```

- 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      返回值为1 : 代表之前未抓取过,需要进行抓取
4      返回值为0 : 代表已经抓取过,无须再次抓取
5
6  【2】代码实现模板
7  import redis
8  from hashlib import md5
9  import sys
10
11 class XxxIncrSpider:
12     def __init__(self):
13         self.r = redis.Redis(host='localhost',port=6379,db=0)
14
15     def url_md5(self,url):
16         """对URL进行md5加密函数"""
17         s = md5()
18         s.update(url.encode())
19         return s.hexdigest()
20
21     def run_spider(self):
22         href_list = ['url1','url2','url3','url4']
23         for href in href_list:
24             href_md5 = self.url_md5(href)
25             if self.r.sadd('spider:urls',href_md5) == 1:
26                 返回值为1表示添加成功,即之前未抓取过,则开始抓取
27             else:
28                 sys.exit()
```

- 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
8
9  【2】响应内容存在特殊字符
10     解码时使用ignore参数
11     html = requests.get(url=url, headers=headers).content.decode(' ', 'ignore')

```

```

1  【1】结果: 节点对象列表
2      1.1) xpath示例: //div、//div[@class="student"]、//div/a[@title="stu"]/span
3
4  【2】结果: 字符串列表
5      2.1) xpath表达式中末尾为: @src、@href、/text()

```

案例1-猫眼电影top100

• 爬虫需求

```

1  【1】确定URL地址
2      百度搜索 - 猫眼电影 - 榜单 - top100榜
3
4  【2】爬取目标
5      所有电影的 电影名称、主演、上映时间

```

• 爬虫实现

```

1  【1】查看网页源码, 确认数据来源
2      响应内容中存在所需抓取数据 - 电影名称、主演、上映时间
3
4  【2】翻页寻找URL地址规律
5      第1页: https://maoyan.com/board/4?offset=0
6      第2页: https://maoyan.com/board/4?offset=10
7      第n页: offset=(n-1)*10
8
9  【3】编写正则表达式
10     <div class="movie-item-info">.*?title="(.*?)".*?class="star">(.*?)</p>.*?
11     releasetime">(.*?)</p>
12
12  【4】开干吧兄弟

```

- 代码实现

```
1  """
2  猫眼电影top100抓取 (电影名称、主演、上映时间)
3  """
4  import requests
5  import re
6  import time
7  import random
8
9  class MaoyanSpider:
10     def __init__(self):
11         self.url = 'https://maoyan.com/board/4?offset={}'
12         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64; Trident/7.0; rv:11.0) like Gecko'}
13
14     def get_html(self, url):
15         html = requests.get(url=url, headers=self.headers).text
16         # 直接调用解析函数
17         self.parse_html(html)
18
19     def parse_html(self, html):
20         """解析提取数据"""
21         regex = '<div class="movie-item-info">.*?title="(.*?)".*?<p class="star">(.*?)</p>.*?<p class="releasetime">(.*?)</p>'
22         pattern = re.compile(regex, re.S)
23         r_list = pattern.findall(html)
24         # r_list: [('活着', '牛犇', '2000-01-01'), (), (), ..., ())
25         self.save_html(r_list)
26
27     def save_html(self, r_list):
28         """数据处理函数"""
29         item = {}
30         for r in r_list:
31             item['name'] = r[0].strip()
32             item['star'] = r[1].strip()
33             item['time'] = r[2].strip()
34             print(item)
35
36     def run(self):
37         """程序入口函数"""
38         for offset in range(0, 91, 10):
39             url = self.url.format(offset)
40             self.get_html(url=url)
41             # 控制数据抓取频率:uniform() 生成指定范围内的浮点数
42             time.sleep(random.uniform(0, 1))
43
44 if __name__ == '__main__':
45     spider = MaoyanSpider()
46     spider.run()
```

案例2-汽车之家

- 领取任务

```
1  【1】爬取地址
2      汽车之家 - 二手车 - 价格从低到高
3      https://www.chel68.com/beijing/a0_0msdgsncngpi11to1csplexx0/
4
5
6  【2】爬取目标
7      所有汽车的 型号、行驶里程、上牌时间、档位、排量、车辆所在地、价格
8
9  【3】爬取分析
10     *****一级页面需抓取*****
11         1、车辆详情页的链接
12
13     *****二级页面需抓取*****
14         1、名称
15         2、行驶里程
16         3、上牌时间
17         4、档位
18         5、排量
19         6、车辆所在地
20         7、价格
```

- 实现步骤

```
1  【1】确定响应内容中是否存在所需抓取数据 - 存在
2
3  【2】找URL地址规律
4      第1页: https://www.chel68.com/beijing/a0_0msdgsncngpi11to1csplexx0/
5      第2页: https://www.chel68.com/beijing/a0_0msdgsncngpi11to1csp2exx0/
6      第n页: https://www.chel68.com/beijing/a0_0msdgsncngpi11to1csp{}exx0/
7
8  【3】 写正则表达式
9      一级页面正则表达式:<li class="cards-li list-photo-li".*?<a href="(.*?)".*?</li>
10     二级页面正则表达式:<div class="car-box">.*?<h3 class="car-brand-name">(.*?)
    </h3>.*?<ul class="brand-unit-item fn-clear">.*?<li>.*?<h4>(.*?)</h4>.*?<h4>(.*?)
    </h4>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<span class="price"
    id="overlayPrice">¥(.*?)<b>
11
12  【4】代码实现
```

- 代码实现

```
1  """
```

```

2 汽车之家二手车信息抓取
3 思路
4     1、一级页面：汽车的链接
5     2、二级页面：具体汽车信息
6
7 建立User-Agent池：防止被网站检测到是爬虫
8     使用fake_useragent模块
9     安装: sudo pip3 install fake_useragent
10    使用:
11        from fake_useragent import UserAgent
12        UserAgent().random
13    """
14    import requests
15    import re
16    import time
17    import random
18    from fake_useragent import UserAgent
19
20    class CarSpider:
21        def __init__(self):
22            self.url =
23            'https://www.chel68.com/beijing/a0_0msdgsncgpilltolcsp{}exx0/'
24
25        def get_html(self, url):
26            """功能函数1 - 获取html"""
27            headers = { 'User-Agent':UserAgent().random }
28            html = requests.get(url=url, headers=headers).text
29
30            return html
31
32        def re_func(self, regex, html):
33            """功能函数2 - 正则解析函数"""
34            pattern = re.compile(regex, re.S)
35            r_list = pattern.findall(html)
36
37            return r_list
38
39        def parse_html(self, one_url):
40            """爬虫逻辑函数"""
41            one_html = self.get_html(url=one_url)
42            one_regex = '<li class="cards-li list-photo-li".*?<a href="(.*?)".*?</li>'
43
44            href_list = self.re_func(regex=one_regex, html=one_html)
45            for href in href_list:
46                two_url = 'https://www.chel68.com' + href
47                # 获取1辆汽车的具体信息
48                self.get_car_info(two_url)
49                # 控制爬取频率
50                time.sleep(random.randint(1,2))
51
52        def get_car_info(self, two_url):
53            """获取1辆汽车的具体信息"""
54
55            two_html = self.get_html(url=two_url)

```



```

53         two_regex = '<div class="car-box">.*?<h3 class="car-brand-name">(.*?)</h3>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<span class="price" id="overlayPrice">¥(.*?)<b>'
54         # car_list: [('福睿斯', '3万公里', '2016年3月', '手动 / 1.5L', '廊坊', '5.60'),]
55         car_list = self.re_func(regex=two_regex, html=two_html)
56         item = {}
57         item['name'] = car_list[0][0].strip()
58         item['km'] = car_list[0][1].strip()
59         item['time'] = car_list[0][2].strip()
60         item['type'] = car_list[0][3].split('/')[0].strip()
61         item['displace'] = car_list[0][3].split('/')[1].strip()
62         item['address'] = car_list[0][4].strip()
63         item['price'] = car_list[0][5].strip()
64         print(item)
65
66     def run(self):
67         for i in range(1,5):
68             url = self.url.format(i)
69             self.parse_html(url)
70
71 if __name__ == '__main__':
72     spider = CarSpider()
73     spider.run()

```

- 练习 - 将数据存入MySQL数据库

```

1  create database cardb charset utf8;
2  use cardb;
3  create table cartab(
4  name varchar(100),
5  km varchar(50),
6  years varchar(50),
7  type varchar(50),
8  displacement varchar(50),
9  city varchar(50),
10 price varchar(50)
11 ) charset=utf8;

```

- 使用redis实现增量爬虫

```

1  """
2      提示：使用redis中的集合,sadd()方法,添加成功返回1,否则返回0
3      请各位大佬忽略掉下面代码,自己独立实现
4  """
5
6  import requests
7  import re
8  import time
9  import random

```

```

10     import pymysql
11     from hashlib import md5
12     import sys
13     import redis
14
15
16     class CarSpider(object):
17         def __init__(self):
18             self.url = 'https://www.chel68.com/beijing/a0_0msdgscncgpillto1csp{}exx0/'
19             self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/535.1 (KHTML, like Gecko) Chrome/14.0.835.163 Safari/535.1'}
20             self.db =
pymysql.connect('localhost', 'root', '123456', 'cardb', charset='utf8')
21             self.cursor = self.db.cursor()
22             # 连接redis去重
23             self.r = redis.Redis(host='localhost', port=6379, db=0)
24
25             # 功能函数1 - 获取响应内容
26             def get_html(self, url):
27                 html = requests.get(url=url, headers=self.headers).text
28
29                 return html
30
31             # 功能函数2 - 正则解析
32             def re_func(self, regex, html):
33                 pattern = re.compile(regex, re.S)
34                 r_list = pattern.findall(html)
35
36                 return r_list
37
38             # 爬虫函数开始
39             def parse_html(self, one_url):
40                 one_html = self.get_html(one_url)
41                 one_regex = '<li class="cards-li list-photo-li".*?<a href="(.*?)".*?</li>'
42                 href_list = self.re_func(one_regex, one_html)
43                 for href in href_list:
44                     # 加密指纹
45                     s = md5()
46                     s.update(href.encode())
47                     finger = s.hexdigest()
48                     # 如果指纹表中不存在
49                     if self.r.sadd('car:urls', finger):
50                         # 每便利一个汽车信息，必须要把此辆汽车所有数据提取完成后再提取下一辆汽车信息
51                         url = 'https://www.chel68.com' + href
52
53                         # 获取一辆汽车的信息
54                         self.get_data(url)
55                         time.sleep(random.randint(1, 2))
56                     else:
57                         sys.exit('抓取结束')
58
59             # 获取一辆汽车信息
60
61             def get_data(self, url):

```

```

61         two_html = self.get_html(url)
62         two_regex = '<div class="car-box">.*?<h3 class="car-brand-name">(.*?)</h3>.*?<ul class="brand-unit-item fn-clear">.*?<li>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<h4>(.*?)</h4>.*?<span class="price" id="overlayPrice">¥(.*?)<b'

63         item = {}
64         car_info_list = self.re_func(two_regex,two_html)
65         item['name'] = car_info_list[0][0]
66         item['km'] = car_info_list[0][1]
67         item['year'] = car_info_list[0][2]
68         item['type'] = car_info_list[0][3].split('/')[0]
69         item['displacement'] = car_info_list[0][3].split('/')[1]
70         item['city'] = car_info_list[0][4]
71         item['price'] = car_info_list[0][5]
72         print(item)
73
74         one_car_list = [
75             item['name'],
76             item['km'],
77             item['year'],
78             item['type'],
79             item['displacement'],
80             item['city'],
81             item['price']
82         ]
83         ins = 'insert into cartab values(%s,%s,%s,%s,%s,%s,%s)'
84         self.cursor.execute(ins,one_car_list)
85         self.db.commit()
86
87     def run(self):
88         for p in range(1,2):
89             url = self.url.format(p)
90             self.parse_html(url)
91
92         # 断开数据库链接
93         self.cursor.close()
94         self.db.close()
95
96     if __name__ == '__main__':
97         spider = CarSpider()
98         spider.run()

```

• 猫眼电影案例-xpath实现

```

1  """
2  猫眼电影top100抓取 (电影名称、主演、上映时间)
3  """
4  import requests
5  import time
6  import random
7  from lxml import etree

```

```

8
9 class MaoyanSpider:
10     def __init__(self):
11         self.url = 'https://maoyan.com/board/4?offset={}'
12         self.headers = {'User-Agent': 'Mozilla/5.0 (compatible; MSIE 9.0; Windows
NT 6.1; Win64; x64; Trident/5.0; .NET CLR 2.0.50727; SLCC2; .NET CLR 3.5.30729;
.NET CLR 3.0.30729; Media Center PC 6.0; InfoPath.3; .NET4.0C; Tablet PC 2.0;
.NET4.0E)'}
13
14     def get_html(self, url):
15         html = requests.get(url=url, headers=self.headers).text
16         # 直接调用解析函数
17         self.parse_html(html)
18
19     def parse_html(self, html):
20         """解析提取数据 - xpath"""
21         p = etree.HTML(html)
22         # 基准xpath: 每个电影信息的节点对象dd列表 [<element dd at xxx>, <element dd at
xxx>, ...]
23         dd_list = p.xpath('//dl[@class="board-wrapper"]/dd')
24         print(dd_list)
25         item = {}
26         for dd in dd_list:
27             item['name'] = dd.xpath('.//p[@class="name"]/a/@title')[0].strip()
28             item['star'] = dd.xpath('.//p[@class="star"]/text()')[0].strip()
29             item['time'] = dd.xpath('.//p[@class="releasetime"]/text()')
[0].strip()
30             print(item)
31
32     def run(self):
33         """程序入口函数"""
34         for offset in range(0, 91, 10):
35             url = self.url.format(offset)
36             self.get_html(url=url)
37             # 控制数据抓取频率: uniform() 生成指定范围内的浮点数
38             time.sleep(random.uniform(0, 1))
39
40 if __name__ == '__main__':
41     spider = MaoyanSpider()
42     spider.run()

```

案例3-豆瓣图书

- 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="pl"] /text()
10         书籍评分 : .//span[@class="rating_nums"] /text()
11         评价人数 : .//span[@class="pl"] /text()
12         书籍类型 : .//span[@class="inq"] /text()

```

• 代码实现

```

1  import requests
2  from lxml import etree
3  from fake_useragent import UserAgent
4  import time
5  import random
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         """使用随机的User-Agent"""
13         headers = {'User-Agent': UserAgent().random}
14         html = requests.get(url=url, headers=headers).text
15         self.parse_html(html)
16
17     def parse_html(self, html):
18         """lxml+xpath进行数据解析"""
19         parse_obj = etree.HTML(html)
20         # 1.基准xpath: 提取每本书的节点对象列表
21         table_list = parse_obj.xpath('//div[@class="indent"]/table')
22         for table in table_list:

```

```

23         item = {}
24         # 书名
25         name_list = table.xpath('..//div[@class="p12"]/a/@title')
26         item['name'] = name_list[0].strip() if name_list else None
27         # 描述
28         content_list = table.xpath('..//p[@class="p1"]/text()')
29         item['content'] = content_list[0].strip() if content_list else None
30         # 评分
31         score_list = table.xpath('..//span[@class="rating_nums"]/text()')
32         item['score'] = score_list[0].strip() if score_list else None
33         # 评价人数
34         nums_list = table.xpath('..//span[@class="p1"]/text()')
35         item['nums'] = nums_list[0][1:-1].strip() if nums_list else None
36         # 类别
37         type_list = table.xpath('..//span[@class="inq"]/text()')
38         item['type'] = type_list[0].strip() if type_list else None
39
40         print(item)
41
42     def run(self):
43         for i in range(5):
44             start = (i - 1) * 25
45             page_url = self.url.format(start)
46             self.get_html(page_url)
47             time.sleep(random.randint(1,2))
48
49 if __name__ == '__main__':
50     spider = DoubanBookSpider()
51     spider.run()

```

案例4-链家二手房

- 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          html =
13          requests.get(url=url, headers=headers, timeout=3).content.decode('utf-8', 'ignore')
14          self.get_data(html)

```

```

1  def get_data(self, html):
2      p = etree.HTML(html)
3      # 基准xpath: [<element li at xxx>, <element li>]
4      li_list = p.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA LOGCLICKDATA"]')
5      # for遍历, 依次提取每个房源信息, 放到字典item中
6      item = {}
7      for li in li_list:
8          # 名称+区域
9          name_list = li.xpath('.//div[@class="positionInfo"]/a[1]/text()')
10         item['name'] = name_list[0].strip() if name_list else None
11         address_list = li.xpath('.//div[@class="positionInfo"]/a[2]/text()')
12         item['address'] = address_list[0].strip() if address_list else None
13         # 户型+面积+方位+是否精装+楼层+年代+类型
14         # h_list: ['']
15         h_list = li.xpath('.//div[@class="houseInfo"]/text()')
16         if h_list:
17             info_list = h_list[0].split('|')
18             if len(info_list) == 7:
19                 item['model'] = info_list[0].strip()
20                 item['area'] = info_list[1].strip()
21                 item['direct'] = info_list[2].strip()
22                 item['perfect'] = info_list[3].strip()
23                 item['floor'] = info_list[4].strip()
24                 item['year'] = info_list[5].strip()[:-2]
25                 item['type'] = info_list[6].strip()
26             else:
27                 item['model'] = item['area'] = item['direct'] = item['perfect'] =
28                 item['floor'] = item['year'] = item['type'] = None
29             else:
30                 item['model'] = item['area'] = item['direct'] = item['perfect'] =

```



```

    item['floor'] = item['year'] = item['type'] = None
30
31     # 总价+单价
32     total_list = li.xpath('..//div[@class="totalPrice"]/span/text()')
33     item['total'] = total_list[0].strip() if total_list else None
34     unit_list = li.xpath('..//div[@class="unitPrice"]/span/text()')
35     item['unit'] = unit_list[0].strip() if unit_list else None
36
37     print(item)
38
39     def run(self):
40         for pg in range(1,101):
41             url = self.url.format(pg)
42             self.parse_html(url)
43             time.sleep(random.randint(1,2))
44
45 if __name__ == '__main__':
46     spider = LianjiaSpider()
47     spider.run()

```

多线程爬虫

- 应用场景

- 1 【1】多进程：CPU密集程序
- 2 【2】多线程：爬虫(网络I/O)、本地磁盘I/O

- 队列

```

1  【1】导入模块
2      from queue import Queue
3
4  【2】使用
5      q = Queue()
6      q.put(url)
7      q.get()    # 当队列为空时，阻塞
8      q.empty() # 判断队列是否为空，True/False
9
10 【3】q.get()解除阻塞方式
11     3.1) q.get(block=False)
12     3.2) q.get(block=True, timeout=3)
13     3.3) if not q.empty():
14         q.get()

```

- 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
10 t_list = []
11
12 for i in range(5):
13     t = Thread(target=函数名)
14     t_list.append(t)
15     t.start()
16
17 for t in t_list:
18     t.join()
```

- 线程锁

```
1  from threading import Lock
2
3  lock = Lock()
4  lock.acquire()
5  lock.release()
6
7  【注意】上锁成功后,再次上锁会阻塞
```

- 多线程爬虫示例代码

```
1  # 抓取豆瓣电影剧情类别下的电影信息
2  """
3  豆瓣电影 - 剧情 - 抓取
4  """
5  import requests
6  from fake_useragent import UserAgent
7  import time
8  import random
9  from threading import Thread, Lock
10 from queue import Queue
11
12 class DoubanSpider:
13     def __init__(self):
14         self.url = 'https://movie.douban.com/j/chart/top_list?
15         type=13&interval_id=100%3A90&action=&start={} &limit=20'
16
17         self.i = 0
```

```

16         # 队列 + 锁
17         self.q = Queue()
18         self.lock = Lock()
19
20     def get_agent(self):
21         """获取随机的User-Agent"""
22         return UserAgent().random
23
24     def url_in(self):
25         """把所有要抓取的URL地址入队列"""
26         for start in range(0, 684, 20):
27             url = self.url.format(start)
28             # url入队列
29             self.q.put(url)
30
31     # 线程事件函数: 请求+解析+数据处理
32     def get_html(self):
33         while True:
34             # 从队列中获取URL地址
35             # 一定要在判断队列是否为空 和 get() 地址 前后加锁,防止队列中只剩一个地址时出现重
36             # 复判断
37             self.lock.acquire()
38             if not self.q.empty():
39                 headers = {'User-Agent': self.get_agent()}
40                 url = self.q.get()
41                 self.lock.release()
42
43                 html = requests.get(url=url, headers=headers).json()
44                 self.parse_html(html)
45             else:
46                 # 如果队列为空,则最终必须释放锁
47                 self.lock.release()
48                 break
49
50     def parse_html(self, html):
51         """解析"""
52         # html: [{},{},{},{}]
53         item = {}
54         for one_film in html:
55             item['rank'] = one_film['rank']
56             item['title'] = one_film['title']
57             item['score'] = one_film['score']
58             print(item)
59             # 加锁 + 释放锁
60             self.lock.acquire()
61             self.i += 1
62             self.lock.release()
63
64     def run(self):
65         # 先让URL地址入队列
66         self.url_in()
67         # 创建多个线程,开干吧
68
69         t_list = []

```

```

68         for i in range(1):
69             t = Thread(target=self.get_html)
70             t_list.append(t)
71             t.start()
72
73         for t in t_list:
74             t.join()
75
76         print('数量:',self.i)
77
78 if __name__ == '__main__':
79     start_time = time.time()
80     spider = DoubanSpider()
81     spider.run()
82     end_time = time.time()
83     print('执行时间:%.2f' % (end_time-start_time))

```

-

```

1  【1】所用到的模块
2      1.1) from threading import Thread
3      1.2) from threading import Lock
4      1.3) from queue import Queue
5
6  【2】整体思路
7      2.1) 创建URL队列: q = Queue()
8      2.2) 产生URL地址,放入队列: q.put(url)
9      2.3) 线程事件函数: 从队列中获取地址,开始抓取: url = q.get()
10     2.4) 创建多线程,并运行
11
12 【3】代码结构
13     def __init__(self):
14         """创建URL队列"""
15         self.q = Queue()
16         self.lock = Lock()
17
18     def url_in(self):
19         """生成待爬取的URL地址,入队列"""
20         pass
21
22     def parse_html(self):
23         """线程事件函数,获取地址,进行数据抓取"""
24         while True:
25             self.lock.acquire()
26             if not self.q.empty():
27                 url = self.q.get()
28                 self.lock.release()
29             else:
30                 self.lock.release()
31                 break
32
33     def run(self):
34         self.url_in()

```

```

35         t_list = []
36         for i in range(3):
37             t = Thread(target=self.parse_html)
38             t_list.append(t)
39             t.start()
40
41         for th in t_list:
42             th.join()
43
44     【4】队列要点: q.get()防止阻塞方式
45     4.1) 方法1: q.get(block=False)
46     4.2) 方法2: q.get(block=True, timeout=3)
47     4.3) 方法3:
48         if not q.empty():
49             q.get()

```

scrapy框架

- 定义

1 异步处理框架,可配置和可扩展程度非常高,Python中使用最广泛的爬虫框架

- 安装

```

1  【1】Ubuntu安装
2      1.1) 安装依赖包
3          a> sudo apt-get install libffi-dev
4          b> sudo apt-get install libssl-dev
5          c> sudo apt-get install libxml2-dev
6          d> sudo apt-get install python3-dev
7          e> sudo apt-get install libxslt1-dev
8          f> sudo apt-get install zlib1g-dev
9          g> sudo pip3 install -I -U service_identity
10
11      1.2) 安装scrapy框架
12          a> sudo pip3 install Scrapy
13
14  【2】Windows安装
15      2.1) cmd命令行(管理员): python -m pip install Scrapy
16      【注意】: 如果安装过程中报如下错误
17          'Error: Microsoft Vistual C++ 14.0 is required xxx'
18          则安装Windows下的Microsoft Vistual C++ 14.0 即可(笔记spiderfiles中有)

```

- Scrapy框架五大组件

```

1  【1】引擎(Engine)      : 整个框架核心
2  【2】调度器(Scheduler) : 维护请求队列
3  【3】下载器(Downloader): 获取响应对象
4  【4】爬虫文件(Spider)  : 数据解析提取
5  【5】项目管道(Pipeline): 数据入库处理
6  *****
7  【中间件1】: 下载器中间件(Downloader Middlewares) : 引擎->下载器,包装请求(随机代理等)
8  【中间件2】: 蜘蛛中间件(Spider Middlewares) : 引擎->爬虫文件,可修改响应对象属性

```

- scrapy爬虫工作流程

```

1  【1】爬虫项目启动,由引擎向爬虫程序索要第一批要爬取的URL,交给调度器去入队列
2  【2】调度器处理请求后出队列,通过下载器中间件交给下载器去下载
3  【3】下载器得到响应对象后,通过蜘蛛中间件交给爬虫程序
4  【4】爬虫程序进行数据提取:
5      4.1) 数据交给管道文件去入库处理
6      4.2) 对于需要继续跟进的URL,再次交给调度器入队列,依次循环

```

- scrapy常用命令

```

1  【1】创建爬虫项目
2      scrapy startproject 项目名
3
4  【2】创建爬虫文件
5      scrapy genspider 爬虫名 域名
6
7  【3】运行爬虫
8      scrapy crawl 爬虫名

```

- scrapy项目目录结构

```

1  Baidu          # 项目文件夹
2  └─ Baidu       # 项目目录
3  │   └─ items.py # 定义数据结构
4  │   └─ middlewares.py # 中间件
5  │   └─ pipelines.py # 数据处理
6  │   └─ settings.py # 全局配置
7  │   └─ spiders
8  │       └─ baidu.py # 爬虫文件
9  └─ scrapy.cfg   # 项目基本配置文件

```

- settings.py常用变量

```

1  【1】USER_AGENT = 'Mozilla/5.0'
2

```

```
3  【2】ROBOTSTXT_OBEY = False
4      是否遵循robots协议,一般我们一定要设置为False
5
6  【3】CONCURRENT_REQUESTS = 32
7      最大并发量,默认为16
8
9  【4】DOWNLOAD_DELAY = 0.5
10     下载延迟时间: 访问相邻页面的间隔时间,降低数据抓取的频率
11
12  【5】COOKIES_ENABLED = False | True
13     Cookie默认是禁用的,取消注释则 启用Cookie, 即: True和False都是启用Cookie
14
15  【6】DEFAULT_REQUEST_HEADERS = {}
16     请求头,相当于requests.get(headers=headers)
```

小试牛刀

```
1  【1】执行3条命令,创建项目基本结构
2      scrapy startproject Baidu
3      cd Baidu
4      scrapy genspider baidu www.baidu.com
5
6  【2】完成爬虫文件: spiders/baidu.py
7      import scrapy
8      class BaiduSpider(scrapy.Spider):
9          name = 'baidu'
10         allowed_domains = ['www.baidu.com']
11         start_urls = ['http://www.baidu.com/']
12
13         def parse(self, response):
14             r_list = response.xpath('/html/head/title/text()').extract()[0]
15             print(r_list)
16
17  【3】完成settings.py配置
18      3.1) ROBOTSTXT_OBEY = False
19      3.2) DEFAULT_REQUEST_HEADERS = {
20          'User-Agent' : 'Mozilla/5.0'
21      }
22
23  【4】运行爬虫
24      4.1) 创建run.py (和scrapy.cfg同路径)
25      4.2) run.py
26          from scrapy import cmdline
27          cmdline.execute('scrapy crawl baidu'.split())
28
29  【5】执行 run.py 运行爬虫
```

瓜子二手车直卖网 - 一级页面

- 目标

```
1  【1】抓取瓜子二手车官网二手车收据（我要买车）
2
3  【2】URL地址: https://www.guazi.com/bj/buy/o{}/#bread
4      URL规律: o1  o2  o3  o4  o5  ...  ...
5
6  【3】所抓数据
7      3.1) 汽车链接
8      3.2) 汽车名称
9      3.3) 汽车价格
```

实现步骤

- 步骤1 - 创建项目和爬虫文件

```
1  scrapy startproject Car
2  cd Car
3  scrapy genspider car www.guazi.com
```

- 步骤2 - 定义要爬取的数据结构

```
1  """items.py"""
2  import scrapy
3
4  class CarItem(scrapy.Item):
5      # 链接、名称、价格
6      url = scrapy.Field()
7      name = scrapy.Field()
8      price = scrapy.Field()
```

- 步骤3 - 编写爬虫文件（代码实现1）

```
1  """
2  此方法其实还是一页一页抓取，效率并没有提升，和单线程一样
3
4  xpath表达式如下：
5  【1】基准xpath，匹配所有汽车节点对象列表
6      li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
7
8  【2】遍历后每辆车信息的xpath表达式
9      汽车链接: './a[1]/@href'
10     汽车名称: './h2[@class="t"]/text()'
11     汽车价格: './div[@class="t-price"]/p/text()'
12 """
13 # -*- coding: utf-8 -*-
14 import scrapy
15 from ..items import CarItem
```



```
class GuaziSpider(scrapy.Spider):
```

```
1     # 爬虫名
2     name = 'car'
3     # 允许爬取的域名
4     allowed_domains = ['www.guazi.com']
5     # 初始的URL地址
6     start_urls = ['https://www.guazi.com/bj/buy/o1/#bread']
7     # 生成URL地址的变量
8     n = 1
9
10    def parse(self, response):
11        # 基准xpath: 匹配所有汽车的节点对象列表
12        li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
13        # 给items.py中的 GuaziItem类 实例化
14        item = CarItem()
15        for li in li_list:
16            item['url'] = li.xpath('./a[1]/@href').get()
17            item['name'] = li.xpath('./a[1]/@title').get()
18            item['price'] = li.xpath('./div[@class="t-price"]/p/text()).get()
19
20            # 把抓取的数据,传递给了管道文件 pipelines.py
21            yield item
22
23        # 1页数据抓取完成,生成下一页的URL地址,交给调度器入队列
24        if self.n < 5:
25            self.n += 1
26            url = 'https://www.guazi.com/bj/buy/o{}/#bread'.format(self.n)
27            # 把url交给调度器入队列
28            yield scrapy.Request(url=url, callback=self.parse)
```

```
1
2 - **步骤3 - 编写爬虫文件 (代码实现2) **
3
4     ``python
5     """
6         重写start_requests()方法, 效率极高
7     """
8     # -*- coding: utf-8 -*-
9     import scrapy
10    from ..items import CarItem
11
12    class GuaziSpider(scrapy.Spider):
13        # 爬虫名
14        name = 'car2'
15        # 允许爬取的域名
16        allowed_domains = ['www.guazi.com']
17        # 1、去掉start_urls变量
18        # 2、重写 start_requests() 方法
19        def start_requests(self):
20            """生成所有要抓取的URL地址,一次性交给调度器入队列"""
21
22            for i in range(1,6):
```

```

22         url = 'https://www.guazi.com/bj/buy/o{}/#bread'.format(i)
23         # scrapy.Request(): 把请求交给调度器入队列
24         yield scrapy.Request(url=url, callback=self.parse)
25
26     def parse(self, response):
27         # 基准xpath: 匹配所有汽车的节点对象列表
28         li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
29         # 给items.py中的 GuaziItem类 实例化
30         item = CarItem()
31         for li in li_list:
32             item['url'] = li.xpath('./a[1]/@href').get()
33             item['name'] = li.xpath('./a[1]/@title').get()
34             item['price'] = li.xpath('./div[@class="t-price"]/p/text()').get()
35
36             # 把抓取的数据,传递给了管道文件 pipelines.py
37             yield item

```

• 步骤4 - 管道文件处理数据

```

1  """
2  pipelines.py处理数据
3  1、mysql数据库建库建表
4  create database cardb charset utf8;
5  use cardb;
6  create table cartab(
7  name varchar(200),
8  price varchar(100),
9  url varchar(500)
10 ) charset=utf8;
11 """
12 # -*- coding: utf-8 -*-
13
14 # 管道1 - 从终端打印输出
15 class CarPipeline(object):
16     def process_item(self, item, spider):
17         print(dict(item))
18         return item
19
20 # 管道2 - 存入MySQL数据库管道
21 import pymysql
22 from .settings import *
23
24 class CarMysqlPipeline(object):
25     def open_spider(self, spider):
26         """爬虫项目启动时只执行1次,一般用于数据库连接"""
27         self.db =
28         pymysql.connect(MYSQL_HOST, MYSQL_USER, MYSQL_PWD, MYSQL_DB, charset=CHARSET)
29         self.cursor = self.db.cursor()
30
31     def process_item(self, item, spider):
32         """处理从爬虫文件传过来的item数据"""
33         ins = 'insert into guazitab values(%s,%s,%s)'
34         car_li = [item['name'], item['price'], item['url']]

```

```

34         self.cursor.execute(ins, car_li)
35         self.db.commit()
36
37         return item
38
39     def close_spider(self, spider):
40         """爬虫程序结束时只执行1次,一般用于数据库断开"""
41         self.cursor.close()
42         self.db.close()

```

瓜子二手车直卖网 - 二级页面

- 目标说明

```

1  【1】在抓取一级页面的代码基础上升级
2  【2】一级页面所抓取数据（和之前一样）：
3      2.1) 汽车链接
4      2.2) 汽车名称
5      2.3) 汽车价格
6  【3】二级页面所抓取数据
7      3.1) 行驶里程：//ul[@class="assort clearfix"]/li[2]/span/text()
8      3.2) 排量：      //ul[@class="assort clearfix"]/li[3]/span/text()
9      3.3) 变速箱：    //ul[@class="assort clearfix"]/li[4]/span/text()

```

在原有项目基础上实现

- 步骤1 - items.py

```

1  # 添加二级页面所需抓取的数据结构
2
3  import scrapy
4
5  class GuaziItem(scrapy.Item):
6      # define the fields for your item here like:
7      # 一级页面：链接、名称、价格
8      url = scrapy.Field()
9      name = scrapy.Field()
10     price = scrapy.Field()
11     # 二级页面：时间、里程、排量、变速箱
12     time = scrapy.Field()
13     km = scrapy.Field()
14     disp = scrapy.Field()
15     trans = scrapy.Field()

```

- 步骤2 - car2.py

```

1  """
2      重写start_requests()方法,效率极高
3  """

```

```

4  # -*- coding: utf-8 -*-
5  import scrapy
6  from ..items import CarItem
7
8  class GuaziSpider(scrapy.Spider):
9      # 爬虫名
10     name = 'car2'
11     # 允许爬取的域名
12     allowed_domains = ['www.guazi.com']
13     # 1、去掉start_urls变量
14     # 2、重写 start_requests() 方法
15     def start_requests(self):
16         """生成所有要抓取的URL地址,一次性交给调度器入队列"""
17         for i in range(1,6):
18             url = 'https://www.guazi.com/bj/buy/o{}/#bread'.format(i)
19             # scrapy.Request(): 把请求交给调度器入队列
20             yield scrapy.Request(url=url,callback=self.parse)
21
22     def parse(self, response):
23         # 基准xpath: 匹配所有汽车的节点对象列表
24         li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
25         # 给items.py中的 GuaziItem类 实例化
26         item = CarItem()
27         for li in li_list:
28             item['url'] = 'https://www.guazi.com' +
li.xpath('./a[1]/@href').get()
29             item['name'] = li.xpath('./a[1]/@title').get()
30             item['price'] = li.xpath('./div[@class="t-price"]/p/text()').get()
31             # Request() 中meta参数: 在不同解析函数之间传递数据,item数据会随着response一起
返回
32             yield scrapy.Request(url=item['url'], meta={'meta_1': item},
callback=self.detail_parse)
33
34     def detail_parse(self, response):
35         """汽车详情页的解析函数"""
36         # 获取上个解析函数传递过来的 meta 数据
37         item = response.meta['meta_1']
38         item['km'] = response.xpath('//ul[@class="assort
clearfix"]/li[2]/span/text()').get()
39         item['disp'] = response.xpath('//ul[@class="assort
clearfix"]/li[3]/span/text()').get()
40         item['trans'] = response.xpath('//ul[@class="assort
clearfix"]/li[4]/span/text()').get()
41
42         # 1条数据最终提取全部完成, 交给管道文件处理
43         yield item

```

• 步骤3 - pipelines.py

```

1  # 将数据存入mongodb数据库,此处我们就不对MySQL表字段进行操作了,如有兴趣可自行完善
2  # MongoDB管道
3  import pymongo
4

```

```

5 class GuaziMongoPipeline(object):
6     def open_spider(self, spider):
7         """爬虫项目启动时只执行1次,用于连接MongoDB数据库"""
8         self.conn = pymongo.MongoClient(MONGO_HOST, MONGO_PORT)
9         self.db = self.conn[MONGO_DB]
10        self.myset = self.db[MONGO_SET]
11
12        def process_item(self, item, spider):
13            car_dict = dict(item)
14            self.myset.insert_one(car_dict)
15            return item

```

• 步骤4 - settings.py

```

1 # 定义MongoDB相关变量
2 MONGO_HOST = 'localhost'
3 MONGO_PORT = 27017
4 MONGO_DB = 'guazidb'
5 MONGO_SET = 'guaziset'

```

盗墓笔记小说抓取 - 三级页面

• 目标

```

1 【1】URL地址 : http://www.daomubiji.com/
2 【2】要求 : 抓取目标网站中盗墓笔记所有章节的所有小说的具体内容, 保存到本地文件
3             ./data/novel/盗墓笔记1:七星鲁王宫/七星鲁王_第一章_血尸.txt
4             ./data/novel/盗墓笔记1:七星鲁王宫/七星鲁王_第二章_五十年后.txt

```

• 准备工作xpath

```

1 【1】一级页面 - 大章节标题、链接:
2     1.1) 基准xpath匹配a节点对象列表: '//li[contains(@id,"menu-item-20")]/a'
3     1.2) 大章节标题: './text()'
4     1.3) 大章节链接: './@href'
5
6 【2】二级页面 - 小章节标题、链接
7     2.1) 基准xpath匹配article节点对象列表: '//article'
8     2.2) 小章节标题: './a/text()'
9     2.3) 小章节链接: './a/@href'
10
11 【3】三级页面 - 小说内容
12     3.1) p节点列表: '//article[@class="article-content"]/p/text()'
13     3.2) 利用join()进行拼接: ' '.join(['p1', 'p2', 'p3', ''])

```

项目实施

• 1、创建项目及爬虫文件

```
1 scrapy startproject Daomu
2 cd Daomu
3 scrapy genspider daomu www.daomubiji.com
```

• 2、定义要爬取的数据结构 - itemspy

```
1 class DaomuItem(scrapy.Item):
2     # 拷问：你的pipelines.py中需要处理哪些数据？ 文件名、路径
3     # 文件名：小标题名称 son_title: 七星鲁王 第一章 血尸
4     son_title = scrapy.Field()
5     directory = scrapy.Field()
6     content = scrapy.Field()
```

• 3、爬虫文件实现数据抓取 - daomu.py

```
1 # -*- coding: utf-8 -*-
2 import scrapy
3 from ..items import DaomuItem
4 import os
5
6 class DaomuSpider(scrapy.Spider):
7     name = 'daomu'
8     allowed_domains = ['www.daomubiji.com']
9     start_urls = ['http://www.daomubiji.com/']
10
11     def parse(self, response):
12         """一级页面解析函数：提取大标题+大链接,并把大链接交给调度器入队列"""
13         a_list = response.xpath('//li[contains(@id,"menu-item-20")]/a')
14         for a in a_list:
15             item = DaomuItem()
16             parent_title = a.xpath('./text()').get()
17             parent_url = a.xpath('./@href').get()
18             item['directory'] = './novel/{}/'.format(parent_title)
19             # 创建对应文件夹
20             if not os.path.exists(item['directory']):
21                 os.makedirs(item['directory'])
22             # 交给调度器入队列
23             yield scrapy.Request(url=parent_url, meta={'meta_1':item},
24                                 callback=self.detail_page)
25
26         # 返回了11个response,调用了这个函数
27     def detail_page(self, response):
28         """二级页面解析函数：提取小标题、小链接"""
29         # 把item接收
30         meta_1 = response.meta['meta_1']
31         art_list = response.xpath('//article')
32         for art in art_list:
33             # 只要有继续交往调度器的请求,就必须新建item对象
34             item = DaomuItem()
35             item['son_title'] = art.xpath('./a/text()').get()
36             son_url = art.xpath('./a/@href').get()
37             item['directory'] = meta_1['directory']
```

```

37         # 再次交给调度器入队列
38         yield scrapy.Request(url=son_url, meta={'item':item},
callback=self.get_content)
39
40     # 盗墓笔记1: 传过来了75个response
41     # 盗墓笔记2: 传过来了 n 个response
42     # ... ...
43     def get_content(self, response):
44         """三级页面解析函数: 提取具体小说内容"""
45         item = response.meta['item']
46         # content_list: ['段落1','段落2','段落3',...]
47         content_list = response.xpath('//article[@class="article-
content"]/p/text()').extract()
48         item['content'] = '\n'.join(content_list)
49
50         # 至此,一条item数据全部提取完成
51         yield item

```

• 4、管道文件实现数据处理 - pipelines.py

```

1 class DaomuPipeline(object):
2     def process_item(self, item, spider):
3         # filename: ./novel/盗墓笔记1:七星鲁王宫/七星鲁王_第一章_血尸.txt
4         filename = '{}{}.txt'.format(item['directory'],
item['son_title'].replace(' ', '_'))
5         with open(filename, 'w') as f:
6             f.write(item['content'])
7
8         return item

```

• 5、全局配置 - setting.py

```

1 ROBOTSTXT_OBEY = False
2 DOWNLOAD_DELAY = 0.5
3 DEFAULT_REQUEST_HEADERS = {
4     'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
5     'Accept-Language': 'en',
6     'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML,
like Gecko) Chrome/80.0.3987.149 Safari/537.36'
7 }
8 ITEM_PIPELINES = {
9     'Daomu.pipelines.DaomuPipeline': 300,
10 }

```

• 完整流程

```

1 【1】 scrapy startproject Tencent
2 【2】 cd Tencent
3 【3】 scrapy genspider tencent tencent.com
4 【4】 items.py(定义爬取数据结构)

```

```

5     import scrapy
6     class TencentItem(scrapy.Item):
7         name = scrapy.Field()
8         address = scrapy.Field()
9
10    【5】 tencent.py (写爬虫文件)
11    import scrapy
12    from ..items import TencentItem
13    class TencentSpider(scrapy.Spider):
14        name = 'tencent'
15        allowed_domains = ['tencent.com']
16        start_urls = ['']
17        def parse(self, response):
18            item = TencentItem()
19            item['name'] = xxxx
20            yield item
21
22    【6】 pipelines.py (数据处理)
23    class TencentPipeline(object):
24        def process_item(self, item, spider):
25            return item
26
27    【7】 settings.py (全局配置)
28
29    【8】 run.py
30    from scrapy import cmdline
31    cmdline.execute('scrapy crawl tencent'.split())

```

我们必须记住

- 熟练记住

```

1    【1】 响应对象response属性及方法
2        1.1) response.text : 获取响应内容 - 字符串
3        1.2) response.body : 获取bytes数据类型
4        1.3) response.xpath('')
5        1.4) response.xpath('').extract() : 提取文本内容, 将列表中所有元素序列化为Unicode字符串
6        1.5) response.xpath('').extract_first() : 序列化提取列表中第1个文本内容
7        1.6) response.xpath('').get() : 提取列表中第1个文本内容 (等同于extract_first())
8
9    【2】 settings.py中常用变量
10       2.1) 设置数据导出编码 (主要针对于json文件)
11           FEED_EXPORT_ENCODING = 'utf-8'
12       2.2) 设置User-Agent
13           USER_AGENT = ''
14       2.3) 设置最大并发数 (默认为16)
15           CONCURRENT_REQUESTS = 32
16       2.4) 下载延迟时间 (每隔多长时间请求一个网页)
17           DOWNLOAD_DELAY = 0.5
18       2.5) 请求头
19           DEFAULT_REQUEST_HEADERS = {'Cookie' : 'xxx'}
20       2.6) 添加项目管道

```



```

21         ITEM_PIPELINES = {'目录名.pipelines.类名' : 优先级}
22     2.7) cookie (默认禁用,取消注释-True|False都为开启)
23         COOKIES_ENABLED = False

```

爬虫项目启动方式

- 启动方式

```

1  【1】方式一:基于start_urls
2      1.1) 从爬虫文件(spider)的start_urls变量中遍历URL地址交给调度器入队列,
3      1.2) 把下载器返回的响应对象(response)交给爬虫文件的parse(self, response)函数处理
4
5  【2】方式二
6      重写start_requests()方法,从此方法中获取URL,交给指定的callback解析函数处理
7      2.1) 去掉start_urls变量
8      2.2) def start_requests(self):
9              # 生成要爬取的URL地址,利用scrapy.Request()方法交给调度器

```

数据持久化存储

- MySQL-MongoDB-Json-csv

```

1  *****存入MySQL、MongoDB*****
2
3  【1】在setting.py中定义相关变量
4  【2】pipelines.py中新建管道类,并导入settings模块
5      def open_spider(self, spider):
6          # 爬虫开始执行1次,用于数据库连接
7
8      def process_item(self, item, spider):
9          # 用于处理抓取的item数据
10         return item
11
12     def close_spider(self, spider):
13         # 爬虫结束时执行1次,用于断开数据库连接
14
15  【3】settings.py中添加此管道
16     ITEM_PIPELINES = {'':200}
17
18  【注意】 process_item() 函数中一定要 return item
19
20  *****存入JSON、CSV文件*****
21  scrapy crawl maoyan -o maoyan.csv
22  scrapy crawl maoyan -o maoyan.json
23  【注意】
24      存入json文件时候需要添加变量(settings.py) : FEED_EXPORT_ENCODING = 'utf-8'

```

分布式爬虫

- 分布式爬虫介绍
- 多台主机共享一个爬取队列

```
1  【1】 原理
2      多台主机共享1个爬取队列
3      scrapy的调度器本身不支持分布式
4
5  【2】 实现
6      2.1) 重写scrapy调度器(scrapy_redis模块)
7      2.2) sudo pip3 install scrapy_redis
```

- 为什么使用redis

```
1  【1】 Redis基于内存,速度快
2  【2】 Redis非关系型数据库,Redis中集合,存储每个request的指纹
```

scrapy_redis详解

- GitHub地址

```
1  https://github.com/rmax/scrapy-redis
```

- settings.py说明

```
1  # 重新指定调度器: 启用Redis调度存储请求队列
2  SCHEDULER = "scrapy_redis.scheduler.Scheduler"
3
4  # 重新指定去重机制: 确保所有的爬虫通过Redis去重
5  DUPEFILTER_CLASS = "scrapy_redis.dupefilter.RFPDupeFilter"
6
7  # 不清除Redis队列: 暂停/恢复/断点续爬(默认清除为False, 设置为True不清除)
8  SCHEDULER_PERSIST = True
9
10 # 优先级队列 (默认)
11 SCHEDULER_QUEUE_CLASS = 'scrapy_redis.queue.PriorityQueue'
12 # 可选用的其它队列
13 # 先进先出
14 SCHEDULER_QUEUE_CLASS = 'scrapy_redis.queue.FifoQueue'
15 # 后进先出
16 SCHEDULER_QUEUE_CLASS = 'scrapy_redis.queue.LifoQueue'
17
18 # redis管道
19 ITEM_PIPELINES = {
20     'scrapy_redis.pipelines.RedisPipeline': 300
21 }
22
23 # 指定连接到redis时使用的端口和地址
24 REDIS_HOST = 'localhost'
25 REDIS_PORT = 6379
```

腾讯招聘分布式改写

- 分布式爬虫完成步骤

- 1 【1】首先完成非分布式scrapy爬虫：正常scrapy爬虫项目抓取
- 2 【2】设置,部署成为分布式爬虫

- 分布式环境说明

- 1 【1】分布式爬虫服务器数量：2（其中1台Windows,1台Ubuntu虚拟机）
- 2 【2】服务器分工：
- 3 2.1) Windows：负责数据抓取
- 4 2.2) Ubuntu：负责URL地址统一管理,同时负责数据抓取

- 腾讯招聘分布式爬虫 - 数据同时存入1个Redis数据库

```
1 【1】完成正常scrapy项目数据抓取（非分布式 - 拷贝之前的Tencent）
2
3 【2】设置settings.py, 完成分布式设置
4 2.1-必须) 使用scrapy_redis的调度器
5     SCHEDULER = "scrapy_redis.scheduler.Scheduler"
6
7 2.2-必须) 使用scrapy_redis的去重机制
8     DUPEFILTER_CLASS = "scrapy_redis.dupefilter.RFPDupeFilter"
9
10 2.3-必须) 定义redis主机地址和端口号
11     REDIS_HOST = '192.168.1.107'
12     REDIS_PORT = 6379
13
14 2.4-非必须) 是否清除请求指纹, True:不清除 False:清除（默认）
15     SCHEDULER_PERSIST = True
16
17 2.5-非必须) 在ITEM_PIPELINES中添加redis管道,数据将会存入redis数据库
18     'scrapy_redis.pipelines.RedisPipeline': 200
19
20 【3】把代码原封不动的拷贝到分布式中的其他爬虫服务器,同时开始运行爬虫
21
22 【结果】：多台机器同时抓取,数据会统一存到Ubuntu的redis中,而且所抓数据不重复
```

- 腾讯招聘分布式爬虫 - 数据存入MySQL数据库

```
1 """和数据存入redis步骤基本一样,只是变更一下管道和MySQL数据库服务器的IP地址"""
2 【1】settings.py
3 1.1) SCHEDULER = 'scrapy_redis.scheduler.Scheduler'
4 1.2) DUPEFILTER_CLASS = 'scrapy_redis.dupefilter.RFPDupeFilter'
5 1.3) SCHEDULER_PERSIST = True
6 1.4) REDIS_HOST = '192.168.1.105'
7 1.5) REDIS_PORT = 6379
8 1.6) ITEM_PIPELINES = {'Tencent.pipelines.TencentMysqlPipeline' : 300}
9 1.7) MYSQL_HOST = '192.168.1.105'
10
```

```
11 【2】将代码拷贝到分布式中所有爬虫服务器
12
13 【3】多台爬虫服务器同时运行scrapy爬虫
14
15 # 赠送腾讯MySQL数据库建库建表语句
16 """
17 create database tencentdb charset utf8;
18 use tencentdb;
19 create table tencenttab(
20 job_name varchar(1000),
21 job_type varchar(200),
22 job_duty varchar(5000),
23 job_require varchar(5000),
24 job_address varchar(200),
25 job_time varchar(200)
26 ) charset=utf8;
27 """
```