

Day01笔记

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, 通用网络爬虫需要遵守robots协议(君子协议)
- 23 示例: <https://www.baidu.com/robots.txt>
- 24 5.2) 聚焦网络爬虫 : 自己写的爬虫程序
- 25
- 26 【6】爬取数据步骤
- 27 6.1) 确定需要爬取的URL地址
- 28 6.2) 由请求模块向URL地址发出请求, 并得到网站的响应
- 29 6.3) 从响应内容中提取所需数据
- 30 a> 所需数据, 保存
- 31 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)
9  AppleWebKit/535.1 (KHTML, like Gecko) Chrome/14.0.835.163 Safari/535.1'}
10 html = requests.get(url=url, headers=headers).content.decode('utf-
11 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
12 # UnicodeDecodeError: utf-8 xxx cannot decode char \xxx in.
13     ignore 可解决
14 # UnicodeEncodeError: gbk code cannot encode char \xxx in,
15     windows 写入文件时常报错误
16     # with open('xxx.txt', 'w', encoding='gb18030') as f:
17 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')
```

正则解析模块re

re模块使用流程

```
1  # 方法一
2  r_list=re.findall('正则表达式',html,re.S)
3  # re.S 让正则的.能够匹配\n换行符
4
5  # 方法二
6  pattern = re.compile('正则表达式',re.S)
7  r_list = pattern.findall(html)
```

- 思考 - 请写出匹配任意一个字符的正则表达式?

```
1 import re
2 # 方法一
3 pattern = re.compile('[\s\S]')
4 result = pattern.findall(html)
5
6 # 方法二
7 pattern = re.compile('.', re.S)
8 result = pattern.findall(html)
```

- 代码示例

```
1 import re
2
3 html = '''
4 <div><p>九霄龙吟惊天变</p></div>
5 <div><p>风云际会潜水游</p></div>
6 '''
7 # 贪婪匹配
8 p = re.compile('<div><p>.*</p></div>', re.S)
9 r_list = p.findall(html)
10 print(r_list)
11
12 # 非贪婪匹配
13 p = re.compile('<div><p>.*?</p></div>', re.S)
14 r_list = p.findall(html)
15 print(r_list)
```

正则表达式分组

- 作用

1 在完整的模式中定义子模式，将每个圆括号中子模式匹配出来的结果提取出来

- 示例代码

```
1 import re
2
3 s = 'A B C D'
4 p1 = re.compile('\w+\s+\w+')
5 print(p1.findall(s))
6 # 分析结果是什么??
7 # ['A B', 'C D']
8
9 p2 = re.compile('(\w+)\s+(\w+)')
10 print(p2.findall(s))
11 # 第一步: ['A B', 'C D']
12 # 第二步: ['A', 'C']
13
```

```

14
15 p3 = re.compile('(\w+)\s+(\w+)')
16 print(p3.findall(s))
17 # 第一步: ['A B', 'C D']
18 # 第二步: [('A','B'), ('C','D')]

```

• 分组总结

- 1、在网页中,想要什么内容,就加()
- 2、先按整体正则匹配,然后再提取分组()中的内容
- 3 如果有2个及以上分组(),则结果中以元组形式显示 [(), (), ()]

• 课堂练习

```

1 # 从如下html代码结构中完成如下内容信息的提取:
2 问题1 : [('Tiger', ' Two...'), ('Rabbit', 'Small..')]
3 问题2 :
4     动物名称 : Tiger
5     动物描述 : Two tigers two tigers run fast
6     *****
7     动物名称 : Rabbit
8     动物描述 : Small white rabbit white and white

```

• 页面结构如下

```

1 <div class="animal">
2     <p class="name">
3         <a title="Tiger"></a>
4     </p>
5     <p class="content">
6         Two tigers two tigers run fast
7     </p>
8 </div>
9
10 <div class="animal">
11     <p class="name">
12         <a title="Rabbit"></a>
13     </p>
14
15     <p class="content">
16         Small white rabbit white and white
17     </p>
18 </div>

```

• 练习答案

```

1 import re
2
3 html = '''<div class="animal">
4     <p class="name">
5         <a title="Tiger"></a>
6     </p>
7
8     <p class="content">
9         Two tigers two tigers run fast

```

```

10     </p>
11 </div>
12
13 <div class="animal">
14     <p class="name">
15         <a title="Rabbit"></a>
16     </p>
17
18     <p class="content">
19         Small white rabbit white and white
20     </p>
21 </div>'''
22
23 p = re.compile('<div class="animal">.*?title="(.*?)".*?content">(.*?)</p>.*?</div>', re.S)
24 r_list = p.findall(html)
25
26 for rt in r_list:
27     print('动物名称:', rt[0].strip())    # strip() 去掉字符串两头的空白包括\n\t
    和空格
28     print('动物描述:', rt[1].strip())
29     print('*' * 50)
30
31
32 # 把想要提取的数据先复制出来，再按需删除，删除的地方加上.*?, 需要提取的地方加(.*?)

```

猫眼电影top100抓取案例

• 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>.*?releasetime">(.*?)</p>
11
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()

```

数据持久化 - MySQL

- pymysql回顾

```

1  import pymysql
2
3  db =
    pymysql.connect('localhost','root','123456','maoyandb',charset='utf8')
4  cursor = db.cursor()
5
6  ins = 'insert into filmtab values(%s,%s,%s)'
7  cursor.execute(ins,['霸王别姬','张国荣','1993'])
8
9  db.commit()
10 cursor.close()
11 db.close()

```

• 练习 - 将电影信息存入MySQL数据库

```

1  【1】 提前建库建表
2  mysql -h127.0.0.1 -uroot -p123456
3  create database maoyandb charset utf8;
4  use maoyandb;
5  create table maoyantab(
6  name varchar(100),
7  star varchar(300),
8  time varchar(100)
9  )charset=utf8;
10
11 【2】 使用excute()方法将数据存入数据库思路
12     2.1) 在 __init__() 中连接数据库并创建游标对象
13     2.2) 在 save_html() 中将所抓取的数据处理成列表，使用execute()方法写入
14     2.3) 在run() 中等数据抓取完成后关闭游标及断开数据库连接

```

• 汽车之家二手车信息抓取

```

1  【1】 URL地址
2      进入汽车之家官网，点击 二手车
3      即: https://www.che168.com/beijing/a0_0msdgsncngpi11to1cspexx0/
4
5  【2】 抓取目标
6      每辆汽车的
7      2.1) 汽车名称
8      2.2) 行驶里程
9      2.3) 城市
10     2.4) 个人还是商家
11     2.5) 价格
12
13 【3】 抓取前5页

```

• 参考答案

```

1  import requests
2  import re
3  import time
4  import random
5

```

```

6 class CarSpider:
7     def __init__(self):
8         self.url =
          'https://www.che168.com/beijing/a0_0msdgsncncgpi1lto1csp{}exx0/?
          pvareaid=102179#currngpotion'
9         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0;
          WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.138
          Safari/537.36'}
10
11     def get_html(self, url):
12         html = requests.get(url=url,
          headers=self.headers).content.decode('gb2312', 'ignore')
13         self.parse_html(html)
14
15     def parse_html(self, html):
16         pattern = re.compile('<li class="cards-li list-photo-li".*?<div
          class="cards-bottom">.*?<h4 class="card-name">(.*)</h4>.*?<p
          class="cards-unit">(.*)</p>.*?<span class="pierce"><em>(.*)</em>',
          re.S)
17         car_list = pattern.findall(html)
18         self.save_html(car_list)
19
20     def save_html(self, car_list):
21         for car in car_list:
22             print(car)
23
24     def run(self):
25         for i in range(1,6):
26             page_url = self.url.format(i)
27             self.get_html(page_url)
28             time.sleep(random.randint(1,2))
29
30 if __name__ == '__main__':
31     spider = CarSpider()
32     spider.run()

```

请求模块(requests)

```

1 html = requests.get(url=url,headers=headers).text
2 html = requests.get(url=url,headers=headers).content.decode('utf-8')
3
4 with open('xxx.txt','w',encoding='utf-8') as f:
5     f.write(html)

```

解析模块(re)

- 使用流程

```

1 p = re.compile('正则表达式',re.S)
2 r_list = p.findall(html)

```

- 贪婪匹配和非贪婪匹配

- 1 贪婪匹配(默认) : `.*`
- 2 非贪婪匹配 : `.*?`

- 正则表达式分组

- 1 【1】想要什么内容在正则表达式中加()
2 【2】多个分组,先按整体正则匹配,然后再提取()中数据。结果: `[((),(),(),(),())]`

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

```
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))
```

day02笔记

数据持久化 - MySQL

- pymysql回顾

```

1  import pymysql
2
3  db =
    pymysql.connect('localhost','root','123456','maoyandb',charset='utf8')
4  cursor = db.cursor()
5
6  ins = 'insert into filmtab values(%s,%s,%s)'
7  cursor.execute(ins,['霸王别姬','张国荣','1993'])
8
9  db.commit()
10 cursor.close()
11 db.close()

```

数据持久化 - MongoDB

• MongoDB特点

```

1  【1】非关系型数据库,数据以键值对方式存储,端口27017
2  【2】MongoDB基于磁盘存储
3  【3】MongoDB数据类型单一,值为JSON文档,而Redis基于内存,
4      3.1> MySQL数据类型: 数值类型、字符类型、日期时间类型、枚举类型
5      3.2> Redis数据类型: 字符串、列表、哈希、集合、有序集合
6      3.3> MongoDB数据类型: 值为JSON文档
7  【4】MongoDB: 库 -> 集合 -> 文档
8      MySQL : 库 -> 表 -> 表记录

```

• MongoDB常用命令

```

1  Linux进入: mongo
2  >show dbs                - 查看所有库
3  >use 库名                 - 切换库
4  >show collections        - 查看当前库中所有集合
5  >db.集合名.find().pretty() - 查看集合中文档
6  >db.集合名.count()       - 统计文档条数
7  >db.集合名.drop()        - 删除集合
8  >db.dropDatabase()       - 删除当前库

```

• pymongo模块使用

```

1  import pymongo
2
3  # 1.连接对象
4  conn = pymongo.MongoClient(host = 'localhost',port = 27017)
5  # 2.库对象
6  db = conn['maoyandb']
7  # 3.集合对象
8  myset = db['maoyanset']
9  # 4.插入数据库
10 myset.insert_one({'name':'赵敏'})

```

• 练习 - 将电影信息存入MongoDB数据库

```

1  """
2  猫眼电影top100抓取（电影名称、主演、上映时间）

```

```

3  存入mongodb数据库中
4  """
5  import requests
6  import re
7  import time
8  import random
9  import pymongo
10
11 class MaoyanSpider:
12     def __init__(self):
13         self.url = 'https://maoyan.com/board/4?offset={}'
14         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0;
WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.113
Safari/537.36'}
15         # 三个对象: 连接对象、库对象、集合对象
16         self.conn = pymongo.MongoClient('127.0.0.1', 27017)
17         self.db = self.conn['maoyandb']
18         self.myset = self.db['maoyanset2']
19
20     def get_html(self, url):
21         html = requests.get(url=url, headers=self.headers).text
22         # 直接调用解析函数
23         self.parse_html(html)
24
25     def parse_html(self, html):
26         """解析提取数据"""
27         regex = '<div class="movie-item-info">.*?title="(.*?)".*?<p
class="star">(.*?)</p>.*?<p class="releasetime">(.*?)</p>'
28         pattern = re.compile(regex, re.S)
29         r_list = pattern.findall(html)
30         # r_list: [('活着', '牛犇', '2000-01-01'), (), (), ..., ())
31         self.save_html(r_list)
32
33     def save_html(self, r_list):
34         """数据处理函数"""
35         for r in r_list:
36             item = {}
37             item['name'] = r[0].strip()
38             item['star'] = r[1].strip()
39             item['time'] = r[2].strip()
40             print(item)
41             # 存入到mongodb数据库
42             self.myset.insert_one(item)
43
44     def run(self):
45         """程序入口函数"""
46         for offset in range(0, 91, 10):
47             url = self.url.format(offset)
48             self.get_html(url=url)
49             # 控制数据抓取频率: uniform()生成指定范围内的浮点数
50             time.sleep(random.uniform(0, 1))
51
52 if __name__ == '__main__':
53     spider = MaoyanSpider()
54     spider.run()

```

- csv描述

```
1  【1】作用
2      将爬取的数据存放到本地的csv文件中
3
4  【2】使用流程
5      2.1> 打开csv文件
6      2.2> 初始化写入对象
7      2.3> 写入数据(参数为列表)
8
9  【3】示例代码
10     import csv
11     with open('sky.csv','w') as f:
12         writer = csv.writer(f)
13         writer.writerow([])
```

- 示例

```
1  【1】题目描述
2      创建 test.csv 文件，在文件中写入数据
3
4  【2】数据写入 - writerow([])方法
5      import csv
6      with open('test.csv','w') as f: # with
open('test.csv','w',newline='') as f:----->windows里面的写法，因为再wiondows
中每条数据会有一个空行
7          writer = csv.writer(f)
8          writer.writerow(['超哥哥','25'])
9
```

- 练习 - 使用 writerow() 方法将猫眼电影数据存入本地 maoyan.csv 文件

```
1  【1】在 __init__() 中打开csv文件，因为csv文件只需要打开和关闭1次即可
2  【2】在 save_html() 中将所抓取的数据处理成列表，使用writerow()方法写入
3  【3】在run() 中等数据抓取完成后关闭文件
```

- 代码实现

```
1  """
2  猫眼电影top100抓取（电影名称、主演、上映时间）
3  存入csv文件,使用writerow()方法
4  """
5  import requests
6  import re
7  import time
8  import random
9  import csv
10
11  class MaoyanSpider:
12      def __init__(self):
13          self.url = 'https://maoyan.com/board/4?offset={}'
14          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)'}

```

```

15         # 打开文件,初始化写入对象
16         self.f = open('maoyan.csv', 'w', newline='', encoding='utf-8')
17         self.writer = csv.writer(self.f)
18
19     def get_html(self, url):
20         html = requests.get(url=url, headers=self.headers).text
21         # 直接调用解析函数
22         self.parse_html(html)
23
24     def parse_html(self, html):
25         """解析提取数据"""
26         regex = '<div class="movie-item-info">.*?title="(.*?)".*?<p
class="star">(.*?)</p>.*?<p class="releasetime">(.*?)</p>'
27         pattern = re.compile(regex, re.S)
28         r_list = pattern.findall(html)
29         # r_list: [('活着', '牛犇', '2000-01-01'), (), (), ..., ())
30         self.save_html(r_list)
31
32     def save_html(self, r_list):
33         """数据处理函数"""
34         for r in r_list:
35             li = [ r[0].strip(), r[1].strip(), r[2].strip() ]
36             self.writer.writerow(li)
37             print(li)
38
39     def run(self):
40         """程序入口函数"""
41         for offset in range(0, 91, 10):
42             url = self.url.format(offset)
43             self.get_html(url=url)
44             # 控制数据抓取频率:uniform()生成指定范围内的浮点数
45             time.sleep(random.uniform(1,2))
46
47         # 所有数据抓取并写入完成后关闭文件
48         self.f.close()
49
50 if __name__ == '__main__':
51     spider = MaoyanSpider()
52     spider.run()

```

useragent 池

```

1  sudo pip3 install fake_useragent          # 爬虫生成useragent的插件, 大概有250个
   UserAgent
2  from fake_useragent import UserAgent
3
4  headers = {'User-Agent': UserAgent().random}

```

汽车之家数据抓取 - 二级页面

• 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.che168.com/beijing/a0_0msdgsncncgpi1lto1csp1exx0/
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.che168.com/beijing/a0_0msdgsncncgpi1lto1csp1exx0/
5      第2页: https://www.che168.com/beijing/a0_0msdgsncncgpi1lto1csp2exx0/
6      第n页: https://www.che168.com/beijing/a0_0msdgsncncgpi1lto1csp{ }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="
10     (.*)" .*?</li>
11      二级页面正则表达式:<div class="car-box">.*?<h3 class="car-brand-name">
12     (.*)"</h3>.*?<ul class="brand-unit-item fn-clear">.*?<li>.*?<h4>(.*)"
13     </h4>.*?<h4>(.*)"</h4>.*?<h4>(.*)"</h4>.*?<h4>(.*)"</h4>.*?<span
14     class="price" id="overlayPrice">¥(.*)"<b>
15
16  【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.che168.com/beijing/a0_0msdgsncncgpi1lto1csp{}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="
43         (.*)".*?</li>'
44         href_list = self.re_func(regex=one_regex, html=one_html)
45         for href in href_list:
46             two_url = 'https://www.che168.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         two_html = self.get_html(url=two_url)
55         two_regex = '<div class="car-box">.*?<h3 class="car-brand-name">
56         (.*)</h3>.*?<h4>(.*)</h4>.*?<h4>(.*)</h4>.*?<h4>(.*)</h4>.*?<h4>
57         (.*)</h4>.*?<span class="price" id="overlayPrice">¥(.*)<b>
58         # car_list: [('福睿斯', '3万公里', '2016年3月', '手动 / 1.5L', '廊坊',
59         '5.60'),]
60         car_list = self.re_func(regex=two_regex, html=two_html)
61         item = {}
62         item['name'] = car_list[0][0].strip()
63         item['km'] = car_list[0][1].strip()
64         item['time'] = car_list[0][2].strip()
65         item['type'] = car_list[0][3].split('/')[0].strip()
66         item['displace'] = car_list[0][3].split('/')[1].strip()
67         item['address'] = car_list[0][4].strip()
68         item['price'] = car_list[0][5].strip()
69         print(item)
70
71     def run(self):
72         for i in range(1,5):
73             url = self.url.format(i)
74             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 =
19         'https://www.che168.com/beijing/a0_0msdgsncnpgi1lto1csp{}exx0/'
20         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 6.1; WOW64)
21         AppleWebKit/535.1 (KHTML, like Gecko) Chrome/14.0.835.163 Safari/535.1'}
22         self.db =
23         pymysql.connect('localhost', 'root', '123456', 'cardb', charset='utf8')
24         self.cursor = self.db.cursor()
25         # 连接redis去重
26         self.r = redis.Redis(host='localhost', port=6379, db=0)
27
28         # 功能函数1 - 获取响应内容
29         def get_html(self, url):
30             html = requests.get(url=url, headers=self.headers).text
31
32             return html
33
34         # 功能函数2 - 正则解析
35         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.che168.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     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()

```

Chrome浏览器安装插件

- 安装方法

```

1  【1】在线安装
2      1.1> 下载插件 - google访问助手
3      1.2> 安装插件 - google访问助手：Chrome浏览器-设置-更多工具-扩展程序-开发者模式-拖拽(解压后的插件)
4      1.3> 在线安装其他插件 - 打开google访问助手 - google应用商店 - 搜索插件 - 添加即可
5
6  【2】离线安装
7      2.1> 网上下载插件 - xxx.crx 重命名为 xxx.zip
8      2.2> Chrome浏览器-设置-更多工具-扩展程序-开发者模式
9      2.3> 拖拽 插件(或者解压后文件夹) 到浏览器中
10     2.4> 重启浏览器，使插件生效

```

- 爬虫常用插件

```

1  【1】google-access-helper : 谷歌访问助手,可访问 谷歌应用商店
2  【2】xpath Helper: 轻松获取HTML元素的xpath路径
3      打开/关闭: Ctrl + Shift + x
4  【3】JsonView: 格式化输出json格式数据
5  【4】Proxy SwitchyOmega: Chrome浏览器中的代理管理扩展程序

```

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d1节点下第2个dd节点的电影节点
6      //d1[@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 [] : //d1[@class="xxx"] //d1/dd[2]
12     2> 只要获取属性值,加 @ : //d1[@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()
```

- 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()

```

• 小作业

1 | 汽车之家案例使用lxml+xpath实现

Day02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

- 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')

```

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
4  # __init__(self):
5      self.conn = pymongo.MongoClient('IP', 27017)
6      self.db = self.conn['cardb']
7      self.myset = self.db['car_set']
8
9  # save_html(self, r_list):
10     self.myset.insert_one(dict)
11
12 # MongoDB - Command - 库->集合->文档
13 mongo
14 >show dbs
15 >use db_name

```

```
16 >show collections
17 >db.集合名.find().pretty()
18 >db.集合名.count()
19 >db.集合名.drop()
20 >db.dropDatabase()
```

xpath表达式

- 匹配规则

```
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()
```

Day03笔记

- xpath匹配汽车之家数据

```
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
5 【2】匹配汽车之家-汽车详情页中,汽车的
6 2.1) 名称: //div[@class="car-box"]/h3/text()
7 2.2) 里程: //ul/li[1]/h4/text()
8 2.3) 时间: //ul/li[2]/h4/text()
9 2.4) 挡位+排量: //ul/li[3]/h4/text()
10 2.5) 所在地: //ul/li[4]/h4/text()
11 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')
7 # 只要调用 xpath 得到的返回值一定是列表: [] 或者 [<element div at xxxx>,
   <xxxxxxxxxx>] 或者 ['字符串','xxx']
```

- 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  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
30         else None
31         # 评分
32         score_list =
33         table.xpath('.//span[@class="rating_nums"]/text()')
34         item['score'] = score_list[0].strip() if score_list else
35         None
36         # 评价人数
37         nums_list = table.xpath('.//span[@class="p1"]/text()')
38         item['nums'] = nums_list[0][1:-1].strip() if nums_list else
39         None
40         # 类别
41         type_list = table.xpath('.//span[@class="inq"]/text()')
42         item['type'] = type_list[0].strip() if type_list else None
43
44         print(item)
45
46 def run(self):
47     for i in range(5):
48         start = (i - 1) * 25
49         page_url = self.url.format(start)
50         self.get_html(page_url)
51         time.sleep(random.randint(1, 2))
52
53 if __name__ == '__main__':
54     spider = DoubanBookSpider()
55     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          html =
requests.get(url=url, headers=headers, timeout=3).content.decode('utf-
8', 'ignore')
13          self.get_data(html)
14
15      def get_data(self, html):
16          p = etree.HTML(html)
17          # 基准xpath: [<element li at xxx>,<element li>]
18          li_list =
p.xpath('//ul[@class="sellListContent"]/li[@class="clear LOGVIEWDATA
LOGCLICKDATA"]')
19
20          # for遍历, 依次提取每个房源信息, 放到字典item中
21          item = {}
22          for li in li_list:
23              # 名称+区域
24              name_list =
li.xpath('..//div[@class="positionInfo"]/a[1]/text()')
25              item['name'] = name_list[0].strip() if name_list else None
26              address_list =
li.xpath('..//div[@class="positionInfo"]/a[2]/text()')
27              item['address'] = address_list[0].strip() if address_list
else None
28              # 户型+面积+方位+是否精装+楼层+年代+类型
29              # h_list: ['']
30              h_list = li.xpath('..//div[@class="houseInfo"]/text()')
31              if h_list:

```

```

32         info_list = h_list[0].split('|')
33         if len(info_list) == 7:
34             item['model'] = info_list[0].strip()
35             item['area'] = info_list[1].strip()
36             item['direct'] = info_list[2].strip()
37             item['perfect'] = info_list[3].strip()
38             item['floor'] = info_list[4].strip()
39             item['year'] = info_list[5].strip()[:-2]
40             item['type'] = info_list[6].strip()
41         else:
42             item['model'] = item['area'] = item['direct'] =
item['perfect'] = item['floor'] = item['year'] = item['type'] = None
43         else:
44             item['model'] = item['area'] = item['direct'] =
item['perfect'] = item['floor'] = item['year'] = item['type'] = None
45
46         # 总价+单价
47         total_list =
li.xpath('..//div[@class="totalPrice"]/span/text()')
48         item['total'] = total_list[0].strip() if total_list else
None
49         unit_list =
li.xpath('..//div[@class="unitPrice"]/span/text()')
50         item['unit'] = unit_list[0].strip() if unit_list else None
51
52         print(item)
53
54     def run(self):
55         for pg in range(1,101):
56             url = self.url.format(pg)
57             self.parse_html(url)
58             time.sleep(random.randint(1,2))
59
60 if __name__ == '__main__':
61     spider = LianjiaSpider()
62     spider.run()

```

- 持久化到数据库中 - 自己实现

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

代理参数-proxies

- 定义及分类

- 【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

```

1  【1】 获取代理IP网站
2      西刺代理、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)

```

- 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)

```

- 建立自己的代理IP池 - 开放代理 | 私密代理

```

1  """

```

```

2  收费代理:
3      建立开放代理的代理IP池
4  思路:
5      1、获取到开放代理
6      2、依次对每个代理IP进行测试,能用的保存到文件中
7  """
8  import requests
9
10 class ProxyPool:
11     def __init__(self):
12         self.url = '代理网站的API链接'
13         self.headers = {'User-Agent': 'Mozilla/5.0 (Windows NT 10.0;
WOW64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/81.0.4044.122
Safari/537.36'}
14         # 打开文件,用来存放可用的代理IP
15         self.f = open('proxy.txt', 'w')
16
17     def get_html(self):
18         html = requests.get(url=self.url, headers=self.headers).text
19         proxy_list = html.split('\r\n')
20         for proxy in proxy_list:
21             # 依次测试每个代理IP是否可用
22             if self.check_proxy(proxy):
23                 self.f.write(proxy + '\n')
24
25     def check_proxy(self, proxy):
26         """测试1个代理IP是否可用,可用返回True,否则返回False"""
27         test_url = 'http://httpbin.org/get'
28         proxies = {
29             'http': 'http://{0}'.format(proxy),
30             'https': 'https://{0}'.format(proxy)
31         }
32         try:
33             res = requests.get(url=test_url, proxies=proxies,
headers=self.headers, timeout=2)
34             if res.status_code == 200:
35                 print(proxy, '\033[31m可用\033[0m')
36                 return True
37             else:
38                 print(proxy, '无效')
39                 return False
40         except:
41             print(proxy, '无效')
42             return False
43
44     def run(self):
45         self.get_html()
46         # 关闭文件
47         self.f.close()
48
49 if __name__ == '__main__':
50     spider = ProxyPool()
51     spider.run()

```

requests.get()

get()方法参数

```
1 【1】url
2 【2】proxies : 字典(普通代理|私密代理+独享代理)
3 【3】headers : 字典(User-Agent, Cookie, Referer, ... ..)
4 【4】timeout : 数字(超过指定时间网站未响应则抛出异常)
5 【5】cookies : 字典
```

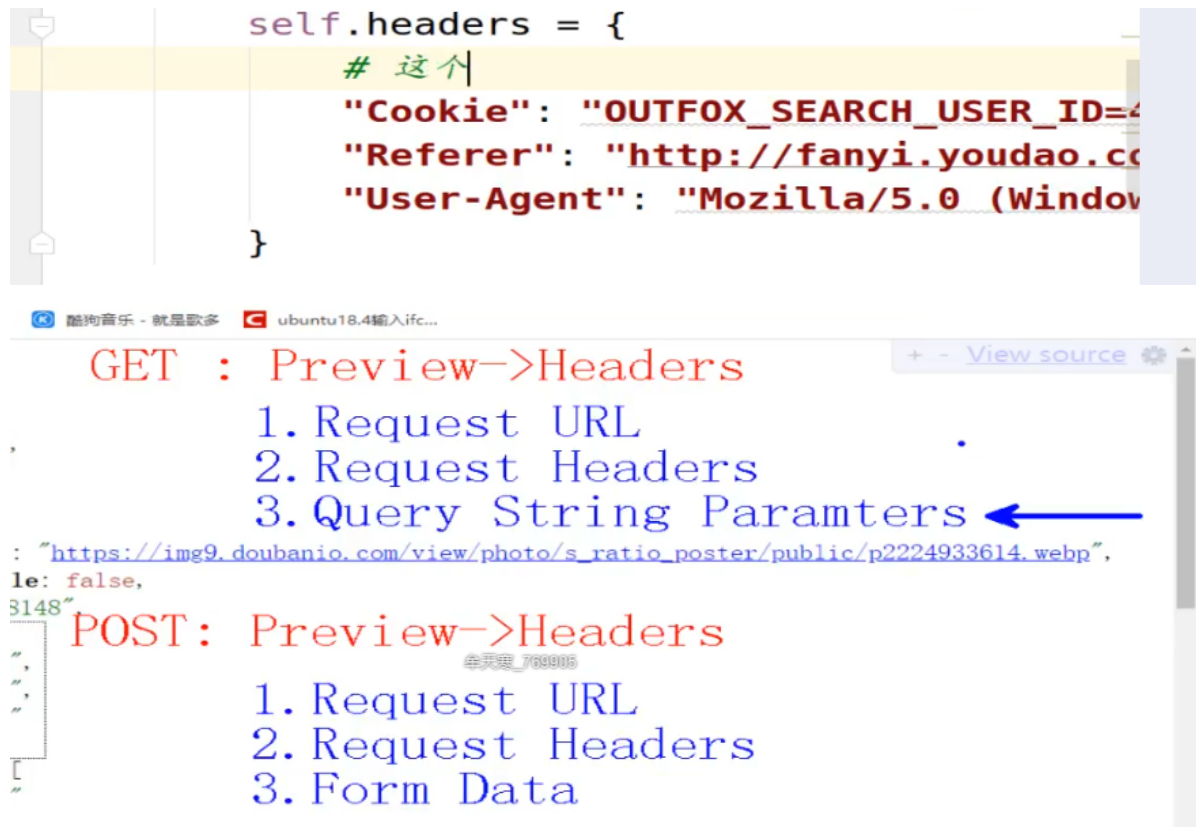
响应对象属性

```
1 【1】text : 字符串
2 【2】content : 获取bytes数据类型
3 【3】status_code : HTTP响应码
4 【4】url : 返回实际数据的URL地址
```

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
8
9
10     # 请求头中最常见的需要检查的三个键
```



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  post请求就找三个数据，url,headers,formdata
4
5  【2】控制台常用选项
6    2.1) Network: 抓取网络数据包
7      a> ALL: 抓取所有的网络数据包
8      b> XHR: 抓取异步加载的网络数据包
9      c> JS : 抓取所有的JS文件
10   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
11   2.3) Console: 交互模式，可对JavaScript中的代码进行测试
12
13  【3】抓取具体网络数据包后
14    3.1) 单击左侧网络数据包地址，进入数据包详情，查看右侧
15    3.2) 右侧:
16      a> Headers: 整个请求信息
17            General、Response Headers、Request Headers、Query String、Form
18            Data
19      b> Preview: 对响应内容进行预览
20      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】浏览器F12开启网络抓包,Network-All,页面翻译单词后找Form表单数据
2  【2】在页面中多翻译几个单词，观察Form表单数据变化（有数据是加密字符串）
3  【3】刷新有道翻译页面，抓取并分析JS代码（本地JS加密）
4  【4】找到JS加密算法，用Python按同样方式加密生成加密数据
5  【5】将Form表单数据处理为字典，通过requests.post()的data参数发送
```

- 具体实现

1、开启F12抓包，找到Form表单数据如下：

```
1  i: 喵喵叫
2  from: AUTO
3  to: AUTO
4  smartresult: dict
5  client: fanyideskweb
6  salt: 15614112641250
7  sign: 94008208919faa19bd531acde36aac5d
8  ts: 1561411264125
```

```
9  bv: f4d62a2579ebb44874d7ef93ba47e822
10 doctype: json
11 version: 2.1
12 keyfrom: fanyi.web
13 action: FY_BY_REALTIME
14
15
16      # 32位是MD5加密      40位是SHA1加密
```

2、在页面中多翻译几个单词，观察Form表单数据变化

```
1  salt: 15614112641250
2  sign: 94008208919faa19bd531acde36aac5d
3  ts: 1561411264125
4  bv: f4d62a2579ebb44874d7ef93ba47e822
5  # 但是bv的值不变
```

3、一般为本地js文件加密，刷新页面，找到js文件并分析JS代码

```
1  【方法1】：Network - JS选项 - 搜索关键词salt
2  【方法2】：控制台右上角 - Search - 搜索salt - 查看文件 - 格式化输出
3
4  【结果】：最终找到相关JS文件：fanyi.min.js
```

4、打开JS文件，分析加密算法，用Python实现

```
1  【ts】经过分析为13位的时间戳，字符串类型
2      js代码实现) "" + (new Date).getTime()
3      python实现) str(int(time.time()*1000))
4
5  【salt】
6      js代码实现) ts + parseInt(10 * Math.random(), 10);
7      python实现) ts + str(random.randint(0,9))
8
9  【sign】('设置断点调试，来查看 e 的值，发现 e 为要翻译的单词')
10     js代码实现) n.md5("fanyideskweb" + e + salt + "n%A-rKaT5fb[Gy?;N5@Tj")
11     python实现)
12     from hashlib import md5
13     string = "fanyideskweb" + e + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
14     s = md5()
15     s.update(string.encode())
16     sign = s.hexdigest()
```

4、pycharm中正则处理headers和formdata

```
1  【1】pycharm进入方法：Ctrl + r，选中 Regex
2  【2】处理headers和formdata
3      (.?): (.*)
4      "$1": "$2",
5  【3】点击 Replace All
```



5、代码实现

```

1 import requests
2 import time
3 import random
4 from hashlib import md5
5
6 class YdSpider(object):
7     def __init__(self):
8         # url一定为F12抓到的 headers -> General -> Request URL
9         self.url = 'http://fanyi.youdao.com/translate_o?
smartresult=dict&smartresult=rule'
10        self.headers = {
11            # 检查频率最高 - 3个
12            "Cookie": "OUTFOX_SEARCH_USER_ID=970246104@10.169.0.83;
OUTFOX_SEARCH_USER_ID_NCOO=570559528.1224236;
_ntes_nnid=96bc13a2f5ce64962adfd6a278467214,1551873108952;
JSESSIONID=aaae9i7p1XPlKaJH_gkYw; td_cookie=18446744072941336803;
SESSION_FROM_COOKIE=unknown; __rl__test__cookies=1565689460872",
13            "Referer": "http://fanyi.youdao.com/",
14            "User-Agent": "Mozilla/5.0 (Windows NT 10.0; win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/76.0.3809.100 Safari/537.36",
15        }
16
17        # 获取salt,sign,ts
18        def get_salt_sign_ts(self,word):
19            # ts
20            ts = str(int(time.time()*1000))
21            # salt
22            salt = ts + str(random.randint(0,9))
23            # sign
24            string = "fanyideskweb" + word + salt + "n%A-rKaT5fb[Gy?;N5@Tj"
25            s = md5()
26            s.update(string.encode())
27            sign = s.hexdigest()
28
29            return salt,sign,ts
30
31        # 主函数
32        def attack_yd(self,word):
33            # 1. 先拿到salt,sign,ts
34            salt,sign,ts = self.get_salt_sign_ts(word)
35            # 2. 定义form表单数据为字典: data={}
36            # 检查了salt sign
37            data = {
38                "i": word,
39                "from": "AUTO",
40                "to": "AUTO",
41                "smartresult": "dict",
42                "client": "fanyideskweb",
43                "salt": salt,
44                "sign": sign,
45                "ts": ts,
46                "bv": "7e3150ecbdf9de52dc355751b074cf60",
47                "doctype": "json",
48                "version": "2.1",
49                "keyfrom": "fanyi.web",
50                "action": "FY_BY_REALTIME",
51            }
52            # 3. 直接发请求:requests.post(url,data=data,headers=xxx)

```

```

53     html = requests.post(
54         url=self.url,
55         data=data,
56         headers=self.headers
57     ).json()
58     # res.json() 将json格式的字符串转为python数据类型
59     result = html['translateResult'][0][0]['tgt']
60
61     print(result)
62
63     # 主函数
64     def run(self):
65         # 输入翻译单词
66         word = input('请输入要翻译的单词:')
67         self.attack_yd(word)
68
69 if __name__ == '__main__':
70     spider = YdSpider()
71     spider.run()

```

今日作业

- 1 【1】将猫眼电影使用 lxml + xpath 实现
- 2 【2】将汽车之家案例使用 lxml + xpath 实现
- 3 【3】完善链家二手房案例，使用 lxml + xpath
- 4 【4】将 周测题 - 电影天堂案例 使用 lxml + xpath 实现
- 5 【5】将 周测题 - 4567TV案例 使用 lxml + xpath 实现
- 6 【6】抓取西刺免费高匿代理，并测试是否可用来建立自己的代理IP池
- 7 <https://www.xicidaili.com/nn/>
- 8 【7】仔细熟悉有道翻译案例抓包及流程分析

Day03回顾

目前反爬总结

- 反爬虫梳理

- 1 【1】Headers反爬虫
- 2 1.1) 检查: Cookie、Referer、User-Agent
- 3 1.2) 解决方案: 通过F12获取headers,传给requests.get()方法
- 4
- 5 【2】IP限制
- 6 2.1) 网站根据IP地址访问频率进行反爬,短时间内限制IP访问
- 7 2.2) 解决方案:
- 8 a) 构造自己IP代理池,每次访问随机选择代理,经常更新代理池
- 9 b) 购买开放代理或私密代理IP
- 10 c) 降低爬取的速度
- 11
- 12 【3】User-Agent限制
- 13 3.1) 类似于IP限制, 检测频率

```
14 3.2) 解决方案：构造自己的User-Agent池,每次访问随机选择
15 a> fake_useragent模块
16 b> 新建py文件,存放大量User-Agent
17
18 【4】对响应内容做处理
19 4.1) 页面结构和响应内容不同
20 4.2) 解决方案：打印并查看响应内容,用xpath或正则做处理
```

Day04笔记

动态加载数据抓取-Ajax

- 特点

```
1 【1】右键 -> 查看网页源码中没有具体数据
2 【2】滚动鼠标滑轮或其他动作时加载,或者页面局部刷新
```

- 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 2.1) 控制台中 XHR : 异步加载的数据包
4 2.2) XHR -> QueryStringParameters(查询参数)
```

豆瓣电影数据抓取案例

- 目标

```
1 【1】地址：豆瓣电影 - 排行榜 - 剧情
2 【2】目标：电影名称、电影评分
```

- F12抓包 (XHR)

```
1 【1】Request URL(基准URL地址) : https://movie.douban.com/j/chart/top_list?
2 【2】Query String(查询参数)
3 # 抓取的查询参数如下:
4 type: 13 # 电影类型
5 interval_id: 100:90
6 action: ''
7 start: 0 # 每次加载电影的起始索引值 0 20 40 60
8 limit: 20 # 每次加载的电影数量
```

- 代码实现 - 全站抓取

```
1 """
2 豆瓣电影 - 全站抓取
3 """
4 import requests
5 from fake_useragent import UserAgent
6 import time
7 import random
8 import re
9 import json
```

```

10
11 class DoubanSpider:
12     def __init__(self):
13         self.url = 'https://movie.douban.com/j/chart/top_list?'
14         self.i = 0
15         # 存入json文件
16         self.f = open('douban.json', 'w', encoding='utf-8')
17         self.all_film_list = []
18
19     def get_agent(self):
20         """获取随机的User-Agent"""
21         return UserAgent().random
22
23     def get_html(self, params):
24         headers = {'User-Agent': self.get_agent()}
25         html = requests.get(url=self.url, params=params,
headers=headers).text
26         # 把json格式的字符串转为python数据类型
27         html = json.loads(html)
28
29         self.parse_html(html)
30
31     def parse_html(self, html):
32         """解析"""
33         # html: [{},{},{},{}]
34         item = {}
35         for one_film in html:
36             item['rank'] = one_film['rank']
37             item['title'] = one_film['title']
38             item['score'] = one_film['score']
39             print(item)
40             self.all_film_list.append(item)
41             self.i += 1
42
43     def run(self):
44         # d: {'剧情': '11', '爱情': '13', '喜剧': '5', ..., ...}
45         d = self.get_d()
46         # 1、给用户提示,让用户选择
47         menu = ''
48         for key in d:
49             menu += key + '|'
50         print(menu)
51         choice = input('请输入电影类别: ')
52         if choice in d:
53             code = d[choice]
54             # 2、total: 电影总数
55             total = self.get_total(code)
56             for start in range(0, total, 20):
57                 params = {
58                     'type': code,
59                     'interval_id': '100:90',
60                     'action': '',
61                     'start': str(start),
62                     'limit': '20'
63                 }
64                 self.get_html(params=params)
65                 time.sleep(random.randint(1, 2))
66

```

```

67         # 把数据存入json文件
68         json.dump(self.all_film_list, self.f, ensure_ascii=False)
69         self.f.close()
70         print('数量:', self.i)
71     else:
72         print('请做出正确的选择')
73
74     def get_d(self):
75         """{'剧情': '11', '爱情': '13', '喜剧': '5', ..., ...}"""
76         url = 'https://movie.douban.com/chart'
77         html = requests.get(url=url, headers={'User-
Agent': self.get_agent()}).text
78         regex = '<span><a href=".*?type_name=(.*?)&type=
(.*?)&interval_id=100:90&action=">'
79         pattern = re.compile(regex, re.S)
80         # r_list: [('剧情', '11'), ('喜剧', '5'), ('爱情', '13')]... ...]
81         r_list = pattern.findall(html)
82         # d: {'剧情': '11', '爱情': '13', '喜剧': '5', ..., ...}
83         d = {}
84         for r in r_list:
85             d[r[0]] = r[1]
86
87         return d
88
89     def get_total(self, code):
90         """获取某个类别下的电影总数"""
91         url = 'https://movie.douban.com/j/chart/top_list_count?type=
{}&interval_id=100%3A90'.format(code)
92         html = requests.get(url=url, headers={'User-
Agent': self.get_agent()}).text
93         html = json.loads(html)
94
95         return html['total']
96
97 if __name__ == '__main__':
98     spider = DoubanSpider()
99     spider.run()

```

json解析模块

- json.loads(json)

```

1  【1】作用 : 把json格式的字符串转为Python数据类型
2
3  【2】示例 : html = json.loads(res.text)

```

- json.dump(python,f,ensure_ascii=False)

```

1  【1】作用
2      把python数据类型 转为 json格式的字符串,一般让你把抓取的数据保存为json文件时使用
3      # 不加 S 是保存到文件之中, dumps 是在程序中使用
4  【2】参数说明
5      2.1) 第1个参数: python类型的数据(字典, 列表等)
6      2.2) 第2个参数: 文件对象
7      2.3) 第3个参数: ensure_ascii=False 序列化时编码,可以显示中文
8

```

```

9  【3】示例代码
10  # 示例1
11  import json
12
13  item = {'name':'QQ','app_id':1}
14  with open('小米.json','a') as f:
15      json.dump(item,f,ensure_ascii=False)
16
17  # 示例2
18  import json
19
20  item_list = []
21  for i in range(3):
22      item = {'name':'QQ','id':i}
23      item_list.append(item)
24
25  with open('xiaomi.json','a') as f:
26      json.dump(item_list,f,ensure_ascii=False)

```

• json模块总结

```

1  # 爬虫最常用
2  【1】数据抓取 - json.loads(html)
3      将响应内容由：json 转为 python
4  【2】数据保存 - json.dump(item_list,f,ensure_ascii=False)
5      将抓取的数据保存到本地 json文件
6
7  # 抓取数据一般处理方式
8  【1】txt文件
9  【2】csv文件
10 【3】json文件
11 【4】MySQL数据库
12 【5】MongoDB数据库
13 【6】Redis数据库

```

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

```

知识点回顾

• 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?
type=13&interval_id=100%3A90&action=&start={}&limit=20'
15         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
地址时出现重复判断
36             self.lock.acquire()
37             if not self.q.empty():
38                 headers = {'User-Agent': self.get_agent()}
39                 url = self.q.get()
40                 self.lock.release()
41
42                 html = requests.get(url=url, headers=headers).json()
43                 self.parse_html(html)
44             else:
45                 # 如果队列为空,则最终必须释放锁
46                 self.lock.release()
47                 break
48
49     def parse_html(self, html):
50         """解析"""
51         # html: [{},{},{},{}]
52         item = {}
53         for one_film in html:
54             item['rank'] = one_film['rank']
55             item['title'] = one_film['title']
56             item['score'] = one_film['score']
57             print(item)
58             # 加锁 + 释放锁
59             self.lock.acquire()
60             self.i += 1
61             self.lock.release()
62
63     def run(self):
64         # 先让URL地址入队列
65         self.url_in()
66         # 创建多个线程,开干吧
67         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))

```

selenium+PhantomJS/Chrome/Firefox

• selenium

```

1  【1】定义
2      1.1) 开源的web自动化测试工具
3
4  【2】用途
5      2.1) 对web系统进行功能性测试,版本迭代时避免重复劳动
6      2.2) 兼容性测试(测试web程序在不同操作系统和不同浏览器中是否运行正常)
7      2.3) 对web系统进行大数量测试
8
9  【3】特点
10     3.1) 可根据指令操控浏览器
11     3.2) 只是工具,必须与第三方浏览器结合使用
12
13 【4】安装
14     4.1) Linux: sudo pip3 install selenium
15     4.2) windows: python -m pip install selenium

```

• PhantomJS浏览器

```

1  【1】定义
2      phantomjs为无界面浏览器(又称无头浏览器),在内存中进行页面加载,高效
3
4  【2】下载地址
5      2.1) chromedriver : 下载对应版本
6          http://npm.taobao.org/mirrors/chromedriver/
7
8      2.2) geckodriver
9          https://github.com/mozilla/geckodriver/releases
10
11     2.3) phantomjs
12         https://phantomjs.org/download.html
13
14 【3】Ubuntu安装
15     3.1) 下载后解压 : tar -zxvf geckodriver.tar.gz
16
17     3.2) 拷贝解压后文件到 /usr/bin/ (添加环境变量)
18         sudo cp geckodriver /usr/bin/
19
20     3.3) 添加可执行权限
21         sudo chmod 777 /usr/bin/geckodriver
22
23 【4】windows安装
24     4.1) 下载对应版本的phantomjs、chromedriver、geckodriver

```

```

25     4.2) 把chromedriver.exe拷贝到python安装目录的Scripts目录下(添加到系统环境
      变量)
26         # 查看python安装路径: where python
27     4.3) 验证
28         cmd命令行: chromedriver
29
30     *****总结*****
31     【1】解压 - 放到用户主目录(chromedriver、geckodriver、phantomjs)
32     【2】拷贝 - sudo cp /home/tarena/chromedriver /usr/bin/
33     【3】权限 - sudo chmod 777 /usr/bin/chromedriver
34
35     # 验证
36     【Ubuntu | windows】
37     ipython3
38     from selenium import webdriver
39     webdriver.Chrome()
40     或者
41     webdriver.Firefox()
42
43     【mac】
44     ipython3
45     from selenium import webdriver
46     webdriver.Chrome(executable_path='/Users/xxx/chromedriver')
47     或者
48     webdriver.Firefox(executable_path='/User/xxx/geckodriver')

```

• 示例代码

```

1  """示例代码一: 使用 selenium+浏览器 打开百度"""
2
3  # 导入selenium的webdriver接口
4  from selenium import webdriver
5  import time
6
7  # 创建浏览器对象
8  browser = webdriver.Chrome()
9  browser.get('http://www.baidu.com/')
10 # 5秒钟后关闭浏览器
11 time.sleep(5)
12 browser.quit()

```

```

1  """示例代码二: 打开百度, 搜索赵丽颖, 点击搜索, 查看"""
2
3  from selenium import webdriver
4  import time
5
6  # 1.创建浏览器对象 - 已经打开了浏览器
7  browser = webdriver.Chrome()
8  # 2.输入: http://www.baidu.com/
9  browser.get('http://www.baidu.com/')
10 # 3.找到搜索框, 向这个节点发送文字: 赵丽颖
11 browser.find_element_by_xpath('//*[@id="kw"]').send_keys('赵丽颖')
12 # 4.找到 百度一下 按钮, 点击一下
13 browser.find_element_by_xpath('//*[@id="su"]').click()

```

• 浏览器对象(browser)方法

```

1 【1】 browser.get(url=url) - 地址栏输入url地址并确认
2 【2】 browser.quit() - 关闭浏览器
3 【3】 browser.close() - 关闭当前页
4 【4】 browser.page_source - HTML结构源码
5 【5】 browser.page_source.find('字符串')
6 从html源码中搜索指定字符串,没有找到返回: -1,经常用于判断是否为最后一页
7 【6】 browser.maximize_window() - 浏览器窗口最大化

```

• 定位节点八种方法

```

1 【1】 单元素查找('结果为1个节点对象')
2 1.1) 【最常用】 browser.find_element_by_id('id属性值')
3 1.2) 【最常用】 browser.find_element_by_name('name属性值')
4 1.3) 【最常用】 browser.find_element_by_class_name('class属性值')
5 1.4) 【最万能】 browser.find_element_by_xpath('xpath表达式')
6 1.5) 【匹配a节点时常用】 browser.find_element_by_link_text('链接文本')
7 1.6) 【匹配a节点时常用】 browser.find_element_by_partial_link_text('部分链接文本')
8 1.7) 【最没用】 browser.find_element_by_tag_name('标记名称')
9 1.8) 【较常用】 browser.find_element_by_css_selector('css表达式')
10
11 【2】 多元素查找('结果为[节点对象列表]')
12 2.1) browser.find_elements_by_id('id属性值')
13 2.2) browser.find_elements_by_name('name属性值')
14 2.3) browser.find_elements_by_class_name('class属性值')
15 2.4) browser.find_elements_by_xpath('xpath表达式')
16 2.5) browser.find_elements_by_link_text('链接文本')
17 2.6) browser.find_elements_by_partial_link_text('部分链接文本')
18 2.7) browser.find_elements_by_tag_name('标记名称')
19 2.8) browser.find_elements_by_css_selector('css表达式')

```

• 猫眼电影示例

```

1 from selenium import webdriver
2 import time
3
4 url = 'https://maoyan.com/board/4'
5 browser = webdriver.Chrome()
6 browser.get(url)
7
8 def get_data():
9     # 基准xpath: [<selenium xxx li at xxx>,<selenium xxx li at>]
10    li_list = browser.find_elements_by_xpath('//*[@id="app"]/div/div/div[1]/dl/dd')
11    for li in li_list:
12        item = {}
13        # info_list: ['1', '霸王别姬', '主演: 张国荣', '上映时间: 1993-01-01', '9.5']
14        info_list = li.text.split('\n')
15        item['number'] = info_list[0]
16        item['name'] = info_list[1]
17        item['star'] = info_list[2]
18        item['time'] = info_list[3]
19        item['score'] = info_list[4]
20
21    print(item)

```

```

22
23 while True:
24     get_data()
25     try:
26         browser.find_element_by_link_text('下一页').click()
27         time.sleep(2)
28     except Exception as e:
29         print('恭喜你!抓取结束')
30         browser.quit()
31         break

```

• 节点对象操作

```

1  【1】文本框操作
2      1.1) node.send_keys('') - 向文本框发送内容
3      1.2) node.clear() - 清空文本
4      1.3) node.get_attribute('value') - 获取文本内容
5
6  【2】按钮操作
7      1.1) node.click() - 点击
8      1.2) node.is_enabled() - 判断按钮是否可用
9      1.3) node.get_attribute('value') - 获取按钮文本

```

chromedriver设置无界面模式

```

1  from selenium import webdriver
2
3  options = webdriver.ChromeOptions()
4  # 添加无界面参数
5  options.add_argument('--headless')
6  browser = webdriver.Chrome(options=options)

```

selenium - 鼠标操作

```

1  from selenium import webdriver
2  # 导入鼠标事件类
3  from selenium.webdriver import ActionChains
4
5  driver = webdriver.Chrome()
6  driver.get('http://www.baidu.com/')
7
8  # 移动到 设置, perform()是真正执行操作, 必须有
9  element = driver.find_element_by_xpath('//*[@id="u1"]/a[8]')
10 ActionChains(driver).move_to_element(element).perform()
11
12 # 单击, 弹出的Ajax元素, 根据链接节点的文本内容查找
13 driver.find_element_by_link_text('高级搜索').click()

```

今日作业

1 |

Day04回顾

- requests.get()参数

```
1  【1】url
2  【2】proxies -> {}
3      proxies = {
4          'http': 'http://1.1.1.1:8888',
5          'https': 'https://1.1.1.1:8888'
6      }
7  【3】timeout
8  【4】headers
```

- requests.post()

```
1  data : 字典, Form表单数据
```

- 常见的反爬机制及处理方式

```
1  【1】Headers反爬虫
2      1.1) 检查: Cookie、Referer、User-Agent
3      1.2) 解决方案: 通过F12获取headers,传给requests.get()方法
4
5  【2】IP限制
6      2.1) 网站根据IP地址访问频率进行反爬,短时间内限制IP访问
7      2.2) 解决方案:
8          a) 构造自己IP代理池,每次访问随机选择代理,经常更新代理池
9          b) 购买开放代理或私密代理IP
10         c) 降低爬取的速度
11
12  【3】User-Agent限制
13      3.1) 类似于IP限制, 检测频率
14      3.2) 解决方案: 构造自己的User-Agent池,每次访问随机选择
15          a> fake_useragent模块
16          b> 新建py文件,存放大量User-Agent
17          c> 程序中定义列表,存放大量的User-Agent
18
19  【4】对响应内容做处理
20      4.1) 页面结构和响应内容不同
21      4.2) 解决方案: 打印并查看响应内容,用xpath或正则做处理
22
23  【5】JS加密
24      5.1) 抓取到对应的JS文件,寻找加密算法
25      5.2) 用Python实现加密算法,生成指定的参数
```

- 有道翻译过程梳理

```
1  【1】打开首页
2
3  【2】准备抓包: F12开启控制台
4
```

```

5  【3】寻找地址
6    3.1) 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
7          F12-Network-XHR-Headers-General-Request URL
8
9  【4】发现规律
10   4.1) 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
11   4.2) 分析对比 Network - All(或者XHR) - Form Data，发现对应的规律
12
13 【5】寻找JS加密文件
14   5.1) 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{}
15
16 【6】查看JS代码
17   6.1) 搜索关键字，找到相关加密方法，用python实现加密算法
18
19 【7】断点调试
20   7.1) JS代码中部分参数不清楚可通过断点调试来分析查看
21
22 【8】完善程序

```

• Ajax动态加载数据抓取流程

```

1  【1】F12打开控制台，执行页面动作抓取网络数据包
2
3  【2】抓取json文件URL地址
4    2.1) 控制台中 XHR：找到异步加载的数据包
5    2.2) GET请求：Network -> XHR -> URL 和 Query String Parameters(查询参数)
6    2.3) POST请求：Network -> XHR -> URL 和 Form Data

```

• json模块

```

1  【1】抓取的json数据转为python数据类型
2    1.1) html = json.loads('[{},{},{}]')
3    1.2) html = requests.get(url=url,headers=headers).json()
4    1.3) html = requests.post(url=url,data=data,headers=headers).json()
5
6  【2】抓取数据保存到json文件
7    import json
8    with open('xxx.json','w') as f:
9        json.dump([{}},{},{}],f,ensure_ascii=False)

```

• 数据抓取最终梳理

```

1  【1】响应内容中存在
2    1.1) 确认抓取数据在响应内容中是否存在
3
4    1.2) 分析页面结构，观察URL地址规律
5      a) 大体查看响应内容结构，查看是否有更改 -- (百度视频案例)
6      b) 查看页面跳转时URL地址变化，查看是否新跳转 -- (民政部案例)
7
8    1.3) 开始写代码进行数据抓取
9
10 【2】响应内容中不存在
11   2.1) 确认抓取数据在响应内容中是否存在
12
13   2.2) F12抓包，开始刷新页面或执行某些行为，主要查看XHR异步加载数据包

```



```
14         a) GET请求: Request URL、Request Headers、Query String Paramters
15         b) POST请求: Request URL、Request Headers、FormData
16
17     2.3) 观察查询参数或者Form表单数据规律,如果需要进行进一步抓包分析处理
18         a) 比如有道翻译的 salt+sign,抓取并分析JS做进一步处理
19         b) 此处注意请求头中的Cookie和Referer以及User-Agent
20
21     2.4) 使用res.json()获取数据,利用列表或者字典的方法获取所需数据
```

• 多线程爬虫梳理

```
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()
```

Day05笔记

selenium+PhantomJS/Chrome/Firefox

- selenium

```
1  【1】定义
2      1.1) 开源的web自动化测试工具
3
4  【2】用途
5      2.1) 对web系统进行功能性测试,版本迭代时避免重复劳动
6      2.2) 兼容性测试(测试web程序在不同操作系统和不同浏览器中是否运行正常)
7      2.3) 对web系统进行大量测试
8
9  【3】特点
10     3.1) 可根据指令操控浏览器
11     3.2) 只是工具,必须与第三方浏览器结合使用
12
13  【4】安装
14     4.1) Linux: sudo pip3 install selenium
15     4.2) windows: python -m pip install selenium
```

- PhantomJS浏览器

```
1 phantomjs为无界面浏览器(又称无头浏览器),在内存中进行页面加载,高效
```

- 环境安装

```
1  【1】下载驱动
2      2.1) chromedriver : 下载对应版本
3          http://npm.taobao.org/mirrors/chromedriver/
4      2.2) geckodriver
5          https://github.com/mozilla/geckodriver/releases
6      2.3) phantomjs
7          https://phantomjs.org/download.html
8
9  【2】添加到系统环境变量
10     2.1) windows: 拷贝到Python安装目录的Scripts目录中
11             windows查看python安装目录(cmd命令行): where python
12     2.2) Linux : 拷贝到/usr/bin目录中 : sudo cp chromedriver /usr/bin/
13
14  【3】Linux中需要修改权限
15     sudo chmod 777 /usr/bin/chromedriver
16
17  【4】验证
18     4.1) ubuntu | windows
19         from selenium import webdriver
20         webdriver.Chrome()
21         webdriver.Firefox()
22
23     4.2) Mac
```

```

24     from selenium import webdriver
25     webdriver.Chrome(executable_path='/Users/xxx/chromedriver')
26     webdriver.Firefox(executable_path='/User/xxx/geckodriver')

```

selenium + PhantomJS
selenium + chromedriver + Chrome
selenium + geckodriver + Firefox

• 示例代码

```

1  """示例代码一：使用 selenium+浏览器 打开百度"""
2
3  # 导入selenium的webdriver接口
4  from selenium import webdriver
5  import time
6
7  # 创建浏览器对象
8  browser = webdriver.Chrome()
9  browser.get('http://www.baidu.com/')
10 # 5秒钟后关闭浏览器
11 time.sleep(5)
12 browser.quit()

```

```

1  """示例代码二：打开百度，搜索赵丽颖，点击搜索，查看"""
2
3  from selenium import webdriver
4  import time
5
6  # 1.创建浏览器对象 - 已经打开了浏览器
7  browser = webdriver.Chrome()
8  # 2.输入：http://www.baidu.com/
9  browser.get('http://www.baidu.com/')
10 # 3.找到搜索框,向这个节点发送文字：赵丽颖
11 browser.find_element_by_xpath('//*[@id="kw"]').send_keys('赵丽颖')
12 # 4.找到 百度一下 按钮,点击一下
13 browser.find_element_by_xpath('//*[@id="su"]').click()

```

• 浏览器对象(browser)方法

```

1  【1】 browser.get(url=url)    - 地址栏输入url地址并确认
2  【2】 browser.quit()          - 关闭浏览器
3  【3】 browser.close()         - 关闭当前页
4  【4】 browser.page_source     - HTML结构源码
5  【5】 browser.page_source.find('字符串')
6      # 从html源码中搜索指定字符串,没有找到返回：-1,经常用于判断是否为最后一页
7  【6】 browser.maximize_window() - 浏览器窗口最大化

```

• 定位节点八种方法

```

1  【1】 单元素查找('结果为1个节点对象')
2      1.1) 【最常用】 browser.find_element_by_id('id属性值')
3      1.2) 【最常用】 browser.find_element_by_name('name属性值')
4      1.3) 【最常用】 browser.find_element_by_class_name('class属性值')

```

```

5     1.4) 【最万能】browser.find_element_by_xpath('xpath表达式')
6     1.5) 【匹配a节点时常用】browser.find_element_by_link_text('链接文本')
7     1.6) 【匹配a节点时常用】browser.find_element_by_partical_link_text('部分链接文本')
8     1.7) 【最没用】browser.find_element_by_tag_name('标记名称')
9     1.8) 【较常用】browser.find_element_by_css_selector('css表达式')
10
11    【2】多元素查找('结果为[节点对象列表]')
12        2.1) browser.find_elements_by_id('id属性值')
13        2.2) browser.find_elements_by_name('name属性值')
14        2.3) browser.find_elements_by_class_name('class属性值')
15        2.4) browser.find_elements_by_xpath('xpath表达式')
16        2.5) browser.find_elements_by_link_text('链接文本')
17        2.6) browser.find_elements_by_partical_link_text('部分链接文本')
18        2.7) browser.find_elements_by_tag_name('标记名称')
19        2.8) browser.find_elements_by_css_selector('css表达式')
20
21    【3】注意!!!
22        当属性值中存在 空格 时,我们要使用 . 去代替空格
23        页面中class属性值为: btn btn-account
24        driver.find_element_by_class_name('btn.btn-account').click()
25
26    【4】如果找不到节点, 会抛出异常
27

```

• 猫眼电影示例

```

1  from selenium import webdriver
2  import time
3
4  url = 'https://maoyan.com/board/4'
5  browser = webdriver.Chrome()
6  browser.get(url)
7
8  def get_data():
9      # 基准xpath: [<selenium xxx li at xxx>,<selenium xxx li at>]
10     li_list = browser.find_elements_by_xpath('//*[
11     @id="app"]/div/div/div[1]/dl/dd')
12     for li in li_list:
13         item = {}
14         # info_list: ['1', '霸王别姬', '主演: 张国荣', '上映时间: 1993-01-
15         01', '9.5']
16         info_list = li.text.split('\n')
17         item['number'] = info_list[0]
18         item['name'] = info_list[1]
19         item['star'] = info_list[2]
20         item['time'] = info_list[3]
21         item['score'] = info_list[4]
22
23         print(item)
24
25  while True:
26     get_data()
27     try:
28         browser.find_element_by_link_text('下一页').click()
29         time.sleep(2)
30     except Exception as e:

```

```
29     print('恭喜你!抓取结束')
30     browser.quit()
31     break
```

- 节点对象操作

```
1  【1】文本框操作
2      1.1) node.send_keys('') - 向文本框发送内容
3      1.2) node.clear()      - 清空文本
4      1.3) node.get_attribute('value') - 获取文本内容
5
6  【2】按钮操作 只针对button按钮
7      1.1) node.click()      - 点击
8      1.2) node.is_enabled() - 判断按钮是否可用
9      1.3) node.get_attribute('value') - 获取按钮文本
```

chromedriver设置无头模式

```
1  from selenium import webdriver
2
3  options = webdriver.ChromeOptions()
4  # 添加无界面参数
5  options.add_argument('--headless')
6  browser = webdriver.Chrome(options=options)
```

selenium - 鼠标操作

```
1  from selenium import webdriver
2  # 导入鼠标事件类
3  from selenium.webdriver import ActionChains
4
5  driver = webdriver.Chrome()
6  driver.get('http://www.baidu.com/')
7
8  # 移动到 设置, perform()是真正执行操作, 必须有
9  element = driver.find_element_by_xpath('//*[@id="u1"]/a[8]')
10 ActionChains(driver).move_to_element(element).perform()
11
12 # 单击, 弹出的Ajax元素, 根据链接节点的文本内容查找
13 driver.find_element_by_link_text('高级搜索').click()
```

selenium - 切换页面

- 适用网站+应对方案

```

1  【1】适用网站类型
2      页面中点开链接出现新的窗口，但是浏览器对象browser还是之前页面的对象，需要切换到不同的窗口进行操作
3
4  【2】应对方案 - browser.switch_to.window()
5
6      # 获取当前所有句柄（窗口）- [handle1,handle2]
7      all_handles = browser.window_handles
8      # 切换browser到新的窗口，获取新窗口的对象
9      browser.switch_to.window(all_handles[1])

```

• 民政部网站案例-selenium

```

1  """
2  适用selenium+Chrome抓取民政部行政区划代码
3  http://www.mca.gov.cn/article/sj/xzqh/2019/
4  """
5  from selenium import webdriver
6
7  class GovSpider(object):
8      def __init__(self):
9          # 设置无界面
10         options = webdriver.ChromeOptions()
11         options.add_argument('--headless')
12         # 添加参数
13         self.browser = webdriver.Chrome(options=options)
14         self.one_url = 'http://www.mca.gov.cn/article/sj/xzqh/2019/'
15
16         def get_incr_url(self):
17             self.browser.get(self.one_url)
18             # 提取最新链接节点对象并点击
19
20         self.browser.find_element_by_xpath('//td[@class="arlisttd"]/a[contains(@title,"代码")]').click()
21         # 切换句柄
22         all_handlers = self.browser.window_handles
23         self.browser.switch_to.window(all_handlers[1])
24         self.get_data()
25
26         def get_data(self):
27             tr_list =
28             self.browser.find_elements_by_xpath('//tr[@height="19"]')
29             for tr in tr_list:
30                 code = tr.find_element_by_xpath('./td[2]').text.strip()
31                 name = tr.find_element_by_xpath('./td[3]').text.strip()
32                 print(name,code)
33
34         def run(self):
35             self.get_incr_url()
36             self.browser.quit()
37
38         if __name__ == '__main__':
39             spider = GovSpider()
40             spider.run()

```

selenium - iframe

• 特点+方法

```
1  【1】特点
2      网页中嵌套了网页，先切换到iframe，然后再执行其他操作
3
4  【2】处理步骤
5      2.1) 切换到要处理的Frame
6      2.2) 在Frame中定位页面元素并进行操作
7      2.3) 返回当前处理的Frame的上一级页面或主页面
8
9  【3】常用方法
10     3.1) 切换到frame - browser.switch_to.frame(frame节点对象)
11     3.2) 返回上一级 - browser.switch_to.parent_frame()
12     3.3) 返回主页面 - browser.switch_to.default_content()
13
14  【4】使用说明
15     4.1) 方法一：默认支持id和name属性值：switch_to.frame(id属性值|name属性
16         值)
17     4.2) 方法二：
18         a> 先找到frame节点：frame_node =
19         browser.find_element_by_xpath('xxx')
20         b> 在切换到frame：browser.switch_to.frame(frame_node)
```

• 示例1 - 登录豆瓣网

```
1  """
2  登录豆瓣网
3  """
4  from selenium import webdriver
5  import time
6
7  # 打开豆瓣官网
8  browser = webdriver.Chrome()
9  browser.get('https://www.douban.com/')
10
11 # 切换到iframe子页面
12 login_frame = browser.find_element_by_xpath('//*[id="anony-reg-
13 new"]/div/div[1]/iframe')
14 browser.switch_to.frame(login_frame)
15
16 # 密码登录 + 用户名 + 密码 + 登录豆瓣
17 browser.find_element_by_xpath('/html/body/div[1]/div[1]/ul[1]/li[2]').click()
18 browser.find_element_by_xpath('//*[id="username"]').send_keys('自己的用
19 户名')
20 browser.find_element_by_xpath('//*[id="password"]').send_keys('自己的密
21 码')
22 browser.find_element_by_xpath('/html/body/div[1]/div[2]/div[1]/div[5]/a')
23 .click()
24 time.sleep(3)
25
26 # 点击我的豆瓣
27 browser.find_element_by_xpath('//*[id="db-nav-
28 sns"]/div/div/div[3]/ul/li[2]/a').click()
```

- selenium+phantomjs|chrome|firefox小总结

```
1  【1】 设置无界面模式
2      options = webdriver.ChromeOptions()
3      options.add_argument('--headless')
4      browser =
webdriver.Chrome(executable_path='/home/tarena/chromedriver',options=options)
5
6  【2】 browser执行JS脚本
7
8      browser.execute_script('window.scrollTo(0,document.body.scrollHeight)')
9
10 【3】 键盘操作
11
12     from selenium.webdriver.common.keys import Keys
13
14 【4】 鼠标操作
15
16     from selenium.webdriver import ActionChains
17     ActionChains(browser).move_to_element('node').perform()
18
19 【5】 切换句柄 - switch_to.frame(handle)
20     all_handles = browser.window_handles
21     browser.switch_to.window(all_handles[1])
22     # 开始进行数据抓取
23     browser.close()
24     browser.switch_to.window(all_handles[0])
25
26 【6】 iframe子页面
27     browser.switch_to.frame(frame_node)
```

- lxml中的xpath 和 selenium中的xpath的区别

```
1  【1】 lxml中的xpath用法 - 推荐自己手写
2      div_list = p.xpath('//div[@class="abc"]/div')
3      item = {}
4      for div in div_list:
5          item['name'] = div.xpath('./a/@href')[0]
6          item['likes'] = div.xpath('./a/text()')[0]
7
8  【2】 selenium中的xpath用法 - 推荐copy - copy xpath
9      div_list = browser.find_elements_by_xpath('//div[@class="abc"]/div')
10     item = {}
11     for div in div_list:
12         item['name'] =
div.find_element_by_xpath('./a').get_attribute('href')
13         item['likes'] = div.find_element_by_xpath('./a').text
```

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 Visual C++ 14.0 is required xxx'
18             则安装Windows下的Microsoft Visual 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/langfang/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
```

今日作业

```
1     【1】使用selenium+浏览器 获取有道翻译结果
2     【2】使用selenium+浏览器 登录网易qq邮箱 : https://mail.qq.com/
3     【3】使用selenium+浏览器 登录网易163邮箱 : https://mail.163.com/
4     【4】熟记scrapy的五大组件,以及工作流程,能够描述的很清楚
```

Day05回顾

selenium+phantomjs/chrome/firefox

- 特点

```
1     【1】简单, 无需去详细抓取分析网络数据包, 使用真实浏览器
2     【2】需要等待页面元素加载, 需要时间, 效率低
```

- 设置无界面模式 (chromedriver | firefox)

```
1     options = webdriver.ChromeOptions()
2     options.add_argument('--headless')
3
4     browser = webdriver.Chrome(options=options)
5     browser.get(url)
```

- browser执行JS脚本

```

1 browser.execute_script(
2     'window.scrollTo(0,document.body.scrollHeight)'
3 )
4 # 把下拉条从上拉到下
5 time.sleep(1)

```

- selenium常用操作

```

1 【1】键盘操作
2     from selenium.webdriver.common.keys import Keys
3     node.send_keys(Keys.SPACE)
4     node.send_keys(Keys.CONTROL, 'a')
5     node.send_keys(Keys.CONTROL, 'c')
6     node.send_keys(Keys.CONTROL, 'v')
7     node.send_keys(Keys.ENTER)
8
9 【2】鼠标操作
10    from selenium.webdriver import ActionChains
11    ActionChains(browser).move_to_element(node).perform()
12
13 【3】切换句柄
14    all_handles = browser.window_handles
15    time.sleep(1)
16    browser.switch_to.window(all_handles[1])
17
18 【4】iframe子框架
19    browser.switch_to.frame(iframe_element)
20    # 写法1 - 任何场景都可以:
21    iframe_node = browser.find_element_by_xpath('')
22    browser.switch_to.frame(iframe_node)
23
24    # 写法2 - 默认支持 id 和 name 两个属性值:
25    browser.switch_to.frame('id属性值|name属性值')

```

Day06笔记

作业概解

作业1 - 有道翻译实现

- 代码实现

```

1 """
2 selenium实现抓取有道翻译结果
3 思路:
4     1、找到输入翻译单词节点,发送文字
5     2、休眠一定时间,等待网站给出响应-翻译结果
6     3、找到翻译结果节点,获取文本内容
7 """
8 from selenium import webdriver
9 import time
10
11 class YdSpider:
12     def __init__(self):
13         self.url = 'http://fanyi.youdao.com/'

```

```

14         # 设置无界面模式
15         self.options = webdriver.ChromeOptions()
16         self.options.add_argument('--headless')
17         self.driver = webdriver.Chrome(options=self.options)
18         # 打开有道翻译官网
19         self.driver.get(self.url)
20
21     def parse_html(self, word):
22         # 发送翻译单词
23         self.driver.find_element_by_id('inputOriginal').send_keys(word)
24         time.sleep(1)
25         # 获取翻译结果
26         result = self.driver.find_element_by_xpath('//*[
[@id="transTarget"]/p/span').text
27
28         return result
29
30     def run(self):
31         word = input('请输入要翻译的单词:')
32         print(self.parse_html(word))
33         self.driver.quit()
34
35 if __name__ == '__main__':
36     spider = YdSpider()
37     spider.run()

```

作业2 - 登录QQ邮箱

- 代码实现

```

1  from selenium import webdriver
2  import time
3
4  driver = webdriver.Chrome()
5  driver.get('https://mail.qq.com/')
6
7  # 切换到iframe子框架
8  driver.switch_to.frame("login_frame")
9
10 # 用户名+密码+登录
11 driver.find_element_by_id('u').send_keys('2621470058')
12 driver.find_element_by_id('p').send_keys('zhanshen001')
13 driver.find_element_by_id('login_button').click()

```

作业3- 163邮箱登陆

- 代码实现

```

1  """
2  selenium模拟登录163邮箱
3  思路:
4      1、密码登录在这里 - 此节点在主页面中,并非iframe内部
5      2、切换iframe - 此处iframe节点中id的值每次都在变化,需要手写xpath,否则会出现
        无法定位iframe
6      3、输入用户名和密码
7      4、点击登录按钮

```

```

8  """
9  from selenium import webdriver
10
11 driver = webdriver.Chrome()
12 driver.get('https://mail.163.com/')
13
14 # 1、切换iframe子页面 - 此处手写xpath,此处iframe中id的值每次都在变化
15 node = driver.find_element_by_xpath('//div[@id="loginDiv"]/iframe[1]')
16 driver.switch_to.frame(node)
17
18 # 2、输入用户名和密码
19 driver.find_element_by_name('email').send_keys('wangweichao_2020')
20 driver.find_element_by_name('password').send_keys('zhanshen001')
21 driver.find_element_by_id('dologin').click()

```

作业4 - 京东爬虫

• 目标

- 1 【1】 目标网址 : <https://www.jd.com/>
- 2 【2】 抓取目标 : 商品名称、商品价格、评价数量、商品商家

• 思路提醒

- 1 【1】 打开京东, 到商品搜索页
- 2 【2】 匹配所有商品节点对象列表
- 3 【3】 把节点对象的文本内容取出来, 查看规律, 是否有更好的处理方法?
- 4 【4】 提取完1页后, 判断如果不是最后1页, 则点击下一页
- 5 '问题: 如何判断是否为最后1页???'

• 实现步骤

```

1  【1】 找节点
2      1.1) 首页搜索框 : //*[@id="key"]
3      2.1) 首页搜索按钮 : //*[@id="search"]/div/div[2]/button
4      2.3) 商品页的 商品信息节点对象列表 : //*[@id="J_goodsList"]/ul/li
5      2.4) for循环遍历后
6          a> 名称: .//div[@class="p-name"]/a/em
7          b> 价格: .//div[@class="p-price"]
8          c> 评论: .//div[@class="p-commit"]/strong
9          d> 商家: .//div[@class="p-shopnum"]
10
11  【2】 执行JS脚本, 获取动态加载数据
12      browser.execute_script(
13          'window.scrollTo(0,document.body.scrollHeight)'
14      )

```

• 代码实现

```

1  from selenium import webdriver
2  import time
3
4  class JdSpider(object):
5      def __init__(self):
6          self.url = 'https://www.jd.com/'

```

```

7         # 设置无界面模式
8         self.options = webdriver.ChromeOptions()
9         self.options.add_argument('--headless')
10        self.browser = webdriver.Chrome(options=self.options)
11
12        def get_html(self):
13            # get():等页面所有元素加载完成后,才会执行后面的代码
14            self.browser.get(self.url)
15            # 搜索框 + 搜索按钮
16            self.browser.find_element_by_xpath('//*[@id="key"]').send_keys('爬虫书')
17            self.browser.find_element_by_xpath('//*[@id="search"]/div/div[2]/button').click()
18
19            # 循环体中的函数: 拉进度条,提取数据
20            def parse_html(self):
21                # 执行js脚本,将进度条拉到最底部
22                self.browser.execute_script(
23                    'window.scrollTo(0,document.body.scrollHeight)'
24                )
25                # 给页面元素加载预留时间
26                time.sleep(3)
27                li_list = self.browser.find_elements_by_xpath('//*[@id="J_goodsList"]/ul/li')
28
29                for li in li_list:
30                    item = {}
31                    item['price'] = li.find_element_by_xpath('.//div[@class="p-price"]').text
32                    item['name'] = li.find_element_by_xpath('.//div[@class="p-name"]/a/em').text
33                    item['commit'] = li.find_element_by_xpath('.//div[@class="p-commit"]/strong').text
34                    item['shop'] = li.find_element_by_xpath('.//div[@class="p-shopnum"]').text
35                    print(item)
36
37            def run(self):
38                self.get_html()
39                while True:
40                    self.parse_html()
41                    if self.browser.page_source.find('pn-next disabled') == -1:
42                        self.browser.find_element_by_xpath('//*[@id="J_bottomPage"]/span[1]/a[9]').click()
43                    else:
44                        self.browser.quit()
45                        break
46
47        if __name__ == '__main__':
48            spider = JdSpider()
49            spider.run()

```

scrapy框架

- Scrapy框架五大组件

```

1  【1】引擎（Engine）-----整个框架核心
2  【2】爬虫程序（Spider）-----数据解析提取
3  【3】调度器（Scheduler）-----维护请求队列
4  【4】下载器（Downloader）-----获取响应对象
5  【5】管道文件（Pipeline）-----数据入库处理
6
7
8  【两个中间件】
9      下载器中间件（Downloader Middlewares）
10         引擎->下载器,包装请求(随机代理等)
11      蜘蛛中间件（Spider Middlewares）
12         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 : scrapy startproject 项目名
2  【2】创建爬虫文件
3      2.1) cd 项目文件夹
4      2.2) scrapy genspider 爬虫名 域名
5  【3】运行爬虫
6      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  【2】ROBOTSTXT_OBEY = False
3      是否遵循robots协议, 一般我们一定要设置为False
4  【3】CONCURRENT_REQUESTS = 32
5      最大并发量, 默认为16
6  【4】DOWNLOAD_DELAY = 0.5
7      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
8  【5】COOKIES_ENABLED = False | True
9      Cookie默认是禁用的, 取消注释则 启用Cookie, 即: True和False都是启用Cookie
10 【6】DEFAULT_REQUEST_HEADERS = {}
11     请求头, 相当于requests.get(headers=headers)

```

• 创建爬虫项目步骤

```

1  【1】新建项目和爬虫文件
2      scrapy startproject 项目名
3      cd 项目文件夹
4      新建爬虫文件: scrapy genspider 文件名 域名
5  【2】明确目标(items.py)
6  【3】写爬虫程序(文件名.py)
7  【4】管道文件(pipelines.py)
8  【5】全局配置(settings.py)
9  【6】运行爬虫
10     8.1) 终端: scrapy crawl 爬虫名
11     8.2) pycharm运行
12         a> 创建run.py(和scrapy.cfg文件同目录)
13             from scrapy import cmdline
14             cmdline.execute('scrapy crawl maoyan'.split())
15         b> 直接运行 run.py 即可

```

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
16
17
18 class GuaziSpider(scrapy.Spider):
19     # 爬虫名
20     name = 'car'
21     # 允许爬取的域名
22     allowed_domains = ['www.guazi.com']
23     # 初始的URL地址
24     start_urls = ['https://www.guazi.com/bj/buy/o1/#bread']
25     # 生成URL地址的变量
26     n = 1
27
28     def parse(self, response):
29         # 基准xpath: 匹配所有汽车的节点对象列表
30         li_list = response.xpath('//ul[@class="carlist clearfix js-
top"]/li')
31         # 给items.py中的 GuaziItem类 实例化
32         item = CarItem()
33         for li in li_list:
34             item['url'] = li.xpath('./a[1]/@href').get()
35             item['name'] = li.xpath('./a[1]/@title').get()
36             item['price'] = li.xpath('./div[@class="t-
price"]/p/text()').get()
37
38             # 把抓取的数据,传递给了管道文件 pipelines.py
39             yield item

```

```

40
41     # 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
42     if self.n < 5:
43         self.n += 1
44         url =
45         'https://www.guazi.com/bj/buy/o{}/#bread'.format(self.n)
46         # 把url交给调度器入队列
47         yield scrapy.Request(url=url, callback=self.parse)

```

• 步骤3 - 编写爬虫文件 (代码实现2)

```

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'] = 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
32             # 把抓取的数据,传递给了管道文件 pipelines.py
33             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 =
pymysql.connect(MYSQL_HOST,MYSQL_USER,MYSQL_PWD,MYSQL_DB,charset=CHARSET
)
28         self.cursor = self.db.cursor()
29
30     def process_item(self, item, spider):
31         """处理从爬虫文件传过来的item数据"""
32         ins = 'insert into guazitab values(%s,%s,%s)'
33         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()
43
44
45 # 管道3 - 存入MongoDB管道
46 import pymongo
47
48 class CarMongoPipeline(object):
49     def open_spider(self, spider):
50         self.conn = pymongo.MongoClient(MONGO_HOST,MONGO_PORT)
51         self.db = self.conn[MONGO_DB]
52         self.myset = self.db[MONGO_SET]
53
54     def process_item(self, item, spider):
55         car_dict = {
56             'name' : item['name'],
57             'price': item['price'],
58             'url' : item['url']
59         }
60         self.myset.insert_one(car_dict)

```

- 步骤5 - 全局配置文件 (settings.py)

```

1  【1】ROBOTSTXT_OBEY = False
2  【2】DOWNLOAD_DELAY = 1
3  【3】COOKIES_ENABLED = False
4  【4】DEFAULT_REQUEST_HEADERS = {
5      "Cookie": "此处填写抓包抓取到的Cookie",
6      "User-Agent": "此处填写自己的User-Agent",
7  }
8
9  # 优先级, 1--1000, 数字越小越高
10 【5】ITEM_PIPELINES = {
11     'Car.pipelines.CarPipeline': 300,
12     'Car.pipelines.CarMySQLPipeline': 400,
13     'Car.pipelines.CarMongoPipeline': 500,
14 }
15
16 【6】定义MySQL相关变量
17 MYSQL_HOST = 'localhost'
18 MYSQL_USER = 'root'
19 MYSQL_PWD = '123456'
20 MYSQL_DB = 'guazidb'
21 CHARSET = 'utf8'
22
23 【7】定义MongoDB相关变量
24 MONGO_HOST = 'localhost'
25 MONGO_PORT = 27017
26 MONGO_DB = 'guazidb'
27 MONGO_SET = 'guaziset'

```

- 步骤6 - 运行爬虫 (run.py)

```

1  """run.py"""
2  from scrapy import cmdline
3  cmdline.execute('scrapy crawl car'.split())

```

知识点汇总

- 数据持久化 - 数据库

```

1  【1】在setting.py中定义相关变量
2  【2】pipelines.py中导入settings模块
3      def open_spider(self, spider):
4          """爬虫开始执行1次,用于数据库连接"""
5
6      def process_item(self, item, spider):
7          """具体处理数据"""
8          return item
9
10     def close_spider(self, spider):
11         """爬虫结束时执行1次,用于断开数据库连接"""
12  【3】settings.py中添加此管道
13     ITEM_PIPELINES = {'':200}
14
15  【注意】: process_item() 函数中一定要 return item ,当前管道的process_item()
    的返回值会作为下一个管道 process_item()的参数

```

- 数据持久化 - csv、json文件

```
1  【1】 存入csv文件
2      scrapy crawl car -o car.csv
3
4  【2】 存入json文件
5      scrapy crawl car -o car.json
6
7  【3】 注意： settings.py中设置导出编码 - 主要针对json文件
8      FEED_EXPORT_ENCODING = 'utf-8'
```

- 节点对象.xpath("")

```
1  【1】 列表,元素为选择器 @
2      [
3          <selector xpath='xxx' data='A'>,
4          <selector xpath='xxx' data='B'>
5      ]
6  【2】 列表.extract() : 序列化列表中所有选择器为Unicode字符串 ['A','B']
7  【3】 列表.extract_first() 或者 get() : 获取列表中第1个序列化的元素(字符串) 'A'
```

- 课堂练习

```
1  【熟悉整个流程】 : 将猫眼电影案例数据抓取, 存入MySQL数据库
```

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

```

• 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: '//*[@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', ''])

```


项目实现

- 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 }

```

今日作业

- 1 **【1】** 腾讯招聘职位信息抓取（二级页面）
- 2 要求：输入职位关键字，抓取该类别下所有职位信息（到职位详情页抓取）
- 3 具体数据如下：
- 4 1.1) 职位名称
- 5 1.2) 职位地点
- 6 1.3) 职位类别
- 7 1.4) 发布时间
- 8 1.5) 工作职责
- 9 1.6) 工作要求

Day06回顾

scrapy框架

- 五大组件+工作流程+常用命令

- 1 **【1】** 五大组件
- 2 1.1) 引擎（Engine）
- 3 1.2) 爬虫程序（Spider）
- 4 1.3) 调度器（Scheduler）
- 5 1.4) 下载器（Downloader）
- 6 1.5) 管道文件（Pipeline）
- 7 1.6) 下载器中间件（Downloader Middlewares）
- 8 1.7) 蜘蛛中间件（Spider Middlewares）
- 9
- 10 **【2】** 工作流程
- 11 2.1) Engine向Spider索要URL,交给Scheduler入队列
- 12 2.2) Scheduler处理后出队列,通过Downloader Middlewares交给Downloader去下
- 13 载
- 14 2.3) Downloader得到响应后,通过Spider Middlewares交给Spider
- 15 2.4) Spider数据提取:
- 16 a) 数据交给Pipeline处理
- 17 b) 需要跟进URL,继续交给Scheduler入队列,依次循环
- 18
- 19 **【3】** 常用命令
- 20 3.1) scrapy startproject 项目名
- 21 3.2) scrapy genspider 爬虫名 域名
- 22 3.3) scrapy crawl 爬虫名

完成scrapy项目完整流程

- 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'
```

Day07笔记

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 """

```

机器视觉与tesseract

• 概述

```

1  【1】作用
2      处理图形验证码
3
4  【2】三个重要概念 - OCR、tesseract-ocr、pytesseract
5      2.1) OCR
6          光学字符识别(Optical Character Recognition),通过扫描等光学输入方式将各
          种票据、报刊、书籍、文稿及其它印刷品的文字转化为图像信息，再利用文字识别技术将图像信息
          转化为电子文本
7
8      2.2) tesseract-ocr
9          OCR的一个底层识别库（不是模块，不能导入），由Google维护的开源OCR识别库
10
11     2.3) pytesseract
12         Python模块,可调用底层识别库，是对tesseract-ocr做的一层Python API封装

```

• 安装tesseract-ocr

```

1  【1】Ubuntu安装
2      sudo apt-get install tesseract-ocr
3
4  【2】Windows安装
5      2.1) 下载安装包
6      2.2) 添加到环境变量(Path)
7
8  【3】测试（终端 | cmd命令行）
9      tesseract xxx.jpg 文件名

```

• 安装pytesseract

```

1  【1】安装
2      sudo pip3 install pytesseract
3
4  【2】使用示例
5      import pytesseract
6      # Python图片处理库
7      from PIL import Image
8
9      # 创建图片对象
10     img = Image.open('test1.jpg')
11     # 图片转字符串
12     result = pytesseract.image_to_string(img)
13     print(result)

```

补充 - 滑块缺口验证码案例

豆瓣网登录

• 案例说明

- 1 【1】URL地址：<https://www.douban.com/>
- 2 【2】先输入几次错误的密码，让登录出现滑块缺口验证，以便于我们破解
- 3 【3】模拟人的行为
- 4 3.1) 先快速滑动
- 5 3.2) 到离重点位置不远的地方开始减速
- 6 【4】详细看代码注释

• 代码实现

```
1  """
2  说明：先输入几次错误的密码，出现滑块缺口验证码
3  """
4  from selenium import webdriver
5  # 导入鼠标事件类
6  from selenium.webdriver import ActionChains
7  import time
8
9  # 加速度函数
10 def get_tracks(distance):
11     """
12     拿到移动轨迹，模仿人的滑动行为，先匀加速后匀减速
13     匀变速运动基本公式：
14     ① $v=v_0+at$ 
15     ② $s=v_0t+\frac{1}{2}at^2$ 
16     """
17     # 初速度
18     v = 0
19     # 单位时间为0.3s来统计轨迹，轨迹即0.3s内的位移
20     t = 0.3
21     # 位置/轨迹列表，列表内的一个元素代表0.3s的位移
22     tracks = []
23     # 当前的位移
24     current = 0
25     # 到达mid值开始减速
26     mid = distance*4/5
27     while current < distance:
28         if current < mid:
29             # 加速度越小,单位时间内的位移越小,模拟的轨迹就越多越详细
30             a = 2
31         else:
32             a = -3
33
34         # 初速度
35         v0 = v
36         # 0.3秒内的位移
37         s = v0*t+0.5*a*(t**2)
38         # 当前的位置
39         current += s
40         # 添加到轨迹列表
41         tracks.append(round(s))
42         # 速度已经达到v, 该速度作为下次的初速度
43         v = v0 + a*t
```

```

44     return tracks
45     # tracks: [第一个0.3秒的移动距离, 第二个0.3秒的移动距离, ...]
46
47
48 # 1、打开豆瓣官网 - 并将窗口最大化
49 browser = webdriver.Chrome()
50 browser.maximize_window()
51 browser.get('https://www.douban.com/')
52
53 # 2、切换到iframe子页面
54 login_frame = browser.find_element_by_xpath('//*[@id="anony-reg-
55 new"]/div/div[1]/iframe')
56 browser.switch_to.frame(login_frame)
57
58 # 3、密码登录 + 用户名 + 密码 + 登录豆瓣
59 browser.find_element_by_xpath('/html/body/div[1]/div[1]/ul[1]/li[2]').click()
60 browser.find_element_by_xpath('//*[@id="username"]').send_keys('15110225726')
61 browser.find_element_by_xpath('//*[@id="password"]').send_keys('zhanshen001')
62 browser.find_element_by_xpath('/html/body/div[1]/div[2]/div[1]/div[5]/a').click()
63 time.sleep(4)
64
65 # 4、切换到新的iframe子页面 - 滑块验证
66 auth_frame = browser.find_element_by_xpath('//*[@id="Tcaptcha"]/iframe')
67 browser.switch_to.frame(auth_frame)
68
69 # 5、按住开始滑动位置按钮 - 先移动180个像素
70 element = browser.find_element_by_xpath('//*[@id="tcaptcha_drag_button"')
71 # click_and_hold(): 按住某个节点并保持
72 ActionChains(browser).click_and_hold(on_element=element).perform()
73 # move_to_element_with_offset(): 移动到距离某个元素(左上角坐标)多少距离的位置
74 ActionChains(browser).move_to_element_with_offset(to_element=element, xoffset=180, yoffset=0).perform()
75
76 # 6、使用加速度函数移动剩下的距离
77 tracks = get_tracks(28)
78 for track in tracks:
79     # move_by_offset(): 鼠标从当前位置移动到某个坐标
80
81     ActionChains(browser).move_by_offset(xoffset=track, yoffset=0).perform()
82
83 # 7、延迟释放鼠标: release()
84 time.sleep(0.5)
85 ActionChains(browser).release().perform()

```

Fiddler抓包工具

- 配置Fiddler

```

1  【1】Tools -> Options -> HTTPS
2      1.1) 添加证书信任： 勾选 Decrypt Https Traffic 后弹出窗口，一路确认
3      1.2) 设置之抓浏览器的包： ...from browsers only
4
5  【2】Tools -> Options -> Connections
6      2.1) 设置监听端口（默认为8888）
7
8  【3】配置完成后重启Fiddler（'重要'）
9      3.1) 关闭Fiddler,再打开Fiddler

```

• 配置浏览器代理

```

1  【1】安装Proxy SwitchyOmega谷歌浏览器插件
2
3  【2】配置代理
4      2.1) 点击浏览器右上角插件SwitchyOmega -> 选项 -> 新建情景模式 ->
myproxy(名字) -> 创建
5      2.2) 输入 HTTP:// 127.0.0.1 8888
6      2.3) 点击：应用选项
7
8  【3】点击右上角SwitchyOmega可切换代理
9
10 【注意】：一旦切换了自己创建的代理,则必须要打开Fiddler才可以上网

```

• Fiddler常用菜单

```

1  【1】Inspector：查看数据包详细内容
2      1.1) 整体分为请求和响应两部分
3
4  【2】Inspector常用菜单
5      2.1) Headers：请求头信息
6      2.2) webForms：POST请求Form表单数据：<body>
7              GET请求查询参数：<QueryString>
8      2.3) Raw：将整个请求显示为纯文本

```

移动端app数据抓取

• 方法1 - 手机 + Fiddler

```

1  设置方法见文件夹 - 移动端抓包配置

```

• 方法2 - F12浏览器工具

有道翻译手机版破解案例

```

1  import requests
2  from lxml import etree
3
4  word = input('请输入要翻译的单词:')
5
6  post_url = 'http://m.youdao.com/translate'
7  post_data = {
8      'inputtext':word,
9      'type':'AUTO'

```

```
10 }  
11  
12 html = requests.post(url=post_url,data=post_data).text  
13 parse_html = etree.HTML(html)  
14 xpath_bds = '//ul[@id="translateResult"]/li/text()'   
15 result = parse_html.xpath(xpath_bds)[0]  
16  
17 print(result)
```