

Day07回顾

selenium+phantomjs/chrome/firefox

■ 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 time.sleep(2)
```

■ 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    mouse_action = ActionChains(browser)
12    mouse_action.move_to_element(node)
13    mouse_action.perform()
14
15 【3】 切换句柄
16    all_handles = browser.window_handles
17    time.sleep(1)
18    browser.switch_to.window(all_handles[1])
19
20 【4】 iframe子框架
21    browser.switch_to.frame(iframe_element)
```

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) ----- 维护请求队列
3 【3】调度器 (Scheduler) ----- 获取响应对象
4 【4】下载器 (Downloader) ----- 数据解析提取
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 蜘蛛中间件 (Spider Middlewares)
```

## ■ 工作流程

- ```
1  【1】 Engine向Spider索要URL,交给Scheduler入队列
2  【2】 Scheduler处理后出队列,通过Downloader Middlewares交给Downloader去下载
3  【3】 Downloader得到响应后,通过Spider Middlewares交给Spider
4  【4】 Spider数据提取:
5      4.1) 数据交给Pipeline处理
6      4.2) 需要跟进URL,继续交给Scheduler入队列, 依次循环
```

■ 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
```

Day08笔记

scrapy框架

■ 全局配置文件settings.py详解

```
1  【1】定义User-Agent
2      USER_AGENT = 'Mozilla/5.0'
3
4  【2】是否遵循robots协议，一般设置为False
5      ROBOTSTXT_OBEY = False
6
7  【3】最大并发量，默认为16
8      CONCURRENT_REQUESTS = 32
9
10 【4】下载延迟时间
11     DOWNLOAD_DELAY = 1
12
13 【5】请求头，此处也可以添加User-Agent
14     DEFAULT_REQUEST_HEADERS={}
```

■ 创建爬虫项目步骤

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

还记得我们抓取的 百度一下,你就知道 吗

■ 步骤跟踪

```
1  【1】创建项目 'Baidu' 和爬虫文件 'baidu'
2      1.1) scrapy startproject Baidu
3      1.2) cd Baidu
4      1.3) scrapy genspider baidu www.baidu.com
5
6  【2】打开爬虫文件: 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         def parse(self, response):
13             r_list = response.xpath('/html/head/title/text()').extract()[0]
14
15     【3】全局配置文件: settings.py
16         ROBOTSTXT_OBEY = False
17         DEFAULT_REQUEST_HEADERS = {'User-Agent': 'Mozilla/5.0'}
18
19     【4】创建文件(和项目目录同路径): run.py
20         from scrapy import cmdline
21         cmdline.execute('scrapy crawl baidu'.split())
22
23     【5】运行 run.py 启动爬虫

```

瓜子二手车直卖网

■ 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

```

实现步骤

■ 1-创建项目和爬虫文件

```

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

```

■ 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()

```

■ 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 class CarSpider(scrapy.Spider):
18     name = 'car'
19     allowed_domains = ['www.guazi.com']
20     i = 1
21     start_urls = ['https://www.guazi.com/langfang/buy/o1/']
22
23     def parse(self, response):
24         # 1.基准xpath,匹配所有汽车节点对象列表
25         li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
26         # 创建item对象,给items.py中定义的数据结构赋值
27         item = CarItem()
28         for li in li_list:
29             item['url'] = 'https://www.guazi.com/' + li.xpath('./a[1]/@href').get()
30             item['name'] = li.xpath('./h2[@class="t"]/text()').get()
31             item['price'] = li.xpath('./div[@class="t-price"]/p/text()').get()
32
33             yield item
34
35         # 生成下一页的链接,继续交给调度器入队列
36         if self.i < 5:
37             self.i += 1
38             url = 'https://www.guazi.com/langfang/buy/o{}/'.format(self.i)
39             # scrapy.Request()是将请求交给调度器入队列的方法
40             yield scrapy.Request(url=url, callback=self.parse)

```

■ 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 CarSpider(scrapy.Spider):
9      name = 'car2'
10     allowed_domains = ['www.guazi.com']
11     # 1、去掉 start_urls
12     # 2、重写start_requests()方法
13     def start_requests(self):

```

```

14     """生成所有待爬取的URL地址,统一交给调度器入队列"""
15     for i in range(1,5):
16         url = 'https://www.guazi.com/langfang/buy/o{}/'.format(i)
17         yield scrapy.Request(url=url,callback=self.parse)
18
19     def parse(self, response):
20         # 1.基准xpath,匹配所有汽车节点对象列表
21         li_list = response.xpath('//ul[@class="carlist clearfix js-top"]/li')
22         # 创建item对象,给items.py中定义的数据结构赋值
23         item = CarItem()
24         for li in li_list:
25             item['url'] = 'https://www.guazi.com/' + li.xpath('./a[1]/@href').get()
26             item['name'] = li.xpath('./h2[@class="t"]/text()').get()
27             item['price'] = li.xpath('./div[@class="t-price"]/p/text()').get()
28
29             yield item

```

■ 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 guazidb charset utf8;
5     use guazidb;
6     create table guaziset(
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(item['name'],item['price'],item['url'])
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)

```

■ 5-全局配置文件 (settings.py)

```

1  【1】ROBOTSTXT_OBEY = False
2  【2】DOWNLOAD_DELAY = 2
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": "Mozilla/5.0 (Windows NT 10.0; WOW64) AppleWebKit/537.36 (KHTML, like
7      Gecko) Chrome/80.0.3987.122 Safari/537.36",
8  }
9  【5】ITEM_PIPELINES = {
10     'Car.pipelines.CarPipeline': 300,
11     #'Car.pipelines.CarMysqlPipeline': 400,
12     #'Car.pipelines.CarMongoPipeline': 500,
13 }
14
15 【6】定义MySQL相关变量
16 MYSQL_HOST = 'localhost'
17 MYSQL_USER = 'root'
18 MYSQL_PWD = '123456'
19 MYSQL_DB = 'guazidb'
20 CHARSET = 'utf8'
21
22 【7】定义MongoDB相关变量
23 MONGO_HOST = 'localhost'
24 MONGO_PORT = 27017
25 MONGO_DB = 'guazidb'
26 MONGO_SET = 'guaziset'

```

■ 6-运行爬虫 (run.py)

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

数据持久化(MySQL)

■ 实现步骤

```
1 【1】 在setting.py中定义相关变量
2
3 【2】 pipelines.py中导入settings模块
4     def open_spider(self, spider):
5         """爬虫开始执行1次,用于数据库连接"""
6
7     def process_item(self, item, spider):
8         """具体处理数据"""
9         return item
10
11     def close_spider(self, spider):
12         """爬虫结束时执行1次,用于断开数据库连接"""
13
14 【3】 settings.py中添加此管道
15     ITEM_PIPELINES = {'':200}
16
17 【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
```

知识点汇总

■ 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
7 【2】 列表.extract() : 序列化列表中所有选择器为Unicode字符串 ['A', 'B']
8
9 【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'
```

■ 日志变量及日志级别(settings.py)


```

1 # 日志相关变量
2 LOG_LEVEL = ''
3 LOG_FILE = '文件名.log'
4
5 # 日志级别
6 5 CRITICAL : 严重错误
7 4 ERROR    : 普通错误
8 3 WARNING  : 警告
9 2 INFO     : 一般信息
10 1 DEBUG    : 调试信息
11 # 注意: 只显示当前级别的日志和比当前级别日志更严重的

```

■ 管道文件使用

```

1 【1】在爬虫文件中为items.py中类做实例化，用爬下来的数据给对象赋值
2     from ..items import MaoyanItem
3     item = MaoyanItem()
4
5 【2】管道文件 (pipelines.py)
6
7 【3】开启管道 (settings.py)
8     ITEM_PIPELINES = { '项目目录名.pipelines.类名':优先级 }

```

保存为csv、json文件

■ 命令格式

```

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

```

■ 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

```

新浪新闻全站抓取

■ 目标

```
1 【1】 抓取新浪新闻下的所有分类的所有新闻，保存到本地
2 【2】 URL： 新浪官网 - 更多 - 导航
3      http://news.sina.com.cn/guide/
4 【3】 要求
5      将信息保存到scrapy项目目录的 data 文件夹中,并按照分类名称创建子文件夹
```

实现步骤

■ 1-创建项目和爬虫文件

```
1 scrapy startproject Sina
2 cd Sina
3 scrapy genspider sina news.sina.com.cn
```

■ 2-定义要抓取的数据结构

```
1 """items.py"""
2 import scrapy
3
4 class SinaItem(scrapy.Item):
5     # 大类标题、url 例：新闻、体育、娱乐、财经... ..
6     parent_name = scrapy.Field()
7     parent_url = scrapy.Field()
8
9     # 小类标题、url 例：体育分类下的 NBA CBA ... ..
10    son_name = scrapy.Field()
11    son_url = scrapy.Field()
12
13    # 小类目录存储路径
14    son_filename = scrapy.Field()
15
16    # 小类下的文章链接、标题、内容
17    article_url = scrapy.Field()
18    article_head = scrapy.Field()
19    article_content = scrapy.Field()
```

■ 3-爬虫文件进行数据解析提取

```
1 """sina.py"""
2 # -*- coding: utf-8 -*-
3 import scrapy
4 import os
5 from ..items import SinaItem
6
7 class SinaSpider(scrapy.Spider):
8     name = 'sina'
9     allowed_domains = ['sina.com.cn']
10    # 起始URL地址为导航页的地址
11    start_urls = ['http://news.sina.com.cn/guide/']
12
13    def parse(self, response):
```

```

14     # 用来存放下一次交给调度器的所有请求,即所有小分类的请求
15     son_items = []
16     # 基准xpath: 获取所有大分类的对象列表
17     div_list = response.xpath('//div[@id="tab01"]/div')
18     for div in div_list:
19         # 1个大分类名称
20         parent_name = div.xpath('./h3/a/text()').get()
21         parent_url = div.xpath('./h3/a/@href').get()
22         if parent_name and parent_url:
23             # 1个大分类下面的所有小分类
24             li_list = div.xpath('./ul/li')
25             for li in li_list:
26                 # 继续交给调度器的item对象,确保每个是独立的
27                 item = SinaItem()
28                 item['son_name'] = li.xpath('./a/text()').get()
29                 item['son_url'] = li.xpath('./a/@href').get()
30                 item['parent_name'] = parent_name
31                 item['parent_url'] = parent_url
32                 # 创建对应的文件夹
33                 directory =
34                 './data/{}/{}/'.format(item['parent_name'], item['son_name'])
35                 item['son_filename'] = directory
36                 if not os.path.exists(directory):
37                     os.makedirs(directory)
38                 son_items.append(item)
39
40     # 大循环结束,items中存放了所有的每个小类请求的 item 对象
41     # 发送每个请求到调度器,得到response连同meta数据一起回调函数parse_son_url方法
42     for item in son_items:
43         yield scrapy.Request(url=item['son_url'], meta=
44         {'meta_1': item}, callback=self.parse_son_url)
45
46     def parse_son_url(self, response):
47         """解析小分类函数"""
48         meta1_item = response.meta['meta_1']
49         # 存放所有新闻页链接的item请求列表
50         meta1_items = []
51         parent_url = meta1_item['parent_url']
52         # 通过观察规律,小分类中的所有新闻连接均以最外层大链接开头
53         news_link_list = response.xpath('//a/@href').extract()
54         for news_link in news_link_list:
55             if news_link.startswith(parent_url) and news_link.endswith('.shtml'):
56                 # 此item对象是需要继续交给调度器入队列的请求
57                 item = SinaItem()
58                 # 把具体的新闻链接发送请求,保存到对应的文件夹下
59                 item['son_name'] = meta1_item['son_name']
60                 item['son_url'] = meta1_item['son_url']
61                 item['parent_name'] = meta1_item['parent_name']
62                 item['parent_url'] = meta1_item['parent_url']
63                 item['son_filename'] = meta1_item['son_filename']
64                 item['article_url'] = news_link
65                 meta1_items.append(item)
66
67     # 发送每个小类下的新闻链接请求,得到response后和meta数据交给get_content()函数处理
68     for item in meta1_items:
69         yield scrapy.Request(url=item['article_url'], meta={'meta_2': item},
70         callback=self.get_content)

```

```

68
69     def get_content(self, response):
70         """获取每个新闻的具体内容"""
71         item = response.meta['meta_2']
72         item['article_head'] = response.xpath('//h1[@class="main-title"]/text() |
//span[@class="location"]/h1/text()).get()
73         item['article_content'] =
'\n'.join(response.xpath('//div[@class="article"]/p/text() | //div[@class="article
clearfix"]/p/text() | //div[@id="artibody"]/p/text()).extract())
74
75         yield item

```

■ 4-数据处理

```

1  """pipelines.py"""
2  # -*- coding: utf-8 -*-
3
4  class SinaPipeline(object):
5      def process_item(self, item, spider):
6          # 文件名使用url地址的中间(即去掉协议和后缀.shtml)
7          filename = item['article_url'][7:-6].replace('/', '-')
8          filename_ = item['son_filename'] + filename + '.txt'
9          # 写入本地文件
10         with open(filename_, 'w', encoding='utf-8') as f:
11             f.write(item['article_content'])
12
13         return item

```

■ 5-全局配置

```

1  """settings.py"""
2  ROBOTSTXT_OBEY = False
3  DOWNLOAD_DELAY = 1
4  DEFAULT_REQUEST_HEADERS = {
5      'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
6      'Accept-Language': 'en',
7      'User-Agent': 'Mozilla/5.0',
8  }
9  ITEM_PIPELINES = {
10     'Sina.pipelines.SinaPipeline': 300,
11 }

```

■ 6-运行爬虫

```

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

```

今日作业

```
1 【1】 scrapy框架有哪几大组件? 以及各个组件之间是如何工作的?
2
3 【2】 腾讯招聘尝试改写为scrapy
4     2.1) response.text : 获取页面响应内容
5     2.2) scrapy中同样可以使用之前学过的模块,比如 json模块 等
6
7 【3】 盗墓笔记小说抓取
8     3.1) 目标: 抓取盗墓笔记1-8中所有章节的所有小说的具体内容, 保存到本地文件
9     3.2) http://www.daomubiji.com/
10    3.3) 保存路径示例: ./novel/盗墓笔记1:七星鲁王宫/七星鲁王_第一章_血尸.txt
11                      ./novel/盗墓笔记1:七星鲁王宫/七星鲁王_第二章_五十年后.txt
```

■ 盗墓笔记作业提示

```
1 【1】 一级页面xpath表达式:
2     a节点: //li[contains(@id,"menu-item-20")]/a
3     title: ./text()
4     link : ./@href
5
6 【2】 二级页面xpath表达式
7     基准xpath : //article
8     for循环遍历后:
9         name=article.xpath('./a/text()').get()
10        link=article.xpath('./a/@href').get()
11
12 【3】 三级页面xpath:
13     response.xpath('//article[@class="article-content"]/p/text()).extract()
14     # 结果: ['p1','p2','p3','']
```