

Day07笔记

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 - 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 - driver.switch_to.frame(frame节点对象)
11    3.2) 返回上一级 - driver.switch_to.parent_frame()
12    3.3) 返回主页面 - driver.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属性值)
16    4.2) 方法二:
17        a> 先找到frame节点 : frame_node = driver.find_element_by_xpath('xxx')
18        b> 在切换到frame : driver.switch_to.frame(frame_node)
```

■ 示例 - 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  driver = webdriver.Chrome()
9  driver.get('https://www.douban.com/')
10
11 # 切换到iframe子页面
12 login_frame = driver.find_element_by_xpath('//*[@id="anony-reg-
new"]/div/div[1]/iframe')
13 driver.switch_to.frame(login_frame)
14
15 # 密码登录 + 用户名 + 密码 + 登录豆瓣
16 driver.find_element_by_xpath('/html/body/div[1]/div[1]/ul[1]/li[2]').click()
17 driver.find_element_by_xpath('//*[@id="username"]').send_keys('自己的用户名')
18 driver.find_element_by_xpath('//*[@id="password"]').send_keys('自己的密码')
19 driver.find_element_by_xpath('/html/body/div[1]/div[2]/div[1]/div[5]/a').click()
20 time.sleep(3)
21
22 # 点击我的豆瓣
23 driver.find_element_by_xpath('//*[@id="db-nav-sns"]/div/div/div[3]/ul/li[2]/a').click()

```

selenium+phantomjs|chrome|firefox总结

```

1  【1】特点
2      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
3      1.2》需要等待页面元素加载，需要时间，效率低
4
5  【2】设置无界面模式
6      options = webdriver.ChromeOptions()
7      options.add_argument('--headless')
8      driver = webdriver.Chrome(executable_path='/home/tarena/chromedriver',options=options)
9
10 【3】鼠标操作
11     from selenium.webdriver import ActionChains
12     ActionChains(driver).move_to_element('node').perform()
13
14 【4】切换句柄 - switch_to.frame(handle)
15     all_handles = driver.window_handles
16     driver.switch_to.window(all_handles[1])
17
18 【5】iframe子页面
19     driver.switch_to.frame(frame_node)
20
21 【6】driver执行JS脚本
22     driver.execute_script('window.scrollTo(0,document.body.scrollHeight)')
23
24 【7】lxml中的xpath 和 selenium中的xpath的区别
25     7.1》lxml中的xpath用法 - 推荐自己手写
26     div_list = p.xpath('//div[@class="abc"]/div')

```

```

27         item = {}
28         for div in div_list:
29             item['name'] = div.xpath('..//a/@href')[0]
30             item['likes'] = div.xpath('..//a/text()')[0]
31
32 7.2》selenium中的xpath用法 - 推荐copy - copy xpath
33     div_list = driver.find_elements_by_xpath('//div[@class="abc"]/div')
34     item = {}
35     for div in div_list:
36         item['name'] = div.find_element_by_xpath('..//a').get_attribute('href')
37         item['likes'] = div.find_element_by_xpath('..//a').text

```

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('//*[
27              @id="transTarget"]/p/span').text
28
29          return result
30
31      def run(self):
32          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- 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
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()

```

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      sudo pip3 install Scrapy
3
4  【2】Windows安装
5      python -m pip install Scrapy
6
7      如果安装过程中报如下错误：'Error: Microsoft Visual C++ 14.0 is required xxx'
8      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】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 = 'https://www.guazi.com/bj/buy/o{}/#bread'.format(self.n)
45         # 把url交给调度器入队列
46         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  【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'
```

■ 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
16
17 【4】 日志级别
18     4.1》5个级别：DEBUG < INFO < WARNING < ERROR < CRITICAL
19         DEBUG：调试信息
20         INFO：一般信息
21         WARNING:警告信息
22         ERROR：错误信息
23         CRITICAL:严重错误
24     4.2》 settings.py中提供了2个变量
25         LOG_LEVEL = 'WARNING' # 只显示WARNING和比WARNING严重信息(WARNING ERROR CRITICAL)
26         LOG_FILE = 'guazi.log' # 把日志存放到日志文件中(guazi.log)

```

■ 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
43         yield item

```

■ 步骤3 - pipelines.py

```

1  # 将数据存入mongodb数据库,此处我们就不对MySQL表字段进行操作了,如有兴趣可自行完善
2  # MongoDB管道
3  import pymongo
4
5  class GuaziMongoPipeline(object):
6      def open_spider(self, spider):
7          """爬虫项目启动时只执行1次,用于连接MongoDB数据库"""
8          self.conn = pymongo.MongoClient(MONGO_HOST, MONGO_PORT)
9          self.db = self.conn[MONGO_DB]
10         self.myset = self.db[MONGO_SET]
11
12     def process_item(self, item, spider):
13         car_dict = dict(item)
14         self.myset.insert_one(car_dict)
15         return item

```

■ 步骤4 - settings.py

```

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

```

今日作业

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 【1】腾讯招聘职位信息抓取（二级页面）

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）

具体数据如下：

 - 1.1) 职位名称
 - 1.2) 职位地点
 - 1.3) 职位类别
 - 1.4) 发布时间
 - 1.5) 工作职责
 - 1.6) 工作要求