

Day10回顾

分布式爬虫原理及实现

```
1  # 原理
2  多台主机共享1个爬取队列
3  # 实现
4  scrapy_redis
```

分布式爬虫配置

```
1  1、完成非分布式scrapy爬虫项目
2  2、settings.py中指向新的调度器、去重机制及redis服务器IP地址及端口号
3  3、设置管道 - MySQL、MongoDB或Redis
```

redis_key分布式爬虫配置

```
1  1、完成非分布式scrapy爬虫项目 - 不能重写 start_requests()
2  2、settings.py中指向新的调度器、去重机制及redis服务器的IP地址及端口号
3  3、设置管道 - MySQL、MongoDB或Redis
4  4、爬虫文件使用RedisSpider类
5      1、去掉start_urls
6      2、添加redis_key : redis_key = "name:spider"
7  5、部署到多台服务器
8  6、连接redis, 执行 LPUSH name:spider First_URL
```

机器视觉

```
1  1、OCR
2  2、tesseract-ocr
3  3、pytesseract
```

Fiddler+Browser配置

```
1  1、Fiddler端 - 安装根证书、只抓浏览器、设置端口
2  2、Browser端 - 安装代理插件(SwithyOmega)并创建切换新代理
```

Fiddler+Phone配置

```
1  1、Fiddler端 - 安装根证书、抓所有进程、设置端口、允许远程
2  2、Phone端 - 修改网络、下载并安装证书
```

爬虫总结

1、什么是爬虫

爬虫是请求网站并提取数据的自动化程序

2、robots协议是什么

爬虫协议或机器人协议,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

3、爬虫的基本流程

1、请求得到响应

2、解析

3、保存数据

4、请求

1、urllib

2、requests

3、scrapy

5、解析

1、re正则表达式

2、lxml+xpath解析

3、json解析模块

6、selenium+browser

7、常见反爬策略

1、Headers : 最基本的反爬手段, 一般被关注的变量是UserAgent和Referer, 可以考虑使用浏览器中

2、UA : 建立User-Agent池, 每次访问页面随机切换

3、拉黑高频访问IP

数据量大用代理IP池伪装成多个访问者, 也可控制爬取速度

4、Cookies

建立有效的cookie池, 每次访问随机切换

5、验证码

验证码数量较少可人工填写

图形验证码可使用tesseract识别

其他情况只能在线打码、人工打码和训练机器学习模型

6、动态生成

一般由js动态生成的数据都是向特定的地址发get请求得到的, 返回的一般是json

7、签名及js加密

一般为本地JS加密, 查找本地JS文件, 分析, 或者使用execjs模块执行JS

8、js调整页面结构

9、js在响应中指向新的地址

8、scrapy框架的运行机制

9、分布式爬虫的原理

多台主机共享一个爬取队列

Day11笔记

移动端app数据抓取 - 浏览器F12

有道翻译手机版破解案例

```
1 import requests
2 from lxml import etree
3
4
5 word = input('请输入要翻译的单词:')
6 url = 'http://m.youdao.com/translate'
7 data = {
8     'inputtext': word,
9     'type': 'AUTO',
10 }
11
12 html = requests.post(url,data=data).text
13 parse_html = etree.HTML(html)
14 result = parse_html.xpath('//ul[@id="translateResult"]/li/text())[0]
15
16 print(result)
```

途牛旅游

目标

- 1 完成途牛旅游爬取系统，输入出发地、目的地，输入时间，抓取热门景点信息及相关评论

地址

- 1、地址: <http://www.tuniu.com/>
- 2、热门 - 搜索
- 3、选择: 相关目的地、出发城市、出游时间 (出发时间和结束时间) 点击确定
- 4、示例地址如下:
- 5 http://s.tuniu.com/search_complex/whole-sh-0-%E7%83%AD%E9%97%A8/list-a{触发时间}_{结束时间}-{出发城市}-{相关目的地}/

项目实现

■ 1、创建项目

```
1 scrapy startproject Tuniu
2 cd Tuniu
3 scrapy genspider tuniu tuniu.com
```

■ 2、定义要抓取的数据结构 - items.py

```
1 # 一级页面
2 # 标题 + 链接 + 价格 + 满意度 + 出游人数 + 点评人数 + 推荐景点 + 供应商
3 title = scrapy.Field()
4 link = scrapy.Field()
5 price = scrapy.Field()
6 satisfaction = scrapy.Field()
7 travelNum = scrapy.Field()
8 reviewNum = scrapy.Field()
9 recommended = scrapy.Field()
10 supplier = scrapy.Field()
11
12 # 二级页面
13 # 优惠券 + 产品评论
14 coupons = scrapy.Field()
15 cp_comments = scrapy.Field()
```

■ 3、爬虫文件数据分析与提取

页面地址分析

```
1 http://s.tuniu.com/search_complex/whole-sh-0-热门/list-a20190828_20190930-1200-m3922/
2 # 分析
3 list-a{出发时间_结束时间-出发城市-相关目的地}/
4 # 如何解决?
5 提取 出发城市及目的地城市的字典, key为城市名称, value为其对应的编码
6
7 # 提取字典, 定义config.py存放
```

代码实现

```
1 # -*- coding: utf-8 -*-
2 import scrapy
3 from ..config import *
4 from ..items import TuniuItem
5 import json
6
7 class TuniuSpider(scrapy.Spider):
8     name = 'tuniu'
9     allowed_domains = ['tuniu.com']
10
11     def start_requests(self):
12         s_city = input('出发城市:')
13         d_city = input('相关目的地:')
14         start_time = input('出发时间(20190828):')
15         end_time = input('结束时间(例如20190830):')
16         s_city = src_citys[s_city]
17         d_city = dst_citys[d_city]
18
```

```

19         url = 'http://s.tuniu.com/search_complex/whole-sh-0-%E7%83%AD%E9%97%A8/list-a{}_{}-
{}-{}'.format(start_time,end_time,s_city, d_city)
20         yield scrapy.Request(url, callback=self.parse)
21
22     def parse(self, response):
23         # 提取所有景点的li节点信息列表
24         items = response.xpath('//ul[@class="thebox clearfix"]/li')
25
26         for item in items:
27             # 此处是否应该在for循环内创建?
28             tuniuItem = TuniuItem()
29             # 景点标题 + 链接 + 价格
30             tuniuItem['title'] = item.xpath('..//span[@class="main-tit"]/@name').get()
31             tuniuItem['link'] = 'http:' + item.xpath('..//div/a/@href').get()
32             tuniuItem['price'] =
int(item.xpath('..//div[@class="tnPrice"]/em/text()').get())
33             # 判断是否为新产品
34             isnews = item.xpath('..//div[@class="new-pro"]').extract()
35             if not len(isnews):
36                 # 满意度 + 出游人数 + 点评人数
37                 tuniuItem['satisfaction'] = item.xpath('..//div[@class="comment-
satNum"]//i/text()').get()
38                 tuniuItem['travelNum'] = item.xpath('..//p[@class="person-
num"]//i/text()').get()
39                 tuniuItem['reviewNum'] = item.xpath('..//p[@class="person-
comment"]//i/text()').get()
40             else:
41                 tuniuItem['satisfaction'] = '新产品'
42                 tuniuItem['travelNum'] = '新产品'
43                 tuniuItem['reviewNum'] = '新产品'
44             # 包含景点+供应商
45             tuniuItem['recommended'] = item.xpath('..//span[@class="overview-
scenery"]/text()').extract()
46             tuniuItem['supplier'] =
item.xpath('..//span[@class="brand"]/span/text()').extract()
47
48             yield scrapy.Request(tuniuItem['link'], callback=self.item_info, meta={'item':
tuniuItem})
49
50         # 解析二级页面
51     def item_info(self, response):
52         tuniuItem = response.meta['item']
53         # 优惠信息
54         coupons = ','.join(response.xpath('//div[@class="detail-favor-coupon-
desc"]/@title').extract())
55         tuniuItem['coupons'] = coupons
56
57         # 想办法获取评论的地址
58         # 产品点评 + 酒店点评 + 景点点评
59         productId = response.url.split('/')[-1]
60         # 产品点评
61         cpdp_url = 'http://www.tuniu.com/papi/tour/comment/product?productId=
{}'.format(productId)
62
63         yield scrapy.Request(cpdp_url, callback=self.cpdp_func, meta={'item': tuniuItem})
64
65         # 解析产品点评

```

```

66     def cpdp_func(self, response):
67         tuniuItem = response.meta['item']
68
69         html = json.loads(response.text)
70         comment = {}
71         for s in html['data']['list']:
72             comment[s['realName']] = s['content']
73
74         tuniuItem['cp_comments'] = comment
75
76         yield tuniuItem

```

■ 4、管道文件处理 - pipelines.py

```

1 | print(dict(item))

```

■ 5、设置settings.py

出发城市和目的地城市的编号如何获取? - tools.py

```

1  # 出发城市
2  # 基准xpath表达式
3  /**/*
   [ @id="niuren_list" ] /div[2]/div[1]/div[2]/div[1]/div/div[1]/dl/dd/ul/li[contains(@class,"filter
   _input")] /a
4  name : ./text()
5  code : ./@href [0].split('/')[ -1].split('-')[ -1]
6
7  # 目的地城市
8  # 基准xpath表达式
9  /**
   [ @id="niuren_list" ] /div[2]/div[1]/div[2]/div[1]/div/div[3]/dl/dd/ul/li[contains(@class,"filter
   _input")] /a
10 name : ./text()
11 code : ./@href [0].split('/')[ -1].split('-')[ -1]

```

代码实现

```

1  import requests
2  from lxml import etree
3
4  url = 'http://s.tuniu.com/search_complex/whole-sh-0-%E7%83%AD%E9%97%A8/'
5  headers = {'User-Agent': 'Mozilla/5.0'}
6
7  html = requests.get(url, headers=headers).text
8  parse_html = etree.HTML(html)
9  # 获取出发地字典
10 # 基准xpath
11 li_list = parse_html.xpath('/**
   [ @id="niuren_list" ] /div[2]/div[1]/div[2]/div[1]/div/div[3]/dl/dd/ul/li[contains(@class,"filter
   _input")] /a')
12 src_citys = {}

```

```
13 dst_citys = {}
14 for li in li_list:
15     city_name_list = li.xpath('./text()')
16     city_code_list = li.xpath('./@href')
17     if city_name_list and city_code_list:
18         city_name = city_name_list[0].strip()
19         city_code = city_code_list[0].split('/')[0].split('-')[0]
20         src_citys[city_name] = city_code
21 print(src_citys)
22
23 # 获取目的地字典
24 li_list = parse_html.xpath('//*
25 [id="niuren_list"]/div[2]/div[1]/div[2]/div[1]/div/div[1]/dl/dd/ul/li[contains(@class,"filter
26 _input")]/a')
27 for li in li_list:
28     city_name_list = li.xpath('./text()')
29     city_code_list = li.xpath('./@href')
30     if city_name_list and city_code_list:
31         city_name = city_name_list[0].strip()
32         city_code = city_code_list[0].split('/')[0].split('-')[0]
33         dst_citys[city_name] = city_code
34 print(dst_citys)
```