

# 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(SwiftyOmega)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

## 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}/](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 |

```

### ■ 4、管道文件处理 - pipelines.py

```

1 |

```

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

```

#### 代码实现

```

1 |

```

