

爬虫期中作业

名称	修改日期	类型	大小
education_project	2019\4\22 星期...	文件夹	
pythonWeb_venv	2019\4\14 星期...	文件夹	

pythonWeb_venv 为虚拟目录，education_project 是项目目录

education_project	2019\4\14 星期...	文件夹	
first_app	2019\4\14 星期...	文件夹	
log	2019\4\14 星期...	文件夹	
static	2019\4\14 星期...	文件夹	
templates	2019\4\22 星期...	文件夹	
爬虫	2019\4\22 星期...	文件夹	
db.sqlite3	2018\6\27 星期...	SQLITE3 文件	196 KB
manage.py	2018\6\27 星期...	Python File	1 KB
新建文本文档.txt	2018\9\26 星期...	文本文档	1 KB

Templates 存放 html，static 存放资源文件，爬虫存放爬虫文件，

名称	修改日期	类型	大小
ad.html	2018\9\9 星期日 ...	360 se HTML Do...	2 KB
article.html	2018\9\26 星期...	360 se HTML Do...	4 KB
articleAdd.html	2018\9\26 星期...	360 se HTML Do...	9 KB
author.html	2018\9\23 星期...	360 se HTML Do...	1 KB
base.html	2018\9\26 星期...	360 se HTML Do...	6 KB
course.html	2018\9\10 星期...	360 se HTML Do...	9 KB
courseVideo.html	2018\9\25 星期...	360 se HTML Do...	10 KB
failure.html	2015\6\18 星期...	360 se HTML Do...	1 KB
index.html	2019\4\22 星期...	360 se HTML Do...	23 KB
jingdong.html	2019\4\22 星期...	360 se HTML Do...	10 KB
login.html	2018\9\26 星期...	360 se HTML Do...	5 KB
my.html	2018\9\17 星期...	360 se HTML Do...	1 KB
ouwang.html	2019\4\22 星期...	360 se HTML Do...	10 KB
pagination.html	2019\4\14 星期...	360 se HTML Do...	2 KB
pagination2.html	2019\4\14 星期...	360 se HTML Do...	3 KB
remen.html	2018\9\24 星期...	360 se HTML Do...	1 KB
userLead.html	2018\9\24 星期...	360 se HTML Do...	1 KB
新建文本文档.txt	2018\9\1 星期日	文本文档	50 KB

nl 修改日期: 2018\9\26 星期三 15:10 创建日期: 2019\4\14 星期日 9:54

名称	修改日期	类型	大小
assets	2019\4\14 星期...	文件夹	
css	2019\4\14 星期...	文件夹	
images	2019\4\14 星期...	文件夹	
js	2019\4\14 星期...	文件夹	
lib	2019\4\14 星期...	文件夹	
mp4	2019\4\14 星期...	文件夹	

Get_ip 从网上爬取免费 ip

名称	修改日期	类型	大小
get_ip	2019\4\22 星期...	文件夹	
jidong	2019\4\22 星期...	文件夹	
scrapy	2019\4\22 星期...	文件夹	

静态网站：瓯网使用 scrapy 爬取，

修改 scrapy 的_requests_to_follow 方法，把 url 写入到 redis.

```

def _requests_to_follow(self, response):
    if not isinstance(response, HtmlResponse):
        return
    reg = r'article\d*show.html'
    seen = set()
    for n, rule in enumerate(self._rules):
        links = [lnk for lnk in rule.link_extractor.extract_links(response)
                  if lnk not in seen]
        if links and rule.process_links:
            links = rule.process_links(links)
        for link in links:
            match = re.findall(reg, str(link.url))
            if match == []:
                if self.redis.is_existArticle(link.url):
                    continue
                else:
                    self.redis.putArticle(link.url)
                    seen.add(link)
                    r = self._build_request(n, link)
                    yield rule.process_request(r)
            else:
                self.redis.put(link.url)

```

解析网页时使用多线程,使用新闻通用 filter。

myfilter.py	2019\4\20 星期...	Python File	6 KB
-------------	-----------------	-------------	------

```

if __name__ == "__main__":
    print('Parent process %s.' % os.getpid())
    p = Pool(10)

    for i in range(10):
        p.apply_async(main, args=())
    print('Waiting for all subprocesses done...')
    p.close()
    p.join()
    print('All subprocesses done.')

```

名称	修改日期	类型	大小
__pycache__	2019\4\22 星期...	文件夹	
myfilter.py	2019\4\20 星期...	Python File	6 KB
myMongodb.py	2018\7\8 星期日 ...	Python File	2 KB
mysqlDb.py	2019\4\20 星期...	Python File	3 KB
ouwang.py	2019\4\22 星期...	Python File	3 KB
ouwang2.py	2018\7\1 星期日 ...	Python File	2 KB
redisDB.py	2018\7\1 星期日 ...	Python File	2 KB
setting.py	2019\4\20 星期...	Python File	1 KB
textExcel.py	2018\7\1 星期日 ...	Python File	1 KB

京东使用 selenium 爬取

名称	修改日期	类型	大小
__pycache__	2019\4\22 星期...	文件夹	
1.py	2019\4\14 星期...	Python File	3 KB
2.py	2019\4\17 星期...	Python File	0 KB
chromedriver.exe	2018\3\20 星期...	应用程序	6,207 KB
db.py	2018\7\23 星期...	Python File	2 KB
jingdong.py	2019\4\19 星期...	Python File	3 KB
mysqlDb.py	2019\4\18 星期...	Python File	3 KB
setting.py	2019\4\10 星期...	Python File	1 KB

```

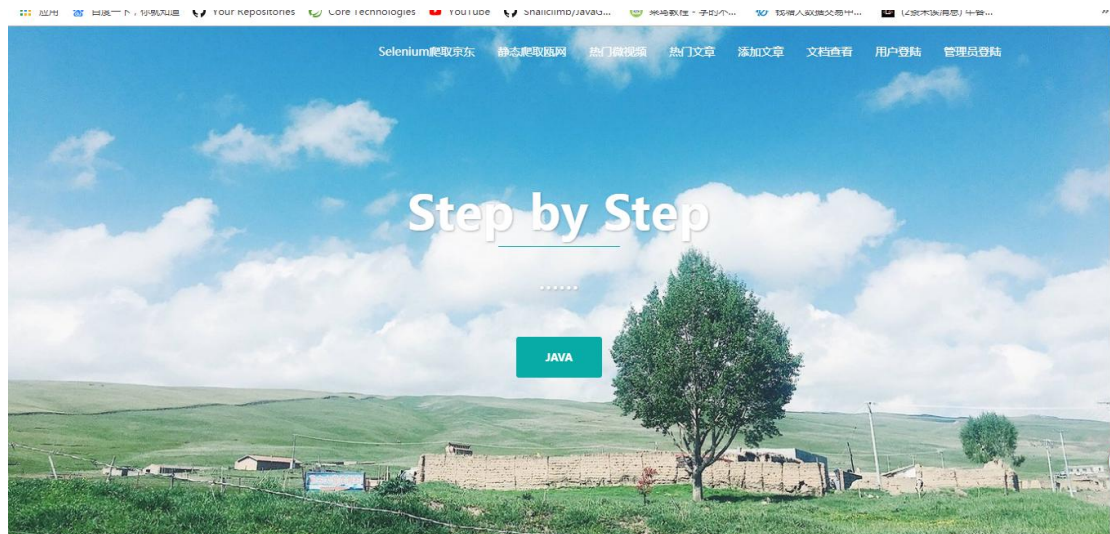
temp += 1

except Exception as e:
    print(e)
finally:
    driver.close()

if __name__ == '__main__':
    keyword = "电脑"
    #设置爬取页面
    page = 30
    page = page * 2
    print("开始爬取京东")
    base_url = 'https://search.jd.com/Search?keyword={key}&enc=utf-8&qrst=1&rt=1&stop=1&vt=2&wg={key}'
    for i in range(1, page, 2):
        url = base_url.format(key=keyword,page=i);
        get_jingdong(url, i)
        time.sleep(5)

```

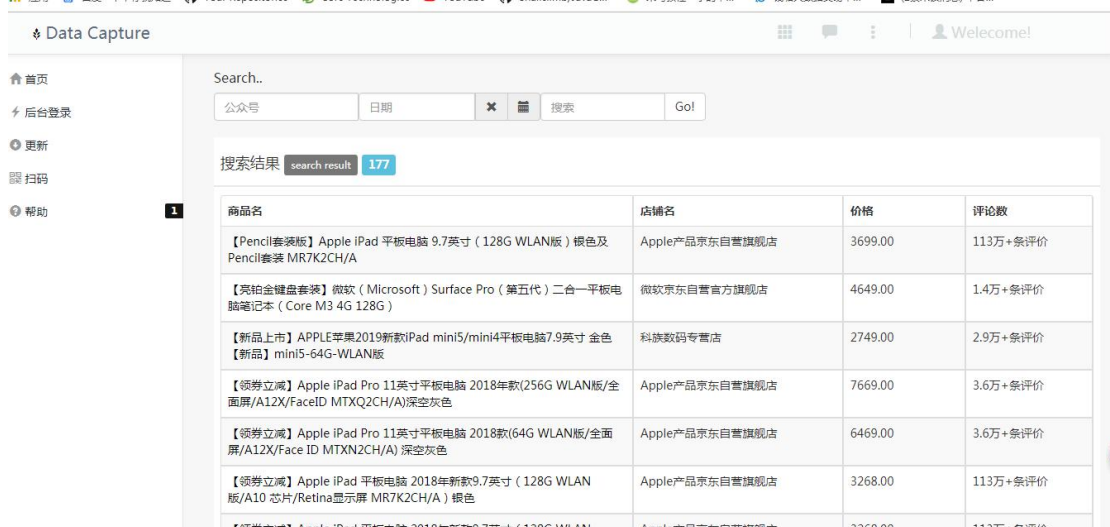
首页



爬取京东

```
def jingdong(request):
    jingdong_list = Jingdong.objects.all().order_by('productname')
    totalNum = 0
    for i in jingdong_list:
        totalNum += 1
    #jingdong_list, thisPage1, thisPage2, thisPage3, thisPage4, thisPage5 = getPage(request, jingdong_list)
    paginator = Paginator(jingdong_list, 25) # 每页显示25条

    page = request.GET.get('page')
    try:
        jingdong_list = paginator.page(page)
    except PageNotAnInteger:
        # 如果请求的页数不是整数，返回第一页。
        jingdong_list = paginator.page(1)
    except EmptyPage:
        # 如果请求的页数不在合法的页数范围内，返回结果的最后一页。
        jingdong_list = paginator.page(paginator.num_pages)
    return render(request, 'jingdong.html', locals())
```



Data Capture

首页

后台登录

更新

扫码

帮助

Search..

公众号

日期

✕

☰

搜索

Go!

搜索结果

search result

1192

文章标题	发布时间	浏览人数	内容(只显示10个)
丹麦警察备战气候大会	2009/12/07	6915	<p>丹麦警察备战气候
【求真 第3期】中共温州市委副书记、市长 张耕加快特色小镇建设 打造局部环境优化的产业平台	2016/12/15	6346	<p>【求真 第3期
【求真 第3期】冯金考:打造品质城市 建设美丽瑞安 ——瑞安品质城市建设的调研与思考	2016/12/15	6645	<p>【求真 第3期
【求真 第3期】冯金考:打造品质城市 建设美丽瑞安 ——瑞安品质城市建设的调研与思考	2016/12/15	6642	<p>【求真 第3期
【求真 第3期】吴育民:巧用微信软件 打造“指尖上的统战部”	2016/12/15	5907	<p>文成海外华侨华
【求真 第3期】周贇:探“问”之道 增“问”之效 ——龙湾区人大常委会开展专题询问的实践与思考	2016/12/15	5226	<p>【求真 第3期
【求真 第3期】姜绍光:打造美丽乡村升级版 开创统筹发展新格局	2016/12/15	6306	<p>近年来,永嘉县
乐清市招商引资情况的调查报告	2016/12/15	8353	<p>招商引资对地方

localhost:7888/suwanan#

Selenium爬取京东

静态爬取网页

热门微视频

热门文章

添加文章

文档查看

用户登陆

管理员登陆

热门文章



集成化信息传播与服务培育

python

推送一些python小资源, 个性化操作, 项目介绍, 实战推荐

爬虫

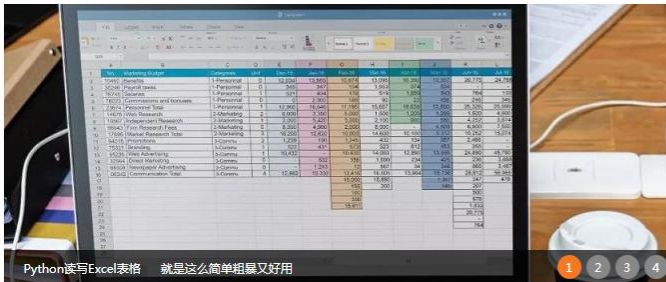
介绍一些各式网站不同的爬取, 全站爬取, 搜索引擎等

java

java SE基础知识, SSM框架, springBoot,springCould微服务等热门技术分享

个人博客

专注Python开发, 欢迎和大家交流

[首页](#)[关于我](#)[文章](#)[添加文章](#)[用户排行榜](#)

用户文章

如何玩转微服务

如何玩转微服务

java 2018-09-26

评论 (5) 浏览 (14)

关注我



CSDN



简书



公众号



邮箱

站长推荐

1 如何玩转微服务

2 Python3 File(文件) 方法

3 (转) JavaEE基础知识回顾

4 Java SE基础知识30问

5 (转)Spring AOP中定义切入点 (PointCut)

[首页](#)[关于我](#)[文章](#)[添加文章](#)[用户排行榜](#)

如何玩转微服务

javaSE 2018-09-26

微服务, 软件应用开发的新纪元

2014年 Martin Fowler 在《MicroServices》论文中首次提出了微服务的概念。近些年, 伴随着互联网的日益发展, 微服务在国内、甚至国际上的发展已达到一个新高潮。

在微服务流行之前, SOA (Service Oriented Architecture) 被广泛熟知与采用。微服务基于 SOA 发展而来, 但与之相比, 微服务更易于理解, 也更利于设计者、开发者的实践落地, 它把“面向服务”的设计思想实现得更加彻底。

关注我



CSDN



简

标签云

Python

spring

评论功能

评论 & comment

您的评论或留言（必填）

评论

#2

万老师(Sept. 26, 2018, 12:29 p.m.)

文章很好！

#2

廖学长(Sept. 26, 2018, 12:30 p.m.)

搭配着项目一起进行就更好了

Selenium爬取京东 静态爬取网页 热门微视频 热门文章 添加文章 文档查看 用户登陆 管理员登陆

微视频



python基础

第一个python程序与数据存储01
print&input与变量和运算符01
字符串与循环中的while01
布尔&list与条件循环语句与turtle01
元组&字符串&字典01



爬虫

爬虫基本原理讲解 Urllib库基本使用
Requests库基本使用 正则表达式基础
BeautifulSoup库详解



Django

Django基本流程走通01
Django中的模型01
Django中的模板01
Django中的视图01
Django的高级使用01

文档简介



python基础



Django



python爬虫

