



官方微博:

加关注

7545

搜索

[首页](#) [Python论坛](#) [Python基础教程](#) [Python高级教程](#) [Python框架](#) [Python函数](#) [GUI教程](#) [Linux教程](#) [PHP教程](#) [在线手册大全](#)百万级python导师
亲身指导

保你120天变身python大牛

有时候,
你需要的只是一句点拨

立即咨询

[首页](#) > [Python高级教程](#) > 正文

python爬虫框架scrapy实例详解

2013-08-14 09:24:47

生成项目scrapy提供一个工具来生成项目，生成的项目中预置了一些文件，用户需要在这些文件中添加自己的代码。打开命令行，执行：scrapy st...

生成项目

scrapy提供一个工具来生成项目，生成的项目中预置了一些文件，用户需要在这些文件中添加自己的代码。

打开命令行，执行：scrapy startproject tutorial，生成的项目类似下面的结构

tutorial/

scrapy.cfg

tutorial/

__init__.py

items.py

pipelines.py

settings.py

spiders/

__init__.py

...

scrapy.cfg是项目的配置文件

用户自己写的spider要放在spiders目录下，一个spider类似

```
1 from scrapy.spider import BaseSpider
2 class DmozSpider(BaseSpider):
3     name = "dmoz"
4     allowed_domains = ["dmoz.org"]
5     start_urls = [
6         "http://www.dmoz.org/Computers/Programming/Languages/Python/Books/",
7         "http://www.dmoz.org/Computers/Programming/Languages/Python/Resources/"
8     ]
9     def parse(self, response):
10         filename = response.url.split("/")[-2]
11         open(filename, 'wb').write(response.body)
```

name属性很重要，不同spider不能使用相同的name

start_urls是spider抓取网页的起始点，可以包括多个url

parse方法是spider抓到一个网页以后默认调用的callback，避免使用这个名字来定义自己的方法。

当spider拿到url的内容以后，会调用parse方法，并且传递一个response参数给它，response包含了抓到的网页的内容，在parse方法里，你可以从抓到的网页里面解析数据。上面的代码只是简单地把

频道总排行

python下redis安装和使用
ImportError: No module named
python 在互联网应用是如此强大
python读取和生成excel文件
python多线程编程2—线程的创建、
分别用python2和python3伪装浏览器
python批量抓取美女图片
Python简单两步实现天气爬虫采集器
成人网站性能提升 20 倍之经验谈
python抓取google搜索结果

频道本月排行

python多线程编程2—线程的创建、
python爬虫框架scrapy实例详解
python读取和生成excel文件
python批量抓取美女图片
分别用python2和python3伪装浏览器
成人网站性能提升 20 倍之经验谈
一步步来用C语言来写python扩展
用python实现的抓取腾讯视频所有电
用python 装饰器打log
python抓取google搜索结果

最新文章

IP地址网段表示法总结
pyOpenSSL版本问题导致Scrapy安装报
Linux开机启动程序或脚本详解
一步步教你理解Python装饰器
字符串的encode与decode解决乱码问题
Python操作cookie之cookielib模块
作为面试官的一些心得
图像库PIL相关笔试题大总结
为什么你会迁移到Python3.5 ?
滴滴悄然提价 - 中国打车市场补贴战或
众筹3万起诉广电总局 - 起因游戏审核条
在Python中添加自定义模块

相关文章

python 在互联网应用是如此强大
sys.argv[] 的使用详解
ImportError: No module named
python下redis安装和使用

网页内容保存到文件。

开始抓取

你可以打开命令行，进入生成的项目根目录tutorial/，执行 `scrapy crawl dmoz`，`dmoz`是spider的name。

解析网页内容

scrapy提供了方便的办法从网页中解析数据，这需要使用到HtmlXPathSelector

```
1 from scrapy.spider import BaseSpider
2 from scrapy.selector import HtmlXPathSelector
3 class DmozSpider(BaseSpider):
4     name = "dmoz"
5     allowed_domains = ["dmoz.org"]
6     start_urls = [
7         "http://www.dmoz.org/Computers/Programming/Languages/Python/Books/",
8         "http://www.dmoz.org/Computers/Programming/Languages/Python/Resources/"
9     ]
10    def parse(self, response):
11        hxs = HtmlXPathSelector(response)
12        sites = hxs.select('//ul/li')
13        for site in sites:
14            title = site.select('a/text()').extract()
15            link = site.select('a/@href').extract()
16            desc = site.select('text()').extract()
17            print title, link, desc
```

HtmlXPathSelector使用了Xpath来解析数据

`//ul/li`表示选择所有的ul标签下的li标签

`a/@href`表示选择所有a标签的href属性

`a/text()`表示选择a标签文本

`a[@href="abc"]`表示选择所有href属性是abc的a标签

我们可以把解析出来的数据保存在一个scrapy可以使用的对象中，然后scrapy可以帮助我们把这些对象保存起来，而不用我们自己把这些数据存到文件中。我们需要在items.py中添加一些类，这些类用来描述我们要保存的数据

```
from scrapy.item import Item, Field
```

```
class DmozItem(Item):
```

```
    title = Field()
```

```
    link = Field()
```

```
    desc = Field()
```

然后在spider的parse方法中，我们把解析出来的数据保存在DmozItem对象中。

```
1 from scrapy.spider import BaseSpider
2 from scrapy.selector import HtmlXPathSelector
3 from tutorial.items import DmozItem
4 class DmozSpider(BaseSpider):
5     name = "dmoz"
6     allowed_domains = ["dmoz.org"]
7     start_urls = [
8         "http://www.dmoz.org/Computers/Programming/Languages/Python/Books/",
9         "http://www.dmoz.org/Computers/Programming/Languages/Python/Resources/"
10    ]
11    def parse(self, response):
12        hxs = HtmlXPathSelector(response)
13        sites = hxs.select('//ul/li')
14        items = []
15        for site in sites:
16            item = DmozItem()
17            item['title'] = site.select('a/text()').extract()
18            item['link'] = site.select('a/@href').extract()
19            item['desc'] = site.select('text()').extract()
20            items.append(item)
21        return items
```

在命令行执行scrapy的时候，我们可以加两个参数，让scrapy把parse方法返回的items输出到json文件中

```
scrapy crawl dmoz -o items.json -t json
```

items.json会被放在项目的根目录

让scrapy自动抓取网页上的所有链接

上面的示例中scrapy只抓取了start_urls里面的两个url的内容，但是通常我们想实现的是scrapy自动发现一个网页上的所有链接，然后再去抓取这些链接的内容。为了实现这一点我们可以在parse方法里面提取我们需要的链接，然后构造一些Request对象，并且把他们返回，scrapy会自动的去抓取这些链接。代码类似：

```

1 class MySpider(BaseSpider):
2     name = 'myspider'
3     start_urls = (
4         'http://example.com/page1',
5         'http://example.com/page2',
6     )
7     def parse(self, response):
8         # collect item_urls
9         for item_url in item_urls:
10             yield Request(url=item_url, callback=self.parse_item)
11     def parse_item(self, response):
12         item = MyItem()
13         # populate `item` fields
14         yield Request(url=item_details_url, meta={'item': item},
15             callback=self.parse_details)
16     def parse_details(self, response):
17         item = response.meta['item']
18         # populate more `item` fields
19         return item

```

parse是默认的callback, 它返回了一个Request列表，scrapy自动的根据这个列表抓取网页，每当抓到一个网页，就会调用parse_item，parse_item也会返回一个列表，scrapy又会根据这个列表去抓网页，并且抓到后调用parse_details

为了让这样的工作更容易，scrapy提供了另一个spider基类，利用它我们可以方便的实现自动抓取链接。我们要用到CrawlSpider

```

1 from scrapy.contrib.linkextractors.sgml import SgmlLinkExtractor
2 class MininovaSpider(CrawlSpider):
3     name = 'mininova.org'
4     allowed_domains = ['mininova.org']
5     start_urls = ['http://www.mininova.org/today']
6     rules = [Rule(SgmlLinkExtractor(allow=['/tor/\d+'])),
7         Rule(SgmlLinkExtractor(allow=['/abc/\d+']), 'parse_torrent')]
8     def parse_torrent(self, response):
9         x = HtmlXPathSelector(response)
10         torrent = TorrentItem()
11         torrent['url'] = response.url
12         torrent['name'] = x.select("//h1/text()").extract()
13         torrent['description'] =
14 x.select("//div[@id='description']").extract()
15         torrent['size'] = x.select("//div[@id='info-left']/p[2]/text()
16 [2]").extract()
17         return torrent

```

相比BaseSpider，新的类多了一个rules属性，这个属性是一个列表，它可以包含多个Rule，每个Rule 描述了哪些链接需要抓取，哪些不需要。这是Rule类的文档<http://doc.scrapy.org/en/latest/topics/spiders.html#scrapy.contrib.spiders.Rule>

这些rule可以有callback，也可以没有，当没有callback的时候，scrapy简单的follow所有这些链接。

pipelines.py的使用

在pipelines.py中我们可以添加一些类来过滤掉我们不想要的item，把item保存到数据库。

```

1 from scrapy.exceptions import DropItem
2 class FilterWordsPipeline(object):
3     """A pipeline for filtering out items which contain certain words in
4     their
5     description"""
6     # put all words in lowercase
7     words_to_filter = ['politics', 'religion']
8     def process_item(self, item, spider):
9         for word in self.words_to_filter:
10             if word in unicode(item['description']).lower():
11                 raise DropItem("Contains forbidden word: %s" % word)
12             else:
13                 return item

```

如果item不符合要求，那么就抛一个异常，这个item不会被输出到json文件中。

要使用pipelines，我们还需要修改settings.py

添加一行

```
ITEM_PIPELINES = ['dirbot.pipelines.FilterWordsPipeline']
```

现在执行scrapy crawl dmoz -o items.json -t json，不符合要求的item就被过滤掉了

相关文章：

- python 在互联网应用是如此强大
- sys.argv[] 的使用详解
- ImportError: No module named setuptools 解决方法
- python下redis安装和使用
- python链接mysql数据库详解
- python+ mysql存储二进制流的方式

关键词搜索： python 爬虫 框架

上一篇： Python 如何将一变量做为函数名？

下一篇： 解决ImportError: libmysqlclient_r.so.16: cannot open shared object file

如对本文内容有疑问，或想进一步交流学习，欢迎通过以下方式：

1. Python论坛
2. python技术互助群(请不要加多个群):
群④：385100854
群③：318130924
群②：333646237
群①：87464755
3. 关注PythonTab微信: Pythontab，公众号: Pythontab中文网

版权声明：本站文章除非注明，均为原创内容，如需转载请务必注明出处，违者本站保留追究其法律责任之权利。



如果你遇到什么问题，请留言回复，我们会及时帮助解决...

发布

还没有评论，沙发等你来抢

[关于我们](#) [联系方式](#) [版权声明](#) [Python入门](#) [Python进阶](#) [Python框架](#) [Python GUI](#) [Python内置函数](#) [Linux教程](#) [MySQL教程](#) [PHP教程](#)

Copyright (C) 2012-2014 PYTHONTAB.COM, All Rights Reserved. PythonTab: Python中文开发者社区门户。

官方唯一联系方式: Email : market@pythontab.com QQ : 20769164