
Scrapy Documentation

Scrapy developers

2019 年 06 月 28 日

1	获得帮助	3
2	第一步	5
2.1	预览 Scrapy	5
2.2	安装指南	8
2.3	Scrapy 教程	12
2.4	例子	27
3	基本概念	29
3.1	命令行工具	29
3.2	爬虫器	41
3.3	选择器	55
3.4	Items	74
3.5	Item 加载器	79
3.6	Scrapy shell	93
3.7	Item 管道	99
3.8	Feed 导出	104
3.9	请求和响应	110
3.10	链接提取	124
3.11	设置	126
3.12	异常	154
4	内置服务	157
4.1	日志	157
4.2	统计数据集合	161
4.3	发送邮件	163
4.4	远程控制台	166
4.5	Web Service	170

5 解决具体问题	171
5.1 常见问答	171
5.2 Debug 爬虫器	177
5.3 爬虫器约束	180
5.4 常见实践	182
5.5 并发爬虫	187
5.6 使用浏览器的开发工具进行抓取	190
5.7 内存泄漏调试	196
5.8 下载并处理文件和图片	202
5.9 部署爬虫器	211
5.10 爬虫器节流	212
5.11 爬虫器硬件性能	214
5.12 作业: 暂停并恢复爬行	217
6 Scrapy 扩展	221
6.1 框架体系	221
6.2 下载器中间件	224
6.3 爬虫器中间件	240
6.4 扩展	247
6.5 核心 API	253
6.6 信号	257
6.7 Item 导出文件	262
7 其他	271
7.1 Scrapy 更新	271
7.2 贡献 Scrapy 代码	330
7.3 版本控制和 API 稳定性	335
Python 模块索引	337
索引	339

Scrapy 是一个快速、高效率的网络爬虫框架，用于抓取 web 站点并从页面中提取结构化的数据。Scrapy 被广泛用于数据挖掘、监测和自动化测试。

CHAPTER 1

获得帮助

遇到困难了？我们乐意帮你解决！

- 试试[FAQ](#) – 这里有一些常见问题的解答。
- 寻找具体信息？试试 `genindex` or `modindex`。
- 在 [StackOverflow using the scrapy tag](#) 里提问或搜索问题。
- 在 [Scrapy subreddit](#) 里提问或搜索问题。
- 在 [scrapy-users mailing list](#) 文档里搜索问题。
- 在 [#scrapy IRC channel](#) 上提问。
- 提交 Scrapy 错误报告请点击 [issue tracker](#)。

2.1 预览 Scrapy

Scrapy 是一个用于抓取和提取 web 站点的结构化数据的应用框架，可用于广泛的有用的应用程序，如数据挖掘、信息处理或历史档案。

尽管 Scrapy 最初是为 web scraping 而设计的，但它也可以用 API 提取数据 (例如 Amazon Associates Web Services) 或者作为一个通用的网络爬虫。

2.1.1 浏览爬虫器例子

为了向你展示 Scrapy 带来了什么，我们将用一个 Scrapy 最简单的方法带你浏览一个 Scrapy 爬虫器的例子。下面是这个爬虫器的代码，它可以从 <http://quotes.toscrape.com> 网站上抓取名人名言，根据以下的页码：

```
import scrapy

class QuotesSpider(scrapy.Spider):
    name = 'quotes'
    start_urls = [
        'http://quotes.toscrape.com/tag/humor/',
    ]
```

(下页继续)

(续上页)

```
def parse(self, response):
    for quote in response.css('div.quote'):
        yield {
            'text': quote.css('span.text::text').get(),
            'author': quote.xpath('span/small/text()').get(),
        }

    next_page = response.css('li.next a::attr("href")').get()
    if next_page is not None:
        yield response.follow(next_page, self.parse)
```

把代码放到一个 Python 文件中，文件名如 `quotes_spider.py` 然后用 `runspider` 命令运行这个爬虫器：

```
scrapy runspider quotes_spider.py -o quotes.json
```

当它执行完毕后，你会发现目录下多了一个名为 `quotes.json` 的文件，文件内容是以 JSON 格式储存的，其中包含名言和作者，格式如下（为了美观进行了排版）：

```
[{
    "author": "Jane Austen",
    "text": "\u201cThe person, be it gentleman or lady, who has not pleasure in a good_
    \u2013novel, must be intolerably stupid.\u201d"
},
{
    "author": "Groucho Marx",
    "text": "\u201cOutside of a dog, a book is man's best friend. Inside of a dog it's_
    \u2013too dark to read.\u201d"
},
{
    "author": "Steve Martin",
    "text": "\u201cA day without sunshine is like, you know, night.\u201d"
},
...]
```

在此期间发生了什么？

当你运行 `scrapy runspider quotes_spider.py` 命令之后，Scrapy 会找到定义的爬虫器，然后通过爬虫引擎运行它们。

爬虫是通过发送请求给 `start_urls` 属性中定义的 url(在本例中，就是 *humor* 分类下的那些 url) 然后调用默认的回调函数 `parse`，并把响应的对象作为参数传给它。在 `parse` 回调函数中，我们使用 CSS 选择器遍

历 quote 元素，并把解析的引用名言和作者生成一个字典通过生成器返回，寻找并请求下一个链接且继续使用 parse 方法作为回调函数。

这里你可以注意到 Scrapy 的一个主要的优势：[请求的调度是异步的](#)。这意味着 Scrapy 不需要等待一个请求被完成处理，它可以同时发送另一个请求或者做一些其他事情。这也同样表明当一些请求失败或者处理发生错误时其他请求可以继续下去。

虽然这使得你可以非常快速地进行爬虫（同时可以并发多个请求，在可承载的压力下），不过 Scrapy 同样提供了优雅的爬虫方法，[简单地配置即可](#)。你也可以设置一些其他的参数，比如给下载器设置一个延迟，限制每个域名或者代理 IP 的并发量，甚至可以让 Scrapy 自动根据请求的响应情况进行[限流](#)。

注解： 本例中[导出](#)了一个 JSON 文件，你可以轻易地更改输出格式（比如 XML 或者 CSV）或者储存到后端（比如 FTP 或者 [Amazon S3](#)）。你同样可以定义一个[item 管道](#)把这些 item 储存到数据库中。

2.1.2 还有什么？

你已经知道了如何用 Scrapy 从一个站点提取和储存 item，但是仅仅是很浅显的了解它。Scrapy 还为爬虫提供了很多强大的功能，比如：

- 用内置的 CSS 选择器和 XPath 语法从 HTML/XML 源中[选择和解析数据](#)。甚至可以在其中使用正则表达式来辅助解析。
- 一个[交互控制台](#) (IPython aware)，可以用 CSS 选择器和 XPath 语法来解析数据，在编写和调试爬虫器的时候十分方便。
- 内置支持[生成输出](#)多种格式文件 (JSON, CSV, XML) 和多种后端存储 (FTP, S3, 本地文件系统)。
- 强大的编码支持和检测，用于处理外部的，不标准的或者损坏的编码。
- [强大的扩展](#)，允许你植入自己写的函数，[信号](#)和定义好的 API(中间件, [扩展](#)，以及[管道](#))。
- 广泛的内置扩展功能和中间件处理：
 - 处理 cookies 和 session
 - HTTP 特性包括压缩、身份验证以及缓存
 - 模拟 user-agent
 - robots.txt 协议
 - 爬虫深度限制
 - 更多
- 一个[远程控制台](#)用于连接运行在 Scrapy 进程的 Python 控制台，以便于自省和调试爬虫程序。
- 还有其他好东西，比如从 [Sitemaps](#) and XML/CSV 源导入可复用的爬虫，一个可以自动下载图片（或者其他和 items 关联的媒体文件）的[媒体中间件](#)，一个 DNS 缓存解析器，以及更多！

2.1.3 下一步?

你下一步要做的就是下载 *Scrapy*, 跟着教程 去创建一个完整的 Scrapy 工程并且 加入社区. 感谢您的关注!

2.2 安装指南

2.2.1 Installing Scrapy

Scrapy runs on Python 2.7 and Python 3.4 or above under CPython (default Python implementation) and PyPy (starting with PyPy 5.9).

If you're using [Anaconda](#) or [Miniconda](#), you can install the package from the [conda-forge](#) channel, which has up-to-date packages for Linux, Windows and OS X.

To install Scrapy using `conda`, run:

```
conda install -c conda-forge scrapy
```

Alternatively, if you're already familiar with installation of Python packages, you can install Scrapy and its dependencies from PyPI with:

```
pip install Scrapy
```

Note that sometimes this may require solving compilation issues for some Scrapy dependencies depending on your operating system, so be sure to check the *Platform specific installation notes*.

We strongly recommend that you install Scrapy in *a dedicated virtualenv*, to avoid conflicting with your system packages.

For more detailed and platform specifics instructions, as well as troubleshooting information, read on.

Things that are good to know

Scrapy is written in pure Python and depends on a few key Python packages (among others):

- [lxml](#), an efficient XML and HTML parser
- [parsel](#), an HTML/XML data extraction library written on top of [lxml](#),
- [w3lib](#), a multi-purpose helper for dealing with URLs and web page encodings
- [twisted](#), an asynchronous networking framework
- [cryptography](#) and [pyOpenSSL](#), to deal with various network-level security needs

The minimal versions which Scrapy is tested against are:

- Twisted 14.0

- lxml 3.4
- pyOpenSSL 0.14

Scrapy may work with older versions of these packages but it is not guaranteed it will continue working because it's not being tested against them.

Some of these packages themselves depends on non-Python packages that might require additional installation steps depending on your platform. Please check *platform-specific guides below*.

In case of any trouble related to these dependencies, please refer to their respective installation instructions:

- [lxml installation](#)
- [cryptography installation](#)

Using a virtual environment (recommended)

TL;DR: We recommend installing Scrapy inside a virtual environment on all platforms.

Python packages can be installed either globally (a.k.a system wide), or in user-space. We do not recommend installing scrapy system wide.

Instead, we recommend that you install scrapy within a so-called “virtual environment” ([virtualenv](#)). Virtualenvs allow you to not conflict with already-installed Python system packages (which could break some of your system tools and scripts), and still install packages normally with `pip` (without `sudo` and the likes).

To get started with virtual environments, see [virtualenv installation instructions](#). To install it globally (having it globally installed actually helps here), it should be a matter of running:

```
$ [sudo] pip install virtualenv
```

Check this [user guide](#) on how to create your virtualenv.

注解: If you use Linux or OS X, [virtualenvwrapper](#) is a handy tool to create virtualenvs.

Once you have created a virtualenv, you can install scrapy inside it with `pip`, just like any other Python package. (See *platform-specific guides below* for non-Python dependencies that you may need to install beforehand).

Python virtualenvs can be created to use Python 2 by default, or Python 3 by default.

- If you want to install scrapy with Python 3, install scrapy within a Python 3 virtualenv.
- And if you want to install scrapy with Python 2, install scrapy within a Python 2 virtualenv.

2.2.2 Platform specific installation notes

Windows

Though it's possible to install Scrapy on Windows using pip, we recommend you to install [Anaconda](#) or [Miniconda](#) and use the package from the [conda-forge](#) channel, which will avoid most installation issues.

Once you've installed [Anaconda](#) or [Miniconda](#), install Scrapy with:

```
conda install -c conda-forge scrapy
```

Ubuntu 14.04 or above

Scrapy is currently tested with recent-enough versions of lxml, twisted and pyOpenSSL, and is compatible with recent Ubuntu distributions. But it should support older versions of Ubuntu too, like Ubuntu 14.04, albeit with potential issues with TLS connections.

Don't use the `python-scrapy` package provided by Ubuntu, they are typically too old and slow to catch up with latest Scrapy.

To install scrapy on Ubuntu (or Ubuntu-based) systems, you need to install these dependencies:

```
sudo apt-get install python-dev python-pip libxml2-dev libxslt1-dev zlib1g-dev libffi-  
dev libssl-dev
```

- `python-dev`, `zlib1g-dev`, `libxml2-dev` and `libxslt1-dev` are required for `lxml`
- `libssl-dev` and `libffi-dev` are required for `cryptography`

If you want to install scrapy on Python 3, you'll also need Python 3 development headers:

```
sudo apt-get install python3 python3-dev
```

Inside a [virtualenv](#), you can install Scrapy with `pip` after that:

```
pip install scrapy
```

注解: The same non-Python dependencies can be used to install Scrapy in Debian Jessie (8.0) and above.

Mac OS X

Building Scrapy's dependencies requires the presence of a C compiler and development headers. On OS X this is typically provided by Apple's Xcode development tools. To install the Xcode command line tools open a terminal window and run:

```
xcode-select --install
```

There's a [known issue](#) that prevents `pip` from updating system packages. This has to be addressed to successfully install Scrapy and its dependencies. Here are some proposed solutions:

- *(Recommended)* **Don't** use system python, install a new, updated version that doesn't conflict with the rest of your system. Here's how to do it using the [homebrew](#) package manager:

- Install [homebrew](#) following the instructions in <https://brew.sh/>
- Update your `PATH` variable to state that homebrew packages should be used before system packages (Change `.bashrc` to `.zshrc` accordingly if you're using `zsh` as default shell):

```
echo "export PATH=/usr/local/bin:/usr/local/sbin:$PATH" >> ~/.bashrc
```

- Reload `.bashrc` to ensure the changes have taken place:

```
source ~/.bashrc
```

- Install python:

```
brew install python
```

- Latest versions of python have `pip` bundled with them so you won't need to install it separately. If this is not the case, upgrade python:

```
brew update; brew upgrade python
```

- *(Optional)* Install Scrapy inside an isolated python environment.

This method is a workaround for the above OS X issue, but it's an overall good practice for managing dependencies and can complement the first method.

[virtualenv](#) is a tool you can use to create virtual environments in python. We recommended reading a tutorial like <http://docs.python-guide.org/en/latest/dev/virtualenvs/> to get started.

After any of these workarounds you should be able to install Scrapy:

```
pip install Scrapy
```

PyPy

We recommend using the latest PyPy version. The version tested is 5.9.0. For PyPy3, only Linux installation was tested.

Most scrapy dependencies now have binary wheels for CPython, but not for PyPy. This means that these dependencies will be built during installation. On OS X, you are likely to face an issue with building

Cryptography dependency, solution to this problem is described [here](#), that is to `brew install openssl` and then export the flags that this command recommends (only needed when installing scrapy). Installing on Linux has no special issues besides installing build dependencies. Installing scrapy with PyPy on Windows is not tested.

You can check that scrapy is installed correctly by running `scrapy bench`. If this command gives errors such as `TypeError: ... got 2 unexpected keyword arguments`, this means that setuptools was unable to pick up one PyPy-specific dependency. To fix this issue, run `pip install 'PyPyDispatcher>=2.1.0'`.

2.2.3 Troubleshooting

AttributeError: 'module' object has no attribute 'OP_NO_TLSv1_1'

After you install or upgrade Scrapy, Twisted or pyOpenSSL, you may get an exception with the following traceback:

```
[...]
File "[...]/site-packages/twisted/protocols/tls.py", line 63, in <module>
    from twisted.internet.sslverify import _setAcceptableProtocols
File "[...]/site-packages/twisted/internet/_sslverify.py", line 38, in <module>
    TLSVersion.TLSv1_1: SSL.OP_NO_TLSv1_1,
AttributeError: 'module' object has no attribute 'OP_NO_TLSv1_1'
```

The reason you get this exception is that your system or virtual environment has a version of pyOpenSSL that your version of Twisted does not support.

To install a version of pyOpenSSL that your version of Twisted supports, reinstall Twisted with the `tls` extra option:

```
pip install twisted[tls]
```

For details, see [Issue #2473](#).

2.3 Scrapy 教程

在本教程中，我们假设你的系统已经安装了 Scrapy，如果没有，点击[Scrapy 安装教程](#)。

我们即将爬取一个列举了名人名言的网站，quotes.toscrape.com。

This tutorial will walk you through these tasks:

1. Creating a new Scrapy project
2. Writing a *spider* to crawl a site and extract data
3. Exporting the scraped data using the command line

4. Changing spider to recursively follow links
5. Using spider arguments

Scrapy is written in [Python](#). If you're new to the language you might want to start by getting an idea of what the language is like, to get the most out of Scrapy.

If you're already familiar with other languages, and want to learn Python quickly, the [Python Tutorial](#) is a good resource.

If you're new to programming and want to start with Python, the following books may be useful to you:

- [Automate the Boring Stuff With Python](#)
- [How To Think Like a Computer Scientist](#)
- [Learn Python 3 The Hard Way](#)

You can also take a look at [this list of Python resources for non-programmers](#), as well as the [suggested resources in the learnpython-subreddit](#).

2.3.1 Creating a project

Before you start scraping, you will have to set up a new Scrapy project. Enter a directory where you'd like to store your code and run:

```
scrapy startproject tutorial
```

This will create a `tutorial` directory with the following contents:

```
tutorial/
  scrapy.cfg          # deploy configuration file

  tutorial/           # project's Python module, you'll import your code from here
    __init__.py

  items.py            # project items definition file

  middlewares.py      # project middlewares file

  pipelines.py        # project pipelines file

  settings.py         # project settings file

  spiders/            # a directory where you'll later put your spiders
    __init__.py
```

2.3.2 Our first Spider

Spiders are classes that you define and that Scrapy uses to scrape information from a website (or a group of websites). They must subclass `scrapy.Spider` and define the initial requests to make, optionally how to follow links in the pages, and how to parse the downloaded page content to extract data.

This is the code for our first Spider. Save it in a file named `quotes_spider.py` under the `tutorial/spiders` directory in your project:

```
import scrapy

class QuotesSpider(scrapy.Spider):
    name = "quotes"

    def start_requests(self):
        urls = [
            'http://quotes.toscrape.com/page/1/',
            'http://quotes.toscrape.com/page/2/',
        ]
        for url in urls:
            yield scrapy.Request(url=url, callback=self.parse)

    def parse(self, response):
        page = response.url.split("/")[-2]
        filename = 'quotes-%s.html' % page
        with open(filename, 'wb') as f:
            f.write(response.body)
        self.log('Saved file %s' % filename)
```

As you can see, our Spider subclasses `scrapy.Spider` and defines some attributes and methods:

- *name*: identifies the Spider. It must be unique within a project, that is, you can't set the same name for different Spiders.
- *start_requests()*: must return an iterable of Requests (you can return a list of requests or write a generator function) which the Spider will begin to crawl from. Subsequent requests will be generated successively from these initial requests.
- *parse()*: a method that will be called to handle the response downloaded for each of the requests made. The response parameter is an instance of `TextResponse` that holds the page content and has further helpful methods to handle it.

The *parse()* method usually parses the response, extracting the scraped data as dicts and also finding new URLs to follow and creating new requests (*Request*) from them.

How to run our spider

To put our spider to work, go to the project's top level directory and run:

```
scrapy crawl quotes
```

This command runs the spider with name `quotes` that we've just added, that will send some requests for the `quotes.toscrape.com` domain. You will get an output similar to this:

```
... (omitted for brevity)
2016-12-16 21:24:05 [scrapy.core.engine] INFO: Spider opened
2016-12-16 21:24:05 [scrapy.extensions.logstats] INFO: Crawled 0 pages (at 0 pages/min),
↳scraped 0 items (at 0 items/min)
2016-12-16 21:24:05 [scrapy.extensions.telnet] DEBUG: Telnet console listening on 127.0.
↳0.1:6023
2016-12-16 21:24:05 [scrapy.core.engine] DEBUG: Crawled (404) <GET http://quotes.
↳toscraper.com/robots.txt> (referer: None)
2016-12-16 21:24:05 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://quotes.
↳toscraper.com/page/1/> (referer: None)
2016-12-16 21:24:05 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://quotes.
↳toscraper.com/page/2/> (referer: None)
2016-12-16 21:24:05 [quotes] DEBUG: Saved file quotes-1.html
2016-12-16 21:24:05 [quotes] DEBUG: Saved file quotes-2.html
2016-12-16 21:24:05 [scrapy.core.engine] INFO: Closing spider (finished)
...
```

Now, check the files in the current directory. You should notice that two new files have been created: `quotes-1.html` and `quotes-2.html`, with the content for the respective URLs, as our `parse` method instructs.

注解: If you are wondering why we haven't parsed the HTML yet, hold on, we will cover that soon.

What just happened under the hood?

Scrapy schedules the `scrapy.Request` objects returned by the `start_requests` method of the Spider. Upon receiving a response for each one, it instantiates `Response` objects and calls the callback method associated with the request (in this case, the `parse` method) passing the response as argument.

A shortcut to the `start_requests` method

Instead of implementing a `start_requests()` method that generates `scrapy.Request` objects from URLs, you can just define a `start_urls` class attribute with a list of URLs. This list will then be used by the

default implementation of `start_requests()` to create the initial requests for your spider:

```
import scrapy

class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
        'http://quotes.toscrape.com/page/2/',
    ]

    def parse(self, response):
        page = response.url.split("/")[-2]
        filename = 'quotes-%s.html' % page
        with open(filename, 'wb') as f:
            f.write(response.body)
```

The `parse()` method will be called to handle each of the requests for those URLs, even though we haven't explicitly told Scrapy to do so. This happens because `parse()` is Scrapy's default callback method, which is called for requests without an explicitly assigned callback.

Extracting data

The best way to learn how to extract data with Scrapy is trying selectors using the *Scrapy shell*. Run:

```
scrapy shell 'http://quotes.toscrape.com/page/1/'
```

注解: Remember to always enclose urls in quotes when running Scrapy shell from command-line, otherwise urls containing arguments (ie. `&` character) will not work.

On Windows, use double quotes instead:

```
scrapy shell "http://quotes.toscrape.com/page/1/"
```

You will see something like:

```
[ ... Scrapy log here ... ]
2016-09-19 12:09:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://quotes.
↪toscrape.com/page/1/> (referer: None)
[s] Available Scrapy objects:
```

(下页继续)

(续上页)

```
[s] scrapy      scrapy module (contains scrapy.Request, scrapy.Selector, etc)
[s] crawler     <scrapy.crawler.Crawler object at 0x7fa91d888c90>
[s] item        {}
[s] request     <GET http://quotes.toscrape.com/page/1/>
[s] response    <200 http://quotes.toscrape.com/page/1/>
[s] settings    <scrapy.settings.Settings object at 0x7fa91d888c10>
[s] spider     <DefaultSpider 'default' at 0x7fa91c8af990>
[s] Useful shortcuts:
[s] shelp()      Shell help (print this help)
[s] fetch(req_or_url) Fetch request (or URL) and update local objects
[s] view(response) View response in a browser
>>>
```

Using the shell, you can try selecting elements using CSS with the response object:

```
>>> response.css('title')
[<Selector xpath='descendant-or-self::title' data='<title>Quotes to Scrape</title>'>]
```

The result of running `response.css('title')` is a list-like object called `SelectorList`, which represents a list of `Selector` objects that wrap around XML/HTML elements and allow you to run further queries to fine-grain the selection or extract the data.

To extract the text from the title above, you can do:

```
>>> response.css('title::text').getall()
['Quotes to Scrape']
```

There are two things to note here: one is that we've added `::text` to the CSS query, to mean we want to select only the text elements directly inside `<title>` element. If we don't specify `::text`, we'd get the full title element, including its tags:

```
>>> response.css('title').getall()
['<title>Quotes to Scrape</title>']
```

The other thing is that the result of calling `.getall()` is a list: it is possible that a selector returns more than one result, so we extract them all. When you know you just want the first result, as in this case, you can do:

```
>>> response.css('title::text').get()
'Quotes to Scrape'
```

As an alternative, you could've written:

```
>>> response.css('title::text')[0].get()
'Quotes to Scrape'
```

However, using `.get()` directly on a `SelectorList` instance avoids an `IndexError` and returns `None` when it doesn't find any element matching the selection.

There's a lesson here: for most scraping code, you want it to be resilient to errors due to things not being found on a page, so that even if some parts fail to be scraped, you can at least get **some** data.

Besides the `getall()` and `get()` methods, you can also use the `re()` method to extract using [regular expressions](#):

```
>>> response.css('title::text').re(r'Quotes.*')
['Quotes to Scrape']
>>> response.css('title::text').re(r'Q\w+')
['Quotes']
>>> response.css('title::text').re(r'(\w+) to (\w+)')
['Quotes', 'Scrape']
```

In order to find the proper CSS selectors to use, you might find useful opening the response page from the shell in your web browser using `view(response)`. You can use your browser's developer tools to inspect the HTML and come up with a selector (see [使用浏览器的开发工具进行抓取](#)).

[Selector Gadget](#) is also a nice tool to quickly find CSS selector for visually selected elements, which works in many browsers.

XPath: a brief intro

Besides [CSS](#), Scrapy selectors also support using [XPath](#) expressions:

```
>>> response.xpath('//title')
[<Selector xpath='//title' data='<title>Quotes to Scrape</title>'>]
>>> response.xpath('//title/text()').get()
'Quotes to Scrape'
```

XPath expressions are very powerful, and are the foundation of Scrapy Selectors. In fact, CSS selectors are converted to XPath under-the-hood. You can see that if you read closely the text representation of the selector objects in the shell.

While perhaps not as popular as CSS selectors, XPath expressions offer more power because besides navigating the structure, it can also look at the content. Using XPath, you're able to select things like: *select the link that contains the text "Next Page"*. This makes XPath very fitting to the task of scraping, and we encourage you to learn XPath even if you already know how to construct CSS selectors, it will make scraping much easier.

We won't cover much of XPath here, but you can read more about *using XPath with Scrapy Selectors here*. To learn more about XPath, we recommend [this tutorial to learn XPath through examples](#), and [this tutorial to learn "how to think in XPath"](#) .

Extracting quotes and authors

Now that you know a bit about selection and extraction, let's complete our spider by writing the code to extract the quotes from the web page.

Each quote in <http://quotes.toscrape.com> is represented by HTML elements that look like this:

```
<div class="quote">
  <span class="text"> "The world as we have created it is a process of our
  thinking. It cannot be changed without changing our thinking." </span>
  <span>
    by <small class="author">Albert Einstein</small>
    <a href="/author/Albert-Einstein">(about)</a>
  </span>
  <div class="tags">
    Tags:
    <a class="tag" href="/tag/change/page/1/">change</a>
    <a class="tag" href="/tag/deep-thoughts/page/1/">deep-thoughts</a>
    <a class="tag" href="/tag/thinking/page/1/">thinking</a>
    <a class="tag" href="/tag/world/page/1/">world</a>
  </div>
</div>
```

Let's open up scrapy shell and play a bit to find out how to extract the data we want:

```
$ scrapy shell 'http://quotes.toscrape.com'
```

We get a list of selectors for the quote HTML elements with:

```
>>> response.css("div.quote")
```

Each of the selectors returned by the query above allows us to run further queries over their sub-elements. Let's assign the first selector to a variable, so that we can run our CSS selectors directly on a particular quote:

```
>>> quote = response.css("div.quote")[0]
```

Now, let's extract `title`, `author` and the `tags` from that quote using the `quote` object we just created:

```
>>> title = quote.css("span.text::text").get()
>>> title
' "The world as we have created it is a process of our thinking. It cannot be changed,
↳without changing our thinking." '
>>> author = quote.css("small.author::text").get()
>>> author
'Albert Einstein'
```

Given that the tags are a list of strings, we can use the `.getall()` method to get all of them:

```
>>> tags = quote.css("div.tags a.tag::text").getall()
>>> tags
['change', 'deep-thoughts', 'thinking', 'world']
```

Having figured out how to extract each bit, we can now iterate over all the quotes elements and put them together into a Python dictionary:

```
>>> for quote in response.css("div.quote"):
...     text = quote.css("span.text::text").get()
...     author = quote.css("small.author::text").get()
...     tags = quote.css("div.tags a.tag::text").getall()
...     print(dict(text=text, author=author, tags=tags))
{'tags': ['change', 'deep-thoughts', 'thinking', 'world'], 'author': 'Albert Einstein',
↳'text': ' "The world as we have created it is a process of our thinking. It cannot be
↳changed without changing our thinking." '}
{'tags': ['abilities', 'choices'], 'author': 'J.K. Rowling', 'text': ' "It is our
↳choices, Harry, that show what we truly are, far more than our abilities." '}
... a few more of these, omitted for brevity
>>>
```

Extracting data in our spider

Let's get back to our spider. Until now, it doesn't extract any data in particular, just saves the whole HTML page to a local file. Let's integrate the extraction logic above into our spider.

A Scrapy spider typically generates many dictionaries containing the data extracted from the page. To do that, we use the `yield` Python keyword in the callback, as you can see below:

```
import scrapy
```

(下页继续)

(续上页)

```
class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
        'http://quotes.toscrape.com/page/2/',
    ]

    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text': quote.css('span.text::text').get(),
                'author': quote.css('small.author::text').get(),
                'tags': quote.css('div.tags a.tag::text').getall(),
            }
```

If you run this spider, it will output the extracted data with the log:

```
2016-09-19 18:57:19 [scrapy.core.scrapers] DEBUG: Scraped from <200 http://quotes.
↳ toscrape.com/page/1/>
{'tags': ['life', 'love'], 'author': 'André Gide', 'text': ' "It is better to be hated
↳ for what you are than to be loved for what you are not." '}
2016-09-19 18:57:19 [scrapy.core.scrapers] DEBUG: Scraped from <200 http://quotes.
↳ toscrape.com/page/1/>
{'tags': ['edison', 'failure', 'inspirational', 'paraphrased'], 'author': 'Thomas A.
↳ Edison', 'text': " "I have not failed. I've just found 10,000 ways that won't work." "}
```

2.3.3 Storing the scraped data

The simplest way to store the scraped data is by using *Feed exports*, with the following command:

```
scrapy crawl quotes -o quotes.json
```

That will generate an `quotes.json` file containing all scraped items, serialized in `JSON`.

For historic reasons, Scrapy appends to a given file instead of overwriting its contents. If you run this command twice without removing the file before the second time, you'll end up with a broken JSON file.

You can also use other formats, like `JSON Lines`:

```
scrapy crawl quotes -o quotes.jl
```

The `JSON Lines` format is useful because it's stream-like, you can easily append new records to it. It

doesn't have the same problem of JSON when you run twice. Also, as each record is a separate line, you can process big files without having to fit everything in memory, there are tools like [JQ](#) to help doing that at the command-line.

In small projects (like the one in this tutorial), that should be enough. However, if you want to perform more complex things with the scraped items, you can write an *Item Pipeline*. A placeholder file for Item Pipelines has been set up for you when the project is created, in `tutorial/pipelines.py`. Though you don't need to implement any item pipelines if you just want to store the scraped items.

2.3.4 Following links

Let's say, instead of just scraping the stuff from the first two pages from <http://quotes.toscrape.com>, you want quotes from all the pages in the website.

Now that you know how to extract data from pages, let's see how to follow links from them.

First thing is to extract the link to the page we want to follow. Examining our page, we can see there is a link to the next page with the following markup:

```
<ul class="pager">
  <li class="next">
    <a href="/page/2/">Next <span aria-hidden="true">&rarr;</span></a>
  </li>
</ul>
```

We can try extracting it in the shell:

```
>>> response.css('li.next a').get()
'<a href="/page/2/">Next <span aria-hidden="true">→</span></a>'
```

This gets the anchor element, but we want the attribute `href`. For that, Scrapy supports a CSS extension that lets you select the attribute contents, like this:

```
>>> response.css('li.next a::attr(href)').get()
'/page/2/'
```

There is also an `attrib` property available (see *Selecting element attributes* for more):

```
>>> response.css('li.next a').attrib['href']
'/page/2/'
```

Let's see now our spider modified to recursively follow the link to the next page, extracting data from it:

```

import scrapy

class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
    ]

    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text': quote.css('span.text::text').get(),
                'author': quote.css('small.author::text').get(),
                'tags': quote.css('div.tags a.tag::text').getall(),
            }

        next_page = response.css('li.next a::attr(href)').get()
        if next_page is not None:
            next_page = response.urljoin(next_page)
            yield scrapy.Request(next_page, callback=self.parse)

```

Now, after extracting the data, the `parse()` method looks for the link to the next page, builds a full absolute URL using the `urljoin()` method (since the links can be relative) and yields a new request to the next page, registering itself as callback to handle the data extraction for the next page and to keep the crawling going through all the pages.

What you see here is Scrapy's mechanism of following links: when you yield a Request in a callback method, Scrapy will schedule that request to be sent and register a callback method to be executed when that request finishes.

Using this, you can build complex crawlers that follow links according to rules you define, and extract different kinds of data depending on the page it's visiting.

In our example, it creates a sort of loop, following all the links to the next page until it doesn't find one – handy for crawling blogs, forums and other sites with pagination.

A shortcut for creating Requests

As a shortcut for creating Request objects you can use `response.follow`:

```
import scrapy
```

(下页继续)

```
class QuotesSpider(scrapy.Spider):
    name = "quotes"
    start_urls = [
        'http://quotes.toscrape.com/page/1/',
    ]

    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text': quote.css('span.text::text').get(),
                'author': quote.css('span small::text').get(),
                'tags': quote.css('div.tags a.tag::text').getall(),
            }

        next_page = response.css('li.next a::attr(href)').get()
        if next_page is not None:
            yield response.follow(next_page, callback=self.parse)
```

Unlike `scrapy.Request`, `response.follow` supports relative URLs directly - no need to call `urljoin`. Note that `response.follow` just returns a `Request` instance; you still have to yield this `Request`.

You can also pass a selector to `response.follow` instead of a string; this selector should extract necessary attributes:

```
for href in response.css('li.next a::attr(href)':
    yield response.follow(href, callback=self.parse)
```

For `<a>` elements there is a shortcut: `response.follow` uses their `href` attribute automatically. So the code can be shortened further:

```
for a in response.css('li.next a'):
    yield response.follow(a, callback=self.parse)
```

注解: `response.follow(response.css('li.next a'))` is not valid because `response.css` returns a list-like object with selectors for all results, not a single selector. A `for` loop like in the example above, or `response.follow(response.css('li.next a')[0])` is fine.

More examples and patterns

Here is another spider that illustrates callbacks and following links, this time for scraping author information:

```
import scrapy

class AuthorSpider(scrapy.Spider):
    name = 'author'

    start_urls = ['http://quotes.toscrape.com/']

    def parse(self, response):
        # follow links to author pages
        for href in response.css('.author + a::attr(href)'):
            yield response.follow(href, self.parse_author)

        # follow pagination links
        for href in response.css('li.next a::attr(href)'):
            yield response.follow(href, self.parse)

    def parse_author(self, response):
        def extract_with_css(query):
            return response.css(query).get(default='').strip()

        yield {
            'name': extract_with_css('h3.author-title::text'),
            'birthdate': extract_with_css('.author-born-date::text'),
            'bio': extract_with_css('.author-description::text'),
        }
```

This spider will start from the main page, it will follow all the links to the authors pages calling the `parse_author` callback for each of them, and also the pagination links with the `parse` callback as we saw before.

Here we're passing callbacks to `response.follow` as positional arguments to make the code shorter; it also works for `scrapy.Request`.

The `parse_author` callback defines a helper function to extract and cleanup the data from a CSS query and yields the Python dict with the author data.

Another interesting thing this spider demonstrates is that, even if there are many quotes from the same author, we don't need to worry about visiting the same author page multiple times. By default, Scrapy filters out duplicated requests to URLs already visited, avoiding the problem of hitting servers too much

because of a programming mistake. This can be configured by the setting `DUPEFILTER_CLASS`.

Hopefully by now you have a good understanding of how to use the mechanism of following links and callbacks with Scrapy.

As yet another example spider that leverages the mechanism of following links, check out the `CrawlSpider` class for a generic spider that implements a small rules engine that you can use to write your crawlers on top of it.

Also, a common pattern is to build an item with data from more than one page, using a *trick to pass additional data to the callbacks*.

2.3.5 Using spider arguments

You can provide command line arguments to your spiders by using the `-a` option when running them:

```
scrapy crawl quotes -o quotes-humor.json -a tag=humor
```

These arguments are passed to the Spider's `__init__` method and become spider attributes by default.

In this example, the value provided for the `tag` argument will be available via `self.tag`. You can use this to make your spider fetch only quotes with a specific tag, building the URL based on the argument:

```
import scrapy

class QuotesSpider(scrapy.Spider):
    name = "quotes"

    def start_requests(self):
        url = 'http://quotes.toscrape.com/'
        tag = getattr(self, 'tag', None)
        if tag is not None:
            url = url + 'tag/' + tag
        yield scrapy.Request(url, self.parse)

    def parse(self, response):
        for quote in response.css('div.quote'):
            yield {
                'text': quote.css('span.text::text').get(),
                'author': quote.css('small.author::text').get(),
            }

        next_page = response.css('li.next a::attr(href)').get()
```

(下页继续)

(续上页)

```
if next_page is not None:
    yield response.follow(next_page, self.parse)
```

If you pass the `tag=humor` argument to this spider, you'll notice that it will only visit URLs from the `humor` tag, such as <http://quotes.toscrape.com/tag/humor>.

You can *learn more about handling spider arguments here*.

2.3.6 Next steps

This tutorial covered only the basics of Scrapy, but there's a lot of other features not mentioned here. Check the [还有什么?](#) section in [预览 Scrapy](#) chapter for a quick overview of the most important ones.

You can continue from the section [基本概念](#) to know more about the command-line tool, spiders, selectors and other things the tutorial hasn't covered like modeling the scraped data. If you prefer to play with an example project, check the [例子](#) section.

2.4 例子

The best way to learn is with examples, and Scrapy is no exception. For this reason, there is an example Scrapy project named `quotesbot`, that you can use to play and learn more about Scrapy. It contains two spiders for <http://quotes.toscrape.com>, one using CSS selectors and another one using XPath expressions.

The `quotesbot` project is available at: <https://github.com/scrapy/quotesbot>. You can find more information about it in the project's README.

If you're familiar with git, you can checkout the code. Otherwise you can download the project as a zip file by clicking [here](#).

[预览 Scrapy](#) 了解什么是 Scrapy, 它是如何帮助你的。

[安装指南](#) 在你的电脑上安装 Scrapy。

[Scrapy 教程](#) 编写你的第一个 Scrapy 项目。

[例子](#) 通过运行一个内置的 Scrapy 例程进一步学习。

3.1 命令行工具

0.10 新版功能.

Scrapy is controlled through the **scrapy** command-line tool, to be referred here as the “Scrapy tool” to differentiate it from the sub-commands, which we just call “commands” or “Scrapy commands” .

The Scrapy tool provides several commands, for multiple purposes, and each one accepts a different set of arguments and options.

(The **scrapy deploy** command has been removed in 1.0 in favor of the standalone **scrapyd-deploy**. See [Deploying your project](#).)

3.1.1 Configuration settings

Scrapy will look for configuration parameters in ini-style **scrapy.cfg** files in standard locations:

1. **/etc/scrapy.cfg** or **c:\scrapy\scrapy.cfg** (system-wide),
2. **~/.config/scrapy.cfg** (**\$XDG_CONFIG_HOME**) and **~/.scrapy.cfg** (**\$HOME**) for global (user-wide) settings, and
3. **scrapy.cfg** inside a scrapy project’ s root (see next section).

Settings from these files are merged in the listed order of preference: user-defined values have higher priority than system-wide defaults and project-wide settings will override all others, when defined.

Scrapy also understands, and can be configured through, a number of environment variables. Currently these are:

- `SCRAPY_SETTINGS_MODULE` (see *Designating the settings*)
- `SCRAPY_PROJECT` (see *Sharing the root directory between projects*)
- `SCRAPY_PYTHON_SHELL` (see *Scrapy shell*)

3.1.2 Default structure of Scrapy projects

Before delving into the command-line tool and its sub-commands, let's first understand the directory structure of a Scrapy project.

Though it can be modified, all Scrapy projects have the same file structure by default, similar to this:

```
scrapy.cfg
myproject/
  __init__.py
  items.py
  middlewares.py
  pipelines.py
  settings.py
  spiders/
    __init__.py
    spider1.py
    spider2.py
    ...
```

The directory where the `scrapy.cfg` file resides is known as the *project root directory*. That file contains the name of the python module that defines the project settings. Here is an example:

```
[settings]
default = myproject.settings
```

3.1.3 Sharing the root directory between projects

A project root directory, the one that contains the `scrapy.cfg`, may be shared by multiple Scrapy projects, each with its own settings module.

In that case, you must define one or more aliases for those settings modules under `[settings]` in your `scrapy.cfg` file:

```
[settings]
default = myproject1.settings
project1 = myproject1.settings
project2 = myproject2.settings
```

By default, the `scrapy` command-line tool will use the `default` settings. Use the `SCRAPY_PROJECT` environment variable to specify a different project for `scrapy` to use:

```
$ scrapy settings --get BOT_NAME
Project 1 Bot
$ export SCRAPY_PROJECT=project2
$ scrapy settings --get BOT_NAME
Project 2 Bot
```

3.1.4 Using the scrapy tool

You can start by running the Scrapy tool with no arguments and it will print some usage help and the available commands:

```
Scrapy X.Y - no active project

Usage:
  scrapy <command> [options] [args]

Available commands:
  crawl          Run a spider
  fetch          Fetch a URL using the Scrapy downloader
  [...]

```

The first line will print the currently active project if you're inside a Scrapy project. In this example it was run from outside a project. If run from inside a project it would have printed something like this:

```
Scrapy X.Y - project: myproject

Usage:
  scrapy <command> [options] [args]

[...]

```

Creating projects

The first thing you typically do with the **scrapy** tool is create your Scrapy project:

```
scrapy startproject myproject [project_dir]
```

That will create a Scrapy project under the `project_dir` directory. If `project_dir` wasn't specified, `project_dir` will be the same as `myproject`.

Next, you go inside the new project directory:

```
cd project_dir
```

And you're ready to use the **scrapy** command to manage and control your project from there.

Controlling projects

You use the **scrapy** tool from inside your projects to control and manage them.

For example, to create a new spider:

```
scrapy genspider mydomain mydomain.com
```

Some Scrapy commands (like *[crawl](#)*) must be run from inside a Scrapy project. See the *[commands reference](#)* below for more information on which commands must be run from inside projects, and which not.

Also keep in mind that some commands may have slightly different behaviours when running them from inside projects. For example, the `fetch` command will use spider-overridden behaviours (such as the `user_agent` attribute to override the user-agent) if the url being fetched is associated with some specific spider. This is intentional, as the `fetch` command is meant to be used to check how spiders are downloading pages.

3.1.5 Available tool commands

This section contains a list of the available built-in commands with a description and some usage examples. Remember, you can always get more info about each command by running:

```
scrapy <command> -h
```

And you can see all available commands with:

```
scrapy -h
```

There are two kinds of commands, those that only work from inside a Scrapy project (Project-specific commands) and those that also work without an active Scrapy project (Global commands), though they

may behave slightly different when running from inside a project (as they would use the project overridden settings).

Global commands:

- *startproject*
- *genspider*
- *settings*
- *runspider*
- *shell*
- *fetch*
- *view*
- *version*

Project-only commands:

- *crawl*
- *check*
- *list*
- *edit*
- *parse*
- *bench*

startproject

- Syntax: `scrapy startproject <project_name> [project_dir]`
- Requires project: *no*

Creates a new Scrapy project named `project_name`, under the `project_dir` directory. If `project_dir` wasn't specified, `project_dir` will be the same as `project_name`.

Usage example:

```
$ scrapy startproject myproject
```

genspider

- Syntax: `scrapy genspider [-t template] <name> <domain>`
- Requires project: *no*

Create a new spider in the current folder or in the current project's `spiders` folder, if called from inside a project. The `<name>` parameter is set as the spider's `name`, while `<domain>` is used to generate the `allowed_domains` and `start_urls` spider's attributes.

Usage example:

```
$ scrapy genspider -l
Available templates:
    basic
    crawl
    csvfeed
    xmlfeed

$ scrapy genspider example example.com
Created spider 'example' using template 'basic'

$ scrapy genspider -t crawl scrapyorg scrapy.org
Created spider 'scrapyorg' using template 'crawl'
```

This is just a convenience shortcut command for creating spiders based on pre-defined templates, but certainly not the only way to create spiders. You can just create the spider source code files yourself, instead of using this command.

crawl

- Syntax: `scrapy crawl <spider>`
- Requires project: *yes*

Start crawling using a spider.

Usage examples:

```
$ scrapy crawl myspider
[ ... myspider starts crawling ... ]
```

check

- Syntax: `scrapy check [-l] <spider>`
- Requires project: *yes*

Run contract checks.

Usage examples:

```
$ scrapy check -l
first_spider
  * parse
  * parse_item
second_spider
  * parse
  * parse_item

$ scrapy check
[FAILED] first_spider:parse_item
>>> 'RetailPricex' field is missing

[FAILED] first_spider:parse
>>> Returned 92 requests, expected 0..4
```

list

- Syntax: `scrapy list`
- Requires project: *yes*

List all available spiders in the current project. The output is one spider per line.

Usage example:

```
$ scrapy list
spider1
spider2
```

edit

- Syntax: `scrapy edit <spider>`
- Requires project: *yes*

Edit the given spider using the editor defined in the `EDITOR` environment variable or (if unset) the *EDITOR* setting.

This command is provided only as a convenience shortcut for the most common case, the developer is of course free to choose any tool or IDE to write and debug spiders.

Usage example:

```
$ scrapy edit spider1
```

fetch

- Syntax: `scrapy fetch <url>`
- Requires project: *no*

Downloads the given URL using the Scrapy downloader and writes the contents to standard output.

The interesting thing about this command is that it fetches the page how the spider would download it. For example, if the spider has a `USER_AGENT` attribute which overrides the User Agent, it will use that one.

So this command can be used to “see” how your spider would fetch a certain page.

If used outside a project, no particular per-spider behaviour would be applied and it will just use the default Scrapy downloader settings.

Supported options:

- `--spider=SPIDER`: bypass spider autodetection and force use of specific spider
- `--headers`: print the response’ s HTTP headers instead of the response’ s body
- `--no-redirect`: do not follow HTTP 3xx redirects (default is to follow them)

Usage examples:

```
$ scrapy fetch --nolog http://www.example.com/some/page.html
[ ... html content here ... ]

$ scrapy fetch --nolog --headers http://www.example.com/
{'Accept-Ranges': ['bytes'],
 'Age': ['1263    '],
 'Connection': ['close    '],
 'Content-Length': ['596'],
 'Content-Type': ['text/html; charset=UTF-8'],
 'Date': ['Wed, 18 Aug 2010 23:59:46 GMT'],
 'Etag': ['"573c1-254-48c9c87349680"'],
 'Last-Modified': ['Fri, 30 Jul 2010 15:30:18 GMT'],
 'Server': ['Apache/2.2.3 (CentOS)']}
```

view

- Syntax: `scrapy view <url>`
- Requires project: *no*

Opens the given URL in a browser, as your Scrapy spider would “see” it. Sometimes spiders see pages differently from regular users, so this can be used to check what the spider “sees” and confirm it’ s what you expect.

Supported options:

- `--spider=SPIDER`: bypass spider autodetection and force use of specific spider
- `--no-redirect`: do not follow HTTP 3xx redirects (default is to follow them)

Usage example:

```
$ scrapy view http://www.example.com/some/page.html
[ ... browser starts ... ]
```

shell

- Syntax: `scrapy shell [url]`
- Requires project: *no*

Starts the Scrapy shell for the given URL (if given) or empty if no URL is given. Also supports UNIX-style local file paths, either relative with `./` or `../` prefixes or absolute file paths. See [Scrapy shell](#) for more info.

Supported options:

- `--spider=SPIDER`: bypass spider autodetection and force use of specific spider
- `-c code`: evaluate the code in the shell, print the result and exit
- `--no-redirect`: do not follow HTTP 3xx redirects (default is to follow them); this only affects the URL you may pass as argument on the command line; once you are inside the shell, `fetch(url)` will still follow HTTP redirects by default.

Usage example:

```
$ scrapy shell http://www.example.com/some/page.html
[ ... scrapy shell starts ... ]

$ scrapy shell --nolog http://www.example.com/ -c '(response.status, response.url)'
(200, 'http://www.example.com/')

# shell follows HTTP redirects by default
$ scrapy shell --nolog http://httpbin.org/redirect-to?url=http%3A%2F%2Fexample.com%2F -c
↪ '(response.status, response.url)'
(200, 'http://example.com/')

# you can disable this with --no-redirect
# (only for the URL passed as command line argument)
```

(下页继续)

(续上页)

```
$ scrapy shell --no-redirect --nolog http://httpbin.org/redirect-to?url=http%3A%2F
%2Fexample.com%2F -c '(response.status, response.url)'
(302, 'http://httpbin.org/redirect-to?url=http%3A%2F%2Fexample.com%2F')
```

parse

- Syntax: `scrapy parse <url> [options]`
- Requires project: *yes*

Fetches the given URL and parses it with the spider that handles it, using the method passed with the `--callback` option, or `parse` if not given.

Supported options:

- `--spider=SPIDER`: bypass spider autodetection and force use of specific spider
- `--a NAME=VALUE`: set spider argument (may be repeated)
- `--callback` or `-c`: spider method to use as callback for parsing the response
- `--meta` or `-m`: additional request meta that will be passed to the callback request. This must be a valid json string. Example: `-meta= ' { "foo" : "bar" } '`
- `--pipelines`: process items through pipelines
- `--rules` or `-r`: use *CrawlSpider* rules to discover the callback (i.e. spider method) to use for parsing the response
- `--noitems`: don't show scraped items
- `--nolinks`: don't show extracted links
- `--nocolour`: avoid using pygments to colorize the output
- `--depth` or `-d`: depth level for which the requests should be followed recursively (default: 1)
- `--verbose` or `-v`: display information for each depth level

Usage example:

```
$ scrapy parse http://www.example.com/ -c parse_item
[ ... scrapy log lines crawling example.com spider ... ]

>>> STATUS DEPTH LEVEL 1 <<<
# Scraped Items -----
[{'name': 'Example item',
  'category': 'Furniture',
  'length': '12 cm'}]
```

(下页继续)

(续上页)

```
# Requests -----  
[]
```

settings

- Syntax: `scrapy settings [options]`
- Requires project: *no*

Get the value of a Scrapy setting.

If used inside a project it'll show the project setting value, otherwise it'll show the default Scrapy value for that setting.

Example usage:

```
$ scrapy settings --get BOT_NAME  
scrapybot  
$ scrapy settings --get DOWNLOAD_DELAY  
0
```

runspider

- Syntax: `scrapy runspider <spider_file.py>`
- Requires project: *no*

Run a spider self-contained in a Python file, without having to create a project.

Example usage:

```
$ scrapy runspider myspider.py  
[ ... spider starts crawling ... ]
```

version

- Syntax: `scrapy version [-v]`
- Requires project: *no*

Prints the Scrapy version. If used with `-v` it also prints Python, Twisted and Platform info, which is useful for bug reports.

bench

0.17 新版功能.

- Syntax: `scrapy bench`
- Requires project: *no*

Run a quick benchmark test. 爬虫器硬件性能.

3.1.6 Custom project commands

You can also add your custom project commands by using the `COMMANDS_MODULE` setting. See the Scrapy commands in `scrapy/commands` for examples on how to implement your commands.

COMMANDS_MODULE

Default: `''` (empty string)

A module to use for looking up custom Scrapy commands. This is used to add custom commands for your Scrapy project.

Example:

```
COMMANDS_MODULE = 'mybot.commands'
```

Register commands via `setup.py` entry points

注解: This is an experimental feature, use with caution.

You can also add Scrapy commands from an external library by adding a `scrapy.commands` section in the entry points of the library `setup.py` file.

The following example adds `my_command` command:

```
from setuptools import setup, find_packages

setup(name='scrapy-mymodule',
      entry_points={
        'scrapy.commands': [
          'my_command=my_scrapy_module.commands.MyCommand',
        ],
      },
```

(下页继续)

(续上页)

```
} ,  
)
```

3.2 爬虫器

Spiders 类定义了如何爬取一个站点（或者一组站点），包括如何执行抓取（即跟踪链接）以及如何从爬取的页面中提取结构化的数据（即爬取 items）。换句话说，Spiders 定义了爬取和解析特定站点（或者一组站点）的行为。

对于爬虫器来说，抓取的过程类似这样运行的：

1. 从放入的一批初始 url 开始抓取，然后指定一个回调函数来处理这个（这些）请求的响应。

第一批请求是调用 `start_requests()` 方法，该方法（默认）为 `start_urls` 中指定的 URLs 生成 *Request* 类，并且 `parse` 方法作为这些请求的回调函数。

2. 在回调函数中，你解析了响应（网页）并返回了提取的（字典化）的数据、*Item* 对象、或者这些对象的迭代器其中之王。这些请求也可以包含一个回调（也许是一个回调）并被 Scrapy 下载然后它们的响应被指定的回调函数处理。
3. 在回调函数中，你解析了页面内容，主要用了 *选择器*（但是你可以用 BeautifulSoup、lxml 或者你喜欢的其它方法）并用这些解析的数据生成 items。
4. 最后，这些从爬虫器返回的 items 会被持久化到数据库（在一些 *Item Pipeline*）或者放入一个用 *Feed 导出* 导出的文件。

尽管这个过程已经适用（或多或少）一些爬虫，不过 Scrapy 还附带了一些其它的爬虫器用于不同的爬取情况。下面我们来了解一下它们。

3.2.1 scrapy.Spider

```
class scrapy.spiders.Spider
```

最简单的爬虫器，其它所有的爬虫器都继承了它（包括 Scrapy 附带的爬虫器，以及你自定义的爬虫器）。它没有提供任何特殊的功能。只有一个默认的 `start_requests()` 启动方法用于发送 `start_urls` 属性中的请求并把返回的响应作为参数调用爬虫器的 `parse` 方法。

name

定义了爬虫器的名字，Scrapy 通过爬虫器的名字来识别不同的爬虫器，所以必需提供。然而，你当然可以定义多个一样的爬虫器，这是一个爬虫器必需提供的属性。

如果爬虫器抓取一个单独的站点，通常的做法是根据域名定义爬虫器 `name`，带或者不带 **TLD**。比如一个爬取 `mywebsite.com` 爬虫器可以被称作 `mywebsite`。

注解: 在 Python2 中只能是 ASCII 码。

`allowed_domains`

一个可选的列表，包含了爬虫器可以抓取的域名。这个列表指定的域名（包含子域名）将不会被请求如果 *OffsiteMiddleware* 被启用。

假设目标 url 是 `https://www.example.com/1.html`，那么就添加 `'example.com'` 到这个列表中。

`start_urls`

一个 URL 列表，如果没有定义特殊的 URL，爬虫器将会从这里开始爬取。所以，第一批页面被从这个列表开始下载。后来的 *Request* 会从这个 url 列表中的数据里相继生成。

`custom_settings`

字典类型的设置配置，当爬虫启动的时候会被项目中的配置重写。它必需作为一个类属性被定义，然后设置会在实例化前更新。

可用的内置设置列表详见: *Built-in settings reference*。

`crawler`

该属性会在 *from_crawler()* 初始化类后被类方法设置，然后与这个爬虫器的实例绑定的 *Crawler* 对象关联。

爬虫程序在项目入口中封装了许多组件 (比如 extensions, middlewares, signals managers, etc). 了解更多详见 *Crawler API* 。

`settings`

运行爬虫器的配置。是一个 *Settings* 实例，从 *设置* 主题了解这个主题的配置。

`logger`

Python 的日志用爬虫器的 *name* 创建。你可以按照 *Logging from Spiders* 里的描述来用它发送日志信息。

`from_crawler(crawler, *args, **kwargs)`

Scrapy 用这个类方法创建你的爬虫器。

你可能不需要直接覆盖它，因为默认的启动行为作为 `__init__()` 方法的一个代理，通过传递 `args` 和 `kwargs` 参数来调用它。

尽管如此，这个方法在新的实例中设置了 *crawler* 和 *settings* 属性以便在后来的爬虫程序中可以被使用。

参数

- **crawler** (*Crawler* 实例) – crawler to which the spider will be bound
- **args** (列表) – `__init__()` 方法传递过来的参数
- **kwargs** (字典) – `__init__()` 方法传递过来的关键字参数。

start_requests()

此方法必需返回一个迭代器，其中包含这个爬虫器要爬取的第一个请求。当爬虫器开始爬虫的时候它被 Scrapy 调用。Scrapy 只会调用它一次，所以作为生成器执行 `start_requests()` 是安全的。

对于每个 `start_urls` 里的 url 默认使用 `Request(url, dont_filter=True)` 执行的。

如果你想要修改最初爬取某个网站的 Request 对象，你可以重写这个方法。比如，你需要在登录时使用 POST 请求，你可以这么写：

```
class MySpider(scrapy.Spider):
    name = 'myspider'

    def start_requests(self):
        return [scrapy.FormRequest("http://www.example.com/login",
                                   formdata={'user': 'john', 'pass': 'secret'},
                                   callback=self.logged_in)]

    def logged_in(self, response):
        # here you would extract links to follow and return Requests for
        # each of them, with another callback
        pass
```

parse(response)

当没有指定回调函数时，Scrapy 默认使用该方法作为回调函数来处理响应。

`parse` 方法负责处理响应并返回抓取到的数据以及/或者跟进的 URL。`Spider` 对其他的 Request 的回调函数也有相同的要求。

此方法以及其他回调函数必需返回一个包含 `Request`、字典或者 `Item` 的可迭代对象。

参数 response (`Response`) – 用于解析 response

log(message[, level, component])

使用爬虫器的 `logger` 发送日志消息的装饰器，保持身后的兼容性。更多内容请参考 *Logging from Spiders*。

closed(reason)

当爬虫器关闭时调用。该方法提供了一个替代调用 `signals.connect()` 来监听 `spider_closed` 信号的快捷方式。

让我们来看一个例子：

```
import scrapy
```

(下页继续)

(续上页)

```
class MySpider(scrapy.Spider):
    name = 'example.com'
    allowed_domains = ['example.com']
    start_urls = [
        'http://www.example.com/1.html',
        'http://www.example.com/2.html',
        'http://www.example.com/3.html',
    ]

    def parse(self, response):
        self.logger.info('A response from %s just arrived!', response.url)
```

R 在单个回调函数中返回多个 Request 以及 Item 的例子:

```
import scrapy

class MySpider(scrapy.Spider):
    name = 'example.com'
    allowed_domains = ['example.com']
    start_urls = [
        'http://www.example.com/1.html',
        'http://www.example.com/2.html',
        'http://www.example.com/3.html',
    ]

    def parse(self, response):
        for h3 in response.xpath('//h3').getall():
            yield {"title": h3}

        for href in response.xpath('//a/@href').getall():
            yield scrapy.Request(response.urljoin(href), self.parse)
```

除了 *start_urls* , 你也可以直接使用 *start_requests()* ; 也可以使用 *Items* 来给予数据更多的结构性 (give data more structure):

```
import scrapy
from myproject.items import MyItem

class MySpider(scrapy.Spider):
    name = 'example.com'
    allowed_domains = ['example.com']
```

(下页继续)

(续上页)

```
def start_requests(self):
    yield scrapy.Request('http://www.example.com/1.html', self.parse)
    yield scrapy.Request('http://www.example.com/2.html', self.parse)
    yield scrapy.Request('http://www.example.com/3.html', self.parse)

def parse(self, response):
    for h3 in response.xpath('//h3').getall():
        yield MyItem(title=h3)

    for href in response.xpath('//a/@href').getall():
        yield scrapy.Request(response.urljoin(href), self.parse)
```

3.2.2 爬虫器参数

爬虫器可以接收一些参数用来改变它的行为。spider 参数的一些常见用途是定义开始 url 或将限制爬取站点的某些部分, 但是它们可以用来配置爬行器的任何功能。

爬虫顺口参数通过 *crawl* 命令使用 *-a* 选项传递。比如:

```
scrapy crawl myspider -a category=electronics
```

爬虫器可以在 `__init__` 方法中获取参数:

```
import scrapy

class MySpider(scrapy.Spider):
    name = 'myspider'

    def __init__(self, category=None, *args, **kwargs):
        super(MySpider, self).__init__(*args, **kwargs)
        self.start_urls = ['http://www.example.com/categories/%s' % category]
        # ...
```

默认的 `__init__` 方法将接收任何爬虫器的参数并自制它们到爬虫器作为它们的属性。上面的例子也可以写成如下:

```
import scrapy

class MySpider(scrapy.Spider):
    name = 'myspider'
```

(下页继续)

(续上页)

```
def start_requests(self):  
    yield scrapy.Request('http://www.example.com/categories/%s' % self.category)
```

请记住, spider 参数只是字符串。爬行器不会自己进行任何解析。如果要从命令行设置 `start_urls` 属性, 你必须使用类似 `ast.literal_eval` 或者 `json.loads` 的方法将其解析为一个列表并将其设置为属性。否则, 你将导致对 `start_urls` 字符串的迭代 (一个非常常见的 python 陷阱) 导致每个字符被视为一个单独的 url。

一个有效的用例是使用 `HttpAuthMiddleware` 设置 http 用户证书。或者使用 `UserAgentMiddleware` 设置用户代理:

```
scrapy crawl myspider -a http_user=myuser -a http_pass=mypassword -a user_agent=mybot
```

爬虫器参数也可以使用 Scrapyd `schedule.json` API 传递。查看 [Scrapyd documentation](#)。

3.2.3 常见的爬虫器

Scrapy 附带了一些有用的通用爬行器, 您可以使用它们对爬行器进行子类化。他们的目标是为一些常见的抓取情况提供方便的功能, 比如根据特定的规则跟进站点上的所有链接,

从 `Sitemaps` 爬取, 或者解析成 XML/CSV feed.

例如使用下面的爬虫器, 我们假设你已经新建了一个项目并在 `myproject.items` 文件明声明了一个 `TestItem`

```
import scrapy  
  
class TestItem(scrapy.Item):  
    id = scrapy.Field()  
    name = scrapy.Field()  
    description = scrapy.Field()
```

CrawlSpider

```
class scrapy.spiders.CrawlSpider
```

爬取一般网站常用的爬虫器。它有一个方便的机制用来跟进一些规则定义的 links。也许该爬虫器并不是完全适合你的特定网站或项目, 但其对很多情况都使用, 因此你可以以其为起点, 根据需求修改部分方法。当然您也可以实现自己的爬虫器。

除了从 Spider 继承过来的 (您必须提供的) 属性外, 其提供了一个新的属性:

rules

一个包含一个 (或多个) `Rule` 对象的集合。每个 `Rule` 都定义了某些行为来爬取网站。Rule 对象在

下边会介绍。如果多个 rule 匹配了相同的链接，则根据他们在本属性中被定义的顺序，第一个会被使用。

该爬虫器也提供了一个可复写的方法：

```
.. method:: parse_start_url(response)
```

当 start_url 的请求返回时，该方法被调用。该方法解析最初的返回值并必须返回一个 *Item* 对象或者一个 *Request* 对象，或者是一个包含它们任何一个的可迭代对象。

爬取规则 (Crawling rules)

```
class scrapy.spiders.Rule(link_extractor, callback=None, cb_kwargs=None, follow=None, process_links=None, process_request=None)
```

link_extractor 是一个 *Link Extractor* 对象。其定义了如何从爬取到的页面提取链接。

callback 是一个 callable 或 string(该 spider 中同名的函数将会被调用)。从 link_extractor 中每获取到链接时将会调用该函数。该回调函数接受一个 response 作为其第一个参数，并返回一个包含 *Item* 以及 (或) *Request* Request 对象 (或者这两者的子类) 的列表 (list)。

警告： 当编写爬虫规则时，请避免使用 parse 作为回调函数。由于 *CrawlSpider* 使用 parse 方法来实现其逻辑，如果你覆盖了 parse 方法，爬虫器将会运行失败。

cb_kwargs 包含传递给回调函数的参数 (keyword argument) 的字典。

follow 是一个布尔 (boolean) 值，指定了根据该规则从 response 提取的链接是否需要跟进。如果 callback 为 None follow 默认为 True，否则默认为 False。

process_links 是一个 callable 或 string(该 spider 中同名的函数将会被调用)。从 link_extractor 中获取链接列表时将会调用该函数。该方法主要用来过滤。

process_request 是一个 callable 或 string(该 spider 中同名的函数将会被调用)。该规则提取到每个 *Request* 时都会调用该函数。该函数必须返回一个 *Request* 对象或者 None。(用来过滤 request)

CrawlSpider 例子

接下来给出配合 rule 使用 CrawlSpider 的例子：

```
import scrapy
from scrapy.spiders import CrawlSpider, Rule
from scrapy.linkextractors import LinkExtractor

class MySpider(CrawlSpider):
```

(下页继续)

```

name = 'example.com'
allowed_domains = ['example.com']
start_urls = ['http://www.example.com']

rules = (
    # Extract links matching 'category.php' (but not matching 'subsection.php')
    # and follow links from them (since no callback means follow=True by default).
    Rule(LinkExtractor(allow=('category\.php', ), deny=('subsection\.php', ))),

    # Extract links matching 'item.php' and parse them with the spider's method parse_item
    Rule(LinkExtractor(allow=('item\.php', )), callback='parse_item'),
)

def parse_item(self, response):
    self.logger.info('Hi, this is an item page! %s', response.url)
    item = scrapy.Item()
    item['id'] = response.xpath('//td[@id="item_id"]/text()').re(r'ID: (\d+)')
    item['name'] = response.xpath('//td[@id="item_name"]/text()').get()
    item['description'] = response.xpath('//td[@id="item_description"]/text()').get()
    return item

```

这个爬虫器将开始爬取 example.com 的首页，获取 category 以及 item 的链接并对后者使用“parse_item”方法。当 item 获得返回 (response) 时，将使用 XPath 处理 HTML 并生成一些数据填入 class: scrapy.item.Item 中。

XMLFeedSpider

```
class scrapy.spiders.XMLFeedSpider
```

XMLFeedSpider 被设计用于通过迭代各个节点来分析 XML feed。迭代器可以从: iternodes, xml 和 html 中选择。鉴于 xml 和 html 迭代器需要先读取所有 DOM 再分析而引起的性能问题。一般还是推荐使用 iternodes。不过使用 html 作为迭代器能有效应对错误的 XML。

您必须定义下列类属性来设置迭代器以及 tag name:

iterator

用于确定使用哪个迭代器的 string。可选项有:

- 'iternodes' 一个高性能的基于正则表达式的迭代器
- 'html' - 使用 Selector 的迭代器。需要注意的是该迭代器使用 DOM 进行分析，其需要将所有的 DOM 载入内存，当数据量大的时候会产生问题。

- 'xml' - 使用 `Selector` 的迭代器。需要注意的是该迭代器使用 DOM 进行分析，其需要将所有的 DOM 载入内存，当数据量大的时候会产生问题。

默认值为: 'iternodes'.

itertag

一个包含开始迭代的节点名的字符串。例如:

```
itertag = 'product'
```

namespaces

一个由 (prefix, uri) 元组所组成的列表。其定义了在该文档中会被爬虫器处理的可用的 namespace。prefix 和 uri 会被自动调用 `register_namespace()` 方法生成 namespace。

您可以通过在 *itertag* 属性中指定节点的 namespace。

例子:

```
class YourSpider(XMLFeedSpider):

    namespaces = [('n', 'http://www.sitemaps.org/schemas/sitemap/0.9')]
    itertag = 'n:url'
    # ...
```

除了这些新的属性之外，该 spider 也有以下可以覆盖 (overrideable) 的方法:

adapt_response(response)

该方法在 spider 分析 response 前被调用。您可以在 response 被分析之前使用该函数来修改内容 (body)。该方法接受一个 response 并返回一个 response(可以相同也可以不同)。

parse_node(response, selector)

当节点符合提供的标签名时 (itertag)。接收到的 response 以及相应的 `Selector` 作为参数传递给该方法。该方法返回一个 *Item* 对象或者 *Request* 对象或者一个包含二者的可迭代对象 (iterable)。

process_results(response, results)

当 spider 返回结果 (item 或 request) 时该方法被调用。设定该方法的目的是在结果返回给框架核心 (framework core) 之前做最后的处理，例如设定 item 的 ID。其接受一个结果的列表 (list of results) 及对应的 response。其结果必须返回一个结果的列表 (list of results)(包含 Item 或者 Request 对象)。

XMLFeedSpider 例子

这些爬虫器用起来都非常简单，让我们看一个例子:

```
from scrapy.spiders import XMLFeedSpider
from myproject.items import TestItem
```

(下页继续)

```

class MySpider(XMLFeedSpider):
    name = 'example.com'
    allowed_domains = ['example.com']
    start_urls = ['http://www.example.com/feed.xml']
    iterator = 'iternodes' # This is actually unnecessary, since it's the default value
    itertag = 'item'

    def parse_node(self, response, node):
        self.logger.info('Hi, this is a <%s> node!: %s', self.itertag, ''.join(node.
↪getall()))

        item = TestItem()
        item['id'] = node.xpath('@id').get()
        item['name'] = node.xpath('name').get()
        item['description'] = node.xpath('description').get()
        return item

```

基本上，我们在上面所做的就是创建一个爬虫器用来从给定的 `start_urls` 下载 feed。并迭代 feed 中每个 `item` 标签。输出，并在 `Item` 中存储有些随机数据。

CSVFeedSpider

```
class scrapy.spiders.CSVFeedSpider
```

该 spider 除了其按行遍历而不是节点之外其他和 XMLFeedSpider 十分类似。而其在每次迭代时调用的是 `parse_row()`。

delimiter

在 CSV 文件中用于区分字段的分隔符。类型为 string 默认是 ',' (逗号)。

quotechar

A string with the enclosure character for each field in the CSV file Defaults to '"' (quotation mark).

headers

在 CSV 文件中包含的用来提取字段的行的列表。

parse_row(response, row)

接收一个 response 对象及一个以提供或检测出来的 header 为键的字典从一个 CSV 文件中。也可以覆盖 `adapt_response` 和 `process_results` 方法来进行预处理 (pre-processing) 及后 (post-processing) 处理。

CSVFeedSpider example

下面的例子和之前的例子很像，但使用了 *CSVFeedSpider*:

```
from scrapy.spiders import CSVFeedSpider
from myproject.items import TestItem

class MySpider(CSVFeedSpider):
    name = 'example.com'
    allowed_domains = ['example.com']
    start_urls = ['http://www.example.com/feed.csv']
    delimiter = ';'
    quotechar = '"'
    headers = ['id', 'name', 'description']

    def parse_row(self, response, row):
        self.logger.info('Hi, this is a row!: %r', row)

        item = TestItem()
        item['id'] = row['id']
        item['name'] = row['name']
        item['description'] = row['description']
        return item
```

SitemapSpider

```
class scrapy.spiders.SitemapSpider
```

SitemapSpider 使您爬取网站时可以通过 Sitemaps 来发现爬取的 URL。

其支持嵌套的 sitemap，并能从 robots.txt 中获取 sitemap 的 url。

sitemap_urls

包含您要爬取的 url 的 sitemap 的 url 列表 (list)。

你也可以指定一个 robots.txt 文件，爬虫器会从中分析并提取 url。

sitemap_rules

A list of tuples (regex, callback) where:

- **regex** 是一个用于匹配从 sitemap 提供的 url 的正则表达式。regex 可以是一个字符串或者编译的正则对象 (compiled regex object)。
- **callback** 指定了匹配正则表达式的 url 的处理函数。callback 可以是一个字符串 (spider 中方法的名字) 或者是 callable。

例子:

```
sitemap_rules = [('/product/', 'parse_product')]
```

规则按顺序进行匹配，只有第一个匹配才会被应用。

如果忽略该属性，sitemap 中发现的所有 url 将会被 parse 处理。

sitemap_follow

一个用于匹配要跟进的 sitemap 的正则表达式的列表 (list)。。其仅仅被应用在使用 [Sitemap index files](#) 来指向其他 sitemap 文件的站点。

默认情况下所有的 sitemap 都会被跟进。

sitemap_alternate_links

当一个 url 有可选的链接，要跟进的时候，指定它。有些非英文网站会在一个 url 块内提供其他语言的网站链接。

例子:

```
<url>
  <loc>http://example.com/</loc>
  <xhtml:link rel="alternate" hreflang="de" href="http://example.com/
↪de"/>
</url>
```

设置了 `sitemap_alternate_links`，将会找到两个 url。禁用了 `sitemap_alternate_links`，只有 `http://example.com/` 被找到。

默认禁用 `sitemap_alternate_links`。

sitemap_filter(entries)

这是一个过滤功能的数据，可以重写它来根据网站标签的值来选择网站标签。

例如:

```
<url>
  <loc>http://example.com/</loc>
  <lastmod>2005-01-01</lastmod>
</url>
```

我们可以定义一个 `sitemap_filter` 方法来根据日期过滤 `entries`

```
from datetime import datetime
from scrapy.spiders import SitemapSpider

class FilteredSitemapSpider(SitemapSpider):
    name = 'filtered_sitemap_spider'
```

(下页继续)

(续上页)

```

allowed_domains = ['example.com']
sitemap_urls = ['http://example.com/sitemap.xml']

def sitemap_filter(self, entries):
    for entry in entries:
        date_time = datetime.strptime(entry['lastmod'], '%Y-%m-%d')
        if date_time.year >= 2005:
            yield entry

```

这将只检索 2005 年及以后年份修改过的 `entries`。

条目是从 sitemap 文档中提取的字典对象。通常，键是标记名称，值是其中的文本。

注意下面的几点:

- 因为 loc 属性是必需的, 没有此标记的条目将被丢弃
- 备用链接使用 alternate 键存储在列表中 (see `sitemap_alternate_links`)
- 命名空间被移除了, 因此命名 lxml tags 的 `{namespace}tagname` 变成了 `tagname`

如果忽略这个方法, 在 sitemaps 找到的所有条目都会被处理, 遵守其他属性及其设置。

SitemapSpider 例子

简单的例子: 使用 `parse` 处理通过 sitemap 发现的所有 url:

```

from scrapy.spiders import SitemapSpider

class MySpider(SitemapSpider):
    sitemap_urls = ['http://www.example.com/sitemap.xml']

    def parse(self, response):
        pass # ... scrape item here ...

```

用特定的函数处理某些 url, 其他的使用另外的回调:

```

from scrapy.spiders import SitemapSpider

class MySpider(SitemapSpider):
    sitemap_urls = ['http://www.example.com/sitemap.xml']
    sitemap_rules = [
        ('/product/', 'parse_product'),
        ('/category/', 'parse_category'),
    ]

```

(下页继续)

(续上页)

```
def parse_product(self, response):  
    pass # ... scrape product ...  
  
def parse_category(self, response):  
    pass # ... scrape category ...
```

跟进 `robots.txt` 定义的文件，并只跟进包含有 `/sitemap_shop` 的 url:

```
from scrapy.spiders import SitemapSpider  
  
class MySpider(SitemapSpider):  
    sitemap_urls = ['http://www.example.com/robots.txt']  
    sitemap_rules = [  
        ('/shop/', 'parse_shop'),  
    ]  
    sitemap_follow = ['/sitemap_shops']  
  
    def parse_shop(self, response):  
        pass # ... scrape shop here ...
```

在 `SitemapSpider` 中使用其他 url:

```
from scrapy.spiders import SitemapSpider  
  
class MySpider(SitemapSpider):  
    sitemap_urls = ['http://www.example.com/robots.txt']  
    sitemap_rules = [  
        ('/shop/', 'parse_shop'),  
    ]  
  
    other_urls = ['http://www.example.com/about']  
  
    def start_requests(self):  
        requests = list(super(MySpider, self).start_requests())  
        requests += [scrapy.Request(x, self.parse_other) for x in self.other_urls]  
        return requests  
  
    def parse_shop(self, response):  
        pass # ... scrape shop here ...
```

(下页继续)

(续上页)

```
def parse_other(self, response):  
    pass # ... scrape other here ...
```

3.3 选择器

When you're scraping web pages, the most common task you need to perform is to extract data from the HTML source. There are several libraries available to achieve this, such as:

- [BeautifulSoup](#) is a very popular web scraping library among Python programmers which constructs a Python object based on the structure of the HTML code and also deals with bad markup reasonably well, but it has one drawback: it's slow.
- [lxml](#) is an XML parsing library (which also parses HTML) with a pythonic API based on [ElementTree](#). (lxml is not part of the Python standard library.)

Scrapy comes with its own mechanism for extracting data. They're called selectors because they "select" certain parts of the HTML document specified either by [XPath](#) or [CSS](#) expressions.

[XPath](#) is a language for selecting nodes in XML documents, which can also be used with HTML. [CSS](#) is a language for applying styles to HTML documents. It defines selectors to associate those styles with specific HTML elements.

注解: Scrapy Selectors is a thin wrapper around [parsel](#) library; the purpose of this wrapper is to provide better integration with Scrapy Response objects.

[parsel](#) is a stand-alone web scraping library which can be used without Scrapy. It uses [lxml](#) library under the hood, and implements an easy API on top of lxml API. It means Scrapy selectors are very similar in speed and parsing accuracy to lxml.

3.3.1 Using selectors

Constructing selectors

Response objects expose a `Selector` instance on `.selector` attribute:

```
>>> response.selector.xpath('//span/text()').get()  
'good'
```

Querying responses using XPath and CSS is so common that responses include two more shortcuts: `response.xpath()` and `response.css()`:

```
>>> response.xpath('//span/text()').get()
'good'
>>> response.css('span::text').get()
'good'
```

Scrapy selectors are instances of `Selector` class constructed by passing either *TextResponse* object or markup as an unicode string (in `text` argument). Usually there is no need to construct Scrapy selectors manually: `response` object is available in Spider callbacks, so in most cases it is more convenient to use `response.css()` and `response.xpath()` shortcuts. By using `response.selector` or one of these shortcuts you can also ensure the response body is parsed only once.

But if required, it is possible to use `Selector` directly. Constructing from text:

```
>>> from scrapy.selector import Selector
>>> body = '<html><body><span>good</span></body></html>'
>>> Selector(text=body).xpath('//span/text()').get()
'good'
```

Constructing from response - *HtmlResponse* is one of *TextResponse* subclasses:

```
>>> from scrapy.selector import Selector
>>> from scrapy.http import HtmlResponse
>>> response = HtmlResponse(url='http://example.com', body=body)
>>> Selector(response=response).xpath('//span/text()').get()
'good'
```

`Selector` automatically chooses the best parsing rules (XML vs HTML) based on input type.

Using selectors

To explain how to use the selectors we'll use the **Scrapy shell** (which provides interactive testing) and an example page located in the Scrapy documentation server:

https://docs.scrapy.org/en/latest/_static/selectors-sample1.html

For the sake of completeness, here's its full HTML code:

```
<html>
<head>
  <base href='http://example.com/' />
  <title>Example website</title>
</head>
<body>
```

(下页继续)

(续上页)

```
<div id='images'>
  <a href='image1.html'>Name: My image 1 <br /><img src='image1_thumb.jpg' /></a>
  <a href='image2.html'>Name: My image 2 <br /><img src='image2_thumb.jpg' /></a>
  <a href='image3.html'>Name: My image 3 <br /><img src='image3_thumb.jpg' /></a>
  <a href='image4.html'>Name: My image 4 <br /><img src='image4_thumb.jpg' /></a>
  <a href='image5.html'>Name: My image 5 <br /><img src='image5_thumb.jpg' /></a>
</div>
</body>
</html>
```

First, let's open the shell:

```
scrapy shell https://docs.scrapy.org/en/latest/_static/selectors-sample1.html
```

Then, after the shell loads, you'll have the response available as `response` shell variable, and its attached selector in `response.selector` attribute.

Since we're dealing with HTML, the selector will automatically use an HTML parser.

So, by looking at the *HTML code* of that page, let's construct an XPath for selecting the text inside the title tag:

```
>>> response.xpath('//title/text()')
[<Selector xpath='//title/text()' data='Example website'>]
```

To actually extract the textual data, you must call the selector `.get()` or `.getall()` methods, as follows:

```
>>> response.xpath('//title/text()').getall()
['Example website']
>>> response.xpath('//title/text()').get()
'Example website'
```

`.get()` always returns a single result; if there are several matches, content of a first match is returned; if there are no matches, `None` is returned. `.getall()` returns a list with all results.

Notice that CSS selectors can select text or attribute nodes using CSS3 pseudo-elements:

```
>>> response.css('title::text').get()
'Example website'
```

As you can see, `.xpath()` and `.css()` methods return a `SelectorList` instance, which is a list of new selectors. This API can be used for quickly selecting nested data:

```
>>> response.css('img').xpath('@src').getall()
['image1_thumb.jpg',
 'image2_thumb.jpg',
 'image3_thumb.jpg',
 'image4_thumb.jpg',
 'image5_thumb.jpg']
```

If you want to extract only the first matched element, you can call the selector `.get()` (or its alias `.extract_first()` commonly used in previous Scrapy versions):

```
>>> response.xpath('//div[@id="images"]/a/text()').get()
'Name: My image 1 '
```

It returns `None` if no element was found:

```
>>> response.xpath('//div[@id="not-exists"]/text()').get() is None
True
```

A default return value can be provided as an argument, to be used instead of `None`:

```
>>> response.xpath('//div[@id="not-exists"]/text()').get(default='not-found')
'not-found'
```

Instead of using e.g. `'@src'` XPath it is possible to query for attributes using `.attrib` property of a Selector:

```
>>> [img.attrib['src'] for img in response.css('img')]
['image1_thumb.jpg',
 'image2_thumb.jpg',
 'image3_thumb.jpg',
 'image4_thumb.jpg',
 'image5_thumb.jpg']
```

As a shortcut, `.attrib` is also available on SelectorList directly; it returns attributes for the first matching element:

```
>>> response.css('img').attrib['src']
'image1_thumb.jpg'
```

This is most useful when only a single result is expected, e.g. when selecting by id, or selecting unique elements on a web page:

```
>>> response.css('base').attrib['href']
'http://example.com/'
```

Now we're going to get the base URL and some image links:

```
>>> response.xpath('//base/@href').get()
'http://example.com/'

>>> response.css('base::attr(href)').get()
'http://example.com/'

>>> response.css('base').attrib['href']
'http://example.com/'

>>> response.xpath('//a[contains(@href, "image")]/@href').getall()
['image1.html',
 'image2.html',
 'image3.html',
 'image4.html',
 'image5.html']

>>> response.css('a[href*=image]::attr(href)').getall()
['image1.html',
 'image2.html',
 'image3.html',
 'image4.html',
 'image5.html']

>>> response.xpath('//a[contains(@href, "image")]/img/@src').getall()
['image1_thumb.jpg',
 'image2_thumb.jpg',
 'image3_thumb.jpg',
 'image4_thumb.jpg',
 'image5_thumb.jpg']

>>> response.css('a[href*=image] img::attr(src)').getall()
['image1_thumb.jpg',
 'image2_thumb.jpg',
 'image3_thumb.jpg',
 'image4_thumb.jpg',
 'image5_thumb.jpg']
```

Extensions to CSS Selectors

Per W3C standards, [CSS selectors](#) do not support selecting text nodes or attribute values. But selecting these is so essential in a web scraping context that Scrapy (parsel) implements a couple of **non-standard pseudo-elements**:

- to select text nodes, use `::text`
- to select attribute values, use `::attr(name)` where *name* is the name of the attribute that you want the value of

警告: These pseudo-elements are Scrapy-/Parsel-specific. They will most probably not work with other libraries like [lxml](#) or [PyQuery](#).

Examples:

- `title::text` selects children text nodes of a descendant `<title>` element:

```
>>> response.css('title::text').get()
'Example website'
```

- `*::text` selects all descendant text nodes of the current selector context:

```
>>> response.css('#images *::text').getall()
['\n ',
 'Name: My image 1 ',
 '\n ',
 'Name: My image 2 ',
 '\n ',
 'Name: My image 3 ',
 '\n ',
 'Name: My image 4 ',
 '\n ',
 'Name: My image 5 ',
 '\n ']
```

- `foo::text` returns no results if `foo` element exists, but contains no text (i.e. text is empty):

```
>>> response.css('img::text').getall()
[]
```

This means `.css('foo::text').get()` could return `None` even if an element exists. Use `default=''` if you always want a string:


```
>>> response.css('img::text').get()
>>> response.css('img::text').get(default='')
''
```

- `a::attr(href)` selects the *href* attribute value of descendant links:

```
>>> response.css('a::attr(href)').getall()
['image1.html',
 'image2.html',
 'image3.html',
 'image4.html',
 'image5.html']
```

注解: See also: *Selecting element attributes*.

注解: You cannot chain these pseudo-elements. But in practice it would not make much sense: text nodes do not have attributes, and attribute values are string values already and do not have children nodes.

Nesting selectors

The selection methods (`.xpath()` or `.css()`) return a list of selectors of the same type, so you can call the selection methods for those selectors too. Here's an example:

```
>>> links = response.xpath('//a[contains(@href, "image")]')
>>> links.getall()
['<a href="image1.html">Name: My image 1 <br></a>',
 '<a href="image2.html">Name: My image 2 <br></a>',
 '<a href="image3.html">Name: My image 3 <br></a>',
 '<a href="image4.html">Name: My image 4 <br></a>',
 '<a href="image5.html">Name: My image 5 <br></a>']

>>> for index, link in enumerate(links):
...     args = (index, link.xpath('@href').get(), link.xpath('img/@src').get())
...     print('Link number %d points to url %r and image %r' % args)

Link number 0 points to url 'image1.html' and image 'image1_thumb.jpg'
Link number 1 points to url 'image2.html' and image 'image2_thumb.jpg'
Link number 2 points to url 'image3.html' and image 'image3_thumb.jpg'
```

(下页继续)

(续上页)

```
Link number 3 points to url 'image4.html' and image 'image4_thumb.jpg'  
Link number 4 points to url 'image5.html' and image 'image5_thumb.jpg'
```

Selecting element attributes

There are several ways to get a value of an attribute. First, one can use XPath syntax:

```
>>> response.xpath("//a/@href").getall()  
['image1.html', 'image2.html', 'image3.html', 'image4.html', 'image5.html']
```

XPath syntax has a few advantages: it is a standard XPath feature, and `@attributes` can be used in other parts of an XPath expression - e.g. it is possible to filter by attribute value.

Scrapy also provides an extension to CSS selectors (`::attr(...)`) which allows to get attribute values:

```
>>> response.css('a::attr(href)').getall()  
['image1.html', 'image2.html', 'image3.html', 'image4.html', 'image5.html']
```

In addition to that, there is a `.attrib` property of Selector. You can use it if you prefer to lookup attributes in Python code, without using XPaths or CSS extensions:

```
>>> [a.attrib['href'] for a in response.css('a')]  
['image1.html', 'image2.html', 'image3.html', 'image4.html', 'image5.html']
```

This property is also available on SelectorList; it returns a dictionary with attributes of a first matching element. It is convenient to use when a selector is expected to give a single result (e.g. when selecting by element ID, or when selecting an unique element on a page):

```
>>> response.css('base').attrib  
{'href': 'http://example.com/'}  
>>> response.css('base').attrib['href']  
'http://example.com/'
```

`.attrib` property of an empty SelectorList is empty:

```
>>> response.css('foo').attrib  
{}
```

Using selectors with regular expressions

Selector also has a `.re()` method for extracting data using regular expressions. However, unlike using `.xpath()` or `.css()` methods, `.re()` returns a list of unicode strings. So you can't construct nested `.re()`

calls.

Here's an example used to extract image names from the *HTML code* above:

```
>>> response.xpath('//a[contains(@href, "image")]/text()').re(r'Name:\s*(.*)')
['My image 1',
 'My image 2',
 'My image 3',
 'My image 4',
 'My image 5']
```

There's an additional helper reciprocating `.get()` (and its alias `.extract_first()`) for `.re()`, named `.re_first()`. Use it to extract just the first matching string:

```
>>> response.xpath('//a[contains(@href, "image")]/text()').re_first(r'Name:\s*(.*)')
'My image 1'
```

`extract()` and `extract_first()`

If you're a long-time Scrapy user, you're probably familiar with `.extract()` and `.extract_first()` selector methods. Many blog posts and tutorials are using them as well. These methods are still supported by Scrapy, there are **no plans** to deprecate them.

However, Scrapy usage docs are now written using `.get()` and `.getall()` methods. We feel that these new methods result in a more concise and readable code.

The following examples show how these methods map to each other.

1. `SelectorList.get()` is the same as `SelectorList.extract_first()`:

```
>>> response.css('a::attr(href)').get()
'image1.html'
>>> response.css('a::attr(href)').extract_first()
'image1.html'
```

2. `SelectorList.getall()` is the same as `SelectorList.extract()`:

```
>>> response.css('a::attr(href)').getall()
['image1.html', 'image2.html', 'image3.html', 'image4.html', 'image5.html']
>>> response.css('a::attr(href)').extract()
['image1.html', 'image2.html', 'image3.html', 'image4.html', 'image5.html']
```

3. `Selector.get()` is the same as `Selector.extract()`:

```
>>> response.css('a::attr(href)')[0].get()
'image1.html'
>>> response.css('a::attr(href)')[0].extract()
'image1.html'
```

4. For consistency, there is also `Selector.getall()`, which returns a list:

```
>>> response.css('a::attr(href)')[0].getall()
['image1.html']
```

So, the main difference is that output of `.get()` and `.getall()` methods is more predictable: `.get()` always returns a single result, `.getall()` always returns a list of all extracted results. With `.extract()` method it was not always obvious if a result is a list or not; to get a single result either `.extract()` or `.extract_first()` should be called.

3.3.2 Working with XPath

Here are some tips which may help you to use XPath with Scrapy selectors effectively. If you are not much familiar with XPath yet, you may want to take a look first at this [XPath tutorial](#).

注解: Some of the tips are based on [this post from ScrapingHub](#)' s blog.

Working with relative XPath

Keep in mind that if you are nesting selectors and use an XPath that starts with `/`, that XPath will be absolute to the document and not relative to the `Selector` you' re calling it from.

For example, suppose you want to extract all `<p>` elements inside `<div>` elements. First, you would get all `<div>` elements:

```
>>> divs = response.xpath('//div')
```

At first, you may be tempted to use the following approach, which is wrong, as it actually extracts all `<p>` elements from the document, not only those inside `<div>` elements:

```
>>> for p in divs.xpath('/p'): # this is wrong - gets all <p> from the whole document
...     print(p.get())
```

This is the proper way to do it (note the dot prefixing the `./p` XPath):

```
>>> for p in divs.xpath('..//p'): # extracts all <p> inside
...     print(p.get())
```

Another common case would be to extract all direct <p> children:

```
>>> for p in divs.xpath('p'):
...     print(p.get())
```

For more details about relative XPath expressions see the [Location Paths](#) section in the XPath specification.

When querying by class, consider using CSS

Because an element can contain multiple CSS classes, the XPath way to select elements by class is the rather verbose:

```
*[contains(concat(' ', normalize-space(@class), ' '), ' someclass ')]
```

If you use @class='someclass' you may end up missing elements that have other classes, and if you just use contains(@class, 'someclass') to make up for that you may end up with more elements that you want, if they have a different class name that shares the string someclass.

As it turns out, Scrapy selectors allow you to chain selectors, so most of the time you can just select by class using CSS and then switch to XPath when needed:

```
>>> from scrapy import Selector
>>> sel = Selector(text='<div class="hero shout"><time datetime="2014-07-23 19:00">
↳Special date</time></div>')
>>> sel.css('.shout').xpath('..//time/@datetime').getall()
['2014-07-23 19:00']
```

This is cleaner than using the verbose XPath trick shown above. Just remember to use the . in the XPath expressions that will follow.

Beware of the difference between //node[1] and (//node)[1]

//node[1] selects all the nodes occurring first under their respective parents.

(//node)[1] selects all the nodes in the document, and then gets only the first of them.

Example:

```
>>> from scrapy import Selector
>>> sel = Selector(text="""
....:     <ul class="list">
```

(下页继续)

(续上页)

```
....:      <li>1</li>
....:      <li>2</li>
....:      <li>3</li>
....:    </ul>
....:    <ul class="list">
....:      <li>4</li>
....:      <li>5</li>
....:      <li>6</li>
....:    </ul>""")
>>> xp = lambda x: sel.xpath(x).getall()
```

This gets all first `` elements under whatever it is its parent:

```
>>> xp("//li[1]")
['<li>1</li>', '<li>4</li>']
```

And this gets the first `` element in the whole document:

```
>>> xp("(//li)[1]")
['<li>1</li>']
```

This gets all first `` elements under an `` parent:

```
>>> xp("//ul/li[1]")
['<li>1</li>', '<li>4</li>']
```

And this gets the first `` element under an `` parent in the whole document:

```
>>> xp("(//ul/li)[1]")
['<li>1</li>']
```

Using text nodes in a condition

When you need to use the text content as argument to an XPath string function, avoid using `./text()` and use just `.` instead.

This is because the expression `./text()` yields a collection of text elements – a *node-set*. And when a node-set is converted to a string, which happens when it is passed as argument to a string function like `contains()` or `starts-with()`, it results in the text for the first element only.

Example:

```
>>> from scrapy import Selector
>>> sel = Selector(text='<a href="#">Click here to go to the <strong>Next Page</strong></a>')
```

Converting a *node-set* to string:

```
>>> sel.xpath('//a//text()').getall() # take a peek at the node-set
['Click here to go to the ', 'Next Page']
>>> sel.xpath("string(//a[1]//text())").getall() # convert it to string
['Click here to go to the ']
```

A *node* converted to a string, however, puts together the text of itself plus of all its descendants:

```
>>> sel.xpath("//a[1]").getall() # select the first node
['<a href="#">Click here to go to the <strong>Next Page</strong></a>']
>>> sel.xpath("string(//a[1])").getall() # convert it to string
['Click here to go to the Next Page']
```

So, using the `./text()` node-set won't select anything in this case:

```
>>> sel.xpath("//a[contains(./text(), 'Next Page')]").getall()
[]
```

But using the `.` to mean the node, works:

```
>>> sel.xpath("//a[contains(., 'Next Page')]").getall()
['<a href="#">Click here to go to the <strong>Next Page</strong></a>']
```

Variables in XPath expressions

XPath allows you to reference variables in your XPath expressions, using the `$somevariable` syntax. This is somewhat similar to parameterized queries or prepared statements in the SQL world where you replace some arguments in your queries with placeholders like `?`, which are then substituted with values passed with the query.

Here's an example to match an element based on its "id" attribute value, without hard-coding it (that was shown previously):

```
>>> # `$val` used in the expression, a `val` argument needs to be passed
>>> response.xpath('//div[@id=$val]/a/text()', val='images').get()
'Name: My image 1 '
```

Here's another example, to find the "id" attribute of a <div> tag containing five <a> children (here we pass the value 5 as an integer):

```
>>> response.xpath('//div[count(a)=$cnt]/@id', cnt=5).get()
'images'
```

All variable references must have a binding value when calling `.xpath()` (otherwise you'll get a `ValueError: XPath error: exception`). This is done by passing as many named arguments as necessary.

`parsel`, the library powering Scrapy selectors, has more details and examples on [XPath variables](#).

Removing namespaces

When dealing with scraping projects, it is often quite convenient to get rid of namespaces altogether and just work with element names, to write more simple/convenient XPath expressions. You can use the `Selector.remove_namespaces()` method for that.

Let's show an example that illustrates this with the Python Insider blog atom feed.

First, we open the shell with the url we want to scrape:

```
$ scrapy shell https://feeds.feedburner.com/PythonInsider
```

This is how the file starts:

```
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet ...
<feed xmlns="http://www.w3.org/2005/Atom"
      xmlns:openSearch="http://a9.com/-/spec/opensearchrss/1.0/"
      xmlns:blogger="http://schemas.google.com/blogger/2008"
      xmlns:georss="http://www.georss.org/georss"
      xmlns:gd="http://schemas.google.com/g/2005"
      xmlns:thr="http://purl.org/syndication/thread/1.0"
      xmlns:feedburner="http://rssnamespace.org/feedburner/ext/1.0">
  ...
```

You can see several namespace declarations including a default "http://www.w3.org/2005/Atom" and another one using the "gd:" prefix for "http://schemas.google.com/g/2005".

Once in the shell we can try selecting all <link> objects and see that it doesn't work (because the Atom XML namespace is obfuscating those nodes):

```
>>> response.xpath("//link")
[]
```


But once we call the `Selector.remove_namespaces()` method, all nodes can be accessed directly by their names:

```
>>> response.selector.remove_namespaces()
>>> response.xpath("//link")
[<Selector xpath='//link' data='<link rel="alternate" type="text/html" h'>,
 <Selector xpath='//link' data='<link rel="next" type="application/atom+'>,
 ...]
```

If you wonder why the namespace removal procedure isn't always called by default instead of having to call it manually, this is because of two reasons, which, in order of relevance, are:

1. Removing namespaces requires to iterate and modify all nodes in the document, which is a reasonably expensive operation to perform by default for all documents crawled by Scrapy
2. There could be some cases where using namespaces is actually required, in case some element names clash between namespaces. These cases are very rare though.

Using EXSLT extensions

Being built atop `lxml`, Scrapy selectors support some [EXSLT](#) extensions and come with these pre-registered namespaces to use in XPath expressions:

prefix	namespace	usage
re	http://exslt.org/regular-expressions	regular expressions
set	http://exslt.org/sets	set manipulation

Regular expressions

The `test()` function, for example, can prove quite useful when XPath's `starts-with()` or `contains()` are not sufficient.

Example selecting links in list item with a “class” attribute ending with a digit:

```
>>> from scrapy import Selector
>>> doc = u"""
... <div>
...     <ul>
...         <li class="item-0"><a href="link1.html">first item</a></li>
...         <li class="item-1"><a href="link2.html">second item</a></li>
...         <li class="item-inactive"><a href="link3.html">third item</a></li>
...         <li class="item-1"><a href="link4.html">fourth item</a></li>
```

(下页继续)

(续上页)

```

...         <li class="item-0"><a href="link5.html">fifth item</a></li>
...     </ul>
... </div>
... """
>>> sel = Selector(text=doc, type="html")
>>> sel.xpath('//li//@href').getall()
['link1.html', 'link2.html', 'link3.html', 'link4.html', 'link5.html']
>>> sel.xpath('//li[re:test(@class, "item-\d$")]/@href').getall()
['link1.html', 'link2.html', 'link4.html', 'link5.html']
>>>

```

警告: C library libxslt doesn't natively support EXSLT regular expressions so lxml's implementation uses hooks to Python's re module. Thus, using regexp functions in your XPath expressions may add a small performance penalty.

Set operations

These can be handy for excluding parts of a document tree before extracting text elements for example.

Example extracting microdata (sample content taken from <http://schema.org/Product>) with groups of item-scopes and corresponding itemprops:

```

>>> doc = u"""
... <div itemscope itemtype="http://schema.org/Product">
...     <span itemprop="name">Kenmore White 17" Microwave</span>
...     
...     <div itemprop="aggregateRating"
...         itemscope itemtype="http://schema.org/AggregateRating">
...         Rated <span itemprop="ratingValue">3.5</span>/5
...         based on <span itemprop="reviewCount">11</span> customer reviews
...     </div>
...
...     <div itemprop="offers" itemscope itemtype="http://schema.org/Offer">
...         <span itemprop="price">$55.00</span>
...         <link itemprop="availability" href="http://schema.org/InStock" />In stock
...     </div>
...
...     Product description:
...     <span itemprop="description">0.7 cubic feet countertop microwave.

```

(下页继续)

(续上页)

```

... Has six preset cooking categories and convenience features like
... Add-A-Minute and Child Lock.</span>
...
... Customer reviews:
...
... <div itemprop="review" itemscope itemtype="http://schema.org/Review">
...   <span itemprop="name">Not a happy camper</span> -
...   by <span itemprop="author">Ellie</span>,
...   <meta itemprop="datePublished" content="2011-04-01">April 1, 2011
...   <div itemprop="reviewRating" itemscope itemtype="http://schema.org/Rating">
...     <meta itemprop="worstRating" content = "1">
...     <span itemprop="ratingValue">1</span>/
...     <span itemprop="bestRating">5</span>stars
...   </div>
...   <span itemprop="description">The lamp burned out and now I have to replace
...   it. </span>
... </div>
...
... <div itemprop="review" itemscope itemtype="http://schema.org/Review">
...   <span itemprop="name">Value purchase</span> -
...   by <span itemprop="author">Lucas</span>,
...   <meta itemprop="datePublished" content="2011-03-25">March 25, 2011
...   <div itemprop="reviewRating" itemscope itemtype="http://schema.org/Rating">
...     <meta itemprop="worstRating" content = "1"/>
...     <span itemprop="ratingValue">4</span>/
...     <span itemprop="bestRating">5</span>stars
...   </div>
...   <span itemprop="description">Great microwave for the price. It is small and
...   fits in my apartment.</span>
... </div>
...   ...
... </div>
... """
>>> sel = Selector(text=doc, type="html")
>>> for scope in sel.xpath('//div[@itemscope]'):
...     print("current scope:", scope.xpath('@itemtype').getall())
...     props = scope.xpath(''
...         set:difference(./descendant::*/@itemprop,
...             .*[@itemscope]/*/@itemprop)''')
...     print("    properties: %s" % (props.getall()))

```

(下页继续)

(续上页)

```
...     print("")

current scope: ['http://schema.org/Product']
    properties: ['name', 'aggregateRating', 'offers', 'description', 'review', 'review']

current scope: ['http://schema.org/AggregateRating']
    properties: ['ratingValue', 'reviewCount']

current scope: ['http://schema.org/Offer']
    properties: ['price', 'availability']

current scope: ['http://schema.org/Review']
    properties: ['name', 'author', 'datePublished', 'reviewRating', 'description']

current scope: ['http://schema.org/Rating']
    properties: ['worstRating', 'ratingValue', 'bestRating']

current scope: ['http://schema.org/Review']
    properties: ['name', 'author', 'datePublished', 'reviewRating', 'description']

current scope: ['http://schema.org/Rating']
    properties: ['worstRating', 'ratingValue', 'bestRating']

>>>
```

Here we first iterate over `itemscope` elements, and for each one, we look for all `itemprops` elements and exclude those that are themselves inside another `itemscope`.

Other XPath extensions

Scrapy selectors also provide a sorely missed XPath extension function `has-class` that returns `True` for nodes that have all of the specified HTML classes.

For the following HTML:

```
<p class="foo bar-baz">First</p>
<p class="foo">Second</p>
<p class="bar">Third</p>
<p>Fourth</p>
```

You can use it like this:

```
>>> response.xpath('//p[has-class("foo")]')
[<Selector xpath='//p[has-class("foo")]' data='<p class="foo bar-baz">First</p>'>,
 <Selector xpath='//p[has-class("foo")]' data='<p class="foo">Second</p>'>]
>>> response.xpath('//p[has-class("foo", "bar-baz")]')
[<Selector xpath='//p[has-class("foo", "bar-baz")]' data='<p class="foo bar-baz">First</
↪p>'>]
>>> response.xpath('//p[has-class("foo", "bar")]')
[]
```

So XPath `//p[has-class("foo", "bar-baz")]` is roughly equivalent to CSS `p.foo.bar-baz`. Please note, that it is slower in most of the cases, because it's a pure-Python function that's invoked for every node in question whereas the CSS lookup is translated into XPath and thus runs more efficiently, so performance-wise its uses are limited to situations that are not easily described with CSS selectors.

Parsel also simplifies adding your own XPath extensions.

3.3.3 Built-in Selectors reference

Selector objects

SelectorList objects

3.3.4 Examples

Selector examples on HTML response

Here are some `Selector` examples to illustrate several concepts. In all cases, we assume there is already a `Selector` instantiated with a `HtmlResponse` object like this:

```
sel = Selector(html_response)
```

1. Select all `<h1>` elements from an HTML response body, returning a list of `Selector` objects (ie. a `SelectorList` object):

```
sel.xpath("//h1")
```

2. Extract the text of all `<h1>` elements from an HTML response body, returning a list of unicode strings:

```
sel.xpath("//h1").getall()           # this includes the h1 tag
sel.xpath("//h1/text()").getall()    # this excludes the h1 tag
```

3. Iterate over all `<p>` tags and print their class attribute:

```
for node in sel.xpath("//p"):
    print(node.attrib['class'])
```

Selector examples on XML response

Here are some examples to illustrate concepts for `Selector` objects instantiated with an *XmlResponse* object:

```
sel = Selector(xml_response)
```

1. Select all `<product>` elements from an XML response body, returning a list of `Selector` objects (ie. a `SelectorList` object):

```
sel.xpath("//product")
```

2. Extract all prices from a *Google Base XML* feed which requires registering a namespace:

```
sel.register_namespace("g", "http://base.google.com/ns/1.0")
sel.xpath("//g:price").getall()
```

3.4 Items

The main goal in scraping is to extract structured data from unstructured sources, typically, web pages. Scrapy spiders can return the extracted data as Python dicts. While convenient and familiar, Python dicts lack structure: it is easy to make a typo in a field name or return inconsistent data, especially in a larger project with many spiders.

To define common output data format Scrapy provides the *Item* class. *Item* objects are simple containers used to collect the scraped data. They provide a dictionary-like API with a convenient syntax for declaring their available fields.

Various Scrapy components use extra information provided by Items: exporters look at declared fields to figure out columns to export, serialization can be customized using Item fields metadata, `trackref` tracks Item instances to help find memory leaks (see *Debugging memory leaks with trackref*), etc.

3.4.1 Declaring Items

Items are declared using a simple class definition syntax and *Field* objects. Here is an example:

```
import scrapy

class Product(scrapy.Item):
```

(下页继续)

(续上页)

```
name = scrapy.Field()
price = scrapy.Field()
stock = scrapy.Field()
tags = scrapy.Field()
last_updated = scrapy.Field(serializer=str)
```

注解: Those familiar with [Django](#) will notice that Scrapy Items are declared similar to [Django Models](#), except that Scrapy Items are much simpler as there is no concept of different field types.

3.4.2 Item Fields

Field objects are used to specify metadata for each field. For example, the serializer function for the `last_updated` field illustrated in the example above.

You can specify any kind of metadata for each field. There is no restriction on the values accepted by *Field* objects. For this same reason, there is no reference list of all available metadata keys. Each key defined in *Field* objects could be used by a different component, and only those components know about it. You can also define and use any other *Field* key in your project too, for your own needs. The main goal of *Field* objects is to provide a way to define all field metadata in one place. Typically, those components whose behaviour depends on each field use certain field keys to configure that behaviour. You must refer to their documentation to see which metadata keys are used by each component.

It's important to note that the *Field* objects used to declare the item do not stay assigned as class attributes. Instead, they can be accessed through the *Item.fields* attribute.

3.4.3 Working with Items

Here are some examples of common tasks performed with items, using the `Product` item *declared above*. You will notice the API is very similar to the `dict` API.

Creating items

```
>>> product = Product(name='Desktop PC', price=1000)
>>> print(product)
Product(name='Desktop PC', price=1000)
```

Getting field values

```
>>> product['name']
Desktop PC
>>> product.get('name')
Desktop PC

>>> product['price']
1000

>>> product['last_updated']
Traceback (most recent call last):
...
KeyError: 'last_updated'

>>> product.get('last_updated', 'not set')
not set

>>> product['lala'] # getting unknown field
Traceback (most recent call last):
...
KeyError: 'lala'

>>> product.get('lala', 'unknown field')
'unknown field'

>>> 'name' in product # is name field populated?
True

>>> 'last_updated' in product # is last_updated populated?
False

>>> 'last_updated' in product.fields # is last_updated a declared field?
True

>>> 'lala' in product.fields # is lala a declared field?
False
```


Setting field values

```
>>> product['last_updated'] = 'today'
>>> product['last_updated']
today

>>> product['lala'] = 'test' # setting unknown field
Traceback (most recent call last):
...
KeyError: 'Product does not support field: lala'
```

Accessing all populated values

To access all populated values, just use the typical [dict API](#):

```
>>> product.keys()
['price', 'name']

>>> product.items()
[('price', 1000), ('name', 'Desktop PC')]
```

Copying items

To copy an item, you must first decide whether you want a shallow copy or a deep copy.

If your item contains [mutable](#) values like lists or dictionaries, a shallow copy will keep references to the same mutable values across all different copies.

For example, if you have an item with a list of tags, and you create a shallow copy of that item, both the original item and the copy have the same list of tags. Adding a tag to the list of one of the items will add the tag to the other item as well.

If that is not the desired behavior, use a deep copy instead.

See the [documentation of the copy module](#) for more information.

To create a shallow copy of an item, you can either call `copy()` on an existing item (`product2 = product.copy()`) or instantiate your item class from an existing item (`product2 = Product(product)`).

To create a deep copy, call `deepcopy()` instead (`product2 = product.deepcopy()`).

Other common tasks

Creating dicts from items:

```
>>> dict(product) # create a dict from all populated values
{'price': 1000, 'name': 'Desktop PC'}
```

Creating items from dicts:

```
>>> Product({'name': 'Laptop PC', 'price': 1500})
Product(price=1500, name='Laptop PC')

>>> Product({'name': 'Laptop PC', 'lala': 1500}) # warning: unknown field in dict
Traceback (most recent call last):
...
KeyError: 'Product does not support field: lala'
```

3.4.4 Extending Items

You can extend Items (to add more fields or to change some metadata for some fields) by declaring a subclass of your original Item.

For example:

```
class DiscountedProduct(Product):
    discount_percent = scrapy.Field(serializer=str)
    discount_expiration_date = scrapy.Field()
```

You can also extend field metadata by using the previous field metadata and appending more values, or changing existing values, like this:

```
class SpecificProduct(Product):
    name = scrapy.Field(Product.fields['name'], serializer=my_serializer)
```

That adds (or replaces) the `serializer` metadata key for the `name` field, keeping all the previously existing metadata values.

3.4.5 Item objects

```
class scrapy.item.Item([arg])
```

Return a new Item optionally initialized from the given argument.

Items replicate the standard `dict` API, including its constructor. The only additional attribute provided by Items is:

fields

A dictionary containing *all declared fields* for this Item, not only those populated. The keys are

the field names and the values are the *Field* objects used in the *Item declaration*.

3.4.6 Field objects

```
class scrapy.item.Field([arg])
```

The *Field* class is just an alias to the built-in *dict* class and doesn't provide any extra functionality or attributes. In other words, *Field* objects are plain-old Python dicts. A separate class is used to support the *item declaration syntax* based on class attributes.

3.5 Item 加载器

Item Loaders provide a convenient mechanism for populating scraped *Items*. Even though Items can be populated using their own dictionary-like API, Item Loaders provide a much more convenient API for populating them from a scraping process, by automating some common tasks like parsing the raw extracted data before assigning it.

In other words, *Items* provide the *container* of scraped data, while Item Loaders provide the mechanism for *populating* that container.

Item Loaders are designed to provide a flexible, efficient and easy mechanism for extending and overriding different field parsing rules, either by spider, or by source format (HTML, XML, etc) without becoming a nightmare to maintain.

3.5.1 Using Item Loaders to populate items

To use an Item Loader, you must first instantiate it. You can either instantiate it with a dict-like object (e.g. Item or dict) or without one, in which case an Item is automatically instantiated in the Item Loader constructor using the Item class specified in the *ItemLoader.default_item_class* attribute.

Then, you start collecting values into the Item Loader, typically using *Selectors*. You can add more than one value to the same item field; the Item Loader will know how to “join” those values later using a proper processing function.

Here is a typical Item Loader usage in a *Spider*, using the *Product item* declared in the *Items chapter*:

```
from scrapy.loader import ItemLoader
from myproject.items import Product

def parse(self, response):
    l = ItemLoader(item=Product(), response=response)
    l.add_xpath('name', '//div[@class="product_name"]')
    l.add_xpath('name', '//div[@class="product_title"]')
```

(下页继续)

```
l.add_xpath('price', '//p[@id="price"]')
l.add_css('stock', 'p#stock')
l.add_value('last_updated', 'today') # you can also use literal values
return l.load_item()
```

By quickly looking at that code, we can see the `name` field is being extracted from two different XPath locations in the page:

1. `//div[@class="product_name"]`
2. `//div[@class="product_title"]`

In other words, data is being collected by extracting it from two XPath locations, using the `add_xpath()` method. This is the data that will be assigned to the `name` field later.

Afterwards, similar calls are used for `price` and `stock` fields (the latter using a CSS selector with the `add_css()` method), and finally the `last_update` field is populated directly with a literal value (`today`) using a different method: `add_value()`.

Finally, when all data is collected, the `ItemLoader.load_item()` method is called which actually returns the item populated with the data previously extracted and collected with the `add_xpath()`, `add_css()`, and `add_value()` calls.

3.5.2 Input and Output processors

An Item Loader contains one input processor and one output processor for each (item) field. The input processor processes the extracted data as soon as it's received (through the `add_xpath()`, `add_css()` or `add_value()` methods) and the result of the input processor is collected and kept inside the ItemLoader. After collecting all data, the `ItemLoader.load_item()` method is called to populate and get the populated `Item` object. That's when the output processor is called with the data previously collected (and processed using the input processor). The result of the output processor is the final value that gets assigned to the item.

Let's see an example to illustrate how the input and output processors are called for a particular field (the same applies for any other field):

```
l = ItemLoader(Product(), some_selector)
l.add_xpath('name', xpath1) # (1)
l.add_xpath('name', xpath2) # (2)
l.add_css('name', css) # (3)
l.add_value('name', 'test') # (4)
return l.load_item() # (5)
```

So what happens is:

1. Data from `xpath1` is extracted, and passed through the *input processor* of the `name` field. The result of the input processor is collected and kept in the Item Loader (but not yet assigned to the item).
2. Data from `xpath2` is extracted, and passed through the same *input processor* used in (1). The result of the input processor is appended to the data collected in (1) (if any).
3. This case is similar to the previous ones, except that the data is extracted from the `css` CSS selector, and passed through the same *input processor* used in (1) and (2). The result of the input processor is appended to the data collected in (1) and (2) (if any).
4. This case is also similar to the previous ones, except that the value to be collected is assigned directly, instead of being extracted from a XPath expression or a CSS selector. However, the value is still passed through the input processors. In this case, since the value is not iterable it is converted to an iterable of a single element before passing it to the input processor, because input processor always receive iterables.
5. The data collected in steps (1), (2), (3) and (4) is passed through the *output processor* of the `name` field. The result of the output processor is the value assigned to the `name` field in the item.

It's worth noticing that processors are just callable objects, which are called with the data to be parsed, and return a parsed value. So you can use any function as input or output processor. The only requirement is that they must accept one (and only one) positional argument, which will be an iterator.

注解: Both input and output processors must receive an iterator as their first argument. The output of those functions can be anything. The result of input processors will be appended to an internal list (in the Loader) containing the collected values (for that field). The result of the output processors is the value that will be finally assigned to the item.

If you want to use a plain function as a processor, make sure it receives `self` as the first argument:

```
def lowercase_processor(self, values):
    for v in values:
        yield v.lower()

class MyItemLoader(ItemLoader):
    name_in = lowercase_processor
```

This is because whenever a function is assigned as a class variable, it becomes a method and would be passed the instance as the the first argument when being called. See [this answer on stackoverflow](#) for more details.

The other thing you need to keep in mind is that the values returned by input processors are collected internally (in lists) and then passed to output processors to populate the fields.

Last, but not least, Scrapy comes with some *commonly used processors* built-in for convenience.

3.5.3 Declaring Item Loaders

Item Loaders are declared like Items, by using a class definition syntax. Here is an example:

```
from scrapy.loader import ItemLoader
from scrapy.loader.processors import TakeFirst, MapCompose, Join

class ProductLoader(ItemLoader):

    default_output_processor = TakeFirst()

    name_in = MapCompose(unicode.title)
    name_out = Join()

    price_in = MapCompose(unicode.strip)

    # ...
```

As you can see, input processors are declared using the `_in` suffix while output processors are declared using the `_out` suffix. And you can also declare a default input/output processors using the *ItemLoader.default_input_processor* and *ItemLoader.default_output_processor* attributes.

3.5.4 Declaring Input and Output Processors

As seen in the previous section, input and output processors can be declared in the Item Loader definition, and it's very common to declare input processors this way. However, there is one more place where you can specify the input and output processors to use: in the *Item Field* metadata. Here is an example:

```
import scrapy
from scrapy.loader.processors import Join, MapCompose, TakeFirst
from w3lib.html import remove_tags

def filter_price(value):
    if value.isdigit():
        return value

class Product(scrapy.Item):
    name = scrapy.Field(
        input_processor=MapCompose(remove_tags),
        output_processor=Join(),
    )
```

(下页继续)

(续上页)

```
price = scrapy.Field(
    input_processor=MapCompose(remove_tags, filter_price),
    output_processor=TakeFirst(),
)
```

```
>>> from scrapy.loader import ItemLoader
>>> il = ItemLoader(item=Product())
>>> il.add_value('name', [u'Welcome to my', u'<strong>website</strong>'])
>>> il.add_value('price', [u'&euro;', u'<span>1000</span>'])
>>> il.load_item()
{'name': u'Welcome to my website', 'price': u'1000'}
```

The precedence order, for both input and output processors, is as follows:

1. Item Loader field-specific attributes: `field_in` and `field_out` (most precedence)
2. Field metadata (`input_processor` and `output_processor` key)
3. Item Loader defaults: `ItemLoader.default_input_processor()` and `ItemLoader.default_output_processor()` (least precedence)

See also: *Reusing and extending Item Loaders*.

3.5.5 Item Loader Context

The Item Loader Context is a dict of arbitrary key/values which is shared among all input and output processors in the Item Loader. It can be passed when declaring, instantiating or using Item Loader. They are used to modify the behaviour of the input/output processors.

For example, suppose you have a function `parse_length` which receives a text value and extracts a length from it:

```
def parse_length(text, loader_context):
    unit = loader_context.get('unit', 'm')
    # ... length parsing code goes here ...
    return parsed_length
```

By accepting a `loader_context` argument the function is explicitly telling the Item Loader that it's able to receive an Item Loader context, so the Item Loader passes the currently active context when calling it, and the processor function (`parse_length` in this case) can thus use them.

There are several ways to modify Item Loader context values:

1. By modifying the currently active Item Loader context (`context` attribute):

```
loader = ItemLoader(product)
loader.context['unit'] = 'cm'
```

2. On Item Loader instantiation (the keyword arguments of Item Loader constructor are stored in the Item Loader context):

```
loader = ItemLoader(product, unit='cm')
```

3. On Item Loader declaration, for those input/output processors that support instantiating them with an Item Loader context. MapCompose is one of them:

```
class ProductLoader(ItemLoader):
    length_out = MapCompose(parse_length, unit='cm')
```

3.5.6 ItemLoader objects

```
class scrapy.loader.ItemLoader([item, selector, response], **kwargs)
```

Return a new Item Loader for populating the given Item. If no item is given, one is instantiated automatically using the class in `default_item_class`.

When instantiated with a `selector` or a `response` parameters the `ItemLoader` class provides convenient mechanisms for extracting data from web pages using *selectors*.

参数

- **item** (*Item* object) – The item instance to populate using subsequent calls to `add_xpath()`, `add_css()`, or `add_value()`.
- **selector** (*Selector* object) – The selector to extract data from, when using the `add_xpath()` (resp. `add_css()`) or `replace_xpath()` (resp. `replace_css()`) method.
- **response** (*Response* object) – The response used to construct the selector using the `default_selector_class`, unless the selector argument is given, in which case this argument is ignored.

The item, selector, response and the remaining keyword arguments are assigned to the Loader context (accessible through the `context` attribute).

`ItemLoader` instances have the following methods:

```
get_value(value, *processors, **kwargs)
```

Process the given `value` by the given `processors` and keyword arguments.

Available keyword arguments:

参数 *re* (*str* or *compiled regex*) – a regular expression to use for extracting data from the given value using `extract_regex()` method, applied before processors

Examples:

```
>>> from scrapy.loader.processors import TakeFirst
>>> loader.get_value(u'name: foo', TakeFirst(), unicode.upper, re='name: (.+)')
'FOO'
```

`add_value(field_name, value, *processors, **kwargs)`

Process and then add the given **value** for the given field.

The value is first passed through `get_value()` by giving the **processors** and **kwargs**, and then passed through the *field input processor* and its result appended to the data collected for that field. If the field already contains collected data, the new data is added.

The given **field_name** can be `None`, in which case values for multiple fields may be added. And the processed value should be a dict with **field_name** mapped to values.

Examples:

```
loader.add_value('name', u'Color TV')
loader.add_value('colours', [u'white', u'blue'])
loader.add_value('length', u'100')
loader.add_value('name', u'name: foo', TakeFirst(), re='name: (.+)')
loader.add_value(None, {'name': u'foo', 'sex': u'male'})
```

`replace_value(field_name, value, *processors, **kwargs)`

Similar to `add_value()` but replaces the collected data with the new value instead of adding it.

`get_xpath(xpath, *processors, **kwargs)`

Similar to `ItemLoader.get_value()` but receives an XPath instead of a value, which is used to extract a list of unicode strings from the selector associated with this *ItemLoader*.

参数

- **xpath** (*str*) – the XPath to extract data from
- **re** (*str* or *compiled regex*) – a regular expression to use for extracting data from the selected XPath region

Examples:

```
# HTML snippet: <p class="product-name">Color TV</p>
loader.get_xpath('//p[@class="product-name"]')
# HTML snippet: <p id="price">the price is $1200</p>
loader.get_xpath('//p[@id="price"]', TakeFirst(), re='the price is (.*)')
```

`add_xpath(field_name, xpath, *processors, **kwargs)`

Similar to `ItemLoader.add_value()` but receives an XPath instead of a value, which is used to extract a list of unicode strings from the selector associated with this `ItemLoader`.

See `get_xpath()` for kwargs.

参数 `xpath (str)` – the XPath to extract data from

Examples:

```
# HTML snippet: <p class="product-name">Color TV</p>
loader.add_xpath('name', '//p[@class="product-name"]')
# HTML snippet: <p id="price">the price is $1200</p>
loader.add_xpath('price', '//p[@id="price"]', re='the price is (.*)')
```

`replace_xpath(field_name, xpath, *processors, **kwargs)`

Similar to `add_xpath()` but replaces collected data instead of adding it.

`get_css(css, *processors, **kwargs)`

Similar to `ItemLoader.get_value()` but receives a CSS selector instead of a value, which is used to extract a list of unicode strings from the selector associated with this `ItemLoader`.

参数

- `css (str)` – the CSS selector to extract data from
- `re (str or compiled regex)` – a regular expression to use for extracting data from the selected CSS region

Examples:

```
# HTML snippet: <p class="product-name">Color TV</p>
loader.get_css('p.product-name')
# HTML snippet: <p id="price">the price is $1200</p>
loader.get_css('p#price', TakeFirst(), re='the price is (.*)')
```

`add_css(field_name, css, *processors, **kwargs)`

Similar to `ItemLoader.add_value()` but receives a CSS selector instead of a value, which is used to extract a list of unicode strings from the selector associated with this `ItemLoader`.

See `get_css()` for kwargs.

参数 `css (str)` – the CSS selector to extract data from

Examples:

```
# HTML snippet: <p class="product-name">Color TV</p>
loader.add_css('name', 'p.product-name')
```

(下页继续)

(续上页)

```
# HTML snippet: <p id="price">the price is $1200</p>
loader.add_css('price', 'p#price', re='the price is (.*)')
```

replace_css(*field_name*, *css*, **processors*, ***kwargs*)

Similar to [add_css\(\)](#) but replaces collected data instead of adding it.

load_item()

Populate the item with the data collected so far, and return it. The data collected is first passed through the [output processors](#) to get the final value to assign to each item field.

nested_xpath(*xpath*)

Create a nested loader with an xpath selector. The supplied selector is applied relative to selector associated with this *ItemLoader*. The nested loader shares the *Item* with the parent *ItemLoader* so calls to [add_xpath\(\)](#), [add_value\(\)](#), [replace_value\(\)](#), etc. will behave as expected.

nested_css(*css*)

Create a nested loader with a css selector. The supplied selector is applied relative to selector associated with this *ItemLoader*. The nested loader shares the *Item* with the parent *ItemLoader* so calls to [add_xpath\(\)](#), [add_value\(\)](#), [replace_value\(\)](#), etc. will behave as expected.

get_collected_values(*field_name*)

Return the collected values for the given field.

get_output_value(*field_name*)

Return the collected values parsed using the output processor, for the given field. This method doesn't populate or modify the item at all.

get_input_processor(*field_name*)

Return the input processor for the given field.

get_output_processor(*field_name*)

Return the output processor for the given field.

ItemLoader instances have the following attributes:

item

The *Item* object being parsed by this Item Loader.

context

The currently active *Context* of this Item Loader.

default_item_class

An Item class (or factory), used to instantiate items when not given in the constructor.

default_input_processor

The default input processor to use for those fields which don't specify one.

default_output_processor

The default output processor to use for those fields which don't specify one.

default_selector_class

The class used to construct the *selector* of this *ItemLoader*, if only a response is given in the constructor. If a selector is given in the constructor this attribute is ignored. This attribute is sometimes overridden in subclasses.

selector

The *Selector* object to extract data from. It's either the selector given in the constructor or one created from the response given in the constructor using the *default_selector_class*. This attribute is meant to be read-only.

3.5.7 Nested Loaders

When parsing related values from a subsection of a document, it can be useful to create nested loaders. Imagine you're extracting details from a footer of a page that looks something like:

Example:

```
<footer>
  <a class="social" href="https://facebook.com/whatever">Like Us</a>
  <a class="social" href="https://twitter.com/whatever">Follow Us</a>
  <a class="email" href="mailto:whatever@example.com">Email Us</a>
</footer>
```

Without nested loaders, you need to specify the full xpath (or css) for each value that you wish to extract.

Example:

```
loader = ItemLoader(item=Item())
# load stuff not in the footer
loader.add_xpath('social', '//footer/a[@class = "social"]/@href')
loader.add_xpath('email', '//footer/a[@class = "email"]/@href')
loader.load_item()
```

Instead, you can create a nested loader with the footer selector and add values relative to the footer. The functionality is the same but you avoid repeating the footer selector.

Example:

```
loader = ItemLoader(item=Item())
# load stuff not in the footer
footer_loader = loader.nested_xpath('//footer')
footer_loader.add_xpath('social', 'a[@class = "social"]/@href')
footer_loader.add_xpath('email', 'a[@class = "email"]/@href')
```

(下页继续)

(续上页)

```
# no need to call footer_loader.load_item()
loader.load_item()
```

You can nest loaders arbitrarily and they work with either xpath or css selectors. As a general guideline, use nested loaders when they make your code simpler but do not go overboard with nesting or your parser can become difficult to read.

3.5.8 Reusing and extending Item Loaders

As your project grows bigger and acquires more and more spiders, maintenance becomes a fundamental problem, especially when you have to deal with many different parsing rules for each spider, having a lot of exceptions, but also wanting to reuse the common processors.

Item Loaders are designed to ease the maintenance burden of parsing rules, without losing flexibility and, at the same time, providing a convenient mechanism for extending and overriding them. For this reason Item Loaders support traditional Python class inheritance for dealing with differences of specific spiders (or groups of spiders).

Suppose, for example, that some particular site encloses their product names in three dashes (e.g. ---Plasma TV---) and you don't want to end up scraping those dashes in the final product names.

Here's how you can remove those dashes by reusing and extending the default Product Item Loader (ProductLoader):

```
from scrapy.loader.processors import MapCompose
from myproject.ItemLoaders import ProductLoader

def strip_dashes(x):
    return x.strip('-')

class SiteSpecificLoader(ProductLoader):
    name_in = MapCompose(strip_dashes, ProductLoader.name_in)
```

Another case where extending Item Loaders can be very helpful is when you have multiple source formats, for example XML and HTML. In the XML version you may want to remove CDATA occurrences. Here's an example of how to do it:

```
from scrapy.loader.processors import MapCompose
from myproject.ItemLoaders import ProductLoader
from myproject.utils.xml import remove_cdata
```

(下页继续)

```
class XmlProductLoader(ProductLoader):
    name_in = MapCompose(remove_cdata, ProductLoader.name_in)
```

And that's how you typically extend input processors.

As for output processors, it is more common to declare them in the field metadata, as they usually depend only on the field and not on each specific site parsing rule (as input processors do). See also: *Declaring Input and Output Processors*.

There are many other possible ways to extend, inherit and override your Item Loaders, and different Item Loaders hierarchies may fit better for different projects. Scrapy only provides the mechanism; it doesn't impose any specific organization of your Loaders collection - that's up to you and your project's needs.

3.5.9 Available built-in processors

Even though you can use any callable function as input and output processors, Scrapy provides some commonly used processors, which are described below. Some of them, like the *MapCompose* (which is typically used as input processor) compose the output of several functions executed in order, to produce the final parsed value.

Here is a list of all built-in processors:

`class scrapy.loader.processors.Identity`

The simplest processor, which doesn't do anything. It returns the original values unchanged. It doesn't receive any constructor arguments, nor does it accept Loader contexts.

Example:

```
>>> from scrapy.loader.processors import Identity
>>> proc = Identity()
>>> proc(['one', 'two', 'three'])
['one', 'two', 'three']
```

`class scrapy.loader.processors.TakeFirst`

Returns the first non-null/non-empty value from the values received, so it's typically used as an output processor to single-valued fields. It doesn't receive any constructor arguments, nor does it accept Loader contexts.

Example:

```
>>> from scrapy.loader.processors import TakeFirst
>>> proc = TakeFirst()
>>> proc(['', 'one', 'two', 'three'])
'one'
```

```
class scrapy.loader.processors.Join(separator=u' ')
```

Returns the values joined with the separator given in the constructor, which defaults to `u' '`. It doesn't accept Loader contexts.

When using the default separator, this processor is equivalent to the function: `u' '.join`

Examples:

```
>>> from scrapy.loader.processors import Join
>>> proc = Join()
>>> proc(['one', 'two', 'three'])
'one two three'
>>> proc = Join('<br>')
>>> proc(['one', 'two', 'three'])
'one<br>two<br>three'
```

```
class scrapy.loader.processors.Compose(*functions, **default_loader_context)
```

A processor which is constructed from the composition of the given functions. This means that each input value of this processor is passed to the first function, and the result of that function is passed to the second function, and so on, until the last function returns the output value of this processor.

By default, stop process on `None` value. This behaviour can be changed by passing keyword argument `stop_on_none=False`.

Example:

```
>>> from scrapy.loader.processors import Compose
>>> proc = Compose(lambda v: v[0], str.upper)
>>> proc(['hello', 'world'])
'HELLO'
```

Each function can optionally receive a `loader_context` parameter. For those which do, this processor will pass the currently active *Loader context* through that parameter.

The keyword arguments passed in the constructor are used as the default Loader context values passed to each function call. However, the final Loader context values passed to functions are overridden with the currently active Loader context accessible through the `ItemLoader.context()` attribute.

```
class scrapy.loader.processors.MapCompose(*functions, **default_loader_context)
```

A processor which is constructed from the composition of the given functions, similar to the *Compose* processor. The difference with this processor is the way internal results are passed among functions, which is as follows:

The input value of this processor is *iterated* and the first function is applied to each element. The results of these function calls (one for each element) are concatenated to construct a new iterable, which is then used to apply the second function, and so on, until the last function is applied to each

value of the list of values collected so far. The output values of the last function are concatenated together to produce the output of this processor.

Each particular function can return a value or a list of values, which is flattened with the list of values returned by the same function applied to the other input values. The functions can also return `None` in which case the output of that function is ignored for further processing over the chain.

This processor provides a convenient way to compose functions that only work with single values (instead of iterables). For this reason the *MapCompose* processor is typically used as input processor, since data is often extracted using the `extract()` method of *selectors*, which returns a list of unicode strings.

The example below should clarify how it works:

```
>>> def filter_world(x):
...     return None if x == 'world' else x
...
>>> from scrapy.loader.processors import MapCompose
>>> proc = MapCompose(filter_world, str.upper)
>>> proc(['hello', 'world', 'this', 'is', 'scrapy'])
['HELLO', 'THIS', 'IS', 'SCRAPY']
```

As with the *Compose* processor, functions can receive Loader contexts, and constructor keyword arguments are used as default context values. See *Compose* processor for more info.

class scrapy.loader.processors.SelectJmes(*json_path*)

Queries the value using the json path provided to the constructor and returns the output. Requires *jmespath* (<https://github.com/jmespath/jmespath.py>) to run. This processor takes only one input at a time.

Example:

```
>>> from scrapy.loader.processors import SelectJmes, Compose, MapCompose
>>> proc = SelectJmes("foo") #for direct use on lists and dictionaries
>>> proc({'foo': 'bar'})
'bar'
>>> proc({'foo': {'bar': 'baz'}})
{'bar': 'baz'}
```

Working with Json:

```
>>> import json
>>> proc_single_json_str = Compose(json.loads, SelectJmes("foo"))
>>> proc_single_json_str('{"foo": "bar"}')
'bar'
```

(下页继续)

(续上页)

```
>>> proc_json_list = Compose(json.loads, MapCompose(SelectJmes('foo')))  
>>> proc_json_list(' [{"foo": "bar"}, {"baz": "tar"} ] ' )  
['bar']
```

3.6 Scrapy shell

The Scrapy shell is an interactive shell where you can try and debug your scraping code very quickly, without having to run the spider. It's meant to be used for testing data extraction code, but you can actually use it for testing any kind of code as it is also a regular Python shell.

The shell is used for testing XPath or CSS expressions and see how they work and what data they extract from the web pages you're trying to scrape. It allows you to interactively test your expressions while you're writing your spider, without having to run the spider to test every change.

Once you get familiarized with the Scrapy shell, you'll see that it's an invaluable tool for developing and debugging your spiders.

3.6.1 Configuring the shell

If you have [IPython](#) installed, the Scrapy shell will use it (instead of the standard Python console). The [IPython](#) console is much more powerful and provides smart auto-completion and colored output, among other things.

We highly recommend you install [IPython](#), specially if you're working on Unix systems (where [IPython](#) excels). See the [IPython installation guide](#) for more info.

Scrapy also has support for [bpython](#), and will try to use it where [IPython](#) is unavailable.

Through scrapy's settings you can configure it to use any one of `ipython`, `bpython` or the standard `python` shell, regardless of which are installed. This is done by setting the `SCRAPY_PYTHON_SHELL` environment variable; or by defining it in your *scrapy.cfg*:

```
[settings]  
shell = bpython
```

3.6.2 Launch the shell

To launch the Scrapy shell you can use the *shell* command like this:

```
scrapy shell <url>
```

Where the `<url>` is the URL you want to scrape.

`shell` also works for local files. This can be handy if you want to play around with a local copy of a web page. `shell` understands the following syntaxes for local files:

```
# UNIX-style
scrapy shell ./path/to/file.html
scrapy shell ../other/path/to/file.html
scrapy shell /absolute/path/to/file.html

# File URI
scrapy shell file:///absolute/path/to/file.html
```

注解: When using relative file paths, be explicit and prepend them with `./` (or `../` when relevant). `scrapy shell index.html` will not work as one might expect (and this is by design, not a bug).

Because `shell` favors HTTP URLs over File URIs, and `index.html` being syntactically similar to `example.com`, `shell` will treat `index.html` as a domain name and trigger a DNS lookup error:

```
$ scrapy shell index.html
[ ... scrapy shell starts ... ]
[ ... traceback ... ]
twisted.internet.error.DNSLookupError: DNS lookup failed:
address 'index.html' not found: [Errno -5] No address associated with hostname.
```

`shell` will not test beforehand if a file called `index.html` exists in the current directory. Again, be explicit.

3.6.3 Using the shell

The Scrapy shell is just a regular Python console (or `IPython` console if you have it available) which provides some additional shortcut functions for convenience.

Available Shortcuts

- `shelp()` - print a help with the list of available objects and shortcuts
- `fetch(url[, redirect=True])` - fetch a new response from the given URL and update all related objects accordingly. You can optionally ask for HTTP 3xx redirections to not be followed by passing `redirect=False`
- `fetch(request)` - fetch a new response from the given request and update all related objects accordingly.

- `view(response)` - open the given response in your local web browser, for inspection. This will add a `<base>` tag to the response body in order for external links (such as images and style sheets) to display properly. Note, however, that this will create a temporary file in your computer, which won't be removed automatically.

Available Scrapy objects

The Scrapy shell automatically creates some convenient objects from the downloaded page, like the *Response* object and the *Selector* objects (for both HTML and XML content).

Those objects are:

- `crawler` - the current *Crawler* object.
- `spider` - the Spider which is known to handle the URL, or a *Spider* object if there is no spider found for the current URL
- `request` - a *Request* object of the last fetched page. You can modify this request using `replace()` or fetch a new request (without leaving the shell) using the `fetch` shortcut.
- `response` - a *Response* object containing the last fetched page
- `settings` - the current *Scrapy settings*

3.6.4 Example of shell session

Here's an example of a typical shell session where we start by scraping the <https://scrapy.org> page, and then proceed to scrape the <https://reddit.com> page. Finally, we modify the (Reddit) request method to POST and re-fetch it getting an error. We end the session by typing Ctrl-D (in Unix systems) or Ctrl-Z in Windows.

Keep in mind that the data extracted here may not be the same when you try it, as those pages are not static and could have changed by the time you test this. The only purpose of this example is to get you familiarized with how the Scrapy shell works.

First, we launch the shell:

```
scrapy shell 'https://scrapy.org' --nolog
```

Then, the shell fetches the URL (using the Scrapy downloader) and prints the list of available objects and useful shortcuts (you'll notice that these lines all start with the `[s]` prefix):

```
[s] Available Scrapy objects:
[s]   scrapy      scrapy module (contains scrapy.Request, scrapy.Selector, etc)
[s]   crawler    <scrapy.crawler.Crawler object at 0x7f07395dd690>
[s]   item       {}
```

(下页继续)

(续上页)

```
[s] request    <GET https://scrapy.org>
[s] response   <200 https://scrapy.org/>
[s] settings   <scrapy.settings.Settings object at 0x7f07395dd710>
[s] spider     <DefaultSpider 'default' at 0x7f0735891690>
[s] Useful shortcuts:
[s] fetch(url[, redirect=True]) Fetch URL and update local objects (by default, ↵
↵ redirects are followed)
[s] fetch(req)                Fetch a scrapy.Request and update local objects
[s] shelp()                   Shell help (print this help)
[s] view(response)           View response in a browser

>>>
```

After that, we can start playing with the objects:

```
>>> response.xpath('//title/text()').get()
'Scrapy | A Fast and Powerful Scraping and Web Crawling Framework'

>>> fetch("https://reddit.com")

>>> response.xpath('//title/text()').get()
'reddit: the front page of the internet'

>>> request = request.replace(method="POST")

>>> fetch(request)

>>> response.status
404

>>> from pprint import pprint

>>> pprint(response.headers)
{'Accept-Ranges': ['bytes'],
 'Cache-Control': ['max-age=0, must-revalidate'],
 'Content-Type': ['text/html; charset=UTF-8'],
 'Date': ['Thu, 08 Dec 2016 16:21:19 GMT'],
 'Server': ['snooserv'],
 'Set-Cookie': ['loid=KqNLou0V9SKMX4qb4n; Domain=reddit.com; Max-Age=63071999; Path=/; ↵
↵ expires=Sat, 08-Dec-2018 16:21:19 GMT; secure'],
```

(下页继续)

(续上页)

```

        'loidcreated=2016-12-08T16%3A21%3A19.445Z; Domain=reddit.com; Max-
↪Age=63071999; Path=/; expires=Sat, 08-Dec-2018 16:21:19 GMT; secure',
        'loid=vi0ZVe4NkxNWd1H7r7; Domain=reddit.com; Max-Age=63071999; Path=/;
↪expires=Sat, 08-Dec-2018 16:21:19 GMT; secure',
        'loidcreated=2016-12-08T16%3A21%3A19.459Z; Domain=reddit.com; Max-
↪Age=63071999; Path=/; expires=Sat, 08-Dec-2018 16:21:19 GMT; secure'],
    'Vary': ['accept-encoding'],
    'Via': ['1.1 varnish'],
    'X-Cache': ['MISS'],
    'X-Cache-Hits': ['0'],
    'X-Content-Type-Options': ['nosniff'],
    'X-Frame-Options': ['SAMEORIGIN'],
    'X-Moose': ['majestic'],
    'X-Served-By': ['cache-cdg8730-CDG'],
    'X-Timer': ['S1481214079.394283,VS0,VE159'],
    'X-Ua-Compatible': ['IE=edge'],
    'X-Xss-Protection': ['1; mode=block']}]
>>>

```

3.6.5 Invoking the shell from spiders to inspect responses

Sometimes you want to inspect the responses that are being processed in a certain point of your spider, if only to check that response you expect is getting there.

This can be achieved by using the `scrapy.shell.inspect_response` function.

Here's an example of how you would call it from your spider:

```

import scrapy

class MySpider(scrapy.Spider):
    name = "myspider"
    start_urls = [
        "http://example.com",
        "http://example.org",
        "http://example.net",
    ]

    def parse(self, response):
        # We want to inspect one specific response.

```

(下页继续)

(续上页)

```

if ".org" in response.url:
    from scrapy.shell import inspect_response
    inspect_response(response, self)

# Rest of parsing code.

```

When you run the spider, you will get something similar to this:

```

2014-01-23 17:48:31-0400 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://example.
↪com> (referer: None)
2014-01-23 17:48:31-0400 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://example.
↪org> (referer: None)
[s] Available Scrapy objects:
[s]   crawler   <scrapy.crawler.Crawler object at 0x1e16b50>
...

>>> response.url
'http://example.org'

```

Then, you can check if the extraction code is working:

```

>>> response.xpath('//h1[@class="fn"]')
[]

```

Nope, it doesn't. So you can open the response in your web browser and see if it's the response you were expecting:

```

>>> view(response)
True

```

Finally you hit Ctrl-D (or Ctrl-Z in Windows) to exit the shell and resume the crawling:

```

>>> ^D
2014-01-23 17:50:03-0400 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://example.
↪net> (referer: None)
...

```

Note that you can't use the `fetch` shortcut here since the Scrapy engine is blocked by the shell. However, after you leave the shell, the spider will continue crawling where it stopped, as shown above.

3.7 Item 管道

After an item has been scraped by a spider, it is sent to the Item Pipeline which processes it through several components that are executed sequentially.

Each item pipeline component (sometimes referred as just “Item Pipeline”) is a Python class that implements a simple method. They receive an item and perform an action over it, also deciding if the item should continue through the pipeline or be dropped and no longer processed.

Typical uses of item pipelines are:

- cleansing HTML data
- validating scraped data (checking that the items contain certain fields)
- checking for duplicates (and dropping them)
- storing the scraped item in a database

3.7.1 Writing your own item pipeline

Each item pipeline component is a Python class that must implement the following method:

`process_item(self, item, spider)`

This method is called for every item pipeline component. `process_item()` must either: return a dict with data, return an *Item* (or any descendant class) object, return a *Twisted Deferred* or raise *DropItem* exception. Dropped items are no longer processed by further pipeline components.

参数

- **item** (*Item* object or a dict) – the item scraped
- **spider** (*Spider* object) – the spider which scraped the item

Additionally, they may also implement the following methods:

`open_spider(self, spider)`

This method is called when the spider is opened.

参数 **spider** (*Spider* object) – the spider which was opened

`close_spider(self, spider)`

This method is called when the spider is closed.

参数 **spider** (*Spider* object) – the spider which was closed

`from_crawler(cls, crawler)`

If present, this classmethod is called to create a pipeline instance from a *Crawler*. It must return a new instance of the pipeline. Crawler object provides access to all Scrapy core components like settings and signals; it is a way for pipeline to access them and hook its functionality into Scrapy.

参数 `crawler` (*Crawler* object) – crawler that uses this pipeline

3.7.2 Item pipeline example

Price validation and dropping items with no prices

Let's take a look at the following hypothetical pipeline that adjusts the `price` attribute for those items that do not include VAT (`price_excludes_vat` attribute), and drops those items which don't contain a price:

```
from scrapy.exceptions import DropItem

class PricePipeline(object):

    vat_factor = 1.15

    def process_item(self, item, spider):
        if item.get('price'):
            if item.get('price_excludes_vat'):
                item['price'] = item['price'] * self.vat_factor
            return item
        else:
            raise DropItem("Missing price in %s" % item)
```

Write items to a JSON file

The following pipeline stores all scraped items (from all spiders) into a single `items.jsonl` file, containing one item per line serialized in JSON format:

```
import json

class JsonWriterPipeline(object):

    def open_spider(self, spider):
        self.file = open('items.jsonl', 'w')

    def close_spider(self, spider):
        self.file.close()

    def process_item(self, item, spider):
        line = json.dumps(dict(item)) + "\n"
```

(下页继续)

(续上页)

```
self.file.write(line)
return item
```

注解: The purpose of `JsonWriterPipeline` is just to introduce how to write item pipelines. If you really want to store all scraped items into a JSON file you should use the *Feed exports*.

Write items to MongoDB

In this example we'll write items to [MongoDB](#) using [pymongo](#). MongoDB address and database name are specified in Scrapy settings; MongoDB collection is named after item class.

The main point of this example is to show how to use *from_crawler()* method and how to clean up the resources properly.:

```
import pymongo

class MongoPipeline(object):

    collection_name = 'scrapy_items'

    def __init__(self, mongo_uri, mongo_db):
        self.mongo_uri = mongo_uri
        self.mongo_db = mongo_db

    @classmethod
    def from_crawler(cls, crawler):
        return cls(
            mongo_uri=crawler.settings.get('MONGO_URI'),
            mongo_db=crawler.settings.get('MONGO_DATABASE', 'items')
        )

    def open_spider(self, spider):
        self.client = pymongo.MongoClient(self.mongo_uri)
        self.db = self.client[self.mongo_db]

    def close_spider(self, spider):
        self.client.close()

    def process_item(self, item, spider):
```

(下页继续)

(续上页)

```
self.db[self.collection_name].insert_one(dict(item))
return item
```

Take screenshot of item

This example demonstrates how to return `Deferred` from `process_item()` method. It uses `Splash` to render screenshot of item url. Pipeline makes request to locally running instance of `Splash`. After request is downloaded and `Deferred` callback fires, it saves item to a file and adds filename to an item.

```
import scrapy
import hashlib
from urllib.parse import quote

class ScreenshotPipeline(object):
    """Pipeline that uses Splash to render screenshot of
    every Scrapy item."""

    SPLASH_URL = "http://localhost:8050/render.png?url={}"

    def process_item(self, item, spider):
        encoded_item_url = quote(item["url"])
        screenshot_url = self.SPLASH_URL.format(encoded_item_url)
        request = scrapy.Request(screenshot_url)
        dfd = spider.crawler.engine.download(request, spider)
        dfd.addBoth(self.return_item, item)
        return dfd

    def return_item(self, response, item):
        if response.status != 200:
            # Error happened, return item.
            return item

        # Save screenshot to file, filename will be hash of url.
        url = item["url"]
        url_hash = hashlib.md5(url.encode("utf8")).hexdigest()
        filename = "{}.png".format(url_hash)
        with open(filename, "wb") as f:
            f.write(response.body)
```

(下页继续)

(续上页)

```
# Store filename in item.
item["screenshot_filename"] = filename
return item
```

Duplicates filter

A filter that looks for duplicate items, and drops those items that were already processed. Let's say that our items have a unique id, but our spider returns multiples items with the same id:

```
from scrapy.exceptions import DropItem

class DuplicatesPipeline(object):

    def __init__(self):
        self.ids_seen = set()

    def process_item(self, item, spider):
        if item['id'] in self.ids_seen:
            raise DropItem("Duplicate item found: %s" % item)
        else:
            self.ids_seen.add(item['id'])
            return item
```

3.7.3 Activating an Item Pipeline component

To activate an Item Pipeline component you must add its class to the `ITEM_PIPELINES` setting, like in the following example:

```
ITEM_PIPELINES = {
    'myproject.pipelines.PricePipeline': 300,
    'myproject.pipelines.JsonWriterPipeline': 800,
}
```

The integer values you assign to classes in this setting determine the order in which they run: items go through from lower valued to higher valued classes. It's customary to define these numbers in the 0-1000 range.

3.8 Feed 导出

0.10 新版功能.

One of the most frequently required features when implementing scrapers is being able to store the scraped data properly and, quite often, that means generating an “export file” with the scraped data (commonly called “export feed”) to be consumed by other systems.

Scrapy provides this functionality out of the box with the Feed Exports, which allows you to generate a feed with the scraped items, using multiple serialization formats and storage backends.

3.8.1 Serialization formats

For serializing the scraped data, the feed exports use the *Item exporters*. These formats are supported out of the box:

- *JSON*
- *JSON lines*
- *CSV*
- *XML*

But you can also extend the supported format through the *FEED_EXPORTERS* setting.

JSON

- *FEED_FORMAT*: `json`
- Exporter used: *JsonItemExporter*
- See *this warning* if you’ re using JSON with large feeds.

JSON lines

- *FEED_FORMAT*: `jsonlines`
- Exporter used: *JsonLinesItemExporter*

CSV

- *FEED_FORMAT*: `csv`
- Exporter used: *CsvItemExporter*

- To specify columns to export and their order use `FEED_EXPORT_FIELDS`. Other feed exporters can also use this option, but it is important for CSV because unlike many other export formats CSV uses a fixed header.

XML

- `FEED_FORMAT`: `xml`
- Exporter used: `XmlItemExporter`

Pickle

- `FEED_FORMAT`: `pickle`
- Exporter used: `PickleItemExporter`

Marshal

- `FEED_FORMAT`: `marshal`
- Exporter used: `MarshalItemExporter`

3.8.2 Storages

When using the feed exports you define where to store the feed using a `URI` (through the `FEED_URI` setting). The feed exports supports multiple storage backend types which are defined by the URI scheme.

The storages backends supported out of the box are:

- *Local filesystem*
- *FTP*
- *S3* (requires `botocore` or `boto`)
- *Standard output*

Some storage backends may be unavailable if the required external libraries are not available. For example, the S3 backend is only available if the `botocore` or `boto` library is installed (Scrapy supports `boto` only on Python 2).

3.8.3 Storage URI parameters

The storage URI can also contain parameters that get replaced when the feed is being created. These parameters are:

- `%(time)s` - gets replaced by a timestamp when the feed is being created

- `%(name)s` - gets replaced by the spider name

Any other named parameter gets replaced by the spider attribute of the same name. For example, `%(site_id)s` would get replaced by the `spider.site_id` attribute the moment the feed is being created.

Here are some examples to illustrate:

- Store in FTP using one directory per spider:
 - `ftp://user:password@ftp.example.com/scraping/feeds/%(name)s/%(time)s.json`
- Store in S3 using one directory per spider:
 - `s3://mybucket/scraping/feeds/%(name)s/%(time)s.json`

3.8.4 Storage backends

Local filesystem

The feeds are stored in the local filesystem.

- URI scheme: `file`
- Example URI: `file:///tmp/export.csv`
- Required external libraries: none

Note that for the local filesystem storage (only) you can omit the scheme if you specify an absolute path like `/tmp/export.csv`. This only works on Unix systems though.

FTP

The feeds are stored in a FTP server.

- URI scheme: `ftp`
- Example URI: `ftp://user:pass@ftp.example.com/path/to/export.csv`
- Required external libraries: none

S3

The feeds are stored on [Amazon S3](#).

- URI scheme: `s3`
- Example URIs:
 - `s3://mybucket/path/to/export.csv`
 - `s3://aws_key:aws_secret@mybucket/path/to/export.csv`

- Required external libraries: `botocore` (Python 2 and Python 3) or `boto` (Python 2 only)

The AWS credentials can be passed as user/password in the URI, or they can be passed through the following settings:

- `AWS_ACCESS_KEY_ID`
- `AWS_SECRET_ACCESS_KEY`

You can also define a custom ACL for exported feeds using this setting:

- `FEED_STORAGE_S3_ACL`

Standard output

The feeds are written to the standard output of the Scrapy process.

- URI scheme: `stdout`
- Example URI: `stdout:`
- Required external libraries: `none`

3.8.5 Settings

These are the settings used for configuring the feed exports:

- `FEED_URI` (mandatory)
- `FEED_FORMAT`
- `FEED_STORAGES`
- `FEED_STORAGE_S3_ACL`
- `FEED_EXPORTERS`
- `FEED_STORE_EMPTY`
- `FEED_EXPORT_ENCODING`
- `FEED_EXPORT_FIELDS`
- `FEED_EXPORT_INDENT`

FEED_URI

Default: `None`

The URI of the export feed. See *Storage backends* for supported URI schemes.

This setting is required for enabling the feed exports.

FEED_FORMAT

The serialization format to be used for the feed. See *Serialization formats* for possible values.

FEED_EXPORT_ENCODING

Default: `None`

The encoding to be used for the feed.

If unset or set to `None` (default) it uses UTF-8 for everything except JSON output, which uses safe numeric encoding (`\uXXXX` sequences) for historic reasons.

Use `utf-8` if you want UTF-8 for JSON too.

FEED_EXPORT_FIELDS

Default: `None`

A list of fields to export, optional. Example: `FEED_EXPORT_FIELDS = ["foo", "bar", "baz"]`.

Use `FEED_EXPORT_FIELDS` option to define fields to export and their order.

When `FEED_EXPORT_FIELDS` is empty or `None` (default), Scrapy uses fields defined in dicts or *Item* subclasses a spider is yielding.

If an exporter requires a fixed set of fields (this is the case for *CSV* export format) and `FEED_EXPORT_FIELDS` is empty or `None`, then Scrapy tries to infer field names from the exported data - currently it uses field names from the first item.

FEED_EXPORT_INDENT

Default: `0`

Amount of spaces used to indent the output on each level. If `FEED_EXPORT_INDENT` is a non-negative integer, then array elements and object members will be pretty-printed with that indent level. An indent level of `0` (the default), or negative, will put each item on a new line. `None` selects the most compact representation.

Currently implemented only by *JsonItemExporter* and *XmlItemExporter*, i.e. when you are exporting to `.json` or `.xml`.

FEED_STORE_EMPTY

Default: `False`

Whether to export empty feeds (ie. feeds with no items).

FEED_STORAGES

Default: {}

A dict containing additional feed storage backends supported by your project. The keys are URI schemes and the values are paths to storage classes.

FEED_STORAGE_S3_ACL

Default: '' (empty string)

A string containing a custom ACL for feeds exported to Amazon S3 by your project.

For a complete list of available values, access the [Canned ACL](#) section on Amazon S3 docs.

FEED_STORAGES_BASE

Default:

```
{
    '': 'scrapy.extensions.feedexport.FileFeedStorage',
    'file': 'scrapy.extensions.feedexport.FileFeedStorage',
    'stdout': 'scrapy.extensions.feedexport.StdoutFeedStorage',
    's3': 'scrapy.extensions.feedexport.S3FeedStorage',
    'ftp': 'scrapy.extensions.feedexport.FTPFeedStorage',
}
```

A dict containing the built-in feed storage backends supported by Scrapy. You can disable any of these backends by assigning `None` to their URI scheme in `FEED_STORAGES`. E.g., to disable the built-in FTP storage backend (without replacement), place this in your `settings.py`:

```
FEED_STORAGES = {
    'ftp': None,
}
```

FEED_EXPORTERS

Default: {}

A dict containing additional exporters supported by your project. The keys are serialization formats and the values are paths to *Item exporter* classes.

FEED_EXPORTERS_BASE

Default:

```
{
    'json': 'scrapy.exporters.JsonItemExporter',
    'jsonlines': 'scrapy.exporters.JsonLinesItemExporter',
    'jl': 'scrapy.exporters.JsonLinesItemExporter',
    'csv': 'scrapy.exporters.CsvItemExporter',
    'xml': 'scrapy.exporters.XmlItemExporter',
    'marshal': 'scrapy.exporters.MarshalItemExporter',
    'pickle': 'scrapy.exporters.PickleItemExporter',
}
```

A dict containing the built-in feed exporters supported by Scrapy. You can disable any of these exporters by assigning `None` to their serialization format in `FEED_EXPORTERS`. E.g., to disable the built-in CSV exporter (without replacement), place this in your `settings.py`:

```
FEED_EXPORTERS = {
    'csv': None,
}
```

3.9 请求和响应

Scrapy uses *Request* and *Response* objects for crawling web sites.

Typically, *Request* objects are generated in the spiders and pass across the system until they reach the Downloader, which executes the request and returns a *Response* object which travels back to the spider that issued the request.

Both *Request* and *Response* classes have subclasses which add functionality not required in the base classes. These are described below in *Request subclasses* and *Response subclasses*.

3.9.1 Request objects

```
class scrapy.http.Request(url[, callback, method='GET', headers, body, cookies, meta,
                           encoding='utf-8', priority=0, dont_filter=False, errback, flags])
```

A *Request* object represents an HTTP request, which is usually generated in the Spider and executed by the Downloader, and thus generating a *Response*.

参数

- `url` (*string*) – the URL of this request

- **callback** (*callable*) – the function that will be called with the response of this request (once its downloaded) as its first parameter. For more information see [Passing additional data to callback functions](#) below. If a Request doesn't specify a callback, the spider's `parse()` method will be used. Note that if exceptions are raised during processing, `errback` is called instead.
- **method** (*string*) – the HTTP method of this request. Defaults to 'GET'.
- **meta** (*dict*) – the initial values for the `Request.meta` attribute. If given, the dict passed in this parameter will be shallow copied.
- **body** (*str or unicode*) – the request body. If a `unicode` is passed, then it's encoded to `str` using the `encoding` passed (which defaults to `utf-8`). If `body` is not given, an empty string is stored. Regardless of the type of this argument, the final value stored will be a `str` (never `unicode` or `None`).
- **headers** (*dict*) – the headers of this request. The dict values can be strings (for single valued headers) or lists (for multi-valued headers). If `None` is passed as value, the HTTP header will not be sent at all.
- **cookies** (*dict or list*) – the request cookies. These can be sent in two forms.

1. Using a dict:

```
request_with_cookies = Request(url="http://www.example.com",
                               cookies={'currency': 'USD', 'country': 'UY'})
```

2. Using a list of dicts:

```
request_with_cookies = Request(url="http://www.example.com",
                               cookies=[{'name': 'currency',
                                         'value': 'USD',
                                         'domain': 'example.com',
                                         'path': '/currency'}])
```

The latter form allows for customizing the `domain` and `path` attributes of the cookie. This is only useful if the cookies are saved for later requests. When some site returns cookies (in a response) those are stored in the cookies for that domain and will be sent again in future requests. That's the typical behaviour of any regular web browser. However, if, for some reason, you want to avoid merging with existing cookies you can instruct Scrapy to do so by setting the `dont_merge_cookies` key to `True` in the `Request.meta`.

Example of request without merging cookies:

```
request_with_cookies = Request(url="http://www.example.com",
                                cookies={'currency': 'USD', 'country': 'UY'},
                                meta={'dont_merge_cookies': True})
```

For more info see [CookiesMiddleware](#).

- **encoding** (*string*) – the encoding of this request (defaults to 'utf-8'). This encoding will be used to percent-encode the URL and to convert the body to `str` (if given as `unicode`).
- **priority** (*int*) – the priority of this request (defaults to 0). The priority is used by the scheduler to define the order used to process requests. Requests with a higher priority value will execute earlier. Negative values are allowed in order to indicate relatively low-priority.
- **dont_filter** (*boolean*) – indicates that this request should not be filtered by the scheduler. This is used when you want to perform an identical request multiple times, to ignore the duplicates filter. Use it with care, or you will get into crawling loops. Default to `False`.
- **errback** (*callable*) – a function that will be called if any exception was raised while processing the request. This includes pages that failed with 404 HTTP errors and such. It receives a `Twisted Failure` instance as first parameter. For more information, see [Using errbacks to catch exceptions in request processing](#) below.
- **flags** (*list*) – Flags sent to the request, can be used for logging or similar purposes.

url

A string containing the URL of this request. Keep in mind that this attribute contains the escaped URL, so it can differ from the URL passed in the constructor.

This attribute is read-only. To change the URL of a Request use [replace\(\)](#).

method

A string representing the HTTP method in the request. This is guaranteed to be uppercase. Example: "GET", "POST", "PUT", etc

headers

A dictionary-like object which contains the request headers.

body

A str that contains the request body.

This attribute is read-only. To change the body of a Request use [replace\(\)](#).

meta

A dict that contains arbitrary metadata for this request. This dict is empty for new Requests,

and is usually populated by different Scrapy components (extensions, middlewares, etc). So the data contained in this dict depends on the extensions you have enabled.

See *Request.meta special keys* for a list of special meta keys recognized by Scrapy.

This dict is *shallow copied* when the request is cloned using the `copy()` or `replace()` methods, and can also be accessed, in your spider, from the `response.meta` attribute.

`copy()`

Return a new Request which is a copy of this Request. See also: *Passing additional data to callback functions*.

`replace([url, method, headers, body, cookies, meta, encoding, dont_filter, callback, errback])`

Return a Request object with the same members, except for those members given new values by whichever keyword arguments are specified. The attribute *Request.meta* is copied by default (unless a new value is given in the `meta` argument). See also *Passing additional data to callback functions*.

Passing additional data to callback functions

The callback of a request is a function that will be called when the response of that request is downloaded. The callback function will be called with the downloaded *Response* object as its first argument.

Example:

```
def parse_page1(self, response):
    return scrapy.Request("http://www.example.com/some_page.html",
                           callback=self.parse_page2)

def parse_page2(self, response):
    # this would log http://www.example.com/some_page.html
    self.logger.info("Visited %s", response.url)
```

In some cases you may be interested in passing arguments to those callback functions so you can receive the arguments later, in the second callback. You can use the *Request.meta* attribute for that.

Here's an example of how to pass an item using this mechanism, to populate different fields from different pages:

```
def parse_page1(self, response):
    item = MyItem()
    item['main_url'] = response.url
    request = scrapy.Request("http://www.example.com/some_page.html",
                              callback=self.parse_page2)
    request.meta['item'] = item
```

(下页继续)

(续上页)

```
yield request

def parse_page2(self, response):
    item = response.meta['item']
    item['other_url'] = response.url
    yield item
```

Using errbacks to catch exceptions in request processing

The errback of a request is a function that will be called when an exception is raised while processing it.

It receives a `Twisted Failure` instance as first parameter and can be used to track connection establishment timeouts, DNS errors etc.

Here's an example spider logging all errors and catching some specific errors if needed:

```
import scrapy

from scrapy.spidermiddlewares.httperror import HttpError
from twisted.internet.error import DNSLookupError
from twisted.internet.error import TimeoutError, TCPTimedOutError

class ErrbackSpider(scrapy.Spider):
    name = "errback_example"
    start_urls = [
        "http://www.httpbin.org/",          # HTTP 200 expected
        "http://www.httpbin.org/status/404", # Not found error
        "http://www.httpbin.org/status/500", # server issue
        "http://www.httpbin.org:12345/",     # non-responding host, timeout expected
        "http://www.httphttpbinbin.org/",    # DNS error expected
    ]

    def start_requests(self):
        for u in self.start_urls:
            yield scrapy.Request(u, callback=self.parse_httpbin,
                                errback=self.errback_httpbin,
                                dont_filter=True)

    def parse_httpbin(self, response):
        self.logger.info('Got successful response from {}'.format(response.url))
        # do something useful here...
```

(下页继续)

(续上页)

```
def errback_httpbin(self, failure):
    # log all failures
    self.logger.error(repr(failure))

    # in case you want to do something special for some errors,
    # you may need the failure's type:

    if failure.check(HttpError):
        # these exceptions come from HttpError spider middleware
        # you can get the non-200 response
        response = failure.value.response
        self.logger.error('HttpError on %s', response.url)

    elif failure.check(DNSLookupError):
        # this is the original request
        request = failure.request
        self.logger.error('DNSLookupError on %s', request.url)

    elif failure.check(TimeoutError, TCPTimedOutError):
        request = failure.request
        self.logger.error('TimeoutError on %s', request.url)
```

3.9.2 Request.meta special keys

The `Request.meta` attribute can contain any arbitrary data, but there are some special keys recognized by Scrapy and its built-in extensions.

Those are:

- `dont_redirect`
- `dont_retry`
- `handle_httpstatus_list`
- `handle_httpstatus_all`
- `dont_merge_cookies`
- `cookiejar`
- `dont_cache`
- `redirect_reasons`

- *redirect_urls*
- *bindaddress*
- *dont_obey_robotstxt*
- *download_timeout*
- *download_maxsize*
- *download_latency*
- *download_fail_on_dataloss*
- *proxy*
- *ftp_user* (See *FTP_USER* for more info)
- *ftp_password* (See *FTP_PASSWORD* for more info)
- *referrer_policy*
- *max_retry_times*

bindaddress

The IP of the outgoing IP address to use for the performing the request.

download_timeout

The amount of time (in secs) that the downloader will wait before timing out. See also: *DOWNLOAD_TIMEOUT*.

download_latency

The amount of time spent to fetch the response, since the request has been started, i.e. HTTP message sent over the network. This meta key only becomes available when the response has been downloaded. While most other meta keys are used to control Scrapy behavior, this one is supposed to be read-only.

download_fail_on_dataloss

Whether or not to fail on broken responses. See: *DOWNLOAD_FAIL_ON_DATALOSS*.

max_retry_times

The meta key is used set retry times per request. When initialized, the *max_retry_times* meta key takes higher precedence over the *RETRY_TIMES* setting.

3.9.3 Request subclasses

Here is the list of built-in *Request* subclasses. You can also subclass it to implement your own custom functionality.

FormRequest objects

The *FormRequest* class extends the base *Request* with functionality for dealing with HTML forms. It uses *lxml.html forms* to pre-populate form fields with form data from *Response* objects.

```
class scrapy.http.FormRequest(url[, formdata, ...])
```

The *FormRequest* class adds a new argument to the constructor. The remaining arguments are the same as for the *Request* class and are not documented here.

参数 *formdata* (*dict or iterable of tuples*) – is a dictionary (or iterable of (key, value) tuples) containing HTML Form data which will be url-encoded and assigned to the body of the request.

The *FormRequest* objects support the following class method in addition to the standard *Request* methods:

```
classmethod from_response(response[, formname=None, formid=None, formnumber=0, formdata=None, formxpath=None, formcss=None, clickdata=None, dont_click=False, ...])
```

Returns a new *FormRequest* object with its form field values pre-populated with those found in the HTML `<form>` element contained in the given response. For an example see *Using FormRequest.from_response() to simulate a user login*.

The policy is to automatically simulate a click, by default, on any form control that looks clickable, like a `<input type="submit">`. Even though this is quite convenient, and often the desired behaviour, sometimes it can cause problems which could be hard to debug. For example, when working with forms that are filled and/or submitted using javascript, the default *from_response()* behaviour may not be the most appropriate. To disable this behaviour you can set the `dont_click` argument to `True`. Also, if you want to change the control clicked (instead of disabling it) you can also use the `clickdata` argument.

警告: Using this method with select elements which have leading or trailing whitespace in the option values will not work due to a [bug in lxml](#), which should be fixed in lxml 3.8 and above.

参数

- **response** (*Response* object) – the response containing a HTML form which will be used to pre-populate the form fields

- **formname** (*string*) – if given, the form with name attribute set to this value will be used.
- **formid** (*string*) – if given, the form with id attribute set to this value will be used.
- **formxpath** (*string*) – if given, the first form that matches the xpath will be used.
- **formcss** (*string*) – if given, the first form that matches the css selector will be used.
- **formnumber** (*integer*) – the number of form to use, when the response contains multiple forms. The first one (and also the default) is 0.
- **formdata** (*dict*) – fields to override in the form data. If a field was already present in the response `<form>` element, its value is overridden by the one passed in this parameter. If a value passed in this parameter is `None`, the field will not be included in the request, even if it was present in the response `<form>` element.
- **clickdata** (*dict*) – attributes to lookup the control clicked. If it's not given, the form data will be submitted simulating a click on the first clickable element. In addition to html attributes, the control can be identified by its zero-based index relative to other submittable inputs inside the form, via the `nr` attribute.
- **dont_click** (*boolean*) – If True, the form data will be submitted without clicking in any element.

The other parameters of this class method are passed directly to the *FormRequest* constructor.

0.10.3 新版功能: The **formname** parameter.

0.17 新版功能: The **formxpath** parameter.

1.1.0 新版功能: The **formcss** parameter.

1.1.0 新版功能: The **formid** parameter.

Request usage examples

Using FormRequest to send data via HTTP POST

If you want to simulate a HTML Form POST in your spider and send a couple of key-value fields, you can return a *FormRequest* object (from your spider) like this:

```
return [FormRequest(url="http://www.example.com/post/action",
                    formdata={'name': 'John Doe', 'age': '27'},
                    callback=self.after_post)]
```

Using `FormRequest.from_response()` to simulate a user login

It is usual for web sites to provide pre-populated form fields through `<input type="hidden">` elements, such as session related data or authentication tokens (for login pages). When scraping, you'll want these fields to be automatically pre-populated and only override a couple of them, such as the user name and password. You can use the `FormRequest.from_response()` method for this job. Here's an example spider which uses it:

```
import scrapy

def authentication_failed(response):
    # TODO: Check the contents of the response and return True if it failed
    # or False if it succeeded.
    pass

class LoginSpider(scrapy.Spider):
    name = 'example.com'
    start_urls = ['http://www.example.com/users/login.php']

    def parse(self, response):
        return scrapy.FormRequest.from_response(
            response,
            formdata={'username': 'john', 'password': 'secret'},
            callback=self.after_login
        )

    def after_login(self, response):
        if authentication_failed(response):
            self.logger.error("Login failed")
            return

        # continue scraping with authenticated session...
```

JSONRequest

The `JSONRequest` class extends the base `Request` class with functionality for dealing with JSON requests.

```
class scrapy.http.JSONRequest(url[, ... data, dumps_kwargs])
```

The `JSONRequest` class adds two new argument to the constructor. The remaining arguments are the same as for the `Request` class and are not documented here.

Using the `JSONRequest` will set the Content-Type header to `application/json` and Accept header

to application/json, text/javascript, */*; q=0.01

参数

- **data** (*JSON serializable object*) – is any JSON serializable object that needs to be JSON encoded and assigned to body. if *Request.body* argument is provided this parameter will be ignored. if *Request.body* argument is not provided and data argument is provided *Request.method* will be set to 'POST' automatically.
- **dumps_kwargs** (*dict*) – Parameters that will be passed to underlying *json.dumps* method which is used to serialize data into JSON format.

JSONRequest usage example

Sending a JSON POST request with a JSON payload:

```
data = {
    'name1': 'value1',
    'name2': 'value2',
}
yield JsonRequest(url='http://www.example.com/post/action', data=data)
```

3.9.4 Response objects

```
class scrapy.http.Response(url[, status=200, headers=None, body=b'', flags=None, request=None])
```

A *Response* object represents an HTTP response, which is usually downloaded (by the Downloader) and fed to the Spiders for processing.

参数

- **url** (*string*) – the URL of this response
- **status** (*integer*) – the HTTP status of the response. Defaults to 200.
- **headers** (*dict*) – the headers of this response. The dict values can be strings (for single valued headers) or lists (for multi-valued headers).
- **body** (*bytes*) – the response body. To access the decoded text as str (unicode in Python 2) you can use *response.text* from an encoding-aware *Response subclass*, such as *TextResponse*.
- **flags** (*list*) – is a list containing the initial values for the *Response.flags* attribute. If given, the list will be shallow copied.
- **request** (*Request object*) – the initial value of the *Response.request* attribute. This represents the *Request* that generated this response.

url

A string containing the URL of the response.

This attribute is read-only. To change the URL of a Response use `replace()`.

status

An integer representing the HTTP status of the response. Example: 200, 404.

headers

A dictionary-like object which contains the response headers. Values can be accessed using `get()` to return the first header value with the specified name or `getlist()` to return all header values with the specified name. For example, this call will give you all cookies in the headers:

```
response.headers.getlist('Set-Cookie')
```

body

The body of this Response. Keep in mind that `Response.body` is always a bytes object. If you want the unicode version use `TextResponse.text` (only available in `TextResponse` and subclasses).

This attribute is read-only. To change the body of a Response use `replace()`.

request

The `Request` object that generated this response. This attribute is assigned in the Scrapy engine, after the response and the request have passed through all `Downloader Middlewares`. In particular, this means that:

- HTTP redirections will cause the original request (to the URL before redirection) to be assigned to the redirected response (with the final URL after redirection).
- `Response.request.url` doesn't always equal `Response.url`
- This attribute is only available in the spider code, and in the `Spider Middlewares`, but not in `Downloader Middlewares` (although you have the `Request` available there by other means) and handlers of the `response_downloaded` signal.

meta

A shortcut to the `Request.meta` attribute of the `Response.request` object (ie. `self.request.meta`).

Unlike the `Response.request` attribute, the `Response.meta` attribute is propagated along redirects and retries, so you will get the original `Request.meta` sent from your spider.

参见:

`Request.meta` attribute

flags

A list that contains flags for this response. Flags are labels used for tagging Responses. For example: `'cached'`, `'redirected'`, etc. And they're shown on the string representation of the Response (`__str__` method) which is used by the engine for logging.

`copy()`

Returns a new Response which is a copy of this Response.

`replace([url, status, headers, body, request, flags, cls])`

Returns a Response object with the same members, except for those members given new values by whichever keyword arguments are specified. The attribute `Response.meta` is copied by default.

`urljoin(url)`

Constructs an absolute url by combining the Response's `url` with a possible relative url.

This is a wrapper over `urlparse.urljoin`, it's merely an alias for making this call:

```
urlparse.urljoin(response.url, url)
```

3.9.5 Response subclasses

Here is the list of available built-in Response subclasses. You can also subclass the Response class to implement your own functionality.

TextResponse objects

`class scrapy.http.TextResponse(url[, encoding[, ...]])`

`TextResponse` objects adds encoding capabilities to the base `Response` class, which is meant to be used only for binary data, such as images, sounds or any media file.

`TextResponse` objects support a new constructor argument, in addition to the base `Response` objects. The remaining functionality is the same as for the `Response` class and is not documented here.

参数 encoding (string) – is a string which contains the encoding to use for this response.

If you create a `TextResponse` object with a unicode body, it will be encoded using this encoding (remember the body attribute is always a string). If `encoding` is `None` (default value), the encoding will be looked up in the response headers and body instead.

`TextResponse` objects support the following attributes in addition to the standard `Response` ones:

text

Response body, as unicode.

The same as `response.body.decode(response.encoding)`, but the result is cached after the first call, so you can access `response.text` multiple times without extra overhead.

注解: `unicode(response.body)` is not a correct way to convert response body to unicode: you would be using the system default encoding (typically `ascii`) instead of the response encoding.

encoding

A string with the encoding of this response. The encoding is resolved by trying the following mechanisms, in order:

1. the encoding passed in the constructor `encoding` argument
2. the encoding declared in the Content-Type HTTP header. If this encoding is not valid (ie. unknown), it is ignored and the next resolution mechanism is tried.
3. the encoding declared in the response body. The `TextResponse` class doesn't provide any special functionality for this. However, the `HtmlResponse` and `XmlResponse` classes do.
4. the encoding inferred by looking at the response body. This is the more fragile method but also the last one tried.

selector

A `Selector` instance using the response as target. The selector is lazily instantiated on first access.

`TextResponse` objects support the following methods in addition to the standard `Response` ones:

xpath(query)

A shortcut to `TextResponse.selector.xpath(query)`:

```
response.xpath('//p')
```

css(query)

A shortcut to `TextResponse.selector.css(query)`:

```
response.css('p')
```

body_as_unicode()

The same as `text`, but available as a method. This method is kept for backward compatibility; please prefer `response.text`.

HtmlResponse objects

```
class scrapy.http.HtmlResponse(url[, ...])
```

The `HtmlResponse` class is a subclass of `TextResponse` which adds encoding auto-discovering support by looking into the HTML meta `http-equiv` attribute. See `TextResponse.encoding`.

XmlResponse objects

```
class scrapy.http.XmlResponse(url[, ...])
```

The `XmlResponse` class is a subclass of `TextResponse` which adds encoding auto-discovering support by looking into the XML declaration line. See `TextResponse.encoding`.

3.10 链接提取

Link extractors are objects whose only purpose is to extract links from web pages (*scrapy.http.Response* objects) which will be eventually followed.

There is `scrapy.linkextractors.LinkExtractor` available in Scrapy, but you can create your own custom Link Extractors to suit your needs by implementing a simple interface.

The only public method that every link extractor has is `extract_links`, which receives a *Response* object and returns a list of `scrapy.link.Link` objects. Link extractors are meant to be instantiated once and their `extract_links` method called several times with different responses to extract links to follow.

Link extractors are used in the *CrawlSpider* class (available in Scrapy), through a set of rules, but you can also use it in your spiders, even if you don't subclass from *CrawlSpider*, as its purpose is very simple: to extract links.

3.10.1 Built-in link extractors reference

Link extractors classes bundled with Scrapy are provided in the *scrapy.linkextractors* module.

The default link extractor is `LinkExtractor`, which is the same as *LxmlLinkExtractor*:

```
from scrapy.linkextractors import LinkExtractor
```

There used to be other link extractor classes in previous Scrapy versions, but they are deprecated now.

LxmlLinkExtractor

```
class scrapy.linkextractors.lxmlhtml.LxmlLinkExtractor(allow=(), deny=(),
                                                         allow_domains=(),
                                                         deny_domains=(),
                                                         deny_extensions=None, re-
strict_xpaths=(), restrict_css=(),
                                                         tags=('a', 'area'), attrs=('href', ),
                                                         canonicalize=False, unique=True,
                                                         process_value=None, strip=True)
```

`LxmlLinkExtractor` is the recommended link extractor with handy filtering options. It is implemented using `lxml`'s robust `HTMLParser`.

参数

- `allow` (*a regular expression (or list of)*) – a single regular expression (or list of regular expressions) that the (absolute) urls must match in order to be extracted. If not given (or empty), it will match all links.

- **deny** (*a regular expression (or list of)*) – a single regular expression (or list of regular expressions) that the (absolute) urls must match in order to be excluded (ie. not extracted). It has precedence over the **allow** parameter. If not given (or empty) it won't exclude any links.
- **allow_domains** (*str or list*) – a single value or a list of string containing domains which will be considered for extracting the links
- **deny_domains** (*str or list*) – a single value or a list of strings containing domains which won't be considered for extracting the links
- **deny_extensions** (*list*) – a single value or list of strings containing extensions that should be ignored when extracting links. If not given, it will default to the `IGNORED_EXTENSIONS` list defined in the `scrapy.linkextractors` package.
- **restrict_xpaths** (*str or list*) – is an XPath (or list of XPath's) which defines regions inside the response where links should be extracted from. If given, only the text selected by those XPath will be scanned for links. See examples below.
- **restrict_css** (*str or list*) – a CSS selector (or list of selectors) which defines regions inside the response where links should be extracted from. Has the same behaviour as **restrict_xpaths**.
- **restrict_text** (*a regular expression (or list of)*) – a single regular expression (or list of regular expressions) that the link's text must match in order to be extracted. If not given (or empty), it will match all links. If a list of regular expressions is given, the link will be extracted if it matches at least one.
- **tags** (*str or list*) – a tag or a list of tags to consider when extracting links. Defaults to ('a', 'area').
- **attrs** (*list*) – an attribute or list of attributes which should be considered when looking for links to extract (only for those tags specified in the **tags** parameter). Defaults to ('href',)
- **canonicalize** (*boolean*) – canonicalize each extracted url (using `w3lib.url.canonicalize_url`). Defaults to **False**. Note that `canonicalize_url` is meant for duplicate checking; it can change the URL visible at server side, so the response can be different for requests with canonicalized and raw URLs. If you're using `LinkExtractor` to follow links it is more robust to keep the default `canonicalize=False`.
- **unique** (*boolean*) – whether duplicate filtering should be applied to extracted links.
- **process_value** (*callable*) – a function which receives each value extracted from the tag and attributes scanned and can modify the value and return a new one, or return `None` to ignore the link altogether. If not given, **process_value** defaults to `lambda x: x`.

For example, to extract links from this code:

```
<a href="javascript:goToPage('../other/page.html'); return false">
↪Link text</a>
```

You can use the following function in `process_value`:

```
def process_value(value):
    m = re.search("javascript:goToPage\('(.*?)'", value)
    if m:
        return m.group(1)
```

- **strip** (*boolean*) – whether to strip whitespaces from extracted attributes. According to HTML5 standard, leading and trailing whitespaces must be stripped from `href` attributes of `<a>`, `<area>` and many other elements, `src` attribute of ``, `<iframe>` elements, etc., so LinkExtractor strips space chars by default. Set `strip=False` to turn it off (e.g. if you're extracting urls from elements or attributes which allow leading/trailing whitespaces).

3.11 设置

The Scrapy settings allows you to customize the behaviour of all Scrapy components, including the core, extensions, pipelines and spiders themselves.

The infrastructure of the settings provides a global namespace of key-value mappings that the code can use to pull configuration values from. The settings can be populated through different mechanisms, which are described below.

The settings are also the mechanism for selecting the currently active Scrapy project (in case you have many).

For a list of available built-in settings see: *Built-in settings reference*.

3.11.1 Designating the settings

When you use Scrapy, you have to tell it which settings you're using. You can do this by using an environment variable, `SCRAPY_SETTINGS_MODULE`.

The value of `SCRAPY_SETTINGS_MODULE` should be in Python path syntax, e.g. `myproject.settings`. Note that the settings module should be on the Python [import search path](#).

3.11.2 Populating the settings

Settings can be populated using different mechanisms, each of which having a different precedence. Here is the list of them in decreasing order of precedence:

1. Command line options (most precedence)
2. Settings per-spider
3. Project settings module
4. Default settings per-command
5. Default global settings (less precedence)

The population of these settings sources is taken care of internally, but a manual handling is possible using API calls. See the *Settings API* topic for reference.

These mechanisms are described in more detail below.

1. Command line options

Arguments provided by the command line are the ones that take most precedence, overriding any other options. You can explicitly override one (or more) settings using the `-s` (or `--set`) command line option.

Example:

```
scrapy crawl myspider -s LOG_FILE=scrapy.log
```

2. Settings per-spider

Spiders (See the [爬虫器](#) chapter for reference) can define their own settings that will take precedence and override the project ones. They can do so by setting their *custom_settings* attribute:

```
class MySpider(scrapy.Spider):
    name = 'myspider'

    custom_settings = {
        'SOME_SETTING': 'some value',
    }
```

3. Project settings module

The project settings module is the standard configuration file for your Scrapy project, it's where most of your custom settings will be populated. For a standard Scrapy project, this means you'll be adding or changing the settings in the `settings.py` file created for your project.

4. Default settings per-command

Each *Scrapy tool* command can have its own default settings, which override the global default settings. Those custom command settings are specified in the `default_settings` attribute of the command class.

5. Default global settings

The global defaults are located in the `scrapy.settings.default_settings` module and documented in the *Built-in settings reference* section.

3.11.3 How to access settings

In a spider, the settings are available through `self.settings`:

```
class MySpider(scrapy.Spider):
    name = 'myspider'
    start_urls = ['http://example.com']

    def parse(self, response):
        print("Existing settings: %s" % self.settings.attributes.keys())
```

注解: The `settings` attribute is set in the base Spider class after the spider is initialized. If you want to use the settings before the initialization (e.g., in your spider's `__init__()` method), you'll need to override the `from_crawler()` method.

Settings can be accessed through the `scrapy.crawler.Crawler.settings` attribute of the Crawler that is passed to `from_crawler` method in extensions, middlewares and item pipelines:

```
class MyExtension(object):
    def __init__(self, log_is_enabled=False):
        if log_is_enabled:
            print("log is enabled!")

    @classmethod
    def from_crawler(cls, crawler):
        settings = crawler.settings
        return cls(settings.getbool('LOG_ENABLED'))
```

The settings object can be used like a dict (e.g., `settings['LOG_ENABLED']`), but it's usually preferred to extract the setting in the format you need it to avoid type errors, using one of the methods provided by the Settings API.

3.11.4 Rationale for setting names

Setting names are usually prefixed with the component that they configure. For example, proper setting names for a fictional robots.txt extension would be `ROBOTSTXT_ENABLED`, `ROBOTSTXT_OBEY`, `ROBOTSTXT_CACHEDIR`, etc.

3.11.5 Built-in settings reference

Here's a list of all available Scrapy settings, in alphabetical order, along with their default values and the scope where they apply.

The scope, where available, shows where the setting is being used, if it's tied to any particular component. In that case the module of that component will be shown, typically an extension, middleware or pipeline. It also means that the component must be enabled in order for the setting to have any effect.

AWS_ACCESS_KEY_ID

Default: `None`

The AWS access key used by code that requires access to [Amazon Web services](#), such as the *S3 feed storage backend*.

AWS_SECRET_ACCESS_KEY

Default: `None`

The AWS secret key used by code that requires access to [Amazon Web services](#), such as the *S3 feed storage backend*.

AWS_ENDPOINT_URL

Default: `None`

Endpoint URL used for S3-like storage, for example Minio or s3.scality. Only supported with `botocore` library.

AWS_USE_SSL

Default: `None`

Use this option if you want to disable SSL connection for communication with S3 or S3-like storage. By default SSL will be used. Only supported with `botocore` library.

AWS_VERIFY

Default: `None`

Verify SSL connection between Scrapy and S3 or S3-like storage. By default SSL verification will occur. Only supported with `botocore` library.

AWS_REGION_NAME

Default: `None`

The name of the region associated with the AWS client. Only supported with `botocore` library.

BOT_NAME

Default: `'scrapybot'`

The name of the bot implemented by this Scrapy project (also known as the project name). This will be used to construct the User-Agent by default, and also for logging.

It's automatically populated with your project name when you create your project with the *startproject* command.

CONCURRENT_ITEMS

Default: `100`

Maximum number of concurrent items (per response) to process in parallel in the Item Processor (also known as the *Item Pipeline*).

CONCURRENT_REQUESTS

Default: `16`

The maximum number of concurrent (ie. simultaneous) requests that will be performed by the Scrapy downloader.

CONCURRENT_REQUESTS_PER_DOMAIN

Default: `8`

The maximum number of concurrent (ie. simultaneous) requests that will be performed to any single domain.

See also: 爬虫器节流 and its *AUTOTHROTTLE_TARGET_CONCURRENCY* option.

CONCURRENT_REQUESTS_PER_IP

Default: 0

The maximum number of concurrent (ie. simultaneous) requests that will be performed to any single IP. If non-zero, the `CONCURRENT_REQUESTS_PER_DOMAIN` setting is ignored, and this one is used instead. In other words, concurrency limits will be applied per IP, not per domain.

This setting also affects `DOWNLOAD_DELAY` and 爬虫器节流: if `CONCURRENT_REQUESTS_PER_IP` is non-zero, download delay is enforced per IP, not per domain.

DEFAULT_ITEM_CLASS

Default: `'scrapy.item.Item'`

The default class that will be used for instantiating items in the *the Scrapy shell*.

DEFAULT_REQUEST_HEADERS

Default:

```
{
    'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
    'Accept-Language': 'en',
}
```

The default headers used for Scrapy HTTP Requests. They're populated in the *DefaultHeadersMiddleware*.

DEPTH_LIMIT

Default: 0

Scope: `scrapy.spidermiddlewares.depth.DepthMiddleware`

The maximum depth that will be allowed to crawl for any site. If zero, no limit will be imposed.

DEPTH_PRIORITY

Default: 0

Scope: `scrapy.spidermiddlewares.depth.DepthMiddleware`

An integer that is used to adjust the priority of a *Request* based on its depth.

The priority of a request is adjusted as follows:

```
request.priority = request.priority - ( depth * DEPTH_PRIORITY )
```

As depth increases, positive values of `DEPTH_PRIORITY` decrease request priority (BFO), while negative values increase request priority (DFO). See also *[Does Scrapy crawl in breadth-first or depth-first order?](#)*.

注解: This setting adjusts priority **in the opposite way** compared to other priority settings *`REDIRECT_PRIORITY_ADJUST`* and *`RETRY_PRIORITY_ADJUST`*.

DEPTH_STATS_VERBOSE

Default: `False`

Scope: `scrapy.spidermiddlewares.depth.DepthMiddleware`

Whether to collect verbose depth stats. If this is enabled, the number of requests for each depth is collected in the stats.

DNSCACHE_ENABLED

Default: `True`

Whether to enable DNS in-memory cache.

DNSCACHE_SIZE

Default: `10000`

DNS in-memory cache size.

DNS_TIMEOUT

Default: `60`

Timeout for processing of DNS queries in seconds. Float is supported.

DOWNLOADER

Default: `'scrapy.core.downloader.Downloader'`

The downloader to use for crawling.

DOWNLOADER_HTTPCLIENTFACTORY

Default: `'scrapy.core.downloader.webclient.ScrapyHTTPClientFactory'`

Defines a `Twisted.protocol.ClientFactory` class to use for HTTP/1.0 connections (for `HTTP10DownloadHandler`).

注解: HTTP/1.0 is rarely used nowadays so you can safely ignore this setting, unless you use `Twisted<11.1`, or if you really want to use HTTP/1.0 and override `DOWNLOAD_HANDLERS_BASE` for `http(s)` scheme accordingly, i.e. to `'scrapy.core.downloader.handlers.http.HTTP10DownloadHandler'`.

DOWNLOADER_CLIENTCONTEXTFACTORY

Default: `'scrapy.core.downloader.contextfactory.ScrapyClientContextFactory'`

Represents the classpath to the `ContextFactory` to use.

Here, “ContextFactory” is a Twisted term for SSL/TLS contexts, defining the TLS/SSL protocol version to use, whether to do certificate verification, or even enable client-side authentication (and various other things).

注解: Scrapy default context factory **does NOT perform remote server certificate verification**. This is usually fine for web scraping.

If you do need remote server certificate verification enabled, Scrapy also has another context factory class that you can set, `'scrapy.core.downloader.contextfactory.BrowserLikeContextFactory'`, which uses the platform's certificates to validate remote endpoints. **This is only available if you use `Twisted>=14.0`.**

If you do use a custom `ContextFactory`, make sure it accepts a `method` parameter at init (this is the `OpenSSL.SSL` method mapping `DOWNLOADER_CLIENT_TLS_METHOD`).

DOWNLOADER_CLIENT_TLS_METHOD

Default: `'TLS'`

Use this setting to customize the TLS/SSL method used by the default HTTP/1.1 downloader.

This setting must be one of these string values:

- `'TLS'`: maps to `OpenSSL's TLS_method()` (a.k.a `SSLv23_method()`), which allows protocol negotiation, starting from the highest supported by the platform; **default, recommended**
- `'TLSv1.0'`: this value forces HTTPS connections to use TLS version 1.0 ; set this if you want the behavior of `Scrapy<1.1`

- 'TLSv1.1': forces TLS version 1.1
- 'TLSv1.2': forces TLS version 1.2
- 'SSLv3': forces SSL version 3 (**not recommended**)

注解: We recommend that you use PyOpenSSL \geq 0.13 and Twisted \geq 0.13 or above (Twisted \geq 14.0 if you can).

DOWNLOADER_MIDDLEWARES

Default:: {}

A dict containing the downloader middlewares enabled in your project, and their orders. For more info see *Activating a downloader middleware*.

DOWNLOADER_MIDDLEWARES_BASE

Default:

```
{
    'scrapy.downloadermiddlewares.robotstxt.RobotsTxtMiddleware': 100,
    'scrapy.downloadermiddlewares.httpauth.HttpAuthMiddleware': 300,
    'scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware': 350,
    'scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware': 400,
    'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware': 500,
    'scrapy.downloadermiddlewares.retry.RetryMiddleware': 550,
    'scrapy.downloadermiddlewares.ajaxcrawl.AjaxCrawlMiddleware': 560,
    'scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware': 580,
    'scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware': 590,
    'scrapy.downloadermiddlewares.redirect.RedirectMiddleware': 600,
    'scrapy.downloadermiddlewares.cookies.CookiesMiddleware': 700,
    'scrapy.downloadermiddlewares.httpproxy.HttpProxyMiddleware': 750,
    'scrapy.downloadermiddlewares.stats.DownloaderStats': 850,
    'scrapy.downloadermiddlewares.httpcache.HttpCacheMiddleware': 900,
}
```

A dict containing the downloader middlewares enabled by default in Scrapy. Low orders are closer to the engine, high orders are closer to the downloader. You should never modify this setting in your project, modify `DOWNLOADER_MIDDLEWARES` instead. For more info see *Activating a downloader middleware*.

DOWNLOADER_STATS

Default: `True`

Whether to enable downloader stats collection.

DOWNLOAD_DELAY

Default: `0`

The amount of time (in secs) that the downloader should wait before downloading consecutive pages from the same website. This can be used to throttle the crawling speed to avoid hitting servers too hard. Decimal numbers are supported. Example:

```
DOWNLOAD_DELAY = 0.25    # 250 ms of delay
```

This setting is also affected by the [RANDOMIZE_DOWNLOAD_DELAY](#) setting (which is enabled by default). By default, Scrapy doesn't wait a fixed amount of time between requests, but uses a random interval between $0.5 * \text{DOWNLOAD_DELAY}$ and $1.5 * \text{DOWNLOAD_DELAY}$.

When [CONCURRENT_REQUESTS_PER_IP](#) is non-zero, delays are enforced per ip address instead of per domain.

You can also change this setting per spider by setting `download_delay` spider attribute.

DOWNLOAD_HANDLERS

Default: `{}`

A dict containing the request downloader handlers enabled in your project. See [DOWNLOAD_HANDLERS_BASE](#) for example format.

DOWNLOAD_HANDLERS_BASE

Default:

```
{
    'file': 'scrapy.core.downloader.handlers.file.FileDownloadHandler',
    'http': 'scrapy.core.downloader.handlers.http.HTTPDownloadHandler',
    'https': 'scrapy.core.downloader.handlers.http.HTTPDownloadHandler',
    's3': 'scrapy.core.downloader.handlers.s3.S3DownloadHandler',
    'ftp': 'scrapy.core.downloader.handlers.ftp.FTPDownloadHandler',
}
```

A dict containing the request download handlers enabled by default in Scrapy. You should never modify this setting in your project, modify [DOWNLOAD_HANDLERS](#) instead.

You can disable any of these download handlers by assigning `None` to their URI scheme in `DOWNLOAD_HANDLERS`. E.g., to disable the built-in FTP handler (without replacement), place this in your `settings.py`:

```
DOWNLOAD_HANDLERS = {
    'ftp': None,
}
```

DOWNLOAD_TIMEOUT

Default: 180

The amount of time (in secs) that the downloader will wait before timing out.

注解: This timeout can be set per spider using `download_timeout` spider attribute and per-request using `download_timeout` Request.meta key.

DOWNLOAD_MAXSIZE

Default: 1073741824 (1024MB)

The maximum response size (in bytes) that downloader will download.

If you want to disable it set to 0.

注解: This size can be set per spider using `download_maxsize` spider attribute and per-request using `download_maxsize` Request.meta key.

This feature needs Twisted >= 11.1.

DOWNLOAD_WARN_SIZE

Default: 33554432 (32MB)

The response size (in bytes) that downloader will start to warn.

If you want to disable it set to 0.

注解: This size can be set per spider using `download_warnsize` spider attribute and per-request using `download_warnsize` Request.meta key.

This feature needs Twisted \geq 11.1.

DOWNLOAD_FAIL_ON_DATALOSS

Default: `True`

Whether or not to fail on broken responses, that is, declared `Content-Length` does not match content sent by the server or chunked response was not properly finish. If `True`, these responses raise a `ResponseFailed([_DataLoss])` error. If `False`, these responses are passed through and the flag `dataloss` is added to the response, i.e.: `'dataloss'` in `response.flags` is `True`.

Optionally, this can be set per-request basis by using the `download_fail_on_dataloss` `Request.meta` key to `False`.

注解: A broken response, or data loss error, may happen under several circumstances, from server misconfiguration to network errors to data corruption. It is up to the user to decide if it makes sense to process broken responses considering they may contain partial or incomplete content. If `RETRY_ENABLED` is `True` and this setting is set to `True`, the `ResponseFailed([_DataLoss])` failure will be retried as usual.

DUPEFILTER_CLASS

Default: `'scrapy.dupefilters.RFPDupeFilter'`

The class used to detect and filter duplicate requests.

The default (`RFPDupeFilter`) filters based on request fingerprint using the `scrapy.utils.request.request_fingerprint` function. In order to change the way duplicates are checked you could subclass `RFPDupeFilter` and override its `request_fingerprint` method. This method should accept scrapy `Request` object and return its fingerprint (a string).

You can disable filtering of duplicate requests by setting `DUPEFILTER_CLASS` to `'scrapy.dupefilters.BaseDupeFilter'`. Be very careful about this however, because you can get into crawling loops. It's usually a better idea to set the `dont_filter` parameter to `True` on the specific `Request` that should not be filtered.

DUPEFILTER_DEBUG

Default: `False`

By default, `RFPDupeFilter` only logs the first duplicate request. Setting `DUPEFILTER_DEBUG` to `True` will make it log all duplicate requests.

EDITOR

Default: `vi` (on Unix systems) or the IDLE editor (on Windows)

The editor to use for editing spiders with the `edit` command. Additionally, if the `EDITOR` environment variable is set, the `edit` command will prefer it over the default setting.

EXTENSIONS

Default:: `{}`

A dict containing the extensions enabled in your project, and their orders.

EXTENSIONS_BASE

Default:

```
{
    'scrapy.extensions.corestats.CoreStats': 0,
    'scrapy.extensions.telnet.TelnetConsole': 0,
    'scrapy.extensions.memusage.MemoryUsage': 0,
    'scrapy.extensions.memdebug.MemoryDebugger': 0,
    'scrapy.extensions.closespider.CloseSpider': 0,
    'scrapy.extensions.feedexport.FeedExporter': 0,
    'scrapy.extensions.logstats.LogStats': 0,
    'scrapy.extensions.spiderstate.SpiderState': 0,
    'scrapy.extensions.throttle.AutoThrottle': 0,
}
```

A dict containing the extensions available by default in Scrapy, and their orders. This setting contains all stable built-in extensions. Keep in mind that some of them need to be enabled through a setting.

For more information See the *extensions user guide* and the *list of available extensions*.

FEED_TEMPDIR

The Feed Temp dir allows you to set a custom folder to save crawler temporary files before uploading with *FTP feed storage* and *Amazon S3*.

FTP_PASSIVE_MODE

Default: `True`

Whether or not to use passive mode when initiating FTP transfers.

FTP_PASSWORD

Default: "guest"

The password to use for FTP connections when there is no "ftp_password" in Request meta.

注解: Paraphrasing [RFC 1635](#), although it is common to use either the password “guest” or one’s e-mail address for anonymous FTP, some FTP servers explicitly ask for the user’s e-mail address and will not allow login with the “guest” password.

FTP_USER

Default: "anonymous"

The username to use for FTP connections when there is no "ftp_user" in Request meta.

ITEM_PIPELINES

Default: {}

A dict containing the item pipelines to use, and their orders. Order values are arbitrary, but it is customary to define them in the 0-1000 range. Lower orders process before higher orders.

Example:

```
ITEM_PIPELINES = {
    'mybot.pipelines.validate.ValidateMyItem': 300,
    'mybot.pipelines.validate.StoreMyItem': 800,
}
```

ITEM_PIPELINES_BASE

Default: {}

A dict containing the pipelines enabled by default in Scrapy. You should never modify this setting in your project, modify *ITEM_PIPELINES* instead.

LOG_ENABLED

Default: True

Whether to enable logging.

LOG_ENCODING

Default: `'utf-8'`

The encoding to use for logging.

LOG_FILE

Default: `None`

File name to use for logging output. If `None`, standard error will be used.

LOG_FORMAT

Default: `'%(asctime)s [%(name)s] %(levelname)s: %(message)s'`

String for formatting log messages. Refer to the [Python logging documentation](#) for the whole list of available placeholders.

LOG_DATEFORMAT

Default: `'%Y-%m-%d %H:%M:%S'`

String for formatting date/time, expansion of the `%(asctime)s` placeholder in `LOG_FORMAT`. Refer to the [Python datetime documentation](#) for the whole list of available directives.

LOG_LEVEL

Default: `'DEBUG'`

Minimum level to log. Available levels are: `CRITICAL`, `ERROR`, `WARNING`, `INFO`, `DEBUG`. For more info see [日志](#).

LOG_STDOUT

Default: `False`

If `True`, all standard output (and error) of your process will be redirected to the log. For example if you `print('hello')` it will appear in the Scrapy log.

LOG_SHORT_NAMES

Default: `False`

If `True`, the logs will just contain the root path. If it is set to `False` then it displays the component responsible for the log output

LOGSTATS_INTERVAL

Default: 60.0

The interval (in seconds) between each logging printout of the stats by `LogStats`.

MEMDEBUG_ENABLED

Default: `False`

Whether to enable memory debugging.

MEMDEBUG_NOTIFY

Default: `[]`

When memory debugging is enabled a memory report will be sent to the specified addresses if this setting is not empty, otherwise the report will be written to the log.

Example:

```
MEMDEBUG_NOTIFY = ['user@example.com']
```

MEMUSAGE_ENABLED

Default: `True`

Scope: `scrapy.extensions.memusage`

Whether to enable the memory usage extension. This extension keeps track of a peak memory used by the process (it writes it to stats). It can also optionally shutdown the Scrapy process when it exceeds a memory limit (see [MEMUSAGE_LIMIT_MB](#)), and notify by email when that happened (see [MEMUSAGE_NOTIFY_MAIL](#)).

See *Memory usage extension*.

MEMUSAGE_LIMIT_MB

Default: 0

Scope: `scrapy.extensions.memusage`

The maximum amount of memory to allow (in megabytes) before shutting down Scrapy (if `MEMUSAGE_ENABLED` is `True`). If zero, no check will be performed.

See *Memory usage extension*.

MEMUSAGE_CHECK_INTERVAL_SECONDS

1.1 新版功能.

Default: 60.0

Scope: `scrapy.extensions.memusage`

The *Memory usage extension* checks the current memory usage, versus the limits set by `MEMUSAGE_LIMIT_MB` and `MEMUSAGE_WARNING_MB`, at fixed time intervals.

This sets the length of these intervals, in seconds.

See *Memory usage extension*.

MEMUSAGE_NOTIFY_MAIL

Default: `False`

Scope: `scrapy.extensions.memusage`

A list of emails to notify if the memory limit has been reached.

Example:

```
MEMUSAGE_NOTIFY_MAIL = ['user@example.com']
```

See *Memory usage extension*.

MEMUSAGE_WARNING_MB

Default: 0

Scope: `scrapy.extensions.memusage`

The maximum amount of memory to allow (in megabytes) before sending a warning email notifying about it. If zero, no warning will be produced.

NEWSPIDER_MODULE

Default: ''

Module where to create new spiders using the *genspider* command.

Example:

```
NEWSPIDER_MODULE = 'mybot.spiders_dev'
```

RANDOMIZE_DOWNLOAD_DELAY

Default: True

If enabled, Scrapy will wait a random amount of time (between $0.5 * \text{DOWNLOAD_DELAY}$ and $1.5 * \text{DOWNLOAD_DELAY}$) while fetching requests from the same website.

This randomization decreases the chance of the crawler being detected (and subsequently blocked) by sites which analyze requests looking for statistically significant similarities in the time between their requests.

The randomization policy is the same used by `wget --random-wait` option.

If `DOWNLOAD_DELAY` is zero (default) this option has no effect.

REACTOR_THREADPOOL_MAXSIZE

Default: 10

The maximum limit for Twisted Reactor thread pool size. This is common multi-purpose thread pool used by various Scrapy components. Threaded DNS Resolver, BlockingFeedStorage, S3FilesStore just to name a few. Increase this value if you're experiencing problems with insufficient blocking IO.

REDIRECT_MAX_TIMES

Default: 20

Defines the maximum times a request can be redirected. After this maximum the request's response is returned as is. We used Firefox default value for the same task.

REDIRECT_PRIORITY_ADJUST

Default: +2

Scope: `scrapy.downloadermiddlewares.redirect.RedirectMiddleware`

Adjust redirect request priority relative to original request:

- a positive priority adjust (default) means higher priority.
- a negative priority adjust means lower priority.

RETRY_PRIORITY_ADJUST

Default: -1

Scope: `scrapy.downloadermiddlewares.retry.RetryMiddleware`

Adjust retry request priority relative to original request:

- a positive priority adjust means higher priority.
- a negative priority adjust (default) means lower priority.

ROBOTSTXT_OBEY

Default: `False`

Scope: `scrapy.downloadermiddlewares.robotstxt`

If enabled, Scrapy will respect robots.txt policies. For more information see [*RobotsTxtMiddleware*](#).

注解: While the default value is `False` for historical reasons, this option is enabled by default in settings.py file generated by `scrapy startproject` command.

SCHEDULER

Default: `'scrapy.core.scheduler.Scheduler'`

The scheduler to use for crawling.

SCHEDULER_DEBUG

Default: `False`

Setting to `True` will log debug information about the requests scheduler. This currently logs (only once) if the requests cannot be serialized to disk. Stats counter (`scheduler/unserializable`) tracks the number of times this happens.

Example entry in logs:

```
1956-01-31 00:00:00+0800 [scrapy.core.scheduler] ERROR: Unable to serialize request:
<GET http://example.com> - reason: cannot serialize <Request at 0x9a7c7ec>
(type Request)> - no more unserializable requests will be logged
(see 'scheduler/unserializable' stats counter)
```

SCHEDULER_DISK_QUEUE

Default: `'scrapy.squeues.PickleLifoDiskQueue'`

Type of disk queue that will be used by scheduler. Other available types are `scrapy.squeues.PickleFifoDiskQueue`, `scrapy.squeues.MarshalFifoDiskQueue`, `scrapy.squeues.MarshalLifoDiskQueue`.

SCHEDULER_MEMORY_QUEUE

Default: `'scrapy.squeues.LifoMemoryQueue'`

Type of in-memory queue used by scheduler. Other available type is: `scrapy.squeues.FifoMemoryQueue`.

SCHEDULER_PRIORITY_QUEUE

Default: `'scrapy.pqueues.ScrapyPriorityQueue'`

Type of priority queue used by the scheduler. Another available type is `scrapy.pqueues.DownloaderAwarePriorityQueue`. `scrapy.pqueues.DownloaderAwarePriorityQueue` works better than `scrapy.pqueues.ScrapyPriorityQueue` when you crawl many different domains in parallel. But currently `scrapy.pqueues.DownloaderAwarePriorityQueue` does not work together with *CONCURRENT_REQUESTS_PER_IP*.

SPIDER_CONTRACTS

Default:: `{}`

A dict containing the spider contracts enabled in your project, used for testing spiders. For more info see [爬虫器约束](#).

SPIDER_CONTRACTS_BASE

Default:

```
{
    'scrapy.contracts.default.UrlContract' : 1,
    'scrapy.contracts.default.ReturnsContract': 2,
    'scrapy.contracts.default.ScrapesContract': 3,
}
```

A dict containing the scrapy contracts enabled by default in Scrapy. You should never modify this setting in your project, modify *SPIDER_CONTRACTS* instead. For more info see [爬虫器约束](#).

You can disable any of these contracts by assigning `None` to their class path in *SPIDER_CONTRACTS*. E.g., to disable the built-in `ScrapesContract`, place this in your `settings.py`:

```
SPIDER_CONTRACTS = {
    'scrapy.contracts.default.ScrapesContract': None,
}
```

SPIDER_LOADER_CLASS

Default: `'scrapy.spiderloader.SpiderLoader'`

The class that will be used for loading spiders, which must implement the *SpiderLoader API*.

SPIDER_LOADER_WARN_ONLY

1.3.3 新版功能.

Default: `False`

By default, when scrapy tries to import spider classes from *SPIDER_MODULES*, it will fail loudly if there is any `ImportError` exception. But you can choose to silence this exception and turn it into a simple warning by setting `SPIDER_LOADER_WARN_ONLY = True`.

注解: Some *scrapy commands* run with this setting to `True` already (i.e. they will only issue a warning and will not fail) since they do not actually need to load spider classes to work: *scrapy runspider*, *scrapy settings*, *scrapy startproject*, *scrapy version*.

SPIDER_MIDDLEWARES

Default:: `{}`

A dict containing the spider middlewares enabled in your project, and their orders. For more info see *Activating a spider middleware*.

SPIDER_MIDDLEWARES_BASE

Default:

```
{
    'scrapy.spidermiddlewares.httperror.HttpErrorMiddleware': 50,
    'scrapy.spidermiddlewares.offsite.OffsiteMiddleware': 500,
    'scrapy.spidermiddlewares.referer.RefererMiddleware': 700,
    'scrapy.spidermiddlewares.urllength.UrlLengthMiddleware': 800,
    'scrapy.spidermiddlewares.depth.DepthMiddleware': 900,
}
```

A dict containing the spider middlewares enabled by default in Scrapy, and their orders. Low orders are closer to the engine, high orders are closer to the spider. For more info see *Activating a spider middleware*.

SPIDER_MODULES

Default: []

A list of modules where Scrapy will look for spiders.

Example:

```
SPIDER_MODULES = ['mybot.spiders_prod', 'mybot.spiders_dev']
```

STATS_CLASS

Default: 'scrapy.statscollectors.MemoryStatsCollector'

The class to use for collecting stats, who must implement the *Stats Collector API*.

STATS_DUMP

Default: True

Dump the *Scrapy stats* (to the Scrapy log) once the spider finishes.

For more info see: [统计数据集合](#).

STATSMAILER_RCPTS

Default: [] (empty list)

Send Scrapy stats after spiders finish scraping. See `StatsMailer` for more info.

TELNETCONSOLE_ENABLED

Default: True

A boolean which specifies if the *telnet console* will be enabled (provided its extension is also enabled).

TELNETCONSOLE_PORT

Default: [6023, 6073]

The port range to use for the telnet console. If set to `None` or `0`, a dynamically assigned port is used. For more info see [远程控制台](#).

TEMPLATES_DIR

Default: `templates` dir inside scrapy module

The directory where to look for templates when creating new projects with *startproject* command and new spiders with *genspider* command.

The project name must not conflict with the name of custom files or directories in the `project` subdirectory.

URLLENGTH_LIMIT

Default: 2083

Scope: `spidermiddlewares.urllength`

The maximum URL length to allow for crawled URLs. For more information about the default value for this setting see: <https://boutell.com/newfaq/misc/urllength.html>

USER_AGENT

Default: "Scrapy/VERSION (+https://scrapy.org)"

The default User-Agent to use when crawling, unless overridden.

Settings documented elsewhere:

The following settings are documented elsewhere, please check each specific case to see how to enable and use them.

- *AJAXCRAWL_ENABLED*
- *AUTOTHROTTLING_DEBUG*
- *AUTOTHROTTLING_ENABLED*
- *AUTOTHROTTLING_MAX_DELAY*
- *AUTOTHROTTLING_START_DELAY*
- *AUTOTHROTTLING_TARGET_CONCURRENCY*
- *AWS_ACCESS_KEY_ID*
- *AWS_ENDPOINT_URL*
- *AWS_REGION_NAME*
- *AWS_SECRET_ACCESS_KEY*
- *AWS_USE_SSL*

- *AWS_VERIFY*
- *BOT_NAME*
- *CLOSESPIDER_ERRORCOUNT*
- *CLOSESPIDER_ITEMCOUNT*
- *CLOSESPIDER_PAGECOUNT*
- *CLOSESPIDER_TIMEOUT*
- *COMMANDS_MODULE*
- *COMPRESSION_ENABLED*
- *CONCURRENT_ITEMS*
- *CONCURRENT_REQUESTS*
- *CONCURRENT_REQUESTS_PER_DOMAIN*
- *CONCURRENT_REQUESTS_PER_IP*
- *COOKIES_DEBUG*
- *COOKIES_ENABLED*
- *DEFAULT_ITEM_CLASS*
- *DEFAULT_REQUEST_HEADERS*
- *DEPTH_LIMIT*
- *DEPTH_PRIORITY*
- *DEPTH_STATS_VERBOSE*
- *DNSCACHE_ENABLED*
- *DNSCACHE_SIZE*
- *DNS_TIMEOUT*
- *DOWNLOADER*
- *DOWNLOADER_CLIENTCONTEXTFACTORY*
- *DOWNLOADER_CLIENT_TLS_METHOD*
- *DOWNLOADER_HTTPCLIENTFACTORY*
- *DOWNLOADER_MIDDLEWARES*
- *DOWNLOADER_MIDDLEWARES_BASE*
- *DOWNLOADER_STATS*
- *DOWNLOAD_DELAY*

- *DOWNLOAD_FAIL_ON_DATALOSS*
- *DOWNLOAD_HANDLERS*
- *DOWNLOAD_HANDLERS_BASE*
- *DOWNLOAD_MAXSIZE*
- *DOWNLOAD_TIMEOUT*
- *DOWNLOAD_WARN_SIZE*
- *DUPEFILTER_CLASS*
- *DUPEFILTER_DEBUG*
- *EDITOR*
- *EXTENSIONS*
- *EXTENSIONS_BASE*
- *FEED_EXPORTERS*
- *FEED_EXPORTERS_BASE*
- *FEED_EXPORT_ENCODING*
- *FEED_EXPORT_FIELDS*
- *FEED_EXPORT_INDENT*
- *FEED_FORMAT*
- *FEED_STORAGES*
- *FEED_STORAGES_BASE*
- *FEED_STORAGE_S3_ACL*
- *FEED_STORE_EMPTY*
- *FEED_TEMPDIR*
- *FEED_URI*
- *FILES_EXPIRES*
- *FILES_RESULT_FIELD*
- *FILES_STORE*
- *FILES_STORE_GCS_ACL*
- *FILES_STORE_S3_ACL*
- *FILES_URLS_FIELD*
- *FTP_PASSIVE_MODE*

- *FTP_PASSWORD*
- *FTP_USER*
- *GCS_PROJECT_ID*
- *HTTPCACHE_ALWAYS_STORE*
- *HTTPCACHE_DBM_MODULE*
- *HTTPCACHE_DIR*
- *HTTPCACHE_ENABLED*
- *HTTPCACHE_EXPIRATION_SECS*
- *HTTPCACHE_GZIP*
- *HTTPCACHE_IGNORE_HTTP_CODES*
- *HTTPCACHE_IGNORE_MISSING*
- *HTTPCACHE_IGNORE_RESPONSE_CACHE_CONTROLS*
- *HTTPCACHE_IGNORE_SCHEMES*
- *HTTPCACHE_POLICY*
- *HTTPCACHE_STORAGE*
- *HTTPERROR_ALLOWED_CODES*
- *HTTPERROR_ALLOW_ALL*
- *HTTPPROXY_AUTH_ENCODING*
- *HTTPPROXY_ENABLED*
- *IMAGES_EXPIRES*
- *IMAGES_MIN_HEIGHT*
- *IMAGES_MIN_WIDTH*
- *IMAGES_RESULT_FIELD*
- *IMAGES_STORE*
- *IMAGES_STORE_GCS_ACL*
- *IMAGES_STORE_S3_ACL*
- *IMAGES_THUMBS*
- *IMAGES_URLS_FIELD*
- *ITEM_PIPELINES*
- *ITEM_PIPELINES_BASE*

- *LOGSTATS_INTERVAL*
- *LOG_DATEFORMAT*
- *LOG_ENABLED*
- *LOG_ENCODING*
- *LOG_FILE*
- *LOG_FORMAT*
- *LOG_LEVEL*
- *LOG_SHORT_NAMES*
- *LOG_STDOUT*
- *MAIL_FROM*
- *MAIL_HOST*
- *MAIL_PASS*
- *MAIL_PORT*
- *MAIL_SSL*
- *MAIL_TLS*
- *MAIL_USER*
- *MEDIA_ALLOW_REDIRECTS*
- *MEMDEBUG_ENABLED*
- *MEMDEBUG_NOTIFY*
- *MEMUSAGE_CHECK_INTERVAL_SECONDS*
- *MEMUSAGE_ENABLED*
- *MEMUSAGE_LIMIT_MB*
- *MEMUSAGE_NOTIFY_MAIL*
- *MEMUSAGE_WARNING_MB*
- *METAREFRESH_ENABLED*
- *METAREFRESH_MAXDELAY*
- *NEWSPIDER_MODULE*
- *RANDOMIZE_DOWNLOAD_DELAY*
- *REACTOR_THREADPOOL_MAXSIZE*
- *REDIRECT_ENABLED*

- *REDIRECT_MAX_TIMES*
- *REDIRECT_MAX_TIMES*
- *REDIRECT_PRIORITY_ADJUST*
- *REFERER_ENABLED*
- *REFERRER_POLICY*
- *RETRY_ENABLED*
- *RETRY_HTTP_CODES*
- *RETRY_PRIORITY_ADJUST*
- *RETRY_TIMES*
- *ROBOTSTXT_OBEY*
- *SCHEDULER*
- *SCHEDULER_DEBUG*
- *SCHEDULER_DISK_QUEUE*
- *SCHEDULER_MEMORY_QUEUE*
- *SCHEDULER_PRIORITY_QUEUE*
- *SPIDER_CONTRACTS*
- *SPIDER_CONTRACTS_BASE*
- *SPIDER_LOADER_CLASS*
- *SPIDER_LOADER_WARN_ONLY*
- *SPIDER_MIDDLEWARES*
- *SPIDER_MIDDLEWARES_BASE*
- *SPIDER_MODULES*
- *STATSMAILER_RCPTS*
- *STATS_CLASS*
- *STATS_DUMP*
- *TELNETCONSOLE_ENABLED*
- *TELNETCONSOLE_HOST*
- *TELNETCONSOLE_PASSWORD*
- *TELNETCONSOLE_PORT*
- *TELNETCONSOLE_PORT*

- `TELNETCONSOLE_USERNAME`
- `TEMPLATES_DIR`
- `URLLENGTH_LIMIT`
- `USER_AGENT`

3.12 异常

3.12.1 Built-in Exceptions reference

Here's a list of all exceptions included in Scrapy and their usage.

DropItem

`exception scrapy.exceptions.DropItem`

The exception that must be raised by item pipeline stages to stop processing an Item. For more information see *Item 管道*.

CloseSpider

`exception scrapy.exceptions.CloseSpider(reason='cancelled')`

This exception can be raised from a spider callback to request the spider to be closed/stopped. Supported arguments:

参数 `reason` (*str*) – the reason for closing

For example:

```
def parse_page(self, response):
    if 'Bandwidth exceeded' in response.body:
        raise CloseSpider('bandwidth_exceeded')
```

DontCloseSpider

`exception scrapy.exceptions.DontCloseSpider`

This exception can be raised in a *spider_idle* signal handler to prevent the spider from being closed.

IgnoreRequest

`exception scrapy.exceptions.IgnoreRequest`

This exception can be raised by the Scheduler or any downloader middleware to indicate that the request should be ignored.

NotConfigured

exception scrapy.exceptions.NotConfigured

This exception can be raised by some components to indicate that they will remain disabled. Those components include:

- Extensions
- Item pipelines
- Downloader middlewares
- Spider middlewares

The exception must be raised in the component's `__init__` method.

NotSupported

exception scrapy.exceptions.NotSupported

This exception is raised to indicate an unsupported feature.

命令行工具 了解如何通过命令行管理 Scrapy 项目。

爬虫器 定义网站爬虫规则。

选择器 使用 Xpath 从网页中提取数据。

Scrapy shell 在交互式环境中测试解析程序。

Items 定义你想要获取的数据。

Item 加载器 将提取的数据填充到项目中。

Item 管道 处理和保存抓取到的数据。

Feed 导出 将你抓取到的数据以不同的方式输出储存。

请求和响应 使用不同的类来实现 HTTP 的请求和响应。

链接提取 便捷的类，用于提取页面中的超链接并继续跟进。

设置 了解如何配置 Scrapy 和查看所有的可用配置。

异常 查看所有可用的异常及其含义。

4.1 日志

注解: `scrapy.log` has been deprecated alongside its functions in favor of explicit calls to the Python standard logging. Keep reading to learn more about the new logging system.

Scrapy uses Python's builtin logging system for event logging. We'll provide some simple examples to get you started, but for more advanced use-cases it's strongly suggested to read thoroughly its documentation.

Logging works out of the box, and can be configured to some extent with the Scrapy settings listed in *Logging settings*.

Scrapy calls `scrapy.utils.log.configure_logging()` to set some reasonable defaults and handle those settings in *Logging settings* when running commands, so it's recommended to manually call it if you're running Scrapy from scripts as described in *Run Scrapy from a script*.

4.1.1 Log levels

Python's builtin logging defines 5 different levels to indicate the severity of a given log message. Here are the standard ones, listed in decreasing order:

1. `logging.CRITICAL` - for critical errors (highest severity)
2. `logging.ERROR` - for regular errors

3. `logging.WARNING` - for warning messages
4. `logging.INFO` - for informational messages
5. `logging.DEBUG` - for debugging messages (lowest severity)

4.1.2 How to log messages

Here' s a quick example of how to log a message using the `logging.WARNING` level:

```
import logging
logging.warning("This is a warning")
```

There are shortcuts for issuing log messages on any of the standard 5 levels, and there' s also a general `logging.log` method which takes a given level as argument. If needed, the last example could be rewritten as:

```
import logging
logging.log(logging.WARNING, "This is a warning")
```

On top of that, you can create different “loggers” to encapsulate messages. (For example, a common practice is to create different loggers for every module). These loggers can be configured independently, and they allow hierarchical constructions.

The previous examples use the root logger behind the scenes, which is a top level logger where all messages are propagated to (unless otherwise specified). Using `logging` helpers is merely a shortcut for getting the root logger explicitly, so this is also an equivalent of the last snippets:

```
import logging
logger = logging.getLogger()
logger.warning("This is a warning")
```

You can use a different logger just by getting its name with the `logging.getLogger` function:

```
import logging
logger = logging.getLogger('mycustomlogger')
logger.warning("This is a warning")
```

Finally, you can ensure having a custom logger for any module you' re working on by using the `__name__` variable, which is populated with current module' s path:

```
import logging
logger = logging.getLogger(__name__)
logger.warning("This is a warning")
```

参见:

Module logging, [HowTo Basic Logging Tutorial](#)

Module logging, [Loggers](#) Further documentation on loggers

4.1.3 Logging from Spiders

Scrapy provides a *logger* within each Spider instance, which can be accessed and used like this:

```
import scrapy

class MySpider(scrapy.Spider):

    name = 'myspider'
    start_urls = ['https://scrapinghub.com']

    def parse(self, response):
        self.logger.info('Parse function called on %s', response.url)
```

That logger is created using the Spider's name, but you can use any custom Python logger you want. For example:

```
import logging
import scrapy

logger = logging.getLogger('mycustomlogger')

class MySpider(scrapy.Spider):

    name = 'myspider'
    start_urls = ['https://scrapinghub.com']

    def parse(self, response):
        logger.info('Parse function called on %s', response.url)
```

4.1.4 Logging configuration

Loggers on their own don't manage how messages sent through them are displayed. For this task, different “handlers” can be attached to any logger instance and they will redirect those messages to appropriate destinations, such as the standard output, files, emails, etc.

By default, Scrapy sets and configures a handler for the root logger, based on the settings below.

Logging settings

These settings can be used to configure the logging:

- `LOG_FILE`
- `LOG_ENABLED`
- `LOG_ENCODING`
- `LOG_LEVEL`
- `LOG_FORMAT`
- `LOG_DATEFORMAT`
- `LOG_STDOUT`
- `LOG_SHORT_NAMES`

The first couple of settings define a destination for log messages. If `LOG_FILE` is set, messages sent through the root logger will be redirected to a file named `LOG_FILE` with encoding `LOG_ENCODING`. If unset and `LOG_ENABLED` is `True`, log messages will be displayed on the standard error. Lastly, if `LOG_ENABLED` is `False`, there won't be any visible log output.

`LOG_LEVEL` determines the minimum level of severity to display, those messages with lower severity will be filtered out. It ranges through the possible levels listed in [Log levels](#).

`LOG_FORMAT` and `LOG_DATEFORMAT` specify formatting strings used as layouts for all messages. Those strings can contain any placeholders listed in [logging's logrecord attributes docs](#) and [datetime's strftime and strptime directives](#) respectively.

If `LOG_SHORT_NAMES` is set, then the logs will not display the scrapy component that prints the log. It is unset by default, hence logs contain the scrapy component responsible for that log output.

Command-line options

There are command-line arguments, available for all commands, that you can use to override some of the Scrapy settings regarding logging.

- `--logfile FILE` Overrides `LOG_FILE`
- `--loglevel/-L LEVEL` Overrides `LOG_LEVEL`
- `--nolog` Sets `LOG_ENABLED` to `False`

参见:

Module [logging.handlers](#) Further documentation on available handlers

Advanced customization

Because Scrapy uses `stdlib` logging module, you can customize logging using all features of `stdlib` logging.

For example, let's say you're scraping a website which returns many HTTP 404 and 500 responses, and you want to hide all messages like this:

```
2016-12-16 22:00:06 [scrapy.spidermiddlewares.httperror] INFO: Ignoring
response <500 http://quotes.toscrape.com/page/1-34/>: HTTP status code
is not handled or not allowed
```

The first thing to note is a logger name - it is in brackets: `[scrapy.spidermiddlewares.httperror]`. If you get just `[scrapy]` then `LOG_SHORT_NAMES` is likely set to `True`; set it to `False` and re-run the crawl.

Next, we can see that the message has `INFO` level. To hide it we should set logging level for `scrapy.spidermiddlewares.httperror` higher than `INFO`; next level after `INFO` is `WARNING`. It could be done e.g. in the spider's `__init__` method:

```
import logging
import scrapy

class MySpider(scrapy.Spider):
    # ...
    def __init__(self, *args, **kwargs):
        logger = logging.getLogger('scrapy.spidermiddlewares.httperror')
        logger.setLevel(logging.WARNING)
        super().__init__(*args, **kwargs)
```

If you run this spider again then `INFO` messages from `scrapy.spidermiddlewares.httperror` logger will be gone.

4.1.5 scrapy.utils.log module

4.2 统计数据集合

Scrapy provides a convenient facility for collecting stats in the form of key/values, where values are often counters. The facility is called the Stats Collector, and can be accessed through the `stats` attribute of the *Crawler API*, as illustrated by the examples in the *Common Stats Collector uses* section below.

However, the Stats Collector is always available, so you can always import it in your module and use its API (to increment or set new stat keys), regardless of whether the stats collection is enabled or not. If it's disabled, the API will still work but it won't collect anything. This is aimed at simplifying the stats

collector usage: you should spend no more than one line of code for collecting stats in your spider, Scrapy extension, or whatever code you're using the Stats Collector from.

Another feature of the Stats Collector is that it's very efficient (when enabled) and extremely efficient (almost unnoticeable) when disabled.

The Stats Collector keeps a stats table per open spider which is automatically opened when the spider is opened, and closed when the spider is closed.

4.2.1 Common Stats Collector uses

Access the stats collector through the `stats` attribute. Here is an example of an extension that access stats:

```
class ExtensionThatAccessStats(object):

    def __init__(self, stats):
        self.stats = stats

    @classmethod
    def from_crawler(cls, crawler):
        return cls(crawler.stats)
```

Set stat value:

```
stats.set_value('hostname', socket.gethostname())
```

Increment stat value:

```
stats.inc_value('custom_count')
```

Set stat value only if greater than previous:

```
stats.max_value('max_items_scraped', value)
```

Set stat value only if lower than previous:

```
stats.min_value('min_free_memory_percent', value)
```

Get stat value:

```
>>> stats.get_value('custom_count')
1
```

Get all stats:

```
>>> stats.get_stats()
{'custom_count': 1, 'start_time': datetime.datetime(2009, 7, 14, 21, 47, 28, 977139)}
```

4.2.2 Available Stats Collectors

Besides the basic `StatsCollector` there are other Stats Collectors available in Scrapy which extend the basic Stats Collector. You can select which Stats Collector to use through the `STATS_CLASS` setting. The default Stats Collector used is the `MemoryStatsCollector`.

MemoryStatsCollector

```
class scrapy.statscollectors.MemoryStatsCollector
```

A simple stats collector that keeps the stats of the last scraping run (for each spider) in memory, after they're closed. The stats can be accessed through the `spider_stats` attribute, which is a dict keyed by spider domain name.

This is the default Stats Collector used in Scrapy.

spider_stats

A dict of dicts (keyed by spider name) containing the stats of the last scraping run for each spider.

DummyStatsCollector

```
class scrapy.statscollectors.DummyStatsCollector
```

A Stats collector which does nothing but is very efficient (because it does nothing). This stats collector can be set via the `STATS_CLASS` setting, to disable stats collect in order to improve performance. However, the performance penalty of stats collection is usually marginal compared to other Scrapy workload like parsing pages.

4.3 发送邮件

Although Python makes sending e-mails relatively easy via the `smtplib` library, Scrapy provides its own facility for sending e-mails which is very easy to use and it's implemented using `Twisted non-blocking IO`, to avoid interfering with the non-blocking IO of the crawler. It also provides a simple API for sending attachments and it's very easy to configure, with a few *settings*.

4.3.1 Quick example

There are two ways to instantiate the mail sender. You can instantiate it using the standard constructor:

```
from scrapy.mail import MailSender
mailer = MailSender()
```

Or you can instantiate it passing a Scrapy settings object, which will respect the *settings*:

```
mailer = MailSender.from_settings(settings)
```

And here is how to use it to send an e-mail (without attachments):

```
mailer.send(to=["someone@example.com"], subject="Some subject", body="Some body", cc=[
    ↪ "another@example.com"])
```

4.3.2 MailSender class reference

MailSender is the preferred class to use for sending emails from Scrapy, as it uses Twisted non-blocking IO, like the rest of the framework.

```
class scrapy.mail.MailSender(smtphost=None, mailfrom=None, smtpuser=None, smtp-
                             pass=None, smtpport=None)
```

参数

- **smtphost** (*str or bytes*) – the SMTP host to use for sending the emails. If omitted, the *MAIL_HOST* setting will be used.
- **mailfrom** (*str*) – the address used to send emails (in the **From:** header). If omitted, the *MAIL_FROM* setting will be used.
- **smtpuser** – the SMTP user. If omitted, the *MAIL_USER* setting will be used. If not given, no SMTP authentication will be performed.
- **smtppass** (*str or bytes*) – the SMTP pass for authentication.
- **smtpport** (*int*) – the SMTP port to connect to
- **smtptls** (*boolean*) – enforce using SMTP STARTTLS
- **smtpssl** (*boolean*) – enforce using a secure SSL connection

```
classmethod from_settings(settings)
```

Instantiate using a Scrapy settings object, which will respect *these Scrapy settings*.

参数 settings (scrapy.settings.Settings object) – the e-mail recipients

```
send(to, subject, body, cc=None, attachs=(), mimetype='text/plain', charset=None)
```

Send email to the given recipients.

参数

- **to** (*str or list of str*) – the e-mail recipients

- **subject** (*str*) – the subject of the e-mail
- **cc** (*str or list of str*) – the e-mails to CC
- **body** (*str*) – the e-mail body
- **attachs** (*iterable*) – an iterable of tuples (**attach_name**, **mimetype**, **file_object**) where **attach_name** is a string with the name that will appear on the e-mail's attachment, **mimetype** is the mimetype of the attachment and **file_object** is a readable file object with the contents of the attachment
- **mimetype** (*str*) – the MIME type of the e-mail
- **charset** (*str*) – the character encoding to use for the e-mail contents

4.3.3 Mail settings

These settings define the default constructor values of the *MailSender* class, and can be used to configure e-mail notifications in your project without writing any code (for those extensions and code that uses *MailSender*).

MAIL_FROM

Default: 'scrapy@localhost'

Sender email to use (From: header) for sending emails.

MAIL_HOST

Default: 'localhost'

SMTP host to use for sending emails.

MAIL_PORT

Default: 25

SMTP port to use for sending emails.

MAIL_USER

Default: None

User to use for SMTP authentication. If disabled no SMTP authentication will be performed.

MAIL_PASS

Default: `None`

Password to use for SMTP authentication, along with *MAIL_USER*.

MAIL_TLS

Default: `False`

Enforce using STARTTLS. STARTTLS is a way to take an existing insecure connection, and upgrade it to a secure connection using SSL/TLS.

MAIL_SSL

Default: `False`

Enforce connecting using an SSL encrypted connection

4.4 远程控制台

Scrapy comes with a built-in telnet console for inspecting and controlling a Scrapy running process. The telnet console is just a regular python shell running inside the Scrapy process, so you can do literally anything from it.

The telnet console is a *built-in Scrapy extension* which comes enabled by default, but you can also disable it if you want. For more information about the extension itself see *Telnet console extension*.

警告: It is not secure to use telnet console via public networks, as telnet doesn't provide any transport-layer security. Having username/password authentication doesn't change that.

Intended usage is connecting to a running Scrapy spider locally (spider process and telnet client are on the same machine) or over a secure connection (VPN, SSH tunnel). Please avoid using telnet console over insecure connections, or disable it completely using *TELNETCONSOLE_ENABLED* option.

4.4.1 How to access the telnet console

The telnet console listens in the TCP port defined in the *TELNETCONSOLE_PORT* setting, which defaults to 6023. To access the console you need to type:

```
telnet localhost 6023
Trying localhost...
Connected to localhost.
Escape character is '^]'.
Username:
Password:
>>>
```

By default Username is **scrapy** and Password is autogenerated. The autogenerated Password can be seen on scrapy logs like the example bellow:

```
2018-10-16 14:35:21 [scrapy.extensions.telnet] INFO: Telnet Password: 16f92501e8a59326
```

Default Username and Password can be overridden by the settings `TELNETCONSOLE_USERNAME` and `TELNETCONSOLE_PASSWORD`.

警告: Username and password provide only a limited protection, as telnet is not using secure transport - by default traffic is not encrypted even if username and password are set.

You need the telnet program which comes installed by default in Windows, and most Linux distros.

4.4.2 Available variables in the telnet console

The telnet console is like a regular Python shell running inside the Scrapy process, so you can do anything from it including importing new modules, etc.

However, the telnet console comes with some default variables defined for convenience:

Shortcut	Description
<code>crawler</code>	the Scrapy Crawler (<i>scrapy.crawler.Crawler</i> object)
<code>engine</code>	Crawler.engine attribute
<code>spider</code>	the active spider
<code>slot</code>	the engine slot
<code>extensions</code>	the Extension Manager (Crawler.extensions attribute)
<code>stats</code>	the Stats Collector (Crawler.stats attribute)
<code>settings</code>	the Scrapy settings object (Crawler.settings attribute)
<code>est</code>	print a report of the engine status
<code>prefs</code>	for memory debugging (see 内存泄漏调试)
<code>p</code>	a shortcut to the <code>pprint.pprint</code> function
<code>hpy</code>	for memory debugging (see 内存泄漏调试)

4.4.3 Telnet console usage examples

Here are some example tasks you can do with the telnet console:

View engine status

You can use the `est()` method of the Scrapy engine to quickly show its state using the telnet console:

```
telnet localhost 6023
>>> est()
Execution engine status

time()-engine.start_time           : 8.62972998619
engine.has_capacity()               : False
len(engine.downloader.active)      : 16
engine.scrapers.is_idle()           : False
engine.spider.name                  : followall
engine.spider_is_idle(engine.spider): False
engine.slot.closing                 : False
len(engine.slot.inprogress)        : 16
len(engine.slot.scheduler.dqs or []) : 0
len(engine.slot.scheduler.mqs)     : 92
len(engine.scrapers.slot.queue)     : 0
len(engine.scrapers.slot.active)    : 0
engine.scrapers.slot.active_size    : 0
engine.scrapers.slot.itemproc_size  : 0
engine.scrapers.slot.needs_backout() : False
```

Pause, resume and stop the Scrapy engine

To pause:

```
telnet localhost 6023
>>> engine.pause()
>>>
```

To resume:

```
telnet localhost 6023
>>> engine.unpause()
>>>
```

To stop:

```
telnet localhost 6023
>>> engine.stop()
Connection closed by foreign host.
```

4.4.4 Telnet Console signals

`scrapy.extensions.telnet.update_telnet_vars(telnet_vars)`

Sent just before the telnet console is opened. You can hook up to this signal to add, remove or update the variables that will be available in the telnet local namespace. In order to do that, you need to update the `telnet_vars` dict in your handler.

参数 `telnet_vars` (*dict*) – the dict of telnet variables

4.4.5 Telnet settings

These are the settings that control the telnet console's behaviour:

TELNETCONSOLE_PORT

Default: [6023, 6073]

The port range to use for the telnet console. If set to `None` or `0`, a dynamically assigned port is used.

TELNETCONSOLE_HOST

Default: '127.0.0.1'

The interface the telnet console should listen on

TELNETCONSOLE_USERNAME

Default: 'scrapy'

The username used for the telnet console

TELNETCONSOLE_PASSWORD

Default: `None`

The password used for the telnet console, default behaviour is to have it autogenerated

4.5 Web Service

webservice has been moved into a separate project.

It is hosted at:

<https://github.com/scrapy-plugins/scrapy-jsonrpc>

日志 了解如何在 Scrapy 上使用 Python 的内置日志。

统计数据集合 收集有关你的抓取爬虫的统计数据。

发送邮件 当某些事件发生时发送电子邮件通知。

远程控制台 使用内置的 Python 控制台检查正在运行的爬虫器。

Web Service 使用 web 服务监视和控制爬虫程序。

5.1 常见问答

5.1.1 How does Scrapy compare to BeautifulSoup or lxml?

`BeautifulSoup` and `lxml` are libraries for parsing HTML and XML. Scrapy is an application framework for writing web spiders that crawl web sites and extract data from them.

Scrapy provides a built-in mechanism for extracting data (called *selectors*) but you can easily use `BeautifulSoup` (or `lxml`) instead, if you feel more comfortable working with them. After all, they're just parsing libraries which can be imported and used from any Python code.

In other words, comparing `BeautifulSoup` (or `lxml`) to Scrapy is like comparing `jinja2` to `Django`.

5.1.2 Can I use Scrapy with BeautifulSoup?

Yes, you can. As mentioned *above*, `BeautifulSoup` can be used for parsing HTML responses in Scrapy callbacks. You just have to feed the response's body into a `BeautifulSoup` object and extract whatever data you need from it.

Here's an example spider using BeautifulSoup API, with `lxml` as the HTML parser:

```
from bs4 import BeautifulSoup
import scrapy
```

(下页继续)

```
class ExampleSpider(scrapy.Spider):
    name = "example"
    allowed_domains = ["example.com"]
    start_urls = (
        'http://www.example.com/',
    )

    def parse(self, response):
        # use lxml to get decent HTML parsing speed
        soup = BeautifulSoup(response.text, 'lxml')
        yield {
            "url": response.url,
            "title": soup.h1.string
        }
```

注解: BeautifulSoup supports several HTML/XML parsers. See [BeautifulSoup's official documentation](#) on which ones are available.

5.1.3 What Python versions does Scrapy support?

Scrapy is supported under Python 2.7 and Python 3.4+ under CPython (default Python implementation) and PyPy (starting with PyPy 5.9). Python 2.6 support was dropped starting at Scrapy 0.20. Python 3 support was added in Scrapy 1.1. PyPy support was added in Scrapy 1.4, PyPy3 support was added in Scrapy 1.5.

注解: For Python 3 support on Windows, it is recommended to use Anaconda/Miniconda as *outlined in the installation guide*.

5.1.4 Did Scrapy “steal” X from Django?

Probably, but we don't like that word. We think Django is a great open source project and an example to follow, so we've used it as an inspiration for Scrapy.

We believe that, if something is already done well, there's no need to reinvent it. This concept, besides being one of the foundations for open source and free software, not only applies to software but also to documentation, procedures, policies, etc. So, instead of going through each problem ourselves, we choose to

copy ideas from those projects that have already solved them properly, and focus on the real problems we need to solve.

We'd be proud if Scrapy serves as an inspiration for other projects. Feel free to steal from us!

5.1.5 Does Scrapy work with HTTP proxies?

Yes. Support for HTTP proxies is provided (since Scrapy 0.8) through the HTTP Proxy downloader middleware. See [HttpProxyMiddleware](#).

5.1.6 How can I scrape an item with attributes in different pages?

See [Passing additional data to callback functions](#).

5.1.7 Scrapy crashes with: ImportError: No module named win32api

You need to install [pywin32](#) because of [this Twisted bug](#).

5.1.8 How can I simulate a user login in my spider?

See [Using FormRequest.from_response\(\) to simulate a user login](#).

5.1.9 Does Scrapy crawl in breadth-first or depth-first order?

By default, Scrapy uses a [LIFO](#) queue for storing pending requests, which basically means that it crawls in [DFO order](#). This order is more convenient in most cases. If you do want to crawl in true [BFO order](#), you can do it by setting the following settings:

```
DEPTH_PRIORITY = 1
SCHEDULER_DISK_QUEUE = 'scrapy.squeues.PickleFifoDiskQueue'
SCHEDULER_MEMORY_QUEUE = 'scrapy.squeues.FifoMemoryQueue'
```

5.1.10 My Scrapy crawler has memory leaks. What can I do?

See [内存泄漏调试](#).

Also, Python has a builtin memory leak issue which is described in [Leaks without leaks](#).

5.1.11 How can I make Scrapy consume less memory?

See previous question.

5.1.12 Can I use Basic HTTP Authentication in my spiders?

Yes, see *HttpAuthMiddleware*.

5.1.13 Why does Scrapy download pages in English instead of my native language?

Try changing the default `Accept-Language` request header by overriding the `DEFAULT_REQUEST_HEADERS` setting.

5.1.14 Where can I find some example Scrapy projects?

See 例子.

5.1.15 Can I run a spider without creating a project?

Yes. You can use the *runspider* command. For example, if you have a spider written in a `my_spider.py` file you can run it with:

```
scrapy runspider my_spider.py
```

See *runspider* command for more info.

5.1.16 I get “Filtered offsite request” messages. How can I fix them?

Those messages (logged with `DEBUG` level) don’ t necessarily mean there is a problem, so you may not need to fix them.

Those messages are thrown by the Offsite Spider Middleware, which is a spider middleware (enabled by default) whose purpose is to filter out requests to domains outside the ones covered by the spider.

For more info see: *OffsiteMiddleware*.

5.1.17 What is the recommended way to deploy a Scrapy crawler in production?

See 部署爬虫器.

5.1.18 Can I use JSON for large exports?

It’ ll depend on how large your output is. See *this warning* in *JsonItemExporter* documentation.

5.1.19 Can I return (Twisted) deferreds from signal handlers?

Some signals support returning deferreds from their handlers, others don't. See the *Built-in signals reference* to know which ones.

5.1.20 What does the response status code 999 means?

999 is a custom response status code used by Yahoo sites to throttle requests. Try slowing down the crawling speed by using a download delay of 2 (or higher) in your spider:

```
class MySpider(CrawlSpider):  
  
    name = 'myspider'  
  
    download_delay = 2  
  
    # [ ... rest of the spider code ... ]
```

Or by setting a global download delay in your project with the `DOWNLOAD_DELAY` setting.

5.1.21 Can I call `pdb.set_trace()` from my spiders to debug them?

Yes, but you can also use the Scrapy shell which allows you to quickly analyze (and even modify) the response being processed by your spider, which is, quite often, more useful than plain old `pdb.set_trace()`.

For more info see *Invoking the shell from spiders to inspect responses*.

5.1.22 Simplest way to dump all my scraped items into a JSON/CSV/XML file?

To dump into a JSON file:

```
scrapy crawl myspider -o items.json
```

To dump into a CSV file:

```
scrapy crawl myspider -o items.csv
```

To dump into a XML file:

```
scrapy crawl myspider -o items.xml
```

For more information see *Feed 导出*

5.1.23 What' s this huge cryptic `__VIEWSTATE` parameter used in some forms?

The `__VIEWSTATE` parameter is used in sites built with ASP.NET/VB.NET. For more info on how it works see [this page](#). Also, here' s an [example spider](#) which scrapes one of these sites.

5.1.24 What' s the best way to parse big XML/CSV data feeds?

Parsing big feeds with XPath selectors can be problematic since they need to build the DOM of the entire feed in memory, and this can be quite slow and consume a lot of memory.

In order to avoid parsing all the entire feed at once in memory, you can use the functions `xmliter` and `csviter` from `scrapy.utils.iterators` module. In fact, this is what the feed spiders (see [爬虫器](#)) use under the cover.

5.1.25 Does Scrapy manage cookies automatically?

Yes, Scrapy receives and keeps track of cookies sent by servers, and sends them back on subsequent requests, like any regular web browser does.

For more info see [请求和响应](#) and `CookiesMiddleware`.

5.1.26 How can I see the cookies being sent and received from Scrapy?

Enable the `COOKIES_DEBUG` setting.

5.1.27 How can I instruct a spider to stop itself?

Raise the `CloseSpider` exception from a callback. For more info see: `CloseSpider`.

5.1.28 How can I prevent my Scrapy bot from getting banned?

See [Avoiding getting banned](#).

5.1.29 Should I use spider arguments or settings to configure my spider?

Both *spider arguments* and *settings* can be used to configure your spider. There is no strict rule that mandates to use one or the other, but settings are more suited for parameters that, once set, don' t change much, while spider arguments are meant to change more often, even on each spider run and sometimes are required for the spider to run at all (for example, to set the start url of a spider).

To illustrate with an example, assuming you have a spider that needs to log into a site to scrape data, and you only want to scrape data from a certain section of the site (which varies each time). In that case, the credentials to log in would be settings, while the url of the section to scrape would be a spider argument.

5.1.30 I' m scraping a XML document and my XPath selector doesn' t return any items

You may need to remove namespaces. See *Removing namespaces*.

5.2 Debug 爬虫器

This document explains the most common techniques for debugging spiders. Consider the following scrapy spider below:

```
import scrapy
from myproject.items import MyItem

class MySpider(scrapy.Spider):
    name = 'myspider'
    start_urls = (
        'http://example.com/page1',
        'http://example.com/page2',
    )

    def parse(self, response):
        # <processing code not shown>
        # collect `item_urls`
        for item_url in item_urls:
            yield scrapy.Request(item_url, self.parse_item)

    def parse_item(self, response):
        # <processing code not shown>
        item = MyItem()
        # populate `item` fields
        # and extract item_details_url
        yield scrapy.Request(item_details_url, self.parse_details, meta={'item': item})

    def parse_details(self, response):
        item = response.meta['item']
```

(下页继续)

(续上页)

```
# populate more `item` fields
return item
```

Basically this is a simple spider which parses two pages of items (the `start_urls`). Items also have a details page with additional information, so we use the `meta` functionality of *Request* to pass a partially populated item.

5.2.1 Parse Command

The most basic way of checking the output of your spider is to use the *parse* command. It allows to check the behaviour of different parts of the spider at the method level. It has the advantage of being flexible and simple to use, but does not allow debugging code inside a method.

In order to see the item scraped from a specific url:

```
$ scrapy parse --spider=mypider -c parse_item -d 2 <item_url>
[ ... scrapy log lines crawling example.com spider ... ]

>>> STATUS DEPTH LEVEL 2 <<<
# Scraped Items -----
[{'url': <item_url>}]

# Requests -----
[]
```

Using the `--verbose` or `-v` option we can see the status at each depth level:

```
$ scrapy parse --spider=mypider -c parse_item -d 2 -v <item_url>
[ ... scrapy log lines crawling example.com spider ... ]

>>> DEPTH LEVEL: 1 <<<
# Scraped Items -----
[]

# Requests -----
[<GET item_details_url>]

>>> DEPTH LEVEL: 2 <<<
# Scraped Items -----
[{'url': <item_url>}]
```

(下页继续)

(续上页)

```
# Requests -----  
[]
```

Checking items scraped from a single start_url, can also be easily achieved using:

```
$ scrapy parse --spider=mypider -d 3 'http://example.com/page1'
```

5.2.2 Scrapy Shell

While the `parse` command is very useful for checking behaviour of a spider, it is of little help to check what happens inside a callback, besides showing the response received and the output. How to debug the situation when `parse_details` sometimes receives no item?

Fortunately, the `shell` is your bread and butter in this case (see *Invoking the shell from spiders to inspect responses*):

```
from scrapy.shell import inspect_response  
  
def parse_details(self, response):  
    item = response.meta.get('item', None)  
    if item:  
        # populate more `item` fields  
        return item  
    else:  
        inspect_response(response, self)
```

See also: *Invoking the shell from spiders to inspect responses*.

5.2.3 Open in browser

Sometimes you just want to see how a certain response looks in a browser, you can use the `open_in_browser` function for that. Here is an example of how you would use it:

```
from scrapy.utils.response import open_in_browser  
  
def parse_details(self, response):  
    if "item name" not in response.body:  
        open_in_browser(response)
```

`open_in_browser` will open a browser with the response received by Scrapy at that point, adjusting the `base` tag so that images and styles are displayed properly.

5.2.4 Logging

Logging is another useful option for getting information about your spider run. Although not as convenient, it comes with the advantage that the logs will be available in all future runs should they be necessary again:

```
def parse_details(self, response):
    item = response.meta.get('item', None)
    if item:
        # populate more `item` fields
        return item
    else:
        self.logger.warning('No item received for %s', response.url)
```

For more information, check the `日志` section.

5.3 爬虫器约束

0.15 新版功能.

注解: This is a new feature (introduced in Scrapy 0.15) and may be subject to minor functionality/API updates. Check the *release notes* to be notified of updates.

Testing spiders can get particularly annoying and while nothing prevents you from writing unit tests the task gets cumbersome quickly. Scrapy offers an integrated way of testing your spiders by the means of contracts.

This allows you to test each callback of your spider by hardcoding a sample url and check various constraints for how the callback processes the response. Each contract is prefixed with an `@` and included in the docstring. See the following example:

```
def parse(self, response):
    """ This function parses a sample response. Some contracts are mingled
    with this docstring.

    @url http://www.amazon.com/s?field-keywords=selfish+gene
    @returns items 1 16
    @returns requests 0 0
    @scrapes Title Author Year Price
    """
```


This callback is tested using three built-in contracts:

class scrapy.contracts.default.UrlContract

This contract (`@url`) sets the sample url used when checking other contract conditions for this spider.

This contract is mandatory. All callbacks lacking this contract are ignored when running the checks:

```
@url url
```

class scrapy.contracts.default.ReturnsContract

This contract (`@returns`) sets lower and upper bounds for the items and requests returned by the spider. The upper bound is optional:

```
@returns item(s)|request(s) [min [max]]
```

class scrapy.contracts.default.ScrapesContract

This contract (`@scrapes`) checks that all the items returned by the callback have the specified fields:

```
@scrapes field_1 field_2 ...
```

Use the `check` command to run the contract checks.

5.3.1 Custom Contracts

If you find you need more power than the built-in scrapy contracts you can create and load your own contracts in the project by using the `SPIDER_CONTRACTS` setting:

```
SPIDER_CONTRACTS = {
    'myproject.contracts.ResponseCheck': 10,
    'myproject.contracts.ItemValidate': 10,
}
```

Each contract must inherit from `scrapy.contracts.Contract` and can override three methods:

class scrapy.contracts.Contract(*method*, **args*)

参数

- **method** (*function*) – callback function to which the contract is associated
- **args** (*list*) – list of arguments passed into the docstring (whitespace separated)

adjust_request_args(*args*)

This receives a `dict` as an argument containing default arguments for request object. `Request` is used by default, but this can be changed with the `request_cls` attribute. If multiple contracts in chain have this attribute defined, the last one is used.

Must return the same or a modified version of it.

`pre_process(response)`

This allows hooking in various checks on the response received from the sample request, before it's being passed to the callback.

`post_process(output)`

This allows processing the output of the callback. Iterators are converted listified before being passed to this hook.

Here is a demo contract which checks the presence of a custom header in the response received. Raise `scrapy.exceptions.ContractFail` in order to get the failures pretty printed:

```
from scrapy.contracts import Contract
from scrapy.exceptions import ContractFail

class HasHeaderContract(Contract):
    """ Demo contract which checks the presence of a custom header
        @has_header X-CustomHeader
    """

    name = 'has_header'

    def pre_process(self, response):
        for header in self.args:
            if header not in response.headers:
                raise ContractFail('X-CustomHeader not present')
```

5.4 常见实践

This section documents common practices when using Scrapy. These are things that cover many topics and don't often fall into any other specific section.

5.4.1 Run Scrapy from a script

You can use the [API](#) to run Scrapy from a script, instead of the typical way of running Scrapy via `scrapy crawl`.

Remember that Scrapy is built on top of the Twisted asynchronous networking library, so you need to run it inside the Twisted reactor.

The first utility you can use to run your spiders is `scrapy.crawler.CrawlerProcess`. This class will start a Twisted reactor for you, configuring the logging and setting shutdown handlers. This class is the one used by all Scrapy commands.

Here's an example showing how to run a single spider with it.

```
import scrapy
from scrapy.crawler import CrawlerProcess

class MySpider(scrapy.Spider):
    # Your spider definition
    ...

process = CrawlerProcess({
    'USER_AGENT': 'Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1)'
})

process.crawl(MySpider)
process.start() # the script will block here until the crawling is finished
```

Make sure to check `CrawlerProcess` documentation to get acquainted with its usage details.

If you are inside a Scrapy project there are some additional helpers you can use to import those components within the project. You can automatically import your spiders passing their name to `CrawlerProcess`, and use `get_project_settings` to get a `Settings` instance with your project settings.

What follows is a working example of how to do that, using the `testspiders` project as example.

```
from scrapy.crawler import CrawlerProcess
from scrapy.utils.project import get_project_settings

process = CrawlerProcess(get_project_settings())

# 'followall' is the name of one of the spiders of the project.
process.crawl('followall', domain='scrapinghub.com')
process.start() # the script will block here until the crawling is finished
```

There's another Scrapy utility that provides more control over the crawling process: `scrapy.crawler.CrawlerRunner`. This class is a thin wrapper that encapsulates some simple helpers to run multiple crawlers, but it won't start or interfere with existing reactors in any way.

Using this class the reactor should be explicitly run after scheduling your spiders. It's recommended you use `CrawlerRunner` instead of `CrawlerProcess` if your application is already using Twisted and you want to run Scrapy in the same reactor.

Note that you will also have to shutdown the Twisted reactor yourself after the spider is finished. This can be achieved by adding callbacks to the deferred returned by the `CrawlerRunner.crawl` method.

Here's an example of its usage, along with a callback to manually stop the reactor after `MySpider` has finished running.

```
from twisted.internet import reactor
import scrapy
from scrapy.crawler import CrawlerRunner
from scrapy.utils.log import configure_logging

class MySpider(scrapy.Spider):
    # Your spider definition
    ...

configure_logging({'LOG_FORMAT': '%(levelname)s: %(message)s'})
runner = CrawlerRunner()

d = runner.crawl(MySpider)
d.addBoth(lambda _: reactor.stop())
reactor.run() # the script will block here until the crawling is finished
```

参见:

[Twisted Reactor Overview](#).

5.4.2 Running multiple spiders in the same process

By default, Scrapy runs a single spider per process when you run `scrapy crawl`. However, Scrapy supports running multiple spiders per process using the *internal API*.

Here is an example that runs multiple spiders simultaneously:

```
import scrapy
from scrapy.crawler import CrawlerProcess

class MySpider1(scrapy.Spider):
    # Your first spider definition
    ...

class MySpider2(scrapy.Spider):
    # Your second spider definition
    ...

process = CrawlerProcess()
process.crawl(MySpider1)
process.crawl(MySpider2)
process.start() # the script will block here until all crawling jobs are finished
```

Same example using CrawlerRunner:

```
import scrapy
from twisted.internet import reactor
from scrapy.crawler import CrawlerRunner
from scrapy.utils.log import configure_logging

class MySpider1(scrapy.Spider):
    # Your first spider definition
    ...

class MySpider2(scrapy.Spider):
    # Your second spider definition
    ...

configure_logging()
runner = CrawlerRunner()
runner.crawl(MySpider1)
runner.crawl(MySpider2)
d = runner.join()
d.addBoth(lambda _: reactor.stop())

reactor.run() # the script will block here until all crawling jobs are finished
```

Same example but running the spiders sequentially by chaining the deferreds:

```
from twisted.internet import reactor, defer
from scrapy.crawler import CrawlerRunner
from scrapy.utils.log import configure_logging

class MySpider1(scrapy.Spider):
    # Your first spider definition
    ...

class MySpider2(scrapy.Spider):
    # Your second spider definition
    ...

configure_logging()
runner = CrawlerRunner()

@defer.inlineCallbacks
```

(下页继续)

(续上页)

```
def crawl():
    yield runner.crawl(MySpider1)
    yield runner.crawl(MySpider2)
    reactor.stop()

crawl()
reactor.run() # the script will block here until the last crawl call is finished
```

参见:

Run Scrapy from a script.

5.4.3 Distributed crawls

Scrapy doesn't provide any built-in facility for running crawls in a distribute (multi-server) manner. However, there are some ways to distribute crawls, which vary depending on how you plan to distribute them.

If you have many spiders, the obvious way to distribute the load is to setup many Scrapy instances and distribute spider runs among those.

If you instead want to run a single (big) spider through many machines, what you usually do is partition the urls to crawl and send them to each separate spider. Here is a concrete example:

First, you prepare the list of urls to crawl and put them into separate files/urls:

```
http://somedomain.com/urls-to-crawl/spider1/part1.list
http://somedomain.com/urls-to-crawl/spider1/part2.list
http://somedomain.com/urls-to-crawl/spider1/part3.list
```

Then you fire a spider run on 3 different Scrapy servers. The spider would receive a (spider) argument `part` with the number of the partition to crawl:

```
curl http://scrapy1.mycompany.com:6800/schedule.json -d project=myproject -d␣
↪spider=spider1 -d part=1
curl http://scrapy2.mycompany.com:6800/schedule.json -d project=myproject -d␣
↪spider=spider1 -d part=2
curl http://scrapy3.mycompany.com:6800/schedule.json -d project=myproject -d␣
↪spider=spider1 -d part=3
```

5.4.4 Avoiding getting banned

Some websites implement certain measures to prevent bots from crawling them, with varying degrees of sophistication. Getting around those measures can be difficult and tricky, and may sometimes require

special infrastructure. Please consider contacting [commercial support](#) if in doubt.

Here are some tips to keep in mind when dealing with these kinds of sites:

- rotate your user agent from a pool of well-known ones from browsers (google around to get a list of them)
- disable cookies (see [COOKIES_ENABLED](#)) as some sites may use cookies to spot bot behaviour
- use download delays (2 or higher). See [DOWNLOAD_DELAY](#) setting.
- if possible, use [Google cache](#) to fetch pages, instead of hitting the sites directly
- use a pool of rotating IPs. For example, the free [Tor project](#) or paid services like [ProxyMesh](#). An open source alternative is [scrapoxy](#), a super proxy that you can attach your own proxies to.
- use a highly distributed downloader that circumvents bans internally, so you can just focus on parsing clean pages. One example of such downloaders is [Crawlera](#)

If you are still unable to prevent your bot getting banned, consider contacting [commercial support](#).

5.5 并发爬虫

Scrapy 默认优化爬取特定的网站。这些站点通常只使用一个爬虫器来爬取, 虽然这不是必需的 (例如, 一些通用爬虫器可以处理任何抛给它们的站点)。

除了这种“集中爬虫”, 还有一些常见的爬虫类型, 包含了大量 (可能是无限) 的 domains , 并且只接受时间或者其他任意的约束, 而不是当 domain 被爬取完或者没有其他请求时停止。这种被称作“广泛爬虫”, 是搜索引擎常用的一种典型爬虫。

这是一些广泛爬虫中常见的属性:

- 有很多 domain (通常无边界) 而不是一组特定的站点
- 它们不一定要爬完整个 domain , 因为可能不会那样做, 而是限制爬虫时间或者爬取的页数
- 它们的逻辑简单 (和一些有着许多解析规则的非常复杂的爬虫器相比) 因为数据通常在一个单独的阶段进行 POST 请求
- 同时爬取多个 domain , 让它们用更快的速度爬取而不受任何站点限制 (出于尊重单独的站点用比较慢的方式爬取, 但是很多站点的时候就要并行了)

正如上面说的, Scrapy 默认设置是优化集中爬虫, 而不是广泛爬虫。不过, 基于它的异步架构, Scrapy 非常适合快速广泛爬虫。本页总结了一些 Scrapy 广泛爬虫时需要注意的内容, 以及一些 Scrapy 配置的建议, 以便于有效地实现广泛爬虫。

5.5.1 如何使用 SCHEDULER_PRIORITY_QUEUE

Scrapy 默认是优先使用 '[scrapy.pqueues.ScrapyPriorityQueue](#)' 队列进行调度。它在单独 domain 的爬取表现是最好的。反之在并行爬取不同的 domains 时表现一般。

优先队列的使用:

```
SCHEDULER_PRIORITY_QUEUE = 'scrapy.pqueues.DownloaderAwarePriorityQueue'
```

5.5.2 增加并发

并发意味着有许多讲求同时被处理。有一个全局限制和单独 domain 的限制。

Scrapy 默认的全局并发限制不适合并行爬虫多个不同 domains , 因此你需要增加。增加多少取决于你的 CPU 能给你的爬虫提供多少资源。100 是个比较好的起点, 但是最好的方法是通过测试来找出你的 Scrapy 在多少并发下会受到 CPU 限制。最佳的情况是, 你的 CPU 利用率在 80-90% 之间。

增加全局并发:

```
CONCURRENT_REQUESTS = 100
```

5.5.3 增加 Twisted IO 线程池的最大量

目前 Scrapy 是通过纯种池以阻塞的方式进行 DNS 解析。如果并发太高爬虫可能会变慢甚至 DNS 解析超时而失败。通过增加处理 DNS 查询的线程数量来解决。DNS 队列将被更快地处理, 加快建立连接和爬虫。

增加最大线程池的大小:

```
REACTOR_THREADPOOL_MAXSIZE = 20
```

5.5.4 设置你自己的 DNS

如果你有多个爬虫进程和单个中心 DNS, 它会像 DoS 攻击 DNS 服务器一样导致整个网络变慢甚至阻塞你的计算机。为了避免这样, 设置你自己的 DNS 服务通过本地缓存和一些上游大型 DNS 如 OpenDNS 或 Verizon。

5.5.5 降低 log 等级

当进行广泛地爬虫时, 你通常只对抓取速率和所有的错误感兴趣。这些统计数据当 Scrapy 的 log 等级是 INFO 时会被记录。为了节省 CPU (和日志存储需求) 你不应该使用 DEBUG log 等级。在开发时使用 DEBUG 等级应该没问题。

设置 log 等级:

```
LOG_LEVEL = 'INFO'
```


5.5.6 禁用 cookies

禁用 cookies 除非你 确实需要。在进行大范围爬取时，通常不需要 cookie（搜索引擎会忽略它们），它们通过节省 CPU 周期和减少内在占用来提高 Scrapy 爬虫的性能。

设置禁用 cookies:

```
COOKIES_ENABLED = False
```

5.5.7 禁用重试

重试失败的 HTTP 请求会显著降低爬虫的速度，尤其是由于站点原因造成的响应慢（或者失败）而造成的超时错误导致过多重试，会阻止爬虫的生产力用到其他 domains。

设置禁用重试:

```
RETRY_ENABLED = False
```

5.5.8 降低下载超时

除非你正从一个非常缓慢的链接爬虫（不该在大范围爬虫出现的情况）减少下载超时以便快速丢弃卡住的情况，释放处理下一个请求的生产力。

设置降低下载超时:

```
DOWNLOAD_TIMEOUT = 15
```

5.5.9 禁止重定向

考虑禁用重定向，除非你很有兴趣跟进它们。当进行广泛爬虫时，通常会保存重定向在以后的爬取时重新访问站点时解析它们。这样也可以保持每批爬取请求的数量不变，否则重定向的循环可能会导致爬虫在任意的 domain 中定向到太多的资源中去。

设置禁用重定向:

```
REDIRECT_ENABLED = False
```

5.5.10 启用爬取 “Ajax 页面爬取”

一些页面 (根据 2013 年的数据统计，高达 1%) 声名它们为 `ajax crawlable`。这意味着他们只能用 AJAX 获得 HTML 中的内容。页面可能用两种方法表示它:

- 1) 在 URL 中使用 `#!` - 这是默认的方法;

2) 使用特殊的 meta 标签 - 这个方法通常用于 “main” 和 “index” 页面。

Scrapy 自动处理 (1); 处理 (2) 使能 *AjaxCrawlMiddleware*:

```
AJAXCRAWL_ENABLED = True
```

当进行广泛的爬虫时经常会爬取很多 “index” 页面; *AjaxCrawlMiddleware* 帮助正确地爬取它们。因为它会有一些性能开销所以默认情况下它是关闭的。而且在集中爬取时启用它也没什么用。

5.6 使用浏览器的开发工具进行抓取

Here is a general guide on how to use your browser’ s Developer Tools to ease the scraping process. Today almost all browsers come with built in *Developer Tools* and although we will use Firefox in this guide, the concepts are applicable to any other browser.

In this guide we’ ll introduce the basic tools to use from a browser’ s Developer Tools by scraping quotes.toscrape.com.

5.6.1 Caveats with inspecting the live browser DOM

Since Developer Tools operate on a live browser DOM, what you’ ll actually see when inspecting the page source is not the original HTML, but a modified one after applying some browser clean up and executing Javascript code. Firefox, in particular, is known for adding `<tbody>` elements to tables. Scrapy, on the other hand, does not modify the original page HTML, so you won’ t be able to extract any data if you use `<tbody>` in your XPath expressions.

Therefore, you should keep in mind the following things:

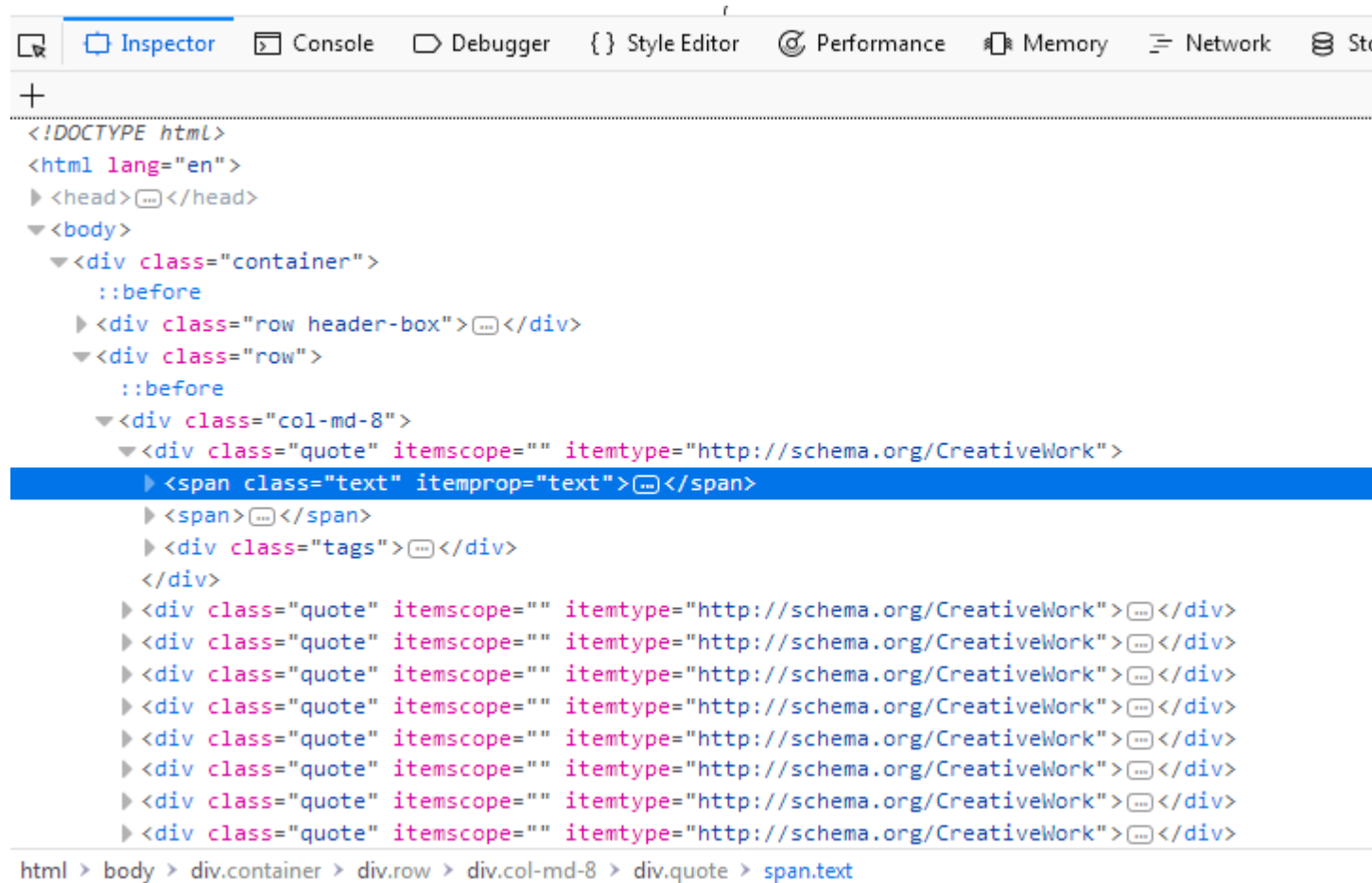
- Disable Javascript while inspecting the DOM looking for XPaths to be used in Scrapy (in the Developer Tools settings click *Disable JavaScript*)
- Never use full XPath paths, use relative and clever ones based on attributes (such as `id`, `class`, `width`, etc) or any identifying features like `contains(@href, 'image')`.
- Never include `<tbody>` elements in your XPath expressions unless you really know what you’ re doing

5.6.2 Inspecting a website

By far the most handy feature of the Developer Tools is the *Inspector* feature, which allows you to inspect the underlying HTML code of any webpage. To demonstrate the Inspector, let’ s look at the quotes.toscrape.com site.

On the site we have a total of ten quotes from various authors with specific tags, as well as the Top Ten Tags. Let’ s say we want to extract all the quotes on this page, without any meta-information about authors, tags, etc.

Instead of viewing the whole source code for the page, we can simply right click on a quote and select **Inspect Element (Q)**, which opens up the *Inspector*. In it you should see something like this:



The interesting part for us is this:

```
<div class="quote" itemscope="" itemtype="http://schema.org/CreativeWork">
  <span class="text" itemprop="text">(...)</span>
  <span>(...)</span>
  <div class="tags">(...)</div>
</div>
```

If you hover over the first `div` directly above the `span` tag highlighted in the screenshot, you'll see that the corresponding section of the webpage gets highlighted as well. So now we have a section, but we can't find our quote text anywhere.

The advantage of the *Inspector* is that it automatically expands and collapses sections and tags of a webpage, which greatly improves readability. You can expand and collapse a tag by clicking on the arrow in front of it or by double clicking directly on the tag. If we expand the `span` tag with the `class= "text"` we will see the quote-text we clicked on. The *Inspector* lets you copy XPaths to selected elements. Let's try it out: Right-click on the `span` tag, select **Copy > XPath** and paste it in the scrapy shell like so:

```
$ scrapy shell "http://quotes.toscrape.com/"
(...)
>>> response.xpath('/html/body/div/div[2]/div[1]/div[1]/span[1]/text()').getall()
['"The world as we have created it is a process of our thinking. It cannot be changed
↳without changing our thinking." ]
```

Adding `text()` at the end we are able to extract the first quote with this basic selector. But this XPath is not really that clever. All it does is go down a desired path in the source code starting from `html`. So let's see if we can refine our XPath a bit:

If we check the *Inspector* again we'll see that directly beneath our expanded `div` tag we have nine identical `div` tags, each with the same attributes as our first. If we expand any of them, we'll see the same structure as with our first quote: Two `span` tags and one `div` tag. We can expand each `span` tag with the `class="text"` inside our `div` tags and see each quote:

```
<div class="quote" itemscope="" itemtype="http://schema.org/CreativeWork">
  <span class="text" itemprop="text">
    "The world as we have created it is a process of our thinking. It cannot be changed
↳without changing our thinking."
  </span>
  <span>(...)</span>
  <div class="tags">(...)</div>
</div>
```

With this knowledge we can refine our XPath: Instead of a path to follow, we'll simply select all `span` tags with the `class="text"` by using the `has-class-extension`:

```
>>> response.xpath('//span[has-class("text")]/text()').getall()
['"The world as we have created it is a process of our thinking. It cannot be changed
↳without changing our thinking." ,
 ' "It is our choices, Harry, that show what we truly are, far more than our abilities."
↳',
 ' "There are only two ways to live your life. One is as though nothing is a miracle.
↳The other is as though everything is a miracle." ',
 (...)]
```

And with one simple, cleverer XPath we are able to extract all quotes from the page. We could have constructed a loop over our first XPath to increase the number of the last `div`, but this would have been unnecessarily complex and by simply constructing an XPath with `has-class("text")` we were able to extract all quotes in one line.

The *Inspector* has a lot of other helpful features, such as searching in the source code or directly scrolling to an element you selected. Let's demonstrate a use case:

Say you want to find the **Next** button on the page. Type **Next** into the search bar on the top right of the *Inspector*. You should get two results. The first is a `li` tag with the `class="text"`, the second the text of an `a` tag. Right click on the `a` tag and select **Scroll into View**. If you hover over the tag, you'll see the button highlighted. From here we could easily create a *Link Extractor* to follow the pagination. On a simple site such as this, there may not be the need to find an element visually but the **Scroll into View** function can be quite useful on complex sites.

Note that the search bar can also be used to search for and test CSS selectors. For example, you could search for `span.text` to find all quote texts. Instead of a full text search, this searches for exactly the `span` tag with the `class="text"` in the page.

5.6.3 The Network-tool

While scraping you may come across dynamic webpages where some parts of the page are loaded dynamically through multiple requests. While this can be quite tricky, the *Network-tool* in the Developer Tools greatly facilitates this task. To demonstrate the Network-tool, let's take a look at the page quotes.toscrape.com/scroll.

The page is quite similar to the basic quotes.toscrape.com-page, but instead of the above-mentioned **Next** button, the page automatically loads new quotes when you scroll to the bottom. We could go ahead and try out different XPath expressions directly, but instead we'll check another quite useful command from the scrapy shell:

```
$ scrapy shell "quotes.toscrape.com/scroll"
(...)
>>> view(response)
```

A browser window should open with the webpage but with one crucial difference: Instead of the quotes we just see a greenish bar with the word **Loading...**



The `view(response)` command lets us view the response our shell or later our spider receives from the

server. Here we see that some basic template is loaded which includes the title, the login-button and the footer, but the quotes are missing. This tells us that the quotes are being loaded from a different request than `quotes.toscrape/scroll`.

If you click on the **Network** tab, you will probably only see two entries. The first thing we do is enable persistent logs by clicking on **Persist Logs**. If this option is disabled, the log is automatically cleared each time you navigate to a different page. Enabling this option is a good default, since it gives us control on when to clear the logs.

If we reload the page now, you'll see the log get populated with six new requests.

Status	Method	File	Domain	Cause	Type	Transferred	Size	0 ms	1.37 min
200	GET	main.css	quotes.toscra...	stylesheet	css	cached	1.34 kB		
200	GET	bootstrap.min.css	quotes.toscra...	stylesheet	css	cached	122.98 kB		
200	GET	scroll	quotes.toscra...	document	html	1.27 kB	2.60 kB		
304	GET	bootstrap.min.css	quotes.toscra...	stylesheet	css	cached	122.98 kB		
304	GET	main.css	quotes.toscra...	stylesheet	css	cached	1.34 kB		
304	GET	jquery.js	quotes.toscra...	script	js	cached	82.37 kB		
200	GET	css?family=Raleway:40...	fonts.googlea...	stylesheet	css	861 B	1.52 kB		
200	GET	quotes?page=1	quotes.toscra...	xhr	json	1.34 kB	2.95 kB		

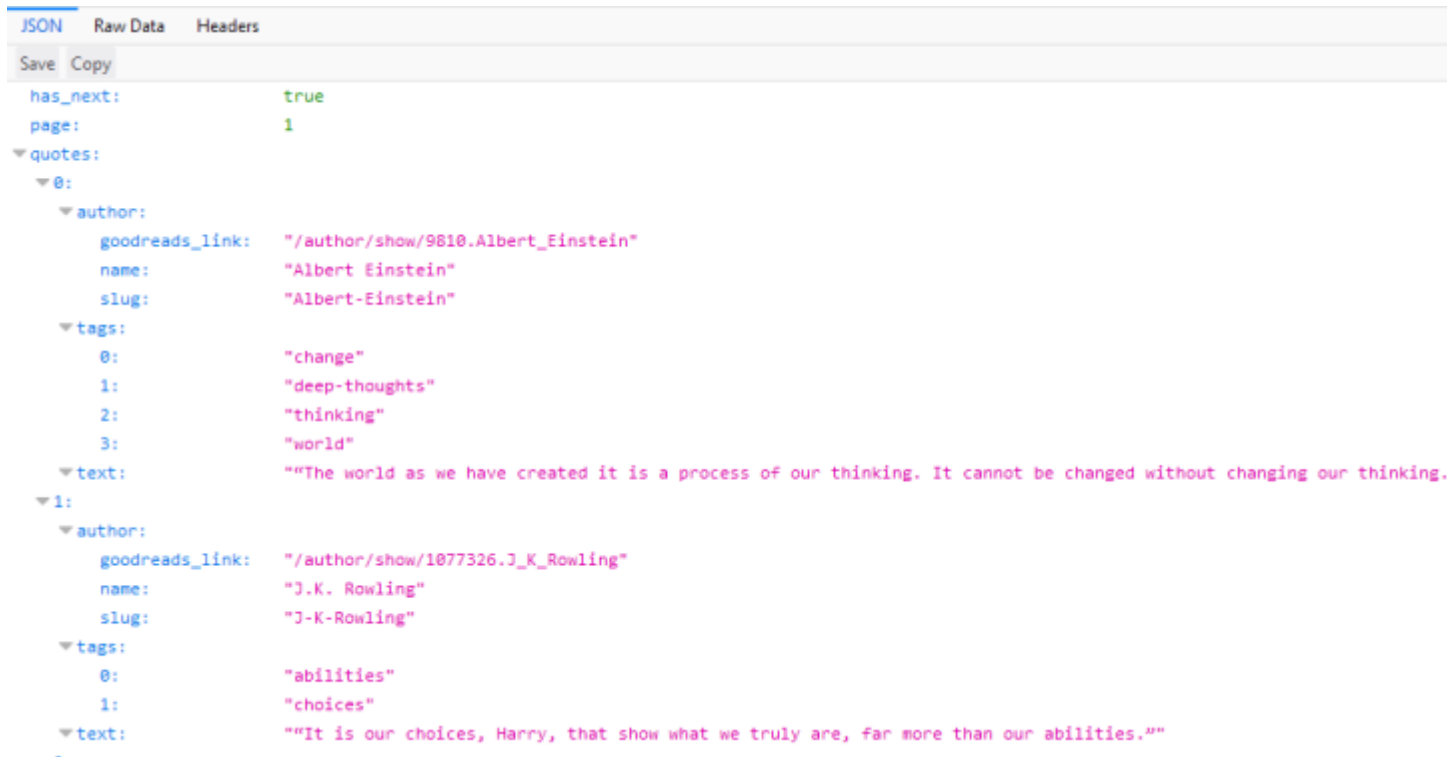
8 requests | 338.07 kB / 64.20 kB transferred | Finish: 3.19 min | DOMContentLoaded: 3.18 min | load: 3.19 min

Here we see every request that has been made when reloading the page and can inspect each request and its response. So let's find out where our quotes are coming from:

First click on the request with the name `scroll`. On the right you can now inspect the request. In **Headers** you'll find details about the request headers, such as the URL, the method, the IP-address, and so on. We'll ignore the other tabs and click directly on **Response**.

What you should see in the **Preview** pane is the rendered HTML-code, that is exactly what we saw when we called `view(response)` in the shell. Accordingly the **type** of the request in the log is `html`. The other requests have types like `css` or `js`, but what interests us is the one request called `quotes?page=1` with the type `json`.

If we click on this request, we see that the request URL is `http://quotes.toscrape.com/api/quotes?page=1` and the response is a JSON-object that contains our quotes. We can also right-click on the request and open **Open in new tab** to get a better overview.



With this response we can now easily parse the JSON-object and also request each page to get every quote on the site:

```

import scrapy
import json

class QuoteSpider(scrapy.Spider):
    name = 'quote'
    allowed_domains = ['quotes.toscrape.com']
    page = 1
    start_urls = ['http://quotes.toscrape.com/api/quotes?page=1']

    def parse(self, response):
        data = json.loads(response.text)
        for quote in data["quotes"]:
            yield {"quote": quote["text"]}
        if data["has_next"]:
            self.page += 1
            url = "http://quotes.toscrape.com/api/quotes?page={}".format(self.page)
            yield scrapy.Request(url=url, callback=self.parse)

```

This spider starts at the first page of the quotes-API. With each response, we parse the `response.text`

and assign it to `data`. This lets us operate on the JSON-object like on a Python dictionary. We iterate through the `quotes` and print out the `quote["text"]`. If the handy `has_next` element is `true` (try loading quotes.toscrape.com/api/quotes?page=10 in your browser or a page-number greater than 10), we increment the `page` attribute and `yield` a new request, inserting the incremented page-number into our `url`.

You can see that with a few inspections in the *Network*-tool we were able to easily replicate the dynamic requests of the scrolling functionality of the page. Crawling dynamic pages can be quite daunting and pages can be very complex, but it (mostly) boils down to identifying the correct request and replicating it in your spider.

5.7 内存泄漏调试

In Scrapy, objects such as Requests, Responses and Items have a finite lifetime: they are created, used for a while, and finally destroyed.

From all those objects, the Request is probably the one with the longest lifetime, as it stays waiting in the Scheduler queue until it's time to process it. For more info see [框架体系](#).

As these Scrapy objects have a (rather long) lifetime, there is always the risk of accumulating them in memory without releasing them properly and thus causing what is known as a “memory leak”.

To help debugging memory leaks, Scrapy provides a built-in mechanism for tracking objects references called *trackref*, and you can also use a third-party library called *Guppy* for more advanced memory debugging (see below for more info). Both mechanisms must be used from the *Telnet Console*.

5.7.1 Common causes of memory leaks

It happens quite often (sometimes by accident, sometimes on purpose) that the Scrapy developer passes objects referenced in Requests (for example, using the *meta* attribute or the request callback function) and that effectively bounds the lifetime of those referenced objects to the lifetime of the Request. This is, by far, the most common cause of memory leaks in Scrapy projects, and a quite difficult one to debug for newcomers.

In big projects, the spiders are typically written by different people and some of those spiders could be “leaking” and thus affecting the rest of the other (well-written) spiders when they get to run concurrently, which, in turn, affects the whole crawling process.

The leak could also come from a custom middleware, pipeline or extension that you have written, if you are not releasing the (previously allocated) resources properly. For example, allocating resources on *spider_opened* but not releasing them on *spider_closed* may cause problems if you're running *multiple spiders per process*.

Too Many Requests?

By default Scrapy keeps the request queue in memory; it includes *Request* objects and all objects referenced in Request attributes (e.g. in *meta*). While not necessarily a leak, this can take a lot of memory. Enabling *persistent job queue* could help keeping memory usage in control.

5.7.2 Debugging memory leaks with trackref

`trackref` is a module provided by Scrapy to debug the most common cases of memory leaks. It basically tracks the references to all live Requests, Responses, Item and Selector objects.

You can enter the telnet console and inspect how many objects (of the classes mentioned above) are currently alive using the `prefs()` function which is an alias to the `print_live_refs()` function:

```
telnet localhost 6023

>>> prefs()
Live References

ExampleSpider          1   oldest: 15s ago
HtmlResponse          10   oldest: 1s ago
Selector               2   oldest: 0s ago
FormRequest           878   oldest: 7s ago
```

As you can see, that report also shows the “age” of the oldest object in each class. If you’re running multiple spiders per process chances are you can figure out which spider is leaking by looking at the oldest request or response. You can get the oldest object of each class using the `get_oldest()` function (from the telnet console).

Which objects are tracked?

The objects tracked by `trackrefs` are all from these classes (and all its subclasses):

- `scrapy.http.Request`
- `scrapy.http.Response`
- `scrapy.item.Item`
- `scrapy.selector.Selector`
- `scrapy.spiders.Spider`

A real example

Let's see a concrete example of a hypothetical case of memory leaks. Suppose we have some spider with a line similar to this one:

```
return Request("http://www.somenastyspider.com/product.php?pid=%d" % product_id,
               callback=self.parse, meta={referer: response})
```

That line is passing a response reference inside a request which effectively ties the response lifetime to the requests' one, and that would definitely cause memory leaks.

Let's see how we can discover the cause (without knowing it a-priori, of course) by using the `trackref` tool.

After the crawler is running for a few minutes and we notice its memory usage has grown a lot, we can enter its telnet console and check the live references:

```
>>> prefs()
Live References

SomenastySpider           1  oldest: 15s ago
HtmlResponse              3890 oldest: 265s ago
Selector                   2  oldest: 0s ago
Request                   3878 oldest: 250s ago
```

The fact that there are so many live responses (and that they're so old) is definitely suspicious, as responses should have a relatively short lifetime compared to Requests. The number of responses is similar to the number of requests, so it looks like they are tied in some way. We can now go and check the code of the spider to discover the nasty line that is generating the leaks (passing response references inside requests).

Sometimes extra information about live objects can be helpful. Let's check the oldest response:

```
>>> from scrapy.utils.trackref import get_oldest
>>> r = get_oldest('HtmlResponse')
>>> r.url
'http://www.somenastyspider.com/product.php?pid=123'
```

If you want to iterate over all objects, instead of getting the oldest one, you can use the `scrapy.utils.trackref.iter_all()` function:

```
>>> from scrapy.utils.trackref import iter_all
>>> [r.url for r in iter_all('HtmlResponse')]
['http://www.somenastyspider.com/product.php?pid=123',
 'http://www.somenastyspider.com/product.php?pid=584',
 ...]
```

Too many spiders?

If your project has too many spiders executed in parallel, the output of `prefs()` can be difficult to read. For this reason, that function has a `ignore` argument which can be used to ignore a particular class (and all its subclasses). For example, this won't show any live references to spiders:

```
>>> from scrapy.spiders import Spider
>>> prefs(ignore=Spider)
```

scrapy.utils.trackref module

Here are the functions available in the `trackref` module.

class scrapy.utils.trackref.object_ref

Inherit from this class (instead of object) if you want to track live instances with the `trackref` module.

scrapy.utils.trackref.print_live_refs(*class_name*, *ignore=NoneType*)

Print a report of live references, grouped by class name.

参数 *ignore* (*class or classes tuple*) – if given, all objects from the specified class (or tuple of classes) will be ignored.

scrapy.utils.trackref.get_oldest(*class_name*)

Return the oldest object alive with the given class name, or `None` if none is found. Use `print_live_refs()` first to get a list of all tracked live objects per class name.

scrapy.utils.trackref.iter_all(*class_name*)

Return an iterator over all objects alive with the given class name, or `None` if none is found. Use `print_live_refs()` first to get a list of all tracked live objects per class name.

5.7.3 Debugging memory leaks with Guppy

`trackref` provides a very convenient mechanism for tracking down memory leaks, but it only keeps track of the objects that are more likely to cause memory leaks (Requests, Responses, Items, and Selectors). However, there are other cases where the memory leaks could come from other (more or less obscure) objects. If this is your case, and you can't find your leaks using `trackref`, you still have another resource: the `Guppy` library. If you're using Python3, see *Debugging memory leaks with muppy*.

If you use `pip`, you can install Guppy with the following command:

```
pip install guppy
```

The telnet console also comes with a built-in shortcut (`hpy`) for accessing Guppy heap objects. Here's an example to view all Python objects available in the heap using Guppy:

```
>>> x = hpy.heap()
>>> x.bytype
```

Partition of a set of 297033 objects. Total size = 52587824 bytes.

Index	Count	%	Size	% Cumulative	% Type
0	22307	8	16423880	31	16423880 31 dict
1	122285	41	12441544	24	28865424 55 str
2	68346	23	5966696	11	34832120 66 tuple
3	227	0	5836528	11	40668648 77 unicode
4	2461	1	2222272	4	42890920 82 type
5	16870	6	2024400	4	44915320 85 function
6	13949	5	1673880	3	46589200 89 types.CodeType
7	13422	5	1653104	3	48242304 92 list
8	3735	1	1173680	2	49415984 94 _sre.SRE_Pattern
9	1209	0	456936	1	49872920 95 scrapy.http.headers.Headers

<1676 more rows. Type e.g. '_.more' to view.>

You can see that most space is used by dicts. Then, if you want to see from which attribute those dicts are referenced, you could do:

```
>>> x.bytype[0].byvia
```

Partition of a set of 22307 objects. Total size = 16423880 bytes.

Index	Count	%	Size	% Cumulative	% Referred Via:
0	10982	49	9416336	57	9416336 57 '.__dict__'
1	1820	8	2681504	16	12097840 74 '.__dict__', '.func_globals'
2	3097	14	1122904	7	13220744 80
3	990	4	277200	2	13497944 82 "['cookies']"
4	987	4	276360	2	13774304 84 "['cache']"
5	985	4	275800	2	14050104 86 "['meta']"
6	897	4	251160	2	14301264 87 '[2]'
7	1	0	196888	1	14498152 88 "['moduleDict']", "['modules']"
8	672	3	188160	1	14686312 89 "['cb_kwargs']"
9	27	0	155016	1	14841328 90 '[1]'

<333 more rows. Type e.g. '_.more' to view.>

As you can see, the Guppy module is very powerful but also requires some deep knowledge about Python internals. For more info about Guppy, refer to the [Guppy documentation](#).

5.7.4 Debugging memory leaks with muppy

If you're using Python 3, you can use muppy from [Pympler](#).

If you use `pip`, you can install muppy with the following command:

```
pip install Pympler
```

Here's an example to view all Python objects available in the heap using muppy:

```
>>> from pympler import muppy
>>> all_objects = muppy.get_objects()
>>> len(all_objects)
28667
>>> from pympler import summary
>>> sum1 = summary.summarize(all_objects)
>>> summary.print_(sum1)
```

types	# objects	total size
=====	=====	=====
<class 'str	9822	1.10 MB
<class 'dict	1658	856.62 KB
<class 'type	436	443.60 KB
<class 'code	2974	419.56 KB
<class '_io.BufferedWriter	2	256.34 KB
<class 'set	420	159.88 KB
<class '_io.BufferedReader	1	128.17 KB
<class 'wrapper_descriptor	1130	88.28 KB
<class 'tuple	1304	86.57 KB
<class 'weakref	1013	79.14 KB
<class 'builtin_function_or_method	958	67.36 KB
<class 'method_descriptor	865	60.82 KB
<class 'abc.ABCMeta	62	59.96 KB
<class 'list	446	58.52 KB
<class 'int	1425	43.20 KB

For more info about muppy, refer to the [muppy documentation](#).

5.7.5 Leaks without leaks

Sometimes, you may notice that the memory usage of your Scrapy process will only increase, but never decrease. Unfortunately, this could happen even though neither Scrapy nor your project are leaking memory. This is due to a (not so well) known problem of Python, which may not return released memory to the operating system in some cases. For more information on this issue see:

- [Python Memory Management](#)
- [Python Memory Management Part 2](#)
- [Python Memory Management Part 3](#)

The improvements proposed by Evan Jones, which are detailed in [this paper](#), got merged in Python 2.5, but this only reduces the problem, it doesn't fix it completely. To quote the paper:

Unfortunately, this patch can only free an arena if there are no more objects allocated in it anymore. This means that fragmentation is a large issue. An application could have many megabytes of free memory, scattered throughout all the arenas, but it will be unable to free any of it. This is a problem experienced by all memory allocators. The only way to solve it is to move to a compacting garbage collector, which is able to move objects in memory. This would require significant changes to the Python interpreter.

To keep memory consumption reasonable you can split the job into several smaller jobs or enable *persistent job queue* and stop/start spider from time to time.

5.8 下载并处理文件和图片

Scrapy provides reusable *item pipelines* for downloading files attached to a particular item (for example, when you scrape products and also want to download their images locally). These pipelines share a bit of functionality and structure (we refer to them as media pipelines), but typically you'll either use the Files Pipeline or the Images Pipeline.

Both pipelines implement these features:

- Avoid re-downloading media that was downloaded recently
- Specifying where to store the media (filesystem directory, Amazon S3 bucket, Google Cloud Storage bucket)

The Images Pipeline has a few extra functions for processing images:

- Convert all downloaded images to a common format (JPG) and mode (RGB)
- Thumbnail generation
- Check images width/height to make sure they meet a minimum constraint

The pipelines also keep an internal queue of those media URLs which are currently being scheduled for download, and connect those responses that arrive containing the same media to that queue. This avoids downloading the same media more than once when it's shared by several items.

5.8.1 Using the Files Pipeline

The typical workflow, when using the `FilesPipeline` goes like this:

1. In a Spider, you scrape an item and put the URLs of the desired into a `file_urls` field.
2. The item is returned from the spider and goes to the item pipeline.

3. When the item reaches the `FilesPipeline`, the URLs in the `file_urls` field are scheduled for download using the standard Scrapy scheduler and downloader (which means the scheduler and downloader middlewares are reused), but with a higher priority, processing them before other pages are scraped. The item remains “locked” at that particular pipeline stage until the files have finish downloading (or fail for some reason).
4. When the files are downloaded, another field (`files`) will be populated with the results. This field will contain a list of dicts with information about the downloaded files, such as the downloaded path, the original scraped url (taken from the `file_urls` field) , and the file checksum. The files in the list of the `files` field will retain the same order of the original `file_urls` field. If some file failed downloading, an error will be logged and the file won’ t be present in the `files` field.

5.8.2 Using the Images Pipeline

Using the *ImagesPipeline* is a lot like using the `FilesPipeline`, except the default field names used are different: you use `image_urls` for the image URLs of an item and it will populate an `images` field for the information about the downloaded images.

The advantage of using the *ImagesPipeline* for image files is that you can configure some extra functions like generating thumbnails and filtering the images based on their size.

The Images Pipeline uses *Pillow* for thumbnailing and normalizing images to JPEG/RGB format, so you need to install this library in order to use it. *Python Imaging Library* (PIL) should also work in most cases, but it is known to cause troubles in some setups, so we recommend to use *Pillow* instead of PIL.

5.8.3 Enabling your Media Pipeline

To enable your media pipeline you must first add it to your project *ITEM_PIPELINES* setting.

For Images Pipeline, use:

```
ITEM_PIPELINES = {'scrapy.pipelines.images.ImagesPipeline': 1}
```

For Files Pipeline, use:

```
ITEM_PIPELINES = {'scrapy.pipelines.files.FilesPipeline': 1}
```

注解: You can also use both the Files and Images Pipeline at the same time.

Then, configure the target storage setting to a valid value that will be used for storing the downloaded images. Otherwise the pipeline will remain disabled, even if you include it in the *ITEM_PIPELINES* setting.

For the Files Pipeline, set the *FILES_STORE* setting:

```
FILES_STORE = '/path/to/valid/dir'
```

For the Images Pipeline, set the *IMAGES_STORE* setting:

```
IMAGES_STORE = '/path/to/valid/dir'
```

5.8.4 Supported Storage

File system is currently the only officially supported storage, but there are also support for storing files in [Amazon S3](#) and [Google Cloud Storage](#).

File system storage

The files are stored using a [SHA1 hash](#) of their URLs for the file names.

For example, the following image URL:

```
http://www.example.com/image.jpg
```

Whose `SHA1 hash` is:

```
3afec3b4765f8f0a07b78f98c07b83f013567a0a
```

Will be downloaded and stored in the following file:

```
<IMAGES_STORE>/full/3afec3b4765f8f0a07b78f98c07b83f013567a0a.jpg
```

Where:

- `<IMAGES_STORE>` is the directory defined in *IMAGES_STORE* setting for the Images Pipeline.
- `full` is a sub-directory to separate full images from thumbnails (if used). For more info see *Thumbnail generation for images*.

Amazon S3 storage

FILES_STORE and *IMAGES_STORE* can represent an Amazon S3 bucket. Scrapy will automatically upload the files to the bucket.

For example, this is a valid *IMAGES_STORE* value:

```
IMAGES_STORE = 's3://bucket/images'
```


You can modify the Access Control List (ACL) policy used for the stored files, which is defined by the `FILES_STORE_S3_ACL` and `IMAGES_STORE_S3_ACL` settings. By default, the ACL is set to `private`. To make the files publicly available use the `public-read` policy:

```
IMAGES_STORE_S3_ACL = 'public-read'
```

For more information, see [canned ACLs](#) in the Amazon S3 Developer Guide.

Because Scrapy uses `boto` / `botocore` internally you can also use other S3-like storages. Storages like self-hosted [Minio](#) or [s3.scality](#). All you need to do is set endpoint option in you Scrapy settings:

```
AWS_ENDPOINT_URL = 'http://minio.example.com:9000'
```

For self-hosting you also might feel the need not to use SSL and not to verify SSL connection:

```
AWS_USE_SSL = False # or True (None by default)
AWS_VERIFY = False # or True (None by default)
```

Google Cloud Storage

`FILES_STORE` and `IMAGES_STORE` can represent a Google Cloud Storage bucket. Scrapy will automatically upload the files to the bucket. (requires [google-cloud-storage](#))

For example, these are valid `IMAGES_STORE` and `GCS_PROJECT_ID` settings:

```
IMAGES_STORE = 'gs://bucket/images/'
GCS_PROJECT_ID = 'project_id'
```

For information about authentication, see [this documentation](#).

You can modify the Access Control List (ACL) policy used for the stored files, which is defined by the `FILES_STORE_GCS_ACL` and `IMAGES_STORE_GCS_ACL` settings. By default, the ACL is set to `''` (empty string) which means that Cloud Storage applies the bucket's default object ACL to the object. To make the files publicly available use the `publicRead` policy:

```
IMAGES_STORE_GCS_ACL = 'publicRead'
```

For more information, see [Predefined ACLs](#) in the Google Cloud Platform Developer Guide.

5.8.5 Usage example

In order to use a media pipeline first, *enable it*.

Then, if a spider returns a dict with the URLs key (`file_urls` or `image_urls`, for the Files or Images Pipeline respectively), the pipeline will put the results under respective key (`files` or `images`).

If you prefer to use *Item*, then define a custom item with the necessary fields, like in this example for Images Pipeline:

```
import scrapy

class MyItem(scrapy.Item):

    # ... other item fields ...
    image_urls = scrapy.Field()
    images = scrapy.Field()
```

If you want to use another field name for the URLs key or for the results key, it is also possible to override it.

For the Files Pipeline, set *FILES_URLS_FIELD* and/or *FILES_RESULT_FIELD* settings:

```
FILES_URLS_FIELD = 'field_name_for_your_files_urls'
FILES_RESULT_FIELD = 'field_name_for_your_processed_files'
```

For the Images Pipeline, set *IMAGES_URLS_FIELD* and/or *IMAGES_RESULT_FIELD* settings:

```
IMAGES_URLS_FIELD = 'field_name_for_your_images_urls'
IMAGES_RESULT_FIELD = 'field_name_for_your_processed_images'
```

If you need something more complex and want to override the custom pipeline behaviour, see *Extending the Media Pipelines*.

If you have multiple image pipelines inheriting from ImagePipeline and you want to have different settings in different pipelines you can set setting keys preceded with uppercase name of your pipeline class. E.g. if your pipeline is called MyPipeline and you want to have custom IMAGES_URLS_FIELD you define setting MYPIPELINE_IMAGES_URLS_FIELD and your custom settings will be used.

5.8.6 Additional features

File expiration

The Image Pipeline avoids downloading files that were downloaded recently. To adjust this retention delay use the *FILES_EXPIRES* setting (or *IMAGES_EXPIRES*, in case of Images Pipeline), which specifies the delay in number of days:

```
# 120 days of delay for files expiration
FILES_EXPIRES = 120
```

(下页继续)

(续上页)

```
# 30 days of delay for images expiration
IMAGES_EXPIRES = 30
```

The default value for both settings is 90 days.

If you have pipeline that subclasses `FilesPipeline` and you'd like to have different setting for it you can set setting keys preceded by uppercase class name. E.g. given pipeline class called `MyPipeline` you can set setting key:

```
MYPIPELINE_FILES_EXPIRES = 180
```

and pipeline class `MyPipeline` will have expiration time set to 180.

Thumbnail generation for images

The Images Pipeline can automatically create thumbnails of the downloaded images. In order to use this feature, you must set `IMAGES_THUMBS` to a dictionary where the keys are the thumbnail names and the values are their dimensions.

For example:

```
IMAGES_THUMBS = {
    'small': (50, 50),
    'big': (270, 270),
}
```

When you use this feature, the Images Pipeline will create thumbnails of the each specified size with this format:

```
<IMAGES_STORE>/thumbs/<size_name>/<image_id>.jpg
```

Where:

- `<size_name>` is the one specified in the `IMAGES_THUMBS` dictionary keys (small, big, etc)
- `<image_id>` is the [SHA1 hash](#) of the image url

Example of image files stored using `small` and `big` thumbnail names:

```
<IMAGES_STORE>/full/63bbfea82b8880ed33cdb762aa11fab722a90a24.jpg
<IMAGES_STORE>/thumbs/small/63bbfea82b8880ed33cdb762aa11fab722a90a24.jpg
<IMAGES_STORE>/thumbs/big/63bbfea82b8880ed33cdb762aa11fab722a90a24.jpg
```

The first one is the full image, as downloaded from the site.

Filtering out small images

When using the Images Pipeline, you can drop images which are too small, by specifying the minimum allowed size in the `IMAGES_MIN_HEIGHT` and `IMAGES_MIN_WIDTH` settings.

For example:

```
IMAGES_MIN_HEIGHT = 110
IMAGES_MIN_WIDTH = 110
```

注解: The size constraints don't affect thumbnail generation at all.

It is possible to set just one size constraint or both. When setting both of them, only images that satisfy both minimum sizes will be saved. For the above example, images of sizes (105 x 105) or (105 x 200) or (200 x 105) will all be dropped because at least one dimension is shorter than the constraint.

By default, there are no size constraints, so all images are processed.

Allowing redirections

By default media pipelines ignore redirects, i.e. an HTTP redirection to a media file URL request will mean the media download is considered failed.

To handle media redirections, set this setting to `True`:

```
MEDIA_ALLOW_REDIRECTS = True
```

5.8.7 Extending the Media Pipelines

See here the methods that you can override in your custom Files Pipeline:

```
class scrapy.pipelines.files.FilesPipeline
```

```
    get_media_requests(item, info)
```

As seen on the workflow, the pipeline will get the URLs of the images to download from the item. In order to do this, you can override the `get_media_requests()` method and return a Request for each file URL:

```
def get_media_requests(self, item, info):
    for file_url in item['file_urls']:
        yield scrapy.Request(file_url)
```

Those requests will be processed by the pipeline and, when they have finished downloading, the results will be sent to the `item_completed()` method, as a list of 2-element tuples. Each tuple will contain (success, file_info_or_error) where:

- **success** is a boolean which is `True` if the image was downloaded successfully or `False` if it failed for some reason
- **file_info_or_error** is a dict containing the following keys (if success is `True`) or a `Twisted Failure` if there was a problem.
 - **url** - the url where the file was downloaded from. This is the url of the request returned from the `get_media_requests()` method.
 - **path** - the path (relative to `FILES_STORE`) where the file was stored
 - **checksum** - a `MD5` hash of the image contents

The list of tuples received by `item_completed()` is guaranteed to retain the same order of the requests returned from the `get_media_requests()` method.

Here's a typical value of the **results** argument:

```
[(True,
  {'checksum': '2b00042f7481c7b056c4b410d28f33cf',
   'path': 'full/0a79c461a4062ac383dc4fade7bc09f1384a3910.jpg',
   'url': 'http://www.example.com/files/product1.pdf'}),
 (False,
  Failure(...))]
```

By default the `get_media_requests()` method returns `None` which means there are no files to download for the item.

`item_completed(results, item, info)`

The `FilesPipeline.item_completed()` method called when all file requests for a single item have completed (either finished downloading, or failed for some reason).

The `item_completed()` method must return the output that will be sent to subsequent item pipeline stages, so you must return (or drop) the item, as you would in any pipeline.

Here is an example of the `item_completed()` method where we store the downloaded file paths (passed in results) in the `file_paths` item field, and we drop the item if it doesn't contain any files:

```
from scrapy.exceptions import DropItem

def item_completed(self, results, item, info):
    file_paths = [x['path'] for ok, x in results if ok]
    if not file_paths:
```

(下页继续)

(续上页)

```

        raise DropItem("Item contains no files")
    item['file_paths'] = file_paths
    return item

```

By default, the `item_completed()` method returns the item.

See here the methods that you can override in your custom Images Pipeline:

```
class scrapy.pipelines.images.ImagesPipeline
```

The *ImagesPipeline* is an extension of the *FilesPipeline*, customizing the field names and adding custom behavior for images.

```
get_media_requests(item, info)
```

Works the same way as *FilesPipeline.get_media_requests()* method, but using a different field name for image urls.

Must return a Request for each image URL.

```
item_completed(results, item, info)
```

The *ImagesPipeline.item_completed()* method is called when all image requests for a single item have completed (either finished downloading, or failed for some reason).

Works the same way as *FilesPipeline.item_completed()* method, but using a different field names for storing image downloading results.

By default, the `item_completed()` method returns the item.

5.8.8 Custom Images pipeline example

Here is a full example of the Images Pipeline whose methods are exemplified above:

```

import scrapy
from scrapy.pipelines.images import ImagesPipeline
from scrapy.exceptions import DropItem

class MyImagesPipeline(ImagesPipeline):

    def get_media_requests(self, item, info):
        for image_url in item['image_urls']:
            yield scrapy.Request(image_url)

    def item_completed(self, results, item, info):
        image_paths = [x['path'] for ok, x in results if ok]
        if not image_paths:

```

(下页继续)

(续上页)

```
raise DropItem("Item contains no images")
item['image_paths'] = image_paths
return item
```

5.9 部署爬虫器

This section describes the different options you have for deploying your Scrapy spiders to run them on a regular basis. Running Scrapy spiders in your local machine is very convenient for the (early) development stage, but not so much when you need to execute long-running spiders or move spiders to run in production continuously. This is where the solutions for deploying Scrapy spiders come in.

Popular choices for deploying Scrapy spiders are:

- *Scrapyd* (open source)
- *Scrapy Cloud* (cloud-based)

5.9.1 Deploying to a Scrapyd Server

Scrapyd is an open source application to run Scrapy spiders. It provides a server with HTTP API, capable of running and monitoring Scrapy spiders.

To deploy spiders to Scrapyd, you can use the *scrapyd-deploy* tool provided by the *scrapyd-client* package. Please refer to the *scrapyd-deploy documentation* for more information.

Scrapyd is maintained by some of the Scrapy developers.

5.9.2 Deploying to Scrapy Cloud

Scrapy Cloud is a hosted, cloud-based service by *Scrapinghub*, the company behind Scrapy.

Scrapy Cloud removes the need to setup and monitor servers and provides a nice UI to manage spiders and review scraped items, logs and stats.

To deploy spiders to Scrapy Cloud you can use the *shub* command line tool. Please refer to the *Scrapy Cloud documentation* for more information.

Scrapy Cloud is compatible with Scrapyd and one can switch between them as needed - the configuration is read from the *scrapy.cfg* file just like *scrapyd-deploy*.

5.10 爬虫器节流

This is an extension for automatically throttling crawling speed based on load of both the Scrapy server and the website you are crawling.

5.10.1 Design goals

1. be nicer to sites instead of using default download delay of zero
2. automatically adjust scrapy to the optimum crawling speed, so the user doesn't have to tune the download delays to find the optimum one. The user only needs to specify the maximum concurrent requests it allows, and the extension does the rest.

5.10.2 How it works

AutoThrottle extension adjusts download delays dynamically to make spider send `AUTOTHROTTLE_TARGET_CONCURRENCY` concurrent requests on average to each remote website.

It uses download latency to compute the delays. The main idea is the following: if a server needs `latency` seconds to respond, a client should send a request each `latency/N` seconds to have `N` requests processed in parallel.

Instead of adjusting the delays one can just set a small fixed download delay and impose hard limits on concurrency using `CONCURRENT_REQUESTS_PER_DOMAIN` or `CONCURRENT_REQUESTS_PER_IP` options. It will provide a similar effect, but there are some important differences:

- because the download delay is small there will be occasional bursts of requests;
- often non-200 (error) responses can be returned faster than regular responses, so with a small download delay and a hard concurrency limit crawler will be sending requests to server faster when server starts to return errors. But this is an opposite of what crawler should do - in case of errors it makes more sense to slow down: these errors may be caused by the high request rate.

AutoThrottle doesn't have these issues.

5.10.3 Throttling algorithm

AutoThrottle algorithm adjusts download delays based on the following rules:

1. spiders always start with a download delay of `AUTOTHROTTLE_START_DELAY`;
2. when a response is received, the target download delay is calculated as `latency / N` where `latency` is a latency of the response, and `N` is `AUTOTHROTTLE_TARGET_CONCURRENCY`.
3. download delay for next requests is set to the average of previous download delay and the target download delay;

4. latencies of non-200 responses are not allowed to decrease the delay;
5. download delay can't become less than `DOWNLOAD_DELAY` or greater than `AUTOTHROTTLE_MAX_DELAY`

注解: The AutoThrottle extension honours the standard Scrapy settings for concurrency and delay. This means that it will respect `CONCURRENT_REQUESTS_PER_DOMAIN` and `CONCURRENT_REQUESTS_PER_IP` options and never set a download delay lower than `DOWNLOAD_DELAY`.

In Scrapy, the download latency is measured as the time elapsed between establishing the TCP connection and receiving the HTTP headers.

Note that these latencies are very hard to measure accurately in a cooperative multitasking environment because Scrapy may be busy processing a spider callback, for example, and unable to attend downloads. However, these latencies should still give a reasonable estimate of how busy Scrapy (and ultimately, the server) is, and this extension builds on that premise.

5.10.4 Settings

The settings used to control the AutoThrottle extension are:

- `AUTOTHROTTLE_ENABLED`
- `AUTOTHROTTLE_START_DELAY`
- `AUTOTHROTTLE_MAX_DELAY`
- `AUTOTHROTTLE_TARGET_CONCURRENCY`
- `AUTOTHROTTLE_DEBUG`
- `CONCURRENT_REQUESTS_PER_DOMAIN`
- `CONCURRENT_REQUESTS_PER_IP`
- `DOWNLOAD_DELAY`

For more information see *How it works*.

AUTOTHROTTLE_ENABLED

Default: `False`

Enables the AutoThrottle extension.

AUTOTHROTTLE_START_DELAY

Default: `5.0`

The initial download delay (in seconds).

AUTOTHROTTLE_MAX_DELAY

Default: 60.0

The maximum download delay (in seconds) to be set in case of high latencies.

AUTOTHROTTLE_TARGET_CONCURRENCY

1.1 新版功能.

Default: 1.0

Average number of requests Scrapy should be sending in parallel to remote websites.

By default, AutoThrottle adjusts the delay to send a single concurrent request to each of the remote websites. Set this option to a higher value (e.g. 2.0) to increase the throughput and the load on remote servers. A lower `AUTOTHROTTLE_TARGET_CONCURRENCY` value (e.g. 0.5) makes the crawler more conservative and polite.

Note that `CONCURRENT_REQUESTS_PER_DOMAIN` and `CONCURRENT_REQUESTS_PER_IP` options are still respected when AutoThrottle extension is enabled. This means that if `AUTOTHROTTLE_TARGET_CONCURRENCY` is set to a value higher than `CONCURRENT_REQUESTS_PER_DOMAIN` or `CONCURRENT_REQUESTS_PER_IP`, the crawler won't reach this number of concurrent requests.

At every given time point Scrapy can be sending more or less concurrent requests than `AUTOTHROTTLE_TARGET_CONCURRENCY`; it is a suggested value the crawler tries to approach, not a hard limit.

AUTOTHROTTLE_DEBUG

Default: `False`

Enable AutoThrottle debug mode which will display stats on every response received, so you can see how the throttling parameters are being adjusted in real time.

5.11 爬虫器硬件性能

0.17 新版功能.

Scrapy comes with a simple benchmarking suite that spawns a local HTTP server and crawls it at the maximum possible speed. The goal of this benchmarking is to get an idea of how Scrapy performs in your hardware, in order to have a common baseline for comparisons. It uses a simple spider that does nothing and just follows links.

To run it use:

```
scrapy bench
```

You should see an output like this:

```
2016-12-16 21:18:48 [scrapy.utils.log] INFO: Scrapy 1.2.2 started (bot: quotesbot)
2016-12-16 21:18:48 [scrapy.utils.log] INFO: Overridden settings: {'CLOSESPIDER_TIMEOUT': 10, 'ROBOTSTXT_OBEY': True, 'SPIDER_MODULES': ['quotesbot.spiders'], 'LOGSTATS_INTERVAL': 1, 'BOT_NAME': 'quotesbot', 'LOG_LEVEL': 'INFO', 'NEWSPIDER_MODULE': 'quotesbot.spiders'}
2016-12-16 21:18:49 [scrapy.middleware] INFO: Enabled extensions:
['scrapy.extensions.closespider.CloseSpider',
 'scrapy.extensions.logstats.LogStats',
 'scrapy.extensions.telnet.TelnetConsole',
 'scrapy.extensions.corestats.CoreStats']
2016-12-16 21:18:49 [scrapy.middleware] INFO: Enabled downloader middlewares:
['scrapy.downloadermiddlewares.robotstxt.RobotsTxtMiddleware',
 'scrapy.downloadermiddlewares.httpauth.HttpAuthMiddleware',
 'scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware',
 'scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware',
 'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware',
 'scrapy.downloadermiddlewares.retry.RetryMiddleware',
 'scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware',
 'scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware',
 'scrapy.downloadermiddlewares.redirect.RedirectMiddleware',
 'scrapy.downloadermiddlewares.cookies.CookiesMiddleware',
 'scrapy.downloadermiddlewares.stats.DownloaderStats']
2016-12-16 21:18:49 [scrapy.middleware] INFO: Enabled spider middlewares:
['scrapy.spidermiddlewares.httperror.HttpErrorMiddleware',
 'scrapy.spidermiddlewares.offsite.OffsiteMiddleware',
 'scrapy.spidermiddlewares.referer.RefererMiddleware',
 'scrapy.spidermiddlewares.urllength.UrlLengthMiddleware',
 'scrapy.spidermiddlewares.depth.DepthMiddleware']
2016-12-16 21:18:49 [scrapy.middleware] INFO: Enabled item pipelines:
[]
2016-12-16 21:18:49 [scrapy.core.engine] INFO: Spider opened
2016-12-16 21:18:49 [scrapy.extensions.logstats] INFO: Crawled 0 pages (at 0 pages/min),
↳scraped 0 items (at 0 items/min)
2016-12-16 21:18:50 [scrapy.extensions.logstats] INFO: Crawled 70 pages (at 4200 pages/
↳min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:51 [scrapy.extensions.logstats] INFO: Crawled 134 pages (at 3840 pages/
↳min), scraped 0 items (at 0 items/min)
```

(下页继续)

(续上页)

```
2016-12-16 21:18:52 [scrapy.extensions.logstats] INFO: Crawled 198 pages (at 3840 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:53 [scrapy.extensions.logstats] INFO: Crawled 254 pages (at 3360 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:54 [scrapy.extensions.logstats] INFO: Crawled 302 pages (at 2880 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:55 [scrapy.extensions.logstats] INFO: Crawled 358 pages (at 3360 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:56 [scrapy.extensions.logstats] INFO: Crawled 406 pages (at 2880 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:57 [scrapy.extensions.logstats] INFO: Crawled 438 pages (at 1920 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:58 [scrapy.extensions.logstats] INFO: Crawled 470 pages (at 1920 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:18:59 [scrapy.core.engine] INFO: Closing spider (closespider_timeout)
2016-12-16 21:18:59 [scrapy.extensions.logstats] INFO: Crawled 518 pages (at 2880 pages/
↪min), scraped 0 items (at 0 items/min)
2016-12-16 21:19:00 [scrapy.statscollectors] INFO: Dumping Scrapy stats:
{'downloader/request_bytes': 229995,
 'downloader/request_count': 534,
 'downloader/request_method_count/GET': 534,
 'downloader/response_bytes': 1565504,
 'downloader/response_count': 534,
 'downloader/response_status_count/200': 534,
 'finish_reason': 'closespider_timeout',
 'finish_time': datetime.datetime(2016, 12, 16, 16, 19, 0, 647725),
 'log_count/INFO': 17,
 'request_depth_max': 19,
 'response_received_count': 534,
 'scheduler/dequeued': 533,
 'scheduler/dequeued/memory': 533,
 'scheduler/enqueued': 10661,
 'scheduler/enqueued/memory': 10661,
 'start_time': datetime.datetime(2016, 12, 16, 16, 18, 49, 799869)}
2016-12-16 21:19:00 [scrapy.core.engine] INFO: Spider closed (closespider_timeout)
```

That tells you that Scrapy is able to crawl about 3000 pages per minute in the hardware where you run it. Note that this is a very simple spider intended to follow links, any custom spider you write will probably do more stuff which results in slower crawl rates. How slower depends on how much your spider does and how well it's written.

In the future, more cases will be added to the benchmarking suite to cover other common scenarios.

5.12 作业: 暂停并恢复爬行

有时, 对于一些大型网站, 需要中途暂停并在某个时间继续爬取。

使用以下工具, Scrapy 可以实现这种开箱即用的这种功能,

- 将请求列表保存到磁盘的一种调度程序
- 在磁盘上保存已经请求过的一种过滤器
- 一个能持续保持爬虫状态 (键/值对) 的扩展

5.12.1 Job 路径

要启动持久化功能你只需通过 `JOBDIR` 设置一个 *job directory* 选项。这个路径会存储所有的请求数据来保持一个单独任务的状态 (如: 一次爬虫)。需要注意的是, 这个路径不能被不同的爬虫之共享, 或者一个爬虫中的不同的 `jobs/runs`。这意味着它只能保存一个 单独任务的状态。

5.12.2 如何使用

启用持久功能来开始一个爬虫, 按照以下方式运行:

```
scrapy crawl somespider -s JOBDIR=crawls/somespider-1
```

然后, 你可以在任何时间使用安全方式停止爬虫 (使用 `Ctrl-C` 或者发送一个信号), 并可以使用相同的命令来继续爬虫:

```
scrapy crawl somespider -s JOBDIR=crawls/somespider-1
```

5.12.3 保持状态

有时你想保持一些爬虫的状态。你可以使用 `spider.state` 属性来实现, 这个属性必需是一个字典。Scrapy 提供了一个内置扩展负责在爬虫开始和停止的时候从 `job` 路径序列化、存储并加载这个属性。

这里有一个例子使用了爬虫状态来进行回调 (简洁起见省略了其它爬虫代码):

```
def parse_item(self, response):  
    # parse item here  
    self.state['items_count'] = self.state.get('items_count', 0) + 1
```

5.12.4 持久化的一些事项

如果你想启用 Scrapy 持久化，以下几条你需要注意：

Cookies 有效期

Cookies 会过期。因此，如果你没有尽快恢复爬虫，请求调度器可能不再工作。如果你的爬虫不依赖 cookies 这也不是问题。

请求序列化

请求需要被 pickle 模块序列化，你需要确保你的请求可以被序列化以确保可以持续运行。

最常见的问题是在 requests 回调函数中使用 lambda 函数，导致无法被序列化。

以下操作将无法工作：

```
def some_callback(self, response):
    somearg = 'test'
    return scrapy.Request('http://www.example.com', callback=lambda r: self.other_
↪callback(r, somearg))

def other_callback(self, response, somearg):
    print("the argument passed is: %s" % somearg)
```

但是可以这样：

```
def some_callback(self, response):
    somearg = 'test'
    return scrapy.Request('http://www.example.com', callback=self.other_callback, meta={
↪'somearg': somearg})

def other_callback(self, response):
    somearg = response.meta['somearg']
    print("the argument passed is: %s" % somearg)
```

如果你想记录这些没有被序列化的请求，你可以在项目中设置 `SCHEDULER_DEBUG` 选项为 `True`。它默认是 `False`。

常见问题 获取最常见问题的答案。

Debug 爬虫器 了解如何 Debug 调试你的 Scrapy 爬虫常见问题。

爬虫器约束 了解如何使用约束条件来测试你的爬虫爬虫器。

常见实践 熟悉一些 Scrapy 常见的实践案例。

并发爬虫 优化 Scrapy 去并行爬取大量的域名。

使用浏览器的开发工具进行抓取 学习如何使用浏览器的开发工具抓取。

内存泄漏调试 学习查找和删除爬虫器中的内存泄漏。

下载并处理文件和图片 从抓取到数据中下载你在 item 中定义过的文件和图片。

部署爬虫器 部署你的 Scrapy 爬虫器并在远程服务器上运行它们。

爬虫器节流 根据负载动态调整爬虫速度。

爬虫器硬件性能 检查一下 Scrapy 在硬件上的性能。

作业：暂停并恢复爬行 学习如何暂停并继续大型的爬虫器。

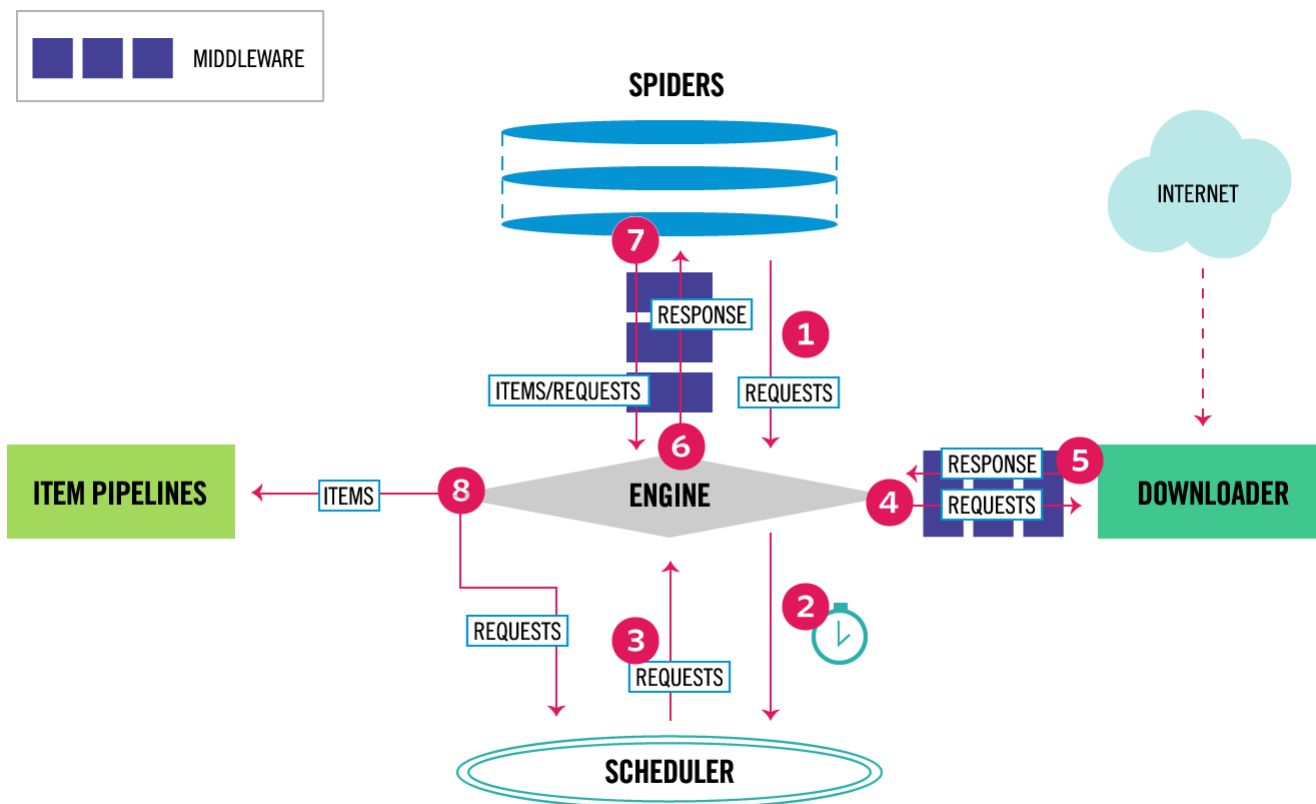
6.1 框架体系

本文档描述了 Scrapy 框架的各部分之间是如何相互联系的。

6.1.1 综述

图示描述了 Scrapy 框架和它各组件，以及内部数据在系统内如何活动（通过红色箭头表示）。下面是对这些组件的描述，以及它们更详细的信息的链接。同时也有关于数据流的描述。

6.1.2 Data flow



Scrapy 中的数据流是由执行引擎控制的，如下所示：

1. 引擎 从爬虫器 中获取爬虫的初始请求对象（一个或者多个）。
2. 引擎 从调度器 中调度请求对象并获取下一个要爬取的请求对象。
3. 调度器 返回下一个请求对象给引擎。
4. 引擎 通过下载器中间件（见`process_request()`）将请求对象传递给下载器。
5. 一旦页面完成下载，下载器 将生成一个响应对象（该页面的）然后通过下载器中间件 传递给引擎（见`process_response()`）。
6. 引擎 从下载器 获得响应对象后把它传递给爬虫器 使用爬虫器中间件 处理（见`process_spider_input()`）。
7. 爬虫器 使用爬虫器中间件 处理响应对象然后返回 `items` 和新的请求对象给引擎（见`process_spider_output()`）。
8. 引擎，然后将处理过的请求对象传递给调度器 并获得后续要爬取的请求对象。
9. 这个过程（从第一步）一直重复直到调度器 中没有请求。

6.1.3 组件

Scrapy 引擎

引擎负责控制系统所有组件之间的数据流，当某些行为发生时触发事件。见[数据流](#)了解更多细节。

调度器

调度器从引擎中接收请求并给它们排序，以便于稍后引擎请求它们时将它们传给引擎。

下载器

下载器负责获取 web 页面并将它们传递给引擎，而引擎又将它们传递给爬虫器。

爬虫器

爬虫器是用户编写的爬虫类，用于解析响应并提取 items 或后续其他的请求，更多细节详见[爬虫器](#)。

Item 管道

Item 管道负责处理爬虫器提取的 items，典型的任务包括清理，验证和持久化（比如把 item 保存到数据库）。更多细节详见[Item 管道](#)。

下载器中间件

下载器中间件是位于引擎和下载器之间的钩子，它处理从引擎到下载器的请求和响应。

如果你要做以下的事情，请用下载器中间件：

- 在请求到下载器之前处理它；（如：请求发送到网站之前）；
- 在响应传递到爬虫器之前改变它；
- 发送一个新的请求，而不是将接收到的响应传递给爬虫器；
- 在不获取网页的情况下将响应传递给爬虫器；
- 默默地减少一些请求。

更多细节详见[下载器中间件](#)。

爬虫器中间件

爬虫器中间件是位于引擎和爬虫器之间的钩子，能够处理爬虫器的输入（响应）和输出（items 和请求）。

如果你要做以下的事情，请用爬虫器中间件：

- post-process output of spider callbacks - change/add/remove requests or items;
- post-process start_requests;
- handle spider exceptions;
- call errback instead of callback for some of the requests based on response content.

更多细节详见[爬虫器中间件](#)。

6.1.4 事件驱动网络

Scrapy 是用 [Twisted](#) 写的, Twisted 是一个著名的 Python 事件驱动网络框架。因此, 它用非阻塞 (即异步) 方式实现并发。

更多关于异步编程和 Twisted 的信息请参考以下链接:

- [Introduction to Deferreds in Twisted](#)
- [Twisted - hello, asynchronous programming](#)
- [Twisted Introduction - Krondo](#)

6.2 下载器中间件

The downloader middleware is a framework of hooks into Scrapy's request/response processing. It's a light, low-level system for globally altering Scrapy's requests and responses.

6.2.1 Activating a downloader middleware

To activate a downloader middleware component, add it to the `DOWNLOADER_MIDDLEWARES` setting, which is a dict whose keys are the middleware class paths and their values are the middleware orders.

Here's an example:

```
DOWNLOADER_MIDDLEWARES = {
    'myproject.middlewares.CustomDownloaderMiddleware': 543,
}
```

The `DOWNLOADER_MIDDLEWARES` setting is merged with the `DOWNLOADER_MIDDLEWARES_BASE` setting defined in Scrapy (and not meant to be overridden) and then sorted by order to get the final sorted list of enabled middlewares: the first middleware is the one closer to the engine and the last is the one closer to the downloader. In other words, the `process_request()` method of each middleware will be invoked in increasing middleware order (100, 200, 300, ...) and the `process_response()` method of each middleware will be invoked in decreasing order.

To decide which order to assign to your middleware see the [DOWNLOADER_MIDDLEWARES_BASE](#) setting and pick a value according to where you want to insert the middleware. The order does matter because each middleware performs a different action and your middleware could depend on some previous (or subsequent) middleware being applied.

If you want to disable a built-in middleware (the ones defined in [DOWNLOADER_MIDDLEWARES_BASE](#) and enabled by default) you must define it in your project's [DOWNLOADER_MIDDLEWARES](#) setting and assign `None` as its value. For example, if you want to disable the user-agent middleware:

```
DOWNLOADER_MIDDLEWARES = {
    'myproject.middlewares.CustomDownloaderMiddleware': 543,
    'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware': None,
}
```

Finally, keep in mind that some middlewares may need to be enabled through a particular setting. See each middleware documentation for more info.

6.2.2 Writing your own downloader middleware

Each downloader middleware is a Python class that defines one or more of the methods defined below.

The main entry point is the `from_crawler` class method, which receives a [Crawler](#) instance. The [Crawler](#) object gives you access, for example, to the [settings](#).

```
class scrapy.downloadermiddlewares.DownloaderMiddleware
```

注解: Any of the downloader middleware methods may also return a deferred.

`process_request(request, spider)`

This method is called for each request that goes through the download middleware.

[process_request\(\)](#) should either: return `None`, return a [Response](#) object, return a [Request](#) object, or raise [IgnoreRequest](#).

If it returns `None`, Scrapy will continue processing this request, executing all other middlewares until, finally, the appropriate downloader handler is called the request performed (and its response downloaded).

If it returns a [Response](#) object, Scrapy won't bother calling *any* other [process_request\(\)](#) or [process_exception\(\)](#) methods, or the appropriate download function; it'll return that response. The [process_response\(\)](#) methods of installed middleware is always called on every response.

If it returns a [Request](#) object, Scrapy will stop calling `process_request` methods and reschedule the returned request. Once the newly returned request is performed, the appropriate middleware

chain will be called on the downloaded response.

If it raises an *IgnoreRequest* exception, the *process_exception()* methods of installed downloader middleware will be called. If none of them handle the exception, the errback function of the request (*Request.errback*) is called. If no code handles the raised exception, it is ignored and not logged (unlike other exceptions).

参数

- **request** (*Request* object) – the request being processed
- **spider** (*Spider* object) – the spider for which this request is intended

process_response(request, response, spider)

process_response() should either: return a *Response* object, return a *Request* object or raise a *IgnoreRequest* exception.

If it returns a *Response* (it could be the same given response, or a brand-new one), that response will continue to be processed with the *process_response()* of the next middleware in the chain.

If it returns a *Request* object, the middleware chain is halted and the returned request is rescheduled to be downloaded in the future. This is the same behavior as if a request is returned from *process_request()*.

If it raises an *IgnoreRequest* exception, the errback function of the request (*Request.errback*) is called. If no code handles the raised exception, it is ignored and not logged (unlike other exceptions).

参数

- **request** (is a *Request* object) – the request that originated the response
- **response** (*Response* object) – the response being processed
- **spider** (*Spider* object) – the spider for which this response is intended

process_exception(request, exception, spider)

Scrapy calls *process_exception()* when a download handler or a *process_request()* (from a downloader middleware) raises an exception (including an *IgnoreRequest* exception)

process_exception() should return: either *None*, a *Response* object, or a *Request* object.

If it returns *None*, Scrapy will continue processing this exception, executing any other *process_exception()* methods of installed middleware, until no middleware is left and the default exception handling kicks in.

If it returns a *Response* object, the *process_response()* method chain of installed middleware is started, and Scrapy won't bother calling any other *process_exception()* methods of middleware.

If it returns a *Request* object, the returned request is rescheduled to be downloaded in the future. This stops the execution of *process_exception()* methods of the middleware the same as returning a response would.

参数

- **request** (is a *Request* object) – the request that generated the exception
- **exception** (an *Exception* object) – the raised exception
- **spider** (*Spider* object) – the spider for which this request is intended

`from_crawler(cls, crawler)`

If present, this classmethod is called to create a middleware instance from a *Crawler*. It must return a new instance of the middleware. Crawler object provides access to all Scrapy core components like settings and signals; it is a way for middleware to access them and hook its functionality into Scrapy.

参数 `crawler` (*Crawler* object) – crawler that uses this middleware

6.2.3 Built-in downloader middleware reference

This page describes all downloader middleware components that come with Scrapy. For information on how to use them and how to write your own downloader middleware, see the *downloader middleware usage guide*.

For a list of the components enabled by default (and their orders) see the *DOWNLOADER_MIDDLEWARES_BASE* setting.

CookiesMiddleware

`class scrapy.downloadermiddlewares.cookies.CookiesMiddleware`

This middleware enables working with sites that require cookies, such as those that use sessions. It keeps track of cookies sent by web servers, and send them back on subsequent requests (from that spider), just like web browsers do.

The following settings can be used to configure the cookie middleware:

- *COOKIES_ENABLED*
- *COOKIES_DEBUG*

Multiple cookie sessions per spider

0.15 新版功能.

There is support for keeping multiple cookie sessions per spider by using the *cookiejar* Request meta key. By default it uses a single cookie jar (session), but you can pass an identifier to use different ones.

For example:

```
for i, url in enumerate(urls):
    yield scrapy.Request(url, meta={'cookiejar': i},
        callback=self.parse_page)
```

Keep in mind that the `cookiejar` meta key is not “sticky”. You need to keep passing it along on subsequent requests. For example:

```
def parse_page(self, response):
    # do some processing
    return scrapy.Request("http://www.example.com/otherpage",
        meta={'cookiejar': response.meta['cookiejar']},
        callback=self.parse_other_page)
```

COOKIES_ENABLED

Default: True

Whether to enable the cookies middleware. If disabled, no cookies will be sent to web servers.

Notice that despite the value of `COOKIES_ENABLED` setting if `Request.meta['dont_merge_cookies']` evaluates to `True` the request cookies will **not** be sent to the web server and received cookies in `Response` will **not** be merged with the existing cookies.

For more detailed information see the `cookies` parameter in `Request`.

COOKIES_DEBUG

Default: False

If enabled, Scrapy will log all cookies sent in requests (ie. `Cookie` header) and all cookies received in responses (ie. `Set-Cookie` header).

Here's an example of a log with `COOKIES_DEBUG` enabled:

```
2011-04-06 14:35:10-0300 [scrapy.core.engine] INFO: Spider opened
2011-04-06 14:35:10-0300 [scrapy.downloadermiddlewares.cookies] DEBUG: Sending cookies
↳to: <GET http://www.diningcity.com/netherlands/index.html>
    Cookie: clientlanguage_nl=en_EN
2011-04-06 14:35:14-0300 [scrapy.downloadermiddlewares.cookies] DEBUG: Received cookies
↳from: <200 http://www.diningcity.com/netherlands/index.html>
    Set-Cookie: JSESSIONID=B~FA4DC0C496C8762AE4F1A620EAB34F38; Path=/
    Set-Cookie: ip_isocode=US
    Set-Cookie: clientlanguage_nl=en_EN; Expires=Thu, 07-Apr-2011 21:21:34 GMT;
↳Path=/
```

(下页继续)

(续上页)

```
2011-04-06 14:49:50-0300 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://www.
↪diningcity.com/netherlands/index.html> (referer: None)
[...]
```

DefaultHeadersMiddleware

```
class scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware
```

This middleware sets all default requests headers specified in the `DEFAULT_REQUEST_HEADERS` setting.

DownloadTimeoutMiddleware

```
class scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware
```

This middleware sets the download timeout for requests specified in the `DOWNLOAD_TIMEOUT` setting or `download_timeout` spider attribute.

注解: You can also set download timeout per-request using `download_timeout` Request.meta key; this is supported even when DownloadTimeoutMiddleware is disabled.

HttpAuthMiddleware

```
class scrapy.downloadermiddlewares.httppauth.HttpAuthMiddleware
```

This middleware authenticates all requests generated from certain spiders using Basic access authentication (aka. HTTP auth).

To enable HTTP authentication from certain spiders, set the `http_user` and `http_pass` attributes of those spiders.

Example:

```
from scrapy.spiders import CrawlSpider

class SomeIntranetSiteSpider(CrawlSpider):

    http_user = 'someuser'
    http_pass = 'somepass'
    name = 'intranet.example.com'

    # .. rest of the spider code omitted ...
```

HttpCacheMiddleware

`class scrapy.downloadermiddlewares.httppcache.HttpCacheMiddleware`

This middleware provides low-level cache to all HTTP requests and responses. It has to be combined with a cache storage backend as well as a cache policy.

Scrapy ships with three HTTP cache storage backends:

- *Filesystem storage backend (default)*
- *DBM storage backend*
- *LevelDB storage backend*

You can change the HTTP cache storage backend with the `HTTPCACHE_STORAGE` setting. Or you can also implement your own storage backend.

Scrapy ships with two HTTP cache policies:

- *RFC2616 policy*
- *Dummy policy (default)*

You can change the HTTP cache policy with the `HTTPCACHE_POLICY` setting. Or you can also implement your own policy. You can also avoid caching a response on every policy using `dont_cache` meta key equals `True`.

Dummy policy (default)

This policy has no awareness of any HTTP Cache-Control directives. Every request and its corresponding response are cached. When the same request is seen again, the response is returned without transferring anything from the Internet.

The Dummy policy is useful for testing spiders faster (without having to wait for downloads every time) and for trying your spider offline, when an Internet connection is not available. The goal is to be able to “replay” a spider run *exactly as it ran before*.

In order to use this policy, set:

- `HTTPCACHE_POLICY` to `scrapy.extensions.httppcache.DummyPolicy`

RFC2616 policy

This policy provides a RFC2616 compliant HTTP cache, i.e. with HTTP Cache-Control awareness, aimed at production and used in continuous runs to avoid downloading unmodified data (to save bandwidth and speed up crawls).

what is implemented:

- Do not attempt to store responses/requests with `no-store` cache-control directive set

- Do not serve responses from cache if `no-cache` cache-control directive is set even for fresh responses
- Compute freshness lifetime from `max-age` cache-control directive
- Compute freshness lifetime from `Expires` response header
- Compute freshness lifetime from `Last-Modified` response header (heuristic used by Firefox)
- Compute current age from `Age` response header
- Compute current age from `Date` header
- Revalidate stale responses based on `Last-Modified` response header
- Revalidate stale responses based on `ETag` response header
- Set `Date` header for any received response missing it
- Support `max-stale` cache-control directive in requests

This allows spiders to be configured with the full RFC2616 cache policy, but avoid revalidation on a request-by-request basis, while remaining conformant with the HTTP spec.

Example:

Add `Cache-Control: max-stale=600` to Request headers to accept responses that have exceeded their expiration time by no more than 600 seconds.

See also: RFC2616, 14.9.3

what is missing:

- **Pragma:** `no-cache` support <https://www.w3.org/Protocols/rfc2616/rfc2616-sec14.html#sec14.9.1>
- **Vary** header support <https://www.w3.org/Protocols/rfc2616/rfc2616-sec13.html#sec13.6>
- Invalidation after updates or deletes <https://www.w3.org/Protocols/rfc2616/rfc2616-sec13.html#sec13.10>
- ...probably others ..

In order to use this policy, set:

- `HTTPCACHE_POLICY` to `scrapy.extensions.httpcache.RFC2616Policy`

Filesystem storage backend (default)

File system storage backend is available for the HTTP cache middleware.

In order to use this storage backend, set:

- `HTTPCACHE_STORAGE` to `scrapy.extensions.httpcache.FilesystemCacheStorage`

Each request/response pair is stored in a different directory containing the following files:

- `request_body` - the plain request body

- `request_headers` - the request headers (in raw HTTP format)
- `response_body` - the plain response body
- `response_headers` - the request headers (in raw HTTP format)
- `meta` - some metadata of this cache resource in Python `repr()` format (grep-friendly format)
- `pickled_meta` - the same metadata in `meta` but pickled for more efficient deserialization

The directory name is made from the request fingerprint (see `scrapy.utils.request.fingerprint`), and one level of subdirectories is used to avoid creating too many files into the same directory (which is inefficient in many file systems). An example directory could be:

```
/path/to/cache/dir/example.com/72/72811f648e718090f041317756c03adb0ada46c7
```

DBM storage backend

0.13 新版功能.

A `DBM` storage backend is also available for the HTTP cache middleware.

By default, it uses the `anydbm` module, but you can change it with the `HTTPCACHE_DBM_MODULE` setting.

In order to use this storage backend, set:

- `HTTPCACHE_STORAGE` to `scrapy.extensions.httpcache.DbmCacheStorage`

LevelDB storage backend

0.23 新版功能.

A `LevelDB` storage backend is also available for the HTTP cache middleware.

This backend is not recommended for development because only one process can access `LevelDB` databases at the same time, so you can't run a crawl and open the scrapy shell in parallel for the same spider.

In order to use this storage backend:

- set `HTTPCACHE_STORAGE` to `scrapy.extensions.httpcache.LevelDbCacheStorage`
- install `LevelDB` python bindings like `pip install leveldb`

HTTPCache middleware settings

The `HttpCacheMiddleware` can be configured through the following settings:

HTTPCACHE_ENABLED

0.11 新版功能.

Default: `False`

Whether the HTTP cache will be enabled.

在 0.11 版更改: Before 0.11, `HTTPCACHE_DIR` was used to enable cache.

HTTPCACHE_EXPIRATION_SECS

Default: `0`

Expiration time for cached requests, in seconds.

Cached requests older than this time will be re-downloaded. If zero, cached requests will never expire.

在 0.11 版更改: Before 0.11, zero meant cached requests always expire.

HTTPCACHE_DIR

Default: `'httpcache'`

The directory to use for storing the (low-level) HTTP cache. If empty, the HTTP cache will be disabled. If a relative path is given, is taken relative to the project data dir. For more info see: *Default structure of Scrapy projects*.

HTTPCACHE_IGNORE_HTTP_CODES

0.10 新版功能.

Default: `[]`

Don't cache response with these HTTP codes.

HTTPCACHE_IGNORE_MISSING

Default: `False`

If enabled, requests not found in the cache will be ignored instead of downloaded.

HTTPCACHE_IGNORE_SCHEMES

0.10 新版功能.

Default: `['file']`

Don't cache responses with these URI schemes.

HTTPCACHE_STORAGE

Default: `'scrapy.extensions.httppcache.FilesystemCacheStorage'`

The class which implements the cache storage backend.

HTTPCACHE_DBM_MODULE

0.13 新版功能.

Default: `'anydbm'`

The database module to use in the *DBM storage backend*. This setting is specific to the DBM backend.

HTTPCACHE_POLICY

0.18 新版功能.

Default: `'scrapy.extensions.httppcache.DummyPolicy'`

The class which implements the cache policy.

HTTPCACHE_GZIP

1.0 新版功能.

Default: `False`

If enabled, will compress all cached data with gzip. This setting is specific to the Filesystem backend.

HTTPCACHE_ALWAYS_STORE

1.1 新版功能.

Default: `False`

If enabled, will cache pages unconditionally.

A spider may wish to have all responses available in the cache, for future use with `Cache-Control: max-stale`, for instance. The `DummyPolicy` caches all responses but never revalidates them, and sometimes a more nuanced policy is desirable.

This setting still respects `Cache-Control: no-store` directives in responses. If you don't want that, filter `no-store` out of the Cache-Control headers in responses you feed to the cache middleware.

HTTPCACHE_IGNORE_RESPONSE_CACHE_CONTROLS

1.1 新版功能.

Default: []

List of Cache-Control directives in responses to be ignored.

Sites often set “no-store” , “no-cache” , “must-revalidate” , etc., but get upset at the traffic a spider can generate if it respects those directives. This allows to selectively ignore Cache-Control directives that are known to be unimportant for the sites being crawled.

We assume that the spider will not issue Cache-Control directives in requests unless it actually needs them, so directives in requests are not filtered.

HttpCompressionMiddleware

```
class scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware
```

This middleware allows compressed (gzip, deflate) traffic to be sent/received from web sites.

This middleware also supports decoding [brotli-compressed](#) responses, provided [brotlipy](#) is installed.

HttpCompressionMiddleware Settings

COMPRESSION_ENABLED

Default: True

Whether the Compression middleware will be enabled.

HttpProxyMiddleware

0.8 新版功能.

```
class scrapy.downloadermiddlewares.httpproxy.HttpProxyMiddleware
```

This middleware sets the HTTP proxy to use for requests, by setting the `proxy` meta value for *Request* objects.

Like the Python standard library modules [urllib](#) and [urllib2](#), it obeys the following environment variables:

- `http_proxy`
- `https_proxy`
- `no_proxy`

You can also set the meta key `proxy` per-request, to a value like `http://some_proxy_server:port` or `http://username:password@some_proxy_server:port`. Keep in mind this value will take precedence over `http_proxy/https_proxy` environment variables, and it will also ignore `no_proxy` environment variable.

RedirectMiddleware

```
class scrapy.downloadermiddlewares.redirect.RedirectMiddleware
```

This middleware handles redirection of requests based on response status.

The urls which the request goes through (while being redirected) can be found in the `redirect_urls` *Request.meta* key. The reason behind each redirect in `redirect_urls` can be found in the `redirect_reasons` *Request.meta* key. For example: `[301, 302, 307, 'meta refresh']`.

The format of a reason depends on the middleware that handled the corresponding redirect. For example, *RedirectMiddleware* indicates the triggering response status code as an integer, while *MetaRefreshMiddleware* always uses the 'meta refresh' string as reason.

The *RedirectMiddleware* can be configured through the following settings (see the settings documentation for more info):

- `REDIRECT_ENABLED`
- `REDIRECT_MAX_TIMES`

If *Request.meta* has `dont_redirect` key set to `True`, the request will be ignored by this middleware.

If you want to handle some redirect status codes in your spider, you can specify these in the `handle_httpstatus_list` spider attribute.

For example, if you want the redirect middleware to ignore 301 and 302 responses (and pass them through to your spider) you can do this:

```
class MySpider(CrawlSpider):
    handle_httpstatus_list = [301, 302]
```

The `handle_httpstatus_list` key of *Request.meta* can also be used to specify which response codes to allow on a per-request basis. You can also set the meta key `handle_httpstatus_all` to `True` if you want to allow any response code for a request.

RedirectMiddleware settings

REDIRECT_ENABLED

0.13 新版功能.

Default: `True`

Whether the Redirect middleware will be enabled.

REDIRECT_MAX_TIMES

Default: 20

The maximum number of redirections that will be followed for a single request.

MetaRefreshMiddleware

```
class scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware
```

This middleware handles redirection of requests based on meta-refresh html tag.

The *MetaRefreshMiddleware* can be configured through the following settings (see the settings documentation for more info):

- *METAREFRESH_ENABLED*
- *METAREFRESH_MAXDELAY*

This middleware obey *REDIRECT_MAX_TIMES* setting, *dont_redirect*, *redirect_urls* and *redirect_reasons* request meta keys as described for *RedirectMiddleware*

MetaRefreshMiddleware settings

METAREFRESH_ENABLED

0.17 新版功能.

Default: True

Whether the Meta Refresh middleware will be enabled.

METAREFRESH_MAXDELAY

Default: 100

The maximum meta-refresh delay (in seconds) to follow the redirection. Some sites use meta-refresh for redirecting to a session expired page, so we restrict automatic redirection to the maximum delay.

RetryMiddleware

```
class scrapy.downloadermiddlewares.retry.RetryMiddleware
```

A middleware to retry failed requests that are potentially caused by temporary problems such as a connection timeout or HTTP 500 error.

Failed pages are collected on the scraping process and rescheduled at the end, once the spider has finished crawling all regular (non failed) pages.

The *RetryMiddleware* can be configured through the following settings (see the settings documentation for more info):

- *RETRY_ENABLED*
- *RETRY_TIMES*
- *RETRY_HTTP_CODES*

If *Request.meta* has *dont_retry* key set to True, the request will be ignored by this middleware.

RetryMiddleware Settings

RETRY_ENABLED

0.13 新版功能.

Default: True

Whether the Retry middleware will be enabled.

RETRY_TIMES

Default: 2

Maximum number of times to retry, in addition to the first download.

Maximum number of retries can also be specified per-request using *max_retry_times* attribute of *Request.meta*. When initialized, the *max_retry_times* meta key takes higher precedence over the *RETRY_TIMES* setting.

RETRY_HTTP_CODES

Default: [500, 502, 503, 504, 522, 524, 408]

Which HTTP response codes to retry. Other errors (DNS lookup issues, connections lost, etc) are always retried.

In some cases you may want to add 400 to *RETRY_HTTP_CODES* because it is a common code used to indicate server overload. It is not included by default because HTTP specs say so.

RobotsTxtMiddleware

```
class scrapy.downloadermiddlewares.robotstxt.RobotsTxtMiddleware
```

This middleware filters out requests forbidden by the robots.txt exclusion standard.

To make sure Scrapy respects robots.txt make sure the middleware is enabled and the `ROBOTSTXT_OBEY` setting is enabled.

If `Request.meta` has `dont_obey_robotstxt` key set to `True` the request will be ignored by this middleware even if `ROBOTSTXT_OBEY` is enabled.

DownloaderStats

```
class scrapy.downloadermiddlewares.stats.DownloaderStats
```

Middleware that stores stats of all requests, responses and exceptions that pass through it.

To use this middleware you must enable the `DOWNLOADER_STATS` setting.

UserAgentMiddleware

```
class scrapy.downloadermiddlewares.useragent.UserAgentMiddleware
```

Middleware that allows spiders to override the default user agent.

In order for a spider to override the default user agent, its `user_agent` attribute must be set.

AjaxCrawlMiddleware

```
class scrapy.downloadermiddlewares.ajaxcrawl.AjaxCrawlMiddleware
```

Middleware that finds ‘AJAX crawlable’ page variants based on meta-fragment html tag. See <https://developers.google.com/webmasters/ajax-crawling/docs/getting-started> for more info.

注解: Scrapy finds ‘AJAX crawlable’ pages for URLs like `'http://example.com/!#foo=bar'` even without this middleware. `AjaxCrawlMiddleware` is necessary when URL doesn't contain `'!#'`. This is often a case for ‘index’ or ‘main’ website pages.

AjaxCrawlMiddleware Settings

AJAXCRAWL_ENABLED

0.21 新版功能.

Default: `False`

Whether the `AjaxCrawlMiddleware` will be enabled. You may want to enable it for *broad crawls*.

HttpProxyMiddleware settings

HTTPPROXY_ENABLED

Default: `True`

Whether or not to enable the `HttpProxyMiddleware`.

HTTPPROXY_AUTH_ENCODING

Default: `"latin-1"`

The default encoding for proxy authentication on `HttpProxyMiddleware`.

6.3 爬虫器中间件

The spider middleware is a framework of hooks into Scrapy's spider processing mechanism where you can plug custom functionality to process the responses that are sent to 爬虫器 for processing and to process the requests and items that are generated from spiders.

6.3.1 Activating a spider middleware

To activate a spider middleware component, add it to the `SPIDER_MIDDLEWARES` setting, which is a dict whose keys are the middleware class path and their values are the middleware orders.

Here's an example:

```
SPIDER_MIDDLEWARES = {  
    'myproject.middlewares.CustomSpiderMiddleware': 543,  
}
```

The `SPIDER_MIDDLEWARES` setting is merged with the `SPIDER_MIDDLEWARES_BASE` setting defined in Scrapy (and not meant to be overridden) and then sorted by order to get the final sorted list of enabled middlewares: the first middleware is the one closer to the engine and the last is the one closer to the spider. In other words, the `process_spider_input()` method of each middleware will be invoked in increasing middleware order (100, 200, 300, ...), and the `process_spider_output()` method of each middleware will be invoked in decreasing order.

To decide which order to assign to your middleware see the `SPIDER_MIDDLEWARES_BASE` setting and pick a value according to where you want to insert the middleware. The order does matter because each middleware performs a different action and your middleware could depend on some previous (or subsequent) middleware being applied.

If you want to disable a builtin middleware (the ones defined in `SPIDER_MIDDLEWARES_BASE`, and enabled by default) you must define it in your project `SPIDER_MIDDLEWARES` setting and assign `None` as its value. For example, if you want to disable the off-site middleware:

```
SPIDER_MIDDLEWARES = {
    'myproject.middlewares.CustomSpiderMiddleware': 543,
    'scrapy.spidermiddlewares.offsite.OffsiteMiddleware': None,
}
```

Finally, keep in mind that some middlewares may need to be enabled through a particular setting. See each middleware documentation for more info.

6.3.2 Writing your own spider middleware

Each spider middleware is a Python class that defines one or more of the methods defined below.

The main entry point is the `from_crawler` class method, which receives a *Crawler* instance. The *Crawler* object gives you access, for example, to the *settings*.

```
class scrapy.spidermiddlewares.SpiderMiddleware
```

```
    process_spider_input(response, spider)
```

This method is called for each response that goes through the spider middleware and into the spider, for processing.

process_spider_input() should return `None` or raise an exception.

If it returns `None`, Scrapy will continue processing this response, executing all other middlewares until, finally, the response is handed to the spider for processing.

If it raises an exception, Scrapy won't bother calling any other spider middleware *process_spider_input()* and will call the request errback if there is one, otherwise it will start the *process_spider_exception()* chain. The output of the errback is chained back in the other direction for *process_spider_output()* to process it, or *process_spider_exception()* if it raised an exception.

参数

- **response** (*Response* object) – the response being processed
- **spider** (*Spider* object) – the spider for which this response is intended

```
    process_spider_output(response, result, spider)
```

This method is called with the results returned from the Spider, after it has processed the response.

process_spider_output() must return an iterable of *Request*, dict or *Item* objects.

参数

- **response** (*Response* object) – the response which generated this output from the spider
- **result** (an iterable of *Request*, dict or *Item* objects) – the result returned by the spider
- **spider** (*Spider* object) – the spider whose result is being processed

process_spider_exception(*response*, *exception*, *spider*)

This method is called when a spider or *process_spider_output()* method (from a previous spider middleware) raises an exception.

process_spider_exception() should return either `None` or an iterable of *Request*, dict or *Item* objects.

If it returns `None`, Scrapy will continue processing this exception, executing any other *process_spider_exception()* in the following middleware components, until no middleware components are left and the exception reaches the engine (where it's logged and discarded).

If it returns an iterable the *process_spider_output()* pipeline kicks in, starting from the next spider middleware, and no other *process_spider_exception()* will be called.

参数

- **response** (*Response* object) – the response being processed when the exception was raised
- **exception** (*Exception* object) – the exception raised
- **spider** (*Spider* object) – the spider which raised the exception

process_start_requests(*start_requests*, *spider*)

0.15 新版功能.

This method is called with the start requests of the spider, and works similarly to the *process_spider_output()* method, except that it doesn't have a response associated and must return only requests (not items).

It receives an iterable (in the **start_requests** parameter) and must return another iterable of *Request* objects.

注解: When implementing this method in your spider middleware, you should always return an iterable (that follows the input one) and not consume all **start_requests** iterator because it can be very large (or even unbounded) and cause a memory overflow. The Scrapy engine is designed to pull start requests while it has capacity to process them, so the start requests iterator can be effectively endless where there is some other condition for stopping the spider (like a time limit or item/page count).

参数

- **start_requests** (an iterable of *Request*) – the start requests
- **spider** (*Spider* object) – the spider to whom the start requests belong

`from_crawler(cls, crawler)`

If present, this classmethod is called to create a middleware instance from a *Crawler*. It must return a new instance of the middleware. Crawler object provides access to all Scrapy core components like settings and signals; it is a way for middleware to access them and hook its functionality into Scrapy.

参数 `crawler` (*Crawler* object) – crawler that uses this middleware

6.3.3 Built-in spider middleware reference

This page describes all spider middleware components that come with Scrapy. For information on how to use them and how to write your own spider middleware, see the *spider middleware usage guide*.

For a list of the components enabled by default (and their orders) see the `SPIDER_MIDDLEWARES_BASE` setting.

DepthMiddleware

`class scrapy.spidermiddlewares.depth.DepthMiddleware`

DepthMiddleware is used for tracking the depth of each Request inside the site being scraped. It works by setting `request.meta['depth'] = 0` whenever there is no value previously set (usually just the first Request) and incrementing it by 1 otherwise.

It can be used to limit the maximum depth to scrape, control Request priority based on their depth, and things like that.

The *DepthMiddleware* can be configured through the following settings (see the settings documentation for more info):

- `DEPTH_LIMIT` - The maximum depth that will be allowed to crawl for any site. If zero, no limit will be imposed.
- `DEPTH_STATS_VERBOSE` - Whether to collect the number of requests for each depth.
- `DEPTH_PRIORITY` - Whether to prioritize the requests based on their depth.

HttpErrorMiddleware

`class scrapy.spidermiddlewares.httperror.HttpErrorMiddleware`

Filter out unsuccessful (erroneous) HTTP responses so that spiders don't have to deal with them, which (most of the time) imposes an overhead, consumes more resources, and makes the spider logic more complex.

According to the [HTTP standard](#), successful responses are those whose status codes are in the 200-300 range.

If you still want to process response codes outside that range, you can specify which response codes the spider is able to handle using the `handle_httpstatus_list` spider attribute or `HTTPERROR_ALLOWED_CODES` setting.

For example, if you want your spider to handle 404 responses you can do this:

```
class MySpider(CrawlSpider):
    handle_httpstatus_list = [404]
```

The `handle_httpstatus_list` key of `Request.meta` can also be used to specify which response codes to allow on a per-request basis. You can also set the meta key `handle_httpstatus_all` to `True` if you want to allow any response code for a request.

Keep in mind, however, that it's usually a bad idea to handle non-200 responses, unless you really know what you're doing.

For more information see: [HTTP Status Code Definitions](#).

HttpErrorMiddleware settings

HTTPERROR_ALLOWED_CODES

Default: `[]`

Pass all responses with non-200 status codes contained in this list.

HTTPERROR_ALLOW_ALL

Default: `False`

Pass all responses, regardless of its status code.

OffsiteMiddleware

```
class scrapy.spidermiddlewares.offsite.OffsiteMiddleware
```

Filters out Requests for URLs outside the domains covered by the spider.

This middleware filters out every request whose host names aren't in the spider's `allowed_domains` attribute. All subdomains of any domain in the list are also allowed. E.g. the rule `www.example.org` will also allow `bob.www.example.org` but not `www2.example.com` nor `example.com`.

When your spider returns a request for a domain not belonging to those covered by the spider, this middleware will log a debug message similar to this one:


```
DEBUG: Filtered offsite request to 'www.othersite.com': <GET http://www.othersite.  
↪com/some/page.html>
```

To avoid filling the log with too much noise, it will only print one of these messages for each new domain filtered. So, for example, if another request for `www.othersite.com` is filtered, no log message will be printed. But if a request for `someothersite.com` is filtered, a message will be printed (but only for the first request filtered).

If the spider doesn't define an `allowed_domains` attribute, or the attribute is empty, the offsite middleware will allow all requests.

If the request has the `dont_filter` attribute set, the offsite middleware will allow the request even if its domain is not listed in allowed domains.

RefererMiddleware

```
class scrapy.spidermiddlewares.referer.RefererMiddleware
```

Populates Request `Referer` header, based on the URL of the Response which generated it.

RefererMiddleware settings

REFERER_ENABLED

0.15 新版功能.

Default: `True`

Whether to enable referer middleware.

REFERRER_POLICY

1.4 新版功能.

Default: `'scrapy.spidermiddlewares.referer.DefaultReferrerPolicy'` [Referrer Policy](#) to apply when populating Request “Referer” header.

注解: You can also set the Referrer Policy per request, using the special `"referrer_policy"` [Request.meta](#) key, with the same acceptable values as for the `REFERRER_POLICY` setting.

Acceptable values for `REFERRER_POLICY`

- either a path to a `scrapy.spidermiddlewares.referer.RefererPolicy` subclass —a custom policy or one of the built-in ones (see classes below),
- or one of the standard W3C-defined string values,
- or the special `"scrapy-default"`.

String value	Class name (as a string)
<code>"scrapy-default"</code> (default)	<code>scrapy.spidermiddlewares.referer.DefaultRefererPolicy</code>
<code>"no-referrer"</code>	<code>scrapy.spidermiddlewares.referer.NoRefererPolicy</code>
<code>"no-referrer-when-downgrade"</code>	<code>scrapy.spidermiddlewares.referer.NoRefererWhenDowngradePolicy</code>
<code>"same-origin"</code>	<code>scrapy.spidermiddlewares.referer.SameOriginPolicy</code>
<code>"origin"</code>	<code>scrapy.spidermiddlewares.referer.OriginPolicy</code>
<code>"strict-origin"</code>	<code>scrapy.spidermiddlewares.referer.StrictOriginPolicy</code>
<code>"origin-when-cross-origin"</code>	<code>scrapy.spidermiddlewares.referer.OriginWhenCrossOriginPolicy</code>
<code>"strict-origin-when-cross-origin"</code>	<code>scrapy.spidermiddlewares.referer.StrictOriginWhenCrossOriginPolicy</code>
<code>"unsafe-url"</code>	<code>scrapy.spidermiddlewares.referer.UnsafeUrlPolicy</code>

警告: Scrapy’s default referrer policy —just like `"no-referrer-when-downgrade"`, the W3C-recommended value for browsers —will send a non-empty `"Referer"` header from any `http(s)://` to any `https://` URL, even if the domain is different.

`"same-origin"` may be a better choice if you want to remove referrer information for cross-domain requests.

注解: `"no-referrer-when-downgrade"` policy is the W3C-recommended default, and is used by major web browsers.

However, it is NOT Scrapy’s default referrer policy (see `DefaultRefererPolicy`).

警告: `"unsafe-url"` policy is NOT recommended.

UrlLengthMiddleware

```
class scrapy.spidermiddlewares.urllength.UrlLengthMiddleware
```

Filters out requests with URLs longer than `URLLENGTH_LIMIT`

The *UrlLengthMiddleware* can be configured through the following settings (see the settings documentation for more info):

- *URLLENGTH_LIMIT* - The maximum URL length to allow for crawled URLs.

6.4 扩展

The extensions framework provides a mechanism for inserting your own custom functionality into Scrapy.

Extensions are just regular classes that are instantiated at Scrapy startup, when extensions are initialized.

6.4.1 Extension settings

Extensions use the *Scrapy settings* to manage their settings, just like any other Scrapy code.

It is customary for extensions to prefix their settings with their own name, to avoid collision with existing (and future) extensions. For example, a hypothetical extension to handle [Google Sitemaps](#) would use settings like `GOOGLESITEMAP_ENABLED`, `GOOGLESITEMAP_DEPTH`, and so on.

6.4.2 Loading & activating extensions

Extensions are loaded and activated at startup by instantiating a single instance of the extension class. Therefore, all the extension initialization code must be performed in the class constructor (`__init__` method).

To make an extension available, add it to the *EXTENSIONS* setting in your Scrapy settings. In *EXTENSIONS*, each extension is represented by a string: the full Python path to the extension's class name. For example:

```
EXTENSIONS = {
    'scrapy.extensions.corestats.CoreStats': 500,
    'scrapy.extensions.telnet.TelnetConsole': 500,
}
```

As you can see, the *EXTENSIONS* setting is a dict where the keys are the extension paths, and their values are the orders, which define the extension *loading* order. The *EXTENSIONS* setting is merged with the *EXTENSIONS_BASE* setting defined in Scrapy (and not meant to be overridden) and then sorted by order to get the final sorted list of enabled extensions.

As extensions typically do not depend on each other, their loading order is irrelevant in most cases. This is why the *EXTENSIONS_BASE* setting defines all extensions with the same order (0). However, this feature can be exploited if you need to add an extension which depends on other extensions already loaded.

6.4.3 Available, enabled and disabled extensions

Not all available extensions will be enabled. Some of them usually depend on a particular setting. For example, the HTTP Cache extension is available by default but disabled unless the `HTTPCACHE_ENABLED` setting is set.

6.4.4 Disabling an extension

In order to disable an extension that comes enabled by default (ie. those included in the `EXTENSIONS_BASE` setting) you must set its order to `None`. For example:

```
EXTENSIONS = {  
    'scrapy.extensions.corestats.CoreStats': None,  
}
```

6.4.5 Writing your own extension

Each extension is a Python class. The main entry point for a Scrapy extension (this also includes middlewares and pipelines) is the `from_crawler` class method which receives a `Crawler` instance. Through the `Crawler` object you can access settings, signals, stats, and also control the crawling behaviour.

Typically, extensions connect to *signals* and perform tasks triggered by them.

Finally, if the `from_crawler` method raises the `NotConfigured` exception, the extension will be disabled. Otherwise, the extension will be enabled.

Sample extension

Here we will implement a simple extension to illustrate the concepts described in the previous section. This extension will log a message every time:

- a spider is opened
- a spider is closed
- a specific number of items are scraped

The extension will be enabled through the `MYEXT_ENABLED` setting and the number of items will be specified through the `MYEXT_ITEMCOUNT` setting.

Here is the code of such extension:

```
import logging  
from scrapy import signals  
from scrapy.exceptions import NotConfigured
```

(下页继续)

(续上页)

```
logger = logging.getLogger(__name__)

class SpiderOpenCloseLogging(object):

    def __init__(self, item_count):
        self.item_count = item_count
        self.items_scraped = 0

    @classmethod
    def from_crawler(cls, crawler):
        # first check if the extension should be enabled and raise
        # NotConfigured otherwise
        if not crawler.settings.getbool('MYEXT_ENABLED'):
            raise NotConfigured

        # get the number of items from settings
        item_count = crawler.settings.getint('MYEXT_ITEMCOUNT', 1000)

        # instantiate the extension object
        ext = cls(item_count)

        # connect the extension object to signals
        crawler.signals.connect(ext.spider_opened, signal=signals.spider_opened)
        crawler.signals.connect(ext.spider_closed, signal=signals.spider_closed)
        crawler.signals.connect(ext.item_scraped, signal=signals.item_scraped)

        # return the extension object
        return ext

    def spider_opened(self, spider):
        logger.info("opened spider %s", spider.name)

    def spider_closed(self, spider):
        logger.info("closed spider %s", spider.name)

    def item_scraped(self, item, spider):
        self.items_scraped += 1
        if self.items_scraped % self.item_count == 0:
            logger.info("scraped %d items", self.items_scraped)
```

6.4.6 Built-in extensions reference

General purpose extensions

Log Stats extension

```
class scrapy.extensions.logstats.LogStats
```

Log basic stats like crawled pages and scraped items.

Core Stats extension

```
class scrapy.extensions.corestats.CoreStats
```

Enable the collection of core statistics, provided the stats collection is enabled (see [统计数据集合](#)).

Telnet console extension

```
class scrapy.extensions.telnet.TelnetConsole
```

Provides a telnet console for getting into a Python interpreter inside the currently running Scrapy process, which can be very useful for debugging.

The telnet console must be enabled by the [TELNETCONSOLE_ENABLED](#) setting, and the server will listen in the port specified in [TELNETCONSOLE_PORT](#).

Memory usage extension

```
class scrapy.extensions.memusage.MemoryUsage
```

注解: This extension does not work in Windows.

Monitors the memory used by the Scrapy process that runs the spider and:

1. sends a notification e-mail when it exceeds a certain value
2. closes the spider when it exceeds a certain value

The notification e-mails can be triggered when a certain warning value is reached ([MEMUSAGE_WARNING_MB](#)) and when the maximum value is reached ([MEMUSAGE_LIMIT_MB](#)) which will also cause the spider to be closed and the Scrapy process to be terminated.

This extension is enabled by the [MEMUSAGE_ENABLED](#) setting and can be configured with the following settings:

- [MEMUSAGE_LIMIT_MB](#)

- `MEMUSAGE_WARNING_MB`
- `MEMUSAGE_NOTIFY_MAIL`
- `MEMUSAGE_CHECK_INTERVAL_SECONDS`

Memory debugger extension

`class scrapy.extensions.memdebug.MemoryDebugger`

An extension for debugging memory usage. It collects information about:

- objects uncollected by the Python garbage collector
- objects left alive that shouldn't. For more info, see *Debugging memory leaks with trackref*

To enable this extension, turn on the `MEMDEBUG_ENABLED` setting. The info will be stored in the stats.

Close spider extension

`class scrapy.extensions.closespider.CloseSpider`

Closes a spider automatically when some conditions are met, using a specific closing reason for each condition.

The conditions for closing a spider can be configured through the following settings:

- `CLOSESPIDER_TIMEOUT`
- `CLOSESPIDER_ITEMCOUNT`
- `CLOSESPIDER_PAGECOUNT`
- `CLOSESPIDER_ERRORCOUNT`

CLOSESPIDER_TIMEOUT

Default: 0

An integer which specifies a number of seconds. If the spider remains open for more than that number of second, it will be automatically closed with the reason `closespider_timeout`. If zero (or non set), spiders won't be closed by timeout.

CLOSESPIDER_ITEMCOUNT

Default: 0

An integer which specifies a number of items. If the spider scrapes more than that amount and those items are passed by the item pipeline, the spider will be closed with the reason `closespider_itemcount`. Requests

which are currently in the downloader queue (up to `CONCURRENT_REQUESTS` requests) are still processed. If zero (or non set), spiders won't be closed by number of passed items.

CLOSESPIDER_PAGECOUNT

0.11 新版功能.

Default: 0

An integer which specifies the maximum number of responses to crawl. If the spider crawls more than that, the spider will be closed with the reason `closespider_pagecount`. If zero (or non set), spiders won't be closed by number of crawled responses.

CLOSESPIDER_ERRORCOUNT

0.11 新版功能.

Default: 0

An integer which specifies the maximum number of errors to receive before closing the spider. If the spider generates more than that number of errors, it will be closed with the reason `closespider_errorcount`. If zero (or non set), spiders won't be closed by number of errors.

StatsMailer extension

```
class scrapy.extensions.statsmailer.StatsMailer
```

This simple extension can be used to send a notification e-mail every time a domain has finished scraping, including the Scrapy stats collected. The email will be sent to all recipients specified in the `STATSMAILER_RCPTS` setting.

Debugging extensions

Stack trace dump extension

```
class scrapy.extensions.debug.StackTraceDump
```

Dumps information about the running process when a `SIGQUIT` or `SIGUSR2` signal is received. The information dumped is the following:

1. engine status (using `scrapy.utils.engine.get_engine_status()`)
2. live references (see *Debugging memory leaks with trackref*)
3. stack trace of all threads

After the stack trace and engine status is dumped, the Scrapy process continues running normally.

This extension only works on POSIX-compliant platforms (ie. not Windows), because the `SIGQUIT` and `SIGUSR2` signals are not available on Windows.

There are at least two ways to send Scrapy the `SIGQUIT` signal:

1. By pressing Ctrl-while a Scrapy process is running (Linux only?)
2. By running this command (assuming `<pid>` is the process id of the Scrapy process):

```
kill -QUIT <pid>
```

Debugger extension

```
class scrapy.extensions.debug.Debugger
```

Invokes a `Python debugger` inside a running Scrapy process when a `SIGUSR2` signal is received. After the debugger is exited, the Scrapy process continues running normally.

For more info see [Debugging in Python](#).

This extension only works on POSIX-compliant platforms (ie. not Windows).

6.5 核心 API

0.15 新版功能.

This section documents the Scrapy core API, and it's intended for developers of extensions and middlewares.

6.5.1 Crawler API

The main entry point to Scrapy API is the `Crawler` object, passed to extensions through the `from_crawler` class method. This object provides access to all Scrapy core components, and it's the only way for extensions to access them and hook their functionality into Scrapy.

The Extension Manager is responsible for loading and keeping track of installed extensions and it's configured through the `EXTENSIONS` setting which contains a dictionary of all available extensions and their order similar to how you *configure the downloader middlewares*.

```
class scrapy.crawler.Crawler(spidercls, settings)
```

The Crawler object must be instantiated with a `scrapy.spiders.Spider` subclass and a `scrapy.settings.Settings` object.

settings

The settings manager of this crawler.

This is used by extensions & middlewares to access the Scrapy settings of this crawler.

For an introduction on Scrapy settings see [设置](#).

For the API see `Settings` class.

signals

The signals manager of this crawler.

This is used by extensions & middlewares to hook themselves into Scrapy functionality.

For an introduction on signals see [信号](#).

For the API see `SignalManager` class.

stats

The stats collector of this crawler.

This is used from extensions & middlewares to record stats of their behaviour, or access stats collected by other extensions.

For an introduction on stats collection see [统计数据集合](#).

For the API see `StatsCollector` class.

extensions

The extension manager that keeps track of enabled extensions.

Most extensions won't need to access this attribute.

For an introduction on extensions and a list of available extensions on Scrapy see [扩展](#).

engine

The execution engine, which coordinates the core crawling logic between the scheduler, downloader and spiders.

Some extension may want to access the Scrapy engine, to inspect or modify the downloader and scheduler behaviour, although this is an advanced use and this API is not yet stable.

spider

Spider currently being crawled. This is an instance of the spider class provided while constructing the crawler, and it is created after the arguments given in the `crawl()` method.

`crawl(*args, **kwargs)`

Starts the crawler by instantiating its spider class with the given `args` and `kwargs` arguments, while setting the execution engine in motion.

Returns a deferred that is fired when the crawl is finished.

6.5.2 Settings API

`scrapy.settings.SETTINGS_PRIORITIES`

Dictionary that sets the key name and priority level of the default settings priorities used in Scrapy.

Each item defines a settings entry point, giving it a code name for identification and an integer priority. Greater priorities take more precedence over lesser ones when setting and retrieving values in the `Settings` class.

```
SETTINGS_PRIORITIES = {
    'default': 0,
    'command': 10,
    'project': 20,
    'spider': 30,
    'cmdline': 40,
}
```

For a detailed explanation on each settings sources, see: [设置](#).

6.5.3 SpiderLoader API

`class scrapy.loader.SpiderLoader`

This class is in charge of retrieving and handling the spider classes defined across the project.

Custom spider loaders can be employed by specifying their path in the `SPIDER_LOADER_CLASS` project setting. They must fully implement the `scrapy.interfaces.ISpiderLoader` interface to guarantee an errorless execution.

`from_settings(settings)`

This class method is used by Scrapy to create an instance of the class. It's called with the current project settings, and it loads the spiders found recursively in the modules of the `SPIDER_MODULES` setting.

参数 settings (`Settings` instance) – project settings

`load(spider_name)`

Get the Spider class with the given name. It'll look into the previously loaded spiders for a spider class with name `spider_name` and will raise a `KeyError` if not found.

参数 spider_name (`str`) – spider class name

`list()`

Get the names of the available spiders in the project.

`find_by_request(request)`

List the spiders' names that can handle the given request. Will try to match the request's url against the domains of the spiders.

参数 `request` (*Request* instance) – queried request

6.5.4 Signals API

6.5.5 Stats Collector API

There are several Stats Collectors available under the `scrapy.statscollectors` module and they all implement the Stats Collector API defined by the *StatsCollector* class (which they all inherit from).

```
class scrapy.statscollectors.StatsCollector
```

```
    get_value(key, default=None)
```

Return the value for the given stats key or default if it doesn't exist.

```
    get_stats()
```

Get all stats from the currently running spider as a dict.

```
    set_value(key, value)
```

Set the given value for the given stats key.

```
    set_stats(stats)
```

Override the current stats with the dict passed in `stats` argument.

```
    inc_value(key, count=1, start=0)
```

Increment the value of the given stats key, by the given count, assuming the start value given (when it's not set).

```
    max_value(key, value)
```

Set the given value for the given key only if current value for the same key is lower than value. If there is no current value for the given key, the value is always set.

```
    min_value(key, value)
```

Set the given value for the given key only if current value for the same key is greater than value. If there is no current value for the given key, the value is always set.

```
    clear_stats()
```

Clear all stats.

The following methods are not part of the stats collection api but instead used when implementing custom stats collectors:

```
    open_spider(spider)
```

Open the given spider for stats collection.

```
    close_spider(spider)
```

Close the given spider. After this is called, no more specific stats can be accessed or collected.

6.6 信号

Scrapy uses signals extensively to notify when certain events occur. You can catch some of those signals in your Scrapy project (using an *extension*, for example) to perform additional tasks or extend Scrapy to add functionality not provided out of the box.

Even though signals provide several arguments, the handlers that catch them don't need to accept all of them - the signal dispatching mechanism will only deliver the arguments that the handler receives.

You can connect to signals (or send your own) through the *Signals API*.

Here is a simple example showing how you can catch signals and perform some action:

```
from scrapy import signals
from scrapy import Spider

class DmozSpider(Spider):
    name = "dmoz"
    allowed_domains = ["dmoz.org"]
    start_urls = [
        "http://www.dmoz.org/Computers/Programming/Languages/Python/Books/",
        "http://www.dmoz.org/Computers/Programming/Languages/Python/Resources/",
    ]

    @classmethod
    def from_crawler(cls, crawler, *args, **kwargs):
        spider = super(DmozSpider, cls).from_crawler(crawler, *args, **kwargs)
        crawler.signals.connect(spider.spider_closed, signal=signals.spider_closed)
        return spider

    def spider_closed(self, spider):
        spider.logger.info('Spider closed: %s', spider.name)

    def parse(self, response):
        pass
```

6.6.1 Deferred signal handlers

Some signals support returning Twisted deferreds from their handlers, see the *Built-in signals reference* below to know which ones.

6.6.2 Built-in signals reference

Here's the list of Scrapy built-in signals and their meaning.

engine_started

`scrapy.signals.engine_started()`

Sent when the Scrapy engine has started crawling.

This signal supports returning deferreds from their handlers.

注解: This signal may be fired *after* the *spider_opened* signal, depending on how the spider was started. So **don't** rely on this signal getting fired before *spider_opened*.

engine_stopped

`scrapy.signals.engine_stopped()`

Sent when the Scrapy engine is stopped (for example, when a crawling process has finished).

This signal supports returning deferreds from their handlers.

item_scraped

`scrapy.signals.item_scraped(item, response, spider)`

Sent when an item has been scraped, after it has passed all the *Item* 管道 stages (without being dropped).

This signal supports returning deferreds from their handlers.

参数

- **item** (dict or *Item* object) – the item scraped
- **spider** (*Spider* object) – the spider which scraped the item
- **response** (*Response* object) – the response from where the item was scraped

item_dropped

`scrapy.signals.item_dropped(item, response, exception, spider)`

Sent after an item has been dropped from the *Item* 管道 when some stage raised a *DropItem* exception.

This signal supports returning deferreds from their handlers.

参数

- **item** (dict or *Item* object) – the item dropped from the *Item* 管道
- **spider** (*Spider* object) – the spider which scraped the item
- **response** (*Response* object) – the response from where the item was dropped
- **exception** (*DropItem* exception) – the exception (which must be a *DropItem* sub-class) which caused the item to be dropped

item_error

`scrapy.signals.item_error(item, response, spider, failure)`

Sent when a *Item* 管道 generates an error (ie. raises an exception), except *DropItem* exception.

This signal supports returning deferreds from their handlers.

参数

- **item** (dict or *Item* object) – the item dropped from the *Item* 管道
- **response** (*Response* object) – the response being processed when the exception was raised
- **spider** (*Spider* object) – the spider which raised the exception
- **failure** (*Failure* object) – the exception raised as a Twisted *Failure* object

spider_closed

`scrapy.signals.spider_closed(spider, reason)`

Sent after a spider has been closed. This can be used to release per-spider resources reserved on *spider_opened*.

This signal supports returning deferreds from their handlers.

参数

- **spider** (*Spider* object) – the spider which has been closed
- **reason** (*str*) – a string which describes the reason why the spider was closed. If it was closed because the spider has completed scraping, the reason is 'finished'. Otherwise, if the spider was manually closed by calling the `close_spider` engine

method, then the reason is the one passed in the `reason` argument of that method (which defaults to `'cancelled'`). If the engine was shutdown (for example, by hitting Ctrl-C to stop it) the reason will be `'shutdown'`.

`spider_opened`

`scrapy.signals.spider_opened(spider)`

Sent after a spider has been opened for crawling. This is typically used to reserve per-spider resources, but can be used for any task that needs to be performed when a spider is opened.

This signal supports returning deferreds from their handlers.

参数 `spider` (*Spider* object) – the spider which has been opened

`spider_idle`

`scrapy.signals.spider_idle(spider)`

Sent when a spider has gone idle, which means the spider has no further:

- requests waiting to be downloaded
- requests scheduled
- items being processed in the item pipeline

If the idle state persists after all handlers of this signal have finished, the engine starts closing the spider. After the spider has finished closing, the `spider_closed` signal is sent.

You may raise a `DontCloseSpider` exception to prevent the spider from being closed.

This signal does not support returning deferreds from their handlers.

参数 `spider` (*Spider* object) – the spider which has gone idle

注解: Scheduling some requests in your `spider_idle` handler does **not** guarantee that it can prevent the spider from being closed, although it sometimes can. That's because the spider may still remain idle if all the scheduled requests are rejected by the scheduler (e.g. filtered due to duplication).

`spider_error`

`scrapy.signals.spider_error(failure, response, spider)`

Sent when a spider callback generates an error (ie. raises an exception).

This signal does not support returning deferreds from their handlers.

参数

- **failure** (*Failure* object) – the exception raised as a Twisted *Failure* object
- **response** (*Response* object) – the response being processed when the exception was raised
- **spider** (*Spider* object) – the spider which raised the exception

request_scheduled

`scrapy.signals.request_scheduled(request, spider)`

Sent when the engine schedules a *Request*, to be downloaded later.

The signal does not support returning deferreds from their handlers.

参数

- **request** (*Request* object) – the request that reached the scheduler
- **spider** (*Spider* object) – the spider that yielded the request

request_dropped

`scrapy.signals.request_dropped(request, spider)`

Sent when a *Request*, scheduled by the engine to be downloaded later, is rejected by the scheduler.

The signal does not support returning deferreds from their handlers.

参数

- **request** (*Request* object) – the request that reached the scheduler
- **spider** (*Spider* object) – the spider that yielded the request

request_reached_downloader

`scrapy.signals.request_reached_downloader(request, spider)`

Sent when a *Request* reached downloader.

The signal does not support returning deferreds from their handlers.

参数

- **request** (*Request* object) – the request that reached downloader
- **spider** (*Spider* object) – the spider that yielded the request

response_received

`scrapy.signals.response_received(response, request, spider)`

Sent when the engine receives a new *Response* from the downloader.

This signal does not support returning deferreds from their handlers.

参数

- **response** (*Response* object) – the response received
- **request** (*Request* object) – the request that generated the response
- **spider** (*Spider* object) – the spider for which the response is intended

response_downloaded

`scrapy.signals.response_downloaded(response, request, spider)`

Sent by the downloader right after a `HTTPResponse` is downloaded.

This signal does not support returning deferreds from their handlers.

参数

- **response** (*Response* object) – the response downloaded
- **request** (*Request* object) – the request that generated the response
- **spider** (*Spider* object) – the spider for which the response is intended

6.7 Item 导出文件

Once you have scraped your items, you often want to persist or export those items, to use the data in some other application. That is, after all, the whole purpose of the scraping process.

For this purpose Scrapy provides a collection of Item Exporters for different output formats, such as XML, CSV or JSON.

6.7.1 Using Item Exporters

If you are in a hurry, and just want to use an Item Exporter to output scraped data see the *Feed 导出*. Otherwise, if you want to know how Item Exporters work or need more custom functionality (not covered by the default exports), continue reading below.

In order to use an Item Exporter, you must instantiate it with its required args. Each Item Exporter requires different arguments, so check each exporter documentation to be sure, in *Built-in Item Exporters reference*. After you have instantiated your exporter, you have to:

1. call the method `start_exporting()` in order to signal the beginning of the exporting process
2. call the `export_item()` method for each item you want to export
3. and finally call the `finish_exporting()` to signal the end of the exporting process

Here you can see an *Item Pipeline* which uses multiple Item Exporters to group scraped items to different files according to the value of one of their fields:

```
from scrapy.exporters import XmlItemExporter

class PerYearXmlExportPipeline(object):
    """Distribute items across multiple XML files according to their 'year' field"""

    def open_spider(self, spider):
        self.year_to_exporter = {}

    def close_spider(self, spider):
        for exporter in self.year_to_exporter.values():
            exporter.finish_exporting()
            exporter.file.close()

    def _exporter_for_item(self, item):
        year = item['year']
        if year not in self.year_to_exporter:
            f = open('{}.xml'.format(year), 'wb')
            exporter = XmlItemExporter(f)
            exporter.start_exporting()
            self.year_to_exporter[year] = exporter
        return self.year_to_exporter[year]

    def process_item(self, item, spider):
        exporter = self._exporter_for_item(item)
        exporter.export_item(item)
        return item
```

6.7.2 Serialization of item fields

By default, the field values are passed unmodified to the underlying serialization library, and the decision of how to serialize them is delegated to each particular serialization library.

However, you can customize how each field value is serialized *before it is passed to the serialization library*.

There are two ways to customize how a field will be serialized, which are described next.

1. Declaring a serializer in the field

If you use *Item* you can declare a serializer in the *field metadata*. The serializer must be a callable which receives a value and returns its serialized form.

Example:

```
import scrapy

def serialize_price(value):
    return '$ %s' % str(value)

class Product(scrapy.Item):
    name = scrapy.Field()
    price = scrapy.Field(serializer=serialize_price)
```

2. Overriding the `serialize_field()` method

You can also override the *serialize_field()* method to customize how your field value will be exported.

Make sure you call the base class *serialize_field()* method after your custom code.

Example:

```
from scrapy.exporter import XmlItemExporter

class ProductXmlExporter(XmlItemExporter):

    def serialize_field(self, field, name, value):
        if field == 'price':
            return '$ %s' % str(value)
        return super(Product, self).serialize_field(field, name, value)
```

6.7.3 Built-in Item Exporters reference

Here is a list of the Item Exporters bundled with Scrapy. Some of them contain output examples, which assume you're exporting these two items:

```
Item(name='Color TV', price='1200')
Item(name='DVD player', price='200')
```

BaseItemExporter

```
class scrapy.exporters.BaseItemExporter(fields_to_export=None, export_empty_fields=False,
                                       encoding='utf-8', indent=0)
```

This is the (abstract) base class for all Item Exporters. It provides support for common features used by all (concrete) Item Exporters, such as defining what fields to export, whether to export empty fields, or which encoding to use.

These features can be configured through the constructor arguments which populate their respective instance attributes: *fields_to_export*, *export_empty_fields*, *encoding*, *indent*.

export_item(item)

Exports the given item. This method must be implemented in subclasses.

serialize_field(field, name, value)

Return the serialized value for the given field. You can override this method (in your custom Item Exporters) if you want to control how a particular field or value will be serialized/exported.

By default, this method looks for a serializer *declared in the item field* and returns the result of applying that serializer to the value. If no serializer is found, it returns the value unchanged except for **unicode** values which are encoded to **str** using the encoding declared in the *encoding* attribute.

参数

- **field** (*Field* object or an empty dict) – the field being serialized. If a raw dict is being exported (not *Item*) *field* value is an empty dict.
- **name** (*str*) – the name of the field being serialized
- **value** – the value being serialized

start_exporting()

Signal the beginning of the exporting process. Some exporters may use this to generate some required header (for example, the *XmlItemExporter*). You must call this method before exporting any items.

finish_exporting()

Signal the end of the exporting process. Some exporters may use this to generate some required footer (for example, the *XmlItemExporter*). You must always call this method after you have no more items to export.

fields_to_export

A list with the name of the fields that will be exported, or None if you want to export all fields. Defaults to None.

Some exporters (like *CsvItemExporter*) respect the order of the fields defined in this attribute.

Some exporters may require *fields_to_export* list in order to export the data properly when spiders return dicts (not *Item* instances).

export_empty_fields

Whether to include empty/unpopulated item fields in the exported data. Defaults to **False**. Some exporters (like *CsvItemExporter*) ignore this attribute and always export all empty fields.

This option is ignored for dict items.

encoding

The encoding that will be used to encode unicode values. This only affects unicode values (which are always serialized to str using this encoding). Other value types are passed unchanged to the specific serialization library.

indent

Amount of spaces used to indent the output on each level. Defaults to 0.

- **indent=None** selects the most compact representation, all items in the same line with no indentation
- **indent<=0** each item on its own line, no indentation
- **indent>0** each item on its own line, indented with the provided numeric value

XmlItemExporter

```
class scrapy.exporters.XmlItemExporter(file, item_element='item', root_element='items',  
                                     **kwargs)
```

Exports Items in XML format to the specified file object.

参数

- **file** – the file-like object to use for exporting the data. Its **write** method should accept **bytes** (a disk file opened in binary mode, a **io.BytesIO** object, etc)
- **root_element** (*str*) – The name of root element in the exported XML.
- **item_element** (*str*) – The name of each item element in the exported XML.

The additional keyword arguments of this constructor are passed to the *BaseItemExporter* constructor.

A typical output of this exporter would be:

```
<?xml version="1.0" encoding="utf-8"?>  
<items>  
  <item>  
    <name>Color TV</name>  
    <price>1200</price>  
  </item>  
  <item>  
    <name>DVD player</name>  
    <price>200</price>
```

(下页继续)

(续上页)

```
</item>
</items>
```

Unless overridden in the `serialize_field()` method, multi-valued fields are exported by serializing each value inside a `<value>` element. This is for convenience, as multi-valued fields are very common.

For example, the item:

```
Item(name=['John', 'Doe'], age='23')
```

Would be serialized as:

```
<?xml version="1.0" encoding="utf-8"?>
<items>
  <item>
    <name>
      <value>John</value>
      <value>Doe</value>
    </name>
    <age>23</age>
  </item>
</items>
```

CsvItemExporter

```
class scrapy.exporters.CsvItemExporter(file, include_headers_line=True, join_multivalued=' ', **kwargs)
```

Exports Items in CSV format to the given file-like object. If the `fields_to_export` attribute is set, it will be used to define the CSV columns and their order. The `export_empty_fields` attribute has no effect on this exporter.

参数

- **file** – the file-like object to use for exporting the data. Its `write` method should accept `bytes` (a disk file opened in binary mode, a `io.BytesIO` object, etc)
- **include_headers_line** (*str*) – If enabled, makes the exporter output a header line with the field names taken from `BaseItemExporter.fields_to_export` or the first exported item fields.
- **join_multivalued** – The char (or chars) that will be used for joining multi-valued fields, if found.

The additional keyword arguments of this constructor are passed to the `BaseItemExporter` constructor, and the leftover arguments to the `csv.writer` constructor, so you can use any `csv.writer` constructor

argument to customize this exporter.

A typical output of this exporter would be:

```
product,price
Color TV,1200
DVD player,200
```

PickleItemExporter

```
class scrapy.exporters.PickleItemExporter(file, protocol=0, **kwargs)
```

Exports Items in pickle format to the given file-like object.

参数

- **file** – the file-like object to use for exporting the data. Its **write** method should accept **bytes** (a disk file opened in binary mode, a `io.BytesIO` object, etc)
- **protocol** (*int*) – The pickle protocol to use.

For more information, refer to the [pickle module documentation](#).

The additional keyword arguments of this constructor are passed to the `BaseItemExporter` constructor.

Pickle isn't a human readable format, so no output examples are provided.

PprintItemExporter

```
class scrapy.exporters.PprintItemExporter(file, **kwargs)
```

Exports Items in pretty print format to the specified file object.

参数 **file** – the file-like object to use for exporting the data. Its **write** method should accept **bytes** (a disk file opened in binary mode, a `io.BytesIO` object, etc)

The additional keyword arguments of this constructor are passed to the `BaseItemExporter` constructor.

A typical output of this exporter would be:

```
{'name': 'Color TV', 'price': '1200'}
{'name': 'DVD player', 'price': '200'}
```

Longer lines (when present) are pretty-formatted.

JsonItemExporter

```
class scrapy.exporters.JsonItemExporter(file, **kwargs)
```

Exports Items in JSON format to the specified file-like object, writing all objects as a list of objects. The additional constructor arguments are passed to the `BaseItemExporter` constructor, and the leftover

arguments to the `JSONEncoder` constructor, so you can use any `JSONEncoder` constructor argument to customize this exporter.

参数 file – the file-like object to use for exporting the data. Its `write` method should accept `bytes` (a disk file opened in binary mode, a `io.BytesIO` object, etc)

A typical output of this exporter would be:

```
[{"name": "Color TV", "price": "1200"},
{"name": "DVD player", "price": "200"}]
```

警告: JSON is very simple and flexible serialization format, but it doesn't scale well for large amounts of data since incremental (aka. stream-mode) parsing is not well supported (if at all) among JSON parsers (on any language), and most of them just parse the entire object in memory. If you want the power and simplicity of JSON with a more stream-friendly format, consider using `JsonLinesItemExporter` instead, or splitting the output in multiple chunks.

JsonLinesItemExporter

`class scrapy.exporters.JsonLinesItemExporter(file, **kwargs)`

Exports Items in JSON format to the specified file-like object, writing one JSON-encoded item per line. The additional constructor arguments are passed to the `BaseItemExporter` constructor, and the leftover arguments to the `JSONEncoder` constructor, so you can use any `JSONEncoder` constructor argument to customize this exporter.

参数 file – the file-like object to use for exporting the data. Its `write` method should accept `bytes` (a disk file opened in binary mode, a `io.BytesIO` object, etc)

A typical output of this exporter would be:

```
{"name": "Color TV", "price": "1200"}
{"name": "DVD player", "price": "200"}
```

Unlike the one produced by `JsonItemExporter`, the format produced by this exporter is well suited for serializing large amounts of data.

框架体系 理解 Scrapy 的架构。

下载器中间件 定制爬虫页面如何请求和下载。

爬虫器中间件 自定义你的爬虫器的输入和输出。

扩展 使用你自定义的函数扩展 Scrapy。

核心 API 在扩展和中间件上使用它来扩展 Scrapy 功能。

信号 查看所有可用的信号以及如何使用它们。

Item **导出文件** 快速导出你的抓取项目到一个文件 (XML, CSV 等)。

7.1 Scrapy 更新

7.1.1 Scrapy 1.6.0 (2019-01-30)

Highlights:

- better Windows support;
- Python 3.7 compatibility;
- big documentation improvements, including a switch from `.extract_first()` + `.extract()` API to `.get()` + `.getall()` API;
- feed exports, FilePipeline and MediaPipeline improvements;
- better extensibility: `item_error` and `request_reached_downloader` signals; `from_crawler` support for feed exporters, feed storages and dupefilters.
- `scrapy.contracts` fixes and new features;
- telnet console security improvements, first released as a backport in *Scrapy 1.5.2 (2019-01-22)*;
- clean-up of the deprecated code;
- various bug fixes, small new features and usability improvements across the codebase.

Selector API changes

While these are not changes in Scrapy itself, but rather in the [parsel](#) library which Scrapy uses for xpath/css selectors, these changes are worth mentioning here. Scrapy now depends on `parsel >= 1.5`, and Scrapy documentation is updated to follow recent `parsel` API conventions.

Most visible change is that `.get()` and `.getall()` selector methods are now preferred over `.extract_first()` and `.extract()`. We feel that these new methods result in a more concise and readable code. See [extract\(\)](#) and [extract_first\(\)](#) for more details.

注解: There are currently **no plans** to deprecate `.extract()` and `.extract_first()` methods.

Another useful new feature is the introduction of `Selector.attrib` and `SelectorList.attrib` properties, which make it easier to get attributes of HTML elements. See [Selecting element attributes](#).

CSS selectors are cached in `parsel >= 1.5`, which makes them faster when the same CSS path is used many times. This is very common in case of Scrapy spiders: callbacks are usually called several times, on different pages.

If you're using custom `Selector` or `SelectorList` subclasses, a **backward incompatible** change in `parsel` may affect your code. See [parsel changelog](#) for a detailed description, as well as for the full list of improvements.

Telnet console

Backward incompatible: Scrapy's telnet console now requires username and password. See [远程控制台](#) for more details. This change fixes a **security issue**; see [Scrapy 1.5.2 \(2019-01-22\)](#) release notes for details.

New extensibility features

- `from_crawler` support is added to feed exporters and feed storages. This, among other things, allows to access Scrapy settings from custom feed storages and exporters ([issue 1605](#), [issue 3348](#)).
- `from_crawler` support is added to dupefilters ([issue 2956](#)); this allows to access e.g. settings or a spider from a dupefilter.
- `item_error` is fired when an error happens in a pipeline ([issue 3256](#));
- `request_reached_downloader` is fired when Downloader gets a new Request; this signal can be useful e.g. for custom Schedulers ([issue 3393](#)).
- new SitemapSpider `sitemap_filter()` method which allows to select sitemap entries based on their attributes in SitemapSpider subclasses ([issue 3512](#)).
- Lazy loading of Downloader Handlers is now optional; this enables better initialization error handling in custom Downloader Handlers ([issue 3394](#)).

New FilePipeline and MediaPipeline features

- Expose more options for S3FilesStore: `AWS_ENDPOINT_URL`, `AWS_USE_SSL`, `AWS_VERIFY`, `AWS_REGION_NAME`. For example, this allows to use alternative or self-hosted AWS-compatible providers (issue 2609, issue 3548).
- ACL support for Google Cloud Storage: `FILES_STORE_GCS_ACL` and `IMAGES_STORE_GCS_ACL` (issue 3199).

scrapy.contracts improvements

- Exceptions in contracts code are handled better (issue 3377);
- `dont_filter=True` is used for contract requests, which allows to test different callbacks with the same URL (issue 3381);
- `request_cls` attribute in Contract subclasses allow to use different Request classes in contracts, for example FormRequest (issue 3383).
- Fixed errback handling in contracts, e.g. for cases where a contract is executed for URL which returns non-200 response (issue 3371).

Usability improvements

- more stats for RobotsTxtMiddleware (issue 3100)
- INFO log level is used to show telnet host/port (issue 3115)
- a message is added to IgnoreRequest in RobotsTxtMiddleware (issue 3113)
- better validation of `url` argument in `Response.follow` (issue 3131)
- non-zero exit code is returned from Scrapy commands when error happens on spider initialization (issue 3226)
- Link extraction improvements: “ftp” is added to scheme list (issue 3152); “flv” is added to common video extensions (issue 3165)
- better error message when an exporter is disabled (issue 3358);
- `scrapy shell --help` mentions syntax required for local files (`./file.html`) - issue 3496.
- Referer header value is added to RFPDufilter log messages (issue 3588)

Bug fixes

- fixed issue with extra blank lines in .csv exports under Windows (issue 3039);
- proper handling of pickling errors in Python 3 when serializing objects for disk queues (issue 3082)

- flags are now preserved when copying Requests (issue 3342);
- `FormRequest.from_response` clickdata shouldn't ignore elements with `input[type=image]` (issue 3153).
- `FormRequest.from_response` should preserve duplicate keys (issue 3247)

Documentation improvements

- Docs are re-written to suggest `.get/.getall` API instead of `.extract/.extract_first`. Also, 选择器 docs are updated and re-structured to match latest parse docs; they now contain more topics, such as *Selecting element attributes* or *Extensions to CSS Selectors* (issue 3390).
- 使用浏览器的开发工具进行抓取 is a new tutorial which replaces old Firefox and Firebug tutorials (issue 3400).
- `SCRAPY_PROJECT` environment variable is documented (issue 3518);
- troubleshooting section is added to install instructions (issue 3517);
- improved links to beginner resources in the tutorial (issue 3367, issue 3468);
- fixed `RETRY_HTTP_CODES` default values in docs (issue 3335);
- remove unused `DEPTH_STATS` option from docs (issue 3245);
- other cleanups (issue 3347, issue 3350, issue 3445, issue 3544, issue 3605).

Deprecation removals

Compatibility shims for pre-1.0 Scrapy module names are removed (issue 3318):

- `scrapy.command`
- `scrapy.contrib` (with all submodules)
- `scrapy.contrib_exp` (with all submodules)
- `scrapy.dupefilter`
- `scrapy.linkextractor`
- `scrapy.project`
- `scrapy.spider`
- `scrapy.spidermanager`
- `scrapy.squeue`
- `scrapy.stats`
- `scrapy.statscol`

- `scrapy.utils.decorator`

See [Module Relocations](#) for more information, or use suggestions from Scrapy 1.5.x deprecation warnings to update your code.

Other deprecation removals:

- Deprecated `scrapy.interfaces.ISpiderManager` is removed; please use `scrapy.interfaces.ISpiderLoader`.
- Deprecated `CrawlerSettings` class is removed ([issue 3327](#)).
- Deprecated `Settings.overrides` and `Settings.defaults` attributes are removed ([issue 3327](#), [issue 3359](#)).

Other improvements, cleanups

- All Scrapy tests now pass on Windows; Scrapy testing suite is executed in a Windows environment on CI ([issue 3315](#)).
- Python 3.7 support ([issue 3326](#), [issue 3150](#), [issue 3547](#)).
- Testing and CI fixes ([issue 3526](#), [issue 3538](#), [issue 3308](#), [issue 3311](#), [issue 3309](#), [issue 3305](#), [issue 3210](#), [issue 3299](#))
- `scrapy.http.cookies.CookieJar.clear` accepts “domain”, “path” and “name” optional arguments ([issue 3231](#)).
- additional files are included to sdist ([issue 3495](#));
- code style fixes ([issue 3405](#), [issue 3304](#));
- unneeded `.strip()` call is removed ([issue 3519](#));
- `collections.deque` is used to store `MiddlewareManager` methods instead of a list ([issue 3476](#))

7.1.2 Scrapy 1.5.2 (2019-01-22)

- *Security bugfix:* Telnet console extension can be easily exploited by rogue websites POSTing content to `http://localhost:6023`, we haven't found a way to exploit it from Scrapy, but it is very easy to trick a browser to do so and elevates the risk for local development environment.

The fix is backward incompatible, it enables telnet user-password authentication by default with a random generated password. If you can't upgrade right away, please consider setting `TELNET_CONSOLE_PORT` out of its default value.

See [telnet console](#) documentation for more info

- Backport CI build failure under GCE environment due to boto import error.

7.1.3 Scrapy 1.5.1 (2018-07-12)

This is a maintenance release with important bug fixes, but no new features:

- 0(N²) gzip decompression issue which affected Python 3 and PyPy is fixed ([issue 3281](#));
- skipping of TLS validation errors is improved ([issue 3166](#));
- Ctrl-C handling is fixed in Python 3.5+ ([issue 3096](#));
- testing fixes ([issue 3092](#), [issue 3263](#));
- documentation improvements ([issue 3058](#), [issue 3059](#), [issue 3089](#), [issue 3123](#), [issue 3127](#), [issue 3189](#), [issue 3224](#), [issue 3280](#), [issue 3279](#), [issue 3201](#), [issue 3260](#), [issue 3284](#), [issue 3298](#), [issue 3294](#)).

7.1.4 Scrapy 1.5.0 (2017-12-29)

This release brings small new features and improvements across the codebase. Some highlights:

- Google Cloud Storage is supported in FilesPipeline and ImagesPipeline.
- Crawling with proxy servers becomes more efficient, as connections to proxies can be reused now.
- Warnings, exception and logging messages are improved to make debugging easier.
- `scrapy parse` command now allows to set custom request meta via `--meta` argument.
- Compatibility with Python 3.6, PyPy and PyPy3 is improved; PyPy and PyPy3 are now supported officially, by running tests on CI.
- Better default handling of HTTP 308, 522 and 524 status codes.
- Documentation is improved, as usual.

Backward Incompatible Changes

- Scrapy 1.5 drops support for Python 3.3.
- Default Scrapy User-Agent now uses https link to scrapy.org ([issue 2983](#)). **This is technically backward-incompatible**; override `USER_AGENT` if you relied on old value.
- Logging of settings overridden by `custom_settings` is fixed; **this is technically backward-incompatible** because the logger changes from `[scrapy.utils.log]` to `[scrapy.crawler]`. If you're parsing Scrapy logs, please update your log parsers ([issue 1343](#)).
- LinkExtractor now ignores `m4v` extension by default, this is change in behavior.
- 522 and 524 status codes are added to `RETRY_HTTP_CODES` ([issue 2851](#))

New features

- Support `<link>` tags in `Response.follow` (issue 2785)
- Support for `ptpython` REPL (issue 2654)
- Google Cloud Storage support for `FilesPipeline` and `ImagesPipeline` (issue 2923).
- New `--meta` option of the “scrapy parse” command allows to pass additional request.meta (issue 2883)
- Populate spider variable when using `shell.inspect_response` (issue 2812)
- Handle HTTP 308 Permanent Redirect (issue 2844)
- Add 522 and 524 to `RETRY_HTTP_CODES` (issue 2851)
- Log versions information at startup (issue 2857)
- `scrapy.mail.MailSender` now works in Python 3 (it requires Twisted 17.9.0)
- Connections to proxy servers are reused (issue 2743)
- Add template for a downloader middleware (issue 2755)
- Explicit message for `NotImplementedError` when parse callback not defined (issue 2831)
- `CrawlerProcess` got an option to disable installation of root log handler (issue 2921)
- `LinkExtractor` now ignores `m4v` extension by default
- Better log messages for responses over `DOWNLOAD_WARN_SIZE` and `DOWNLOAD_MAX_SIZE` limits (issue 2927)
- Show warning when a URL is put to `Spider.allowed_domains` instead of a domain (issue 2250).

Bug fixes

- Fix logging of settings overridden by `custom_settings`; **this is technically backward-incompatible** because the logger changes from `[scrapy.utils.log]` to `[scrapy.crawler]`, so please update your log parsers if needed (issue 1343)
- Default Scrapy User-Agent now uses https link to scrapy.org (issue 2983). **This is technically backward-incompatible**; override `USER_AGENT` if you relied on old value.
- Fix PyPy and PyPy3 test failures, support them officially (issue 2793, issue 2935, issue 2990, issue 3050, issue 2213, issue 3048)
- Fix DNS resolver when `DNSCACHE_ENABLED=False` (issue 2811)
- Add `cryptography` for Debian Jessie tox test env (issue 2848)
- Add verification to check if Request callback is callable (issue 2766)
- Port `extras/qpsclient.py` to Python 3 (issue 2849)
- Use `getfullargspec` under the scenes for Python 3 to stop `DeprecationWarning` (issue 2862)

- Update deprecated test aliases ([issue 2876](#))
- Fix `SitemapSpider` support for alternate links ([issue 2853](#))

Docs

- Added missing bullet point for the `AUTOTHROTTLE_TARGET_CONCURRENCY` setting. ([issue 2756](#))
- Update Contributing docs, document new support channels ([issue 2762](#), [issue:3038](#))
- Include references to Scrapy subreddit in the docs
- Fix broken links; use <https://> for external links ([issue 2978](#), [issue 2982](#), [issue 2958](#))
- Document `CloseSpider` extension better ([issue 2759](#))
- Use `pymongo.collection.Collection.insert_one()` in MongoDB example ([issue 2781](#))
- Spelling mistake and typos ([issue 2828](#), [issue 2837](#), [issue 2884](#), [issue 2924](#))
- Clarify `CSVFeedSpider.headers` documentation ([issue 2826](#))
- Document `DontCloseSpider` exception and clarify `spider_idle` ([issue 2791](#))
- Update “Releases” section in README ([issue 2764](#))
- Fix rst syntax in `DOWNLOAD_FAIL_ON_DATALOSS` docs ([issue 2763](#))
- Small fix in description of `startproject` arguments ([issue 2866](#))
- Clarify data types in `Response.body` docs ([issue 2922](#))
- Add a note about `request.meta['depth']` to `DepthMiddleware` docs ([issue 2374](#))
- Add a note about `request.meta['dont_merge_cookies']` to `CookiesMiddleware` docs ([issue 2999](#))
- Up-to-date example of project structure ([issue 2964](#), [issue 2976](#))
- A better example of `ItemExporters` usage ([issue 2989](#))
- Document `from_crawler` methods for spider and downloader middlewares ([issue 3019](#))

7.1.5 Scrapy 1.4.0 (2017-05-18)

Scrapy 1.4 does not bring that many breathtaking new features but quite a few handy improvements nonetheless.

Scrapy now supports anonymous FTP sessions with customizable user and password via the new `FTP_USER` and `FTP_PASSWORD` settings. And if you’re using Twisted version 17.1.0 or above, FTP is now available with Python 3.

There’s a new `response.follow` method for creating requests; **it is now a recommended way to create Requests in Scrapy spiders**. This method makes it easier to write correct spiders; `response.follow` has several advantages over creating `scrapy.Request` objects directly:

- it handles relative URLs;
- it works properly with non-ascii URLs on non-UTF8 pages;
- in addition to absolute and relative URLs it supports Selectors; for `<a>` elements it can also extract their href values.

For example, instead of this:

```
for href in response.css('li.page a::attr(href)').extract():
    url = response.urljoin(href)
    yield scrapy.Request(url, self.parse, encoding=response.encoding)
```

One can now write this:

```
for a in response.css('li.page a'):
    yield response.follow(a, self.parse)
```

Link extractors are also improved. They work similarly to what a regular modern browser would do: leading and trailing whitespace are removed from attributes (think `href=" http://example.com"`) when building Link objects. This whitespace-stripping also happens for `action` attributes with `FormRequest`.

Please also note that link extractors do not canonicalize URLs by default anymore. This was puzzling users every now and then, and it's not what browsers do in fact, so we removed that extra transformation on extracted links.

For those of you wanting more control on the `Referer`: header that Scrapy sends when following links, you can set your own `Referrer Policy`. Prior to Scrapy 1.4, the default `RefererMiddleware` would simply and blindly set it to the URL of the response that generated the HTTP request (which could leak information on your URL seeds). By default, Scrapy now behaves much like your regular browser does. And this policy is fully customizable with W3C standard values (or with something really custom of your own if you wish). See [REFERRER_POLICY](#) for details.

To make Scrapy spiders easier to debug, Scrapy logs more stats by default in 1.4: memory usage stats, detailed retry stats, detailed HTTP error code stats. A similar change is that HTTP cache path is also visible in logs now.

Last but not least, Scrapy now has the option to make JSON and XML items more human-readable, with newlines between items and even custom indenting offset, using the new `FEED_EXPORT_INDENT` setting.

Enjoy! (Or read on for the rest of changes in this release.)

Deprecations and Backward Incompatible Changes

- Default to `canonicalize=False` in `scrapy.linkextractors.LinkExtractor` ([issue 2537](#), fixes [issue 1941](#) and [issue 1982](#)): **warning, this is technically backward-incompatible**

- Enable memusage extension by default (issue 2539, fixes issue 2187); **this is technically backward-incompatible** so please check if you have any non-default `MEMUSAGE_***` options set.
- `EDITOR` environment variable now takes precedence over `EDITOR` option defined in `settings.py` (issue 1829); Scrapy default settings no longer depend on environment variables. **This is technically a backward incompatible change.**
- `Spider.make_requests_from_url` is deprecated (issue 1728, fixes issue 1495).

New Features

- Accept proxy credentials in *proxy* request meta key (issue 2526)
- Support brotli-compressed content; requires optional `brotlipy` (issue 2535)
- New *response.follow* shortcut for creating requests (issue 1940)
- Added `flags` argument and attribute to *Request* objects (issue 2047)
- Support Anonymous FTP (issue 2342)
- Added `retry/count`, `retry/max_reached` and `retry/reason_count/<reason>` stats to *RetryMiddleware* (issue 2543)
- Added `httperror/response_ignored_count` and `httperror/response_ignored_status_count/<status>` stats to *HttpErrorMiddleware* (issue 2566)
- Customizable *Referrer policy* in *ReferrerMiddleware* (issue 2306)
- New `data`: URI download handler (issue 2334, fixes issue 2156)
- Log cache directory when HTTP Cache is used (issue 2611, fixes issue 2604)
- Warn users when project contains duplicate spider names (fixes issue 2181)
- `CaselessDict` now accepts `Mapping` instances and not only dicts (issue 2646)
- *Media downloads*, with `FilesPipelines` or `ImagesPipelines`, can now optionally handle HTTP redirects using the new `MEDIA_ALLOW_REDIRECTS` setting (issue 2616, fixes issue 2004)
- Accept non-complete responses from websites using a new `DOWNLOAD_FAIL_ON_DATALOSS` setting (issue 2590, fixes issue 2586)
- Optional pretty-printing of JSON and XML items via `FEED_EXPORT_INDENT` setting (issue 2456, fixes issue 1327)
- Allow dropping fields in `FormRequest.from_response` formdata when `None` value is passed (issue 667)
- Per-request retry times with the new *max_retry_times* meta key (issue 2642)
- `python -m scrapy` as a more explicit alternative to `scrapy` command (issue 2740)

Bug fixes

- LinkExtractor now strips leading and trailing whitespaces from attributes ([issue 2547](#), fixes [issue 1614](#))
- Properly handle whitespaces in action attribute in `FormRequest` ([issue 2548](#))
- Buffer CONNECT response bytes from proxy until all HTTP headers are received ([issue 2495](#), fixes [issue 2491](#))
- FTP downloader now works on Python 3, provided you use Twisted ≥ 17.1 ([issue 2599](#))
- Use body to choose response type after decompressing content ([issue 2393](#), fixes [issue 2145](#))
- Always decompress `Content-Encoding: gzip` at `HttpCompressionMiddleware` stage ([issue 2391](#))
- Respect custom log level in `Spider.custom_settings` ([issue 2581](#), fixes [issue 1612](#))
- ‘make htmlview’ fix for macOS ([issue 2661](#))
- Remove “commands” from the command list ([issue 2695](#))
- Fix duplicate Content-Length header for POST requests with empty body ([issue 2677](#))
- Properly cancel large downloads, i.e. above `DOWNLOAD_MAXSIZE` ([issue 1616](#))
- ImagesPipeline: fixed processing of transparent PNG images with palette ([issue 2675](#))

Cleanups & Refactoring

- Tests: remove temp files and folders ([issue 2570](#)), fixed ProjectUtilsTest on OS X ([issue 2569](#)), use portable pypy for Linux on Travis CI ([issue 2710](#))
- Separate building request from `_requests_to_follow` in CrawlSpider ([issue 2562](#))
- Remove “Python 3 progress” badge ([issue 2567](#))
- Add a couple more lines to `.gitignore` ([issue 2557](#))
- Remove bumpversion prerelease configuration ([issue 2159](#))
- Add codecov.yml file ([issue 2750](#))
- Set context factory implementation based on Twisted version ([issue 2577](#), fixes [issue 2560](#))
- Add omitted `self` arguments in default project middleware template ([issue 2595](#))
- Remove redundant `slot.add_request()` call in ExecutionEngine ([issue 2617](#))
- Catch more specific `os.error` exception in FSFilesStore ([issue 2644](#))
- Change “localhost” test server certificate ([issue 2720](#))
- Remove unused `MEMUSAGE_REPORT` setting ([issue 2576](#))

Documentation

- Binary mode is required for exporters ([issue 2564](#), fixes [issue 2553](#))
- Mention issue with `FormRequest.from_response` due to bug in lxml ([issue 2572](#))
- Use single quotes uniformly in templates ([issue 2596](#))
- Document `ftp_user` and `ftp_password` meta keys ([issue 2587](#))
- Removed section on deprecated `contrib/` ([issue 2636](#))
- Recommend Anaconda when installing Scrapy on Windows ([issue 2477](#), fixes [issue 2475](#))
- FAQ: rewrite note on Python 3 support on Windows ([issue 2690](#))
- Rearrange selector sections ([issue 2705](#))
- Remove `__nonzero__` from `SelectorList` docs ([issue 2683](#))
- Mention how to disable request filtering in documentation of `DUPEFILTER_CLASS` setting ([issue 2714](#))
- Add sphinx_rtd_theme to docs setup readme ([issue 2668](#))
- Open file in text mode in JSON item writer example ([issue 2729](#))
- Clarify `allowed_domains` example ([issue 2670](#))

7.1.6 Scrapy 1.3.3 (2017-03-10)

Bug fixes

- Make `SpiderLoader` raise `ImportError` again by default for missing dependencies and wrong `SPIDER_MODULES`. These exceptions were silenced as warnings since 1.3.0. A new setting is introduced to toggle between warning or exception if needed ; see `SPIDER_LOADER_WARN_ONLY` for details.

7.1.7 Scrapy 1.3.2 (2017-02-13)

Bug fixes

- Preserve request class when converting to/from dicts (`utils.reqser`) ([issue 2510](#)).
- Use consistent selectors for author field in tutorial ([issue 2551](#)).
- Fix TLS compatibility in Twisted 17+ ([issue 2558](#))

7.1.8 Scrapy 1.3.1 (2017-02-08)

New features

- Support 'True' and 'False' string values for boolean settings ([issue 2519](#)); you can now do something like `scrapy crawl myspider -s REDIRECT_ENABLED=False`.
- Support kwargs with `response.xpath()` to use *XPath variables* and ad-hoc namespaces declarations ; this requires at least Parsel v1.1 ([issue 2457](#)).
- Add support for Python 3.6 ([issue 2485](#)).
- Run tests on PyPy (warning: some tests still fail, so PyPy is not supported yet).

Bug fixes

- Enforce `DNS_TIMEOUT` setting ([issue 2496](#)).
- Fix *view* command ; it was a regression in v1.3.0 ([issue 2503](#)).
- Fix tests regarding `*EXPIRES` settings with Files/Images pipelines ([issue 2460](#)).
- Fix name of generated pipeline class when using basic project template ([issue 2466](#)).
- Fix compatibility with Twisted 17+ ([issue 2496](#), [issue 2528](#)).
- Fix `scrapy.Item` inheritance on Python 3.6 ([issue 2511](#)).
- Enforce numeric values for components order in `SPIDER_MIDDLEWARES`, `DOWNLOADER_MIDDLEWARES`, `EXTENSIONS` and `SPIDER_CONTRACTS` ([issue 2420](#)).

Documentation

- Reword Code of Conduct section and upgrade to Contributor Covenant v1.4 ([issue 2469](#)).
- Clarify that passing spider arguments converts them to spider attributes ([issue 2483](#)).
- Document `formid` argument on `FormRequest.from_response()` ([issue 2497](#)).
- Add `.rst` extension to README files ([issue 2507](#)).
- Mention LevelDB cache storage backend ([issue 2525](#)).
- Use `yield` in sample callback code ([issue 2533](#)).
- Add note about HTML entities decoding with `.re()/.re_first()` ([issue 1704](#)).
- Typos ([issue 2512](#), [issue 2534](#), [issue 2531](#)).

Cleanups

- Remove redundant check in `MetaRefreshMiddleware` ([issue 2542](#)).
- Faster checks in `LinkExtractor` for allow/deny patterns ([issue 2538](#)).

- Remove dead code supporting old Twisted versions ([issue 2544](#)).

7.1.9 Scrapy 1.3.0 (2016-12-21)

This release comes rather soon after 1.2.2 for one main reason: it was found out that releases since 0.18 up to 1.2.2 (included) use some backported code from Twisted (`scrapy.xlib.tx.*`), even if newer Twisted modules are available. Scrapy now uses `twisted.web.client` and `twisted.internet.endpoints` directly. (See also cleanups below.)

As it is a major change, we wanted to get the bug fix out quickly while not breaking any projects using the 1.2 series.

New Features

- `MailSender` now accepts single strings as values for `to` and `cc` arguments ([issue 2272](#))
- `scrapy fetch url`, `scrapy shell url` and `fetch(url)` inside scrapy shell now follow HTTP redirections by default ([issue 2290](#)); See *fetch* and *shell* for details.
- `HttpErrorMiddleware` now logs errors with `INFO` level instead of `DEBUG`; this is technically **backward incompatible** so please check your log parsers.
- By default, logger names now use a long-form path, e.g. `[scrapy.extensions.logstats]`, instead of the shorter “top-level” variant of prior releases (e.g. `[scrapy]`); this is **backward incompatible** if you have log parsers expecting the short logger name part. You can switch back to short logger names using `LOG_SHORT_NAMES` set to `True`.

Dependencies & Cleanups

- Scrapy now requires Twisted `>= 13.1` which is the case for many Linux distributions already.
- As a consequence, we got rid of `scrapy.xlib.tx.*` modules, which copied some of Twisted code for users stuck with an “old” Twisted version
- `ChunkedTransferMiddleware` is deprecated and removed from the default downloader middlewares.

7.1.10 Scrapy 1.2.3 (2017-03-03)

- Packaging fix: disallow unsupported Twisted versions in `setup.py`

7.1.11 Scrapy 1.2.2 (2016-12-06)

Bug fixes

- Fix a cryptic traceback when a pipeline fails on `open_spider()` ([issue 2011](#))

- Fix embedded IPython shell variables (fixing [issue 396](#) that re-appeared in 1.2.0, fixed in [issue 2418](#))
- A couple of patches when dealing with robots.txt:
 - handle (non-standard) relative sitemap URLs ([issue 2390](#))
 - handle non-ASCII URLs and User-Agents in Python 2 ([issue 2373](#))

Documentation

- Document "download_latency" key in Request's meta dict ([issue 2033](#))
- Remove page on (deprecated & unsupported) Ubuntu packages from ToC ([issue 2335](#))
- A few fixed typos ([issue 2346](#), [issue 2369](#), [issue 2369](#), [issue 2380](#)) and clarifications ([issue 2354](#), [issue 2325](#), [issue 2414](#))

Other changes

- Advertize [conda-forge](#) as Scrapy's official conda channel ([issue 2387](#))
- More helpful error messages when trying to use `.css()` or `.xpath()` on non-Text Responses ([issue 2264](#))
- `startproject` command now generates a sample `middlewares.py` file ([issue 2335](#))
- Add more dependencies' version info in `scrapy version` verbose output ([issue 2404](#))
- Remove all `*.pyc` files from source distribution ([issue 2386](#))

7.1.12 Scrapy 1.2.1 (2016-10-21)

Bug fixes

- Include OpenSSL's more permissive default ciphers when establishing TLS/SSL connections ([issue 2314](#)).
- Fix "Location" HTTP header decoding on non-ASCII URL redirects ([issue 2321](#)).

Documentation

- Fix `JsonWriterPipeline` example ([issue 2302](#)).
- Various notes: [issue 2330](#) on spider names, [issue 2329](#) on middleware methods processing order, [issue 2327](#) on getting multi-valued HTTP headers as lists.

Other changes

- Removed `www.` from `start_urls` in built-in spider templates (issue 2299).

7.1.13 Scrapy 1.2.0 (2016-10-03)

New Features

- New `FEED_EXPORT_ENCODING` setting to customize the encoding used when writing items to a file. This can be used to turn off `\uXXXX` escapes in JSON output. This is also useful for those wanting something else than UTF-8 for XML or CSV output (issue 2034).
- `startproject` command now supports an optional destination directory to override the default one based on the project name (issue 2005).
- New `SCHEDULER_DEBUG` setting to log requests serialization failures (issue 1610).
- JSON encoder now supports serialization of `set` instances (issue 2058).
- Interpret `application/json-amazonui-streaming` as `TextResponse` (issue 1503).
- `scrapy` is imported by default when using shell tools (*`shell`*, *`inspect_response`*) (issue 2248).

Bug fixes

- `DefaultRequestHeaders` middleware now runs before `UserAgent` middleware (issue 2088). **Warning: this is technically backward incompatible**, though we consider this a bug fix.
- HTTP cache extension and plugins that use the `.scrapy` data directory now work outside projects (issue 1581). **Warning: this is technically backward incompatible**, though we consider this a bug fix.
- `Selector` does not allow passing both `response` and `text` anymore (issue 2153).
- Fixed logging of wrong callback name with `scrapy parse` (issue 2169).
- Fix for an odd gzip decompression bug (issue 1606).
- Fix for selected callbacks when using `CrawlSpider` with *`scrapy parse`* (issue 2225).
- Fix for invalid JSON and XML files when spider yields no items (issue 872).
- Implement `flush()` for `StreamLogger` avoiding a warning in logs (issue 2125).

Refactoring

- `canonicalize_url` has been moved to `w3lib.url` (issue 2168).

Tests & Requirements

Scrapy's new requirements baseline is Debian 8 "Jessie". It was previously Ubuntu 12.04 Precise. What this means in practice is that we run continuous integration tests with these (main) packages versions at a minimum: Twisted 14.0, pyOpenSSL 0.14, lxml 3.4.

Scrapy may very well work with older versions of these packages (the code base still has switches for older Twisted versions for example) but it is not guaranteed (because it's not tested anymore).

Documentation

- Grammar fixes: [issue 2128](#), [issue 1566](#).
- Download stats badge removed from README ([issue 2160](#)).
- New scrapy *architecture diagram* ([issue 2165](#)).
- Updated **Response** parameters documentation ([issue 2197](#)).
- Reworded misleading *RANDOMIZE_DOWNLOAD_DELAY* description ([issue 2190](#)).
- Add StackOverflow as a support channel ([issue 2257](#)).

7.1.14 Scrapy 1.1.4 (2017-03-03)

- Packaging fix: disallow unsupported Twisted versions in setup.py

7.1.15 Scrapy 1.1.3 (2016-09-22)

Bug fixes

- Class attributes for subclasses of `ImagesPipeline` and `FilesPipeline` work as they did before 1.1.1 ([issue 2243](#), fixes [issue 2198](#))

Documentation

- *Overview* and *tutorial* rewritten to use <http://toscrape.com> websites ([issue 2236](#), [issue 2249](#), [issue 2252](#)).

7.1.16 Scrapy 1.1.2 (2016-08-18)

Bug fixes

- Introduce a missing `IMAGES_STORE_S3_ACL` setting to override the default ACL policy in `ImagesPipeline` when uploading images to S3 (note that default ACL policy is “private” – instead of “public-read” – since Scrapy 1.1.0)
- `IMAGES_EXPIRES` default value set back to 90 (the regression was introduced in 1.1.1)

7.1.17 Scrapy 1.1.1 (2016-07-13)

Bug fixes

- Add “Host” header in CONNECT requests to HTTPS proxies (issue 2069)
- Use response body when choosing response class (issue 2001, fixes issue 2000)
- Do not fail on canonicalizing URLs with wrong netlocs (issue 2038, fixes issue 2010)
- a few fixes for `HttpCompressionMiddleware` (and `SitemapSpider`):
 - Do not decode HEAD responses (issue 2008, fixes issue 1899)
 - Handle charset parameter in gzip Content-Type header (issue 2050, fixes issue 2049)
 - Do not decompress gzip octet-stream responses (issue 2065, fixes issue 2063)
- Catch (and ignore with a warning) exception when verifying certificate against IP-address hosts (issue 2094, fixes issue 2092)
- Make `FilesPipeline` and `ImagesPipeline` backward compatible again regarding the use of legacy class attributes for customization (issue 1989, fixes issue 1985)

New features

- Enable genspider command outside project folder (issue 2052)
- Retry HTTPS CONNECT `TunnelError` by default (issue 1974)

Documentation

- `FEED_TEMPDIR` setting at lexicographical position (commit 9b3c72c)
- Use idiomatic `.extract_first()` in overview (issue 1994)
- Update years in copyright notice (commit c2c8036)
- Add information and example on errbacks (issue 1995)
- Use “url” variable in downloader middleware example (issue 2015)
- Grammar fixes (issue 2054, issue 2120)

- New FAQ entry on using BeautifulSoup in spider callbacks ([issue 2048](#))
- Add notes about scrapy not working on Windows with Python 3 ([issue 2060](#))
- Encourage complete titles in pull requests ([issue 2026](#))

Tests

- Upgrade py.test requirement on Travis CI and Pin pytest-cov to 2.2.1 ([issue 2095](#))

7.1.18 Scrapy 1.1.0 (2016-05-11)

This 1.1 release brings a lot of interesting features and bug fixes:

- Scrapy 1.1 has beta Python 3 support (requires Twisted ≥ 15.5). See *Beta Python 3 Support* for more details and some limitations.
- Hot new features:
 - Item loaders now support nested loaders ([issue 1467](#)).
 - `FormRequest.from_response` improvements ([issue 1382](#), [issue 1137](#)).
 - Added setting `AUTOTHROTTLLE_TARGET_CONCURRENCY` and improved AutoThrottle docs ([issue 1324](#)).
 - Added `response.text` to get body as unicode ([issue 1730](#)).
 - Anonymous S3 connections ([issue 1358](#)).
 - Deferreds in downloader middlewares ([issue 1473](#)). This enables better robots.txt handling ([issue 1471](#)).
 - HTTP caching now follows RFC2616 more closely, added settings `HTTPCACHE_ALWAYS_STORE` and `HTTPCACHE_IGNORE_RESPONSE_CACHE_CONTROLS` ([issue 1151](#)).
 - Selectors were extracted to the `parsel` library ([issue 1409](#)). This means you can use Scrapy Selectors without Scrapy and also upgrade the selectors engine without needing to upgrade Scrapy.
 - HTTPS downloader now does TLS protocol negotiation by default, instead of forcing TLS 1.0. You can also set the SSL/TLS method using the new `DOWNLOADER_CLIENT_TLS_METHOD`.
- These bug fixes may require your attention:
 - Don't retry bad requests (HTTP 400) by default ([issue 1289](#)). If you need the old behavior, add 400 to `RETRY_HTTP_CODES`.
 - Fix shell files argument handling ([issue 1710](#), [issue 1550](#)). If you try `scrapy shell index.html` it will try to load the URL `http://index.html`, use `scrapy shell ./index.html` to load a local file.

- Robots.txt compliance is now enabled by default for newly-created projects ([issue 1724](#)). Scrapy will also wait for robots.txt to be downloaded before proceeding with the crawl ([issue 1735](#)). If you want to disable this behavior, update `ROBOTSTXT_OBEY` in `settings.py` file after creating a new project.
- Exporters now work on unicode, instead of bytes by default ([issue 1080](#)). If you use `PythonItemExporter`, you may want to update your code to disable binary mode which is now deprecated.
- Accept XML node names containing dots as valid ([issue 1533](#)).
- When uploading files or images to S3 (with `FilesPipeline` or `ImagesPipeline`), the default ACL policy is now “private” instead of “public” **Warning: backward incompatible!**. You can use `FILES_STORE_S3_ACL` to change it.
- We’ve reimplemented `canonicalize_url()` for more correct output, especially for URLs with non-ASCII characters ([issue 1947](#)). This could change link extractors output compared to previous scrapy versions. This may also invalidate some cache entries you could still have from pre-1.1 runs.
Warning: backward incompatible!

Keep reading for more details on other improvements and bug fixes.

Beta Python 3 Support

We have been [hard at work](#) to make Scrapy run on Python 3. As a result, now you can run spiders on Python 3.3, 3.4 and 3.5 (Twisted ≥ 15.5 required). Some features are still missing (and some may never be ported).

Almost all builtin extensions/middlewares are expected to work. However, we are aware of some limitations in Python 3:

- Scrapy does not work on Windows with Python 3
- Sending emails is not supported
- FTP download handler is not supported
- Telnet console is not supported

Additional New Features and Enhancements

- Scrapy now has a [Code of Conduct](#) ([issue 1681](#)).
- Command line tool now has completion for zsh ([issue 934](#)).
- Improvements to `scrapy shell`:
 - Support for bpython and configure preferred Python shell via `SCRAPY_PYTHON_SHELL` ([issue 1100](#), [issue 1444](#)).

- Support URLs without scheme ([issue 1498](#)) **Warning: backward incompatible!**
 - Bring back support for relative file path ([issue 1710](#), [issue 1550](#)).
- Added `MEMUSAGE_CHECK_INTERVAL_SECONDS` setting to change default check interval ([issue 1282](#)).
- Download handlers are now lazy-loaded on first request using their scheme ([issue 1390](#), [issue 1421](#)).
- HTTPS download handlers do not force TLS 1.0 anymore; instead, OpenSSL's `SSLv23_method()`/`TLS_method()` is used allowing to try negotiating with the remote hosts the highest TLS protocol version it can ([issue 1794](#), [issue 1629](#)).
- `RedirectMiddleware` now skips the status codes from `handle_httpstatus_list` on spider attribute or in `Request`'s `meta` key ([issue 1334](#), [issue 1364](#), [issue 1447](#)).
- Form submission:
 - now works with `<button>` elements too ([issue 1469](#)).
 - an empty string is now used for submit buttons without a value ([issue 1472](#))
- Dict-like settings now have per-key priorities ([issue 1135](#), [issue 1149](#) and [issue 1586](#)).
- Sending non-ASCII emails ([issue 1662](#))
- `CloseSpider` and `SpiderState` extensions now get disabled if no relevant setting is set ([issue 1723](#), [issue 1725](#)).
- Added method `ExecutionEngine.close` ([issue 1423](#)).
- Added method `CrawlerRunner.create_crawler` ([issue 1528](#)).
- Scheduler priority queue can now be customized via `SCHEDULER_PRIORITY_QUEUE` ([issue 1822](#)).
- `.pps` links are now ignored by default in link extractors ([issue 1835](#)).
- temporary data folder for FTP and S3 feed storages can be customized using a new `FEED_TEMPDIR` setting ([issue 1847](#)).
- `FilesPipeline` and `ImagesPipeline` settings are now instance attributes instead of class attributes, enabling spider-specific behaviors ([issue 1891](#)).
- `JsonItemExporter` now formats opening and closing square brackets on their own line (first and last lines of output file) ([issue 1950](#)).
- If available, `botocore` is used for `S3FeedStorage`, `S3DownloadHandler` and `S3FilesStore` ([issue 1761](#), [issue 1883](#)).
- Tons of documentation updates and related fixes ([issue 1291](#), [issue 1302](#), [issue 1335](#), [issue 1683](#), [issue 1660](#), [issue 1642](#), [issue 1721](#), [issue 1727](#), [issue 1879](#)).
- Other refactoring, optimizations and cleanup ([issue 1476](#), [issue 1481](#), [issue 1477](#), [issue 1315](#), [issue 1290](#), [issue 1750](#), [issue 1881](#)).

Deprecations and Removals

- Added `to_bytes` and `to_unicode`, deprecated `str_to_unicode` and `unicode_to_str` functions (issue 778).
- `binary_is_text` is introduced, to replace use of `isbinarytext` (but with inverse return value) (issue 1851)
- The `optional_features` set has been removed (issue 1359).
- The `--lsprof` command line option has been removed (issue 1689). **Warning: backward incompatible**, but doesn't break user code.
- The following datatypes were deprecated (issue 1720):
 - `scrapy.utils.datatypes.MultiValueDictKeyError`
 - `scrapy.utils.datatypes.MultiValueDict`
 - `scrapy.utils.datatypes.SiteNode`
- The previously bundled `scrapy.xlib.pydispatch` library was deprecated and replaced by `pydispatcher`.

Relocations

- `telnetconsole` was relocated to `extensions/` (issue 1524).
 - Note: `telnet` is not enabled on Python 3 (<https://github.com/scrapy/scrapy/pull/1524#issuecomment-146985595>)

Bugfixes

- Scrapy does not retry requests that got a HTTP 400 Bad Request response anymore (issue 1289). **Warning: backward incompatible!**
- Support empty password for `http_proxy` config (issue 1274).
- Interpret `application/x-json` as `TextResponse` (issue 1333).
- Support `link rel` attribute with multiple values (issue 1201).
- Fixed `scrapy.http.FormRequest.from_response` when there is a `<base>` tag (issue 1564).
- Fixed `TEMPLATES_DIR` handling (issue 1575).
- Various `FormRequest` fixes (issue 1595, issue 1596, issue 1597).
- Makes `_monkeypatches` more robust (issue 1634).
- Fixed bug on `XMLItemExporter` with non-string fields in items (issue 1738).
- Fixed `startproject` command in OS X (issue 1635).

- Fixed PythonItemExporter and CSVExporter for non-string item types ([issue 1737](#)).
- Various logging related fixes ([issue 1294](#), [issue 1419](#), [issue 1263](#), [issue 1624](#), [issue 1654](#), [issue 1722](#), [issue 1726](#) and [issue 1303](#)).
- Fixed bug in `utils.template.render_templatefile()` ([issue 1212](#)).
- sitemaps extraction from `robots.txt` is now case-insensitive ([issue 1902](#)).
- HTTPS+CONNECT tunnels could get mixed up when using multiple proxies to same remote host ([issue 1912](#)).

7.1.19 Scrapy 1.0.7 (2017-03-03)

- Packaging fix: disallow unsupported Twisted versions in setup.py

7.1.20 Scrapy 1.0.6 (2016-05-04)

- FIX: RetryMiddleware is now robust to non-standard HTTP status codes ([issue 1857](#))
- FIX: Filestorage HTTP cache was checking wrong modified time ([issue 1875](#))
- DOC: Support for Sphinx 1.4+ ([issue 1893](#))
- DOC: Consistency in selectors examples ([issue 1869](#))

7.1.21 Scrapy 1.0.5 (2016-02-04)

- FIX: [Backport] Ignore bogus links in LinkExtractors (fixes [issue 907](#), [commit 108195e](#))
- TST: Changed buildbot makefile to use ‘pytest’ ([commit 1f3d90a](#))
- DOC: Fixed typos in tutorial and media-pipeline ([commit 808a9ea](#) and [commit 803bd87](#))
- DOC: Add AjaxCrawlMiddleware to DOWNLOADER_MIDDLEWARES_BASE in settings docs ([commit aa94121](#))

7.1.22 Scrapy 1.0.4 (2015-12-30)

- Ignoring xlib/tx folder, depending on Twisted version. ([commit 7dfa979](#))
- Run on new travis-ci infra ([commit 6e42f0b](#))
- Spelling fixes ([commit 823a1cc](#))
- escape nodename in xmliter regex ([commit da3c155](#))
- test xml nodename with dots ([commit 4418fc3](#))
- TST don’ t use broken Pillow version in tests ([commit a55078c](#))

- disable log on version command. closes #1426 (commit 86fc330)
- disable log on startproject command (commit db4c9fe)
- Add PyPI download stats badge (commit df2b944)
- don't run tests twice on Travis if a PR is made from a scrapy/scrapy branch (commit a83ab41)
- Add Python 3 porting status badge to the README (commit 73ac80d)
- fixed RFPDupeFilter persistence (commit 97d080e)
- TST a test to show that dupefilter persistence is not working (commit 97f2fb3)
- explicit close file on file:// scheme handler (commit d9b4850)
- Disable dupefilter in shell (commit c0d0734)
- DOC: Add captions to toctrees which appear in sidebar (commit aa239ad)
- DOC Removed pywin32 from install instructions as it's already declared as dependency. (commit 10eb400)
- Added installation notes about using Conda for Windows and other OSes. (commit 1c3600a)
- Fixed minor grammar issues. (commit 7f4ddd5)
- fixed a typo in the documentation. (commit b71f677)
- Version 1 now exists (commit 5456c0e)
- fix another invalid xpath error (commit 0a1366e)
- fix ValueError: Invalid XPath: //div/[id=" not-exists"]/text() on selectors.rst (commit ca8d60f)
- Typos corrections (commit 7067117)
- fix typos in downloader-middleware.rst and exceptions.rst, middleware -> middleware (commit 32f115c)
- Add note to ubuntu install section about debian compatibility (commit 23fda69)
- Replace alternative OSX install workaround with virtualenv (commit 98b63ee)
- Reference Homebrew's homepage for installation instructions (commit 1925db1)
- Add oldest supported tox version to contributing docs (commit 5d10d6d)
- Note in install docs about pip being already included in python>=2.7.9 (commit 85c980e)
- Add non-python dependencies to Ubuntu install section in the docs (commit fbd010d)
- Add OS X installation section to docs (commit d8f4cba)
- DOC(ENH): specify path to rtd theme explicitly (commit de73b1a)
- minor: scrapy.Spider docs grammar (commit 1ddcc7b)
- Make common practices sample code match the comments (commit 1b85bcf)

- nextcall repetitive calls (heartbeats). ([commit 55f7104](#))
- Backport fix compatibility with Twisted 15.4.0 ([commit b262411](#))
- pin pytest to 2.7.3 ([commit a6535c2](#))
- Merge pull request #1512 from mgedmin/patch-1 ([commit 8876111](#))
- Merge pull request #1513 from mgedmin/patch-2 ([commit 5d4daf8](#))
- Typo ([commit f8d0682](#))
- Fix list formatting ([commit 5f83a93](#))
- fix scrapy queue tests after recent changes to queuelib ([commit 3365c01](#))
- Merge pull request #1475 from rweindl/patch-1 ([commit 2d688cd](#))
- Update tutorial.rst ([commit fbc1f25](#))
- Merge pull request #1449 from rhoekman/patch-1 ([commit 7d6538c](#))
- Small grammatical change ([commit 8752294](#))
- Add openssl version to version command ([commit 13c45ac](#))

7.1.23 Scrapy 1.0.3 (2015-08-11)

- add service_identity to scrapy install_requires ([commit cbc2501](#))
- Workaround for travis#296 ([commit 66af9cd](#))

7.1.24 Scrapy 1.0.2 (2015-08-06)

- Twisted 15.3.0 does not raises PicklingError serializing lambda functions ([commit b04dd7d](#))
- Minor method name fix ([commit 6f85c7f](#))
- minor: scrapy.Spider grammar and clarity ([commit 9c9d2e0](#))
- Put a blurb about support channels in CONTRIBUTING ([commit c63882b](#))
- Fixed typos ([commit a9ae7b0](#))
- Fix doc reference. ([commit 7c8a4fe](#))

7.1.25 Scrapy 1.0.1 (2015-07-01)

- Unquote request path before passing to FTPClient, it already escape paths ([commit cc00ad2](#))
- include tests/ to source distribution in MANIFEST.in ([commit eca227e](#))
- DOC Fix SelectJmes documentation ([commit b8567bc](#))

- DOC Bring Ubuntu and Archlinux outside of Windows subsection ([commit 392233f](#))
- DOC remove version suffix from ubuntu package ([commit 5303c66](#))
- DOC Update release date for 1.0 ([commit c89fa29](#))

7.1.26 Scrapy 1.0.0 (2015-06-19)

You will find a lot of new features and bugfixes in this major release. Make sure to check our updated [overview](#) to get a glance of some of the changes, along with our brushed [tutorial](#).

Support for returning dictionaries in spiders

Declaring and returning Scrapy Items is no longer necessary to collect the scraped data from your spider, you can now return explicit dictionaries instead.

Classic version

```
class MyItem(scrapy.Item):
    url = scrapy.Field()

class MySpider(scrapy.Spider):
    def parse(self, response):
        return MyItem(url=response.url)
```

New version

```
class MySpider(scrapy.Spider):
    def parse(self, response):
        return {'url': response.url}
```

Per-spider settings (GSoC 2014)

Last Google Summer of Code project accomplished an important redesign of the mechanism used for populating settings, introducing explicit priorities to override any given setting. As an extension of that goal, we included a new level of priority for settings that act exclusively for a single spider, allowing them to redefine project settings.

Start using it by defining a `custom_settings` class variable in your spider:

```
class MySpider(scrapy.Spider):
    custom_settings = {
        "DOWNLOAD_DELAY": 5.0,
```

(下页继续)

(续上页)

```
"RETRY_ENABLED": False,  
}
```

Read more about settings population: [设置](#)

Python Logging

Scrapy 1.0 has moved away from Twisted logging to support Python built in's as default logging system. We're maintaining backward compatibility for most of the old custom interface to call logging functions, but you'll get warnings to switch to the Python logging API entirely.

Old version

```
from scrapy import log  
log.msg('MESSAGE', log.INFO)
```

New version

```
import logging  
logging.info('MESSAGE')
```

Logging with spiders remains the same, but on top of the `log()` method you'll have access to a custom *logger* created for the spider to issue log events:

```
class MySpider(scrapy.Spider):  
    def parse(self, response):  
        self.logger.info('Response received')
```

Read more in the logging documentation: [日志](#)

Crawler API refactoring (GSoC 2014)

Another milestone for last Google Summer of Code was a refactoring of the internal API, seeking a simpler and easier usage. Check new core interface in: [核心 API](#)

A common situation where you will face these changes is while running Scrapy from scripts. Here's a quick example of how to run a Spider manually with the new API:

```
from scrapy.crawler import CrawlerProcess  
  
process = CrawlerProcess({  
    'USER_AGENT': 'Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 5.1)'
```

(下页继续)

(续上页)

```
}  
process.crawl(MySpider)  
process.start()
```

Bear in mind this feature is still under development and its API may change until it reaches a stable status.

See more examples for scripts running Scrapy: [常见实践](#)

Module Relocations

There' s been a large rearrangement of modules trying to improve the general structure of Scrapy. Main changes were separating various subpackages into new projects and dissolving both `scrapy.contrib` and `scrapy.contrib_exp` into top level packages. Backward compatibility was kept among internal relocations, while importing deprecated modules expect warnings indicating their new place.

Full list of relocations

Outsourced packages

注解: These extensions went through some minor changes, e.g. some setting names were changed. Please check the documentation in each new repository to get familiar with the new usage.

Old location	New location
<code>scrapy.commands.deploy</code>	<code>scrapyd-client</code> (See other alternatives here: 部署爬虫器)
<code>scrapy.contrib.djangoitem</code>	<code>scrapy-djangoitem</code>
<code>scrapy.webservice</code>	<code>scrapy-jsonrpc</code>

`scrapy.contrib_exp` and `scrapy.contrib` dissolutions

Old location	New location
scrapy.contrib_exp.downloadermiddleware.decompression	scrapy.downloadermiddlewares.decompression
scrapy.contrib_exp.iterators	scrapy.utils.iterators
scrapy.contrib.downloadermiddleware	scrapy.downloadermiddlewares
scrapy.contrib.exporter	scrapy.exporters
scrapy.contrib.linkextractors	scrapy.linkextractors
scrapy.contrib.loader	scrapy.loader
scrapy.contrib.loader.processor	scrapy.loader.processors
scrapy.contrib.pipeline	scrapy.pipelines
scrapy.contrib.spidermiddleware	scrapy.spidermiddlewares
scrapy.contrib.spiders	scrapy.spiders
<ul style="list-style-type: none"> • scrapy.contrib.closespider • scrapy.contrib.corestats • scrapy.contrib.debug • scrapy.contrib.feedexport • scrapy.contrib.httpcache • scrapy.contrib.logstats • scrapy.contrib.memdebug • scrapy.contrib.memusage • scrapy.contrib.spiderstate • scrapy.contrib.statsmailer • scrapy.contrib.throttle 	scrapy.extensions.*

Plural renames and Modules unification

Old location	New location
scrapy.command	scrapy.commands
scrapy.dupefilter	scrapy.dupefilters
scrapy.linkextractor	scrapy.linkextractors
scrapy.spider	scrapy.spiders
scrapy.squeue	scrapy.squeues
scrapy.statscol	scrapy.statscollectors
scrapy.utils.decorator	scrapy.utils.decorators

Class renames

Old location	New location
scrapy.spidermanager.SpiderManager	scrapy.spiderloader.SpiderLoader

Settings renames

Old location	New location
SPIDER_MANAGER_CLASS	SPIDER_LOADER_CLASS

Changelog

New Features and Enhancements

- Python logging ([issue 1060](#), [issue 1235](#), [issue 1236](#), [issue 1240](#), [issue 1259](#), [issue 1278](#), [issue 1286](#))
- FEED_EXPORT_FIELDS option ([issue 1159](#), [issue 1224](#))
- Dns cache size and timeout options ([issue 1132](#))
- support namespace prefix in xmliter_lxml ([issue 963](#))
- Reactor threadpool max size setting ([issue 1123](#))
- Allow spiders to return dicts. ([issue 1081](#))
- Add Response.urljoin() helper ([issue 1086](#))
- look in ~/.config/scrapy.cfg for user config ([issue 1098](#))
- handle TLS SNI ([issue 1101](#))
- Selectorlist extract first ([issue 624](#), [issue 1145](#))
- Added JmesSelect ([issue 1016](#))
- add gzip compression to filesystem http cache backend ([issue 1020](#))
- CSS support in link extractors ([issue 983](#))
- httpcache dont_cache meta #19 #689 ([issue 821](#))
- add signal to be sent when request is dropped by the scheduler ([issue 961](#))
- avoid download large response ([issue 946](#))
- Allow to specify the quotechar in CSVFeedSpider ([issue 882](#))
- Add referer to “Spider error processing” log message ([issue 795](#))
- process robots.txt once ([issue 896](#))
- GSoC Per-spider settings ([issue 854](#))
- Add project name validation ([issue 817](#))
- GSoC API cleanup ([issue 816](#), [issue 1128](#), [issue 1147](#), [issue 1148](#), [issue 1156](#), [issue 1185](#), [issue 1187](#), [issue 1258](#), [issue 1268](#), [issue 1276](#), [issue 1285](#), [issue 1284](#))
- Be more responsive with IO operations ([issue 1074](#) and [issue 1075](#))

- Do leveldb compaction for httpcache on closing ([issue 1297](#))

Deprecations and Removals

- Deprecate htmlparser link extractor ([issue 1205](#))
- remove deprecated code from FeedExporter ([issue 1155](#))
- a leftover for .15 compatibility ([issue 925](#))
- drop support for CONCURRENT_REQUESTS_PER_SPIDER ([issue 895](#))
- Drop old engine code ([issue 911](#))
- Deprecate SgmlLinkExtractor ([issue 777](#))

Relocations

- Move exporters/___init___py to exporters.py ([issue 1242](#))
- Move base classes to their packages ([issue 1218](#), [issue 1233](#))
- Module relocation ([issue 1181](#), [issue 1210](#))
- rename SpiderManager to SpiderLoader ([issue 1166](#))
- Remove djangoitem ([issue 1177](#))
- remove scrapy deploy command ([issue 1102](#))
- dissolve contrib_exp ([issue 1134](#))
- Deleted bin folder from root, fixes #913 ([issue 914](#))
- Remove jsonrpc based webservice ([issue 859](#))
- Move Test cases under project root dir ([issue 827](#), [issue 841](#))
- Fix backward incompatibility for relocated paths in settings ([issue 1267](#))

Documentation

- CrawlerProcess documentation ([issue 1190](#))
- Favoring web scraping over screen scraping in the descriptions ([issue 1188](#))
- Some improvements for Scrapy tutorial ([issue 1180](#))
- Documenting Files Pipeline together with Images Pipeline ([issue 1150](#))
- deployment docs tweaks ([issue 1164](#))
- Added deployment section covering scrapyd-deploy and shub ([issue 1124](#))
- Adding more settings to project template ([issue 1073](#))
- some improvements to overview page ([issue 1106](#))
- Updated link in docs/topics/architecture.rst ([issue 647](#))

- DOC reorder topics ([issue 1022](#))
- updating list of Request.meta special keys ([issue 1071](#))
- DOC document download_timeout ([issue 898](#))
- DOC simplify extension docs ([issue 893](#))
- Leaks docs ([issue 894](#))
- DOC document from_crawler method for item pipelines ([issue 904](#))
- Spider_error doesn't support deferreds ([issue 1292](#))
- Corrections & Sphinx related fixes ([issue 1220](#), [issue 1219](#), [issue 1196](#), [issue 1172](#), [issue 1171](#), [issue 1169](#), [issue 1160](#), [issue 1154](#), [issue 1127](#), [issue 1112](#), [issue 1105](#), [issue 1041](#), [issue 1082](#), [issue 1033](#), [issue 944](#), [issue 866](#), [issue 864](#), [issue 796](#), [issue 1260](#), [issue 1271](#), [issue 1293](#), [issue 1298](#))

Bugfixes

- Item multi inheritance fix ([issue 353](#), [issue 1228](#))
- ItemLoader.load_item: iterate over copy of fields ([issue 722](#))
- Fix Unhandled error in Deferred (RobotsTxtMiddleware) ([issue 1131](#), [issue 1197](#))
- Force to read DOWNLOAD_TIMEOUT as int ([issue 954](#))
- scrapy.utils.misc.load_object should print full traceback ([issue 902](#))
- Fix bug for “.local” host name ([issue 878](#))
- Fix for Enabled extensions, middlewares, pipelines info not printed anymore ([issue 879](#))
- fix dont_merge_cookies bad behaviour when set to false on meta ([issue 846](#))

Python 3 In Progress Support

- disable scrapy.telnet if twisted.conch is not available ([issue 1161](#))
- fix Python 3 syntax errors in ajaxcrawl.py ([issue 1162](#))
- more python3 compatibility changes for urllib ([issue 1121](#))
- assertItemsEqual was renamed to assertCountEqual in Python 3. ([issue 1070](#))
- Import unittest.mock if available. ([issue 1066](#))
- updated deprecated cgi.parse_qs to use six's parse_qs ([issue 909](#))
- Prevent Python 3 port regressions ([issue 830](#))
- PY3: use MutableMapping for python 3 ([issue 810](#))
- PY3: use six.BytesIO and six.moves.cStringIO ([issue 803](#))
- PY3: fix xmlrpclib and email imports ([issue 801](#))
- PY3: use six for robotparser and urlparse ([issue 800](#))

- PY3: use `six.iterkeys`, `six.iteritems`, and `tempfile` (issue 799)
- PY3: fix `has_key` and use `six.moves.configparser` (issue 798)
- PY3: use `six.moves.cPickle` (issue 797)
- PY3 make it possible to run some tests in Python3 (issue 776)

Tests

- remove unnecessary lines from `py3-ignores` (issue 1243)
- Fix remaining warnings from `pytest` while collecting tests (issue 1206)
- Add docs build to travis (issue 1234)
- TST don't collect tests from deprecated modules. (issue 1165)
- install `service_identity` package in tests to prevent warnings (issue 1168)
- Fix deprecated settings API in tests (issue 1152)
- Add test for webclient with POST method and no body given (issue 1089)
- `py3-ignores.txt` supports comments (issue 1044)
- modernize some of the asserts (issue 835)
- `selector.__repr__` test (issue 779)

Code refactoring

- `CSVFeedSpider` cleanup: use `iterate_spider_output` (issue 1079)
- remove unnecessary check from `scrapy.utils.spider.iter_spider_output` (issue 1078)
- `Pydispatch` pep8 (issue 992)
- Removed unused `'load=False'` parameter from `walk_modules()` (issue 871)
- For consistency, use `job_dir` helper in `SpiderState` extension. (issue 805)
- rename `"sflo"` local variables to less cryptic `"log_observer"` (issue 775)

7.1.27 Scrapy 0.24.6 (2015-04-20)

- encode invalid xpath with `unicode_escape` under PY2 (commit 07cb3e5)
- fix IPython shell scope issue and load IPython user config (commit 2c8e573)
- Fix small typo in the docs (commit d694019)
- Fix small typo (commit f92fa83)
- Converted `sel.xpath()` calls to `response.xpath()` in `Extracting the data` (commit c2c6d15)

7.1.28 Scrapy 0.24.5 (2015-02-25)

- Support new `__getEndpoint` Agent signatures on Twisted 15.0.0 ([commit 540b9bc](#))
- DOC a couple more references are fixed ([commit b4c454b](#))
- DOC fix a reference ([commit e3c1260](#))
- `t.i.b.ThreadedResolver` is now a new-style class ([commit 9e13f42](#))
- `S3DownloadHandler`: fix auth for requests with quoted paths/query params ([commit cdb9a0b](#))
- fixed the variable types in mailsender documentation ([commit bb3a848](#))
- Reset `items_scraped` instead of `item_count` ([commit edb07a4](#))
- Tentative attention message about what document to read for contributions ([commit 7ee6f7a](#))
- `mitmproxy` 0.10.1 needs `netlib` 0.10.1 too ([commit 874fcdd](#))
- pin `mitmproxy` 0.10.1 as >0.11 does not work with tests ([commit c6b21f0](#))
- Test the parse command locally instead of against an external url ([commit c3a6628](#))
- Patches Twisted issue while closing the connection pool on `HTTPDownloadHandler` ([commit d0bf957](#))
- Updates documentation on dynamic item classes. ([commit eeb589a](#))
- Merge pull request #943 from Lazar-T/patch-3 ([commit 5fdab02](#))
- typo ([commit b0ae199](#))
- `pywin32` is required by Twisted. closes #937 ([commit 5cb0cfb](#))
- Update `install.rst` ([commit 781286b](#))
- Merge pull request #928 from Lazar-T/patch-1 ([commit b415d04](#))
- comma instead of fullstop ([commit 627b9ba](#))
- Merge pull request #885 from jsma/patch-1 ([commit de909ad](#))
- Update `request-response.rst` ([commit 3f3263d](#))
- `SgmlLinkExtractor` - fix for parsing `<area>` tag with Unicode present ([commit 49b40f0](#))

7.1.29 Scrapy 0.24.4 (2014-08-09)

- pem file is used by mockserver and required by scrapy bench ([commit 5eddc68](#))
- scrapy bench needs `scrapy.tests*` ([commit d6cb999](#))

7.1.30 Scrapy 0.24.3 (2014-08-09)

- no need to waste travis-ci time on py3 for 0.24 ([commit 8e080c1](#))
- Update installation docs ([commit 1d0c096](#))
- There is a trove classifier for Scrapy framework! ([commit 4c701d7](#))
- update other places where w3lib version is mentioned ([commit d109c13](#))
- Update w3lib requirement to 1.8.0 ([commit 39d2ce5](#))
- Use `w3lib.html.replace_entities()` (`remove_entities()` is deprecated) ([commit 180d3ad](#))
- set `zip_safe=False` ([commit a51ee8b](#))
- do not ship tests package ([commit ee3b371](#))
- scrapy.bat is not needed anymore ([commit c3861cf](#))
- Modernize setup.py ([commit 362e322](#))
- headers can not handle non-string values ([commit 94a5c65](#))
- fix ftp test cases ([commit a274a7f](#))
- The sum up of travis-ci builds are taking like 50min to complete ([commit ae1e2cc](#))
- Update shell.rst typo ([commit e49c96a](#))
- removes weird indentation in the shell results ([commit 1ca489d](#))
- improved explanations, clarified blog post as source, added link for XPath string functions in the spec ([commit 65c8f05](#))
- renamed `UserTimeoutError` and `ServerTimeoutError` #583 ([commit 037f6ab](#))
- adding some xpath tips to selectors docs ([commit 2d103e0](#))
- fix tests to account for <https://github.com/scrapy/w3lib/pull/23> ([commit f8d366a](#))
- `get_func_args` maximum recursion fix #728 ([commit 81344ea](#))
- Updated input/output processor example according to #560. ([commit f7c4ea8](#))
- Fixed Python syntax in tutorial. ([commit db59ed9](#))
- Add test case for tunneling proxy ([commit f090260](#))
- Bugfix for leaking Proxy-Authorization header to remote host when using tunneling ([commit d8793af](#))
- Extract links from XHTML documents with MIME-Type “application/xml” ([commit ed1f376](#))
- Merge pull request #793 from roysc/patch-1 ([commit 91a1106](#))
- Fix typo in commands.rst ([commit 743e1e2](#))
- better testcase for settings.overrides.setdefault ([commit e22daaf](#))

- Using CRLF as line marker according to http 1.1 definition ([commit 5ec430b](#))

7.1.31 Scrapy 0.24.2 (2014-07-08)

- Use a mutable mapping to proxy deprecated settings.overrides and settings.defaults attribute ([commit e5e8133](#))
- there is not support for python3 yet ([commit 3cd6146](#))
- Update python compatible version set to debian packages ([commit fa5d76b](#))
- DOC fix formatting in release notes ([commit c6a9e20](#))

7.1.32 Scrapy 0.24.1 (2014-06-27)

- Fix deprecated CrawlerSettings and increase backward compatibility with .defaults attribute ([commit 8e3f20a](#))

7.1.33 Scrapy 0.24.0 (2014-06-26)

Enhancements

- Improve Scrapy top-level namespace ([issue 494](#), [issue 684](#))
- Add selector shortcuts to responses ([issue 554](#), [issue 690](#))
- Add new lxml based LinkExtractor to replace unmaintained SgmlLinkExtractor ([issue 559](#), [issue 761](#), [issue 763](#))
- Cleanup settings API - part of per-spider settings **GSoC project** ([issue 737](#))
- Add UTF8 encoding header to templates ([issue 688](#), [issue 762](#))
- Telnet console now binds to 127.0.0.1 by default ([issue 699](#))
- Update debian/ubuntu install instructions ([issue 509](#), [issue 549](#))
- Disable smart strings in lxml XPath evaluations ([issue 535](#))
- Restore filesystem based cache as default for http cache middleware ([issue 541](#), [issue 500](#), [issue 571](#))
- Expose current crawler in Scrapy shell ([issue 557](#))
- Improve testsuite comparing CSV and XML exporters ([issue 570](#))
- New `offsite/filtered` and `offsite/domains` stats ([issue 566](#))
- Support `process_links` as generator in CrawlSpider ([issue 555](#))
- Verbose logging and new stats counters for DupeFilter ([issue 553](#))

- Add a `mimetype` parameter to `MailSender.send()` (issue 602)
- Generalize file pipeline log messages (issue 622)
- Replace unencodeable codepoints with html entities in `SGMLLinkExtractor` (issue 565)
- Converted SEP documents to rst format (issue 629, issue 630, issue 638, issue 632, issue 636, issue 640, issue 635, issue 634, issue 639, issue 637, issue 631, issue 633, issue 641, issue 642)
- Tests and docs for `clickdata`'s `nr_index` in `FormRequest` (issue 646, issue 645)
- Allow to disable a downloader handler just like any other component (issue 650)
- Log when a request is discarded after too many redirections (issue 654)
- Log error responses if they are not handled by spider callbacks (issue 612, issue 656)
- Add content-type check to http compression mw (issue 193, issue 660)
- Run pypy tests using latest pypi from ppa (issue 674)
- Run test suite using pytest instead of trial (issue 679)
- Build docs and check for dead links in tox environment (issue 687)
- Make `scrapy.version_info` a tuple of integers (issue 681, issue 692)
- Infer exporter's output format from filename extensions (issue 546, issue 659, issue 760)
- Support case-insensitive domains in `url_is_from_any_domain()` (issue 693)
- Remove pep8 warnings in project and spider templates (issue 698)
- Tests and docs for `request_fingerprint` function (issue 597)
- Update SEP-19 for GSoC project `per-spider settings` (issue 705)
- Set exit code to non-zero when contracts fails (issue 727)
- Add a setting to control what class is instantiated as Downloader component (issue 738)
- Pass response in `item_dropped` signal (issue 724)
- Improve `scrapy check` contracts command (issue 733, issue 752)
- Document `spider.closed()` shortcut (issue 719)
- Document `request_scheduled` signal (issue 746)
- Add a note about reporting security issues (issue 697)
- Add LevelDB http cache storage backend (issue 626, issue 500)
- Sort spider list output of `scrapy list` command (issue 742)
- Multiple documentation enhancements and fixes (issue 575, issue 587, issue 590, issue 596, issue 610, issue 617, issue 618, issue 627, issue 613, issue 643, issue 654, issue 675, issue 663, issue 711, issue 714)

Bugfixes

- Encode unicode URL value when creating Links in RegexLinkExtractor ([issue 561](#))
- Ignore None values in ItemLoader processors ([issue 556](#))
- Fix link text when there is an inner tag in SGMLLinkExtractor and HtmlParserLinkExtractor ([issue 485](#), [issue 574](#))
- Fix wrong checks on subclassing of deprecated classes ([issue 581](#), [issue 584](#))
- Handle errors caused by inspect.stack() failures ([issue 582](#))
- Fix a reference to nonexistent engine attribute ([issue 593](#), [issue 594](#))
- Fix dynamic itemclass example usage of type() ([issue 603](#))
- Use lucasdemarchi/codespell to fix typos ([issue 628](#))
- Fix default value of attrs argument in SgmlLinkExtractor to be tuple ([issue 661](#))
- Fix XXE flaw in sitemap reader ([issue 676](#))
- Fix engine to support filtered start requests ([issue 707](#))
- Fix offsite middleware case on urls with no hostnames ([issue 745](#))
- Testsuite doesn't require PIL anymore ([issue 585](#))

7.1.34 Scrapy 0.22.2 (released 2014-02-14)

- fix a reference to nonexistent engine.slots. closes [#593](#) ([commit 13c099a](#))
- downloaderMW doc typo (spiderMW doc copy remnant) ([commit 8ae11bf](#))
- Correct typos ([commit 1346037](#))

7.1.35 Scrapy 0.22.1 (released 2014-02-08)

- localhost666 can resolve under certain circumstances ([commit 2ec2279](#))
- test inspect.stack failure ([commit cc3eda3](#))
- Handle cases when inspect.stack() fails ([commit 8cb44f9](#))
- Fix wrong checks on subclassing of deprecated classes. closes [#581](#) ([commit 46d98d6](#))
- Docs: 4-space indent for final spider example ([commit 13846de](#))
- Fix HtmlParserLinkExtractor and tests after [#485](#) merge ([commit 368a946](#))
- BaseSgmlLinkExtractor: Fixed the missing space when the link has an inner tag ([commit b566388](#))
- BaseSgmlLinkExtractor: Added unit test of a link with an inner tag ([commit c1cb418](#))

- BaseSgmlLinkExtractor: Fixed `unknown_endtag()` so that it only set `current_link=None` when the end tag match the opening tag ([commit 7e4d627](#))
- Fix tests for Travis-CI build ([commit 76c7e20](#))
- replace unencodeable codepoints with html entities. fixes #562 and #285 ([commit 5f87b17](#))
- RegexLinkExtractor: encode URL unicode value when creating Links ([commit d0ee545](#))
- Updated the tutorial crawl output with latest output. ([commit 8da65de](#))
- Updated shell docs with the crawler reference and fixed the actual shell output. ([commit 875b9ab](#))
- PEP8 minor edits. ([commit f89efaf](#))
- Expose current crawler in the scrapy shell. ([commit 5349cec](#))
- Unused re import and PEP8 minor edits. ([commit 387f414](#))
- Ignore None' s values when using the ItemLoader. ([commit 0632546](#))
- DOC Fixed HTTPCACHE_STORAGE typo in the default value which is now Filesystem instead Dbm. ([commit cde9a8c](#))
- show ubuntu setup instructions as literal code ([commit fb5c9c5](#))
- Update Ubuntu installation instructions ([commit 70fb105](#))
- Merge pull request #550 from stray-leone/patch-1 ([commit 6f70b6a](#))
- modify the version of scrapy ubuntu package ([commit 725900d](#))
- fix 0.22.0 release date ([commit af0219a](#))
- fix typos in news.rst and remove (not released yet) header ([commit b7f58f4](#))

7.1.36 Scrapy 0.22.0 (released 2014-01-17)

Enhancements

- **[Backward incompatible]** Switched HTTPCacheMiddleware backend to filesystem ([issue 541](#)) To restore old backend set `HTTPCACHE_STORAGE` to `scrapy.contrib.httpcache.DbmCacheStorage`
- Proxy `https://` urls using `CONNECT` method ([issue 392](#), [issue 397](#))
- Add a middleware to crawl ajax crawlable pages as defined by google ([issue 343](#))
- Rename `scrapy.spider.BaseSpider` to `scrapy.spider.Spider` ([issue 510](#), [issue 519](#))
- Selectors register EXSLT namespaces by default ([issue 472](#))
- Unify item loaders similar to selectors renaming ([issue 461](#))
- Make `RFPDupeFilter` class easily subclassable ([issue 533](#))
- Improve test coverage and forthcoming Python 3 support ([issue 525](#))

- Promote startup info on settings and middleware to INFO level (issue 520)
- Support partials in `get_func_args` util (issue 506, issue:504)
- Allow running individual tests via tox (issue 503)
- Update extensions ignored by link extractors (issue 498)
- Add middleware methods to get files/images/thumbs paths (issue 490)
- Improve offsite middleware tests (issue 478)
- Add a way to skip default Referer header set by RefererMiddleware (issue 475)
- Do not send `x-gzip` in default `Accept-Encoding` header (issue 469)
- Support defining http error handling using settings (issue 466)
- Use modern python idioms wherever you find legacies (issue 497)
- Improve and correct documentation (issue 527, issue 524, issue 521, issue 517, issue 512, issue 505, issue 502, issue 489, issue 465, issue 460, issue 425, issue 536)

Fixes

- Update Selector class imports in CrawlSpider template (issue 484)
- Fix nonexistent reference to `engine.slots` (issue 464)
- Do not try to call `body_as_unicode()` on a non-TextResponse instance (issue 462)
- Warn when subclassing XPathItemLoader, previously it only warned on instantiation. (issue 523)
- Warn when subclassing XPathSelector, previously it only warned on instantiation. (issue 537)
- Multiple fixes to memory stats (issue 531, issue 530, issue 529)
- Fix overriding url in `FormRequest.from_response()` (issue 507)
- Fix tests runner under pip 1.5 (issue 513)
- Fix logging error when spider name is unicode (issue 479)

7.1.37 Scrapy 0.20.2 (released 2013-12-09)

- Update CrawlSpider Template with Selector changes (commit 6d1457d)
- fix method name in tutorial. closes GH-480 (commit b4fc359)

7.1.38 Scrapy 0.20.1 (released 2013-11-28)

- `include_package_data` is required to build wheels from published sources (commit 5ba1ad5)

- `process_parallel` was leaking the failures on its internal deferreds. closes #458 ([commit 419a780](#))

7.1.39 Scrapy 0.20.0 (released 2013-11-08)

Enhancements

- New Selector's API including CSS selectors ([issue 395](#) and [issue 426](#)),
- Request/Response url/body attributes are now immutable (modifying them had been deprecated for a long time)
- `ITEM_PIPELINES` is now defined as a dict (instead of a list)
- Sitemap spider can fetch alternate URLs ([issue 360](#))
- `Selector.remove_namespaces()` now remove namespaces from element's attributes. ([issue 416](#))
- Paved the road for Python 3.3+ ([issue 435](#), [issue 436](#), [issue 431](#), [issue 452](#))
- New item exporter using native python types with nesting support ([issue 366](#))
- Tune HTTP1.1 pool size so it matches concurrency defined by settings ([commit b43b5f575](#))
- `scrapy.mail.MailSender` now can connect over TLS or upgrade using STARTTLS ([issue 327](#))
- New FilesPipeline with functionality factored out from ImagesPipeline ([issue 370](#), [issue 409](#))
- Recommend Pillow instead of PIL for image handling ([issue 317](#))
- Added debian packages for Ubuntu quantal and raring ([commit 86230c0](#))
- Mock server (used for tests) can listen for HTTPS requests ([issue 410](#))
- Remove multi spider support from multiple core components ([issue 422](#), [issue 421](#), [issue 420](#), [issue 419](#), [issue 423](#), [issue 418](#))
- Travis-CI now tests Scrapy changes against development versions of `w3lib` and `queuelib` python packages.
- Add pypy 2.1 to continuous integration tests ([commit ecfa7431](#))
- Pylint, pep8 and removed old-style exceptions from source ([issue 430](#), [issue 432](#))
- Use importlib for parametric imports ([issue 445](#))
- Handle a regression introduced in Python 2.7.5 that affects `XmlItemExporter` ([issue 372](#))
- Bugfix crawling shutdown on SIGINT ([issue 450](#))
- Do not submit `reset` type inputs in `FormRequest.from_response` ([commit b326b87](#))
- Do not silence download errors when request errback raises an exception ([commit 684cfc0](#))

Bugfixes

- Fix tests under Django 1.6 (commit [b6bed44c](#))
- Lot of bugfixes to retry middleware under disconnections using HTTP 1.1 download handler
- Fix inconsistencies among Twisted releases (issue [406](#))
- Fix scrapy shell bugs (issue [418](#), issue [407](#))
- Fix invalid variable name in setup.py (issue [429](#))
- Fix tutorial references (issue [387](#))
- Improve request-response docs (issue [391](#))
- Improve best practices docs (issue [399](#), issue [400](#), issue [401](#), issue [402](#))
- Improve django integration docs (issue [404](#))
- Document `bindaddress` request meta (commit [37c24e01d7](#))
- Improve `Request` class documentation (issue [226](#))

Other

- Dropped Python 2.6 support (issue [448](#))
- Add `cssselect` python package as install dependency
- Drop libxml2 and multi selector's backend support, `lxml` is required from now on.
- Minimum Twisted version increased to 10.0.0, dropped Twisted 8.0 support.
- Running test suite now requires `mock` python library (issue [390](#))

Thanks

Thanks to everyone who contribute to this release!

List of contributors sorted by number of commits:

```
69 Daniel Graña <dangra@...>
37 Pablo Hoffman <pablo@...>
13 Mikhail Korobov <kmike84@...>
 9 Alex Cepoi <alex.cepoi@...>
 9 alexanderlukanin13 <alexander.lukanin.13@...>
 8 Rolando Espinoza La fuente <darkrho@...>
 8 Lukasz Biedrycki <lukasz.biedrycki@...>
 6 Nicolas Ramirez <namirez.uy@...>
```

(下页继续)

(续上页)

```
3 Paul Tremberth <paul.tremberth@...>
2 Martin Olveyra <molveyra@...>
2 Stefan <misc@...>
2 Rolando Espinoza <darkrho@...>
2 Loren Davie <loren@...>
2 irgmedeiros <irgmedeiros@...>
1 Stefan Koch <taikano@...>
1 Stefan <cct@...>
1 scraperdragon <dragon@...>
1 Kumara Tharmalingam <ktharmal@...>
1 Francesco Piccinno <stack.box@...>
1 Marcos Campal <duendex@...>
1 Dragon Dave <dragon@...>
1 Capi Etheriel <barraponto@...>
1 cacovsky <amarquesferraz@...>
1 Berend Iwema <berend@...>
```

7.1.40 Scrapy 0.18.4 (released 2013-10-10)

- IPython refuses to update the namespace. fix #396 (commit 3d32c4f)
- Fix AlreadyCalledError replacing a request in shell command. closes #407 (commit b1d8919)
- Fix start_requests laziness and early hangs (commit 89faf52)

7.1.41 Scrapy 0.18.3 (released 2013-10-03)

- fix regression on lazy evaluation of start requests (commit 12693a5)
- forms: do not submit reset inputs (commit e429f63)
- increase unittest timeouts to decrease travis false positive failures (commit 912202e)
- backport master fixes to json exporter (commit cfc2d46)
- Fix permission and set umask before generating sdist tarball (commit 06149e0)

7.1.42 Scrapy 0.18.2 (released 2013-09-03)

- Backport scrapy check command fixes and backward compatible multi crawler process(issue 339)

7.1.43 Scrapy 0.18.1 (released 2013-08-27)

- remove extra import added by cherry picked changes ([commit d20304e](#))
- fix crawling tests under twisted pre 11.0.0 ([commit 1994f38](#))
- py26 can not format zero length fields {} ([commit abf756f](#))
- test PotentiaDataLoss errors on unbound responses ([commit b15470d](#))
- Treat responses without content-length or Transfer-Encoding as good responses ([commit c4bf324](#))
- do not include ResponseFailed if http11 handler is not enabled ([commit 6cbe684](#))
- New HTTP client wraps connection losses in ResponseFailed exception. fix #373 ([commit 1a20bba](#))
- limit travis-ci build matrix ([commit 3b01bb8](#))
- Merge pull request #375 from peterarenot/patch-1 ([commit fa766d7](#))
- Fixed so it refers to the correct folder ([commit 3283809](#))
- added quantal & raring to support ubuntu releases ([commit 1411923](#))
- fix retry middleware which didn't retry certain connection errors after the upgrade to http1 client, closes GH-373 ([commit bb35ed0](#))
- fix XmlItemExporter in Python 2.7.4 and 2.7.5 ([commit de3e451](#))
- minor updates to 0.18 release notes ([commit c45e5f1](#))
- fix contributters list format ([commit 0b60031](#))

7.1.44 Scrapy 0.18.0 (released 2013-08-09)

- Lot of improvements to testsuite run using Tox, including a way to test on pypi
- Handle GET parameters for AJAX crawlable urls ([commit 3fe2a32](#))
- Use lxml recover option to parse sitemaps ([issue 347](#))
- Bugfix cookie merging by hostname and not by netloc ([issue 352](#))
- Support disabling `HttpCompressionMiddleware` using a flag setting ([issue 359](#))
- Support xml namespaces using `iternodes` parser in `XMLFeedSpider` ([issue 12](#))
- Support `dont_cache` request meta flag ([issue 19](#))
- Bugfix `scrapy.utils.gz.gunzip` broken by changes in python 2.7.4 ([commit 4dc76e](#))
- Bugfix url encoding on `SgmlLinkExtractor` ([issue 24](#))
- Bugfix `TakeFirst` processor shouldn't discard zero (0) value ([issue 59](#))
- Support nested items in xml exporter ([issue 66](#))

- Improve cookies handling performance ([issue 77](#))
- Log dupe filtered requests once ([issue 105](#))
- Split redirection middleware into status and meta based middlewares ([issue 78](#))
- Use HTTP1.1 as default downloader handler ([issue 109](#) and [issue 318](#))
- Support xpath form selection on `FormRequest.from_response` ([issue 185](#))
- Bugfix unicode decoding error on `SgmlLinkExtractor` ([issue 199](#))
- Bugfix signal dispatching on pypi interpreter ([issue 205](#))
- Improve request delay and concurrency handling ([issue 206](#))
- Add RFC2616 cache policy to `HttpCacheMiddleware` ([issue 212](#))
- Allow customization of messages logged by engine ([issue 214](#))
- Multiples improvements to `DjangoItem` ([issue 217](#), [issue 218](#), [issue 221](#))
- Extend Scrapy commands using setuptools entry points ([issue 260](#))
- Allow spider `allowed_domains` value to be set/tuple ([issue 261](#))
- Support `settings.getdict` ([issue 269](#))
- Simplify internal `scrapy.core.scrapers` slot handling ([issue 271](#))
- Added `Item.copy` ([issue 290](#))
- Collect idle downloader slots ([issue 297](#))
- Add `ftp://` scheme downloader handler ([issue 329](#))
- Added downloader benchmark webserver and spider tools 爬虫器硬件性能
- Moved persistent (on disk) queues to a separate project (`queuelib`) which scrapy now depends on
- Add scrapy commands using external libraries ([issue 260](#))
- Added `--pdb` option to `scrapy` command line tool
- Added `XPathSelector.remove_namespaces()` which allows to remove all namespaces from XML documents for convenience (to work with namespace-less XPath). Documented in [选择器](#).
- Several improvements to spider contracts
- New default middleware named `MetaRefreshMiddleware` that handles meta-refresh html tag redirections,
- `MetaRefreshMiddleware` and `RedirectMiddleware` have different priorities to address [#62](#)
- added `from_crawler` method to spiders
- added system tests with mock server
- more improvements to Mac OS compatibility (thanks Alex Cepoi)

- several more cleanups to singletons and multi-spider support (thanks Nicolas Ramirez)
- support custom download slots
- added `--spider` option to “shell” command.
- log overridden settings when scrapy starts

Thanks to everyone who contribute to this release. Here is a list of contributors sorted by number of commits:

```
130 Pablo Hoffman <pablo@...>
97 Daniel Graña <dangra@...>
20 Nicolás Ramírez <namirez.uy@...>
13 Mikhail Korobov <kmike84@...>
12 Pedro Faustino <pedrobandim@...>
11 Steven Almeroth <sroth77@...>
5 Rolando Espinoza La fuente <darkrho@...>
4 Michal Danilak <mimino.coder@...>
4 Alex Cepoi <alex.cepoi@...>
4 Alexandr N Zamaraev (aka tonal) <tonal@...>
3 paul <paul.tremberth@...>
3 Martin Olveyra <molveyra@...>
3 Jordi Llonch <llonchj@...>
3 arijitchakraborty <myself.arijit@...>
2 Shane Evans <shane.evans@...>
2 joehillen <joehillen@...>
2 Hart <HartSimha@...>
2 Dan <ellis23@...>
1 Zuhao Wan <wanzuhao@...>
1 whodatninja <blake@...>
1 vkrest <v.krestiannykov@...>
1 tpeng <pengtaoo@...>
1 Tom Mortimer-Jones <tom@...>
1 Rocio Aramberri <roschegel@...>
1 Pedro <pedro@...>
1 notsobad <wangxiaohugg@...>
1 Natan L <kuyanatan.nlao@...>
1 Mark Grey <mark.grey@...>
1 Luan <luanpab@...>
1 Libor Nenadál <libor.nenadal@...>
1 Juan M Uys <opyate@...>
1 Jonas Brunsgaard <jonas.brunsgaard@...>
1 Ilya Baryshev <baryshev@...>
1 Hasnain Lakhani <m.hasnain.lakhani@...>
```

(下页继续)

(续上页)

```
1 Emanuel Schorsch <emschorsch@...>
1 Chris Tilden <chris.tilden@...>
1 Capi Etheriel <barraponto@...>
1 cacovsky <amarquesferraz@...>
1 Berend Iwema <berend@...>
```

7.1.45 Scrapy 0.16.5 (released 2013-05-30)

- obey request method when scrapy deploy is redirected to a new endpoint (commit 8c4fcee)
- fix inaccurate downloader middleware documentation. refs #280 (commit 40667cb)
- doc: remove links to diveintopython.org, which is no longer available. closes #246 (commit bd58bfa)
- Find form nodes in invalid html5 documents (commit e3d6945)
- Fix typo labeling attrs type bool instead of list (commit a274276)

7.1.46 Scrapy 0.16.4 (released 2013-01-23)

- fixes spelling errors in documentation (commit 6d2b3aa)
- add doc about disabling an extension. refs #132 (commit c90de33)
- Fixed error message formatting. log.err() doesn't support cool formatting and when error occurred, the message was: "ERROR: Error processing %(item)s" (commit c16150c)
- lint and improve images pipeline error logging (commit 56b45fc)
- fixed doc typos (commit 243be84)
- add documentation topics: Broad Crawls & Common Practices (commit 1fbb715)
- fix bug in scrapy parse command when spider is not specified explicitly. closes #209 (commit c72e682)
- Update docs/topics/commands.rst (commit 28eac7a)

7.1.47 Scrapy 0.16.3 (released 2012-12-07)

- Remove concurrency limitation when using download delays and still ensure inter-request delays are enforced (commit 487b9b5)
- add error details when image pipeline fails (commit 8232569)
- improve mac os compatibility (commit 8dcf8aa)
- setup.py: use README.rst to populate long_description (commit 7b5310d)

- doc: removed obsolete references to ClientForm ([commit 80f9bb6](#))
- correct docs for default storage backend ([commit 2aa491b](#))
- doc: removed broken proxyhub link from FAQ ([commit bdf61c4](#))
- Fixed docs typo in SpiderOpenCloseLogging example ([commit 7184094](#))

7.1.48 Scrapy 0.16.2 (released 2012-11-09)

- scrapy contracts: python2.6 compat ([commit a4a9199](#))
- scrapy contracts verbose option ([commit ec41673](#))
- proper unittest-like output for scrapy contracts ([commit 86635e4](#))
- added open_in_browser to debugging doc ([commit c9b690d](#))
- removed reference to global scrapy stats from settings doc ([commit dd55067](#))
- Fix SpiderState bug in Windows platforms ([commit 58998f4](#))

7.1.49 Scrapy 0.16.1 (released 2012-10-26)

- fixed LogStats extension, which got broken after a wrong merge before the 0.16 release ([commit 8c780fd](#))
- better backward compatibility for scrapy.conf.settings ([commit 3403089](#))
- extended documentation on how to access crawler stats from extensions ([commit c4da0b5](#))
- removed .hgtags (no longer needed now that scrapy uses git) ([commit d52c188](#))
- fix dashes under rst headers ([commit fa4f7f9](#))
- set release date for 0.16.0 in news ([commit e292246](#))

7.1.50 Scrapy 0.16.0 (released 2012-10-18)

Scrapy changes:

- added 爬虫器约束, a mechanism for testing spiders in a formal/reproducible way
- added options `-o` and `-t` to the *runspider* command
- documented 爬虫器节流 and added to extensions installed by default. You still need to enable it with `AUTOTHROTTLE_ENABLED`
- major Stats Collection refactoring: removed separation of global/per-spider stats, removed stats-related signals (`stats_spider_opened`, etc). Stats are much simpler now, backward compatibility is kept on the Stats Collector API and signals.
- added `process_start_requests()` method to spider middlewares

- dropped Signals singleton. Signals should now be accessed through the Crawler.signals attribute. See the signals documentation for more info.
- dropped Signals singleton. Signals should now be accessed through the Crawler.signals attribute. See the signals documentation for more info.
- dropped Stats Collector singleton. Stats can now be accessed through the Crawler.stats attribute. See the stats collection documentation for more info.
- documented 核心 *API*
- lxml is now the default selectors backend instead of libxml2
- ported FormRequest.from_response() to use lxml instead of ClientForm
- removed modules: scrapy.xlib.BeautifulSoup and scrapy.xlib.ClientForm
- SitemapSpider: added support for sitemap urls ending in .xml and .xml.gz, even if they advertise a wrong content type (commit 10ed28b)
- StackTraceDump extension: also dump trackref live references (commit fe2ce93)
- nested items now fully supported in JSON and JSONLines exporters
- added cookiejar Request meta key to support multiple cookie sessions per spider
- decoupled encoding detection code to w3lib.encoding, and ported Scrapy code to use that module
- dropped support for Python 2.5. See <https://blog.scrapinghub.com/2012/02/27/scrapy-0-15-dropping-support-for-python-2-5/>
- dropped support for Twisted 2.5
- added REFERER_ENABLED setting, to control referer middleware
- changed default user agent to: Scrapy/VERSION (+http://scrapy.org)
- removed (undocumented) HTMLImageLinkExtractor class from scrapy.contrib.linkextractors.image
- removed per-spider settings (to be replaced by instantiating multiple crawler objects)
- USER_AGENT spider attribute will no longer work, use user_agent attribute instead
- DOWNLOAD_TIMEOUT spider attribute will no longer work, use download_timeout attribute instead
- removed ENCODING_ALIASES setting, as encoding auto-detection has been moved to the w3lib library
- promoted topics-djangoitem to main contrib
- LogFormatter method now return dicts (instead of strings) to support lazy formatting (issue 164, commit dcef7b0)
- downloader handlers (DOWNLOAD_HANDLERS setting) now receive settings as the first argument of the constructor

- replaced memory usage accounting with (more portable) `resource` module, removed `scrapy.utils.memory` module
- removed signal: `scrapy.mail.mail_sent`
- removed `TRACK_REFS` setting, now `trackrefs` is always enabled
- DBM is now the default storage backend for HTTP cache middleware
- number of log messages (per level) are now tracked through Scrapy stats (stat name: `log_count/LEVEL`)
- number received responses are now tracked through Scrapy stats (stat name: `response_received_count`)
- removed `scrapy.log.started` attribute

7.1.51 Scrapy 0.14.4

- added precise to supported ubuntu distros (commit [b7e46df](#))
- fixed bug in json-rpc webservice reported in <https://groups.google.com/forum/#!topic/scrapy-users/qgVBmFybNAQ/discussion>. also removed no longer supported ‘run’ command from extras/scrapy-ws.py (commit [340fbdb](#))
- meta tag attributes for content-type http equiv can be in any order. #123 (commit [0cb68af](#))
- replace “import Image” by more standard “from PIL import Image” . closes #88 (commit [4d17048](#))
- return trial status as bin/runtests.sh exit value. #118 (commit [b7b2e7f](#))

7.1.52 Scrapy 0.14.3

- forgot to include pydispatch license. #118 (commit [fd85f9c](#))
- include egg files used by testsuite in source distribution. #118 (commit [c897793](#))
- update docstring in project template to avoid confusion with genspider command, which may be considered as an advanced feature. refs #107 (commit [2548dcc](#))
- added note to docs/topics/firebug.rst about google directory being shut down (commit [668e352](#))
- dont discard slot when empty, just save in another dict in order to recycle if needed again. (commit [8e9f607](#))
- do not fail handling unicode xpathes in libxml2 backed selectors (commit [b830e95](#))
- fixed minor mistake in Request objects documentation (commit [bf3c9ee](#))
- fixed minor defect in link extractors documentation (commit [ba14f38](#))
- removed some obsolete remaining code related to sqlite support in scrapy (commit [0665175](#))

7.1.53 Scrapy 0.14.2

- move buffer pointing to start of file before computing checksum. refs #92 (commit 6a5bef2)
- Compute image checksum before persisting images. closes #92 (commit 9817df1)
- remove leaking references in cached failures (commit 673a120)
- fixed bug in MemoryUsage extension: get_engine_status() takes exactly 1 argument (0 given) (commit 11133e9)
- fixed struct.error on http compression middleware. closes #87 (commit 1423140)
- ajax crawling wasn't expanding for unicode urls (commit 0de3fb4)
- Catch start_requests iterator errors. refs #83 (commit 454a21d)
- Speed-up libxml2 XPathSelector (commit 2fbd662)
- updated versioning doc according to recent changes (commit 0a070f5)
- scrapyd: fixed documentation link (commit 2b4e4c3)
- extras/makedeb.py: no longer obtaining version from git (commit caff0e)

7.1.54 Scrapy 0.14.1

- extras/makedeb.py: no longer obtaining version from git (commit caff0e)
- bumped version to 0.14.1 (commit 6cb9e1c)
- fixed reference to tutorial directory (commit 4b86bd6)
- doc: removed duplicated callback argument from Request.replace() (commit 1aecdd)
- fixed formatting of scrapyd doc (commit 8bf19e6)
- Dump stacks for all running threads and fix engine status dumped by StackTraceDump extension (commit 14a8e6e)
- added comment about why we disable ssl on boto images upload (commit 5223575)
- SSL handshaking hangs when doing too many parallel connections to S3 (commit 63d583d)
- change tutorial to follow changes on dmoz site (commit bcb3198)
- Avoid _disconnectedDeferred AttributeError exception in Twisted>=11.1.0 (commit 98f3f87)
- allow spider to set autothrottle max concurrency (commit 175a4b5)

7.1.55 Scrapy 0.14

New features and settings

- Support for [AJAX](#) [crawleable](#) [urls](#)
- New persistent scheduler that stores requests on disk, allowing to suspend and resume crawls ([r2737](#))
- added `-o` option to `scrapy crawl`, a shortcut for dumping scraped items into a file (or standard output using `-`)
- Added support for passing custom settings to Scrapyd `schedule.json` api ([r2779](#), [r2783](#))
- New `ChunkedTransferMiddleware` (enabled by default) to support [chunked transfer encoding](#) ([r2769](#))
- Add boto 2.0 support for S3 downloader handler ([r2763](#))
- Added [marshal](#) to formats supported by feed exports ([r2744](#))
- In request errbacks, offending requests are now received in `failure.request` attribute ([r2738](#))
- **Big downloader refactoring to support per domain/ip concurrency limits ([r2732](#))**
 - `CONCURRENT_REQUESTS_PER_SPIDER` setting has been deprecated and replaced by:
 - * `CONCURRENT_REQUESTS`, `CONCURRENT_REQUESTS_PER_DOMAIN`,
`CONCURRENT_REQUESTS_PER_IP`
 - check the documentation for more details
- Added builtin caching DNS resolver ([r2728](#))
- Moved Amazon AWS-related components/extensions (SQS spider queue, SimpleDB stats collector) to a separate project: [scaws](<https://github.com/scrapinghub/scaws>) ([r2706](#), [r2714](#))
- Moved spider queues to scrapyd: `scrapy.spiderqueue` -> `scrapyd.spiderqueue` ([r2708](#))
- Moved sqlite utils to scrapyd: `scrapy.utils.sqlite` -> `scrapyd.sqlite` ([r2781](#))
- Real support for returning iterators on `start_requests()` method. The iterator is now consumed during the crawl when the spider is getting idle ([r2704](#))
- Added `REDIRECT_ENABLED` setting to quickly enable/disable the redirect middleware ([r2697](#))
- Added `RETRY_ENABLED` setting to quickly enable/disable the retry middleware ([r2694](#))
- Added `CloseSpider` exception to manually close spiders ([r2691](#))
- Improved encoding detection by adding support for HTML5 meta charset declaration ([r2690](#))
- Refactored close spider behavior to wait for all downloads to finish and be processed by spiders, before closing the spider ([r2688](#))
- Added `SitemapSpider` (see documentation in Spiders page) ([r2658](#))
- Added `LogStats` extension for periodically logging basic stats (like crawled pages and scraped items) ([r2657](#))

- Make handling of gzipped responses more robust (#319, r2643). Now Scrapy will try and decompress as much as possible from a gzipped response, instead of failing with an `IOError`.
- Simplified `!MemoryDebugger` extension to use stats for dumping memory debugging info (r2639)
- Added new command to edit spiders: `scrapy edit` (r2636) and `-e` flag to `genspider` command that uses it (r2653)
- Changed default representation of items to pretty-printed dicts. (r2631). This improves default logging by making log more readable in the default case, for both `Scraped` and `Dropped` lines.
- Added `spider_error` signal (r2628)
- Added `COOKIES_ENABLED` setting (r2625)
- Stats are now dumped to Scrapy log (default value of `STATS_DUMP` setting has been changed to `True`). This is to make Scrapy users more aware of Scrapy stats and the data that is collected there.
- Added support for dynamically adjusting download delay and maximum concurrent requests (r2599)
- Added new DBM HTTP cache storage backend (r2576)
- Added `listjobs.json` API to Scrapyd (r2571)
- `CsvItemExporter`: added `join_multivalued` parameter (r2578)
- Added namespace support to `xmliter_lxml` (r2552)
- Improved cookies middleware by making `COOKIES_DEBUG` nicer and documenting it (r2579)
- Several improvements to Scrapyd and Link extractors

Code rearranged and removed

- **Merged `item passed` and `item scraped` concepts, as they have often proved confusing in the past. This**
 - original `item_scraped` signal was removed
 - original `item_passed` signal was renamed to `item_scraped`
 - old log lines `Scraped Item...` were removed
 - old log lines `Passed Item...` were renamed to `Scraped Item...` lines and downgraded to `DEBUG` level
- **Reduced Scrapy codebase by striping part of Scrapy code into two new libraries:**
 - `w3lib` (several functions from `scrapy.utils.{http,markup,multipart,response,url}`, done in r2584)
 - `scrapely` (was `scrapy.contrib.ibl`, done in r2586)
- Removed unused function: `scrapy.utils.request.request_info()` (r2577)

- Removed googledir project from `examples/googledir`. There's now a new example project called `dirbot` available on github: <https://github.com/scrapy/dirbot>
- Removed support for default field values in Scrapy items (r2616)
- Removed experimental crawls spider v2 (r2632)
- Removed scheduler middleware to simplify architecture. Duplicates filter is now done in the scheduler itself, using the same dupe filtering class as before (`DUPEFILTER_CLASS` setting) (r2640)
- Removed support for passing urls to `scrapy crawl` command (use `scrapy parse` instead) (r2704)
- Removed deprecated Execution Queue (r2704)
- Removed (undocumented) spider context extension (from `scrapy.contrib.spidercontext`) (r2780)
- removed `CONCURRENT_SPIDERS` setting (use `scrapyd maxproc` instead) (r2789)
- Renamed attributes of core components: `downloader.sites` -> `downloader.slots`, `scraper.sites` -> `scraper.slots` (r2717, r2718)
- Renamed setting `CLOSESPIDER_ITEMPASSED` to `CLOSESPIDER_ITEMCOUNT` (r2655). Backward compatibility kept.

7.1.56 Scrapy 0.12

The numbers like #NNN reference tickets in the old issue tracker (Trac) which is no longer available.

New features and improvements

- Passed item is now sent in the `item` argument of the `item_passed` (#273)
- Added verbose option to `scrapy version` command, useful for bug reports (#298)
- HTTP cache now stored by default in the project data dir (#279)
- Added project data storage directory (#276, #277)
- Documented file structure of Scrapy projects (see command-line tool doc)
- New lxml backend for XPath selectors (#147)
- Per-spider settings (#245)
- Support exit codes to signal errors in Scrapy commands (#248)
- Added `-c` argument to `scrapy shell` command
- Made `libxml2` optional (#260)
- New `deploy` command (#261)
- Added `CLOSESPIDER_PAGECOUNT` setting (#253)

- Added `CLOSESPIDER_ERRORCOUNT` setting (#254)

Scrapyd changes

- Scrapyd now uses one process per spider
- It stores one log file per spider run, and rotate them keeping the latest 5 logs per spider (by default)
- A minimal web ui was added, available at <http://localhost:6800> by default
- There is now a `scrapy server` command to start a Scrapyd server of the current project

Changes to settings

- added `HTTPCACHE_ENABLED` setting (False by default) to enable HTTP cache middleware
- changed `HTTPCACHE_EXPIRATION_SECS` semantics: now zero means “never expire” .

Deprecated/obsoleted functionality

- Deprecated `runserver` command in favor of `server` command which starts a Scrapyd server. See also: Scrapyd changes
- Deprecated `queue` command in favor of using Scrapyd `schedule.json` API. See also: Scrapyd changes
- Removed the `!XmlItemLoader` (experimental contrib which never graduated to main contrib)

7.1.57 Scrapy 0.10

The numbers like #NNN reference tickets in the old issue tracker (Trac) which is no longer available.

New features and improvements

- New Scrapy service called `scrapyd` for deploying Scrapy crawlers in production (#218) (documentation available)
- Simplified Images pipeline usage which doesn't require subclassing your own images pipeline now (#217)
- Scrapy shell now shows the Scrapy log by default (#206)
- Refactored execution queue in a common base code and pluggable backends called “spider queues” (#220)
- New persistent spider queue (based on SQLite) (#198), available by default, which allows to start Scrapy in server mode and then schedule spiders to run.

- Added documentation for Scrapy command-line tool and all its available sub-commands. (documentation available)
- Feed exporters with pluggable backends (#197) (documentation available)
- Deferred signals (#193)
- Added two new methods to item pipeline `open_spider()`, `close_spider()` with deferred support (#195)
- Support for overriding default request headers per spider (#181)
- Replaced default Spider Manager with one with similar functionality but not depending on Twisted Plugins (#186)
- Splitted Debian package into two packages - the library and the service (#187)
- Scrapy log refactoring (#188)
- New extension for keeping persistent spider contexts among different runs (#203)
- Added `dont_redirect` request.meta key for avoiding redirects (#233)
- Added `dont_retry` request.meta key for avoiding retries (#234)

Command-line tool changes

- New `scrapy` command which replaces the old `scrapy-ctl.py` (#199) - there is only one global `scrapy` command now, instead of one `scrapy-ctl.py` per project - Added `scrapy.bat` script for running more conveniently from Windows
- Added bash completion to command-line tool (#210)
- Renamed command `start` to `runserver` (#209)

API changes

- `url` and `body` attributes of Request objects are now read-only (#230)
- `Request.copy()` and `Request.replace()` now also copies their `callback` and `errback` attributes (#231)
- Removed `UrlFilterMiddleware` from `scrapy.contrib` (already disabled by default)
- Offsite middleware doesn't filter out any request coming from a spider that doesn't have a `allowed_domains` attribute (#225)
- Removed Spider Manager `load()` method. Now spiders are loaded in the constructor itself.
- **Changes to Scrapy Manager (now called "Crawler"):**
 - `scrapy.core.manager.ScrapyManager` class renamed to `scrapy.crawler.Crawler`
 - `scrapy.core.manager.scrapymanager` singleton moved to `scrapy.project.crawler`

- Moved module: `scrapy.contrib.spidermanager` to `scrapy.spidermanager`
- Spider Manager singleton moved from `scrapy.spider.spiders` to the `spiders`` attribute of `scrapy.project.crawler` singleton.
- **moved Stats Collector classes: (#204)**
 - `scrapy.stats.collector.StatsCollector` to `scrapy.statscol.StatsCollector`
 - `scrapy.stats.collector.SimplesdbStatsCollector` to `scrapy.contrib.statscol.SimplesdbStatsCollector`
- default per-command settings are now specified in the `default_settings` attribute of command object class (#201)
- **changed arguments of Item pipeline `process_item()` method from `(spider, item)` to `(item, spider)`**
 - backward compatibility kept (with deprecation warning)
- **moved `scrapy.core.signals` module to `scrapy.signals`**
 - backward compatibility kept (with deprecation warning)
- **moved `scrapy.core.exceptions` module to `scrapy.exceptions`**
 - backward compatibility kept (with deprecation warning)
- added `handles_request()` class method to `BaseSpider`
- dropped `scrapy.log.exc()` function (use `scrapy.log.err()` instead)
- dropped `component` argument of `scrapy.log.msg()` function
- dropped `scrapy.log.log_level` attribute
- Added `from_settings()` class methods to Spider Manager, and Item Pipeline Manager

Changes to settings

- Added `HTTPCACHE_IGNORE_SCHEMES` setting to ignore certain schemes on `!HttpCacheMiddleware` (#225)
- Added `SPIDER_QUEUE_CLASS` setting which defines the spider queue to use (#220)
- Added `KEEP_ALIVE` setting (#220)
- Removed `SERVICE_QUEUE` setting (#220)
- Removed `COMMANDS_SETTINGS_MODULE` setting (#201)
- Renamed `REQUEST_HANDLERS` to `DOWNLOAD_HANDLERS` and make download handlers classes (instead of functions)

7.1.58 Scrapy 0.9

The numbers like #NNN reference tickets in the old issue tracker (Trac) which is no longer available.

New features and improvements

- Added SMTP-AUTH support to scrapy.mail
- New settings added: MAIL_USER, MAIL_PASS (r2065 | #149)
- Added new scrapy-ctl view command - To view URL in the browser, as seen by Scrapy (r2039)
- Added web service for controlling Scrapy process (this also deprecates the web console. (r2053 | #167)
- Support for running Scrapy as a service, for production systems (r1988, r2054, r2055, r2056, r2057 | #168)
- Added wrapper induction library (documentation only available in source code for now). (r2011)
- Simplified and improved response encoding support (r1961, r1969)
- Added LOG_ENCODING setting (r1956, documentation available)
- Added RANDOMIZE_DOWNLOAD_DELAY setting (enabled by default) (r1923, doc available)
- MailSender is no longer IO-blocking (r1955 | #146)
- Linkextractors and new Crawlspider now handle relative base tag urls (r1960 | #148)
- Several improvements to Item Loaders and processors (r2022, r2023, r2024, r2025, r2026, r2027, r2028, r2029, r2030)
- Added support for adding variables to telnet console (r2047 | #165)
- Support for requests without callbacks (r2050 | #166)

API changes

- Change Spider.domain_name to Spider.name (SEP-012, r1975)
- Response.encoding is now the detected encoding (r1961)
- HttpErrorMiddleware now returns None or raises an exception (r2006 | #157)
- scrapy.command modules relocation (r2035, r2036, r2037)
- Added ExecutionQueue for feeding spiders to scrape (r2034)
- Removed ExecutionEngine singleton (r2039)
- Ported S3ImagesStore (images pipeline) to use boto and threads (r2033)
- Moved module: scrapy.management.telnet to scrapy.telnet (r2047)

Changes to default settings

- Changed default `SCHEDULER_ORDER` to `DFO` ([r1939](#))

7.1.59 Scrapy 0.8

The numbers like `#NNN` reference tickets in the old issue tracker (Trac) which is no longer available.

New features

- Added `DEFAULT_RESPONSE_ENCODING` setting ([r1809](#))
- Added `dont_click` argument to `FormRequest.from_response()` method ([r1813](#), [r1816](#))
- Added `clickdata` argument to `FormRequest.from_response()` method ([r1802](#), [r1803](#))
- Added support for HTTP proxies (`HttpProxyMiddleware`) ([r1781](#), [r1785](#))
- Offsite spider middleware now logs messages when filtering out requests ([r1841](#))

Backward-incompatible changes

- Changed `scrapy.utils.response.get_meta_refresh()` signature ([r1804](#))
- Removed deprecated `scrapy.item.ScrapedItem` class - use `scrapy.item.Item` instead ([r1838](#))
- Removed deprecated `scrapy.xpath` module - use `scrapy.selector` instead. ([r1836](#))
- Removed deprecated `core.signals.domain_open` signal - use `core.signals.domain_opened` instead ([r1822](#))
- **`log.msg()` now receives a spider argument** ([r1822](#))
 - Old domain argument has been deprecated and will be removed in 0.9. For spiders, you should always use the `spider` argument and pass spider references. If you really want to pass a string, use the `component` argument instead.
- Changed core signals `domain_opened`, `domain_closed`, `domain_idle`
- **Changed Item pipeline to use spiders instead of domains**
 - The `domain` argument of `process_item()` item pipeline method was changed to `spider`, the new signature is: `process_item(spider, item)` ([r1827](#) | [#105](#))
 - To quickly port your code (to work with Scrapy 0.8) just use `spider.domain_name` where you previously used `domain`.
- **Changed Stats API to use spiders instead of domains** ([r1849](#) | [#113](#))
 - `StatsCollector` was changed to receive spider references (instead of domains) in its methods (`set_value`, `inc_value`, etc).

- added `StatsCollector.iter_spider_stats()` method
- removed `StatsCollector.list_domains()` method
- Also, Stats signals were renamed and now pass around spider references (instead of domains). Here's a summary of the changes:
 - To quickly port your code (to work with Scrapy 0.8) just use `spider.domain_name` where you previously used `domain`. `spider_stats` contains exactly the same data as `domain_stats`.
- **CloseDomain extension moved to `scrapy.contrib.closespider.CloseSpider` (r1833)**
 - Its settings were also renamed:
 - * `CLOSEDOMAIN_TIMEOUT` to `CLOSESPIDER_TIMEOUT`
 - * `CLOSEDOMAIN_ITEMCOUNT` to `CLOSESPIDER_ITEMCOUNT`
- Removed deprecated `SCRAPYSETTINGS_MODULE` environment variable - use `SCRAPY_SETTINGS_MODULE` instead (r1840)
- Renamed setting: `REQUESTS_PER_DOMAIN` to `CONCURRENT_REQUESTS_PER_SPIDER` (r1830, r1844)
- Renamed setting: `CONCURRENT_DOMAINS` to `CONCURRENT_SPIDERS` (r1830)
- Refactored HTTP Cache middleware
- HTTP Cache middleware has been heavily refactored, retaining the same functionality except for the domain sectorization which was removed. (r1843)
- Renamed exception: `DontCloseDomain` to `DontCloseSpider` (r1859 | #120)
- Renamed extension: `DelayedCloseDomain` to `SpiderCloseDelay` (r1861 | #121)
- Removed obsolete `scrapy.utils.markup.remove_escape_chars` function - use `scrapy.utils.markup.replace_escape_chars` instead (r1865)

7.1.60 Scrapy 0.7

First release of Scrapy.

7.2 贡献 Scrapy 代码

重要: Double check that you are reading the most recent version of this document at <https://docs.scrapy.org/en/master/contributing.html>

There are many ways to contribute to Scrapy. Here are some of them:

- Blog about Scrapy. Tell the world how you’ re using Scrapy. This will help newcomers with more examples and will help the Scrapy project to increase its visibility.
- Report bugs and request features in the [issue tracker](#), trying to follow the guidelines detailed in *Reporting bugs* below.
- Submit patches for new functionalities and/or bug fixes. Please read *Writing patches* and *Submitting patches* below for details on how to write and submit a patch.
- Join the [Scrapy subreddit](#) and share your ideas on how to improve Scrapy. We’ re always open to suggestions.
- Answer Scrapy questions at [Stack Overflow](#).

7.2.1 Reporting bugs

注解: Please report security issues **only** to scrapy-security@googlegroups.com. This is a private list only open to trusted Scrapy developers, and its archives are not public.

Well-written bug reports are very helpful, so keep in mind the following guidelines when you’ re going to report a new bug.

- check the *FAQ* first to see if your issue is addressed in a well-known question
- if you have a general question about scrapy usage, please ask it at [Stack Overflow](#) (use “scrapy” tag).
- check the [open issues](#) to see if the issue has already been reported. If it has, don’ t dismiss the report, but check the ticket history and comments. If you have additional useful information, please leave a comment, or consider *sending a pull request* with a fix.
- search the [scrapy-users](#) list and [Scrapy subreddit](#) to see if it has been discussed there, or if you’ re not sure if what you’ re seeing is a bug. You can also ask in the [#scrapy](#) IRC channel.
- write **complete, reproducible, specific bug reports**. The smaller the test case, the better. Remember that other developers won’ t have your project to reproduce the bug, so please include all relevant files required to reproduce it. See for example StackOverflow’ s guide on creating a [Minimal, Complete, and Verifiable example](#) exhibiting the issue.
- the most awesome way to provide a complete reproducible example is to send a pull request which adds a failing test case to the Scrapy testing suite (see *Submitting patches*). This is helpful even if you don’ t have an intention to fix the issue yourselves.
- include the output of `scrapy version -v` so developers working on your bug know exactly which version and platform it occurred on, which is often very helpful for reproducing it, or knowing if it was already fixed.

7.2.2 Writing patches

The better a patch is written, the higher the chances that it'll get accepted and the sooner it will be merged.

Well-written patches should:

- contain the minimum amount of code required for the specific change. Small patches are easier to review and merge. So, if you're doing more than one change (or bug fix), please consider submitting one patch per change. Do not collapse multiple changes into a single patch. For big changes consider using a patch queue.
- pass all unit-tests. See *Running tests* below.
- include one (or more) test cases that check the bug fixed or the new functionality added. See *Writing tests* below.
- if you're adding or changing a public (documented) API, please include the documentation changes in the same patch. See *Documentation policies* below.

7.2.3 Submitting patches

The best way to submit a patch is to issue a [pull request](#) on GitHub, optionally creating a new issue first.

Remember to explain what was fixed or the new functionality (what it is, why it's needed, etc). The more info you include, the easier will be for core developers to understand and accept your patch.

You can also discuss the new functionality (or bug fix) before creating the patch, but it's always good to have a patch ready to illustrate your arguments and show that you have put some additional thought into the subject. A good starting point is to send a pull request on GitHub. It can be simple enough to illustrate your idea, and leave documentation/tests for later, after the idea has been validated and proven useful. Alternatively, you can start a conversation in the [Scrapy subreddit](#) to discuss your idea first.

Sometimes there is an existing pull request for the problem you'd like to solve, which is stalled for some reason. Often the pull request is in a right direction, but changes are requested by Scrapy maintainers, and the original pull request author hasn't had time to address them. In this case consider picking up this pull request: open a new pull request with all commits from the original pull request, as well as additional changes to address the raised issues. Doing so helps a lot; it is not considered rude as soon as the original author is acknowledged by keeping his/her commits.

You can pull an existing pull request to a local branch by running `git fetch upstream pull/$PR_NUMBER/head:$BRANCH_NAME_TO_CREATE` (replace 'upstream' with a remote name for scrapy repository, \$PR_NUMBER with an ID of the pull request, and \$BRANCH_NAME_TO_CREATE with a name of the branch you want to create locally). See also: <https://help.github.com/articles/checking-out-pull-requests-locally/#modifying-an-inactive-pull-request-locally>.

When writing GitHub pull requests, try to keep titles short but descriptive. E.g. For bug #411: "Scrapy hangs if an exception raises in start_requests" prefer "Fix hanging when exception occurs in start_requests (#411)" instead of "Fix for #411". Complete titles make it easy to skim through the issue tracker.

Finally, try to keep aesthetic changes (**PEP 8** compliance, unused imports removal, etc) in separate commits from functional changes. This will make pull requests easier to review and more likely to get merged.

7.2.4 Coding style

Please follow these coding conventions when writing code for inclusion in Scrapy:

- Unless otherwise specified, follow **PEP 8**.
- It's OK to use lines longer than 80 chars if it improves the code readability.
- Don't put your name in the code you contribute; git provides enough metadata to identify author of the code. See <https://help.github.com/articles/setting-your-username-in-git/> for setup instructions.

7.2.5 Documentation policies

For reference documentation of API members (classes, methods, etc.) use docstrings and make sure that the Sphinx documentation uses the `autodoc` extension to pull the docstrings. API reference documentation should follow docstring conventions (**PEP 257**) and be IDE-friendly: short, to the point, and it may provide short examples.

Other types of documentation, such as tutorials or topics, should be covered in files within the `docs/` directory. This includes documentation that is specific to an API member, but goes beyond API reference documentation.

In any case, if something is covered in a docstring, use the `autodoc` extension to pull the docstring into the documentation instead of duplicating the docstring in files within the `docs/` directory.

7.2.6 Tests

Tests are implemented using the Twisted `unit-testing framework`, running tests requires `tox`.

Running tests

Make sure you have a recent enough `tox` installation:

```
tox --version
```

If your version is older than 1.7.0, please update it first:

```
pip install -U tox
```

To run all tests go to the root directory of Scrapy source code and run:

```
tox
```

To run a specific test (say `tests/test_loader.py`) use:

```
tox -- tests/test_loader.py
```

To run the tests on a specific `tox` environment, use `-e <name>` with an environment name from `tox.ini`. For example, to run the tests with Python 3.6 use:

```
tox -e py36
```

You can also specify a comma-separated list of environments, and use `tox`'s `parallel mode` to run the tests on multiple environments in parallel:

```
tox -e py27,py36 -p auto
```

To pass command-line options to `pytest`, add them after `--` in your call to `tox`. Using `--` overrides the default positional arguments defined in `tox.ini`, so you must include those default positional arguments (`scrapy tests`) after `--` as well:

```
tox -- scrapy tests -x # stop after first failure
```

You can also use the `pytest-xdist` plugin. For example, to run all tests on the Python 3.6 `tox` environment using all your CPU cores:

```
tox -e py36 -- scrapy tests -n auto
```

To see coverage report install `coverage` (`pip install coverage`) and run:

```
coverage report
```

see output of `coverage --help` for more options like `html` or `xml` report.

Writing tests

All functionality (including new features and bug fixes) must include a test case to check that it works as expected, so please include tests for your patches if you want them to get accepted sooner.

Scrapy uses unit-tests, which are located in the `tests/` directory. Their module name typically resembles the full path of the module they're testing. For example, the item loaders code is in:

```
scrapy.loader
```

And their unit-tests are in:

```
tests/test_loader.py
```

7.3 版本控制和 API 稳定性

7.3.1 Versioning

There are 3 numbers in a Scrapy version: *A.B.C*

- *A* is the major version. This will rarely change and will signify very large changes.
- *B* is the release number. This will include many changes including features and things that possibly break backward compatibility, although we strive to keep these cases at a minimum.
- *C* is the bugfix release number.

Backward-incompatibilities are explicitly mentioned in the [release notes](#), and may require special attention before upgrading.

Development releases do not follow 3-numbers version and are generally released as `dev` suffixed versions, e.g. `1.3dev`.

注解: With Scrapy 0.* series, Scrapy used [odd-numbered versions for development releases](#). This is not the case anymore from Scrapy 1.0 onwards.

Starting with Scrapy 1.0, all releases should be considered production-ready.

For example:

- `1.1.1` is the first bugfix release of the `1.1` series (safe to use in production)

7.3.2 API Stability

API stability was one of the major goals for the `1.0` release.

Methods or functions that start with a single dash (`_`) are private and should never be relied as stable.

Also, keep in mind that stable doesn't mean complete: stable APIs could grow new methods or functionality but the existing methods should keep working the same way.

Scrapy 更新 看看在最近的 Scrapy 版本中发生了什么变化。

贡献 Scrapy 代码 了解如何为 Scrapy 仓库贡献代码。

版本控制和 API 稳定性 了解 Scrapy 的版本控制和 API 的稳定性。

S

`scrapy.contracts`, 181
`scrapy.contracts.default`, 181
`scrapy.crawler`, 253
`scrapy.downloadermiddlewares`, 225
`scrapy.downloadermiddlewares.ajaxcrawl`, 239
`scrapy.downloadermiddlewares.cookies`, 227
`scrapy.downloadermiddlewares.defaultheaders`, 229
`scrapy.downloadermiddlewares.downloadtimeout`, 229
`scrapy.downloadermiddlewares.httppauth`, 229
`scrapy.downloadermiddlewares.httpcache`, 230
`scrapy.downloadermiddlewares.httpcompression`, 235
`scrapy.downloadermiddlewares.httpproxy`, 235
`scrapy.downloadermiddlewares.redirect`, 236
`scrapy.downloadermiddlewares.retry`, 237
`scrapy.downloadermiddlewares.robotstxt`, 239
`scrapy.downloadermiddlewares.stats`, 239
`scrapy.downloadermiddlewares.useragent`, 239
`scrapy.exceptions`, 154
`scrapy.exporters`, 262
`scrapy.extensions.closespider`, 251
`scrapy.extensions.corestats`, 250
`scrapy.extensions.debug`, 252
`scrapy.extensions.logstats`, 250
`scrapy.extensions.memdebug`, 251
`scrapy.extensions.memusage`, 250
`scrapy.extensions.statsmailer`, 252
`scrapy.extensions.telnet`, 166
`scrapy.http`, 110
`scrapy.item`, 74
`scrapy.linkextractors`, 124
`scrapy.linkextractors.lxmlhtml`, 124
`scrapy.loader`, 79
`scrapy.loader.processors`, 90
`scrapy.mail`, 163
`scrapy.pipelines.files`, 208
`scrapy.pipelines.images`, 210
`scrapy.selector`, 73
`scrapy.settings`, 255
`scrapy.signals`, 258
`scrapy.spidermiddlewares`, 241
`scrapy.spidermiddlewares.depth`, 243
`scrapy.spidermiddlewares.httperror`, 243
`scrapy.spidermiddlewares.offsite`, 244
`scrapy.spidermiddlewares.referer`, 245
`scrapy.spidermiddlewares.urllength`, 247
`scrapy.spiders`, 41
`scrapy.statscollectors`, 163
`scrapy.utils.log`, 161
`scrapy.utils.trackref`, 199

A

`adapt_response()` (*scrapy.spiders.XMLFeedSpider* 方法), 49

`add_css()` (*scrapy.loader.ItemLoader* 方法), 86

`add_value()` (*scrapy.loader.ItemLoader* 方法), 85

`add_xpath()` (*scrapy.loader.ItemLoader* 方法), 85

`adjust_request_args()` (*scrapy.contracts.Contract* 方法), 181

`AJAXCRAWL_ENABLED`
setting, 239

`AjaxCrawlMiddleware`
(*scrapy.downloadermiddlewares.ajaxcrawl* 中的类), 239

`allowed_domains` (*scrapy.spiders.Spider* 属性), 42

`AUTOTHROTTLING_DEBUG`
setting, 214

`AUTOTHROTTLING_ENABLED`
setting, 213

`AUTOTHROTTLING_MAX_DELAY`
setting, 213

`AUTOTHROTTLING_START_DELAY`
setting, 213

`AUTOTHROTTLING_TARGET_CONCURRENCY`
setting, 214

`AWS_ACCESS_KEY_ID`
setting, 129

`AWS_ENDPOINT_URL`
setting, 129

`AWS_REGION_NAME`
setting, 130

`AWS_SECRET_ACCESS_KEY`

setting, 129

`AWS_USE_SSL`

setting, 129

`AWS_VERIFY`

setting, 129

B

`BaseItemExporter` (*scrapy.exporters* 中的类), 265

`bench`

command, 39

`bindaddress`

reqmeta, 116

`body` (*scrapy.http.Request* 属性), 112

`body` (*scrapy.http.Response* 属性), 121

`body_as_unicode()` (*scrapy.http.TextResponse* 方法), 123

`BOT_NAME`

setting, 130

C

`check`

command, 34

`clear_stats()` (*scrapy.statscollectors.StatsCollector* 方法), 256

`close_spider()`, 99

`close_spider()` (*scrapy.statscollectors.StatsCollector* 方法), 256

`closed()` (*scrapy.spiders.Spider* 方法), 43

`CloseSpider`, 154

CLOSESPIDER_ERRORCOUNT

setting, 252

CLOSESPIDER_ITEMCOUNT

setting, 251

CLOSESPIDER_PAGECOUNT

setting, 252

CLOSESPIDER_TIMEOUT

setting, 251

command

bench, 39

check, 34

crawl, 34

edit, 35

fetch, 35

genspider, 33

list, 35

parse, 38

runspider, 39

settings, 39

shell, 37

startproject, 33

version, 39

view, 36

COMMANDS_MODULE

setting, 40

Compose (*scrapy.loader.processors* 中的类), 91

COMPRESSION_ENABLED

setting, 235

CONCURRENT_ITEMS

setting, 130

CONCURRENT_REQUESTS

setting, 130

CONCURRENT_REQUESTS_PER_DOMAIN

setting, 130

CONCURRENT_REQUESTS_PER_IP

setting, 130

context (*scrapy.loader.ItemLoader* 属性), 87

Contract (*scrapy.contracts* 中的类), 181

cookiejar

reqmeta, 227

COOKIES_DEBUG

setting, 228

COOKIES_ENABLED

setting, 228

CookiesMiddleware (*scrapy.downloadermiddlewares.cookies* 中的类), 227

copy() (*scrapy.http.Request* 方法), 113

copy() (*scrapy.http.Response* 方法), 121

CoreStats (*scrapy.extensions.corestats* 中的类), 250

crawl

command, 34

crawl() (*scrapy.crawler.Crawler* 方法), 254

Crawler (*scrapy.crawler* 中的类), 253

crawler (*scrapy.spiders.Spider* 属性), 42

CrawlSpider (*scrapy.spiders* 中的类), 46

css() (*scrapy.http.TextResponse* 方法), 123

CSVFeedSpider (*scrapy.spiders* 中的类), 50

CsvItemExporter (*scrapy.exporters* 中的类), 267

custom_settings (*scrapy.spiders.Spider* 属性), 42

D

default_input_processor

(*scrapy.loader.ItemLoader* 属性), 87

DEFAULT_ITEM_CLASS

setting, 131

default_item_class (*scrapy.loader.ItemLoader* 属性), 87

default_output_processor

(*scrapy.loader.ItemLoader* 属性), 87

DEFAULT_REQUEST_HEADERS

setting, 131

default_selector_class

(*scrapy.loader.ItemLoader* 属性), 88

DefaultHeadersMiddleware

(*scrapy.downloadermiddlewares.defaultheaders* 中的类), 229

delimiter (*scrapy.spiders.CSVFeedSpider* 属性), 50

DEPTH_LIMIT

setting, 131

DEPTH_PRIORITY

setting, 131

DEPTH_STATS_VERBOSE

setting, 132

- DepthMiddleware (*scrapy.spidermiddlewares.depth* 中的类), 243
- DNS_TIMEOUT
 - setting, 132
- DNSCACHE_ENABLED
 - setting, 132
- DNSCACHE_SIZE
 - setting, 132
- dont_cache
 - reqmeta, 230
- dont_merge_cookies
 - reqmeta, 111
- dont_obey_robotstxt
 - reqmeta, 239
- dont_redirect
 - reqmeta, 236
- dont_retry
 - reqmeta, 238
- DontCloseSpider, 154
- DOWNLOAD_DELAY
 - setting, 135
- DOWNLOAD_FAIL_ON_DATALOSS
 - setting, 137
- download_fail_on_dataloss
 - reqmeta, 116
- DOWNLOAD_HANDLERS
 - setting, 135
- DOWNLOAD_HANDLERS_BASE
 - setting, 135
- download_latency
 - reqmeta, 116
- DOWNLOAD_MAXSIZE
 - setting, 136
- download_maxsize
 - reqmeta, 136
- DOWNLOAD_TIMEOUT
 - setting, 136
- download_timeout
 - reqmeta, 116
- DOWNLOAD_WARN_SIZE
 - setting, 136
- DOWNLOADER
 - setting, 132
- DOWNLOADER_CLIENT_TLS_METHOD
 - setting, 133
- DOWNLOADER_CLIENTCONTEXTFACTORY
 - setting, 133
- DOWNLOADER_HTTPCLIENTFACTORY
 - setting, 132
- DOWNLOADER_MIDDLEWARES
 - setting, 134
- DOWNLOADER_MIDDLEWARES_BASE
 - setting, 134
- DOWNLOADER_STATS
 - setting, 134
- DownloaderMiddleware
 - (*scrapy.downloadermiddlewares* 中的类), 225
- DownloaderStats (*scrapy.downloadermiddlewares.stats* 中的类), 239
- DownloadTimeoutMiddleware
 - (*scrapy.downloadermiddlewares.downloadtimeout* 中的类), 229
- DropItem, 154
- DummyStatsCollector (*scrapy.statscollectors* 中的类), 163
- DUPEFILTER_CLASS
 - setting, 137
- DUPEFILTER_DEBUG
 - setting, 137
- E
- edit
 - command, 35
- EDITOR
 - setting, 137
- encoding (*scrapy.exporters.BaseItemExporter* 属性), 266
- encoding (*scrapy.http.TextResponse* 属性), 122
- engine (*scrapy.crawler.Crawler* 属性), 254
- engine_started
 - signal, 258
- engine_started() (在 *scrapy.signals* 模块中), 258
- engine_stopped

signal, 258
engine_stopped() (在 *scrapy.signals* 模块中), 258
export_empty_fields
(*scrapy.exporters.BaseItemExporter* 属性), 265
export_item() (*scrapy.exporters.BaseItemExporter* 方法), 265

EXTENSIONS

setting, 138

extensions (*scrapy.crawler.Crawler* 属性), 254

EXTENSIONS_BASE

setting, 138

F

FEED_EXPORT_ENCODING

setting, 108

FEED_EXPORT_FIELDS

setting, 108

FEED_EXPORT_INDENT

setting, 108

FEED_EXPORTERS

setting, 109

FEED_EXPORTERS_BASE

setting, 109

FEED_FORMAT

setting, 107

FEED_STORAGE_S3_ACL

setting, 109

FEED_STORAGES

setting, 108

FEED_STORAGES_BASE

setting, 109

FEED_STORE_EMPTY

setting, 108

FEED_TEMPDIR

setting, 138

FEED_URI

setting, 107

fetch

command, 35

Field (*scrapy.item* 中的类), 79

fields (*scrapy.item.Item* 属性), 78

fields_to_export (*scrapy.exporters.BaseItemExporter* 属性), 265

FILES_EXPIRES

setting, 206

FILES_RESULT_FIELD

setting, 205

FILES_STORE

setting, 203

FILES_STORE_GCS_ACL

setting, 205

FILES_STORE_S3_ACL

setting, 204

FILES_URLS_FIELD

setting, 205

FilesPipeline (*scrapy.pipelines.files* 中的类), 208

find_by_request() (*scrapy.loader.SpiderLoader* 方法), 255

finish_exporting() (*scrapy.exporters.BaseItemExporter* 方法), 265

flags (*scrapy.http.Response* 属性), 121

FormRequest (*scrapy.http* 中的类), 117

from_crawler(), 99

from_crawler() (*scrapy.downloadermiddlewares.DownloaderMiddleware* 方法), 227

from_crawler() (*scrapy.spidermiddlewares.SpiderMiddleware* 方法), 243

from_crawler() (*scrapy.spiders.Spider* 方法), 42

from_response() (*scrapy.http.FormRequest* 类方法), 117

from_settings() (*scrapy.loader.SpiderLoader* 方法), 255

from_settings() (*scrapy.mail.MailSender* 类方法), 164

FTP_PASSIVE_MODE

setting, 138

FTP_PASSWORD

setting, 138

FTP_USER

setting, 139

G

GCS_PROJECT_ID

- setting, 205
- genspider
 - command, 33
- get_collected_values()
 - (*scrapy.loader.ItemLoader* 方法), 87
- get_css() (*scrapy.loader.ItemLoader* 方法), 86
- get_input_processor() (*scrapy.loader.ItemLoader* 方法), 87
- get_media_requests()
 - (*scrapy.pipelines.files.FilesPipeline* 方法), 208
- get_media_requests()
 - (*scrapy.pipelines.images.ImagesPipeline* 方法), 210
- get_oldest() (在 *scrapy.utils.trackref* 模块中), 199
- get_output_processor()
 - (*scrapy.loader.ItemLoader* 方法), 87
- get_output_value() (*scrapy.loader.ItemLoader* 方法), 87
- get_stats() (*scrapy.statscollectors.StatsCollector* 方法), 256
- get_value() (*scrapy.loader.ItemLoader* 方法), 84
- get_value() (*scrapy.statscollectors.StatsCollector* 方法), 256
- get_xpath() (*scrapy.loader.ItemLoader* 方法), 85
- H
- handle_httpstatus_all
 - reqmeta, 244
- handle_httpstatus_list
 - reqmeta, 244
- headers (*scrapy.http.Request* 属性), 112
- headers (*scrapy.http.Response* 属性), 121
- headers (*scrapy.spiders.CSVFeedSpider* 属性), 50
- HtmlResponse (*scrapy.http* 中的类), 123
- HttpAuthMiddleware (*scrapy.downloadermiddlewares.httpauth* 中的类), 229
- HTTPCACHE_ALWAYS_STORE
 - setting, 234
- HTTPCACHE_DBM_MODULE
 - setting, 234
- HTTPCACHE_DIR
 - setting, 233
- HTTPCACHE_ENABLED
 - setting, 232
- HTTPCACHE_EXPIRATION_SECS
 - setting, 233
- HTTPCACHE_GZIP
 - setting, 234
- HTTPCACHE_IGNORE_HTTP_CODES
 - setting, 233
- HTTPCACHE_IGNORE_MISSING
 - setting, 233
- HTTPCACHE_IGNORE_RESPONSE_CACHE_CONTROLS
 - setting, 234
- HTTPCACHE_IGNORE_SCHEMES
 - setting, 233
- HTTPCACHE_POLICY
 - setting, 234
- HTTPCACHE_STORAGE
 - setting, 234
- HttpCacheMiddleware
 - (*scrapy.downloadermiddlewares.httppache* 中的类), 230
- HttpCompressionMiddleware
 - (*scrapy.downloadermiddlewares.httpcompression* 中的类), 235
- HTTPERROR_ALLOW_ALL
 - setting, 244
- HTTPERROR_ALLOWED_CODES
 - setting, 244
- HttpErrorMiddleware
 - (*scrapy.spidermiddlewares.httperror* 中的类), 243
- HTTPPROXY_AUTH_ENCODING
 - setting, 240
- HTTPPROXY_ENABLED
 - setting, 240
- HttpProxyMiddleware
 - (*scrapy.downloadermiddlewares.httpproxy* 中的类), 235
- I
- Identity (*scrapy.loader.processors* 中的类), 90

- IgnoreRequest, 154
- IMAGES_EXPIRES
 - setting, 206
- IMAGES_MIN_HEIGHT
 - setting, 208
- IMAGES_MIN_WIDTH
 - setting, 208
- IMAGES_RESULT_FIELD
 - setting, 205
- IMAGES_STORE
 - setting, 203
- IMAGES_STORE_GCS_ACL
 - setting, 205
- IMAGES_STORE_S3_ACL
 - setting, 204
- IMAGES_THUMBS
 - setting, 207
- IMAGES_URLS_FIELD
 - setting, 205
- ImagesPipeline (*scrapy.pipelines.images* 中的类), 210
- inc_value() (*scrapy.statscollectors.StatsCollector* 方法), 256
- indent (*scrapy.exporters.BaseItemExporter* 属性), 266
- Item (*scrapy.item* 中的类), 78
- item (*scrapy.loader.ItemLoader* 属性), 87
- item_completed() (*scrapy.pipelines.files.FilesPipeline* 方法), 209
- item_completed() (*scrapy.pipelines.images.ImagesPipeline* 方法), 210
- item_dropped
 - signal, 259
- item_dropped() (在 *scrapy.signals* 模块中), 259
- item_error
 - signal, 259
- item_error() (在 *scrapy.signals* 模块中), 259
- ITEM_PIPELINES
 - setting, 139
- ITEM_PIPELINES_BASE
 - setting, 139
- item_scraped
 - signal, 258
- item_scraped() (在 *scrapy.signals* 模块中), 258
- ItemLoader (*scrapy.loader* 中的类), 84
- iter_all() (在 *scrapy.utils.trackref* 模块中), 199
- iterator (*scrapy.spiders.XMLFeedSpider* 属性), 48
- itertag (*scrapy.spiders.XMLFeedSpider* 属性), 49
- J
- Join (*scrapy.loader.processors* 中的类), 90
- JsonItemExporter (*scrapy.exporters* 中的类), 268
- JsonLinesItemExporter (*scrapy.exporters* 中的类), 269
- JSONRequest (*scrapy.http* 中的类), 119
- L
- list
 - command, 35
- list() (*scrapy.loader.SpiderLoader* 方法), 255
- load() (*scrapy.loader.SpiderLoader* 方法), 255
- load_item() (*scrapy.loader.ItemLoader* 方法), 87
- log() (*scrapy.spiders.Spider* 方法), 43
- LOG_DATEFORMAT
 - setting, 140
- LOG_ENABLED
 - setting, 139
- LOG_ENCODING
 - setting, 139
- LOG_FILE
 - setting, 140
- LOG_FORMAT
 - setting, 140
- LOG_LEVEL
 - setting, 140
- LOG_SHORT_NAMES
 - setting, 140
- LOG_STDOUT
 - setting, 140
- logger (*scrapy.spiders.Spider* 属性), 42
- LogStats (*scrapy.extensions.logstats* 中的类), 250
- LOGSTATS_INTERVAL
 - setting, 140

LxmlLinkExtractor (*scrapy.linkextractors.lxmlhtml* 中的类), 124

M

MAIL_FROM
setting, 165

MAIL_HOST
setting, 165

MAIL_PASS
setting, 165

MAIL_PORT
setting, 165

MAIL_SSL
setting, 166

MAIL_TLS
setting, 166

MAIL_USER
setting, 165

MailSender (*scrapy.mail* 中的类), 164

MapCompose (*scrapy.loader.processors* 中的类), 91

max_retry_times
reqmeta, 116

max_value() (*scrapy.statscollectors.StatsCollector* 方法), 256

MEDIA_ALLOW_REDIRECTS
setting, 208

MEMDEBUG_ENABLED
setting, 141

MEMDEBUG_NOTIFY
setting, 141

MemoryStatsCollector (*scrapy.statscollectors* 中的类), 163

MEMUSAGE_CHECK_INTERVAL_SECONDS
setting, 141

MEMUSAGE_ENABLED
setting, 141

MEMUSAGE_LIMIT_MB
setting, 141

MEMUSAGE_NOTIFY_MAIL
setting, 142

MEMUSAGE_WARNING_MB
setting, 142

meta (*scrapy.http.Request* 属性), 112

meta (*scrapy.http.Response* 属性), 121

METAREFRESH_ENABLED
setting, 237

METAREFRESH_MAXDELAY
setting, 237

MetaRefreshMiddleware
(*scrapy.downloadermiddlewares.redirect* 中的类), 237

method (*scrapy.http.Request* 属性), 112

min_value() (*scrapy.statscollectors.StatsCollector* 方法), 256

N

name (*scrapy.spiders.Spider* 属性), 41

namespaces (*scrapy.spiders.XMLFeedSpider* 属性), 49

nested_css() (*scrapy.loader.ItemLoader* 方法), 87

nested_xpath() (*scrapy.loader.ItemLoader* 方法), 87

NEWSPIDER_MODULE
setting, 142

NotConfigured, 155

NotSupported, 155

O

object_ref (*scrapy.utils.trackref* 中的类), 199

OffsiteMiddleware (*scrapy.spidermiddlewares.offsite* 中的类), 244

open_spider(), 99

open_spider() (*scrapy.statscollectors.StatsCollector* 方法), 256

P

parse
command, 38

parse() (*scrapy.spiders.Spider* 方法), 43

parse_node() (*scrapy.spiders.XMLFeedSpider* 方法), 49

parse_row() (*scrapy.spiders.CSVFeedSpider* 方法), 50

PickleItemExporter (*scrapy.exporters* 中的类), 268

- `post_process()` (*scrapy.contracts.Contract* 方法), 182
`PprintItemExporter` (*scrapy.exporters* 中的类), 268
`pre_process()` (*scrapy.contracts.Contract* 方法), 181
`print_live_refs()` (在 *scrapy.utils.trackref* 模块中), 199
`process_exception()`
 (*scrapy.downloadermiddlewares.DownloaderMiddleware* 方法), 226
`process_item()`, 99
`process_request()` (*scrapy.downloadermiddlewares.DownloaderMiddleware* 方法), 225
`process_response()` (*scrapy.downloadermiddlewares.DownloaderMiddleware* 方法), 226
`process_results()` (*scrapy.spiders.XMLFeedSpider* 方法), 49
`process_spider_exception()`
 (*scrapy.spidermiddlewares.SpiderMiddleware* 方法), 242
`process_spider_input()`
 (*scrapy.spidermiddlewares.SpiderMiddleware* 方法), 241
`process_spider_output()`
 (*scrapy.spidermiddlewares.SpiderMiddleware* 方法), 241
`process_start_requests()`
 (*scrapy.spidermiddlewares.SpiderMiddleware* 方法), 242
`proxy`
 reqmeta, 235
 Python 提高建议
 PEP 8, 333
- ## Q
- `quotechar` (*scrapy.spiders.CSVFeedSpider* 属性), 50
- ## R
- `RANDOMIZE_DOWNLOAD_DELAY`
 setting, 142
`REACTOR_THREADPOOL_MAXSIZE`
 setting, 143
`REDIRECT_ENABLED`
 setting, 236
`REDIRECT_MAX_TIMES`
 setting, 143, 237
`REDIRECT_PRIORITY_ADJUST`
 setting, 143
`redirect_reasons`
 reqmeta, 236
`redirect_urls`
 reqmeta, 236
`RedirectMiddleware` (*scrapy.downloadermiddlewares.redirect* 中的类), 226
`REFERER_ENABLED`
 setting, 245
`RefererMiddleware` (*scrapy.spidermiddlewares.referer* 中的类), 245
`REFERRER_POLICY`
 setting, 245
`referrer_policy`
 reqmeta, 245
`replace()` (*scrapy.http.Request* 方法), 113
`replace()` (*scrapy.http.Response* 方法), 122
`replace_css()` (*scrapy.loader.ItemLoader* 方法), 87
`replace_value()` (*scrapy.loader.ItemLoader* 方法), 85
`replace_xpath()` (*scrapy.loader.ItemLoader* 方法), 86
`reqmeta`
 bindaddress, 116
 cookiejar, 227
 dont_cache, 230
 dont_merge_cookies, 111
 dont_obey_robotstxt, 239
 dont_redirect, 236
 dont_retry, 238
 download_fail_on_datatloss, 116
 download_latency, 116
 download_maxsize, 136
 download_timeout, 116
 handle_httpstatus_all, 244
 handle_httpstatus_list, 244
 max_retry_times, 116

- proxy, 235
 - redirect_reasons, 236
 - redirect_urls, 236
 - referrer_policy, 245
 - Request (*scrapy.http* 中的类), 110
 - request (*scrapy.http.Response* 属性), 121
 - request_dropped
 - signal, 261
 - request_dropped() (在 *scrapy.signals* 模块中), 261
 - request_reached_downloader
 - signal, 261
 - request_reached_downloader() (在 *scrapy.signals* 模块中), 261
 - request_scheduled
 - signal, 261
 - request_scheduled() (在 *scrapy.signals* 模块中), 261
 - Response (*scrapy.http* 中的类), 120
 - response_downloaded
 - signal, 262
 - response_downloaded() (在 *scrapy.signals* 模块中), 262
 - response_received
 - signal, 261
 - response_received() (在 *scrapy.signals* 模块中), 261
 - RETRY_ENABLED
 - setting, 238
 - RETRY_HTTP_CODES
 - setting, 238
 - RETRY_PRIORITY_ADJUST
 - setting, 143
 - RETRY_TIMES
 - setting, 238
 - RetryMiddleware (*scrapy.downloadermiddlewares.retry* 中的类), 237
 - ReturnsContract (*scrapy.contracts.default* 中的类), 181
 - ROBOTSTXT_OBEY
 - setting, 144
 - RobotsTxtMiddleware
 - (*scrapy.downloadermiddlewares.robotstxt* 中的类), 239
 - Rule (*scrapy.spiders* 中的类), 47
 - rules (*scrapy.spiders.CrawlSpider* 属性), 46
 - runspider
 - command, 39
- ## S
- SCHEDULER
 - setting, 144
 - SCHEDULER_DEBUG
 - setting, 144
 - SCHEDULER_DISK_QUEUE
 - setting, 144
 - SCHEDULER_MEMORY_QUEUE
 - setting, 144
 - SCHEDULER_PRIORITY_QUEUE
 - setting, 145
 - ScrapesContract (*scrapy.contracts.default* 中的类), 181
 - scrapy.contracts (模块), 181
 - scrapy.contracts.default (模块), 181
 - scrapy.crawler (模块), 253
 - scrapy.downloadermiddlewares (模块), 225
 - scrapy.downloadermiddlewares.ajaxcrawl (模块), 239
 - scrapy.downloadermiddlewares.cookies (模块), 227
 - scrapy.downloadermiddlewares.defaultheaders (模块), 229
 - scrapy.downloadermiddlewares.downloadtimeout (模块), 229
 - scrapy.downloadermiddlewares.httpauth (模块), 229
 - scrapy.downloadermiddlewares.httppcache (模块), 230
 - scrapy.downloadermiddlewares.httpcompression (模块), 235
 - scrapy.downloadermiddlewares.httpproxy (模块), 235
 - scrapy.downloadermiddlewares.redirect (模块), 236
 - scrapy.downloadermiddlewares.retry (模块), 237

`scrapy.downloadermiddlewares.robotstxt` (模块), 239

`scrapy.downloadermiddlewares.stats` (模块), 239

`scrapy.downloadermiddlewares.useragent` (模块), 239

`scrapy.exceptions` (模块), 154

`scrapy.exporters` (模块), 262

`scrapy.extensions.clospider` (模块), 251

`scrapy.extensions.clospider.CloseSpider` (*scrapy.extensions.clospider* 中的类), 251

`scrapy.extensions.corestats` (模块), 250

`scrapy.extensions.debug` (模块), 252

`scrapy.extensions.debug.Debugger` (*scrapy.extensions.debug* 中的类), 253

`scrapy.extensions.debug.StackTraceDump` (*scrapy.extensions.debug* 中的类), 252

`scrapy.extensions.logstats` (模块), 250

`scrapy.extensions.memdebug` (模块), 251

`scrapy.extensions.memdebug.MemoryDebugger` (*scrapy.extensions.memdebug* 中的类), 251

`scrapy.extensions.memusage` (模块), 250

`scrapy.extensions.memusage.MemoryUsage` (*scrapy.extensions.memusage* 中的类), 250

`scrapy.extensions.statsmailer` (模块), 252

`scrapy.extensions.statsmailer.StatsMailer` (*scrapy.extensions.statsmailer* 中的类), 252

`scrapy.extensions.telnet` (模块), 166, 250

`scrapy.extensions.telnet.TelnetConsole` (*scrapy.extensions.telnet* 中的类), 250

`scrapy.http` (模块), 110

`scrapy.item` (模块), 74

`scrapy.linkextractors` (模块), 124

`scrapy.linkextractors.lxmlhtml` (模块), 124

`scrapy.loader` (模块), 79, 255

`scrapy.loader.processors` (模块), 90

`scrapy.mail` (模块), 163

`scrapy.pipelines.files` (模块), 208

`scrapy.pipelines.images` (模块), 210

`scrapy.selector` (模块), 73

`scrapy.settings` (模块), 255

`scrapy.signals` (模块), 258

`scrapy.spidermiddlewares` (模块), 241

`scrapy.spidermiddlewares.depth` (模块), 243

`scrapy.spidermiddlewares.httperror` (模块), 243

`scrapy.spidermiddlewares.offsite` (模块), 244

`scrapy.spidermiddlewares.referer` (模块), 245

`scrapy.spidermiddlewares.urllength` (模块), 247

`scrapy.spiders` (模块), 41

`scrapy.statscollectors` (模块), 163, 256

`scrapy.utils.log` (模块), 161

`scrapy.utils.trackref` (模块), 199

`SelectJmes` (*scrapy.loader.processors* 中的类), 92

`selector` (*scrapy.http.TextResponse* 属性), 123

`selector` (*scrapy.loader.ItemLoader* 属性), 88

`send()` (*scrapy.mail.MailSender* 方法), 164

`serialize_field()` (*scrapy.exporters.BaseItemExporter* 方法), 265

`set_stats()` (*scrapy.statscollectors.StatsCollector* 方法), 256

`set_value()` (*scrapy.statscollectors.StatsCollector* 方法), 256

setting

- `AJAXCRAWL_ENABLED`, 239
- `AUTOTHROTTLING_DEBUG`, 214
- `AUTOTHROTTLING_ENABLED`, 213
- `AUTOTHROTTLING_MAX_DELAY`, 213
- `AUTOTHROTTLING_START_DELAY`, 213
- `AUTOTHROTTLING_TARGET_CONCURRENCY`, 214
- `AWS_ACCESS_KEY_ID`, 129
- `AWS_ENDPOINT_URL`, 129
- `AWS_REGION_NAME`, 130
- `AWS_SECRET_ACCESS_KEY`, 129
- `AWS_USE_SSL`, 129
- `AWS_VERIFY`, 129
- `BOT_NAME`, 130
- `CLOSESPIDER_ERRORCOUNT`, 252
- `CLOSESPIDER_ITEMCOUNT`, 251
- `CLOSESPIDER_PAGECOUNT`, 252
- `CLOSESPIDER_TIMEOUT`, 251
- `COMMANDS_MODULE`, 40
- `COMPRESSION_ENABLED`, 235
- `CONCURRENT_ITEMS`, 130
- `CONCURRENT_REQUESTS`, 130

CONCURRENT_REQUESTS_PER_DOMAIN, 130
CONCURRENT_REQUESTS_PER_IP, 130
COOKIES_DEBUG, 228
COOKIES_ENABLED, 228
DEFAULT_ITEM_CLASS, 131
DEFAULT_REQUEST_HEADERS, 131
DEPTH_LIMIT, 131
DEPTH_PRIORITY, 131
DEPTH_STATS_VERBOSE, 132
DNS_TIMEOUT, 132
DNSCACHE_ENABLED, 132
DNSCACHE_SIZE, 132
DOWNLOAD_DELAY, 135
DOWNLOAD_FAIL_ON_DATALOSS, 137
DOWNLOAD_HANDLERS, 135
DOWNLOAD_HANDLERS_BASE, 135
DOWNLOAD_MAXSIZE, 136
DOWNLOAD_TIMEOUT, 136
DOWNLOAD_WARN_SIZE, 136
DOWNLOADER, 132
DOWNLOADER_CLIENT_TLS_METHOD, 133
DOWNLOADER_CLIENTCONTEXTFACTORY, 133
DOWNLOADER_HTTPCLIENTFACTORY, 132
DOWNLOADER_MIDDLEWARES, 134
DOWNLOADER_MIDDLEWARES_BASE, 134
DOWNLOADER_STATS, 134
DUPEFILTER_CLASS, 137
DUPEFILTER_DEBUG, 137
EDITOR, 137
EXTENSIONS, 138
EXTENSIONS_BASE, 138
FEED_EXPORT_ENCODING, 108
FEED_EXPORT_FIELDS, 108
FEED_EXPORT_INDENT, 108
FEED_EXPORTERS, 109
FEED_EXPORTERS_BASE, 109
FEED_FORMAT, 107
FEED_STORAGE_S3_ACL, 109
FEED_STORAGES, 108
FEED_STORAGES_BASE, 109
FEED_STORE_EMPTY, 108
FEED_TEMPDIR, 138
FEED_URI, 107
FILES_EXPIRES, 206
FILES_RESULT_FIELD, 205
FILES_STORE, 203
FILES_STORE_GCS_ACL, 205
FILES_STORE_S3_ACL, 204
FILES_URLS_FIELD, 205
FTP_PASSIVE_MODE, 138
FTP_PASSWORD, 138
FTP_USER, 139
GCS_PROJECT_ID, 205
HTTPCACHE_ALWAYS_STORE, 234
HTTPCACHE_DBM_MODULE, 234
HTTPCACHE_DIR, 233
HTTPCACHE_ENABLED, 232
HTTPCACHE_EXPIRATION_SECS, 233
HTTPCACHE_GZIP, 234
HTTPCACHE_IGNORE_HTTP_CODES, 233
HTTPCACHE_IGNORE_MISSING, 233
HTTPCACHE_IGNORE_RESPONSE_CACHE_CONTROLS, 234
HTTPCACHE_IGNORE_SCHEMES, 233
HTTPCACHE_POLICY, 234
HTTPCACHE_STORAGE, 234
HTTPERROR_ALLOW_ALL, 244
HTTPERROR_ALLOWED_CODES, 244
HTTPPROXY_AUTH_ENCODING, 240
HTTPPROXY_ENABLED, 240
IMAGES_EXPIRES, 206
IMAGES_MIN_HEIGHT, 208
IMAGES_MIN_WIDTH, 208
IMAGES_RESULT_FIELD, 205
IMAGES_STORE, 203
IMAGES_STORE_GCS_ACL, 205
IMAGES_STORE_S3_ACL, 204
IMAGES_THUMBS, 207
IMAGES_URLS_FIELD, 205
ITEM_PIPELINES, 139
ITEM_PIPELINES_BASE, 139
LOG_DATEFORMAT, 140
LOG_ENABLED, 139
LOG_ENCODING, 139

LOG_FILE, 140
LOG_FORMAT, 140
LOG_LEVEL, 140
LOG_SHORT_NAMES, 140
LOG_STDOUT, 140
LOGSTATS_INTERVAL, 140
MAIL_FROM, 165
MAIL_HOST, 165
MAIL_PASS, 165
MAIL_PORT, 165
MAIL_SSL, 166
MAIL_TLS, 166
MAIL_USER, 165
MEDIA_ALLOW_REDIRECTS, 208
MEMDEBUG_ENABLED, 141
MEMDEBUG_NOTIFY, 141
MEMUSAGE_CHECK_INTERVAL_SECONDS, 141
MEMUSAGE_ENABLED, 141
MEMUSAGE_LIMIT_MB, 141
MEMUSAGE_NOTIFY_MAIL, 142
MEMUSAGE_WARNING_MB, 142
METAREFRESH_ENABLED, 237
METAREFRESH_MAXDELAY, 237
NEWSPIDER_MODULE, 142
RANDOMIZE_DOWNLOAD_DELAY, 142
REACTOR_THREADPOOL_MAXSIZE, 143
REDIRECT_ENABLED, 236
REDIRECT_MAX_TIMES, 143, 237
REDIRECT_PRIORITY_ADJUST, 143
REFERER_ENABLED, 245
REFERRER_POLICY, 245
RETRY_ENABLED, 238
RETRY_HTTP_CODES, 238
RETRY_PRIORITY_ADJUST, 143
RETRY_TIMES, 238
ROBOTSTXT_OBEY, 144
SCHEDULER, 144
SCHEDULER_DEBUG, 144
SCHEDULER_DISK_QUEUE, 144
SCHEDULER_MEMORY_QUEUE, 144
SCHEDULER_PRIORITY_QUEUE, 145
SPIDER_CONTRACTS, 145

SPIDER_CONTRACTS_BASE, 145
SPIDER_LOADER_CLASS, 145
SPIDER_LOADER_WARN_ONLY, 146
SPIDER_MIDDLEWARES, 146
SPIDER_MIDDLEWARES_BASE, 146
SPIDER_MODULES, 146
STATS_CLASS, 147
STATS_DUMP, 147
STATSMAILER_RCPTS, 147
TELNETCONSOLE_ENABLED, 147
TELNETCONSOLE_HOST, 169
TELNETCONSOLE_PASSWORD, 169
TELNETCONSOLE_PORT, 147, 169
TELNETCONSOLE_USERNAME, 169
TEMPLATES_DIR, 147
URLLENGTH_LIMIT, 148
USER_AGENT, 148

settings
 command, 39

settings (*scrapy.crawler.Crawler* 属性), 253
settings (*scrapy.spiders.Spider* 属性), 42
SETTINGS_PRIORITIES() (在 *scrapy.settings* 模块中), 255

shell
 command, 37

signal
 engine_started, 258
 engine_stopped, 258
 item_dropped, 259
 item_error, 259
 item_scraped, 258
 request_dropped, 261
 request_reached_downloader, 261
 request_scheduled, 261
 response_downloaded, 262
 response_received, 261
 spider_closed, 259
 spider_error, 260
 spider_idle, 260
 spider_opened, 260
 update_telnet_vars, 169

signals (*scrapy.crawler.Crawler* 属性), 254

- sitemap_alternate_links
(*scrapy.spiders.SitemapSpider* 属性), 52
- sitemap_filter() (*scrapy.spiders.SitemapSpider* 方法), 52
- sitemap_follow (*scrapy.spiders.SitemapSpider* 属性), 52
- sitemap_rules (*scrapy.spiders.SitemapSpider* 属性), 51
- sitemap_urls (*scrapy.spiders.SitemapSpider* 属性), 51
- SitemapSpider (*scrapy.spiders* 中的类), 51
- spider (*scrapy.crawler.Crawler* 属性), 254
- Spider (*scrapy.spiders* 中的类), 41
- spider_closed
 signal, 259
- spider_closed() (在 *scrapy.signals* 模块中), 259
- SPIDER_CONTRACTS
 setting, 145
- SPIDER_CONTRACTS_BASE
 setting, 145
- spider_error
 signal, 260
- spider_error() (在 *scrapy.signals* 模块中), 260
- spider_idle
 signal, 260
- spider_idle() (在 *scrapy.signals* 模块中), 260
- SPIDER_LOADER_CLASS
 setting, 145
- SPIDER_LOADER_WARN_ONLY
 setting, 146
- SPIDER_MIDDLEWARES
 setting, 146
- SPIDER_MIDDLEWARES_BASE
 setting, 146
- SPIDER_MODULES
 setting, 146
- spider_opened
 signal, 260
- spider_opened() (在 *scrapy.signals* 模块中), 260
- spider_stats (*scrapy.statscollectors.MemoryStatsCollector* 属性), 163
- SpiderLoader (*scrapy.loader* 中的类), 255
- SpiderMiddleware (*scrapy.spidermiddlewares* 中的类), 241
- start_exporting() (*scrapy.exporters.BaseItemExporter* 方法), 265
- start_requests() (*scrapy.spiders.Spider* 方法), 42
- start_urls (*scrapy.spiders.Spider* 属性), 42
- startproject
 command, 33
- stats (*scrapy.crawler.Crawler* 属性), 254
- STATS_CLASS
 setting, 147
- STATS_DUMP
 setting, 147
- StatsCollector (*scrapy.statscollectors* 中的类), 256
- STATSMAILER_RCPTS
 setting, 147
- status (*scrapy.http.Response* 属性), 121
- T**
- TakeFirst (*scrapy.loader.processors* 中的类), 90
- TELNETCONSOLE_ENABLED
 setting, 147
- TELNETCONSOLE_HOST
 setting, 169
- TELNETCONSOLE_PASSWORD
 setting, 169
- TELNETCONSOLE_PORT
 setting, 147, 169
- TELNETCONSOLE_USERNAME
 setting, 169
- TEMPLATES_DIR
 setting, 147
- text (*scrapy.http.TextResponse* 属性), 122
- TextResponse (*scrapy.http* 中的类), 122
- U**
- update_telnet_vars
 signal, 169
- update_telnet_vars() (在 *scrapy.extensions.telnet* 模块中), 169
- url (*scrapy.http.Request* 属性), 112

`url` (*scrapy.http.Response* 属性), 120

`UrlContract` (*scrapy.contracts.default* 中的类), 181

`urljoin()` (*scrapy.http.Response* 方法), 122

`URLLENGTH_LIMIT`

setting, 148

`UrlLengthMiddleware`

(*scrapy.spidermiddlewares.urllength* 中的类), 247

`USER_AGENT`

setting, 148

`UserAgentMiddleware`

(*scrapy.downloadermiddlewares.useragent* 中的类), 239

V

`version`

command, 39

`view`

command, 36

X

`XMLFeedSpider` (*scrapy.spiders* 中的类), 48

`XmlItemExporter` (*scrapy.exporters* 中的类), 266

`XmlResponse` (*scrapy.http* 中的类), 123

`xpath()` (*scrapy.http.TextResponse* 方法), 123