



链滴

# Python3 网络爬虫实战：10、爬虫框架的安装：PySpider、Scrapy

作者：[zhaolixiang](#)

原文链接：<https://ld246.com/article/1566383608762>

来源网站：[链滴](#)

许可协议：[署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

<p>我们直接用 Requests、Selenium 等库写爬虫，如果爬取量不是太大，速度要求不高，是完全以满足需求的。但是写多了会发现其内部许多代码和组件是可以复用的，如果我们把这些组件抽离出，将各个功能模块化，就慢慢会形成一个框架雏形，久而久之，爬虫框架就诞生了。</p>

<p>利用框架我们可以不用再去关心某些功能的具体实现，只需要去关心爬取逻辑即可。有了它们，以大大简化代码量，而且架构也会变得清晰，爬取效率也会高许多。所以如果对爬虫有一定基础，上框架是一种好的选择。</p>

<p>本书主要介绍的爬虫框架有 PySpider 和 Scrapy，本节我们来介绍一下 PySpider、Scrapy 以它们的一些扩展库的安装方式。</p>

## 

<p>PySpider 是国人 binux 编写的强大的网络爬虫框架，它带有强大的 WebUI、脚本编辑器、任务监控器、项目管理器以及结果处理器，同时它支持多种数据库后端、多种消息队列，另外它还支持 JavaScript 渲染页面的爬取，使用起来非常方便，本节介绍一下它的安装过程。</p>

## 

<ul>

<li>官方文档：<a href="https://ld246.com/forward?goto=http%3A%2F%2Fdocs.pyspider.org%2F" target="\_blank" rel="nofollow ugc">http://docs.pyspider.org/</a> </li>

<li>PyPi：<a href="https://ld246.com/forward?goto=https%3A%2F%2Fpypi.python.org%2Fpypi%2Fpyspider" target="\_blank" rel="nofollow ugc">https://pypi.python.org/pypi/pyspider</a> </li>

<li>GitHub：<a href="https://ld246.com/forward?goto=https%3A%2F%2Fgithub.com%2Fbinux%2Fpyspider" target="\_blank" rel="nofollow ugc">https://github.com/binux/pyspider</a> </li>

<li>官方教程：<a href="https://ld246.com/forward?goto=http%3A%2F%2Fdocs.pyspider.org%2Fen%2Flatest%2Ftutorial" target="\_blank" rel="nofollow ugc">http://docs.pyspider.org/en/latest/tutorial</a> </li>

<li>在线实例：<a href="https://ld246.com/forward?goto=http%3A%2F%2Fdemo.pyspider.org" target="\_blank" rel="nofollow ugc">http://demo.pyspider.org</a> </li>

</ul>

## 

<p>PySpider 是支持 JavaScript 渲染的，而这个过程是依赖于 PhantomJS 的，所以还需要安装 PhantomJS，所以在安装之前请安装好 PhantomJS，安装方式在前文有介绍。</p>

## 

<p>推荐使用 Pip 安装，命令如下：</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install pyspider</span></span></code></pre>
```

<p>命令执行完毕即可完成安装。</p>

## 

<p>Windows 下可能会出现这样的错误提示：Command "python setup.py egg\_info" failed with error code 1 in /tmp/pip-build-vXo1W3/pycurl<br>

这个是 PyCurl 安装错误，一般会出现在 Windows 下，需要安装 PyCurl 库，下载链接为：<a href="https://ld246.com/forward?goto=http%3A%2F%2Fwww.lfd.uci.edu%2F%7Egohlke%2Fpythonlibs%2F%23pycurl" target="\_blank" rel="nofollow ugc">http://www.lfd.uci.edu/~gohlke/pythonlibs/#pycurl</a>，找到对应 Python 版本然后下载相应的 Wheel 文件。<br>

如 Windows 64 位，Python3.6 则下载 pycurl-7.43.0-cp36-cp36m-win\_amd64.whl，随后用 Pip 安装即可，命令如下：</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install pycurl-7.43.0-cp36-cp36m-win_amd64.whl</span></span></code></pre>
```

<p>Linux 下如果遇到 PyCurl 的错误可以参考本文：<a href="https://ld246.com/forward?goto=https%3A%2F%2Fimlonghao.com%2F19.html" target="\_blank" rel="nofollow ugc">https://imlonghao.com/19.html</a> </p>

<p>Mac 遇到这种情况执行下面操作：</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl"></span></span></code></pre>
```

```
cl">brew install openssl
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">openssl version
</span></span><span class="highlight-line"><span class="highlight-cl">查看版本
</span></span><span class="highlight-line"><span class="highlight-cl">find /usr/local -n
me ssl.h
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">可以看到形如:
</span></span><span class="highlight-line"><span class="highlight-cl">usr/local/Cellar/o
enssl/1.0.2s/include/openssl/ssl.h
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">添加环境变量
</span></span><span class="highlight-line"><span class="highlight-cl">export PYCURL_S
L_LIBRARY=openssl
</span></span><span class="highlight-line"><span class="highlight-cl">export LDFLAGS=
L/usr/local/Cellar/openssl/1.0.2s/lib
</span></span><span class="highlight-line"><span class="highlight-cl">export CPPFLAGS
-l/usr/local/Cellar/openssl/1.0.2s/include
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl"> pip3 install pyspi
er
</span></span></code></pre>
```

## 5. 验证安装

安装完成之后，可以直接在命令行下启动 PySpider: </p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">pyspider all
</span></span></code></pre>
```

控制台会有类似如下输出，如图 1-75 所示: <br>

</p>

图 1-75 控制台<br>

这时 PySpider 的 Web 服务就会在本地 5000 端口运行，直接在浏览器打开: http://localhost:5000/ 即可进入 PySpider 的 WebUI 管理页面，如图 1-76 所示: <br>

</p>

图 1-76 管理页面<br>

如果出现类似页面那证明 PySpider 已经安装成功了。<br>

在后文会介绍 PySpider 的详细用法。</p>

这里有一个深坑，PySpider 在 Python3.7 上运行时会报错</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">File "/usr/local/lib/python3.7/site-packages/pyspider/run.py", line 231
</span></span><span class="highlight-line"><span class="highlight-cl">    async=True, get
object=False, no_input=False):
</span></span><span class="highlight-line"><span class="highlight-cl">        ^
</span></span><span class="highlight-line"><span class="highlight-cl">SyntaxError: invali
syntax
</span></span></code></pre>
```

原因是 python3.7 中 async 已经变成了关键字。因此出现这个错误。<br>

修改方式是手动替换一下</p>

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight
cl">下面位置的async改为mark_async
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">/usr/local/lib/pyt
on3.7/site-packages/pyspider/run.py 的231行、245行 (两个)、365行
</span></span></code></pre>
```

```
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">/usr/local/lib/pyt
on3.7/site-packages/pyspider/webui/app.py 的95行
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span><span class="highlight-line"><span class="highlight-cl">/usr/local/lib/pyt
on3.7/site-packages/pyspider/fetcher/tornado_fetcher.py 的81行、89行 (两个) 、95行、117行
</span></span><span class="highlight-line"><span class="highlight-cl">
</span></span></code></pre>
```

## Scrapy的安装

Scrapy 是一个十分强大的爬虫框架，依赖的库比较多，至少需要依赖库有 Twisted 14.0, lxml 3.4, pyOpenSSL 0.14。而在不同平台环境又各不相同，所以在安装之前最好确保把一些基本库安装好。本节介绍一下 Scrapy 在不同平台的安装方法。

### 1. 相关链接

<ul>

<li>官方网站: <a href="https://ld246.com/forward?goto=https%3A%2F%2Fscrapy.org" target="\_blank" rel="nofollow ugc">https://scrapy.org</a></li>

<li>官方文档: <a href="https://ld246.com/forward?goto=https%3A%2F%2Fdocs.scrapy.org" target="\_blank" rel="nofollow ugc">https://docs.scrapy.org</a></li>

<li>PyPi: <a href="https://ld246.com/forward?goto=https%3A%2F%2Fpypi.python.org%2Fpackage%2FScrapy" target="\_blank" rel="nofollow ugc">https://pypi.python.org/pypi/Scrapy</a></li>

<li>GitHub: <a href="https://ld246.com/forward?goto=https%3A%2F%2Fgithub.com%2Fscrapy%2Fscrapy" target="\_blank" rel="nofollow ugc">https://github.com/scrapy/scrapy</a></li>

<li>中文文档: <a href="https://ld246.com/forward?goto=http%3A%2F%2Fscrapy-chs.readthedocs.io" target="\_blank" rel="nofollow ugc">http://scrapy-chs.readthedocs.io</a></li>

</ul>

### 3. Mac 下的安装

在 Mac 上构建 Scrapy 的依赖库需要 C 编译器以及开发头文件，它一般由 Xcode 提供，运行如命令安装即可：

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">xcode-select --install
</span></span></code></pre>
```

随后利用 Pip 安装 Scrapy 即可，运行如下命令：

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install Scrapy
</span></span></code></pre>
```

运行完毕之后即可完成 Scrapy 的安装。

### 4. 验证安装

安装之后，在命令行下输入 scrapy，如果出现类似下方的结果，就证明 Scrapy 安装成功，如图 1-80 所示：

</p>

图 1-80 验证安装

### 5. 常见错误

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pkg_resources.VersionConflict: (six 1.5.2 (/usr/lib/python3/dist-packages), Requirement.p
rse('six&gt;=1.6.0'))
</span></span></code></pre>
```

six 包版本过低，six 包是一个提供兼容 Python2 和 Python3 的库，升级 six 包即可：

```
<pre><code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">sudo pip3 install -U six
</span></span></code></pre>
```

---

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">c/_cffi_backend.c:15:17: fatal error: ffi.h: No such file or directory</span></span></code></pre>
```

<p>这是在 Linux 下常出现的错误，缺少 Libffi 这个库。什么是 libffi？“FFI”的全名是 Foreign Function Interface，通常指的是允许以一种语言编写的代码调用另一种语言的代码。而 Libffi 库提供最底层的、与架构相关的、完整的“FFI”。<br>

安装相应的库即可。<br>

Ubuntu、Debian: </p>

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">sudo apt-get install build-essential libssl-dev libffi-dev python3-dev</span></span></code></pre>
```

<p>CentOS、RedHat:</p>

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">sudo yum install gcc libffi-devel python-devel openssl-devel</span></span></code></pre>
```

---

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">Command "python setup.py egg_info" failed with error code 1 in /tmp/pip-build/cryptography/</span></span></code></pre>
```

</span></span></code></pre>

<p>这是缺少加密的相关组件，利用 Pip 安装即可。</p>

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install cryptography</span></span></code></pre>
```

---

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">ImportError: No module named 'packaging'</span></span></code></pre>
```

<p>缺少 packaging 这个包，它提供了 Python 包的核心功能，利用 Pip 安装即可。</p>

---

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">ImportError: No module named '_cffi_backend'</span></span></code></pre>
```

<p>缺少 cffi 包，使用 Pip 安装即可：</p>

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install cffi</span></span></code></pre>
```

---

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">ImportError: No module named 'pyparsing'</span></span></code></pre>
```

<p>缺少 pyparsing 包，使用 Pip 安装即可：</p>

```
<code class="highlight-chroma"><span class="highlight-line"><span class="highlight-cl">pip3 install pyparsing appdirs</span></span></code></pre>
```