

资深爬虫工程师专业奉献

黄永祥 著

基于Python 3编撰 ◀
从零基础到项目实战 ◀
提供技术交流与服务群 ◀

玩转Python 网络爬虫



本书实例
源代码

清华大学出版社



玩转 Python 网络爬虫



黄永祥 著

清华大学出版社
北京

内 容 简 介

本书站在初学者的角度，从原理到实践，循序渐进地讲述了使用Python开发网络爬虫的核心技术。全书从逻辑上可分为基础篇、实战篇和爬虫框架篇三部分。基础篇主要介绍了编写网络爬虫所需的基础知识，分别是网站分析、数据抓取、数据清洗和数据入库。网站分析讲述如何使用Chrome和Fiddler抓包工具对网络做全面分析；数据抓取介绍了Python爬虫模块Urllib和Requests的基础知识；数据清洗主要介绍字符串操作、正则和Beautiful Soup的使用；数据入库分别讲述了MySQL和MongoDB的操作，通过ORM框架SQLAlchemy实现数据持久化，实现企业级开发。实战篇深入讲解了分布式爬虫、爬虫软件开发与应用、12306抢票程序和微博爬取，所举示例均来自于开发实践，可帮助读者快速提升技能，开发实际项目。框架篇主要讲述Scrapy的基础知识，并通过爬取QQ音乐为实例，让读者深层次了解Scrapy的使用。

本书内容丰富，注重实战，适用于从零开始学习网络爬虫的初学者，或者是已经有一些网络爬虫编写经验，但希望更加全面、深入理解Python爬虫的开发人员。

本书封面贴有清华大学出版社防伪标签，无标签者不得销售。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目（CIP）数据

玩转Python网络爬虫/黄永祥著. —北京：清华大学出版社，2018（2018.11重印）
ISBN 978-7-302-50328-6

I. ①玩… II. ①黄… III. ①软件工具—程序设计 IV. ①TP311.56

中国版本图书馆CIP数据核字（2018）第114988号

责任编辑：王金柱

封面设计：王 翔

责任校对：闫秀华

责任印制：沈 露

出版发行：清华大学出版社

网 址：<http://www.tup.com.cn>，<http://www.wqbook.com>

地 址：北京清华大学学研大厦A座 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62776969，c-service@tup.tsinghua.edu.cn

质量反馈：010-62772015，zhiliang@tup.tsinghua.edu.cn

印 刷 者：北京富博印刷有限公司

装 订 者：北京市密云县京文制本装订厂

经 销：全国新华书店

开 本：180mm×230mm 印 张：20.25 字 数：454千字

版 次：2018年8月第1版 印 次：2018年11月第2次印刷

定 价：69.00元

产品编号：077679-01

前言

随着大数据和人工智能的普及，Python 的地位也变得水涨船高，许多技术人员投身于 Python 开发，其中网络爬虫是 Python 最为热门的应用领域之一。在爬虫领域，Python 可以说是处于霸主地位，Python 能解决爬虫开发过程中所遇到的难题，开发速度快且支持异步编程，大大缩短了开发周期。因此，Python 爬虫编程已成为爬虫工程师的必备技能。

本书结构

本书共分 18 章，各章内容概述如下：

第 1 章介绍什么是网络爬虫、爬虫的类型和原理、爬虫搜索策略和反爬虫技术以及解决方案。

第 2 章讲解爬虫开发的基础知识，包括 HTTP 协议、请求头和 Cookies 的作用、HTML 的布局结构、JavaScript 的介绍、JSON 的数据格式和 Ajax 的原理。

第 3 章介绍使用 Chrome 开发工具分析爬取网站，重点介绍开发者工具的 Elements 和 Network 标签的功能和使用方式，并通过开发者工具分析 QQ 网站。

第 4 章主要介绍 Fiddler 抓包工具的原理和安装配置，Fiddler 用户界面的各个功能及使用方法。

第 5 章讲述了 Urllib 在 Python 2 和 Python 3 的变化及使用，包括发送请求、使用代理 IP、Cookies 的读写、HTTP 证书验收和数据处理。

第 6 章介绍 Python 第三方库 Requests 的安装和使用，包括发送请求、使用代理 IP、Cookies 的读写、HTTP 证书验收和文件下载与上传。

第 7 章介绍验证码的种类和识别方法，包括 OCR 的安装和使用、验证码图片处理和使用第三方平台识别验证码。

第 8 章讲述数据清洗的三种方法，包括字符串操作（截取、查找、分割和替换）、正则表达式的使用和第三方库 Beautiful Soup 的安装以及使用。

第 9 章讲述如何将数据存储到文件，分别介绍了 CSV、Excel 和 Word 的读写方法及数据存储。

第 10 章介绍 ORM 框架 SQLAlchemy 的安装及使用，实现关系型数据库持久化存储数据，这是企业级的关系型数据库操作。

第 11 章讲述非关系型数据库 MongoDB 的操作，介绍 MongoDB 的安装、原理结构和 Python 实现 MongoDB 读写。

第 12 章介绍爬取淘宝商品信息实例，包括网站分析、数据抓取、数据清洗以及存储在 CSV 文件中，读者应掌握爬虫的开发流程。

第 13 章介绍爬取 QQ 音乐全站歌曲实例，包括网站分析、数据抓取和实现 SQLAlchemy 存储歌曲信息并下载文件，使用异步编程实现分布式开发，提高爬取效率。

第 14 章是在第 12 章的基础上实现爬虫软件开发，包括 PyQt5 的安装、使用 Qt Designer 设计软件界面、搭建 MVC 开发架构。

第 15 章实现 12306 抢票爬虫开发，包括用户登录、查询车次、预订车票、提交订单和生成订单的分析以及功能实现。

第 16 章介绍微博爬虫开发，包括微博登录、采集热门微博、发布微博、关注微博用户和转发评论的分析以及功能实现。

第 17 章介绍 Scrapy 爬虫框架，包括 Scrapy 的运行机制、安装、项目创建以及各个组件的编写（Setting、Items、Item Pipelines 和 Spider）和文件下载。

第 18 章介绍 Scrapy 爬取 QQ 音乐全站歌曲实例，包括编写 Spider 实现数据抓取、Item Pipelines 实现歌曲信息存储和歌曲下载、Items 定义数据存储对象和 Setting 配置项目设置。

本书特色

循序渐进，知识全面：本书站在初学者的角度，围绕 Python 网络爬虫开发展开讲解，从初学者必备基础知识着手，循序渐进地介绍了使用 Python 3 开发网络爬虫的各种知识，内容难度适中，由浅入深，实用性强，覆盖面广，条理清晰，且具有较强的逻辑性和系统性。

实例丰富，扩展性强：本书采用大量的实例进行讲解，力求通过实际操作使读者更容易地掌握爬虫开发。本书实例都经过作者精心设计和挑选，根据作者的实际开发经验总结而来，涵盖了在实际开发中所遇到的各种问题。对于精选案例，都尽可能做到步骤详尽、结构清晰、分析深入浅出，而且案例的扩展性强，读者可根据实际需求扩展开发。

基于理论，注重实践：在讲解过程中，不仅介绍理论知识，而且安排了综合应用实例或小型应用程序，将理论应用到实践中，加强读者的实际开发能力，巩固开发技能和相关知识。

源码提供

本书的实例源码可以在百度网盘下载。

链接地址 1: <https://pan.baidu.com/s/1htxBpic> 密码: aesy

链接地址 2: <https://pan.baidu.com/s/1E7axRN9rC0i9ASM1124IAw>

也可以扫描以下二维码下载。



如果你在下载过程中遇到问题，可发送邮件至 booksaga@126.com 获得帮助，邮件标题为“玩转 Python 网络爬虫下载资源”。

技术服务

读者在学习或者工作的过程中，如果遇到实际问题，可以加入 QQ 群 93314951 或 657341423 与笔者联系，笔者会在第一时间给予回复。

读者对象

本书主要适合以下读者阅读：

- Python 网络爬虫初学者及在校学生。
- Python 初级爬虫工程师。
- 从事数据抓取的技术人员。
- 其他学习 Python 网络爬虫的开发人员。

虽然笔者力求本书更臻完美，但由于水平所限，难免会出现错误，特别是实例中爬取的网站可能随时更新，导致源码在运行过程中出现问题，欢迎广大读者和高手专家给予指正，笔者将十分感谢。

编者

2018 年 1 月

目 录

第 1 章 理解网络爬虫	1
1.1 爬虫的定义	1
1.2 爬虫的类型	2
1.3 爬虫的原理	3
1.4 爬虫的搜索策略	5
1.5 反爬虫技术及解决方案	6
1.6 本章小结	8
第 2 章 爬虫开发基础	9
2.1 HTTP 与 HTTPS	9
2.2 请求头	11
2.3 Cookies	13
2.4 HTML	14
2.5 JavaScript	16
2.6 JSON	18
2.7 Ajax	19
2.8 本章小结	20
第 3 章 Chrome 分析网站	21
3.1 Chrome 开发工具	21
3.2 Elements 标签	22
3.3 Network 标签	23
3.4 分析 QQ 音乐	27
3.5 本章小结	29

第 4 章 Fiddler 抓包工具	30
4.1 Fiddler 介绍	30
4.2 Fiddler 安装配置	31
4.3 Fiddler 抓取手机应用	33
4.4 Toolbar 工具栏	36
4.5 Web Session 列表	37
4.6 View 选项视图	40
4.7 Quickexec 命令行	41
4.8 本章小结	42
第 5 章 Urllib 数据抓取	43
5.1 Urllib 简介	43
5.2 发送请求	44
5.3 复杂的请求	46
5.4 代理 IP	47
5.5 使用 Cookies	48
5.6 证书验证	50
5.7 数据处理	51
5.8 本章小结	52
第 6 章 Requests 数据抓取	54
6.1 Requests 简介及安装	54
6.2 请求方式	55
6.3 复杂的请求方式	57
6.4 下载与上传	60
6.5 本章小结	63
第 7 章 验证码识别	64
7.1 验证码类型	64
7.2 OCR 技术	66
7.3 第三方平台	69
7.4 本章小结	72

第 8 章 数据清洗	74
8.1 字符串操作.....	74
8.2 正则表达式.....	78
8.3 Beautiful Soup 介绍及安装.....	84
8.4 Beautiful Soup 的使用.....	86
8.5 本章小结.....	90
第 9 章 文档数据存储	92
9.1 CSV 数据写入和读取.....	92
9.2 Excel 数据写入和读取.....	94
9.3 Word 数据写入和读取.....	99
9.4 本章小结.....	101
第 10 章 ORM 框架	104
10.1 SQLAlchemy 介绍.....	104
10.2 安装 SQLAlchemy.....	105
10.3 连接数据库.....	106
10.4 创建数据表.....	108
10.5 添加数据.....	111
10.6 更新数据.....	112
10.7 查询数据.....	114
10.8 本章小结.....	116
第 11 章 MongoDB 数据库操作	118
11.1 MongoDB 介绍.....	118
11.2 安装及使用.....	120
11.2.1 MongoDB.....	120
11.2.2 MongoDB 可视化工具.....	121
11.2.3 PyMongo.....	123
11.3 连接数据库.....	123
11.4 添加文档.....	125

11.5 更新文档	126
11.6 查询文档	127
11.7 本章小结	130
第 12 章 项目实战：爬取淘宝商品信息	131
12.1 分析说明	131
12.2 功能实现	134
12.3 数据存储	136
12.4 本章小结	138
第 13 章 项目实战：分布式爬虫——QQ 音乐	139
13.1 分析说明	139
13.2 歌曲下载	140
13.3 歌手和歌曲信息	145
13.4 分类歌手列表	148
13.5 全站歌手列表	150
13.6 数据存储	152
13.7 分布式概念	154
13.7.1 GIL 是什么	154
13.7.2 为什么会有 GIL	154
13.8 并发库 concurrent.futures	155
13.9 分布式爬虫	157
13.10 本章小结	159
第 14 章 项目实战：爬虫软件——淘宝商品信息	161
14.1 分析说明	161
14.2 GUI 库介绍	162
14.3 PyQt5 安装及环境搭建	162
14.4 软件界面开发	165
14.5 MVC——视图	169
14.6 MVC——控制器	171
14.7 MVC——模型	172

14.8 扩展思路	173
14.9 本章小结	174
第 15 章 项目实战：12306 抢票	176
15.1 分析说明	176
15.2 验证码验证	177
15.3 用户登录与验证	181
15.4 查询车次	187
15.5 预订车票	193
15.6 提交订单	196
15.7 生成订单	204
15.8 本章小结	209
第 16 章 项目实战：玩转微博	219
16.1 分析说明	219
16.2 用户登录	220
16.3 用户登录（带验证码）	232
16.4 关键字搜索热门微博	240
16.5 发布微博	247
16.6 关注用户	253
16.7 点赞和转发评论	257
16.8 本章小结	263
第 17 章 Scrapy 爬虫框架	265
17.1 爬虫框架	265
17.2 Scrapy 的运行机制	267
17.3 安装 Scrapy	268
17.4 爬虫开发快速入门	270
17.5 Spiders 介绍	277
17.6 Spider 的编写	278
17.7 Items 的编写	282

17.8 Item Pipeline 的编写	284
17.9 Selectors 的编写	288
17.10 文件下载	291
17.11 本章小结	296
第 18 章 项目实战: Scrapy 爬取 QQ 音乐	298
18.1 分析说明	298
18.2 创建项目	299
18.3 编写 setting	300
18.4 编写 Items	301
18.5 编写 Item Pipelines	302
18.6 编写 Spider	305
18.7 本章小结	310

第 1 章

理解网络爬虫

1.1 爬虫的定义

网络爬虫是一种按照一定的规则自动地抓取网络信息的程序或者脚本。简单来说，网络爬虫就是根据一定的算法实现编程开发，主要通过 URL 实现数据的抓取和发掘。

随着大数据时代的发展，数据规模越来越庞大，数据类型繁多，但是数据价值普遍较低。为了从庞大的数据体系里获取有价值的信息，从而延伸了网络爬虫、数据分析等多个职位。近几年，网络爬虫的需求更是井喷式地爆发，在招聘的供求市场上往往是供不应求，造成这个现状的主要原因就是求职者的专业水平低于需求企业的要求。

传统的爬虫有百度、Google、必应等搜索引擎，这类通用的搜索引擎都有自己的核心算法。但是，这类通用的搜索引擎也存在着一一定的局限性：

(1) 不同的搜索引擎对于同一个搜索会有不同的结果，搜索出来的结果未必是用户需要的信息。

(2) 通用的引擎扩大网络覆盖率，但有限的搜索引擎服务器资源与无限的网络数据资源之间的矛盾将进一步加深。

(3) 随着网络上数据形式繁多和网络技术不断发展，图片、数据库、音频、视频多媒体等不同数据大量出现，通用搜索引擎往往对这些信息含量密集且具有一定结构的数据无能为力，不能很好地发现和获取。

因此，为了得到准确的数据，定向抓取相关网页资源的聚焦爬虫应运而生。聚焦爬虫是一个自动下载网页的程序，根据设定的抓取目标有目的地访问互联网上的网页与相关的 URL，从而获取所需要的信息。与通用爬虫不同，聚焦爬虫并不追求全面的覆盖率，而是抓取与某一特定内容相关的网页，为面向特定的用户提供准备数据资源。

1.2 爬虫的类型

网络爬虫根据系统结构和开发技术大致可以分为 4 种类型：通用网络爬虫、聚焦网络爬虫、增量式网络爬虫和深层网络爬虫。

通用网络爬虫又称全网爬虫，常见的有百度、Google、必应等搜索引擎，爬行对象从一些初始 URL 扩充到整个网站，主要为门户网站搜索引擎和大型网站服务采集数据，具有以下特点：

(1) 由于商业原因，引擎的算法是不会对外公布的。

(2) 这类网络爬虫的爬行范围和数量巨大，对于爬行速度和存储空间要求较高，爬行页面的顺序要求相对较低。

(3) 待刷新的页面太多，通常采用并行工作方式，但需要较长时间才能刷新一次页面。

(4) 存在一定缺陷，通用网络爬虫适用于为搜索引擎搜索广泛的需求。

聚焦网络爬虫又称主题网络爬虫，是选择性地爬行根据需求的主题相关页面的网络爬虫。与通用网络爬虫相比，聚焦爬虫只需要爬行与主题相关的页面，不需要广泛地覆盖无关的网页，很好地满足一些特定人群对特定领域信息的需求。

增量式网络爬虫是指对已下载网页采取增量式更新和只爬行新产生或者已经发生变化的网页的爬虫，它能够在一定程度上保证所爬行的页面尽可能是新的页面。只会在需要的时候爬行新产生或发生更新的页面，并不重新下载没有发生变化的页面，可有效减少数据下载量，及时更新已爬行的网页，减小时间和空间上的耗费，但是增加了爬行算法的复杂度和实现难度，基本上这类爬虫在实际开发中不太普及。

深层网络爬虫是大部分内容不能通过静态 URL 获取的、隐藏在搜索表单后的、只有用户提交一些关键词才能获得的网络页面。例如某些网站需要用户登录或者通过提交表单实现提交数据。这类爬虫也是本书主要讲述的重点之一。

这 4 种类型的爬虫大致上又可以分为两类，就是通用爬虫和聚焦爬虫，其中聚焦网络爬虫、增量式网络爬虫和深层网络爬虫可以通俗地归纳为一类，因为这类爬虫都是定向爬取数据。相比于通用爬虫，这类爬虫比较有目的性，也就是网络上经常说的网络爬虫，而通用爬虫在网络上通常称为搜索引擎。

1.3 爬虫的原理

通用网络爬虫的实现原理及过程如图 1-1 所示。

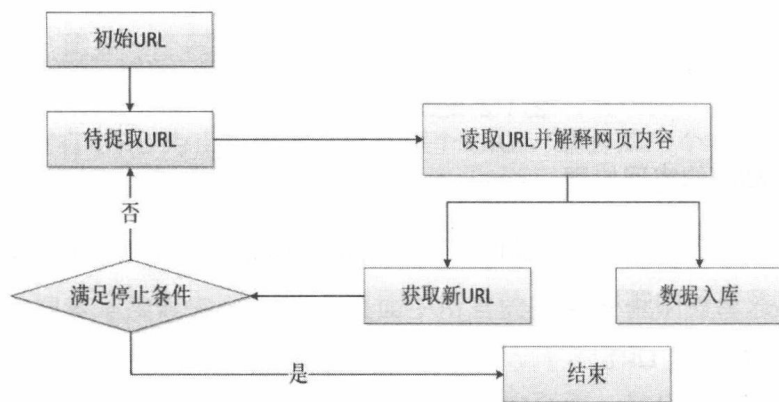


图 1-1 通用爬虫实现的原理及过程

通用网络爬虫的实现原理：

(1) 获取初始的 URL。初始的 URL 地址可以人为地指定，也可以由用户指定的某个或某几个初始爬取网页决定。

(2) 根据初始的 URL 爬取页面并获得新的 URL。获得初始的 URL 地址之后，先爬取当前 URL 地址中的网页信息，然后解析网页信息内容，将网页存储到原始数据库中，并且在当前获得的网页信息里发现新的 URL 地址，存放于一个 URL 队列里面。

(3) 从 URL 队列中读取新的 URL，从而获得新的网页信息，同时在新网页中获取新 URL，并重复上述的爬取过程。

(4) 满足爬虫系统设置的停止条件时，停止爬取。在编写爬虫的时候，一般会设置相应的停止条件，爬虫则会在停止条件满足时停止爬取。如果没有设置停止条件，爬虫就会一直爬取下去，一直到无法获取新的 URL 地址为止。

聚焦网络爬虫的执行原理和过程与通用爬虫大致相同，在通用爬虫的基础上增加两个步骤：定义爬取目标和筛选过滤 URL，原理如图 1-2 所示。

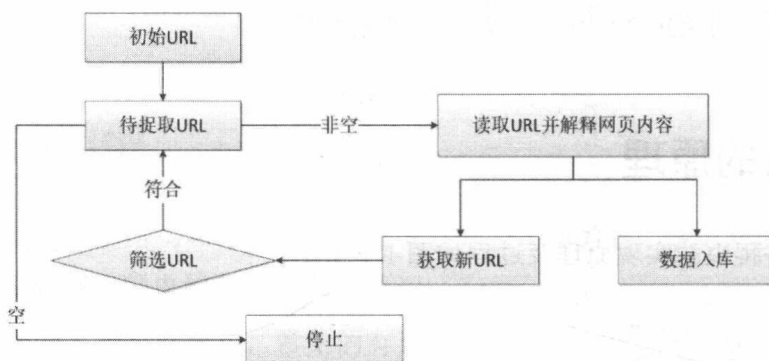


图 1-2 聚焦网络爬虫的原理

聚焦网络爬虫的实现原理：

(1) 制定爬取的方案。在聚焦网络爬虫中，首先要依据需求定义聚焦网络爬虫爬取的目标以及整体的爬取方案。

(2) 设定初始的 URL。

(3) 根据初始的 URL 抓取页面，并获得新的 URL。

(4) 从新的 URL 中过滤掉与需求无关的 URL，将过滤后的 URL 放到 URL 队列中。

(5) 在 URL 队列中，根据搜索算法确定 URL 的优先级，并确定下一步要爬取的 URL 地址。因为聚焦网络爬虫具有目的性，所以 URL 的爬取顺序不同会导致爬虫的执行效率不同。