



# 实战 Python 网络爬虫

从爬虫软件开发到自己动手开发爬虫框架

黄永祥 / 著

从原理到实践，深入浅出，热门爬虫核心技术全掌握  
涵盖丰富的爬虫工具、库、框架，十余个实战项目  
资深爬虫工程师倾力奉献，入门、进阶、求职必备



清华大学出版社



实战  
**Python**  
网络爬虫

黄永祥 / 著

清华大学出版社  
北京

## 内 容 简 介

本书从原理到实践,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,包括网站分析、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 抓包工具对网站做全面分析;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 BeautifulSoup 的使用;数据入库讲述了 MySQL 和 MongoDB 的操作,通过 ORM 框架 SQLAlchemy 实现数据持久化,进行企业级开发。实战篇深入讲解了分布式爬虫、爬虫软件的开发、12306 抢票程序和微博爬取等。框架篇主要讲述流行的爬虫框架 Scrapy,并以 Scrapy 与 Selenium、Splash、Redis 结合的项目案例,让读者深层次了解 Scrapy 的使用。此外,本书还介绍了爬虫的上线部署、如何自己动手开发一款爬虫框架、反爬虫技术的解决方案等内容。

本书使用 Python 3.X 编写,技术先进,项目丰富,适合欲从事爬虫工程师和数据分析师岗位的初学者、大学生和研究生使用,也很适合有一些网络爬虫编写经验,但希望更加全面、深入理解 Python 爬虫的开发人员使用。

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

实战 Python 网络爬虫 / 黄永祥著. —北京:清华大学出版社, 2019

ISBN 978-7-302-52489-2

I. ①实… II. ①黄… III. ①软件工具—程序设计 IV. ①TP311.561

中国版本图书馆 CIP 数据核字(2019)第 043080 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](mailto:c-service@tup.tsinghua.edu.cn)

质量反馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印装者:清华大学印刷厂

经 销:全国新华书店

开 本:190mm×260mm 印 张:30.25 字 数:774千字

版 次:2019年6月第1版 印 次:2019年6月第1次印刷

定 价:99.00元

---

产品编号:082567-01

# 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。此外，从事数据分析的工程师，为获取数据，很多时候也会用到网络爬虫的相关技术，因此，Python 爬虫编程已成为爬虫工程师和数据分析师的必备技能。

## 本书结构

本书共分 28 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、爬虫搜索策略和爬虫的合法性及开发流程。

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，重点介绍开发工具的 Elements 和 Network 标签的功能和使用方式，并通过开发工具分析 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~第 8 章介绍 Python 第三方库 Requests、Requests-Cache 爬虫缓存和 Requests-HTML，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、复杂的请求方式、缓存的存储机制、数据清洗以及 Ajax 动态数据爬取等内容。

第 9 章介绍网页操控和数据爬取，重点讲解 Selenium 的安装与使用，并通过实战项目“百度自动答题”，讲解了 Selenium 的使用。

第 10 章介绍手机 App 数据爬取，包括 Appium 的原理与开发环境搭建、连接 Android 系统，并通过实战项目“淘宝商品采集”，介绍了 App 数据的爬取技巧。

第 11 章介绍 Splash、Mitmproxy 与 Aiohttp 的安装和使用，包括 Splash 动态数据抓取、Mitmproxy 抓包和 Aiohttp 高并发抓取。

第 12 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 13 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 BeautifulSoup 的安装以及使用。

第 14 章讲述如何将数据存储到文件，包括 CSV、Excel 和 Word 文件的读取和写入方法。

第 15 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。

第 16 章讲述非关系型数据库 MongoDB 的操作，包括 MongoDB 的安装、原理和 Python 实现

MongoDB 的读写。

第 17 章至第 21 章介绍了 5 个实战项目，分别是：爬取 51Job 招聘信息、分布式爬虫——QQ 音乐、12306 抢票爬虫、微博爬取和微博爬虫软件的开发。

第 22 章至第 25 章介绍了 Scrapy 爬虫框架，包括 Scrapy 的运行机制、项目创建、各个组件的编写（Setting、Items、Item Pipelines 和 Spider）和文件下载及 Scrapy 中间件，并通过实战项目“Scrapy+Selenium 爬取豆瓣电影评论”、“Scrapy+Splash 爬取 B 站动漫信息”和“Scrapy+Redis 分布式爬取猫眼排行榜”、“爬取链家楼盘信息”和“QQ 音乐全站爬取”，深入讲解了 Scrapy 的应用和分布式爬虫的编写技巧。

第 26 章介绍爬虫的上线部署，包括非框架式爬虫和框架式爬虫的部署技巧。

第 27 章介绍常见的反爬虫技术，并给出了可行的反爬虫解决方案。

第 28 章介绍爬虫框架的编写，学习如何自己动手编写一款爬虫框架，以满足特定业务场景的需求。

### 本书特色

**循序渐进，涉及面广：**本书站在初学者的角度，循序渐进地介绍了使用 Python 开发网络爬虫的各种知识，内容由浅入深，几乎涵盖了目前网络爬虫开发的各种热门工具和前瞻性技术。

**实战项目丰富，扩展性强：**本书采用大量的实战项目进行讲解，力求通过实际应用使读者更容易地掌握爬虫开发技术，以应对业务需求。本书项目经过编者精心设计和挑选，根据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。对于精选项目，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。

**从理论到实践，注重培养爬虫开发思维：**在讲解过程中，不仅介绍理论知识，注重培养读者的爬虫开发思维，而且安排了综合应用实例或小型应用程序，使读者能顺利地将理论应用到实践中。

**特色干货，倾情分享：**本书大部分内容都来自作者多年来的编程实践，操作性很强。值得关注的是，本书还介绍了爬虫软件和爬虫框架的开发，供学有余力的读者扩展知识结构，提升开发技能。

### 源代码下载

本书所有程序代码均在 Python 3.6 下调试通过，源代码 Github 下载地址：

<https://github.com/xyjw/python-Reptile>

你也可以扫描下面的二维码下载。



如果你在下载过程中遇到问题，可发送邮件至 [554301449@qq.com](mailto:554301449@qq.com) 获得帮助，邮件标题为“实

战 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 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。
- 从事数据抓取和分析的技术人员。
- 学习 Python 程序设计的开发人员。

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。

黄永祥

2019 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 .....	2
1.4 爬虫的搜索策略 .....	4
1.5 爬虫的合法性与开发流程 .....	5
1.6 本章小结 .....	6
第 2 章 爬虫开发基础 .....	7
2.1 HTTP 与 HTTPS .....	7
2.2 请求头 .....	9
2.3 Cookies .....	10
2.4 HTML .....	11
2.5 JavaScript .....	12
2.6 JSON .....	14
2.7 Ajax .....	14
2.8 本章小结 .....	15
第 3 章 Chrome 分析网站 .....	16
3.1 Chrome 开发工具 .....	16
3.2 Elements 标签 .....	17
3.3 Network 标签 .....	18
3.4 分析 QQ 音乐 .....	20
3.5 本章小结 .....	23
第 4 章 Fiddler 抓包 .....	24
4.1 Fiddler 介绍 .....	24
4.2 Fiddler 安装配置 .....	24
4.3 Fiddler 抓取手机应用 .....	26
4.4 Toolbar 工具栏 .....	29
4.5 Web Session 列表 .....	30
4.6 View 选项视图 .....	32
4.7 Quickexec 命令行 .....	33



4.8	本章小结 .....	34
<b>第 5 章</b>	<b>爬虫库 Urllib.....</b>	<b>35</b>
5.1	Urllib 简介 .....	35
5.2	发送请求 .....	36
5.3	复杂的请求 .....	37
5.4	代理 IP .....	38
5.5	使用 Cookies .....	39
5.6	证书验证 .....	40
5.7	数据处理 .....	41
5.8	本章小结 .....	42
<b>第 6 章</b>	<b>爬虫库 Requests.....</b>	<b>43</b>
6.1	Requests 简介及安装 .....	43
6.2	请求方式 .....	44
6.3	复杂的请求方式 .....	45
6.4	下载与上传 .....	47
6.5	本章小结 .....	49
<b>第 7 章</b>	<b>Requests-Cache 爬虫缓存 .....</b>	<b>50</b>
7.1	简介及安装 .....	50
7.2	在 Requests 中使用缓存 .....	50
7.3	缓存的存储机制 .....	53
7.4	本章小结 .....	54
<b>第 8 章</b>	<b>爬虫库 Requests-HTML.....</b>	<b>55</b>
8.1	简介及安装 .....	55
8.2	请求方式 .....	56
8.3	数据清洗 .....	56
8.4	Ajax 动态数据抓取 .....	59
8.5	本章小结 .....	61
<b>第 9 章</b>	<b>网页操控与数据爬取 .....</b>	<b>62</b>
9.1	了解 Selenium .....	62
9.2	安装 Selenium .....	63
9.3	网页元素定位 .....	66
9.4	网页元素操控 .....	70
9.5	常用功能 .....	73
9.6	实战：百度自动答题 .....	80
9.7	本章小结 .....	85



第 10 章 手机 App 数据爬取.....	86
10.1 Appium 简介及原理 .....	86
10.2 搭建开发环境 .....	87
10.3 连接 Android 系统.....	92
10.4 App 的元素定位.....	97
10.5 App 的元素操控.....	99
10.6 实战：淘宝商品采集 .....	102
10.7 本章小结 .....	107
第 11 章 Splash、Mitmproxy 与 Aiohttp .....	109
11.1 Splash 动态数据抓取.....	109
11.1.1 简介及安装.....	109
11.1.2 使用 Splash 的 API 接口 .....	112
11.2 Mitmproxy 抓包 .....	116
11.2.1 简介及安装.....	116
11.2.2 用 Mitmdump 抓取爱奇艺视频 .....	116
11.3 Aiohttp 高并发抓取 .....	119
11.3.1 简介及使用 .....	119
11.3.2 Aiohttp 异步爬取小说排行榜 .....	123
11.4 本章小结 .....	126
第 12 章 验证码识别 .....	128
12.1 验证码的类型 .....	128
12.2 OCR 技术 .....	129
12.3 第三方平台 .....	131
12.4 本章小结 .....	134
第 13 章 数据清洗 .....	136
13.1 字符串操作 .....	136
13.1.1 截取.....	136
13.1.2 替换.....	137
13.1.3 查找.....	137
13.1.4 分割.....	138
13.2 正则表达式 .....	139
13.2.1 正则语法 .....	140
13.2.2 正则处理函数.....	141
13.3 BeautifulSoup 数据清洗 .....	144
13.3.1 BeautifulSoup 介绍与安装.....	144
13.3.2 BeautifulSoup 的使用示例 .....	146

13.4	本章小结 .....	149
<b>第 14 章</b>	<b>文档数据存储 .....</b>	<b>150</b>
14.1	CSV 数据的写入和读取 .....	150
14.2	Excel 数据的写入和读取 .....	151
14.3	Word 数据的写入和读取 .....	154
14.4	本章小结 .....	156
<b>第 15 章</b>	<b>ORM 框架 .....</b>	<b>158</b>
15.1	SQLAlchemy 介绍与安装 .....	158
15.1.1	操作数据库的方法 .....	158
15.1.2	SQLAlchemy 框架介绍 .....	158
15.1.3	SQLAlchemy 的安装 .....	159
15.2	连接数据库 .....	160
15.3	创建数据表 .....	162
15.4	添加数据 .....	164
15.5	更新数据 .....	165
15.6	查询数据 .....	166
15.7	本章小结 .....	168
<b>第 16 章</b>	<b>MongoDB 数据库操作 .....</b>	<b>169</b>
16.1	MongoDB 介绍 .....	169
16.2	MogoDB 的安装及使用 .....	170
16.2.1	MongoDB 的安装与配置 .....	170
16.2.2	MongoDB 可视化工具 .....	172
16.2.3	PyMongo 的安装 .....	173
16.3	连接 MongoDB 数据库 .....	173
16.4	添加文档 .....	174
16.5	更新文档 .....	175
16.6	查询文档 .....	176
16.7	本章小结 .....	178
<b>第 17 章</b>	<b>实战：爬取 51Job 招聘信息 .....</b>	<b>180</b>
17.1	项目分析 .....	180
17.2	获取城市编号 .....	180
17.3	获取招聘职位总页数 .....	182
17.4	爬取每个职位信息 .....	184
17.5	数据存储 .....	188
17.6	爬虫配置文件 .....	190
17.7	本章小结 .....	191

第 18 章 实战：分布式爬虫——QQ 音乐	193
18.1 项目分析	193
18.2 歌曲下载	194
18.3 歌手的歌曲信息	198
18.4 分类歌手列表	201
18.5 全站歌手列表	203
18.6 数据存储	204
18.7 分布式爬虫	205
18.7.1 分布式概念	205
18.7.2 并发库 concurrent.futures	206
18.7.3 分布式策略	207
18.8 本章小结	209
第 19 章 实战：12306 抢票爬虫	211
19.1 项目分析	211
19.2 验证码验证	211
19.3 用户登录与验证	214
19.4 查询车次	219
19.5 预订车票	225
19.6 提交订单	227
19.7 生成订单	233
19.8 本章小结	236
第 20 章 实战：玩转微博	244
20.1 项目分析	244
20.2 用户登录	244
20.3 用户登录（带验证码）	253
20.4 关键词搜索热门微博	259
20.5 发布微博	264
20.6 关注用户	268
20.7 点赞和转发评论	271
20.8 本章小结	277
第 21 章 实战：微博爬虫软件开发	278
21.1 GUI 库及 PyQt5 的安装与配置	278
21.1.1 GUI 库	278
21.1.2 PyQt5 安装及环境搭建	279
21.2 项目分析	281
21.3 软件主界面	284

21.4	相关服务界面 .....	288
21.5	微博采集界面 .....	292
21.6	微博发布界面 .....	297
21.7	微博爬虫功能 .....	308
21.8	本章小结 .....	315
第 22 章	Scrapy 爬虫开发 .....	317
22.1	认识与安装 Scrapy .....	317
22.1.1	常见爬虫框架介绍 .....	317
22.1.2	Scrapy 的运行机制 .....	318
22.1.3	安装 Scrapy .....	319
22.2	Scrapy 爬虫开发示例 .....	320
22.3	Spider 的编写 .....	326
22.4	Items 的编写 .....	329
22.5	Item Pipeline 的编写 .....	330
22.5.1	用 MongoDB 实现数据入库 .....	330
22.5.2	用 SQLAlchemy 实现数据入库 .....	332
22.6	Selectors 的编写 .....	333
22.7	文件下载 .....	336
22.8	本章小结 .....	339
第 23 章	Scrapy 扩展开发 .....	341
23.1	剖析 Scrapy 中间件 .....	341
23.1.1	SpiderMiddleware 中间件 .....	342
23.1.2	DownloaderMiddleware 中间件 .....	344
23.2	自定义中间件 .....	347
23.2.1	设置代理 IP 服务 .....	347
23.2.2	动态设置请求头 .....	350
23.2.3	设置随机 Cookies .....	353
23.3	实战: Scrapy+Selenium 爬取豆瓣电影评论 .....	355
23.3.1	网站分析 .....	355
23.3.2	项目设计与实现 .....	357
23.3.3	定义 Selenium 中间件 .....	359
23.3.4	开发 Spider 程序 .....	360
23.4	实战: Scrapy+Splash 爬取 B 站动漫信息 .....	362
23.4.1	Scrapy_Splash 实现原理 .....	363
23.4.2	网站分析 .....	363
23.4.3	项目设计与实现 .....	365
23.4.4	开发 Spider 程序 .....	367

23.5 实战: Scrapy+Redis 分布式爬取猫眼排行榜 .....	369
23.5.1 Scrapy_Redis 实现原理 .....	369
23.5.2 安装 Redis 数据库 .....	371
23.5.3 网站分析 .....	372
23.5.4 项目设计与实现 .....	373
23.5.5 开发 Spider 程序 .....	375
23.6 分布式爬虫与增量式爬虫 .....	377
23.6.1 基于管道实现增量式 .....	378
23.6.2 基于中间件实现增量式 .....	381
23.7 本章小结 .....	384
<b>第 24 章 实战: 爬取链家楼盘信息 .....</b>	<b>386</b>
24.1 项目分析 .....	386
24.2 创建项目 .....	389
24.3 项目配置 .....	389
24.4 定义存储字段 .....	391
24.5 定义管道类 .....	392
24.6 编写爬虫规则 .....	396
24.7 本章小结 .....	400
<b>第 25 章 实战: QQ 音乐全站爬取 .....</b>	<b>402</b>
25.1 项目分析 .....	402
25.2 项目创建与配置 .....	403
25.2.1 项目创建 .....	403
25.2.2 项目配置 .....	403
25.3 定义存储字段和管道类 .....	405
25.3.1 定义存储字段 .....	405
25.3.2 定义管道类 .....	405
25.4 编写爬虫规则 .....	408
25.5 本章小结 .....	413
<b>第 26 章 爬虫的上线部署 .....</b>	<b>415</b>
26.1 非框架式爬虫部署 .....	415
26.1.1 创建可执行程序 .....	415
26.1.2 制定任务计划程序 .....	417
26.1.3 创建服务程序 .....	421
26.2 框架式爬虫部署 .....	424
26.2.1 Scrapy 部署爬虫服务 .....	424
26.2.2 Gerapy 爬虫管理框架 .....	429
26.3 本章小结 .....	434

---

第 27 章 反爬虫的解决方案.....	435
27.1 常见的反爬虫技术 .....	435
27.2 基于验证码的反爬虫 .....	436
27.2.1 验证码出现的情况 .....	437
27.2.2 解决方案 .....	438
27.3 基于请求参数的反爬虫.....	439
27.3.1 请求参数的数据来源 .....	439
27.3.2 请求参数的查找 .....	440
27.4 基于请求头的反爬虫 .....	441
27.5 基于 Cookies 的反爬虫 .....	443
27.6 本章小结 .....	447
第 28 章 自己动手开发爬虫框架 .....	449
28.1 框架设计说明 .....	449
28.2 异步爬取方式 .....	450
28.3 数据清洗机制 .....	455
28.4 数据存储机制 .....	457
28.5 实战：用自制框架爬取豆瓣电影.....	463
28.6 本章小结 .....	468

# 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，通用的搜索引擎存在着一定的局限性：

- (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) 通用的搜索引擎扩大了网络覆盖率，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的不断发展，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，可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) 这类网络爬虫的爬取范围和数量巨大，对于爬取速度和存储空间要求较高，爬取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，是选择性地爬取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，聚焦爬虫只需要爬取与主题相关的页面，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取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取的页面尽可能是新的页面。只会在需要的时候爬取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，及时更新已爬取的网页，减小时间和空间上的耗费，但是增加了爬取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。这类爬虫也是本书讲述的重点之一。

实际上，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。

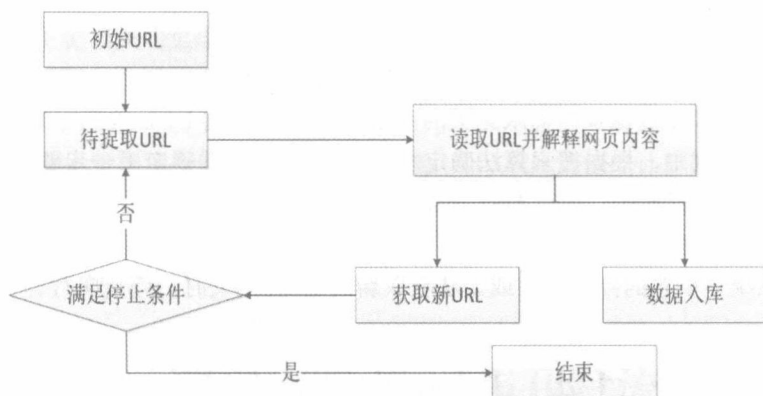


图 1-1 通用爬虫实现的原理及过程

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。

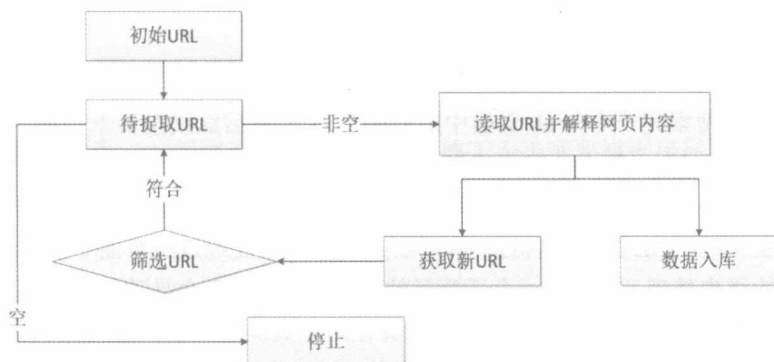


图 1-2 聚焦网络爬虫的原理

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。