



JetBrains大中华区市场部经理赵磊作序

300个实战案例 | 10万行源代码 | 22个综合实战项目

宁哥  
大讲堂

# Python Crawler Programming

Principles, Technologies and Development

# Python爬虫技术

## 深入理解原理、技术与开发

李宁◎编著

Li Ning

```
for i in range(1, tim + 1): if i == 1: tour.append(he) else: tour.  
#user/bin/python #-*- coding: UTF-8 -*- tour.append(2*he)  
tour.append(he) else:  
if i == 1: tour.append(he) else: tour.  
tour.append(he) else:  
if i == 1: tour.append(he) else: tour.  
tour.append(he) else:
```



清华大学出版社

海量资源

宁哥  
大讲堂

Python Crawler Programming

Principles, Technologies and Development

# Python爬虫技术

深入理解原理、技术与开发

李宁◎编著  
Li Ning

清华大学出版社  
北京

## 内 容 简 介

本书从实战角度系统讲解 Python 爬虫的核心知识点,并通过大量的真实项目让读者熟练掌握 Python 爬虫技术。本书用 20 多个实战案例,完美演绎了使用各种技术编写 Python 爬虫的方式,读者可以任意组合这些技术,完成非常复杂的爬虫应用。

全书共 20 章,分为 5 篇。第 1 篇基础知识(第 1、2 章),主要包括 Python 运行环境的搭建、HTTP 基础、网页基础(HTML、CSS、JavaScript 等)、爬虫的基本原理、Session 与 Cookie。第 2 篇网络库(第 3~6 章),主要包括网络库 urllib、urllib3、requests 和 Twisted 的核心使用方法,如发送 HTTP 请求、处理超时、设置 HTTP 请求头、搭建和使用代理、解析链接、Robots 协议等。第 3 篇解析库(第 7~10 章),主要包括 3 个常用解析库(lxml、Beautiful Soup 和 pyquery)的使用方法,同时介绍多种用于分析 HTML 代码的技术,如正则表达式、XPath、CSS 选择器、方法选择器等。第 4 篇数据存储(第 11、12 章),主要包括 Python 中数据存储的解决方案,如文件存储和数据库存储,其中数据库存储包括多种数据库,如本地数据库 SQLite、网络数据库 MySQL 以及文档数据库 MongoDB。第 5 篇爬虫高级应用(第 13~20 章),主要包括 Python 爬虫的一些高级技术,如抓取异步数据、Selenium、Splash、抓取移动 App 数据、Appium、多线程爬虫、爬虫框架 Scrapy,最后给出一个综合的实战案例,综合了 Python 爬虫、数据存储、PyQt5、多线程、数据可视化、Web 等多种技术实现一个可视化爬虫。

本书可以作为广大计算机软件技术开发者、互联网技术人员学习“爬虫技术”的参考用书。也可以作为高等院校计算机科学与技术、软件工程、人工智能等专业的教学参考用书。

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

Python 爬虫技术:深入理解原理、技术与开发/李宁编著. —北京:清华大学出版社,2020.1(2020.5重印)  
(宁哥大讲堂)

ISBN 978-7-302-53568-3

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

中国版本图书馆 CIP 数据核字(2019)第 173106 号

责任编辑:盛东亮 钟志芳

封面设计:李召霞

责任校对:白蕾

责任印制:丛怀宇

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:三河市铭诚印务有限公司

经 销:全国新华书店

开 本:203mm×260mm

印 张:31.25 字 数:855千字

版 次:2020年1月第1版

印 次:2020年5月第2次印刷

定 价:89.00元

产品编号:082638-01

## P R E F A C E 前 言

Python 现在非常火爆。但 Python 就和英语一样，如果只会 Python 语言，就相当于只能用英语进行日常会话。然而，真正的英语高手是可以作为专业领域翻译的，如 IT、金融、数学等专业领域。Python 也是一样，光学习 Python 语言是不行的，要想找到更好的工作，或得到更高的薪水，需要学会用 Python 做某一领域的应用。

现在 Python 应用的热门领域比较广，例如人工智能，不过人工智能不光涉及 Python 语言本身的技术，还涉及数学领域的知识，虽然比较火爆，但绝对不是短时间可以掌握的。然后有一个领域与人工智能的火爆程度相当，但不像人工智能那样难入门，这就是爬虫领域。

为什么爬虫领域如此火爆呢？其实爬虫的基本功能就是从网上下载各种类型的数据（如 HTML、图像文件等）。但不要小瞧这些下载的数据，因为这些数据将成为很多应用的数据源。例如，著名的 Google 搜索引擎，每天都有数以亿计的查询请求，而搜索引擎为这些请求返回的数据，都是来源于强大的爬虫。编写搜索引擎的第一步就是通过爬虫抓取整个互联网的数据，然后将这些数据库保存到本地（以特定的数据格式），接下来就是对这些数据进行分析整理。然后才可以通过搜索引擎进行查询。虽然搜索引擎的实现技术非常多，也非常复杂，但爬虫是 1，其他的所有技术都是 0，如果没有爬虫搜集数据，再强大的分析程序也毫无用武之地。

除了搜索引擎外，人工智能中的重要分支深度学习也需要爬虫抓取的数据来训练模型。例如，要想训练一个识别金字塔的深度学习模型，就需要大量与金字塔相关的图片进行训练。最简单的方式，就是使用百度或谷歌搜索金字塔图片，然后用爬虫抓取这些图片到本地。这是利用了搜索引擎通过关键字分类的特性，并且重新利用了这些分类的图片。

通过这些例子可以了解到，学习爬虫是进入其他更高端领域的钥匙，所以学习 Python 爬虫将成为第一个需要选择的热门领域。

尽管爬虫的基本功能是下载文件，但一个复杂的爬虫应用，可不光涉及网络技术。将数据下载后，还需要对数据进行分析，提取需要的信息，以及进行数据可视化，甚至需要一个基于 UI 的可视化爬虫。所以与爬虫有关的技术还是很多的。

由于 Python 爬虫涉及的技术很多，学习资料过于分散。所以，笔者觉得很有必要编写一本全面介绍 Python 爬虫实战类的书籍，在书中分享笔者对 Python 爬虫以及相关技术的理解和经验，帮助同行和感兴趣的朋友快速入门，并利用 Python 语言编写各种复杂的爬虫应用。笔者希望本书能起到抛砖引玉的作用，使读者对 Python 爬虫以及相关技术产生浓厚的兴趣，并能成功进入 Python 爬虫领域。加油！高薪的工作在等着你们！

本书使用最新的 Python 3 编写，并在书中探讨了关于 Python 爬虫的核心技术。全书分 5 篇，共 20 章。内容涵盖 Python 爬虫的基础知识、常用网络库、常用分析库、数据存储技术、异步数据处理、可见即可爬技术、抓取移动 App、Scrapy 等。本书还包含 20 多个真实的项目，以便让读者身临其境

地体验 Python 爬虫的魅力。

限于篇幅，本书无法囊括 Python 爬虫以及相关技术的方方面面，只能尽自己所能，与大家分享尽可能多的知识和经验。相信通过本书的学习，读者可以拥有进一步深入学习的能力，达到 Python 爬虫高手的程度也只是时间问题。

为方便读者学习，本书配套提供了程序代码，并录制了一集视频，扫码即可下载程序代码或观看视频。



最后，笔者希望本书能为国内的 Python 爬虫以及相关技术的普及，为广大从业者提供有价值的实践经验并帮助他们快速上手贡献绵薄之力。

编著者

2019年10月

# C O N T E N T S 目 录

前言

## 第 1 篇 基础知识

### 第 1 章 开发环境配置 ..... 2

- 1.1 安装官方的 Python 运行环境..... 2
- 1.2 配置 PATH 环境变量..... 5
- 1.3 安装 Anaconda Python 开发环境..... 6
- 1.4 安装 PyCharm..... 7
- 1.5 配置 PyCharm..... 8
- 1.6 小结..... 10

### 第 2 章 爬虫基础..... 11

- 2.1 HTTP 基础..... 11
  - 2.1.1 URI 和 URL..... 11
  - 2.1.2 超文本..... 12
  - 2.1.3 HTTP 与 HTTPS..... 12
  - 2.1.4 HTTP 的请求过程..... 15
  - 2.1.5 请求..... 17
  - 2.1.6 响应..... 20
- 2.2 网页基础..... 23
  - 2.2.1 HTML..... 23
  - 2.2.2 CSS..... 24
  - 2.2.3 CSS 选择器..... 25
  - 2.2.4 JavaScript..... 27
- 2.3 爬虫的基本原理..... 27
  - 2.3.1 爬虫的分类..... 27
  - 2.3.2 爬虫抓取数据的方式和手段..... 28
- 2.4 Session 与 Cookie..... 28
  - 2.4.1 静态页面和动态页面..... 29
  - 2.4.2 无状态 HTTP 与 Cookie..... 30
  - 2.4.3 利用 Session 和 Cookie 保持状态..... 30
  - 2.4.4 查看网站的 Cookie..... 31
  - 2.4.5 HTTP 状态何时会失效..... 32

- 2.5 实战案例：抓取所有的网络资源..... 33
- 2.6 实战案例：抓取博客文章列表..... 37
- 2.7 小结..... 40

## 第 2 篇 网络库

### 第 3 章 网络库 urllib..... 42

- 3.1 urllib 简介..... 42
- 3.2 发送请求与获得响应..... 43
  - 3.2.1 用 urlopen 函数发送 HTTP GET 请求..... 43
  - 3.2.2 用 urlopen 函数发送 HTTP POST 请求..... 44
  - 3.2.3 请求超时..... 45
  - 3.2.4 设置 HTTP 请求头..... 46
  - 3.2.5 设置中文 HTTP 请求头..... 48
  - 3.2.6 请求基础验证页面..... 50
  - 3.2.7 搭建代理与使用代理..... 54
  - 3.2.8 读取和设置 Cookie..... 56
- 3.3 异常处理..... 60
  - 3.3.1 URLError..... 60
  - 3.3.2 HTTPError..... 61
- 3.4 解析链接..... 62
  - 3.4.1 拆分与合并 URL (urlparse 与 urlunparse)..... 62
  - 3.4.2 另一种拆分与合并 URL 的方式 (urlsplit 与 urlunsplit)..... 63
  - 3.4.3 连接 URL (urljoin)..... 65
  - 3.4.4 URL 编码 (urlencode)..... 65
  - 3.4.5 编码与解码 (quote 与 unquote)..... 66
  - 3.4.6 参数转换 (parse\_qs 与 parse\_qsl)..... 66
- 3.5 Robots 协议..... 67
  - 3.5.1 Robots 协议简介..... 67
  - 3.5.2 分析 Robots 协议..... 68

3.6 小结 .....	69
<b>第4章 网络库 urllib3 .....</b>	<b>70</b>
4.1 urllib3 简介 .....	70
4.2 urllib3 模块 .....	70
4.3 发送 HTTP GET 请求 .....	71
4.4 发送 HTTP POST 请求 .....	72
4.5 HTTP 请求头 .....	74
4.6 HTTP 响应头 .....	76
4.7 上传文件 .....	76
4.8 超时 .....	78
4.9 小结 .....	79
<b>第5章 网络库 requests .....</b>	<b>80</b>
5.1 基本用法 .....	80
5.1.1 requests 的 HelloWorld .....	81
5.1.2 GET 请求 .....	81
5.1.3 添加 HTTP 请求头 .....	82
5.1.4 抓取二进制数据 .....	83
5.1.5 POST 请求 .....	84
5.1.6 响应数据 .....	85
5.2 高级用法 .....	87
5.2.1 上传文件 .....	88
5.2.2 处理 Cookie .....	89
5.2.3 使用同一个会话 (Session) .....	90
5.2.4 SSL 证书验证 .....	91
5.2.5 使用代理 .....	94
5.2.6 超时 .....	95
5.2.7 身份验证 .....	97
5.2.8 将请求打包 .....	97
5.3 小结 .....	98
<b>第6章 Twisted 网络框架 .....</b>	<b>99</b>
6.1 异步编程模型 .....	99
6.2 Reactor (反应堆) 模式 .....	101
6.3 HelloWorld, Twisted 框架 .....	101
6.4 用 Twisted 实现时间戳客户端 .....	103
6.5 用 Twisted 实现时间戳服务端 .....	104
6.6 小结 .....	105

## 第3篇 解析库

<b>第7章 正则表达式 .....</b>	<b>108</b>
7.1 使用正则表达式 .....	108
7.1.1 使用 match 方法匹配字符串 .....	108
7.1.2 使用 search 方法在一个字符串中 查找模式 .....	109
7.1.3 匹配多个字符串 .....	110
7.1.4 匹配任何单个字符 .....	111
7.1.5 使用字符集 .....	112
7.1.6 重复、可选和特殊字符 .....	114
7.1.7 分组 .....	117
7.1.8 匹配字符串的起始和结尾以及 单词边界 .....	118
7.1.9 使用 findall 和 finditer 查找 每一次出现的位置 .....	120
7.1.10 用 sub 和 subn 搜索与替换 .....	121
7.1.11 使用 split 分隔字符串 .....	122
7.2 一些常用的正则表达式 .....	123
7.3 项目实战: 抓取小说目录和全文 .....	124
7.4 项目实战: 抓取猫眼电影 Top100 榜单 .....	128
7.5 项目实战: 抓取糗事百科网的段子 .....	133
7.6 小结 .....	136
<b>第8章 lxml 与 XPath .....</b>	<b>137</b>
8.1 lxml 基础 .....	137
8.1.1 安装 lxml .....	137
8.1.2 操作 XML .....	138
8.1.3 操作 HTML .....	140
8.2 XPath .....	141
8.2.1 XPath 概述 .....	141
8.2.2 使用 XPath .....	141
8.2.3 选取所有节点 .....	143
8.2.4 选取子节点 .....	145
8.2.5 选取父节点 .....	146
8.2.6 属性匹配与获取 .....	146
8.2.7 多属性匹配 .....	147
8.2.8 按序选择节点 .....	148
8.2.9 节点轴选择 .....	149

8.2.10 在 Chrome 中自动获得 XPath 代码..... 151

8.2.11 使用 Chrome 验证 XPath ..... 153

8.3 项目实战：抓取豆瓣 Top250 图书榜单 ..... 154

8.4 项目实战：抓取起点中文网的小说信息 ..... 158

8.5 小结 ..... 161

**第 9 章 Beautiful Soup 库..... 162**

9.1 Beautiful Soup 简介 ..... 162

9.2 Beautiful Soup 基础 ..... 162

9.2.1 安装 Beautiful Soup ..... 163

9.2.2 选择解析器 ..... 164

9.2.3 编写第一个 Beautiful Soup 程序 ..... 164

9.3 节点选择器 ..... 165

9.3.1 选择节点 ..... 165

9.3.2 嵌套选择节点 ..... 167

9.3.3 选择子节点 ..... 168

9.3.4 选择父节点 ..... 171

9.3.5 选择兄弟节点 ..... 172

9.4 方法选择器 ..... 174

9.4.1 find\_all 方法 ..... 174

9.4.2 find 方法 ..... 177

9.5 CSS 选择器 ..... 178

9.5.1 基本用法 ..... 179

9.5.2 嵌套选择节点 ..... 180

9.5.3 获取属性值与文本 ..... 181

9.5.4 通过浏览器获取 CSS 选择器代码 ..... 182

9.6 实战案例：抓取租房信息 ..... 184

9.7 实战案例：抓取酷狗网络红歌榜 ..... 188

9.8 小结 ..... 191

**第 10 章 pyquery 库..... 192**

10.1 pyquery 简介 ..... 192

10.2 pyquery 基础 ..... 192

10.2.1 安装 pyquery ..... 193

10.2.2 pyquery 的基本用法 ..... 193

10.3 CSS 选择器 ..... 194

10.4 查找节点 ..... 196

10.4.1 查找子节点 ..... 196

10.4.2 查找父节点 ..... 197

10.4.3 查找兄弟节点 ..... 198

10.4.4 获取节点信息 ..... 199

10.5 修改节点 ..... 203

10.5.1 添加和移除节点的样式  
(addClass 和 removeClass)..... 204

10.5.2 修改节点属性和文本内容  
(attr、removeAttr、text 和 html)..... 205

10.5.3 删除节点 (remove) ..... 207

10.6 伪类选择器 ..... 208

10.7 项目实战：抓取当当图书排行榜 ..... 210

10.8 项目实战：抓取京东商城手机销售排行榜... 213

10.9 小结 ..... 219

## 第 4 篇 数据存储

**第 11 章 文件存储 ..... 222**

11.1 打开文件 ..... 222

11.2 操作文件的基本方法 ..... 224

11.2.1 读文件和写文件 ..... 224

11.2.2 读行和写行 ..... 226

11.3 使用 FileInput 对象读取文件 ..... 227

11.4 处理 XML 格式的数据 ..... 228

11.4.1 读取与搜索 XML 文件 ..... 228

11.4.2 字典转换为 XML 字符串 ..... 229

11.4.3 XML 字符串转换为字典 ..... 231

11.5 处理 JSON 格式的数据 ..... 232

11.5.1 JSON 字符串与字典互相转换 ..... 233

11.5.2 将 JSON 字符串转换为类实例 ..... 234

11.5.3 将类实例转换为 JSON 字符串 ..... 236

11.5.4 类实例列表与 JSON 字符串互相  
转换 ..... 236

11.6 将 JSON 字符串转换为 XML 字符串 ..... 237

11.7 CSV 文件存储 ..... 238

11.7.1 写入 CSV 文件 ..... 238

11.7.2 读取 CSV 文件 ..... 241

11.8 小结 ..... 241

**第 12 章 数据库存储 ..... 242**

12.1 SQLite 数据库 ..... 242

12.1.1 管理 SQLite 数据库 ..... 243

12.1.2 用 Python 操作 SQLite 数据库 ..... 245

12.2	MySQL 数据库 .....	247
12.2.1	安装 MySQL .....	247
12.2.2	在 Python 中使用 MySQL .....	250
12.3	非关系型数据库 .....	253
12.3.1	NoSQL 简介 .....	253
12.3.2	MongoDB 数据库 .....	253
12.3.3	pymongo 模块 .....	255
12.4	项目实战: 抓取豆瓣音乐排行榜 .....	256
12.5	项目实战: 抓取豆瓣电影排行榜 .....	260
12.6	小结 .....	264

## 第 5 篇 爬虫高级应用

### 第 13 章 抓取异步数据 .....

266

13.1	异步加载与 AJAX .....	266
13.2	基本原理 .....	267
13.3	逆向工程 .....	270
13.4	提取结果 .....	274
13.5	项目实战: 支持搜索功能的图片爬虫 .....	274
13.6	项目实战: 抓取京东图书评价 .....	279
13.7	小结 .....	284

### 第 14 章 可见即可爬: Selenium .....

285

14.1	安装 Selenium .....	286
14.2	安装 WebDriver .....	286
14.2.1	安装 ChromeDriver .....	287
14.2.2	装 Edge WebDriver .....	288
14.2.3	安装其他浏览器的 WebDriver .....	289
14.3	Selenium 的基本使用方法 .....	289
14.4	查找节点 .....	293
14.4.1	查找单个节点 .....	293
14.4.2	查找多个节点 .....	295
14.5	节点交互 .....	297
14.6	动作链 .....	298
14.7	执行 JavaScript 代码 .....	301
14.8	获取节点信息 .....	302
14.9	管理 Cookies .....	303
14.10	改变节点的属性值 .....	304
14.11	项目实战: 抓取 QQ 空间说说的内容 .....	306
14.12	小结 .....	308

### 第 15 章 基于 Splash 的爬虫 .....

309

15.1	Splash 基础 .....	309
15.1.1	Splash 功能简介 .....	309
15.1.2	安装 Docker .....	310
15.1.3	安装 Splash .....	310
15.2	Splash Lua 脚本 .....	312
15.2.1	第一个 Lua 脚本 .....	312
15.2.2	异步处理 .....	313
15.2.3	Splash 对象属性 .....	314
15.2.4	go 方法 .....	318
15.2.5	wait 方法 .....	319
15.2.6	jsfunc 方法 .....	320
15.2.7	evaljs 方法 .....	320
15.2.8	runjs 方法 .....	320
15.2.9	autoload 方法 .....	321
15.2.10	call_later 方法 .....	322
15.2.11	http_get 方法 .....	323
15.2.12	http_post 方法 .....	324
15.2.13	set_content 方法 .....	325
15.2.14	html 方法 .....	325
15.2.15	png 方法 .....	326
15.2.16	jpeg 方法 .....	326
15.2.17	har 方法 .....	326
15.2.18	其他方法 .....	327
15.3	使用 CSS 选择器 .....	331
15.3.1	select 方法 .....	331
15.3.2	select_all 方法 .....	332
15.4	模拟鼠标和键盘的动作 .....	333
15.5	Splash HTTP API .....	334
15.6	项目实战: 使用 Splash Lua 抓取京东 搜索结果 .....	338
15.7	小结 .....	340

### 第 16 章 抓取移动 App 的数据 .....

341

16.1	使用 Charles .....	341
16.1.1	抓取 HTTP 数据包 .....	342
16.1.2	安装 PC 端证书 .....	344
16.1.3	在手机端安装证书 .....	345
16.1.4	监听 HTTPS 数据包 .....	346
16.2	使用 mitmproxy .....	348

16.2.1	安装 mitmproxy.....	348	18.1.1	进程.....	389
16.2.2	在 PC 端安装 mitmproxy 证书.....	349	18.1.2	线程.....	390
16.2.3	在移动端安装 mitmproxy 证书.....	352	18.2	Python 与线程.....	390
16.2.4	mitmproxy 有哪些功能.....	353	18.2.1	使用单线程执行程序.....	390
16.2.5	设置手机的代理.....	353	18.2.2	使用多线程执行程序.....	391
16.2.6	用 mitmproxy 监听 App 的请求与 响应数据.....	354	18.2.3	为线程函数传递参数.....	393
16.2.7	使用 mitmproxy 编辑请求信息.....	356	18.2.4	线程和锁.....	394
16.2.8	mitmdump 与 Python 对接.....	357	18.3	高级线程模块 (threading).....	395
16.2.9	使用 mitmweb 监听请求与响应.....	361	18.3.1	Thread 类与线程函数.....	395
16.3	项目实战: 实时抓取“得到” App 在线课程.....	363	18.3.2	Thread 类与线程对象.....	396
16.4	小结.....	367	18.3.3	从 Thread 类继承.....	398
<b>第 17 章 使用 Appium 在移动端抓取数据... 368</b>			18.4	线程同步.....	399
17.1	安装 Appium.....	368	18.4.1	线程锁.....	400
17.1.1	安装 Appium 桌面端.....	368	18.4.2	信号量.....	402
17.1.2	配置 Android 开发环境.....	370	18.5	生产者—消费者问题与 queue 模块.....	405
17.1.3	配置 iOS 开发环境.....	371	18.6	多进程.....	407
17.2	Appium 的基本使用方法.....	372	18.7	项目实战: 抓取豆瓣音乐 Top250 排行榜 (多线程版).....	408
17.2.1	启动 Appium 服务.....	372	18.8	项目实战: 抓取豆瓣音乐 Top250 排行榜 (多进程版).....	411
17.2.2	查找 Android App 的 Package 和 入口 Activity.....	374	18.9	小结.....	412
17.2.3	控制 App.....	376	<b>第 19 章 网络爬虫框架: Scrapy..... 413</b>		
17.3	使用 Python 控制手机 App.....	379	19.1	Scrapy 基础知识.....	413
17.4	AppiumPythonClient API.....	380	19.1.1	Scrapy 简介.....	413
17.4.1	初始化 (Remote 类).....	380	19.1.2	Scrapy 安装.....	414
17.4.2	查找元素.....	381	19.1.3	Scrapy Shell 抓取 Web 资源.....	415
17.4.3	单击元素.....	381	19.2	用 Scrapy 编写网络爬虫.....	417
17.4.4	屏幕拖动.....	382	19.2.1	创建和使用 Scrapy 工程.....	417
17.4.5	屏幕滑动.....	382	19.2.2	在 PyCharm 中使用 Scrapy.....	419
17.4.6	拖曳操作.....	383	19.2.3	在 PyCharm 中使用扩展工具运行 Scrapy 程序.....	421
17.4.7	文本输入.....	383	19.2.4	使用 Scrapy 抓取数据, 并通过 XPath 指定解析规则.....	423
17.4.8	动作链.....	383	19.2.5	将抓取到的数据保存为多种格式的 文件.....	424
17.5	项目实战: 利用 Appium 抓取微信朋友圈 信息.....	384	19.2.6	使用 ItemLoader 保存单条抓取的 数据.....	426
17.6	小结.....	388	19.2.7	使用 ItemLoader 保存多条抓取的 数据.....	428
<b>第 18 章 多线程和多进程爬虫..... 389</b>					
18.1	线程与进程.....	389			

19.2.8	抓取多个 URL	430
19.3	Scrapy 的高级应用	431
19.3.1	处理登录页面	431
19.3.2	处理带隐藏文本框的登录页面	434
19.3.3	通过 API 抓取天气预报数据	436
19.3.4	从 CSV 格式转换到 JSON 格式	443
19.3.5	下载器中间件	447
19.3.6	爬虫中间件	452
19.3.7	Item 管道	455
19.3.8	通用爬虫	465
19.4	小结	474

第 20 章	综合爬虫项目：可视化爬虫	475
20.1	项目简介	475
20.2	主界面设计和实现	477
20.3	获取商品页数和每页商品数	478
20.4	并发抓取商品列表	479
20.5	数据库操作类	481
20.6	情感分析	484
20.7	抓取和分析商品评论数据	485
20.8	可视化评论数据	486
20.9	小结	488

# 第1篇

# 基础知识

- 
- ◎ 第1章 开发环境配置
  - ◎ 第2章 爬虫基础
-

## 第 1 章

## 开发环境配置

“工欲善其事，必先利其器。”，由于本书涉及的爬虫使用 Python 语言编写，所以在学习编写爬虫之前，必须先搭建好 Python 开发环境。Python 程序可以直接使用记事本开发，也可以使用 IDE (Integrated Development Environment, 集成开发环境) 开发。不过大多数项目都会使用 IDE 开发。因为 IDE 支持代码高亮、智能提示和可视化等功能，这些功能可以让开发效率大大提升。

本章主要介绍以下内容：

- (1) 安装 Python 标准环境；
- (2) 安装 Anaconda Python 环境；
- (3) 设置 PATH 环境变量；
- (4) 安装 PyCharm；
- (5) 配置 PyCharm。



## 1.1 安装官方的 Python 运行环境

不管用什么工具开发 Python 程序，都必须先安装 Python 的运行环境。由于 Python 是跨平台的，所以在安装之前，先要确定在哪个操作系统平台上安装，目前最常用的是 Windows、macOS 和 Linux 三大平台。由于目前国内使用 Windows 和 macOS 的程序员比较多，所以本章主要以 Windows 和 macOS 为例介绍如何安装和使用 Python 运行环境。

读者可以直接到 Python 的官网 (<https://www.python.org/downloads>) 下载相应操作系统平台的 Python 安装包。

进入下载页面，浏览器会根据不同的操作系统显示不同的 Python 安装包下载链接。如果读者使用的是 Windows 平台，会显示如图 1-1 所示的 Python 下载页面。

如果读者使用的是 macOS 平台，会显示如图 1-2 所示的 Python 下载页面。

不管是哪个操作系统平台的下载页面，都会出现“Download Python 3.7.2”按钮（随着时间的推移，可能版本号略有不同）。直接单击“Download Python3.7.2”按钮下载相应平台的 Python 安装

包即可。如果是 Windows 平台，下载的是 exe 文件，如果是 macOS 平台，下载的是 pkg 文件，这是 macOS 上的安装程序，直接安装即可。

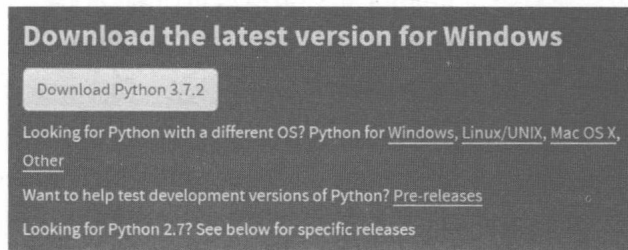


图 1-1 Windows 平台的 Python 下载页面

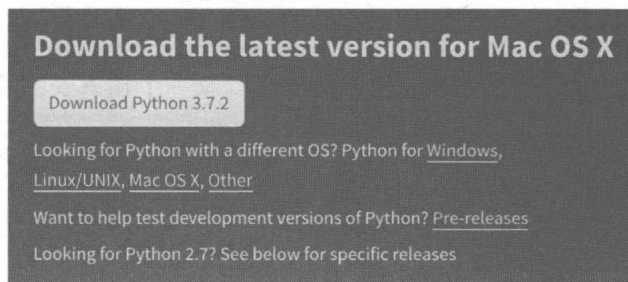


图 1-2 macOS 平台的 Python 下载页面

现在主要介绍在 Windows 平台如何安装 Python 运行环境。首先运行下载的 exe 文件，会显示 Python 安装界面。如果读者的机器已经安装了 Python 环境，那么会显示如图 1-3 所示的升级或重新安装的界面。如果读者的机器没有安装过 Python 环境，显示的安装界面与图 1-3 的界面类似，只是 Upgrade Now 变成了 Install Now。建议读者选中界面下方的“Add Python 3.7 to PATH”复选框，这样安装程序就会自动将 Python 的路径加到 PATH 环境变量中。

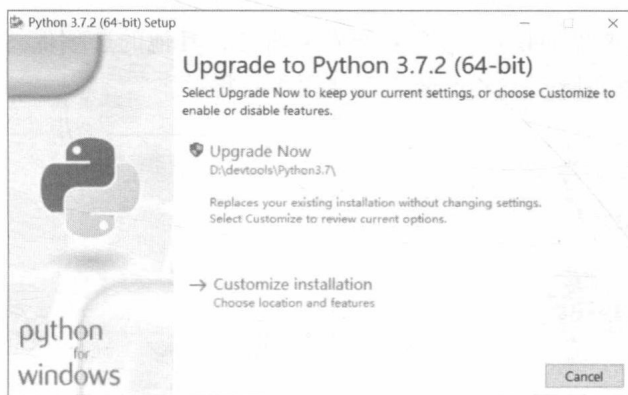


图 1-3 Windows 版 Python 3.7 安装程序的第一个界面

在图 1-3 所示的界面中出现两个安装选项，Upgrade Now 和 Customize installation，如果读者要升级 Python，可以单击 Upgrade Now 按钮，如果读者想定制安装选项，可以单击 Customize installation。现在我们定制安装 Python 环境。单击 Customize installation 按钮，会显示如图 1-4 所示的界面，默认所有的复选框全部选中，保留默认设置，然后单击 Next 按钮进入下一个安装界面。

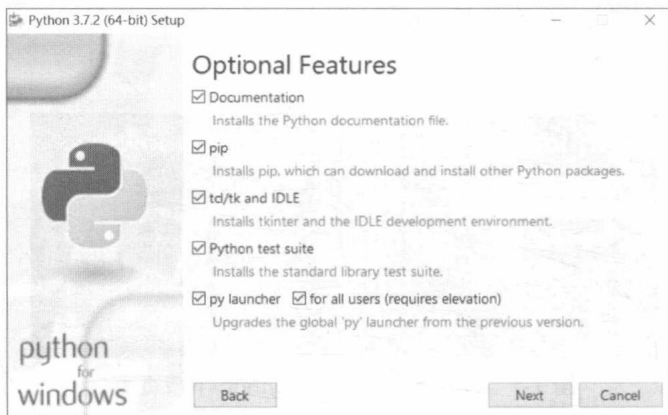


图 1-4 Python 选项界面

进入下一个安装界面后，会看到如图 1-5 所示的选项以及一个输入安装路径的文本框。

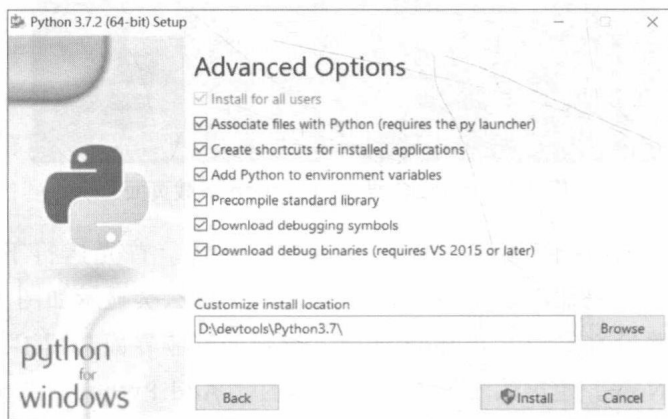


图 1-5 Python 安装程序的高级选项

读者可以在这个高级选项界面指定 Python 的安装路径，其他的选项保持默认设置即可。最后单击 Install 按钮安装 Python 开发环境，安装进度界面如图 1-6 所示。

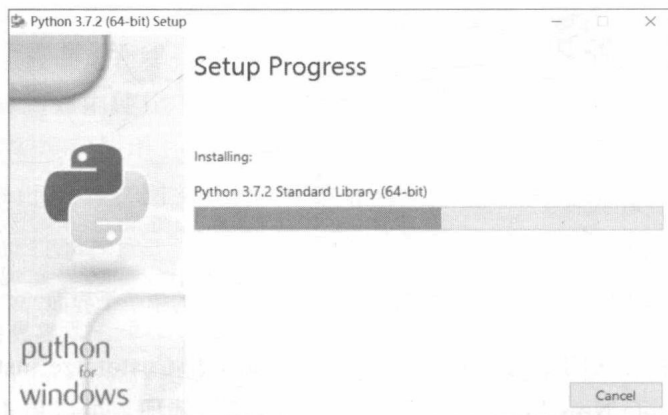


图 1-6 Python 安装进度界面

安装完成后，关闭安装界面。macOS 下安装 Python 的过程与 Windows 类似，读者可以自行安装。

## 1.2 配置 PATH 环境变量

在安装完 Python 运行环境后，我们可以测试一下 Python 运行环境，如果在安装 Python 的过程中忘记了选中“Add Python 3.7 to PATH”复选框，那么默认情况下，Python 安装程序是不会将 Python 安装目录添加到 PATH 环境变量中的。这样一来，我们就无法在 Windows 命令行工具中的任何目录执行 python 命令了，必须进入 Python 的安装目录才可以使用 python 命令。

为了更方便地执行 python 命令，建议将 Python 安装目录添加到 PATH 环境变量中。在 Windows 平台配置 PATH 环境变量的步骤如下。

回到桌面，右击“此电脑”，在弹出菜单中单击“属性”菜单项，会显示如图 1-7 所示的“系统”窗口。



图 1-7 “系统”窗口

单击“系统”窗口左侧的“高级系统设置”，弹出如图 1-8 所示的“系统属性”窗口。

单击“系统属性”窗口下方的“环境变量(N)...”按钮，弹出如图 1-9 所示的“环境变量”窗口。



图 1-8 “系统属性”窗口

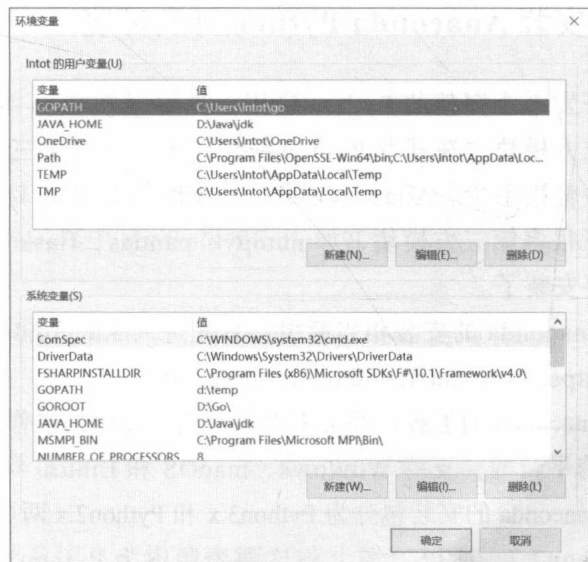


图 1-9 “环境变量”窗口

“环境变量”窗口有两个列表，上面的列表是为 Windows 当前登录用户设置环境变量，在这里设

置的环境变量只对当前登录用户有效。下面的列表是对所有用户设置的环境变量，这些变量对所有的用户都有效。读者在哪里设置 PATH 环境变量都可以，本书在上面的列表中设置了 PATH 环境变量。如果在列表中没有 PATH 环境变量，可单击“新建(N)...”按钮添加一个新的 PATH 环境变量。如果已经有了 PATH 环境变量，双击 PATH，就会弹出如图 1-10 所示的“编辑环境变量”对话框。单击“新建”按钮，添加 Python 的安装路径即可。注意，这里要填写目录，而不是 python.exe 文件的路径。

在 Path 环境变量中添加 Python 的路径后，打开 Windows 命令行工具，执行 `python --version` 命令，如果输出如图 1-11 所示的内容，说明 Python 已经安装成功。

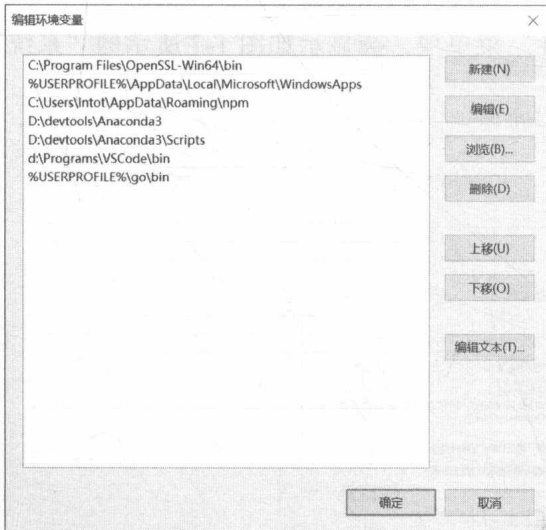


图 1-10 “编辑环境变量”对话框



图 1-11 测试 Python 运行环境是否安装成功

### 1.3 安装 Anaconda Python 开发环境

开发一个完整的 Python 应用，光使用 Python 本身提供的模块是远远不够的，因此，需要使用大量第三方模块。在开发 Python 应用时安装这些第三方模块是一件令人头痛的事，不过 Anaconda 会让这件事轻松不少。Anaconda 是一个集成的 Python 运行环境。除了包含 Python 本身的运行环境外，还集成了很多第三方模块，如 numpy、pandas、flask 等。也就是说，安装了 Anaconda 后，这些模块都不需要安装了。

Anaconda 的安装相当简单，首先进入 Anaconda 的下载页面，地址为：

<https://www.anaconda.com/distribution>

Anaconda 的下载页面会根据用户当前使用的操作系统自动切换到相应的 Anaconda 安装包。Anaconda 是跨平台的，支持 Windows、macOS 和 Linux。

Anaconda 的安装包分为 Python3.x 和 Python2.x 两个版本，目前 Anaconda 稳定版本分别支持 Python3.6 和 Python2.7，所以习惯上称这两个版本为 Python3.6 版和 Python2.7 版，尽管目前 Python 的最新版本是 3.7，但使用 Python3.6 仍然可以完美地运行本书的案例，所以读者只需要选择 Python3.6 或其以上版本即可。