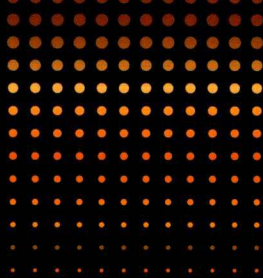


大数据及人工智能产教融合系列丛书



# Python

## 网络爬虫从入门到实践

庄培杰◎编 著

### 前沿

全书基于  
Python 3.7.0  
讲解爬虫

### 全面

涵盖抓包、请  
求模拟、数据  
解析、反爬虫  
应对等众多知  
识点

### 实战

20多个经典实  
战案例详解，  
随学随用

### 超值

附赠22个  
Python爬虫项  
目源码，帮助  
新手轻松入门

 中国工信出版集团

 电子工业出版社  
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY  
<http://www.phei.com.cn>

大数据及人工智能产教融合系列丛书

# Python 网络爬虫从入门到实践

庄培杰 编著

Publishing House of Electronics Industry

北京 · BEIJING

## 内 容 简 介

本书讲解了如何使用 Python 编写网络爬虫，涵盖爬虫的概念、Web 基础、Chrome、Charles 和 Packet Capture 抓包、urllib、Requests 请求库、lxml、Beautiful Soup、正则表达式解析数据、CSV、Excel、MySQL、Redis、MongoDB 保存数据、反爬虫策略应对、爬虫框架 Scrapy 的使用与部署，以及应用案例。

本书结构清晰、内容精练，代码示例典型实用，附带实践过程中遇到问题的解决方案，非常适合 Python 初学者和进阶读者阅读。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

### 图书在版编目 (CIP) 数据

Python 网络爬虫从入门到实践 / 庄培杰编著. —北京: 电子工业出版社, 2019.8

(大数据及人工智能产教融合系列丛书)

ISBN 978-7-121-37105-9

I. ①P… II. ①庄… 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 144186 号

责任编辑: 李 冰 特约编辑: 王 纲

印 刷: 三河市鑫金马印装有限公司

装 订: 三河市鑫金马印装有限公司

出版发行: 电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本: 787×1092 1/16 印张: 19.5 字数: 499 千字

版 次: 2019 年 8 月第 1 版

印 次: 2019 年 8 月第 1 次印刷

定 价: 79.00 元

凡所购买电子工业出版社图书有缺损问题, 请向购买书店调换。若书店售缺, 请与本社发行部联系, 联系及邮购电话: (010) 88254888, 88258888。

质量投诉请发邮件至 [zllts@phei.com.cn](mailto:zllts@phei.com.cn), 盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

本书咨询联系方式: [libing@phei.com.cn](mailto:libing@phei.com.cn)。

# 前 言

---

笔者是一名 Android 开发工程师，在接触 Python 之前，每天的工作流程基本上都是接到新版本的需求→写新页面→修改接口和业务逻辑，非常乏味。

持续性的重复劳动，让笔者意识到一个问题：如果只会 Android 开发，能做的事情非常有限！例如，自己写一个 App，如果没有可供调用的 API，那么只能得到一个单机的 App。因为自己对后台相关的技术一窍不通，平时根本不用去了解这方面的知识，只要给后台发出请求，然后解析数据，显示到页面上就好。

笔者开始琢磨花点时间去学习后台开发的知识，候选方案有 Java、Kotlin、Python、PHP 和 Go。因为之前接触过 Python，加上笔者只是想写一个给自己的 App 调用的 API，所以最后选择了 Python 这门对初学者非常友好的编程语言。

笔者花了一周多的时间把 Python 的基本语法研究了一遍，发现 Python 语法简单、代码简洁。正当笔者准备去看 Flask 这个轻量级 Python Web 框架的文档时，问题来了：没有数据源。数据源都没有，写什么 API？于是，笔者把学习重心转移到 Python 爬虫编写上。

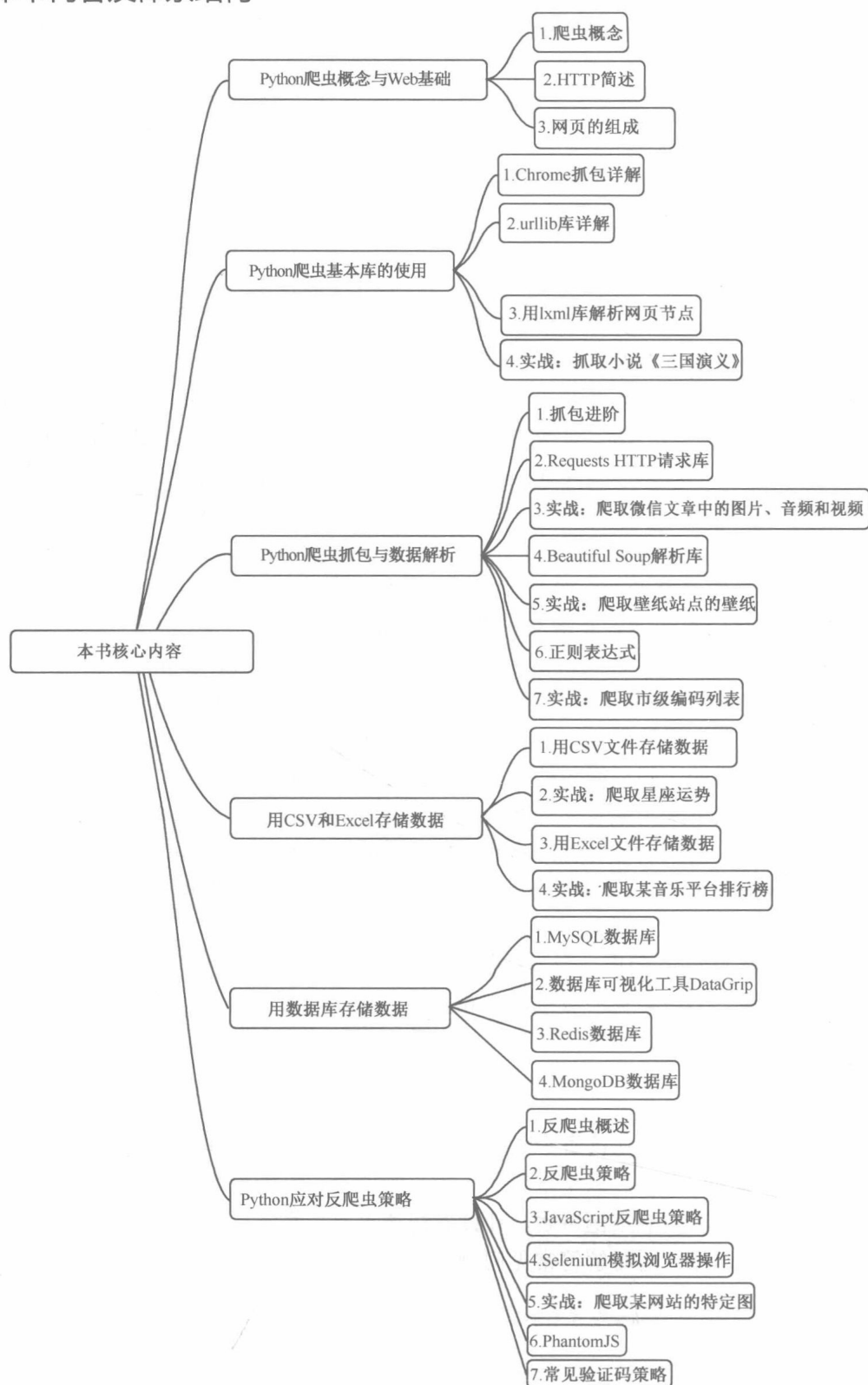
这仿佛打开了新的大门，简单的几行代码就完成了站点模拟请求、页面解析，还把图片下载到了本地。接下来，笔者又开始研究如何提高爬取效率、解析效率，以及采用多样化的数据存储形式。后来发现有些站点有反爬虫的策略，得不到数据，于是笔者又对此进行研究。爬虫学得差不多了，笔者又开始研究数据分析，分析爬取的结果，得出一些有用的结论，如分析招聘网站上某个工作岗位的行情等。

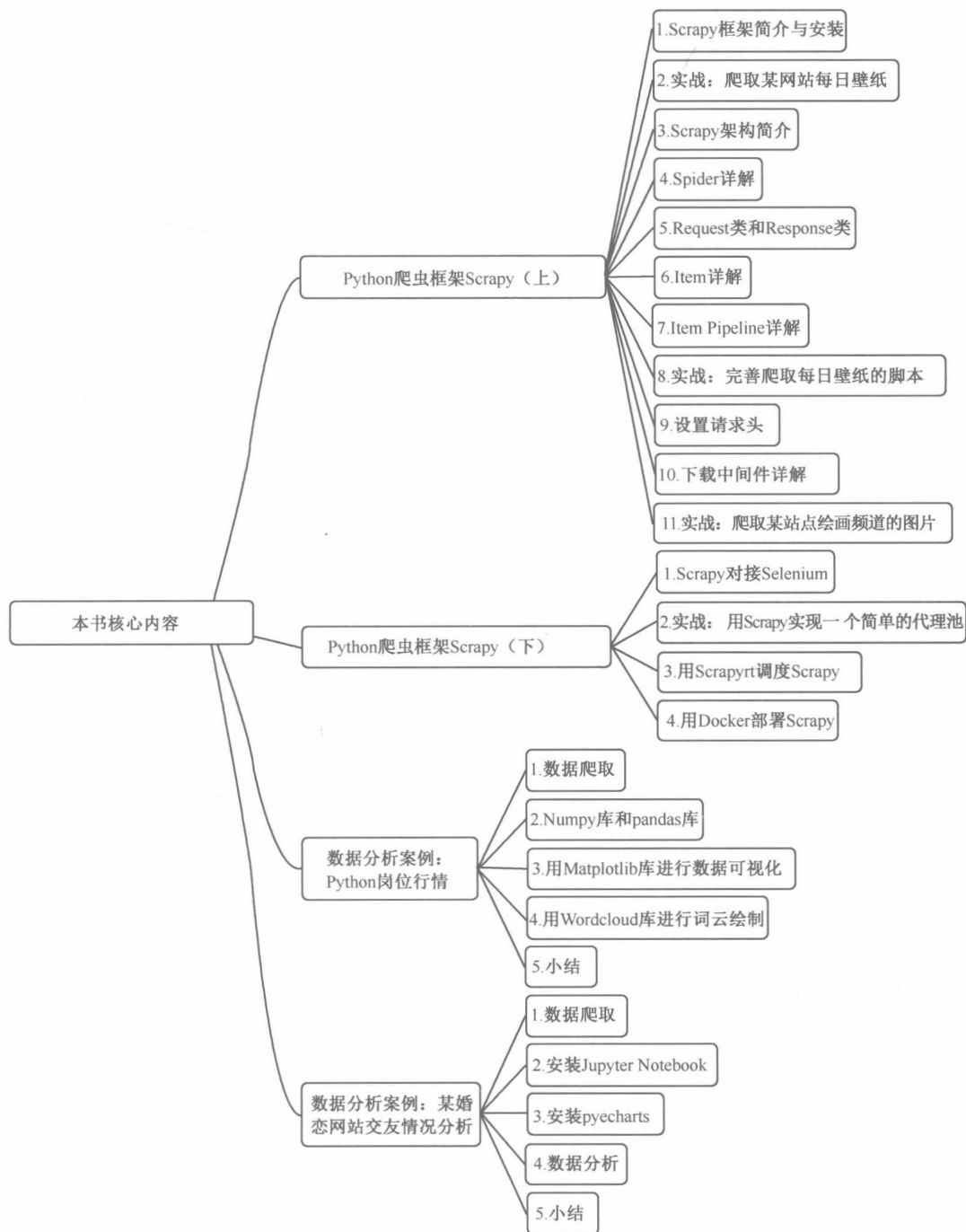
笔者相信掌握 Python 爬虫，会为你的工作、生活带来一些意想不到的便利。因笔者能力和成书时间所限，书中难免会有纰漏，相关代码也不是尽善尽美，希望读者批评指正，笔者将感激不尽。

## 本书特色

- (1) 清晰的学习路线，让读者少走弯路。
- (2) 精练的知识点讲解，只学重要的、核心的内容，不拖泥带水。
- (3) 简单而经典的实战案例讲解，让读者快速上手，举一反三。
- (4) 代码注释详尽，适时解答常见问题。

## 本书内容及体系结构





## 本书读者对象

- 想要学习 Python 编程的人员。
- 有 Python 编程基础、想进阶的人员。
- 各计算机、软件专业的在校学生。
- 其他对 Python 爬虫感兴趣的读者。

# 目 录

---

第 1 章 Python 爬虫概念与 Web 基础	1
1.1 爬虫概念	1
1.1.1 什么是爬虫	1
1.1.2 爬虫使用场景的引入	2
1.1.3 爬虫的组成部分	3
1.1.4 模拟请求	3
1.1.5 数据解析	4
1.1.6 数据保存	5
1.1.7 爬虫的学习路线	5
1.2 HTTP 简述	6
1.2.1 简述一次网络请求过程	6
1.2.2 URI 和 URL	7
1.2.3 HTTP 请求报文	8
1.2.4 HTTP 响应报文	10
1.3 网页的组成	13
1.3.1 HTML 简介	13
1.3.2 CSS 选择器简介	16
1.3.3 JavaScript 简介	17
第 2 章 Python 爬虫基本库的使用	18
2.1 Chrome 抓包详解	18
2.1.1 Controls	20
2.1.2 Filter	21
2.1.3 Request Table	21

2.2	urllib 库详解	23
2.2.1	发送请求	23
2.2.2	抓取二进制文件	24
2.2.3	模拟 GET 和 POST 请求	25
2.2.4	修改请求头	26
2.2.5	设置连接超时	27
2.2.6	延迟提交数据	27
2.2.7	设置代理	27
2.2.8	Cookie	28
2.2.9	urllib.parse 模块	29
2.2.10	urllib.error 异常处理模块	31
2.2.11	urllib.robotparser 模块	32
2.3	用 lxml 库解析网页节点	34
2.3.1	安装库	34
2.3.2	XPath 语法速成	34
2.4	实战: 爬取小说《三国演义》	36
第 3 章	Python 爬虫抓包与数据解析	41
3.1	抓包进阶	41
3.1.1	HTTPS 介绍	42
3.1.2	HTTPS 的工作流程	43
3.1.3	Charles 抓包	43
3.1.4	Packet Capture 抓包	49
3.2	Requests HTTP 请求库	52
3.2.1	Requests 库简介	53
3.2.2	Requests HTTP 基本请求	53
3.2.3	Requests 请求常用设置	54
3.2.4	Requests 处理返回结果	54
3.2.5	Requests 处理 Cookie	55
3.2.6	Requests 重定向与请求历史	55
3.2.7	Requests 错误与异常处理	55
3.2.8	Requests Session 会话对象	55
3.2.9	Requests SSL 证书验证	56
3.3	实战: 爬取微信文章中的图片、音频和视频	56
3.3.1	爬取标题	56
3.3.2	爬取图片	57

3.3.3	爬取音频	58
3.3.4	爬取视频	60
3.3.5	代码整理	64
3.4	Beautiful Soup 解析库	67
3.4.1	Beautiful Soup 简介	67
3.4.2	Beautiful Soup 对象实例化	67
3.4.3	Beautiful Soup 的四大对象	68
3.4.4	Beautiful Soup 的各种节点	69
3.4.5	Beautiful Soup 文档树搜索	69
3.4.6	Beautiful Soup 使用 CSS 选择器	70
3.5	实战：爬取壁纸站点的壁纸	70
3.6	正则表达式	74
3.6.1	re 模块	74
3.6.2	正则规则详解	75
3.6.3	正则练习	77
3.7	实战：爬取市级编码列表	79
3.7.1	获取所有市级的跳转链接列表	80
3.7.2	解析表格获得所有市级天气链接	81
3.7.3	提取市级编码	82
3.7.4	整合调整代码	83
第 4 章	用 CSV 和 Excel 存储数据	85
4.1	用 CSV 文件存储数据	85
4.1.1	CSV 写入	86
4.1.2	CSV 读取	87
4.2	实战：爬取星座运势	88
4.3	用 Excel 文件存储数据	89
4.3.1	Excel 写入	89
4.3.2	Excel 读取	90
4.4	实战：爬取某音乐平台排行榜	91
第 5 章	用数据库存储数据	99
5.1	MySQL 数据库	99
5.1.1	安装 MySQL	100
5.1.2	在 Windows 环境下安装 MySQL	100
5.1.3	在 Windows 环境下配置 MYSQL_HOME 环境变量	101

5.1.4	在 Windows 环境下设置 MySQL 登录密码	101
5.1.5	在 Windows 环境下启动或关闭 MySQL 服务	102
5.1.6	Mac 环境	103
5.1.7	Ubuntu 环境	103
5.1.8	MySQL 的基本操作	104
5.1.9	MySQL 数据库语法速成	106
5.1.10	Python 连接 MySQL 数据库	110
5.1.11	MySQL 特殊符号和表情问题	114
5.1.12	实战：抓取某技术网站数据	115
5.2	数据库可视化工具 DataGrip	122
5.2.1	建立数据库关联	122
5.2.2	编写 SQL 语句	123
5.2.3	常见问题：连接远程主机	124
5.3	Redis 数据库	125
5.3.1	安装 Redis	126
5.3.2	redis-py 库的安装	130
5.3.3	redis-py 基本操作示例	130
5.3.4	实战：爬取视频弹幕并保存到 Redis	134
5.4	MongoDB 数据库	137
5.4.1	安装 MongoDB	137
5.4.2	安装 PyMongo 库	140
5.4.3	PyMongo 基本操作示例	140
5.4.4	实战：爬取某电商网站关键字搜索结果并保存到 MongoDB	144
第 6 章	Python 应对反爬虫策略	148
6.1	反爬虫概述	148
6.1.1	为什么会出现反爬虫	149
6.1.2	常见的爬虫与反爬虫大战	149
6.2	反爬虫策略	150
6.2.1	User-Agent 限制	150
6.2.2	302 重定向	151
6.2.3	IP 限制	151
6.2.4	什么是网络代理	151
6.2.5	如何获取代理 IP	151
6.2.6	ADSL 拨号代理	152
6.2.7	Squid 配置代理缓存服务器	156

6.2.8	TinyProxy 配置代理缓存服务器	158
6.2.9	Cookie 限制	159
6.3	JavaScript 反爬虫策略	159
6.3.1	Ajax 动态加载数据	159
6.3.2	实战：爬取某素材网内容分析	159
6.3.3	数据请求分析	160
6.3.4	编写代码	163
6.4	Selenium 模拟浏览器操作	166
6.4.1	Selenium 简介	166
6.4.2	安装 Selenium	167
6.4.3	Selenium 常用函数	168
6.5	实战：爬取某网站的特定图	172
6.6	PhantomJS	175
6.6.1	在 Windows 上安装 PhantomJS	175
6.6.2	在 Mac 上安装 PhantomJS	175
6.6.3	在 Ubuntu 上安装 PhantomJS	176
6.6.4	关于 PhantomJS 的重要说明	176
6.7	常见验证码策略	176
6.7.1	图片验证码	177
6.7.2	实战：实现图片验证码自动登录	178
6.7.3	实战：实现滑动验证码自动登录	185
第 7 章	Python 爬虫框架 Scrapy (上)	196
7.1	Scrapy 框架简介与安装	197
7.1.1	Scrapy 相关信息	197
7.1.2	Scrapy 的安装	197
7.2	实战：爬取某网站每日壁纸	199
7.2.1	抓取目标分析	199
7.2.2	创建爬虫脚本	201
7.2.3	编写爬虫脚本	202
7.2.4	运行爬虫脚本	203
7.2.5	解析数据	203
7.3	Scrapy 架构简介	204
7.3.1	Scrapy 架构图	204
7.3.2	各个模块间的协作流程	205
7.3.3	协作流程拟人化对话版	206

7.4	Spider 详解	207
7.4.1	Spider 的主要属性和函数	207
7.4.2	Spider 运行流程	207
7.5	Request 类和 Response 类	209
7.5.1	Request 详解	209
7.5.2	Response 类常用参数、方法与子类	210
7.5.3	选择器	211
7.5.4	Scrapy Shell	212
7.6	Item 详解	213
7.7	Item Pipeline 详解	213
7.7.1	自定义 Item Pipeline 类	213
7.7.2	启用 Item Pipeline	214
7.8	实战：完善爬取每日壁纸的脚本	214
7.8.1	定义 BingItem	215
7.8.2	使用 ImagesPipeline	215
7.8.3	修改 Spider 代码	216
7.8.4	运行爬虫脚本	216
7.9	设置请求头	217
7.9.1	构造 Request 时传入	217
7.9.2	修改 settings.py 文件	217
7.9.3	为爬虫添加 custom_settings 字段	218
7.10	下载中间件详解	218
7.10.1	自定义 Downloader Middleware 类	218
7.10.2	启用自定义的代理下载中间件	219
7.11	实战：爬取某站点绘画频道的图片	219
7.11.1	分析爬取的站点	219
7.11.2	新建项目与明确爬取目标	221
7.11.3	创建爬虫爬取网页	221
7.11.4	设置代理	223
7.11.5	解析数据	223
7.11.6	存储数据	224
7.11.7	完善代码	226
第 8 章	Python 爬虫框架 Scrapy (下)	228
8.1	Scrapy 对接 Selenium	228
8.1.1	如何对接	228

8.1.2	对接示例：爬取某网站首页文章	229
8.2	实战：用 Scrapy 实现一个简单的代理池	232
8.2.1	代理池的设计	232
8.2.2	创建项目	232
8.2.3	编写获取 IP 的爬虫	233
8.2.4	编写检测 IP 的爬虫	238
8.2.5	编写调度程序	240
8.2.6	编写获取代理 IP 的接口	241
8.2.7	使用代理	243
8.3	用 Scrapyrt 调度 Scrapy	243
8.3.1	相关文档与安装 Scrapyrt	243
8.3.2	Scrapyrt GET 请求相关参数	244
8.3.3	Scrapyrt POST 请求相关参数	246
8.4	用 Docker 部署 Scrapy	246
8.4.1	Docker 简介	246
8.4.2	下载并安装 Docker	247
8.4.3	创建 Dockerfile	249
8.4.4	构建 Docker 镜像	250
8.4.5	把生成的 Docker 镜像推送到 Docker Hub	251
8.4.6	在云服务器上运行 Docker 镜像	253
第 9 章	数据分析案例：Python 岗位行情	254
9.1	数据爬取	254
9.2	NumPy 库和 pandas 库	258
9.2.1	ndarray 数组	259
9.2.2	ndarray 数组的常用操作	260
9.2.3	pandas 库	263
9.3	用 Matplotlib 实现数据可视化	268
9.3.1	Matplotlib 中文乱码问题	269
9.3.2	Matplotlib 绘制显示不全	270
9.3.3	用 Matplotlib 生成图表并进行分析	271
9.4	用 Wordcloud 库进行词云绘制	275
9.4.1	Wordcloud 简介	275
9.4.2	Wordcloud 构造函数与常用方法	276
9.4.3	词云绘制	277
9.5	小结	280

第 10 章 数据分析案例：某婚恋网站交友情况分析 .....	281
10.1 数据爬取 .....	281
10.2 安装 Jupyter Notebook .....	287
10.3 安装 pyecharts .....	288
10.4 数据分析 .....	289
10.4.1 读取 CSV 文件里的数据 .....	289
10.4.2 分析身高 .....	290
10.4.3 分析学历 .....	292
10.4.4 分析年龄 .....	292
10.4.5 分析城市 .....	294
10.4.6 分析交友宣言 .....	294
10.5 小结 .....	296

# 第 1 章

## Python 爬虫概念与 Web 基础

---

在开始学习如何编写 Python 爬虫之前，我们先学习一些和爬虫相关的概念和 Web 基础知识。

本章主要学习内容：

- 爬虫概念。
- HTTP 简述。
- 网页的组成（HTML，CSS，JavaScript）。

### 1.1 爬虫概念

---

本节将详细介绍与爬虫相关的一些概念，这将有助于我们后面的学习。

#### 1.1.1 什么是爬虫

爬虫，即网络爬虫，又称网络蜘蛛（Web Spider），是一种按照一定规则，用来自动浏览或抓取万维网数据的程序。可以把爬虫程序看成一个机器人，它的功能就是模拟人的行为去访问各种站点，或者带回一些与站点相关的信息。它可以 24 小时不间断地做一些重复性的工作，还可以自动提取一些数据。

有一点要注意，本章标题虽然写明 Python 爬虫，但并不是只有 Python 才能编写爬虫程序。其他的编程语言也可以用来编写爬虫程序，如 PHP 有 `phpspider` 爬虫框架，Java 有 `WebMagic` 爬虫框架，C# 有 `DotnetSpider` 爬虫框架等。

## 1.1.2 爬虫使用场景的引入

爬虫到底能干些什么呢？我们通过下面几个场景引入。

### 场景一：

大部分读者应该都有看网络小说的习惯，而去正站点看小说一般是需要付费的，所以衍生了很多盗版小说站点。这种盗版小说站点一般通过挂广告的形式来盈利，在浏览器底部会嵌入各种类型的广告，用户单击这些广告会打开其对应的网址，盗版小说站点以此获利。

这些盗版小说站点的广告一般都是误导读者打开它。比如，广告上有个 X 按钮，读者通常以为单击 X 按钮可以关闭这个广告，殊不知是打开了广告，而且有些广告中含有不堪入目的内容。想一想，在拥挤的地铁里，旁边的人看到你打开了这样的网页，会有多尴尬。有没有办法既能看到小说又不用看这些烦人的广告呢？

答：通过爬虫可以解决这个问题，让爬虫只解析小说部分的内容并显示出来，甚至可以把整本小说的内容解析完，保存到本地，以便离线阅读。

### 场景二：

有些读者有一种类似于强迫症的行为，比如，快递预计今天会到，每隔一段时间你就会不由自主地输入单号看看快递到了没。有没有办法能摆脱这种频繁而又枯燥的工作呢？

答：对于这种轮询（每隔一段时间查询一次）的任务，我们可以编写一个定时爬虫，每隔一段时间自动去请求相应的站点，然后处理结果，判断其是否符合我们的预期，与人工刷新相比，爬虫的频率更快、效率更高。

### 场景三：

有一些站点，会通过签到或以做每日任务的形式来提高用户的活跃度，当坚持到一定天数后会发放一些小奖励。比如，坚持十五天就能获得一个小礼物。因为忙碌或其他原因，你中断了签到，使得之前的努力就都白费了。而且有些每日任务非常枯燥，但每天要为此花上好几分钟。有没有办法把这种任务交给程序来做，然后坐享其成呢？

答：可以把这些任务交给爬虫。对于签到这种操作，通过抓包获取签到所需的接口、参数和请求规则，接着让脚本每天定时执行即可。每日任务一般通过模拟单击的方式完成，通过 Selenium 自动化框架来模拟。对于 App 签到，则可以通过 Appium 移动端自动化测试框架来完成。

### 场景四：

如果你有志于从事 Python 开发相关的工作，想了解这个行业的一些情况，比如薪资、年限和相关要求等，而你身边又没有从事相关工作的朋友，怎样才能获取到这些信息呢？

答：可以编写爬虫爬取一些招聘网站，比如拉勾、前程无忧这类站点，把 Python 岗位相关的信息都抓取下来，然后通过数据分析基础三件套（NumPy、pandas、Matplotlib）进行基本的数据分析，以此获取和这个行业相关的一些信息。

### 1.1.3 爬虫的组成部分

相信在看完爬虫应用的四个场景后，读者对爬虫能做什么有了大致了解，接下来我们来了解爬虫由哪几部分组成。爬虫的三个组成部分如图 1.1 所示。

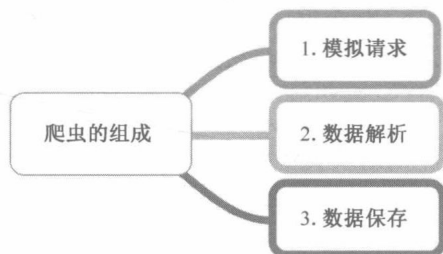


图 1.1 爬虫的三个组成部分

### 1.1.4 模拟请求

模拟请求就是如何把我们的爬虫伪装得像一个人一样去访问互联网站点。

- 最简单的站点，什么都不处理，只要发送请求，就会给出相应的结果。
- 稍微复杂一点的站点，会判断请求头中的 User-Agent 是否为浏览器请求，Host 字段是否为正确的服务器域名，以及 Referer 字段的地址是否合法。
- 再复杂一点的站点，需要登录后才能访问。登录后会持有有一个 Cookie 或 Session 会话，你需要带着这个东西才能执行一些请求，否则都会跳到登录页。
- 更复杂一点的站点，登录很复杂，需要五花八门的验证码、最简单的数字图片加噪点、滑动验证码、点触验证码。除此之外，还有一些其他特立独行的验证方式，如最经典的微博宫格验证码、极验证码的行为验证等。
- 还有更复杂的站点，其链接和请求参数都是加密的，需要研究、破解加密规则才能够模拟访问。
- 除此之外，还有一些反爬虫套路，如限制 IP 访问频次、JavaScript 动态加载数据等。

在模拟请求之前，先要了解请求规则，一般通过抓包工具来完成。

(1) 对最简单的浏览器请求（以 Chrome 谷歌浏览器为例），在网页空白处右键单击并选择检查，或者依次单击如图 1.2 所示的“更多工具”→“开发者工具”。在 Windows 中打开开发者工具的快捷键 F12。

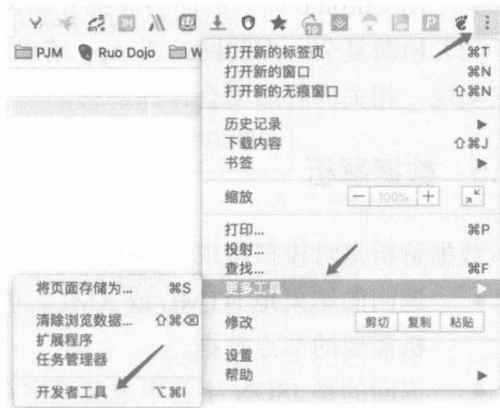


图 1.2 打开 Chrome 开发者工具