

O'REILLY®



华章 IT

原书第2版

利用Python 进行数据分析

Python for Data Analysis: Data Wrangling with Pandas, NumPy,
and IPython



powered by



Wes McKinney 著

徐敬一 译



机械工业出版社
China Machine Press

原书第2版

利用 Python 进行数据分析

韦斯·麦金尼 (Wes McKinney) 著
徐敬一 译

Beijing • Boston • Farnham • Sebastopol • Tokyo

O'REILLY®

O'Reilly Media, Inc. 授权机械工业出版社出版

机械工业出版社

图书在版编目 (CIP) 数据

利用 Python 进行数据分析 (原书第 2 版)/(美) 韦斯·麦金尼 (Wes McKinney) 著;
徐敬一译. - 北京: 机械工业出版社, 2018.6 (2018.9 重印)
(O'Reilly 精品图书系列)

书名原文: Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython, 2nd Edition

ISBN 978-7-111-60370-2

I. 利… II. ①韦… ②徐… III. 统计分析—应用软件 IV. C819

中国版本图书馆 CIP 数据核字 (2018) 第 150175 号

北京市版权局著作权合同登记

图字: 01-2017-8013 号

Copyright © 2018 William McKinney. All rights reserved.

Simplified Chinese Edition, jointly published by O'Reilly Media, Inc. and China Machine Press, 2018.
Authorized translation of the English edition, 2018 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 2018。

简体中文版由机械工业出版社出版 2018。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。

封底无防伪标均为盗版

本书法律顾问

北京大成律师事务所 韩光 / 邹晓东

书 名 / 利用 Python 进行数据分析 (原书第 2 版)

书 号 / ISBN 978-7-111-60370-2

责任编辑 / 冯秀泳

封面设计 / Karen Montgomery, 张健

出版发行 / 机械工业出版社

地 址 / 北京市西城区百万庄大街 22 号 (邮政编码 100037)

印 刷 / 三河市宏图印务有限公司

开 本 / 178 毫米 × 233 毫米 16 开本 30.75 印张

版 次 / 2018 年 7 月第 1 版 2018 年 9 月第 2 次印刷

定 价 / 119.00 元 (册)

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010)88379426, 88361066

购书热线: (010)68326294, 88379649, 68995259

投稿热线: (010)88379604

读者信箱: hzit@hzbook.com

译者序

2014年我刚走上工作岗位时，Python数据分析技术开始在中国流行，越来越多的人开始学习和使用这门技术。当时，Python的主流版本是2.7, 3.4版本刚刚发布，但认可度和流行度并不高。很幸运，我阅读了本书的第1版。第1版基于Python 2.7，介绍了大量的实用技术，包括pandas、NumPy的使用等。在实践中，书中的技术被我大量应用，帮助我更快更好地处理了各类数据。

时过境迁，从本书英文版第1版2012年出版至今，已经过去了6年。这6年中，Python的主流版本从2.7升级到了3.6。无论是否情愿，大部分Pythoner都不得不学会适应新版本；而pandas则从0.1.0版本迭代到如今的0.22.0版本。版本号的持续增加意味着新技术、新特性的不断丰富。举例来说，将带有多层索引的数据透视表写入Excel在2015年之前是无法使用pandas完成的，在0.17版本后该功能被加入。因此，本书的第一版内容已经略显落后。

2017年10月下旬，本书作者Wes McKinney先生更新了本书的第2版。在第2版中，他将Python版本更新到3.6，介绍了pandas的一些新接口和功能，并新增了大量现实世界的数据分析实例，以确保本书的可实践性。我在2017年11月18日接到出版社的翻译邀请后便开始了翻译工作。McKinney先生的写作风格朴实、形象，因而我的翻译过程较为顺畅，但书中部分口语化的叙述也因为中英文表达方式的差异而增加了意译的难度。

在本书翻译过程中，我得到了很多帮助。首先要感谢华章公司的王春华编辑和冯秀泳编辑，他们在翻译过程中给了我耐心指导。在审稿时，来自国内Python圈的朋友们对本书进行了仔细而全面的审核，他们是网易数据工程师马喆诚、早稻田大学研究生梁婷、大疆工程师王波。感谢他们对本书的付梓而做出的贡献。此外，我还要感谢我的女朋友易

慧娟，感谢你在生活中对我的各种好，我爱你！

为了让国内读者在第一时间读到这本畅销国外技术著作，出版社和我都加快了工作进度，但时间紧、任务重，再加之本人水平有限，翻译工作中难免会出现一些失误。欢迎读者将阅读过程中发现的问题发送至我的邮箱（xujingyi46@163.com）。

本书英文版的副书名是“Data Wrangling with Pandas, NumPy, and IPython”，其中Wrangling是一个很难直译的词汇，它的原意是争执、争论，但在书中它描述的是将数据进行规整、处理的意思。希望读者读完本书后，可以使用好pandas、NumPy和IPython这些工具，更好地完成数据处理、分析的学习和工作。Enjoy it！

徐敬一

2018年5月于中国工商银行后台中心

目录

前言	1
第 1 章 准备工作	7
1.1 本书内容	7
1.1.1 什么类型的数据	7
1.2 为何利用 Python 进行数据分析	8
1.2.1 Python 作为胶水	8
1.2.2 解决“双语言”难题	8
1.2.3 为何不使用 Python	9
1.3 重要的 Python 库	9
1.3.1 NumPy	9
1.3.2 pandas	10
1.3.3 matplotlib	11
1.3.4 IPython 与 Jupyter	11
1.3.5 SciPy	12
1.3.6 scikit-learn	12
1.3.7 statsmodels	13
1.4 安装与设置	13
1.4.1 Windows	14
1.4.2 Apple (OS X 和 macOS)	14
1.4.3 GNU/Linux	14
1.4.4 安装及更新 Python 包	15
1.4.5 Python 2 和 Python 3	16
1.4.6 集成开发环境和文本编辑器	16
1.5 社区和会议	17
1.6 快速浏览本书	17
1.6.1 代码示例	18
1.6.2 示例数据	18

1.6.3 导入约定	18
1.6.4 术语	19
第 2 章 Python 语言基础、IPython 及 Jupyter notebook ...	20
2.1 Python 解释器	21
2.2 IPython 基础	22
2.2.1 运行 IPython 命令行	22
2.2.2 运行 Jupyter notebook	23
2.2.3 Tab 补全	25
2.2.4 内省	27
2.2.5 %run 命令	28
2.2.6 执行剪贴板中的程序	30
2.2.7 终端快捷键	30
2.2.8 关于魔术命令	31
2.2.9 matplotlib 集成	33
2.3 Python 语言基础	34
2.3.1 语言语义	34
2.3.2 标量类型	42
2.3.3 控制流	49
第 3 章 内建数据结构、函数及文件 ...	54
3.1 数据结构和序列	54
3.1.1 元组	54
3.1.2 列表	57
3.1.3 内建序列函数	61
3.1.4 字典	64
3.1.5 集合	67
3.1.6 列表、集合和字典的推导式	69
3.2 函数	72
3.2.1 命名空间、作用域和本地函数	72
3.2.2 返回多个值	73
3.2.3 函数是对象	74
3.2.4 匿名 (Lambda) 函数	75
3.2.5 柯里化：部分参数应用	76
3.2.6 生成器	77
3.2.7 错误和异常处理	79

3.3 文件与操作系统	82
3.3.1 字节与 Unicode 文件.....	85
3.4 本章小结.....	86

第 4 章 NumPy 基础：数组与向量化计算 87

4.1 NumPy ndarray：多维数组对象	89
4.1.1 生成 ndarray	90
4.1.2 ndarray 的数据类型.....	92
4.1.3 NumPy 数组算术.....	94
4.1.4 基础索引与切片	95
4.1.5 布尔索引	100
4.1.6 神奇索引	103
4.1.7 数组转置和换轴	104
4.2 通用函数：快速的逐元素数组函数	106
4.3 使用数组进行面向数组编程.....	109
4.3.1 将条件逻辑作为数组操作	110
4.3.2 数学和统计方法	111
4.3.3 布尔值数组的方法.....	113
4.3.4 排序	114
4.3.5 唯一值与其他集合逻辑.....	115
4.4 使用数组进行文件输入和输出	115
4.5 线性代数.....	116
4.6 伪随机数生成.....	118
4.7 示例：随机漫步	120
4.7.1 一次性模拟多次随机漫步	121
4.8 本章小结.....	122

第 5 章 pandas 入门 123

5.1 pandas 数据结构介绍	123
5.1.1 Series	123
5.1.2 DataFrame.....	128
5.1.3 索引对象	134
5.2 基本功能	135
5.2.1 重建索引	136
5.2.2 轴向上删除条目	138
5.2.3 索引、选择与过滤	140

5.2.4 整数索引	144
5.2.5 算术和数据对齐	145
5.2.6 函数应用和映射	150
5.2.7 排序和排名	152
5.2.8 含有重复标签的轴索引	155
5.3 描述性统计的概述与计算	157
5.3.1 相关性和协方差	159
5.3.2 唯一值、计数和成员属性	161
5.4 本章小结	164

第 6 章 数据载入、存储及文件格式 165

6.1 文本格式数据的读写	165
6.1.1 分块读入文本文件	171
6.1.2 将数据写入文本格式	172
6.1.3 使用分隔格式	174
6.1.4 JSON 数据	176
6.1.5 XML 和 HTML：网络抓取	177
6.2 二进制格式	180
6.2.1 使用 HDF5 格式	181
6.2.2 读取 Microsoft Excel 文件	183
6.3 与 Web API 交互	184
6.4 与数据库交互	186
6.5 本章小结	187

第 7 章 数据清洗与准备 188

7.1 处理缺失值	188
7.1.1 过滤缺失值	189
7.1.2 补全缺失值	191
7.2 数据转换	194
7.2.1 删除重复值	194
7.2.2 使用函数或映射进行数据转换	195
7.2.3 替代值	197
7.2.4 重命名轴索引	198
7.2.5 离散化和分箱	199
7.2.6 检测和过滤异常值	202
7.2.7 置换和随机抽样	203
7.2.8 计算指标 / 虚拟变量	204

7.3 字符串操作	207
7.3.1 字符串对象方法	208
7.3.2 正则表达式	210
7.3.3 pandas 中的向量化字符串函数	213
7.4 本章小结	215

第 8 章 数据规整：连接、联合与重塑 216

8.1 分层索引	216
8.1.1 重排序和层级排序	219
8.1.2 按层级进行汇总统计	220
8.1.3 使用 DataFrame 的列进行索引	220
8.2 联合与合并数据集	221
8.2.1 数据库风格的 DataFrame 连接	222
8.2.2 根据索引合并	226
8.2.3 沿轴向连接	230
8.2.4 联合重叠数据	235
8.3 重塑和透视	236
8.3.1 使用多层索引进行重塑	236
8.3.2 将“长”透视为“宽”	240
8.3.3 将“宽”透视为“长”	242
8.4 本章小结	244

第 9 章 绘图与可视化 245

9.1 简明 matplotlib API 入门	245
9.1.1 图片与子图	246
9.1.2 颜色、标记和线类型	250
9.1.3 刻度、标签和图例	252
9.1.4 注释与子图加工	255
9.1.5 将图片保存到文件	258
9.1.6 matplotlib 设置	258
9.2 使用 pandas 和 seaborn 绘图	259
9.2.1 折线图	259
9.2.2 柱状图	262
9.2.3 直方图和密度图	266
9.2.4 散点图或点图	269
9.2.5 分面网格和分类数据	270

9.3 其他 Python 可视化工具	271
9.4 本章小结	272
第 10 章 数据聚合与分组操作	274
10.1 GroupBy 机制	274
10.1.1 遍历各分组	278
10.1.2 选择一列或所有列的子集	279
10.1.3 使用字典和 Series 分组	280
10.1.4 使用函数分组	281
10.1.5 根据索引层级分组	282
10.2 数据聚合	282
10.2.1 逐列及多函数应用	284
10.2.2 返回不含行索引的聚合数据	287
10.3 应用：通用拆分 - 应用 - 联合	288
10.3.1 压缩分组键	290
10.3.2 分位数与桶分析	291
10.3.3 示例：使用指定分组值填充缺失值	292
10.3.4 示例：随机采样与排列	294
10.3.5 示例：分组加权平均和相关性	296
10.3.6 示例：逐组线性回归	298
10.4 数据透视表与交叉表	298
10.4.1 交叉表：crosstab	301
10.5 本章小结	302
第 11 章 时间序列	303
11.1 日期和时间数据的类型及工具	303
11.1.1 字符串与 datetime 互相转换	305
11.2 时间序列基础	307
11.2.1 索引、选择、子集	308
11.2.2 含有重复索引的时间序列	311
11.3 日期范围、频率和移位	312
11.3.1 生成日期范围	313
11.3.2 频率和日期偏置	316
11.3.3 移位（前向和后向）日期	317
11.4 时区处理	320
11.4.1 时区的本地化和转换	320

11.4.2 时区感知时间戳对象的操作	323
11.4.3 不同时区间的操作	324
11.5 时间区间和区间算术	324
11.5.1 区间频率转换	326
11.5.2 季度区间频率	327
11.5.3 将时间戳转换为区间（以及逆转换）.....	329
11.5.4 从数组生成 PeriodIndex	330
11.6 重新采样与频率转换	332
11.6.1 向下采样	334
11.6.2 向上采样与插值	336
11.6.3 使用区间进行重新采样	337
11.7 移动窗口函数	339
11.7.1 指数加权函数	342
11.7.2 二元移动窗口函数	343
11.7.3 用户自定义的移动窗口函数	344
11.8 本章小结	344

第 12 章 高阶 pandas 346

12.1 分类数据	346
12.1.1 背景和目标	346
12.1.2 pandas 中的 Categorical 类型	348
12.1.3 使用 Categorical 对象进行计算	350
12.1.4 分类方法	352
12.2 高阶 GroupBy 应用	355
12.2.1 分组转换和“展开” GroupBy	355
12.2.2 分组的时间重新采样	359
12.3 方法链技术	361
12.3.1 pipe 方法	362
12.4 本章小结	363

第 13 章 Python 建模库介绍 364

13.1 pandas 与建模代码的结合	364
13.2 使用 Patsy 创建模型描述	367
13.2.1 Patsy 公式中的数据转换	369
13.2.2 分类数据与 Patsy	371
13.3 statsmodels 介绍	373

13.3.1 评估线性模型	374
13.3.2 评估时间序列处理	377
13.4 scikit-learn 介绍	377
13.5 继续你的教育	381
第 14 章 数据分析示例	382
14.1 从 Bitly 获取 1.USA.gov 数据	382
14.1.1 纯 Python 时区计数	383
14.1.2 使用 pandas 进行时区计数	385
14.2 MovieLens 1M 数据集	392
14.2.1 测量评价分歧	396
14.3 美国 1880 ~ 2010 年的婴儿名字	397
14.3.1 分析名字趋势	402
14.4 美国农业部食品数据库	410
14.5 2012 年联邦选举委员会数据库	416
14.5.1 按职业和雇主的捐献统计	419
14.5.2 捐赠金额分桶	421
14.5.3 按州进行捐赠统计	423
14.6 本章小结	424
附录 A 高阶 NumPy	425
附录 B 更多 IPython 系统相关内容	457

前言

第 2 版新内容

本书第 1 版出版于 2012 年，彼时基于 Python 的开源数据分析库（例如 pandas）仍然是一个发展迅速的新事物。在本次更新、拓展的第 2 版中，我在一些章节内进行了修改，以解释过去 5 年中发生的不兼容的变更、弃用和一些新特性。此外，我还添加了新内容，用以介绍在 2012 年还不存在或者不成熟的工具。最后，我会避免把一些新兴的或者不太可能走向成熟的开源项目写入本书。我希望本版的读者能够发现本书内容在 2020 年或者 2021 年仍然几乎像在 2017 年一样适用。

第 2 版中的主要更新包括：

- 所有的代码，包括把 Python 的教程更新到了 Python 3.6 版本（第 1 版中使用的是 Python 2.7）
- 更新了 Python 第三方发布版 Anaconda 和其他所需 Python 包的安装指引
- 更新 pandas 库到 2017 年的最新版
- 新增一章，关于更多高级 pandas 工具和一些使用提示
- 新增 statsmodels 和 scikit-learn 的简明使用介绍

除了以上更新内容，我还重新组织了第 1 版的部分重要内容，使本书对新手来说更易于理解。

本书约定

以下印刷约定将在本书中使用：

斜体 (*Italic*)

表示新的术语、URL、email 地址、文件名和文件扩展名。

等宽字体 (**Constant width**)

用于程序清单以及段落中的程序元素，例如变量名、函数名、数据库、数据类型、环境变量、表达式和关键字等。

等宽粗体 (**Constant width bold**)

表示命令或其他应当由用户键入的文本。

等宽斜体 (*Constant width italic*)

表示应当由用户提供值来替代的文本，或者其他由上下文决定的值。

本符号表示提示或建议。



本符号表示一般性说明。



本符号表示警告。



使用代码示例

可以通过本书的 GitHub 仓库获取本书每一章中的数据文件和相关材料。GitHub 仓库地址：<http://github.com/wesm/pydata-book>。

本书的目的在于帮助你完成工作。一般来说，本书提供的示例代码，你可以在你的程序或文档中使用而无须联系我们获取许可，除非你需要重造大量代码。举例来说，使用本书中的代码段编写程序无须授权许可，但销售或发行 O'Reilly 图书的 CD-ROM 代码示例则需要许可。引用本书代码回答问题不需要许可，但在你的产品文档中大量使用本书示例代码则需要许可。

我们鼓励注明资料来源的行为，但这并不是必需的。来源注明通常包括书名、作者、出版社及 ISBN，例如：“*Python for Data Analysis* by Wes McKinney(O'Reilly). Copyright 2017 Wes McKinney, 978-1-491-95766-0”。

如果你认为你对本书示例代码的使用超过了正常使用范围或者需要以上介绍的授权许可,请联系 permissions@oreilly.com。

O'Reilly Safari

 **Safari**[®] Safari (前身为 Safari Books Online) 是一个会员制的培训、参考网站,服务于企业、政府、教育者和个人。

会员可以访问数千书籍、培训视频、学习路径、交互教程和超过 250 家出版商的企划列表,包括 O'Reilly Media、Harvard Business Review、Prentice Hall Professional、Addison-Wesley Professional、Microsoft Press、Sams、Que、Peachpit Press、Adobe、Focal Press、Cisco Press、John Wiley & Sons、Syngress、Morgan Kaufmann、IBM Redbooks、Packt、Adobe Press、FT Press、Apress、Manning、New Riders、McGraw-Hill、Jones & Bartlett、Course Technology 等。

更多信息,请访问 <http://oreilly.com/safari>。

如何联系我们

对于本书如果有任何意见或疑问,请按照以下地址联系本书出版商。

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

我们为本书准备了一个网页,用于陈列勘误、示例和其他附加信息。访问地址是: http://bit.ly/python_data_analysis_2e。

针对本书评论或提出技术问题,请发送邮件至: bookquestions@oreilly.com。

关于本书的更多信息、课程、会议及新闻,请访问我们的网站: <http://www.oreilly.com>。

Facebook 联系我们: <http://facebook.com/oreilly>

Twitter 联系我们: <http://twitter.com/oreillymedia>

YouTube 观看我们的视频: <http://www.youtube.com/oreillymedia>

致谢

本书是全世界很多人多年来富有成效的讨论、协作和支持的成果。我想对他们中的一些代表致以谢意。

怀念：John D. Hunter (1968—2012)

我们亲爱的朋友和同行 John D. Hunter 在经历了一场与结肠癌的战斗后，于 2012 年 8 月 28 日离开了世界。那时正是我完成本书第 1 版最终手稿后不久。

John 对 Python 科学计算和数据社区的影响之大难以估量，他给我们留下的遗产价值非凡。除了在 2000 年初期开发 matplotlib 之外（那时 Python 还没有当下如此流行），他还帮助塑造了一代核心开源开发者的文化，如今这些开发者已经成为 Python 生态系统的顶梁柱，而 Python 生态系统对于现如今的我们来说似乎是理所当然的。

在 2010 年 1 月，我开源生涯的早期，那时候 pandas 刚刚发布了 0.1 版本，我便有幸结识了 John。即便在最黑暗的时期，他的才华和指导仍在帮助我推动 pandas 前进，实现 Python 成为数据分析第一语言的愿景。

John 与 IPython、Jupyter 项目的先锋 Fernando Pérez、Brian Granger 及其他很多 Python 社区的倡议人联系紧密。我们四人曾经希望共同写作一本书，但只有我个人时间最为自由，所以这个想法被搁置了。我非常确信他会为过去 5 年中我们个人及我们社区所取得的成就感到骄傲。

第 2 版致谢 (2017)

距离我在 2012 年 7 月完成第 1 版手稿已经 5 年了。很多事情都发生了变化。Python 社区获得了极大的成长，围绕 Python 的开源软件生态系统也十分繁荣。pandas 核心开发者孜孜不倦的付出，使得 pandas 项目高速成长，也使得 pandas 的用户群体遍布 Python 数据科学生态系统的各个角落，没有他们本书将不会存在。pandas 的核心开发者包括但不限于：Tom Augspurger、Joris van den Bossche、Chris Bartak、Phillip Cloud、gfyoung、Andy Hayden、Masaaki Horikoshi、Stephan Hoyer、Adam Klein、Wouter、Overmeire、Jeff Reback、Chang She、Skipper Seabold、Jeff Tratner 和 y-p。

在第 2 版的实际写作过程中，非常感谢 O'Reilly 的工作人员在写作进程中给予的耐心帮助。他们是 Marie Beaugureau、Ben Lorica 和 Colleen Toporek。我再次得到了优秀技术审阅人的支持，他们是 Tom Augspurger、Paul Barry、Hugh Brown、Jonathan Coe 和 Andreas Müller。感谢你们。