

Python

统计分析

[奥] 托马斯·哈斯尔万特(Thomas Haslwanter) 著 李锐 译 张志杰 审

An Introduction
to Statistics
with Python

With Applications in the Life Sciences

非
外
借



中国工信出版集团



人民邮电出版社
POSTS & TELECOM PRESS

Python

统计分析

[奥] 托马斯·哈斯尔万特(Thomas Haslwanter) 著 李锐 译 张志杰 审



人民邮电出版社
北京

图书在版编目(CIP)数据

Python统计分析 / (奥) 托马斯·哈斯尔万特
(Thomas Haslwanter) 著 ; 李锐译. — 北京 : 人民邮
电出版社, 2018. 12
ISBN 978-7-115-49384-2

I. ①P… II. ①托… ②李… III. ①统计分析—应用
软件 IV. ①C819

中国版本图书馆CIP数据核字(2018)第214278号

版权声明

Translation from English language edition

An Introduction to Statistics with Python: With Applications in the Life Sciences

by Thomas Haslwanter

Copyright © Springer International Publishing Switzerland 2016

This Springer imprint is published by Springer Nature

The registered company is Springer International Publishing AG

All Rights Reserved

本书中文简体字版由施普林格出版社授权人民邮电出版社出版。未经出版者书面许可, 不得以
任何方式复制或抄袭本书任何部分。

版权所有, 侵权必究。

◆ 著 [奥] 托马斯·哈斯尔万特(Thomas Haslwanter)

译 李 锐

审 张志杰

责任编辑 王峰松

责任印制 焦志炜

◆ 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 <http://www.ptpress.com.cn>

北京缤索印刷有限公司印刷

◆ 开本: 720×960 1/16

印张: 15.25

字数: 295 千字

2018 年 12 月第 1 版

印数: 1—3 000 册

2018 年 12 月北京第 1 次印刷

著作权合同登记号 图字: 01-2017-5035 号

定价: 79.00 元

读者服务热线: (010)81055410 印装质量热线: (010)81055316

反盗版热线: (010)81055315

广告经营许可证: 京东工商广登字 20170147 号

内容提要

本书以基础的统计学知识和假设检验为重点，简明扼要地讲述了 Python 在数据分析、可视化和统计建模中的应用。本书主要包括 Python 的简单介绍、研究设计、数据管理、概率分布、不同数据类型的假设检验、广义线性模型、生存分析和贝叶斯统计学等从入门到高级的内容。

本书利用 Python 这门开源语言，不仅在直观上对数据分析和统计检验提供了很好的讲解，而且在相关数学公式的讲解上也能够做到深入浅出。本书所讲内容的可操作性很强，配套提供相关的代码和数据，方便读者动手练习。

本书适合对统计学和 Python 有兴趣的读者，特别是在实验学科中需要利用 Python 的强大功能来进行数据处理和统计分析的学生和研究人员。

作者简介

托马斯·哈斯尔万特 (Thomas Haslwanter) 在学术机构中有超过 10 年的教学经验, 是林茨上奥地利州应用科学大学 (University of Applied Sciences Upper Austria in Linz) 医学工程系的教授, 瑞士苏黎世联邦理工学院讲师, 并曾在澳大利亚悉尼大学和德国图宾根大学担任过研究员。他在医学研究方面经验丰富, 专注于眩晕症的诊断、治疗和康复。在深入使用 Matlab 软件 15 年后, 他发现 Python 非常强大, 并将其用于统计分析、声音和图像处理以及生物仿真应用。

译者简介

李锐, 复旦大学公共卫生学院流行病与生物统计专业博士生, Python、R 和 Lisp 语言的爱好者, 主要研究方向为统计学习和机器学习建模以及组学数据的数据挖掘。先后以第一作者身份发表学术论文 6 篇, 其中 SCI 论文 4 篇。参编中文专著 2 本。

审校者简介

张志杰, 复旦大学公共卫生学院副教授, 多本国际医疗卫生期刊的特邀编辑, 研究方向为统计建模和医学领域的统计分析方法。参与并完成国家重大科技专项、“863”计划项目、国家“十五”科技攻关课题、自然科学基金重大项目等多项国家级课题的研究, 研究成果先后获 2010 年全国百篇优秀博士学位论文、2012 年上海市医学奖二等奖、上海市科技进步奖二等奖以及中华医学奖三等奖, 2011 年入选复旦大学首批“卓学人才计划”, 2013 年入选上海市第二批新优青人才计划。

前言

在我自己的研究工作中，当我进行数据分析时，有两件事情经常让我深陷泥沼：(1) 我了解的统计学知识不够多；(2) 市面上的书籍大多是理论性的，缺少一些实践上的帮助。你手上拿着的（或者在你的平板电脑或笔记本电脑上的）这本书就是要解决这两个关键问题。这本书将为你提供足够的统计学领域知识，这样你就不会迷失，与此同时，这本书也会教你使用统计学分析所需要的工具。我认为 Python 提供的解决方案能够解决生物学家、物理学家、医学博士们在他们的工作中遇到的 90% 的问题。如果你正在攻读研究生学位，或者是一个正在分析最近实验数据的医学研究者，那么你将会在本书中找到你所需要的工具和它们的使用说明及源代码。

基于上述原因，本书将重点讲解统计学的基础知识和假设检验，并简单地介绍其他的统计学方法。我明白本书中介绍的大多数统计学检验也可以使用统计学建模的方法来完成。但是在大多数情况下，统计学建模并不是生命科学领域的期刊所使用的方法论。高级的统计分析超出了本书的范围，并且坦率地说，也超出了我对统计学的了解。

本书主要使用 Python 语言进行统计学检验和数据分析，主要是基于下面两点原因：首先，我希望这些方法能够被所有人使用。尽管市面上有一些商业的解决方案，比如 Matlab、SPSS、Minitab 等，它们提供了强大的分析工具，但是这些软件大多数只能在学术环境中合法使用¹，而 Python 则是完全免费的²（“像啤酒一样免费”经常能够在 Python 社区中听到）。另一个原因是，Python 是我见过的最优美的编程语言之一，大约在 2010 年，Python 及其文档就发展得较为成熟了，这使得一个业余的编码人员也能够轻松地使用它。配合这本书一起使用，Python 和 Python 生态系统提供的优美又免费的工具包，将覆盖大多数研究者一辈子所需要了解的所有统计学分析。

这本书是为谁写的

这本书的前提条件如下：

- 你有一些基本的编程经验。如果你之前从来没有接触过编程，你最好先从学习 Python 语言开始，在书中我提供了一些非常好的 Python 学习参考链

1 应该是免费的合法使用。——译者注。

2 free 是双关，自由软件也叫 free software。——译者注。

接给你。同时学习编程和学习统计学可能有点拔苗助长了。

- 你不必是一个统计学专家。如果你有高级的统计学分析的经验，那么借助 Python 和 Python 包的在线帮助文档，你马上就能够进行大部分的数据分析。尽管本书会让你熟悉 Python 编程，但主要还是聚焦在统计学的基本概念和假设检验上，只有在本书的最后一个部分才会涉及线性回归建模和贝叶斯统计学等内容。

本书旨在提供所有（至少大部分）你需要的统计学分析工具。我在本书中会提供足够的理论背景知识帮助你明白你正在做什么。如非必要，我不会证明任何的定理或应用数学。对于本书中提到的所有统计学检验，我都会提供一个能够正常运行的 Python 程序。总的来说，你只需要定义好你的问题，选择合适的程序，稍微修改一下程序让它符合你的需求。这样的话，就算你没有太多的 Python 编程经验，你也能快速上手。我并没有提供给你单独的 Python 包，因为我希望你能够根据你自己的需求对每个程序进行修改以适应你的设置（数据类型、自定义的绘图的标签和返回值等）。

全书共分为 3 个部分。

第一部分 简单介绍 Python：如何安装和配置 Python 运行环境，运行一些简单的 Python 程序；为了防止你犯一些常见错误，我们也提供了一些小建议。这部分也会介绍如何从不同的数据源读取数据到 Python 中，并对数据进行可视化。

第二部分 介绍统计学分析：如何进行研究设计，如何分析数据，概率分布的基本知识，概述常见的重要假设检验方法。尽管现代统计学扎根于统计学建模，但是假设检验仍然占据着生命科学领域的主导地位。对于每一个检验方法，我们都会提供一个 Python 程序来展示该检验是如何用 Python 语言完成的。

第三部分 介绍统计建模的知识并简单介绍高级统计分析的步骤。因为 logistic 回归等离散型数据检验方法使用了一种叫作“广义线性模型”的高级统计学方法，所以在这个部分中，也会包含这些内容。在本书的最后，将会展示贝叶斯统计学中的基本概念。

补充材料

随本书出版的还有大量的 Python 程序和示例数据，均可以在网上获取。这些程序包括：所有书中的程序，每章末尾示例的答案，每个检验方法的示例代码。此外还包括书中插图的绘图代码，运行代码所需要的数据，等等。

本书附带的 Python 程序和数据集可以在 Github 代码库上下载（https://github.com/thomas-haslwanter/statsintro_python），所有的材料都可以在 <http://www.springer.com/de/book/9783319283159> 下载。

致谢

Python 是由许多用户社区做出的贡献组成的，本书中的一些章节也基于互联网上的优秀信息（已获得作者同意在本书中引用他们的内容）。

我要特别感谢下面这些人。

- Paul E Johnson 阅读了全书书稿，并对全书的内容组织和统计细节方面提出了很有价值的反馈建议。
- Connor Johnson 写了篇博客解释了 statsmodels OLS（最小二乘法）命令的结果，这篇博客是本书“统计学模型”的基础。
- Cam Davidson Pilon 写了一本名叫《Probabilistic-Programming-and-Bayesian-Methods-for-Hackers》的开源电子书，我从中借鉴了 Challenger disaster 的例子用来阐述贝叶斯统计学。
- 多亏了 Fabian Pedregosa 的一篇关于有序 logistic 回归的博客，让我在本书中加入了相关的内容，因为我对这部分内容并不熟悉。

我还想感谢 Carolyn Mayer，他阅读了我的手稿并将一些口语化的语言润色为正式的书面语。此外我还要特别感谢我的妻子，她不仅对全书的组织结构提出了重要的建议，还提出了很多编程教学的小技巧，并且本书中和茶有关的主题都有她的支持和协助。

如果你有任何建议或者勘误，请给我的工作邮箱（thomas.haslwanger@fh-linz.at）发电子邮件。除非另有通知，如果我基于你的反馈建议对本书做出了改变，我会将你加入到贡献者名单中。如果你愿意附上错误出现位置的句子，哪怕只有一部分，也能让我更容易定位错误。页码和章节名也很好，但不如句子容易定位。非常感谢！

Thomas Haslwanger
奥地利林茨
2015 年 12 月

缩 写

ANOVA	方差分析
CDF	累积分布函数
CI	置信区间
DF/DOF	自由度
EOL	行末
GLM	广义线性模型
HTML	超文本标记语言
IDE	集成开发环境
IQR	四分位数间距
ISF	逆生存函数
KDE	核密度估计
MCMC	马尔可夫链蒙特卡洛
NAN	不是数字
OLS	普通最小二乘法
PDF	概率密度函数
PPF	百分比点函数
QQ-Plot	分位数—分位数图
ROC	受试者操作特征
RVS	随机变数样本
SD	标准差
SE/SEM	(均值的)标准误
SF	生存函数
SQL	结构化查询语言
SS	平方和
Tukey HSD	Tukey 显著差异检验

资源与支持

本书由异步社区出品，社区（<https://www.epubit.com/>）为您提供相关资源和后续服务。

配套资源

本书提供如下资源：

- 书中彩图文件。

要获得以上配套资源，请在异步社区本书页面中单击 **配套资源**，跳转到下载界面，按提示进行操作即可。注意：为保证购书读者的权益，该操作会给出相关提示，要求输入提取码进行验证。

提交勘误

作者和编辑尽最大努力来确保书中内容的准确性，但难免会存在疏漏。欢迎您将发现的问题反馈给我们，帮助我们提升图书的质量。

当您发现错误时，请登录异步社区，按书名搜索，进入本书页面，单击“提交勘误”，输入勘误信息，单击“提交”按钮即可。本书的作者和编辑会对您提交的勘误进行审核，确认并接受后，您将获赠异步社区的 100 积分。积分可用于在异步社区兑换优惠券、样书或奖品。



扫码关注本书

扫描下方二维码，您将会在异步社区微信服务号中看到本书信息及相关的服务提示。



与我们联系

我们的联系邮箱是 contact@epubit.com.cn。

如果您对本书有任何疑问或建议，请您发邮件给我们，并在邮件标题中注明本书书名，以便我们更高效地做出反馈。

如果您有兴趣出版图书、录制教学视频，或者参与图书翻译、技术审校等工作，可以发邮件给我们；有意出版图书的作者也可以到异步社区在线提交投稿（直接访问 www.epubit.com/selfpublish/submission 即可）。

如果您是学校、培训机构或企业，想批量购买本书或异步社区出版的其他图书，也可以发邮件给我们。

如果您在网上发现有针对异步社区出品图书的各种形式的盗版行为，包括对图书全部或部分内容的非授权传播，请您将怀疑有侵权行为的链接发邮件给我们。您的这一举动是对作者权益的保护，也是我们持续为您提供有价值的内容的动力之源。

关于异步社区和异步图书

“异步社区”是人民邮电出版社旗下 IT 专业图书社区，致力于出版精品 IT 技术图书和相关学习产品，为作译者提供优质出版服务。异步社区创办于 2015 年 8 月，提供大量精品 IT 技术图书和电子书，以及高品质技术文章和视频课程。更多详情请访问异步社区官网 <https://www.epubit.com>。

“异步图书”是由异步社区编辑团队策划出版的精品 IT 专业图书的品牌，依托于人民邮电出版社近 30 年的计算机图书出版积累和专业编辑团队，相关图书在封面上印有异步图书的 LOGO。异步图书的出版领域包括软件开发、大数据、AI、测试、前端、网络技术 etc。



异步社区



微信服务号

目 录

第一部分 Python 和统计学

第1章 为什么学习统计学 2

第2章 Python 4

- 2.1 开始 4
 - 2.1.1 惯例 4
 - 2.1.2 发行版和包 5
 - 2.1.3 安装 Python 7
 - 2.1.4 安装 R 和 rpy2 8
 - 2.1.5 个性化 IPython/Jupyter 9
 - 2.1.6 Python 资源 12
 - 2.1.7 第一个 Python 程序 13
- 2.2 Python 数据结构 14
 - 2.2.1 Python 数据类型 14
 - 2.2.2 索引和切片 16
 - 2.2.3 向量和数组 17
- 2.3 IPython/Jupyter：一个交互式的 Python 编程环境 18
 - 2.3.1 Qt 控制台的第一个会话 19
 - 2.3.2 Notebook 和 rpy2 21
 - 2.3.3 IPython 小贴士 23
- 2.4 开发 Python 程序 24
 - 2.4.1 将交互式命令转化为一个 Python 程序 24
 - 2.4.2 函数、模块和包 26
 - 2.4.3 Python 小贴士 30
 - 2.4.4 代码版本控制 31
- 2.5 Pandas：用于统计学的数据结构 31
 - 2.5.1 数据处理 31
 - 2.5.2 分组 (Grouping) 33

2 目 录

- 2.6 Statsmodels : 统计建模的工具 34
- 2.7 Seaborn : 数据可视化 35
- 2.8 一般惯例 36
- 2.9 练习 36

第3章 数据输入 38

- 3.1 从文本文件中输入 38
 - 3.1.1 目视检查 38
 - 3.1.2 读入 ASCII 数据到 Python 中 38
- 3.2 从 MS Excel 中导入 42
- 3.3 从其他格式导入数据 43

第4章 统计数据的展示 45

- 4.1 数据类型 45
 - 4.1.1 分类数据 45
 - 4.1.2 数值型 46
- 4.2 在 Python 中作图 46
 - 4.2.1 函数式和面向对象式的绘图方法 47
 - 4.2.2 交互式绘图 48
- 4.3 展示统计学数据集 52
 - 4.3.1 单变量数据 53
 - 4.3.2 二元变量和多元变量绘图 59
- 4.4 练习 61

第二部分 分布和假设检验

第5章 背景 63

- 5.1 总体和样本 63
- 5.2 概率分布 64
 - 5.2.1 离散分布 64
 - 5.2.2 连续分布 65
 - 5.2.3 期望值和方差 65
- 5.3 自由度 66
- 5.4 研究设计 66
 - 5.4.1 术语 67
 - 5.4.2 概述 67

- 5.4.3 研究类型 68
- 5.4.4 实验设计 69
- 5.4.5 个人建议 72
- 5.4.6 临床研究计划 73

第6章 单变量的分布 74

- 6.1 分布的特征描述 74
 - 6.1.1 分布中心 74
 - 6.1.2 量化变异度 76
 - 6.1.3 分布形状的参数描述 79
 - 6.1.4 概率密度的重要展示 81
- 6.2 离散分布 82
 - 6.2.1 伯努利分布 82
 - 6.2.2 二项分布 83
 - 6.2.3 泊松分布 85
- 6.3 正态分布 86
 - 6.3.1 正态分布的例子 88
 - 6.3.2 中心极限定理 88
 - 6.3.3 分布和假设检验 89
- 6.4 来自正态分布的连续型分布 90
 - 6.4.1 t 分布 90
 - 6.4.2 卡方分布 92
 - 6.4.3 F 分布 94
- 6.5 其他连续型分布 95
 - 6.5.1 对数正态分布 96
 - 6.5.2 韦伯分布 96
 - 6.5.3 指数分布 97
 - 6.5.4 均匀分布 98
- 6.6 练习 98

第7章 假设检验 100

- 7.1 典型分析步骤 100
 - 7.1.1 数据筛选和离群值 100
 - 7.1.2 正态性检验 101
 - 7.1.3 转换 104

- 7.2 假设概念、错误、 p 值和样本量 104
 - 7.2.1 一个例子 104
 - 7.2.2 推广和应用 105
 - 7.2.3 p 值的解释 106
 - 7.2.4 错误的类型 107
 - 7.2.5 样本量 108
- 7.3 灵敏度和特异度 110
- 7.4 受试者操作特征(ROC)曲线 113

第8章 数值型数据的均值检验 114

- 8.1 样本均值的分布 114
 - 8.1.1 单样本均值的 t 检验 114
 - 8.1.2 Wilcoxon符号秩和检验 116
- 8.2 两组之间的比较 117
 - 8.2.1 配对 t 检验 117
 - 8.2.2 独立组别之间的 t 检验 118
 - 8.2.3 两组之间的非参数比较:
Mann-Whitney检验 118
 - 8.2.4 统计学假设检验与统计学建模 118
- 8.3 多组比较 120
 - 8.3.1 方差分析(ANOVA) 120
 - 8.3.2 多重比较 123
 - 8.3.3 Kruskal-Wallis检验 125
 - 8.3.4 两因素方差分析 126
 - 8.3.5 三因素方差分析 126
- 8.4 总结:选择正确的检验方法进行组间比较 127
 - 8.4.1 典型的检验 127
 - 8.4.2 假设的例子 128
- 8.5 练习 129

第9章 分类数据的检验 131

- 9.1 单个率 131
 - 9.1.1 置信区间 131
 - 9.1.2 解释 132

- 9.1.3 例子 132
- 9.2 频数表 133
 - 9.2.1 单因素卡方检验 133
 - 9.2.2 卡方列联表检验 134
 - 9.2.3 Fisher 精确检验 136
 - 9.2.4 McNemar 检验 139
 - 9.2.5 Cochran's Q 检验 140
- 9.3 练习 141
- 第 10 章 生存时间分析 144
 - 10.1 生存分布 144
 - 10.2 生存概率 145
 - 10.2.1 删失 145
 - 10.2.2 Kaplan - Meier 生存曲线 146
 - 10.3 在两组间比较生存曲线 148

第三部分 统计建模

- 第 11 章 线性回归模型 150
 - 11.1 线性相关 150
 - 11.1.1 相关系数 150
 - 11.1.2 秩相关 151
 - 11.2 一般线性回归模型 152
 - 11.2.1 例子 1: 简单线性回归 153
 - 11.2.2 例子 2: 二次方拟合 153
 - 11.2.3 决定系数 154
 - 11.3 Patsy: 公式的语言 155
 - 11.4 用 Python 进行线性回归分析 158
 - 11.4.1 例子 1: 拟合带置信区间的直线 158
 - 11.4.2 例子 2: 嘈杂的二次多项式 159
 - 11.5 线性回归模型的结果 162
 - 11.5.1 例子: 英国的烟草和酒精 162
 - 11.5.2 带有截距的回归的定义 165
 - 11.5.3 R^2 值 165

- 11.5.4 \bar{R}^2 : 调整后的 R^2 值 165
- 11.5.5 模型的系数和它们的解释 168
- 11.5.6 残差分析 171
- 11.5.7 异常值 174
- 11.5.8 用 Sklearn 进行回归 175
- 11.5.9 结论 176

- 11.6 线性回归模型的假设 177
- 11.7 线性回归模型结果的解释 180
- 11.8 Bootstrapping 180
- 11.9 练习 181

第 12 章 多元数据分析 182

- 12.1 可视化多元相关 182
 - 12.1.1 散点图矩阵 182
 - 12.1.2 相关性矩阵 182
- 12.2 多重线性回归 184

第 13 章 离散数据的检验 185

- 13.1 等级资料的组间比较 185
- 13.2 Logistic 回归 186
- 13.3 广义线性模型 188
 - 13.3.1 指数族分布 189
 - 13.3.2 线性预测器和连接函数 189
- 13.4 有序 Logistic 回归 189
 - 13.4.1 问题定义 189
 - 13.4.2 优化 191
 - 13.4.3 代码 191
 - 13.4.4 性能 191

第 14 章 贝叶斯统计学 193

- 14.1 贝叶斯学派与频率学派的解释 193
- 14.2 计算机时代的贝叶斯方法 195
- 14.3 例子: 用马尔可夫链蒙特卡罗模拟分析挑战者号灾难 195