



“十三五”全国统计规划教材

(第2版)

Python 数据分析基础

阮 敬 编著



“十三五”全国统计规划教材

内容简介

(第2版)

Python

数据分析基础

阮 敬 编著

内容简介

本书通过真实案例,全面介绍 python3 编程基础及其数据分析工具的应用,培养读者通过数据提出问题、分析问题、解决问题以及对分析结果评价的能力。全书内容包括:python3 基本配置和编程基础、面向对象编程及并行处理、数据预处理、数据描述与可视化、统计推断、相关分析、关联分析、回归分析、主成分和因子分析、聚类、判别与分类、列联分析、对应分析、定性数据分析、时间序列分析、神经网络与深度学习等,将数据分析工作中的基本理论、方法和应用进行深入剖析。

图书在版编目(CIP)数据

Python 数据分析基础 / 阮敬编著. —— 2 版. —— 北京：
中国统计出版社，2018.8

“十三五”全国统计规划教材

ISBN 978-7-5037-8614-3

I. ①P… II. ①阮… III. ①软件工具—程序设计—
高等学校—教材 IV. ①TP311.561

中国版本图书馆 CIP 数据核字(2018)第 190456 号

Python 数据分析基础(第 2 版)

作 者/阮 敬

责任编辑/姜 洋

封面设计/张 冰

出版发行/中国统计出版社

通信地址/北京市丰台区西三环南路甲 6 号 邮政编码/100073

电 话/邮购(010)63376909 书店(010)68783171

网 址/<http://www.zgtjcb.com>

印 刷/河北鑫兆源印刷有限公司

经 销/新华书店

开 本/787×1092mm 1/16

字 数/550 千字

印 张/31

版 别/2018 年 8 月第 2 版

版 次/2018 年 8 月第 1 次印刷

定 价/128.00 元

版权所有。未经许可,本书的任何部分不得以任何方式在世界任何地区
以任何文字翻印、仿制或转载。

中国统计出版社,如有印装错误,本社发行部负责调换。

国家统计局

全国统计教材编审委员会第七届委员会

主任委员：宋跃征

副主任委员：叶植材 许亦频 赵彦云 邱东 徐勇勇 肖红叶
耿直

常务委员（按姓氏笔划排序）：

万东华 叶植材 许亦频 李金昌 杨映霜 肖红叶
邱东 宋跃征 陈峰 周勇 赵彦云 耿直
徐勇勇 徐辉 郭建华 程维虎 曾五一

学术委员（按姓氏笔划排序）：

方积乾 冯士雍 刘扬 杨灿 肖红叶 吴喜之
何书元 汪荣明 金勇进 郑京平 赵彦云 柯惠新
贺铿 耿直 徐一帆 徐勇勇 蒋萍 曾五一

专业委员（按姓氏笔划排序）：

万崇华 马骏 王汉生 王兆军 王志电 王彤
王学钦 王振龙 王震 尹建鑫 石玉峰 石磊
史代敏 冯兴东 朱启贵 朱建平 朱胜 向书坚
刘玉秀 刘立丰 刘立新 米子川 苏为华 杜金柱
李元 李金昌 李勇 李晓松 李萍 李朝鲜
杨仲山 杨军 杨汭华 杨贵军 吴海山 吴德胜
余华银 宋旭光 张波 张宝学 陈峰 林华珍
罗良清 周勇 郑明 房祥忠 郝元涛 胡太忠
洪永森 夏结来 徐国祥 郭建华 唐年胜 程维虎
傅德印 虞文武 薛付忠

出版说明

全国统计教材编审委员会成立于1988年，是国家统计局领导下的全国统计教材建设工作的最高指导机构和咨询机构。自编审委员会成立以来，分别制定并实施了“七五”至“十三五”全国统计教材建设规划，共组织编写和出版了“七五”至“十二五”六轮“全国统计教材编审委员会规划教材”，这些规划教材被全国各院校师生广泛使用，对中国的统计和教育事业作出了积极贡献。自本轮规划教材起，“全国统计教材编审委员会规划教材”更名为“全国统计规划教材”，将以全新的面貌和更积极的精神，继续服务全国院校师生。

《国家教育事业发展“十三五”规划》指出，要实行产学研用协同育人，探索通识教育和专业教育相结合的人才培养方式，推动高校针对不同层次、不同类型人才培养的特点，改进专业培养方案，构建科学的课程体系和学习支持体系。强化课程研发、教材编写、教学成果推广，及时将最新科研成果、企业先进技术等转化为教学内容。加快培养能够解决一线实际问题、宽口径的高层次复合型人才。提高应用型、技术技能型和复合型人才培养比重。

《“十三五”时期统计改革发展规划纲要》指出，“十三五”时期，统计改革发展的总体目标是：形成依靠创新驱动、坚持依法治统、更加公开透明的统计工作格局，逐步实现统计调查的科学规范，统计管理的严谨高效，统计服务的普惠优质，统计运作效率、数据质量和服务水平明显提升，建立适应全面建成小康社会要求的现代统计调查体系，保障统计数据真实准确完整及时，为实现统计现代化打下坚实基础。

围绕新时代中国特色社会主义教育事业和统计事业新特点，全国统计教材编审委员会将组织编写和出版适应新时代特色、高质量、高水平的优秀统计规划教材，以培养出应用型、复合型、高素质、创新型的统计人才。

2015年9月，经李克强总理签批，国务院印发了《促进大数据发展行动纲要》系统部署大数据发展工作，我国各项工作进入大数据时代，拉开了统计教育和统计教材建设的大数据新时代。因此，在完成以往传统统计专业规划教材的编写和出版外，本轮规划教材要把编写大数据内容统计规划教材作为重点工作，以培养新一代适应大数据时代需要的统计人才。

为了适应新时代对统计人才的需要，组织编写出版高质量、高水平教材，本轮规划教材在组织编写和出版中，将坚持以下原则：

1. 坚持质量第一的原则。本轮规划教材将从内容编写、装帧设计、排版印刷等各环节把好质量关，组织编写和出版高质量的统计规划教材。
2. 坚持高水平原则。本轮规划教材将在作者选定、选题编写内容确定、编辑加工等环节上严格把关，确保规划教材在专业内容和写作水平等各方面，保证高水平高标准，坚决杜绝在低水平上重复编写。
3. 坚持创新的原则。无论是对以往规划教材进行修订改版，还是组织编写新编教材，本轮规划教材将把统计工作、统计科研、统计教学以及教学方法、方式的新内容融合在教材中，从规划教材的内容和传播方式上，实行创新。
4. 坚持多层次、多样性规划的原则。本轮规划教材将组织编写出版专科类、本科类、研究生和职业教育类等不同层次的统计教材，并可以考虑根据需要组织编写社会培训类教材；对于同一门课程，鼓励教师编写若干不同风格和适应不同专业培养对象的教材。
5. 坚持教材编写与教材研讨并重的原则。本轮规划教材将注重帮助院校师生学习和使用这些教材，使他们对教材中一些重要概念进一步理解，使教材内容的安排与学生的认知规律相符，发挥教材对统计教学的指导作用，进一步加强统计教材研讨工作，对教材进行分课程的研讨，以促进统计教材的向前发展。
6. 坚持创品牌、出精品、育经典的原则。本轮规划教材将继续修订改版已经出版的优秀规划教材，使它们成为精品，乃至经典，与此同时，将有意识地培养优秀的新作者和新内容规划教材，为以后培养新的精品教材打下基础，把“全国统计规划教材”打造成国内具有巨大影响力的统计教材品牌。
7. 坚持向国际优秀统计教材学习和看齐的原则。不论是修订改版教材还是新编教材，本轮规划教材将坚持与国际接轨，积极吸收国内外统计科学的新成果和统计教学改革的新成就，把这些优秀内容融进去。
8. 坚持积极利用新的教学方式和教学科技成果的原则。本轮规划教材将积极利用数据和互联网发展成果，适应院校教学方式、教学方法以及教材编写方式的重大变化，立体发展纸介质和利用数据、互联网传播方式的统计规划教材内容，适应新时代发展需要。

总之，全国统计教材编审委员会将不忘初心，牢记使命，积极组织各院校统计专家学者参与编写和评审本轮规划教材，虚心听取读者的积极建议，努力组织编写出版好本轮规划教材，使本轮规划教材能够在以往的基础上，百尺竿头，更进一步，为我国的统计和教育事业作出更大贡献。

国家统计局

全国统计教材编审委员会

第2版前言

为适应数据科学与大数据技术领域的飞速发展，《Python 数据分析基础》第2版经过将近1年时间的广泛教学实践和市场检验，终于面世并顺利入选“‘十三五’全国统计规划教材”。在保留第1版全部优点和特色的基础上，第2版做了许多优化、改进和创新，具体内容如下：

1. 全书基于 python 3.6.4 对全部内容进行了更新。
2. 将第1版的编程基础部分根据教学难度和教学要求调整为两个章节。即，第1章强调编程基础，第2章强调编程的高级技能，并补充了类特性、异常捕获与容错处理、并行计算等编程的进阶内容。
3. 增加了“神经网络与深度学习”章节。深度学习是当前数据科学、人工智能领域较为热门的研究内容，第2版增加了对神经网络和深度学习基本思想、基本框架以及基本步骤的介绍，以及如何利用 python 提供的 tensorflow 框架工具进行解决实际问题的案例，帮助读者理解深度学习的理论基础和基本算法。
4. 可读和易用性进一步提高。本书第1版在去年9月份正式出版之后，被全国几十所高等院校采纳为基础课、专业课和选修课的教材。经过多次与授课教师和学生的沟通交流及意见反馈，第2版针对教学过程中的突出问题进行了仔细斟酌和调整，力争使得本书内容更加生动、深入浅出和言简意赅。
5. 第2版除了继续提供书中案例数据（请访问中国统计出版社官方网站 www.zgtjcbs.com 下载），同时还提供 python 编程基础和编程进阶章节的课程课件 PPT，请教师读者联系作者索取。

Email: ruanjing@msn.com

读者交流群：



阮 敬

2018年7月22日于洛杉矶

第1版前言

数据分析是科学研究中的重要环节，随着大数据时代的迅猛发展，其越来越受社会和市场的重视，是科学研究、经营管理、预测与决策等过程中必不可少的基础工作。Python 是当今大数据时代下最为流行的编程工具之一，在大数据领域有着十分广泛的应用，可以实现从数据收集和数据管理到数据分析和挖掘的完整过程，其高效的编程和程序执行过程，能够完全胜任日常数据分析工作的需求。

随着数据分析作用的日益凸显，如何对现有数据进行整理、加工、处理和分析，以期得到所谓的结论，作为人们进行决策的依据进而实现数据的价值？如何利用现有数据对将来可能出现的数据结果或结论进行判断或预测？不管是针对企事业单位的管理者或决策者还是从事具体数据分析的工作人员而言，都需要进行合理数据分析流程的规划，区分数据类型，利用适合的数据分析方法，使用方便、快捷、可靠的统计软件作为工具，对特定数据进行分析与预测，从而洞察市场动向，观测人心所在，把握商机，提升竞争力。

而具有深厚数学背景的统计分析和数据分析方法往往成为相关人员继续深入学习的门槛，甚至成为枯燥乏味的代名词，无法体验到数据分析成果带来的成效。本书就是要力求降低学习难度，通过编者积累的大量真实案例和数据，主要以文字阐述替代复杂公式推导，深入浅出剖析数据分析方法的基本原理和步骤，重点在于厘清数据分析的基本思路，合理得到恰当的分析结果。在分析过程中，本书基于 python 2.7，从基础编程入手，主要通过调用 python 基本库和常用工具库的方式，用大量的实例来展示数据分析每一步骤的细节，带领读者走入数据分析的奇妙世界。

本书的第 1 章和第 2 章主要介绍 python 的基本环境、编程基础和数据预处理方面的内容，具体内容包括 python 数据类型及数据结构、语句与控制流、基本库、函数和面向对象编程的基础，以及数据分析最为常用的基本分析工具库 numpy 和 pandas 基础等；

第 3 章和第 4 章主要介绍利用 python 进行描述分析的基本过程和方法，涵盖了各种常用数据分析图形的绘制和解读以及统计量和统计表等具体内容；

第 5、6、7 章主要介绍利用 python 如何进行总体推断。在大数据时代即使数据量再大，但也离不开利用统计思想对总体特征进行推测和判断，这些具体内容包括参数估计、假设检验和非参数分析；

第 8 章主要介绍如何用 python 来分析数据之间的关系，具体涵盖了简单相关分析、非参数相关分析、偏相关分析、点二列相关分析以及数据挖掘中常用的关联分析等内容；

第 9 章和第 10 章主要介绍如何利用 python 来进行回归分析。回归模型可以说是大部分统计分析和数据挖掘方法的基础，本书介绍的具体内容有线性回归、非线性回归、多项式回归、分位数回归、自变量含有定性变量的回归以及因变量含有定性变量的广义线性回归分析；

第 11 章和第 12 章主要就日常数据分析中所使用的多元统计分析方法进行介绍，具体内容包括主成分分析、因子分析、列联分析以及对应分析等；

第 13 章和第 14 章主要介绍在 python 中进行数据挖掘所使用的聚类和分类方法。内容涵盖系统聚类、 k -means 聚类、DBSCAN 聚类、距离判别和线性判别、贝叶斯判别以及数据挖掘中的 k -近邻、决策树、支持向量机和随机森林等分类方法；

第 15 章主要介绍 python 中使用 ARIMA 建模进行时间序列分析的基本方法和思路。

本书以实用为主要目的，因此上述大部分的数据分析过程均会调用现有常用且公认的结果较为合理的工具库（如 numpy、pandas、matplotlib、scipy、statsmodels、scikit-learn 等）。对于本书提及的数据分析方法无法通过调用现成工具库实现的，本书在相应章节中使用 python 编制了相应的函数或类，以供读者在分析实际问题时调用和复用。读者在复用这些函数或类时，也可根据自身需要对它们进行进一步优化。

全书采用 macOS Sierra 操作系统下的 python 2.7.13 和 Anaconda 4.3.1 的 jupyter notebook 作为分析环境，希望读者参考本书的内容边做边学习。为了提高学习效果，读者应该自行把本书全部代码在 python 中一字一句的敲一遍并运行之，故本书不提供电子版程序代码。但为了提高学习效率，本书附送随书案例的全部数据（下载地址：www.zgtjcbs.com）。

本书由本人在原书《实用 SAS 统计分析教程》（中国统计出版社 2013 年版）基础上亲自编写完成。开源软件的显著特点大家都懂的。因此，读者可在阅读本书时对照原书进行实际操作，认真体会商业软件和开源软件分析流程和分析结果

的异同。此外，我的研究生杨磊磊和王禹提供了部分分析程序并对全书所编制的程序进行了运行验证。同时感谢中国统计出版社的支持。尽管作者已经投入了大量时间和精力来编写此书，但由于水平有限，如有不足之处，敬请专家与同行批评指正。同时也欢迎广大读者与作者积极联系，共同探讨数据分析方面的心得与体会。

Email: ruanjing@msn.com

读者交流群：



阮 敬

2017年8月23日

目 录

第1章 Python 编程基础	1
1.1 Python 系统配置	1
1.2 Python 基础知识	6
1.2.1 帮助	6
1.2.2 标识符	7
1.2.3 行与缩进	7
1.2.4 变量与对象	8
1.2.5 数字与表达式	10
1.2.6 运算符	11
1.2.7 字符串	12
1.2.7.1 转义字符	12
1.2.7.2 字符串格式化	13
1.2.7.3 字符串的内置方法	14
1.2.8 日期和时间	19
1.3 数据结构与序列	20
1.3.1 列表	21
1.3.1.1 列表索引和切片	21
1.3.1.2 列表操作	22
1.3.1.3 内置列表函数	23
1.3.1.4 列表方法	23
1.3.2 元组	24
1.3.3 字典	25
1.3.4 集合	27
1.3.5 推导式	28
1.4 语句与控制流	29
1.4.1 条件语句	29

• 1 •

1. 4. 2 循环语句	31
1. 4. 2. 1 while 循环	31
1. 4. 2. 2 for 循环	31
1. 4. 2. 3 循环控制	33
1. 5 函数	34
1. 5. 1 函数的参数	35
1. 5. 2 全局变量与局部变量	36
1. 5. 3 匿名函数	37
1. 5. 4 递归和闭包	38
1. 5. 5 柯里化与反柯里化	39
1. 5. 6 常用的内置高阶函数	40
1. 5. 6. 1 filter 函数	40
1. 5. 6. 2 map 函数	40
1. 5. 6. 3 reduce 函数	40
1. 6 迭代器、生成器和装饰器	41
1. 6. 1 迭代器	41
1. 6. 2 生成器	42
1. 6. 3 装饰器	44
第 2 章 Python 编程进阶	47
2. 1 类	47
2. 1. 1 声明类	47
2. 1. 2 方法	49
2. 1. 2. 1 实例方法	49
2. 1. 2. 2 类方法	50
2. 1. 2. 3 静态方法	51
2. 1. 3 属性	52
2. 1. 3. 1 实例属性和类属性	53
2. 1. 3. 2 私有属性和公有属性	53
2. 1. 4 继承	54
2. 1. 4. 1 隐式继承	54

2.1.4.2 显式覆盖	56
2.1.4.3 super 继承	56
2.1.4.4 多态	57
2.1.4.5 多重继承	59
2.1.5 特性	60
2.2 异常捕获与容错处理	64
2.2.1 错误和异常	64
2.2.2 异常处理	66
2.2.2.1 触发异常	66
2.2.2.2 捕获异常	67
2.2.2.3 其他处理	68
2.3 模块	69
2.4 包	70
2.4.1 包的组成与调用	71
2.4.2 常用数据分析工具库	71
2.4.2.1 scipy	71
2.4.2.2 statsmodels	72
2.4.2.3 sklearn	73
2.4.2.4 TensorFlow	73
2.5 文件 I/O	74
2.6 多核并行计算	77
2.6.1 多进程	78
2.6.2 并行	81
第3章 数据预处理	84
3.1 numpy 基础	84
3.1.1 向量	86
3.1.2 数组	88
3.1.2.1 数据类型与结构数组	88
3.1.2.2 索引与切片	91
3.1.2.3 数组的属性	94

3.1.2.4	数组排序	95
3.1.2.5	数组维度	96
3.1.2.6	数组组合	98
3.1.2.7	数组分拆	101
3.1.2.8	ufunc 运算	102
3.1.3	矩阵	107
3.1.4	文件读写	107
3.2	pandas 基础	109
3.2.1	pandas 的数据结构	109
3.2.1.1	Series	109
3.2.1.2	DataFrame	113
3.2.2	pandas 的数据操作	123
3.2.2.1	排序	123
3.2.2.2	排名	125
3.2.2.3	运算	126
3.2.2.4	函数应用与映射	127
3.2.2.5	分组	129
3.2.2.6	合并	129
3.2.2.7	分类数据	132
3.2.2.8	时间序列	133
3.2.2.9	缺失值处理	142
第 4 章	数据描述	148
4.1	统计量	148
4.1.1	集中趋势	148
4.1.1.1	均值	148
4.1.1.2	中位数	150
4.1.1.3	分位数	151
4.1.1.4	众数	151
4.1.2	离散程度	152
4.1.2.1	极差	152

4.1.2.2 四分位差	153
4.1.2.3 方差和标准差	153
4.1.2.4 协方差	154
4.1.2.5 变异系数	154
4.1.3 分布形状	154
4.1.3.1 偏度	154
4.1.3.2 峰度	155
4.2 统计表	156
4.2.1 统计表的基本要素	156
4.2.2 统计表的编制	157

第5章 统计图形与可视化 161

5.1 matplotlib 基本绘图	161
5.1.1 函数绘图	161
5.1.2 图形基本设置	166
5.1.2.1 创建图例	166
5.1.2.2 刻度设置	167
5.1.2.3 图像注解	168
5.1.2.4 图像大小	169
5.1.2.5 创建子图	170
5.1.2.6 其他绘图函数	171
5.1.3 面向对象绘图	172
5.1.4 绘图样式	174
5.2 pandas 基本绘图	174
5.3 基本统计图形	176
5.3.1 折线图	177
5.3.2 面积图	179
5.3.3 直方图	179
5.3.4 条形图	181
5.3.5 龙卷风图	184
5.3.6 饼图	185

5.3.7 阶梯图	186
5.3.8 盒须图	187
5.3.9 小提琴图	189
5.3.10 散点图	190
5.3.11 气泡图	192
5.3.12 六边形箱图	193
5.3.13 雷达坐标图	194
5.3.14 轮廓图	195
5.3.15 调和曲线图	195
5.3.16 等高线图	196
5.3.17 极坐标图	196
5.3.18 词云图	197
5.3.19 数据地图	200
5.4 其他绘图工具	202

第6章 简单统计推断	204
6.1 简单统计推断的基本原理	204
6.1.1 数据分布	204
6.1.1.1 总体分布	205
6.1.1.2 样本分布	205
6.1.1.3 抽样分布	205
6.1.2 参数估计	207
6.1.2.1 点估计	208
6.1.2.2 区间估计	208
6.1.3 假设检验	209
6.1.3.1 假设检验的基本思想	209
6.1.3.2 假设检验基本步骤	210
6.1.3.3 假设检验中总体的几种不同情况	211
6.2 单总体参数的估计及假设检验	213
6.2.1 单总体的参数估计	213
6.2.1.1 单总体均值的参数估计	213

6.2.1.2 单总体方差、标准差的参数估计	214
6.2.1.3 单总体比例的参数估计	215
6.2.2 单总体参数的假设检验	215
6.2.2.1 总体均值的假设检验	215
6.2.2.2 总体比例的假设检验	218
6.3 两总体参数的假设检验	218
6.3.1 独立样本的假设检验	219
6.3.1.1 独立样本均值之差的假设检验	219
6.3.1.2 独立样本比例之差的假设检验	221
6.3.2 成对样本的假设检验	222
 第 7 章 方差分析	225
7.1 方差分析的基本原理	225
7.2 一元方差分析	229
7.2.1 一元单因素方差分析	229
7.2.1.1 方差同质性检验	230
7.2.1.2 方差来源分解及检验过程	230
7.2.1.3 多重比较检验	231
7.2.1.4 方差分析模型的参数估计和预测	232
7.2.1.5 方差分析模型的预测	234
7.2.2 一元多因素方差分析	234
7.2.2.1 只考虑主效应的多因素方差分析	235
7.2.2.2 存在交互效应的多因素方差分析	239
7.3 协方差分析	241
 第 8 章 非参数检验	244
8.1 非参数检验的基本问题	244
8.2 单样本非参数检验	245
8.2.1 中位数（均值）的检验	245
8.2.2 分布的检验	247
8.2.3 游程检验	248