

WILEY

大数据应用与技术丛书

Statistical Data Cleaning with Applications in R

R统计数据 清洗及应用

[荷兰] 马克·范德鲁(Mark van der Loo)
埃德温·德荣格(Edwin de Jonge)
杨小冬

著
译

清华大学出版社



大数据应用与技术丛书

R 统计数据清洗 及应用

[荷兰] 马克·范德鲁(Mark van der Loo) 著
埃德温·德荣格(Edwin de Jong) 译
杨小冬



清华大学出版社

北 京

Mark van der Loo, Edwin de Jonge
Statistical Data Cleaning with Applications in R
EISBN: 978-1-118-89715-7

Trademarks: Wiley, the Wiley logo, Wrox, the Wrox logo, Programmer to Programmer, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. Visual Studio is a registered trademark of Microsoft Corporation. All other trademarks are the property of their respective owners. John Wiley & Sons, Inc., is not associated with any product or vendor mentioned in this book.

本书中文简体字版由 Wiley Publishing, Inc. 授权清华大学出版社出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

北京市版权局著作权合同登记号 图字：01-2018-5172

本书封面贴有 Wiley 公司防伪标签，无标签者不得销售。
版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目(CIP)数据

R 统计数据清洗及应用 / (荷)马克·范德鲁(Mark van der Loo), (荷)埃德温·德荣格(Edwin de Jonge) 著; 杨小冬 译. —北京: 清华大学出版社, 2019
(大数据应用与技术丛书)
书名原文: Statistical Data Cleaning with Applications in R
ISBN 978-7-302-52662-9

I. ①R… II. ①马… ②埃… ③杨… III. ①统计分析—统计程序 IV. ①C819

中国版本图书馆 CIP 数据核字(2019)第 053590 号

责任编辑: 王 军
封面设计: 孔祥峰
版式设计: 思创景点
责任校对: 牛艳敏
责任印制: 沈 露

出版发行: 清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址: 北京清华大学学研大厦 A 座 邮 编: 100084

社 总 机: 010-62770175 邮 购: 010-62786544

投稿与读者服务: 010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈: 010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者: 三河市少明印务有限公司

经 销: 全国新华书店

开 本: 170mm×240mm 印 张: 18.5 字 数: 374 千字

版 次: 2019 年 6 月第 1 版 印 次: 2019 年 6 月第 1 次印刷

定 价: 79.80 元

产品编号: 080801-01

译者序

我们生活在一个信息大爆炸的时代，通过各种媒介获取的数据不仅数量巨大，而且纷繁复杂、形式各异。其中既包含我们需要的有用信息，也充斥着很多冗余的、无用的甚至错误的信息。这种情况为数据分析带来了非常大的困难。面对这种困境，数据清洗适时出现。

数据清洗是对数据进行重新审查和校验的过程，目的在于删除重复信息、校正存在的错误、处理无效值和缺失值，并提供数据一致性。顾名思义，数据清洗就是把“脏”的东西“洗掉”，它是发现并校正数据文件中可识别的错误的最后一道程序。其原理是利用数理统计、数据挖掘或预定义的清洗规则等技术将脏数据转换为满足相应质量要求的数据。

本书既包含丰富的理论知识，同时又列举了大量示例，让我们可以轻松地了解如何在实践中应用相关知识。当然，书中涉及很多专业性很强的概念和术语，要求读者对数学和统计学有基本的了解，同时还应该具备一定的编程经验。其中期望值、方差、基础微积分和线性代数方面的知识是必不可少的。此外，如果具备一定的 R 语言知识，对理解本书中的内容会有很大的帮助。除了理论知识介绍和示例之外，书中还提供了一些练习，可以进一步加强读者对软件和相应方法的理解和运用。如果正好面临数据分析的难题，相信本书一定能够为你提供必要的帮助，助你顺利完成工作。

本书的两位作者均来自荷兰统计局统计方法部，在数据清洗、统计和分析方面具备丰富的理论和实践经验，同时还融入了很多相关领域专家的真知灼见，也得到了 R 核心团队、软件包开发人员的大力帮助以及 R 社区的大力支持，这使得本书在所述领域具备很高的权威性。学完本书，一定会让你受益匪浅。

在这里我要对清华大学出版社的李阳编辑及其他相关人员表示感谢，感谢你们的辛勤付出，没有你们就没有本书的顺利出版。

在本书的翻译过程中，译者本着“信、达、雅”的原则，力求能够准确表达出作者的原意，让读者能够轻松理解所讲述的内容。当然，由于译者本身的水平有限，书中难免会出现一些错误，欢迎读者不吝指正，在此向你表示感谢。本书内容主要由杨小冬翻译，参与翻译工作的还有彭明珍、王亚、王浩、程建福、彭林、王欢、彭征肖、武彦龙、杜娟、汪春燕、白德、王亚卫。

最后，要对我的家人表示深深的感谢，谢谢你们的理解和支持，没有你们做坚强的后盾，本书就无法按时翻译完成。

前 言

在数据分析中，数据清洗往往是最为耗时的部分。在“官方统计”(Official Statistics)社区，很久以前就已经开始将数据清洗作为一个单独的学科进行研究(在该社区，数据清洗被称为“数据编辑”)。此外，在研究中还引入了数据库的相关知识，尽管如此，针对大型统计社区的文献著作还是非常有限。正是因为这个原因，当出版方邀请我们对之前为 *useR!*2013 大会编写的《R 数据清洗简介》教程进行扩展，进而编纂成一本书时，我们毫不犹豫就答应了，这也是我们的心声。一方面，我们认为，过去 50 年中在“官方统计”社区发布的一些方法应该为更多的用户所了解和使用，而本书或许可以为此助一臂之力。另一方面，我们正在从基于调查的数据源过渡到管理型“大”数据源，希望本书能够帮助为“官方统计”社区增加一些(通常是预先存在的)相关技术。

对于我们来说，通过编写本书也可以帮助我们系统地梳理相关知识，进一步完善之前针对这一主题所编写的软件。回过头来看，我们最终不仅成功完成了本书的编写，还重新开发并普及了很多之前编写的数据库清洗 R 软件包。为什么要这样做呢？其中一个原因是，我们发现了一些很好的方法，能够普及和扩展我们的软件和方法；另一个原因就是，我们希望将最近出现的“tidyverse”接口风格与 R 功能联系起来。

本书包含的内容

本书包含一系列精选的主题，我们认为这些主题对于开发数据库清洗(也称为数据编辑)系统非常有用。主题范围非常广泛，与计算机科学、数字方法、技术标准、统计以及数据库建模和编程等相关的主题，全部涵盖其中。

本书涵盖“技术数据库清洗”方面的主题，包括数字、文本和日期类型的转换和解释。同时对与这些数据类型相关的技术标准也做了较为详细的介绍。在对象的数据内容方面，相关主题包括数据验证(数据检查)、错误定位、各种错误校正方法以及缺失值插补方法。

对于本书中讨论的理论知识，为了便于用户理解，我们会尽可能地提供可执行的 R 代码进行举例说明。此外，我们还提供了相应的练习，希望可以指导读者进一步加

强对软件 and 对应方法的理解。

广泛的主题既反映出这一课题涵盖范围之广，同时也体现了作者广博的专业知识。当然，还有很多主题并未在书中进行介绍，其中，最重要的主题可能要算清洗时间序列对象和离群值检测。

本书面向的读者

本书的读者应该对数学和统计学有基本的了解，同时还应该具备一定的编程经验。我们假定读者已经了解期望值、方差、基础微积分和线性代数方面的知识。如果具备一定的 R 语言知识，那么对理解本书中的内容会有很大的帮助，因为本书就是使用 R 语言进行说明介绍的。不过，为了便于读者理解和参考，我们还是利用一章内容简要介绍了相关的基础知识。

致谢

本书最终能够顺利出版，离不开很多人的辛勤工作。在这里，我们要感谢荷兰统计局的同事们，他们在百忙之中抽出时间与我们就数据验证、插补和错误定位进行了卓有成效的讨论，为我们提供了很多真知灼见。本书中的部分章节参考了合著者所发表的论文和报告。我们要感谢 Jeroen Pannekoek、Sander Scholtus 和 Jacco Daalmans 的帮助，没有你们的密切合作，就没有本书的成功出版。此外，R 核心团队、软件包开发人员也为我们提供了非常大的帮助，当然，还有 R 社区的大力支持，在此，一并表示感谢。

最后，还要感谢我们的家人，感谢他们的关爱与支持。

2017 年 6 月
Mark 和 Edwin

关于本书配套网站

可以通过以下网址访问本书的配套网站：

www.data-cleaning.org

在本书的配套网站中，可以找到很多极具价值的参考资料，从而更好地学习本书的相关内容，其中包括一些补充资料。

目 录

第 1 章 数据清洗	1	2.7 本书中使用的软件包	24
1.1 统计价值链	1	第 3 章 数据的技术表示	27
1.1.1 原始数据	2	3.1 数值数据	28
1.1.2 输入数据	2	3.1.1 整数	28
1.1.3 有效数据	3	3.1.2 R 中的整数	30
1.1.4 统计数据	3	3.1.3 实数	31
1.1.5 输出	3	3.1.4 双精度数	31
1.2 本书使用的表示法和约定	3	3.1.5 机器精度的概念	33
第 2 章 R 语言简介	5	3.1.6 处理浮点数的不良结果	34
2.1 命令行中的 R 语言	5	3.1.7 处理不良结果	35
2.2 向量	7	3.1.8 R 中的数值数据	37
2.2.1 向量计算	9	3.2 文本数据	38
2.2.2 数组和矩阵	10	3.2.1 术语和编码	38
2.3 数据帧	11	3.2.2 Unicode	39
2.3.1 公式-数据接口	12	3.2.3 一些常见的编码方案	40
2.3.2 选择行和列, 布尔运算符	13	3.2.4 R 中的文本数据: character 类的对象	43
2.3.3 使用索引进行选择	13	3.2.5 R 中的编码方案	45
2.3.4 数据帧操纵: dplyr 软件包	15	3.2.6 使用非本地编码方案进行数据的读取和写入	46
2.4 特殊值	16	3.2.7 检测编码方案	48
2.5 在 R 中导入和导出数据	19	3.2.8 排序规则和排序	49
2.5.1 R 中的文件路径	20	3.3 时间和日期	51
2.5.2 软件包提供的格式	20	3.3.1 TAI、UTC 以及 POSIX 从 Epoch 开始的秒数	51
2.5.3 从数据库读取数据	21	3.3.2 时间和日期表示法	52
2.5.4 处理 R 外部的数据	21	3.3.3 R 中的时间和日期存储	54
2.6 函数	22		
2.6.1 使用函数	22		
2.6.2 编写函数	23		

3.3.4	R 中的时间和日期 转换	55	5.4.1	字符串指标	101
3.3.5	闰日、时区和夏令时	57	5.4.2	R 中的字符串指标和 近似文本匹配	110
3.4	区域设置注意事项	58	第 6 章	数据验证	121
第 4 章	数据结构	61	6.1	简介	121
4.1	简介	61	6.2	初识 validate 软件包	122
4.2	表格数据	61	6.2.1	使用 check_that 快速 检查	122
4.2.1	data.frame 对象	62	6.2.2	基本工作流程: validator 和 confront	124
4.2.2	数据库	62	6.2.3	validate 和 DSL 背景 简介	126
4.2.3	dplyr	64	6.3	定义数据验证	127
4.3	矩阵数据	65	6.3.1	数据验证的正式 定义	128
4.4	时间序列	66	6.3.2	验证函数的运算	130
4.5	图表数据	68	6.3.3	验证和缺失值	132
4.6	Web 数据	70	6.3.4	验证函数的结构	133
4.6.1	网页爬取	70	6.3.5	界定 validate 中的验证 规则	134
4.6.2	Web API	70	6.4	数据验证函数的形式 类型	135
4.7	其他数据	73	6.4.1	深入了解测量	135
4.8	整理表格数据	73	6.4.2	验证规则的分类	137
4.8.1	每列变量	75	6.5	使用 validate 软件包验证 数据	139
4.8.2	单个观测值存储在多个 表中	75	6.5.1	控制台和 validator 对象 中的验证规则	139
第 5 章	清洗文本数据	77	6.5.2	在管道中验证	141
5.1	字符规范化	78	6.5.3	抛出错误或警告	141
5.1.1	编码转换和 Unicode 规范化	78	6.5.4	测试线性方程式的 公差	142
5.1.2	字符转换和音译	80	6.5.5	设置和重置选项	143
5.2	使用正则表达式进行模式 匹配	82	6.5.6	从文件导入验证规则/将 验证规则导出到文件	144
5.2.1	基本正则表达式	82			
5.2.2	实用的正则表达式	85			
5.2.3	在 R 中生成正则 表达式	93			
5.3	R 中的常见字符串处理 任务	94			
5.4	近似文本匹配	99			

6.5.7	检查变量类型和元数据	146	8.1.1	完备性	185
6.5.8	检查值范围和代码列表	147	8.1.2	多余的规则和不可行性	186
6.5.9	检查记录中一致性规则	148	8.2	以逻辑语言表述规则	186
6.5.10	检查跨记录验证规则	150	8.3	规则集问题	188
6.5.11	检查函数依赖	151	8.3.1	不可行规则集	188
6.5.12	跨数据集验证	152	8.3.2	固定值	190
6.5.13	宏、变量组、键	153	8.3.3	冗余规则	191
6.5.14	分析输出: validation 对象	154	8.3.4	非松弛子句	191
6.5.15	输出维度和输出选择	156	8.3.5	非约束子句	191
第 7 章	在数据记录中定位错误	159	8.4	检测和简化过程	192
7.1	错误定位	159	8.4.1	混合整数规划	193
7.2	使用 R 进行错误定位	162	8.4.2	检测可行性	193
7.3	以 MIP 问题的形式进行错误定位	164	8.4.3	查找导致不可行的规则	193
7.3.1	错误定位和混合整数规划	165	8.4.4	检测冲突规则	194
7.3.2	线性限制	166	8.4.5	检测部分不可行性	194
7.3.3	分类限制	167	8.4.6	检测固定值	194
7.3.4	混合类型限制	169	8.4.7	检测非松弛子句	195
7.4	数值稳定性问题	171	8.4.8	检测非约束子句	195
7.4.1	解决 MIP 问题	172	8.4.9	检测冗余规则	195
7.4.2	缩放数值记录	174	8.5	结论	196
7.4.3	设置数值阈值	174	第 9 章	基于领域知识模型的方法	197
7.5	实际问题	176	9.1	使用数据修改规则进行校正	197
7.5.1	设置可靠性权重	176	9.1.1	修改函数	198
7.5.2	简化条件验证规则	177	9.1.2	针对数值数据的一类修改函数	202
7.6	结论	181	9.2	使用 dcmofidy 进行基于规则的校正	206
第 8 章	规则集的维护和简化	185	9.2.1	从文件中读取规则	207
8.1	验证规则的质量	185	9.2.2	修改规则语法	208
			9.2.3	缺失值	209
			9.2.4	顺序执行和与顺序无关的执行	209

9.2.5 选项设置管理	210	10.7 插补下的抽样方差	245
9.3 演绎校正	210	10.8 多重插补	246
9.3.1 校正数值数据中的键入 错误	211	10.8.1 基于 EM 算法的多重 插补	249
9.3.2 使用线性限制进行演绎 插补	214	10.8.2 Amelia 软件包	249
第 10 章 插补和调整	221	10.8.3 基于链式方程的多 变量插补	253
10.1 缺失数据	221	10.8.4 使用 mice 软件包进行 插补	254
10.1.1 缺失数据机制	221	10.9 用于估计插补方差的分析 方法	257
10.1.2 使用 R 可视化和测试 缺失数据中的 模式	222	10.10 选择插补方法	257
10.2 基于模型的插补	226	10.11 约束值调整	260
10.3 R 中基于模型的插补	228	10.11.1 形式化描述	260
10.3.1 使用 <code>simputation</code> 指定 插补方法	228	10.11.2 对插补数据的 应用	263
10.3.2 基于线性回归的 插补	229	10.11.3 使用 <code>rspa</code> 软件包调整 插补值	263
10.3.3 M 估计	231	第 11 章 示例：一个小型数据清洗 系统	265
10.3.4 Lasso 回归、岭回归和 弹性网络回归	233	11.1 设置	266
10.3.5 分类和回归树	233	11.1.1 确定性方法	267
10.3.6 随机森林	236	11.1.2 错误定位	268
10.4 使用 R 进行赋值元素 插补	237	11.1.3 插补	269
10.4.1 随机和顺序热卡 插补	238	11.1.4 调整插补数据	271
10.4.2 k 最近邻和预测均值 匹配	239	11.2 监控数据更改	273
10.5 <code>simputation</code> 软件包中的其他 方法	240	11.2.1 数据差异(Daff)	273
10.6 基于 EM 算法的插补	241	11.2.2 汇总单元格更改	275
10.6.1 EM 算法	242	11.2.3 按照验证规则汇总 更改	276
10.6.2 假定多变量正态分布情 况下的 EM 插补	244	11.2.4 使用 <code>lumberjack</code> 自动 跟踪数据更改	278
		11.3 集成和自动化	282
		11.3.1 使用 RScript	282
		11.3.2 <code>docopt</code> 软件包	283
		11.3.3 自动化数据清洗	283

第1章

数据清洗

1.1 统计价值链

之所以要进行数据清洗，目的在于提升数据质量，使其能够可靠地用于生成统计模型或统计报表。创建部分统计输出必须达到的质量要求由下面这个简单的成本效益问题决定：统计输出什么时候适合使用，为使数据达到这一质量水平需要付出多大的努力？

要想准确回答这个问题，一种非常有用的方法就是按照价值链进行数据分析。粗略地说，价值链由能够逐步增加产品价值的一系列活动组成。在过去 20 年中，统计价值链已经成为“官方统计”社区的常用术语，不过，似乎并未对此形成一个通用的定义。大致来讲，统计价值链通过定义一些有意义的中间数据产品构造而成，针对这些中间数据产品，使用一组精选的质量属性对其进行描述¹。构造统计价值链的方法有很多，但对于这些作者来说，图 1.1 中所示的内容已经被证明是非常通用的，并且有助于围绕统计生产过程组织整理各种想法。

图 1.1 中的架构有一个显著的特点，那就是它很自然地将通常分类为“数据清洗”的活动引入统计生产过程。从左侧开始，首先是原始数据。必须对其进行处理，使之满足(足够的)技术标准，从而可以作为一致性检查、数据校正和插补过程的输入。实现这一目标后，数据便可以被认为是有效的，能够(足以)用于生成统计参数。之后，仍要对这些数据进行格式化处理，以准备生成相应的输出。

应该认识到，尽管此架构能够很好地组织数据分析活动，但实际上，此过程基本上不可能是线性的。更常见的过程是清洗数据、创建一些聚合、发现问题，然后返回。统计价值链的目的更多的是简要了解各个活动发生的位置(例如，通过将它们放在单独脚本中)，而不在于规定实际工作流程的线性顺序。实际上，在统计价值链的后续阶段，工作流程会循环多次，直到生成质量足够理想的输出。在接下来的内容中，我们将对每个阶段进行较为详细的讨论。

¹ 最早引入这一定义的似乎是 Willeboordse(2000)。

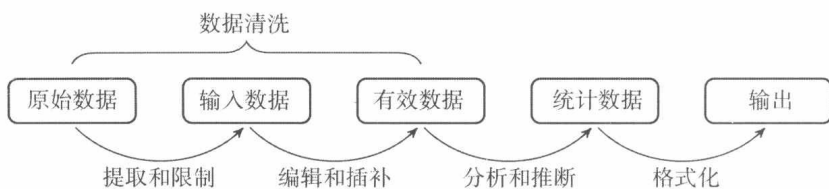


图 1.1 统计价值链的一部分，从原始数据到统计产品，显示 5 个不同的统计价值级别

1.1.1 原始数据

所谓原始数据，指的是最初提供给分析人员的数据。当然，根据数据源的不同，此类数据的状态可能会有非常大的差异。不管是哪种情况，我们都假定一个前提，那就是负责分析的人员对数据的收集方式没有任何影响，或者影响微乎其微，可以忽略。第一步包括使数据易于访问和理解。为了在完成最初的处理步骤后使数据变得更精确，我们要求数据中的每个值都由它们所代表的实际对象(人员、公司等)来标识。对于每个值，已知它所代表的变量(年龄、收入等)，并且值以对应的技术格式(数字和字符串)进行存储。

根据技术上的原始数据格式，实现所需的技术格式所必需的活动通常包括文件转换、字符串规范化(例如编码转换)以及数值的标准化和转换。针对已知统计对象的主干布局注册表连接数据(可能使用能够处理键不精确匹配的过程)也被认为是此类过程的一部分。这些过程将在后面的第 3~5 章中进行更为详细的介绍。

1.1.2 输入数据

输入数据指的是每个值都以正确的类型存储，并以它所代表的变量和所属的统计实体标识的数据。在很多情况下，此类数据集可以通过表格格式来表示，行表示实体，列表示变量。在 R 社区中，这被称为整洁数据(Wickham, 2014b)。在这里，我们将格式置于开放状态。很多数据都可以按照树或图表的形式表示(例如，网页和 XML 结构)。只要所有元素都能轻松地标识并且具有正确的格式，那么此数据就可以作为输入数据。

当数据集达到输入数据的级别以后，必须处理缺失值、非置信值以及非置信值组合。这一过程通常称为数据编辑和插补。它与前面各个步骤的不同之处在于，它重点关注领域知识的数据一致性。此类领域知识通常可以表示为一组规则，例如 `age >= 0`、`mean(profit) > 0` 或 `if (age < 15) has_job = FALSE`。我们将在本书后面的部分(第 6~8 章)详细介绍如何定义、应用和维护此类规则，以便可以实现数据清洗的自动化，从而能够以可重复的方式执行。此外，在第 7 章中，我们将介绍一些方法，使得可以在一条记录中拾取最小数量的字段，通过对这些字段进行更改或插补，能够使所有规则都得到满足。

在后面的第 9 章中，我们将正式介绍如何利用知识规则安全地自动执行数据修改，而第 10 章将介绍如何进行缺失值插补。

1.1.3 有效数据

如果数据能够如实地表示它们所代表的变量和对象，它们就是有效的。如果能够确保数据满足通过一组验证规则的形式表达的领域知识，便可以可重复的方式执行此操作。通常情况下，还会在此基础上补充某些形式的专家审阅。例如，根据各种可视化表示形式进行审阅，或者由领域专家对聚合值进行审阅。

数据被认为有效以后，便可以通过已知的建模和推理技术生成统计数据。这些技术可能需要考虑一些过程，例如，在特定插补过程之后估计方差时，具体取决于前面所采用的数据清洗过程。

1.1.4 统计数据

统计数据只是对应输出变量的估计值。通常情况下，这些数据是简单的聚合值(总计和平均值)。但是，从原则上来说，它们可以由更为复杂的参数(例如，回归模型系数或训练好的机器学习模型)组成。

1.1.5 输出

输出是分析过程的终点。接收统计数据并使其做好传播准备，即可创建输出。这可能涉及技术格式化，例如，使数字可以通过(Web) API 或布局格式化获取(例如，通过准备报告或可视化表示形式)。对于技术格式化，技术验证可能再次成为必不可少的步骤，例如，通过针对某些(JSON 或 XML)架构检查输出格式。通常情况下，一个分析人员的输出往往是另一个分析人员的原始数据。

1.2 本书使用的表示法和约定

本书讨论的主题与数学中的各种子域、逻辑、统计学、计算机科学和编程有关。由于涵盖的领域如此之广，因此很难对各种不同的变量类型和概念使用一种统一的表示法。不过，我们尽量在整本书中使用一种统一的表示法。

基础数学和逻辑

我们遵循传统的表示法，并使用 \mathbb{N} 、 \mathbb{Z} 和 \mathbb{R} 分别表示自然数、整数和实数。符号 \vee 、 \wedge 、 \neg 分别表示逻辑或、逻辑与和逻辑非。在讲到逻辑的情况下，我们使用 \oplus 来表示“异或”。有时候，区分定义和等式是非常有用的。在这种情况下，定义将使用 \equiv 来表示。

线性代数

向量使用粗体的小写字母表示，通常为 \mathbf{x} 、 \mathbf{y} 等。除非另有说明，否则本书所说

的向量一般都是指列向量。符号 $\mathbf{1}$ 和 $\mathbf{0}$ 分别表示系数全部为 1 或 0 的向量。矩阵使用粗体的大写字母表示，通常为 \mathbf{A} 、 \mathbf{B} 等。单位矩阵使用 \mathbf{I} 表示。转置通过上标 T 表示，而矩阵相乘通过并置表示，例如， $\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x}$ 。向量的标准欧几里得范数通过 $\|\mathbf{x}\|$ 表示，如果使用了另一个 L_k 范数，则使用下标表示。例如， $\|\mathbf{x}\|_1$ 表示 \mathbf{x} 的 L_1 范数。在线性代数中， \otimes 和 \oplus 分别表示直积(张量积)和直和。

概率和统计

随机变量使用大写字母表示，通常为 X 、 Y 等。概率和概率密度都使用 $P(X)$ 表示。期望值、方差和协方差分别表示为 $E(X)$ 、 $V(X)$ 和 $\text{cov}(X, Y)$ 。估计值通过 $\hat{\cdot}$ 符号表示，例如， $\hat{E}(X)$ 表示 X 的估计期望值。

代码和变量

R 代码使用固定字体宽度的代码段表示。输出内容以两个注释符号为前缀，如下所示：

```
age <- sample(100,25,replace=TRUE)
mean(age)
## [1] 52.44
```

有时区分代码中的变量与逻辑概念是非常有用的。在这种情况下，代码中的变量将表示为 `age`，而逻辑概念将表示为 *age*。

第2章

R语言简介

接下来，我们将简要介绍 R 语言的一些核心功能。除了安装 R 语言程序以外，我们建议同时安装适用于 R 语言的一种集成开发环境(Integrated Development Environment, IDE)。选择恰当的 IDE 不仅可以为 R 语言编程提供良好的界面，其帮助系统还可以帮助组织项目、代码和数据。

为了充分理解并合理应用本书中的内容，建议亲自运行提供的代码示例，也可以进行一些自定义调整并对生成的结果进行解释。

2.1 命令行中的R语言

在启动 R 程序或适用于 R 语言的 IDE 之后，即可访问一个交互式控制台或命令行界面。第一个应用是替换袖珍计算器。可以键入一个算式，R 程序会返回计算结果(以[1]为前缀)，如下所示：

```
1 + 1
## [1] 2
```

要开始使用 R 语言，可以先试着键入下面的语句。当然，可以使用不同的数字或运算。R 语言支持所有常用的数学函数。

```
1 + 1
3^2
sin(pi/2)
(1 + 4) * 3
exp(1)
sqrt(16)
```

如果想要重复使用特定的结果或值，可以使用<-运算符将它们存储起来。

```
x <- 10
y <- 20
```

R 现在已经记住 10 和 20 这两个值，并将它们命名为 x 和 y 。实际上， x 和 y 现在已经正式成为 R 对象。R 语言非常灵活，除了上面介绍的方法以外，还可以通过多种