

资深数据科学家多维度详解R语言数据可视化技术的精髓
通过500多个可视化示例，全面展现30多类统计图形的绘制

The R logo consists of a red 'R' with a grey circular background element behind it.

R语言

A large, stylized graphic of an eye with radiating lines, rendered in black and yellow, positioned behind the main title.

**数据可视化
实战**

米霖◎编著



机械工业出版社
China Machine Press

R语言

**数据可视化
实战**

米霖◎编著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

R语言数据可视化实战 / 米霖编著. —北京: 机械工业出版社, 2020.11

ISBN 978-7-111-66791-9

I. R… II. 米… III. 统计分析—应用软件 IV. C819

中国版本图书馆CIP数据核字 (2020) 第199849号

R 语言数据可视化实战

出版发行: 机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码: 100037)

责任编辑: 陈佳媛

责任校对: 姚志娟

印刷: 中国电影出版社印刷厂

版次: 2020 年 11 月第 1 版第 1 次印刷

开本: 186mm × 240mm 1/16

印张: 35.75

书号: ISBN 978-7-111-66791-9

定价: 169.00 元

客服电话: (010) 88361066 88379833 68326294

投稿热线: (010) 88379604

华章网站: www.hzbook.com

读者信箱: hzit@hzbook.com

版权所有·侵权必究

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

数据可视化是一种非常重要的技能，能够帮助人们快速理解信息，分析数据中存在的趋势，识别关系中的模式。当人们有了新的见解后，往往需要将这些见解传达给其他人，在传达的过程中，使用图表、图形或其他对视觉有影响的表现形式非常重要，因为这些表现形式吸引人并可以迅速传达信息。

R 是由统计学家设计的专门用于统计计算的语言，它也是一种非常好的数据可视化工具。随着技术的进步，数据公司或研究机构对数据的收集变得越来越复杂，许多人已经采用 R 语言作为分析数据的首选语言。R 语言适用于机器学习、数据分析、数据可视化及科学计算等领域。它有多个主题的软件包，如计量经济学、金融学和时间序列等，也拥有一流的可视化、报告和交互性工具。这些工具对于企业和科学研究都非常重要，被科学家、工程师和商业人士广泛使用。

笔者多年来一直以 R 语言为工具从事数据可视化、数据分析、统计建模和机器学习等数据科学工作，在工作中积累了大量的经验，对 R 语言的原理和应用有比较全面和深入的认识，尤其在数据可视化方面更是有独到的见解。R 语言提供了强大的数据可视化功能，可以生成高质量的图形，如条形图、直方图、散点图、动态图和数学符号，甚至可以用很少的代码来实现自己设计的全新图形。R 语言有很多数据可视化包，如 `ggplot2`、`ggvis` 和 `plotly` 等，使用这些包可以非常轻松地绘制出令人印象深刻的数据可视化图表。其中，`plotly` 包提供了一套绘制交互式图形的工具，所绘制的图形非常有表现力。另外，R 语言的文档资源很丰富，而且质量很高，这对学习 R 语言有很大的帮助。

为了帮助广大相关从业人员和数据技术爱好者快速掌握 R 语言数据可视化技术，笔者编写了本书。本书主要介绍如何使用 R 语言绘制常见的统计图形，如直方图、散点图和箱线图，另外也介绍了如何绘制一些不常见但很实用的统计图形，如桑基图、和弦图和时间序列图等。

本书不但介绍了普通的统计图形的绘制，而且介绍了交互式图形及动画图形的绘制，另外还介绍了如何使用 Shiny 工具包进行数据探索与可视化。相信通过阅读本书，读者可以在较短时间内比较系统地掌握 R 语言数据可视化技术。

本书特色

1. 内容全面，理论结合实践

本书全面介绍如何利用 R 语言绘制各种统计图形，涵盖普通的统计图形绘制、高级统

计图形绘制、交互式图形绘制、动画图形绘制，以及统计图形的细节调整。另外，本书还介绍了如何利用 Shiny 进行数据探索与可视化，是一部理论与实践紧密结合的数据可视化学习宝典。

2. 零基础入门，学习门槛低

阅读本书不需要读者具备太多的预备知识，只要有基本的数学与统计学知识即可，学习门槛很低，零基础即可入门。

3. 实例丰富，实用性强

本书在讲解的过程中给出了大量的 R 语言绘图实例，这些实例涵盖单变量图形绘制、两个同类型变量的图形绘制、分类变量和连续变量的图形绘制、高维图形绘制等。这些例子都非常实用，可以比较容易地迁移到实际工作中。

4. 绘图代码详细，效果精美

本书中的绘图实例都给出了详细的实现代码，读者可以按照代码和操作步骤亲自动手实现每一个实例的效果，而且这些绘图实例的效果非常精美，让人赏心悦目。

本书内容

本书共分为 13 章，各章内容简单介绍如下：

第 1 章主要介绍 R 语言的基本概念、Rstudio 跨平台集成开发环境及常见的统计图形等。

第 2 章主要介绍一些数据处理与数据探索的方法，如数据转换和数据重塑等，因为在数据可视化之前需要先对数据进行处理。

第 3 章主要介绍在不进行更多细节调整的情况下如何快速地进行数据可视化。

第 4 章主要介绍面积图、密度图、直方图和频率图这几种单变量图形的绘制。这类图形往往只涉及数据集的一个变量。

第 5 章主要介绍两个同类型变量的图形绘制，包括散点图、抖动点图、连续二维分布图和线图。

第 6 章主要介绍分类变量和连续变量的图形绘制，包括箱线图、小提琴图、棒棒糖图、条形图、圆形条形图、饼图和甜甜圈图。

第 7 章主要介绍高维图形的绘制，包括气泡图、三维散点图、流型图、相关矩阵图、树状图、圆形包装图和树形图。

第 8 章主要介绍其他类型的统计图形绘制，包括和弦图、桑基图、网络图、旭日图、雷达图、词云、平行图、时间序列图、交互式图形及动画图。这些图形并不是很常见，但是非常有用，使用这些高级图形能够让数据可视化的效果更加引人注目。

第 9~11 章主要介绍图形的细节调整，如添加图形元素、图形的颜色调整、线条类型

调整、坐标轴范围调整、删除面板边框和网格线、合并多幅图形等。一幅优秀的统计图形往往需要许多细节上的调整，通过调整细节，可以让图形更具表现力。

第 12、13 章主要介绍一些扩展内容，包括 `ggfortify` 绘图包和 `Shiny` 工具包，它们可以实现数据可视化的一些高级功能，如交互式图形绘制和动画图形绘制等。

读者对象

- 数据可视化从业人员；
- 统计学、数学、经济学、计算机和财经等专业的本科生和研究生；
- 互联网从业人员，如产品经理；
- R 语言数据可视化初学者与进阶者；
- 对数据可视化感兴趣的人员；
- 相关培训机构的学员。

配套资源

本书的所有实例源代码文件及彩色效果图等相关资源需要读者自行下载。方法是：在华章公司官网 www.hzbook.com 上搜索到本书，然后单击“资料下载”按钮，即可在本书页面上找到“配书资源”下载链接。

售后支持

本书涉及的内容比较庞杂，R 数据可视化技术也是日新月异，加之作者水平和成书时间所限，书中难免有一些疏漏和不当之处，敬请读者指正。阅读过程中若有疑问，请发 E-mail 至 hzbook2017@163.com，以获得帮助。

前言

第 1 章 R 语言数据可视化简介	1
1.1 R 语言介绍	1
1.1.1 向量	3
1.1.2 列表	3
1.1.3 矩阵	4
1.1.4 数组	5
1.1.5 因子	5
1.1.6 数据框	6
1.1.7 for 循环	7
1.1.8 条件判断	8
1.1.9 函数	9
1.2 Rstudio 介绍	12
1.3 R 包介绍	13
1.4 R 语言数据读取	14
1.4.1 读取 Excel 数据	15
1.4.2 读取 SPSS、SAS 和 STATA 数据	17
1.5 ggplot2 介绍	18
1.5.1 使用 qplot 函数快速绘图	19
1.5.2 使用 ggplot 函数绘图	20
1.6 统计图形	22
1.6.1 散点图	22
1.6.2 箱线图	24
1.6.3 小提琴图	25
1.6.4 条形图	27
1.6.5 和弦图	28
1.6.6 桑基图	30
1.6.7 棒棒糖图	31
1.6.8 克利夫兰点图	32
1.6.9 艺术图	34
1.7 tidyverse 介绍	38

1.8 总结	41
第 2 章 数据处理与探索	42
2.1 数据转换	42
2.1.1 筛选数据集的行	42
2.1.2 筛选数据集的列	46
2.1.3 数据排序及新变量生成	48
2.1.4 数据分组汇总	49
2.1.5 数据合并	50
2.2 数据重塑	54
2.2.1 数据聚合	55
2.2.2 数据分散	56
2.2.3 数据切割	57
2.2.4 数据合并	58
2.3 总结	59
第 3 章 数据可视化	60
3.1 ggplot2 核心概念	60
3.1.1 散点图	62
3.1.2 折线图	64
3.1.3 条形图	66
3.1.4 直方图	68
3.1.5 密度图	72
3.1.6 箱线图	75
3.2 总结	79
第 4 章 单变量图形绘制	80
4.1 面积图	80
4.1.1 面积图的绘制方式	82
4.1.2 绘制堆叠的面积图	85
4.1.3 绘制比例堆叠面积图	88
4.2 密度图	90
4.2.1 基础密度图	95
4.2.2 绘制少量分组的密度图	96
4.2.3 绘制大量分组的密度图	101
4.2.4 密度图的其他调整	103
4.3 直方图	105
4.3.1 基础直方图	107
4.3.2 分组直方图的绘制	109

4.3.3	合并直方图与密度图	111
4.4	频率图	113
4.5	总结	114
第 5 章	两个同类型变量的图形绘制	115
5.1	散点图	115
5.1.1	绘制基础散点图	118
5.1.2	绘制分组散点图	121
5.1.3	添加拟合曲线	127
5.1.4	在散点图中添加地毯图	130
5.1.5	在散点图中添加文本	136
5.2	抖动点图	140
5.3	连续二维分布图	142
5.3.1	绘制二维直方图	143
5.3.2	绘制六角直方图	144
5.3.3	绘制二维密度直方图	145
5.3.4	调整图形配色	147
5.4	线图	148
5.4.1	绘制基础线图	155
5.4.2	绘制连线图	160
第 6 章	分类变量和连续变量的图形绘制	163
6.1	箱线图	163
6.1.1	绘制基础箱线图	174
6.1.2	调整参数	175
6.1.3	调整箱线图组别的顺序	176
6.1.4	调整颜色	183
6.1.5	构建分组箱线图	191
6.1.6	调整箱线图的宽度	193
6.1.7	构建连续变量的箱线图	194
6.1.8	添加平均值	195
6.1.9	添加抖动点	196
6.2	小提琴图	197
6.2.1	绘制基础的小提琴图	200
6.2.2	绘制水平的小提琴图	201
6.2.3	在小提琴图中添加箱线图	203
6.3	棒棒糖图	204
6.3.1	绘制分组的棒棒糖图	209

6.3.2	绘制基础棒棒糖图	217
6.3.3	棒棒糖图参数的调节	219
6.3.4	添加标注	224
6.4	条形图	226
6.4.1	绘制基础条形图	229
6.4.2	改变条形图宽度	235
6.4.3	添加误差棒	235
6.5	圆形条形图	241
6.5.1	绘制基础圆形条形图	250
6.5.2	添加标签	252
6.5.3	圆形条形图的更多调整	253
6.6	饼图	259
6.6.1	绘制基础饼图	263
6.6.2	调整细节	264
6.6.3	添加标签	265
6.7	甜甜圈图	266
第 7 章	高维图形绘制	270
7.1	气泡图	270
7.1.1	绘制基础气泡图	274
7.1.2	控制气泡的大小	275
7.1.3	设置颜色	276
7.1.4	调整更多的细节	277
7.1.5	绘制动态图	279
7.2	三维散点图	280
7.3	流型图	282
7.3.1	绘制基础流型图	285
7.3.2	调整流型图的偏移	286
7.3.3	调整流型图的形状与颜色	287
7.4	相关矩阵图	288
7.5	树状图	291
7.5.1	绘制基础树状图	295
7.5.2	绘制圆形树状图	297
7.5.3	绘制聚类结果的树状图	298
7.5.4	更多调整	302
7.6	圆形包装图	308
7.6.1	具有一个层次的圆形包装图	310
7.6.2	调整颜色	311

7.6.3	调整圆形之间的距离	315
7.6.4	绘制多层次的圆形包装图	315
7.6.5	调整细节	317
7.6.6	隐藏第一级	321
7.7	树形图	325
7.7.1	绘制基础树形图	326
7.7.2	绘制带有多个级别的树形图	327
7.7.3	自定义树形图	328
第 8 章	其他图形绘制	332
8.1	和弦图	332
8.1.1	绘制圆形图	334
8.1.2	绘制基础和弦图	337
8.1.3	调整细节	340
8.2	桑基图	343
8.3	网络图	347
8.3.1	绘制基础网络图	356
8.3.2	调整网络图的参数	358
8.3.3	网络图布局	361
8.3.4	将变量映射到节点和链接特征	362
8.3.5	使用网络图可视化聚类结果	364
8.4	旭日图	366
8.5	雷达图	368
8.5.1	绘制雷达图	374
8.5.2	绘制多组雷达图	375
8.6	词云	376
8.6.1	绘制词云	378
8.6.2	调整颜色和背景颜色	379
8.6.3	调整形状	381
8.6.4	调整单词方向	382
8.7	平行图	383
8.7.1	绘制基础平行图	389
8.7.2	自定义颜色、主题和外观	390
8.8	时间序列图	391
8.8.1	时间序列包 dygraphs	396
8.8.2	时间序列热图	397
8.9	交互式图形	399

8.9.1	散点图	400
8.9.2	气泡图	401
8.9.3	面积图	402
8.9.4	条形图	404
8.9.5	饼图	405
8.9.6	桑基图	406
8.9.7	误差棒图	408
8.9.8	箱线图	409
8.9.9	直方图	411
8.9.10	二维直方图	413
8.9.11	二维轮廓直方图	414
8.9.12	小提琴图	415
8.9.13	雷达图	416
8.9.14	热图	418
8.9.15	三维散点图	418
8.9.16	动画图	420
8.9.17	调整图形图例	421
8.9.18	修改交互文本	422
8.10	动画图	423
8.10.1	绘制基础动画图	424
8.10.2	使用分面	425
8.10.3	动态变化图形	426
第 9 章	图形元素、标题和图例绘制	429
9.1	添加图形元素	429
9.2	主标题、轴标签和图例标题	432
9.2.1	改变标签的外观	434
9.2.2	修改图例	436
9.2.3	修改图例的位置和外貌	436
9.2.4	使用 <code>guides</code> 函数修改图例	440
第 10 章	颜色等参数的调整	445
10.1	图形颜色调整	445
10.1.1	使用单个颜色调整图形	446
10.1.2	通过分组调整颜色	448
10.1.3	渐变或连续颜色	455
10.2	点的形状、颜色和大小调整	457
10.3	线条类型调整	460

10.4	坐标轴范围调整	462
10.5	坐标轴转换	465
10.6	时间数据坐标轴	468
10.7	自定义标签	471
10.8	图形主题和背景颜色	477
10.9	自定义图形的背景	480
10.10	删除面板边框和网格线	481
10.11	ggthemes 包	482
10.12	文本注释	483
10.13	ggrepel 包	485
10.14	添加直线	488
10.15	图形翻转和反向	490
10.16	分面	491
第 11 章	合并多幅图形	499
11.1	合并多幅图形到一张图中	499
11.2	gridExtra 包	502
11.3	添加边际分布图	505
11.4	在 ggplot 中插入一个外部图形元素	506
第 12 章	R 语言绘图包	509
12.1	ggstatsplot 包	509
12.2	ggfortify 包	520
12.2.1	生存分析	520
12.2.2	时间序列图	521
12.2.3	密度图	523
12.2.4	时间序列预测图	524
12.2.5	聚类图	527
12.2.6	热力图	530
12.2.7	主成分分析可视化	532
12.3	quantmod 包	535
第 13 章	Shiny 工具包	544
13.1	Shiny 工具包简介	544
13.2	Shiny App 的基础部分	548
13.3	Shiny 示例	550
13.4	Shiny 总结	553
13.5	制作一个 Shiny 程序	554
13.6	Shiny 部署	556

第 1 章 R 语言数据可视化简介

本章主要介绍 R 语言数据可视化的一些基本内容，包括什么是 R、什么是 Rstudio、什么是 R 包等，另外还会介绍 R 语言基本语法的相关内容。由于本书的重点在于数据可视化，因此不会深入讲解 R 语言的语法。本章旨在帮助读者了解 R 语言的一些核心内容，以尽可能简单的方式帮助读者了解如何使用 R 语言，最后介绍一些常见的统计图形。

1.1 R 语言介绍

R 是一种编程语言和自由软件，常用于统计计算和图形绘制。R 语言广泛应用于统计学领域和数据科学领域，用于开发统计软件 and 进行数据分析。R 是世界上最强大的统计计算、机器学习和图形编程语言，拥有蓬勃发展的全球用户、开发人员和贡献者社区。

R 语言是罗斯·伊哈卡（Ross Ihaka）和罗伯特·杰特曼（Robert Gentleman）在新西兰的奥克兰大学所创建的，以两个作者名字的第一个字母命名。R 项目于 1992 年构思，1995 年发布初始版本，2000 年发布稳定版本。R 可从 CRAN 网页进行下载，链接为 <http://cran.r-project.org/>。

R 及其库实现了各种统计和图形方法，包括线性和非线性建模、经典统计测试、时间序列分析、分类、聚类等。一般而言，最新的统计方法都有对应的 R 包，R 用户可以非常方便地学习、应用最新的统计方法。

R 可以通过函数和扩展包轻松扩展其功能，R 社区因在软件包方面的积极贡献而闻名。任何人都可以在 R 社区贡献出自己的包。到目前为止，R 中有超过 10 000 个包可供下载。R 开发人员社区非常活跃。R 在统计计算、数据科学及数据可视化方面有着无与伦比的优势，几乎你面临的所有问题或者你关心的问题都有人实现了 R 包。

R 是数据科学领域最流行的语言，并且 R 语言是完全面向数据的，更加注重从数据的角度去思考问题，在这一点上与其他的编程语言有很大的区别。

另外，R 是用于统计研究的主要工具。因此当新的方法被开发的时候，它不仅仅被作为论文发表，而且往往会被开发成为一个 R 包。这让 R 永远成为新算法的前沿，而 R 用户可以非常方便地使用这些新的算法。

CRAN 本身是一个非常有效的共享 R 扩展平台，具有用于包创作、构建、测试和分发的成熟系统。R 核心团队，特别是 CRAN 维护者，为 R 包创建了这样一个充满活力的生

态系统。图 1.1 所示为 R 包数量在不同年份的变化。



图 1.1 R 包数量的变化

从图 1.1 中可以看到，R 包的开发速度越来越快。但是过多的包导致人们在寻找自己所需要的包时困难重重，可能需要借助一些搜索工具来完成。目前有很多包的搜索与汇总工具，例如：

- CRAN 提供包任务视图 (<https://cran.r-project.org/web/views/>)，按主题区域（例如财务或临床试验）提供包目录。
- MRAN (Microsoft R 应用程序网络) 为 CRAN 上的 R 软件包提供搜索工具。MRAN 的链接为 <https://mran.microsoft.com/taskview>。
- 为了找到最受欢迎的软件包，Rdocumentation.org 按下载次数提供了软件包的排行榜，并提供了新发布和最近更新的包的列表。RDocumentation.org (<https://www.rdocumentation.org/taskviews#Bayesian>) 还提供了基于 CRAN 任务视图的可搜索版本。通过上面的几个方法可以很方便地找到自己所需要的包。

R 有许多标准函数，使用这些函数可以解决非常多的统计机器学习任务。对于计算密集型任务，可以在运行时连接和调用 C、C++ 和 Fortran 代码。高级用户可以编写 C、C++、Java、.NET 或 Python 代码直接操作 R 对象，还可以通过使用用户提交的包来执行特定功能或特定研究领域的任务。软件包使 R 具有高度可扩展性。

R 的语法非常简单，介绍 R 的语法首先需要从数据的角度入手，在使用任何编程语言进行编程时，需要使用各种变量来存储各种信息。与其他编程语言（如 C 和 Java）不同，在 R 中使用变量不需要声明类型，并且 R 中的数据由 R 语言的对象存储。这些对象包括

向量、列表、矩阵、数组、因子、数据框。

数据是数据科学的基础，也是数据可视化的基础。下面开始介绍 R 的基本数据结构。

1.1.1 向量

这里所说的向量并不是物理学中有方向的一个量，而是一组数据。例如，1~10，这10个数字就可以表示为一个向量。或者26个字母，也可以表示为一个向量。向量是 R 中的一个基本数据结构。当你想用多个元素创建向量时，应该使用 `c` 函数，这意味着将元素组合成一个向量。下面的代码则是创建向量的一个例子。

```
# 创建一个向量
letter <- c('a','b',"c")
letter
## [1] "a" "b" "c"
# 查看数据类型
class(letter)
## [1] "character"
```

上面的代码使用 `c` 函数创建了一个向量，包含3个元素，分别为 a、b、c 这3个字母。然后 `class` 函数可以用于查看数据的类型。从代码的输出结果中可以看出，这个向量中元素的数据类型是字符类型 (`character`)。上面创建的向量有3个元素，如果希望提取向量中的某个元素，则可以通过 `[]` 进行提取，代码如下：

```
letter[1]
## a
```

上面的代码选取了向量中的第一个元素，希望提取第几个元素只需要传入对应数字即可。

1.1.2 列表

列表 (`list`) 这个数据类型常用于存储不同数据类型的数据。列表里面可以存储向量、列表、矩阵、数组、因子、数据框等所有数据类型。另外，列表还可以嵌套多层列表，下面的代码创建了一个列表，列表中分别包含了向量、列表、数组、矩阵、数据框这几种 R 语言中的数据结构。

```
l <- list(c(1,2,3),list(1),array(c(1,2,3,3,1,2),dim = c(1,2,3)),matrix
(c(1,2,3,4),nrow = 2),factor(c(1,2,3)),data.frame(nu = c(1:3),id =c("a",
"b","c")))
l
## [[1]]
## [1] 1 2 3
##
## [[2]]
## [[2]][[1]]
## [1] 1
```

```
##
##
## [[3]]
## , , 1
##
##      [,1] [,2]
## [1,]  1   2
##
## , , 2
##
##      [,1] [,2]
## [1,]  3   3
##
## , , 3
##
##      [,1] [,2]
## [1,]  1   2
##
##
## [[4]]
##      [,1] [,2]
## [1,]  1   3
## [2,]  2   4
##
## [[5]]
## [1] 1 2 3
## Levels: 1 2 3
##
## [[6]]
##   nu id
## 1 1 a
## 2 2 b
## 3 3 c
class(l)
## [1] "list"
```

上面的代码创建了 6 种类型的数据，然后将它们全部放到一个 list 里面。如果希望提取列表中的元素，则需要使用[[]]这个符号。l[[1]]表示提取列表的第一个元素，也就是整个向量。如果希望进一步提取向量中的元素，提取方式与普通向量的提取方式相同。例如，代码 l[[1]][1]提取了列表中向量的第一个元素。

1.1.3 矩阵

矩阵可以理解为二维数组，它可以使用矩阵函数（matrix）来创建。执行下面的代码，将创建一个矩阵。

```
# Create a matrix.
M = matrix( c('a','a','b','c','b','a'), nrow = 2, ncol = 3, byrow = TRUE)
print(M)
##      [,1] [,2] [,3]
## [1,] "a"  "a"  "b"
## [2,] "c"  "b"  "a"
```