



R Data Analysis:
Methods and
Application

R 数据分析 方法与案例详解

方匡南 朱建平 姜叶飞 编著



R 数据分析 方法与案例详解

方匡南 朱建平 姜叶飞 编著

电子工业出版社

Publishing House of Electronics Industry

北京·BEIJING

内 容 简 介

本书是一本 R 语言和数据分析的入门教材,循序渐进、深入浅出,每个知识点尽量从实际的应用案例出发,以问题为导向,在解决问题中学习统计方法、R 语言的基本使用以及编程技巧。

本书内容涵盖 R 数据结构、函数与优化、抽样模拟、统计分析、假设检验、回归分析、统计绘图和 R 包制作等内容。

本书的定位是为业界数据分析人员、经济管理类、医学的学生提供方法和程序上的参考,在写作过程中尽量删去比较理论的数学原理,这样能够帮助读者轻松上手学习。

未经许可,不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有,侵权必究。

图书在版编目(CIP)数据

R 数据分析:方法与案例详解 / 方匡南,朱建平,姜叶飞编著. —北京:电子工业出版社,2015.2
ISBN 978-7-121-25290-7

I. ①R… II. ①方… ②朱… ③姜… III. ①程序设计—教材 IV. ①TP312

中国版本图书馆 CIP 数据核字(2015)第 30474 号



责任编辑:张月萍

印 刷:三河市双峰印刷装订有限公司

装 订:三河市双峰印刷装订有限公司

出版发行:电子工业出版社

北京市海淀区万寿路173信箱

邮编:100036

开 本:787×980 1/16 印张:24.5

字数:583千字

版 次:2015年2月第1版

印 次:2015年2月第1次印刷

定 价:69.00元

凡所购买电子工业出版社图书有缺损问题,请向购买书店调换。若书店售缺,请与本社发行部联系,联系及邮购电话:(010)88254888。

质量投诉请发邮件至zlt@phei.com.cn,盗版侵权举报请发邮件至dbqq@phei.com.cn。

服务热线:(010)88258888。

前言

R语言是由新西兰奥克兰大学的Ross Ihaka与Robert Gentleman一起开发的一个面向对象的编程语言，因两人的名字都是以R开头，所以命名为“R”。R是“GNU S”，一个免费开源、能够自由有效地用于统计计算和绘图的语言和环境，可以在UNIX、Windows和Mac OS系统中运行，它提供了广泛的统计分析和绘图技术，包括回归分析、时间序列、分类和聚类等方法。2009年，《纽约时报》发表了题为“Data Analysts Captivated by R's Power”的社评，集中讨论了R语言在数据分析领域的发展，并引发了SAS和R用户广泛而激烈的争论。文章认为，让R变得如此有用和如此快地广受欢迎的原因是统计学家、工程师和科学家们在不断精炼代码或编写各种特有而具体的包。而且现在R软件增添了很多高级算法、作图颜色和文本注释，以及为与数据库链接等提供挖掘技术。文中引用了几位科学家对R做的高度评价，如Google首席经济学家Hal Varian讲了一句很好的话：“R最优美的地方是它能够修改很多前人编写的包的代码做各种你所需要的事情，实际上你是站在巨人的肩膀上。”

2010年，美国统计协会（American Statistical Association）将第一届“统计计算及图形奖”授予R语言，用于表彰其在统计应用和统计研究方面广泛的影响。

笔者自接触R语言至今已近10年，记得最早听到R语言是2004年日本同志社文化大学情报研究所金明哲教授前来讲座，他给我们讲了一个关于利用文本挖掘方法破了人身意外保险诈骗杀人案。第一次听到数据挖掘竟然可以如此神奇的破案，便被深深地吸引住了，当初金教授提到其主要分析工具就是R语言。这是笔者第一次接触R语言。

其后，笔者的硕士导师王斌会教授要求我们学习R语言，坦诚地讲，R语言对于刚入门的人来讲确实有点困难，而且那个时候R语言的参考书少得可怜，当时市面上还没有R语言中文教材，用的人也很少。我参与了王斌会老师R语言中文教材的编写小组，一边学习R语言英文经典教材，一边整理学习心得和相关资料。一年之后，第一本R语言中文教材《R语言统计分析软件教程》面世了。通过编写教材，笔者从中体会到了学习R语言的乐趣，自从学会了R语言，笔者之后的数据分析主要使用的都是它。

后来，人大经济论坛邀请笔者为他们录制R语言视频课程，笔者利用闲暇时间整理讲义，先后录制了《R语言初级》《R语言高级》《R金融时间序列初级》和《R金融时间序列高级》等视频课程，反响很好，很多学员要求笔者写相应的系列教材，由于繁重的科研压力以及个人

的惰性，虽有此想法已久，但却迟迟未动笔。后来成为厦门大学教师，系里让教《计算机在统计中的应用》一课，笔者觉得讲R语言是最恰当不过了，因此边上课，边整理讲义，并不断地完善，教该课程已经4年，讲义也先后修改了4年。

再后来，人大经济论坛又邀请笔者为他们在北京、上海等地开设暑期和节假日现场公开课，前来听课的学生既有国内外著名高校的教师和研究生，又有医药公司的数据分析人员以及互联网公司的数据分析人员等，他们对笔者的讲义提了很多有用的建议。笔者在这些公开课的讲课中，不断地完善讲义，最终形成了此书。关于如何取一个恰当而又响亮的书名，确实是一件挺难的事，笔者想了好几个书名，但最终还是认为取名为《R数据分析——方法与案例详解》可能比较合适。

该书是一本R语言和数据分析的入门教材，内容循序渐进、深入浅出，每个知识点都尽量从实际的应用案例出发，以问题为导向，在解决问题中学习统计方法、R语言的基本使用以及编程技巧。本书的定位是为业界数据分析人员、经济管理类和医学类的学生提供方法和程序上的参考，在写作过程中尽量删去比较理论的数学原理，当然有些原理部分无法跳过，所以在学习时需要有一定高等数学和概率论的基础，但如果想真正掌握某个统计方法，那么学习其背后的原理还是非常有必要的。当然，如果对方法的原理确实不感兴趣，只是为了用R程序实现某种方法，或者分析某个有意义的数据，则可以跳过原理，只看案例和程序即可。

该书的姐妹书，会讲解更为高深的统计方法，涉及聚类分析、分类分析、关联规则、时间序列等问题，更着重在互联网、金融、企业营销、基因分析等领域的应用，目前暂定书名为《R数据挖掘——方法与案例详解》。

该书适合作为高校数据分析相关专业的教科书，也适合供医学、市场调查、金融以及互联网等企业的相关数据分析人员阅读。

本书第1章和第2章由方匡南、姜叶飞撰写，第3章和第4章由方匡南撰写，第5章由方匡南、易焯迪撰写，第6章到第9章由方匡南撰写，第10章到第15章由朱建平、方匡南撰写，第16章由方匡南撰写，第17章由方匡南、姜叶飞撰写。第18章由方匡南、张声威撰写。最后由方匡南、朱建平进行加工总纂、修改定稿。

王秉权参与了第4章到第12章的校对、修改等工作。欧阳汉参与了第7章到第9章的校对、修改。

感谢成都道然科技有限责任公司专业意见和建议。感谢四川大学严珂玮为本书配的精美插画。感谢夫人张晶在我的写作过程中给予的支持和帮助。再次感谢为本书提供了直接或者间接帮助的各位朋友，没有你们的帮助，本书的出版没有这么顺利。

该书的数据可以从网址zhiliaobang.com或者kuangnanfang.com下载。

由于作者水平有限，书中难免有错误和不足之处，恳请读者批评指正！

方匡南
2014.8.1

目录

第 1 章 初识R语言.....	1	3.3.3 repeat语句.....	41
1.1 什么是R语言.....	1	3.4 编写自己的函数.....	41
1.2 为什么用R语言.....	2	3.4.1 函数名.....	42
1.3 安装R.....	4	3.4.2 关键词function.....	42
1.4 R扩展包.....	4	3.4.3 参数.....	42
1.4.1 R扩展包的安装与载入.....	5	3.4.4 函数体和函数返回值.....	44
1.4.2 R包的使用.....	6	3.5 程序调试.....	45
1.5 R编辑器.....	7	3.6 程序运行时间与效率.....	46
1.6 工作空间.....	11	3.7 用R做优化求解.....	47
第 2 章 数据结构与基本运算.....	13	3.7.1 一元函数优化求解.....	48
2.1 数据类型.....	13	3.7.2 多元函数优化求解.....	48
2.2 数据对象.....	14	3.7.3 约束条件下的优化求解.....	50
2.2.1 向量.....	15	3.8 习题.....	52
2.2.2 矩阵.....	21	第 4 章 随机数与抽样模拟.....	54
2.2.3 数组.....	31	4.1 一元随机数的产生.....	54
2.2.4 因子.....	32	4.1.1 均匀分布随机数.....	54
2.2.5 列表.....	33	4.1.2 正态分布随机数.....	56
2.2.6 数据框.....	34	4.1.3 指数分布随机数.....	57
2.3 习题.....	36	4.1.4 离散分布随机数的生成.....	58
第 3 章 函数与优化.....	38	4.1.5 常见分布函数表.....	59
3.1 常用的R内置函数.....	38	4.2 多元随机数的生成.....	61
3.2 条件控制语句.....	38	4.2.1 多元正态分布随机数.....	61
3.2.1 if/else语句.....	38	4.2.2 多元正态分布密度函数、	
3.2.2 ifelse 语句.....	39	分位数与累积概率.....	63
3.2.3 switch语句.....	39	4.2.3 多元t分布随机数.....	64
3.3 循环语句.....	40	4.3 随机抽样.....	65
3.3.1 for循环.....	40	4.3.1 放回与无放回抽样.....	65
3.3.2 while循环.....	40	4.3.2 bootstrap重抽样.....	66
		4.4 统计模拟.....	67

4.4.1 几种常见的模拟方法	67	6.4.1 访问数据框数据.....	115
4.4.2 模拟函数的建立方法	70	6.4.2 多变量数据的分析.....	118
4.5 习题.....	73	6.5 习题.....	124
第5章 数据读写与预处理	74	第7章 参数假设检验	126
5.1 数据的读入	74	7.1 假设检验的思想与步骤	126
5.1.1 直接输入数据.....	74	7.1.1 假设检验的基本思想	126
5.1.2 读R包中的数据.....	75	7.1.2 假设检验的基本步骤	128
5.1.3 从外部文件读入数据	75	7.2 正态总体单样本参数假设检验	129
5.2 写出数据	79	7.2.1 均值的检验.....	130
5.3 数据预处理	80	7.2.2 方差检验.....	132
5.3.1 变量预处理.....	81	7.3 正态总体双样本参数假设检验	134
5.3.2 变量重编码.....	82	7.3.1 双样本方差的检验	
5.3.3 变量重命名.....	84	(方差齐性检验)	134
5.3.4 变量类型的转换.....	85	7.3.2 两样本均值检验.....	135
5.3.5 日期变量的变换.....	86	7.4 比例假设检验	139
5.4 缺失数据处理	87	7.4.1 单样本比例检验.....	139
5.4.1 缺失数据的识别.....	87	7.4.2 两样本比例检验.....	141
5.4.2 缺失数据的探索与检验	88	7.5 习题.....	142
5.4.3 缺失数据的处理.....	89	第8章 非参数假设检验	144
5.5 数据集的合并与拆分	90	8.1 图示法.....	144
5.5.1 数据框的合并与拆分	90	8.2 卡方检验	146
5.5.2 数据集的合并.....	92	8.2.1 卡方分布(χ^2 distribution)	147
5.5.3 数据集的抽取.....	92	8.2.2 卡方拟合优度检验.....	148
5.6 习题.....	93	8.2.3 卡方独立性检验.....	151
第6章 探索性数据分析	94	8.2.4 卡方两样本同质性检验	151
6.1 主要分析工具	94	8.3 秩和检验	152
6.1.1 探索性数据分析的工具	94	8.3.1 秩的概念.....	153
6.1.2 数据的类型.....	98	8.3.2 单样本符号秩检验.....	153
6.2 单变量数据分析	99	8.3.3 两独立秩和检验.....	154
6.2.1 分类型数据.....	99	8.3.4 多个独立样本的秩和检验... ..	155
6.2.2 数值型数据.....	101	8.3.5 多个相关样本的秩和检验... ..	158
6.2.3 离群值探索.....	106	8.4 K-S检验.....	160
6.3 双变量数据分析	109	8.4.1 K-S单样本总体分布验证... ..	160
6.3.1 分类数据对分类数据	109	8.4.2 K-S两独立样本同质检验... ..	160
6.3.2 分类数据对数值型数据	111	8.5 常用正态性检验	162
6.3.3 数值型数据对数值型数据... ..	112	8.5.1 偏度、峰度检验法.....	162
6.4 多变量数据分析	115	8.5.2 Shapiro-Wilk (W检验)	163

8.5.3 其他常用正态检验.....	165	11.3 序列相关性	240
8.6 习题.....	167	11.3.1 问题的提出.....	241
第 9 章 方差分析	169	11.3.2 序列相关性定义及后果....	243
9.1 单因素方差分析	170	11.3.3 序列相关性检验.....	245
9.2 双因素方差分析	174	11.3.4 序列相关性克服.....	248
9.2.1 不考虑交互作用的双因		11.4 习题.....	251
素方差分析.....	174	第 12 章 非线性回归分析	254
9.2.2 考虑交互作用的双因素		12.1 问题的提出	254
分析.....	178	12.2 可线性化的非线性回归	255
9.3 习题.....	183	12.2.1 Cobb-Douglas生产函数 ...	255
第 10 章 线性回归模型	184	12.2.2 多项式方程模型.....	257
10.1 问题提出	184	12.2.3 指数函数模型.....	259
10.2 一元线性回归	185	12.3 不可线性化的非线性回归	260
10.2.1 一元线性回归概述.....	186	12.3.1 非线性模型的参数估计	
10.2.2 一元线性回归的参数估计 ...	188	与迭代算法.....	262
10.2.3 一元线性回归模型的检验 ...	195	12.3.2 初始值选取.....	269
10.2.4 一元线性回归的预测	197	12.3.3 收敛性.....	270
10.2.5 一元线性回归综合案例 ...	201	12.4 非线性回归评价和假设检验	271
10.3 多元线性回归分析	205	12.4.1 可决系数.....	271
10.3.1 多元线性回归模型及假定... 206		12.4.2 参数显著性的 F 检验	271
10.3.2 参数估计.....	207	12.4.3 似然比检验.....	272
10.3.3 模型检验.....	209	12.5 习题.....	274
10.3.4 预测.....	211	第 13 章 二元选择模型	275
10.3.5 多元线性回归综合案例 ... 213		13.1 问题的提出	276
10.4 习题.....	218	13.2 线性概率 (LP) 模型原理.....	277
第 11 章 线性回归模型的扩展	220	13.3 Probit模型原理.....	279
11.1 多重共线性	220	13.4 Logit模型原理	280
11.1.1 问题的提出.....	220	13.5 边际效应分析	281
11.1.2 多重共线性定义及后果....	222	13.6 最大似然估计 (MLE)	282
11.1.3 多重共线性检验.....	222	13.7 似然比检验和拟合优度	282
11.1.4 多重共线性克服.....	225	13.8 案例分析: 经济学教学新方法	
11.2 异方差性	229	的效果.....	284
11.2.1 问题的提出.....	229	13.9 扩展案例: 信用卡违约预测分析... 289	
11.2.2 异方差性定义及后果	231	13.9.1 描述性统计.....	290
11.2.3 异方差性检验.....	232	13.9.2 模型建立与参数估计	291
11.2.4 异方差性克服.....	236	13.9.3 系数意义与边际分析	295
		13.9.4 拟合与预测.....	296

13.9.5 结论与建议.....	297	第 16 章 分位数回归.....	330
13.10 习题.....	297	16.1 问题的提出.....	330
第 14 章 多元选择模型.....	299	16.2 总体分位数和总体中位数.....	332
14.1 有序选择模型.....	299	16.3 经验分位数估计.....	333
14.1.1 问题的提出: 本科生申 请研究生的影响因素.....	300	16.4 分位数回归原理.....	334
14.1.2 有序选择模型.....	300	16.5 扩展案例: 社会保障与城乡 家庭消费.....	339
14.1.3 案例分析: 本科生申请 研究生的影响因素.....	302	16.5.1 问题的提出.....	339
14.2 多元无序Logit模型.....	304	16.5.2 数据说明.....	339
14.2.1 问题的提出: 关于钓鱼 模式的选择.....	304	16.5.3 实证分析.....	342
14.2.2 多元无序Logit模型.....	305	16.5.4 结论与建议.....	345
14.2.3 案例分析: 关于钓鱼模 式的选择.....	307	16.6 习题.....	345
14.3 嵌套Logit模型.....	309	第 17 章 高级统计绘图.....	346
14.3.1 问题的提出: 旅行交通 方式选择.....	309	17.1 绘制地图.....	346
14.3.2 嵌套Logit模型原理.....	310	17.2 高阶绘图工具——ggplot2.....	355
14.3.3 案例分析: 旅行交通方 式选择.....	311	17.2.1 散点图.....	355
14.4 习题.....	313	17.2.2 散点图上添加平滑曲线... ..	358
第 15 章 计数模型与受限因变量模型.....	314	17.2.3 条形图和箱线图.....	360
15.1 计数模型.....	314	17.2.4 直方图和密度曲线图.....	362
15.1.1 问题的提出: 轮船事故 的计数数据模型.....	314	17.2.5 时间序列图.....	364
15.1.2 计数数据模型的设定.....	316	17.2.6 图形标注.....	365
15.1.3 计数数据模型的估计.....	317	17.3 三维图形与等高线图.....	366
15.2 受限因变量模型.....	319	17.3.1 三维图形.....	366
15.2.1 截断模型的问题提出.....	319	17.3.2 等高图/等高线.....	368
15.2.2 截断模型原理.....	319	17.4 词云.....	369
15.2.3 审查模型问题的提出.....	321	17.5 散点图矩阵与关系矩阵图.....	370
15.2.4 审查模型原理.....	322	17.6 马赛克图.....	372
15.2.5 最大似然估计 (MLE).....	323	17.7 习题.....	374
15.3 习题.....	328	第 18 章 如何制作自己的R包.....	375
		18.1 R包基础.....	376
		18.2 在Windows中制作R包.....	377
		18.3 在RStudio中制作R包.....	381
		18.4 习题.....	383
		参考文献.....	384

第 1 章

初识R语言

本章主要对R语言的基础内容做讲解，使读者可以尽快熟悉这款流行的数据分析工具。本章内容涵盖：R语言的优势、安装流程、扩展包、常用编辑器和工作空间。

1.1 什么是R语言

你现在最常用的统计软件是什么？SAS，MATLAB，还是Eviews？如果你是一个经常和数据分析打交道，需要运用或自己编写各种统计方法的人，但还没用上R，那么你已经脱离现在的主流了。

R语言是由新西兰奥克兰大学的Ross Ihaka与Robert Gentleman一起开发的一个面向对象的编程语言，因两人的名字都是以R开头，所以命名为“R”。R是“GNU S”，一个免费开源、能够自由有效地用于统计计算和绘图的语言和环境，可以在UNIX、Windows和Mac OS系统中运行，它提供了广泛的统计分析和绘图技术，包括回归分析、时间序列、分类和聚类等方法。R语言的前身是S语言，S语言是贝尔实验室（Bell Laboratories）的Rick Becker、John Chambers和Allan Wilks共同开发的，提供了一系列统计和图形显示工具，这个语言过去一度是数据分析领域里面的标准语言，但现在正逐步被R语言取代。

R是一套完整的数据处理、计算和制图软件系统，它是一套开源的数据分析解决方案，由一个庞大而活跃的全国性社区维护。其功能包括：数据存储和处理系统；数组运算工具（其向量、矩阵运算方面功能尤其强大）；完整连贯的统计分析工具；优秀的统计制图功能；简便而强大的编程语言——可操纵数据的输入和输出，可实现分支、循环，用户可自定义功能。与其说R是一种统计软件，还不如说R是一种统计分析与计算的环境，因为R不仅提供若干统计程序，而且还可进行统计分析，只需使用者指定数据库和若干参数即可。R的思想是：它可以提

供一些集成的统计工具，更重要的是，它还可以提供各种数学计算、统计计算的函数，从而令使用者（用户）能够灵活地进行数据分析，甚至创造出符合需要的新的统计计算方法。随着R功能的不断完善，也提供了如图1-1所示的许多高级功能。



图1-1 R能够提供的高级功能

1.2 为什么用R语言

2009年，《纽约时报》发表了题为“Data Analysts Captivated by R’s Power”的社评，集中讨论了R语言在数据分析领域的发展，并引发了SAS和R用户广泛而激烈的争论。文章认为，让R变得如此有用和如此快地广受欢迎的原因是统计学家、工程师和科学家们在不断精炼代码或编写各种特有、具体的包。而且现在R软件增添了很多高级算法、作图颜色、文本注释以及与数据库等链接的挖掘技术。

KDnuggets网站每年都会做一些数据分析和数据挖掘软件使用的专题问卷调查。据KDnuggets网站2011年对570个数据挖掘和数据分析的工作者关于过去12个月数据挖掘和数据分析所使用的编程语言的调查显示（<http://www.kdnuggets.com/2011/08/poll-languages-for-data-mining-analytics.html>），R语言名列榜首，如图1-2所示，占接近半壁江山（45%），而紧随其后的SQL、Python和Java则在某一领域具有各自独到的优势，而SAS和MATLAB分别名列第5和第6，被R远远甩在后面。

除了R之外，现在还有许多人使用SPSS、SAS或MATLAB，但大都用的是盗版软件。免费是R流行开来的最大的一个因素。不过，R的最大优点是出色的可视化图形、丰富的统计方法及高效的更新速度。其主要特点如下。

- ① 高效的数据处理和保存机制。
- ② 完整的数组和矩阵操作运算符，以及完整的数据分析工具。

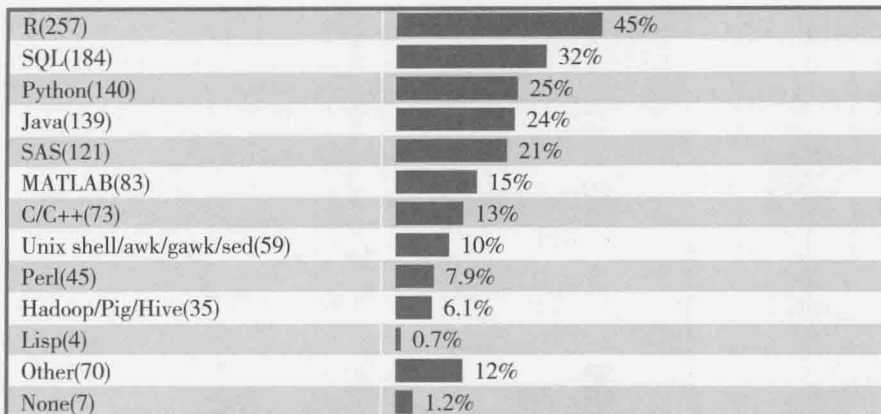


图1-2 数据挖掘和数据分析编程语言使用调查结果

③ 出色的图形统计功能。

④ 简单高效的建模工具。

⑤ 作为一个开源软件，R具有丰富的扩展包（packages），可以自由加载其他开发者提供的函数和数据包，直接利用可以节省很多重新编写算法的精力。

⑥ 提供了很多高级功能。除统计外，还可以使用R来给电脑关机、发微博、发校内状态、下五子棋，或是配合LaTeX撰写动态统计报告以及自动生成概率统计的试卷和答案，等等。

⑦ 几乎兼容全平台。除了支持OS X、Linux和Windows之外，甚至可以在iOS设备上编辑和运行R的程序，还可以在iPhone等移动设备上安装R程序。

⑧ 逐渐支持多国语言。作为一个开源软件，R在其主页上提供了可供大家添加自己国家语言的文件(<http://developer.r-project.org/Translations.html>)，其中，中文的翻译在<https://github.com/r-cn/r-cntrans>上面得到了国内众多的R爱好者的支持。

⑨ 更新速度快。R几乎囊括了所有统计方法，当其他软件还不能完成一些最新的统计方法时，在R中几乎都可以完成。R的更新速度是以周来计算的。

R功能强大、内容丰富，各种帮助文档很多。目前R已有近5000个扩展包，这些包由独立贡献者编写，质量参差不齐，文档零散而且很难找到。不得不承认R的学习曲线比较陡峭。要掌握R的所有功能是非常困难的，因此必须有一本能带你从入门到精通的参考书，这样会让你事半功倍！本书总结作者10多年R的使用经验，深入浅出地讲解了R语言如何进行数据分析和作图。

1.3 安装R

R的官方网站是<http://www.r-project.org>，打开该网站，界面如图1-3所示。

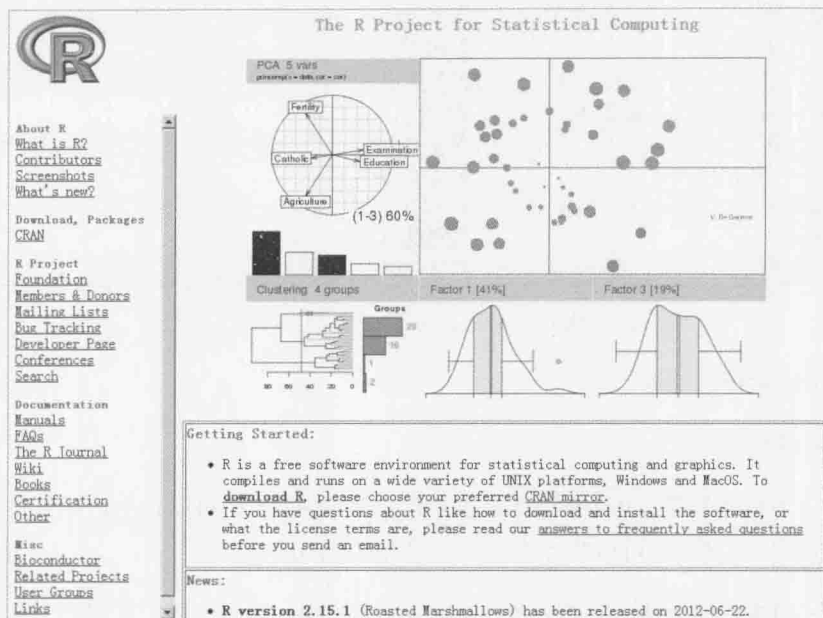


图1-3 R官方网站首页

R安装软件和安装包（package）的下载链接在首页左侧的CRAN。R除了自己的主程序外，还有5000多个用户贡献的扩展（附加）包，以及各种文档，这一大箩筐东西被复制到世界各地几十台服务器上供用户下载。

作为中国用户，一般都会选择一个国内的镜像。目前，中国大陆的镜像有如下4个。

http://ftp.ctex.org/mirrors/CRAN/	CTEX.ORG
http://mirror.bjtu.edu.cn/cran	Beijing Jiaotong University, Beijing
http://mirrors.ustc.edu.cn/CRAN/	University of Science and Technology of China
http://mirrors.xmu.edu.cn/CRAN/	Xiamen University

在R官方主页单击download R，选择对应的镜像后，在右侧下载和安装包栏目里会出现三种操作系统的R版本（Linux、OS X和Windows），选择相应的操作系统后，再选择base就会进入R的下载页面，在页面上会出现R的最新版本和安装说明等文件。

1.4 R扩展包

整个R系统主要是由一系列程序包（package）组成的。R包是函数、数据和预编译代码以

一种定义完善的格式组成的集合。这些包默认储存在library里。R安装好后，默认自带了base、datasets、utils、stats、grDevices、graphics和methods包，除这些包外，CRAN上还有数千附加包，由R核心开发团队之外的用户自行提交，如此多的程序包基本上保证了R可以实现目前几乎所有的统计分析。扩展包需要下载安装并载入到会话中才能使用，接下来介绍如何安装和载入R扩展包。

1.4.1 R扩展包的安装与载入

R有多种安装方法，这里主要介绍其中两种。

① 在线安装。例如，需要安装“class”这个扩展包，则输入命令 `install.packages("class")` 执行即可。也可以同时安装多个包，例如，需同时安装“class”和“cluster”两个包，则输入命令 `install.packages(c("class","cluster"))` 执行即可。

② 利用RStudio安装。RStudio在1.5节会详细介绍。在RStudio右下角栏目里单击“packages”，然后再单击“install packages”，会出现如图1-4所示对话框，在packages对话框里输入需要安装的包名，如输入class，就可以安装包。

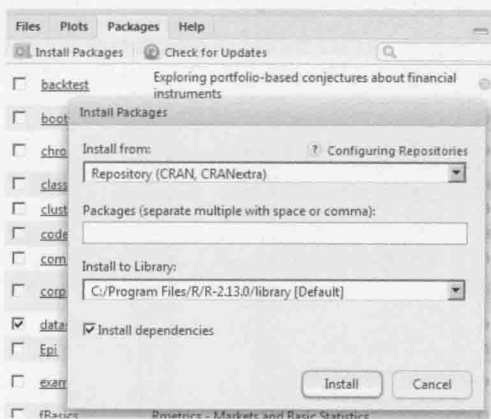


图1-4 用RStudio安装包

安装完R扩展包后，并不能马上使用，需要载入到当前会话中才能使用。R包的载入有两种方法。

① 在R Console中输入命令 `library()` 载入。如需要载入class包，则输入 `library(class)` 即可。

② 利用RStudio载入。安装好后的包会出现在RStudio右下方的栏目里，即如图1-5所示对话框，单击左侧的方框，即可载入相应的包。

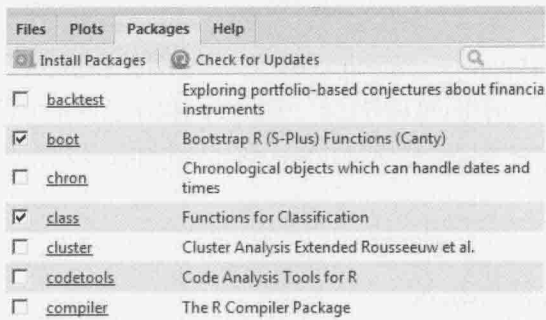


图1-5 RStudio载入包

一个包仅需安装一次即可，但是R的扩展包经常被更新，要想使用更新后的扩展包，则需要对已安装的包进行更新。下面介绍两种常用方法。

① 在R Console中输入命令`update.packages()`，则会出现需要更新的包和相关信息，并可以选择是否需要更新。如果选择“y”，则会在线更新。

```
> update.packages()
boot :
Version 1.2-43 installed in C:/Program Files/R/R-2.13.0/library
Version 1.3-2 available at http://mirrors.xmu.edu.cn/CRAN
Update (y/N/c)?
```

② 利用RStudio更新R包。单击“Check for Updates”按钮，会出现如图1-6所示对话框。在左侧单击需要更新的包，即可进行更新。

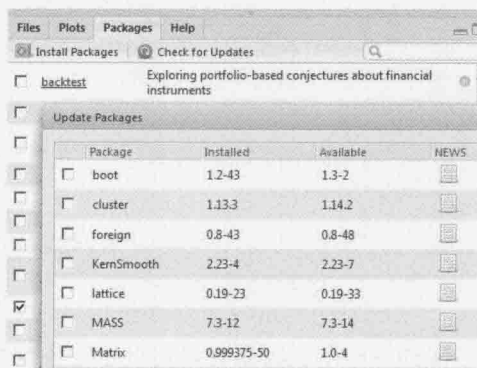


图1-6 利用RStudio更新R包

1.4.2 R包的使用

载入R包后，可以使用该扩展包的函数和数据。例如，我们载入“class”包后，想要查看该包里的函数，一种方法是在R Console里输入命令`help(package="class")`，会返回一个包含了“class”包里所有函数和数据集名称的帮助文档，如图1-7所示；另一种方法是在RStudio右侧

单击“class”，则会返回如图1-8所示的帮助文档，如需进一步了解该函数的使用方法，则继续单击相应的函数，就会返回该函数的详细帮助文档。

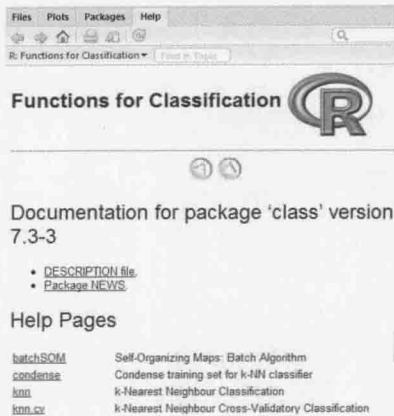


图1-7 class帮助文档

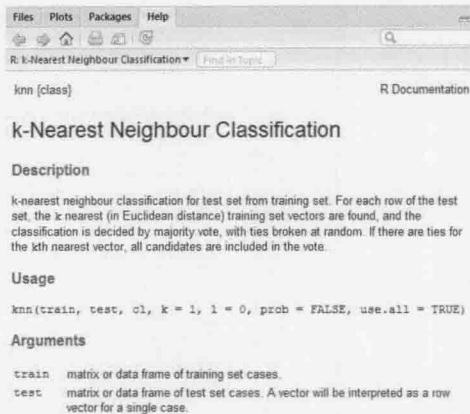


图1-8 knn函数帮助文档

1.5 R编辑器

Windows版本下载安装完成后，在桌面会出现图标^①。双击该快捷图标，打开R的操作平台（R Console），如图1-9所示。虽然可以直接在此操作平台上输入命令，但是R本身默认的IDE不是很友好，甚至可以用粗糙来形容。通常情况下，建议不要直接在操作平台上输入命令，而是使用R的编辑软件。

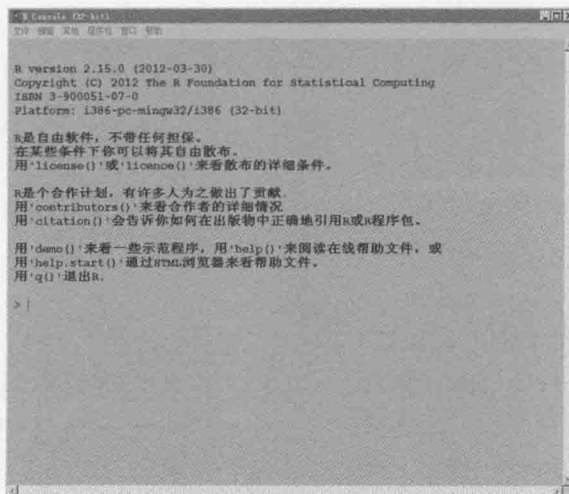


图 1-9 R的操作平台

作为一个开源软件，R不乏各种优秀的IDE。本书主要使用其中一款优秀的跨平台开源IDE——RStudio。RStudio将常用的窗口都整合在一起，使开发者无须在命令行和绘图窗口之间跳

来跳去,更便于控制。不过,要想使用RStudio,除了需要安装R开发环境外,还需要到<http://www.rstudio.org/download/desktop/> 上下载适合你电脑的RStudio版本并安装。

打开RStudio,如图1-10所示,左上角是脚本编辑窗口,左下角是命令窗口,右上角是工作空间和历史记录窗口,右下角是作图、帮助和包管理窗口。

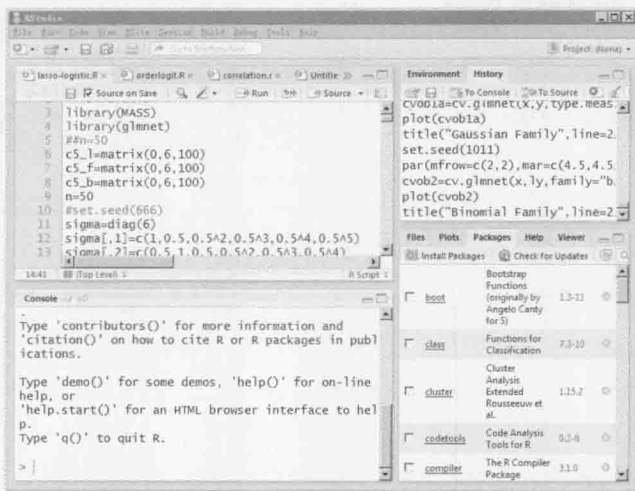


图1-10 RStudio用户界面

如果需要编写R语言代码,则可以用这里提供的两种途径运行。

① 在左下角的命令窗口输入并回车运行。

② 使用左上角的脚本编辑窗口编写好代码之后,根据你的需要选择运行代码。

如果需要使用脚本编辑代码,则需要在RStudio中执行File→New File→RScript命令(快捷键Ctrl+Shift+N)来新建一个脚本编辑窗口。写好代码之后,可以通过单击如图1-11所示的Run按钮来运行程序。如果直接单击Run按钮,则只是运行当前行代码。如果先用鼠标选好要运行的代码,然后再单击Run按钮,则会运行所选代码。

在使用脚本编辑窗口进行编辑时,RStudio还提供了一些高级功能。



图1-11 运行R脚本程序