

Official Study Kit for
Wrox Certified Big Data Analyst Program

大数据分析师权威教程

大数据分析 与 预测建模

Wrox 国际 IT 认证项目组 / 编 姚军 / 译

Official Study Kit for
Wrox Certified Big Data Analyst Program

大数据分析师权威教程

大数据分析 与 预测建模



Wrox 国际 IT 认证项目组 / 编 姚军 / 译

人民邮电出版社
北京

图书在版编目 (C I P) 数据

大数据分析权威教程. 大数据分析与预测建模 /
Wrox国际IT认证项目组编; 姚军译. — 北京: 人民邮
电出版社, 2017. 11

书名原文: Official Study Kit for Wrox
Certified Big Data Analyst Program
ISBN 978-7-115-46366-1

I. ①大… II. ①W… ②姚… III. ①数据处理—教材
IV. ①TP274

中国版本图书馆CIP数据核字(2017)第202599号

内 容 提 要

“大数据”已连年入选 IT 领域的热点话题, 人们每天都会通过互联网、移动设备等生产大量数据。如何从海量数据中洞悉出隐藏其后的见解是当今社会各领域人士极为关注的话题。本系列图书以“大数据分析师”应掌握的 IT 技术为主线, 共分两卷, 以 7 个模块(第 1 卷包括 4 个模块, 第 2 卷包括 3 个模块)分别介绍大数据入门, 分析和 R 编程入门, 使用 R 进行数据分析, 用 R 进行高级分析, 机器学习的概念, 社交媒体、移动分析和可视化, 大数据分析的行业应用等核心内容, 全面且详尽地涵盖了大数据分析的各个领域。

本书为第 1 卷, 首先提供大数据的概览, 介绍大数据概念及其在商业中的应用、处理大数据的技术、Hadoop 生态系统和 MapReduce 的相关内容, 然后介绍如何理解分析、分析方法与工具, 重点讲解流行分析工具 R, 介绍如何将数据集导入 R 和从 R 导出数据、在 R 中如何操纵和处理数据, 最后详细介绍 R 中的函数和包、R 的描述性统计、R 中的图形分析、R 中的假设检验、R 中的线性回归、非线性回归、聚类分析、决策树、R 和 Hadoop 的集成及 Hive, 通过这些实战内容, 使读者掌握 R 语言在数据分析中的全面应用。通过本书, 读者能对大数据概念、重要性及其应用有全面的了解, 熟悉各种大数据分析工具。

本书适用于想成为大数据分析师的人员以及所有对大数据分析感兴趣的技术人员和决策者阅读。

-
- ◆ 编 者 Wrox 国际 IT 认证项目组
 - 译 者 姚 军
 - 责任编辑 杨海玲
 - 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>
 - 北京鑫正大印刷有限公司印刷

 - ◆ 开本: 800×1000 1/16
 - 印张: 32
 - 字数: 715 千字 2017 年 11 月第 1 版
 - 印数: 1—2 400 册 2017 年 11 月北京第 1 次印刷
- 著作权合同登记号 图字: 01-2015-2395 号
-

定价: 108.00 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

版权声明

Copyright © 2014 by respective authors:

Module 1	Session 1	Kogent Learning Solutions Inc.
Module 1	Session 2	Bill Franks
Module 1	Session 3~5	Judith Hurwitz, Alan Nugent, Dr. Fern Halper, Marcia Kaufman
Module 2	Session 1~2	Bill Franks
Module 2	Session 3~4	Joris Meys, Andrie de Vries, Mark Gardener
Module 2	Session 5	Joris Meys, Andrie de Vries
Module 3	Session 1	Joris Meys, Andrie de Vries
Module 3	Session 2	Mark Gardener
Module 3	Session 3	Dr Murray Logan
Module 3	Session 4	Michael J. Crawley
Module 3	Session 5	Mark Gardener
Module 4	Session 1	Michael J. Crawley, Deborah J. Rumsey, Mark Gardener
Module 4	Session 2	Michael J. Crawley
Module 4	Session 3	Johannes Ledolter, Stephane Tuffery
Module 4	Session 4	Dean Abbott
Module 4	Session 5	Kogent Learning Solutions Inc.
Appendix		Joris Meys, Andrie de Vries

All rights reserved and the following credit: Authorized translation from English language edition published by Wiley India Pvt. Ltd.

前言

欢迎阅读《大数据分析师权威教程》和《大数据开发者权威教程》!

信息技术蓬勃发展,每天都有新产品问世,同时不断地形成新的趋势。这种不断的变化使得信息技术和软件专业人员、开发人员、科学家以及投资者都不敢怠慢,并引发了新的职业机会和有趣的工作。然而,竞争是激烈的,与最新的技术和趋势保持同步是永恒的要求。对于专业人士来说,在全球 IT 行业中,入行、生存和成长都变得日益复杂。

想在 IT 这样一个充满活力的行业中高效地学习,就必须做到:

- 对核心技术概念和设计通则有很好的理解;
- 具备适应各种平台和应用的敏捷性;
- 对当前和即将到来的行业趋势和标准有充分的认识。

鉴于以上几点,我们很高兴地为大家介绍《大数据分析师权威教程》(两卷)和《大数据开发者权威教程》(两卷)系列。

这两个系列共 4 本书,旨在培育新一代年轻 IT 专业人士——他们能够灵活地在多个平台之间切换,并能胜任核心职位。这两个系列是在对技术、IT 市场需求以及当今就业培训方面的全球行业标准进行了广泛并严格的调研之后才开发出来的。这些计划的构思目标是成为理想的就业能力培训项目,为那些有志于在国际 IT 行业取得事业成功的人提供服务。这一系列目前已经包含了一些最为热门的 IT 领域中的认证项目,如大数据、云、移动和网络应用程序、网络安全、数据库和网络、计算机操作、软件测试等。根据我们的全球质量标准加以调整之后,这些项目还能帮助你识别和评估职业机会,并为符合全球最著名企业的招聘流程做好最佳的准备。

这两个系列是学习和培训资源的知识库,为在重要领域和信息技术行业中培养厂商中立和平台独立的专业能力而设立。这些资源有效地利用了创新的学习手段和以成果为导向的学习工具,培养富有抱负的 IT 专业人士。同时也为开设大数据分析师和大数据开发者相关培训课程的讲师提供了全面综合的教学和指导方案。

《大数据分析师权威教程》系列概览

大数据可能是今天的科技行业中最受欢迎的流行语之一。全世界的企业都已经意识到了可用的大量数据的价值,并尽最大努力来管理和分析数据、发挥其作用,以建立战略和发展竞争优势。与此同时,这项技术的出现,导致了各种新的和增强的工作角色的演变。

《大数据分析师权威教程》系列的目标是培养新一代的国际化全能大数据分析师,使他们精通数据挖掘、数据操纵和数据分析方面的基本及高级分析技术,熟悉大数据平台以及业务和行业

需求，能够高效地参与大数据分析项目。

本系列旨在：

- 使参与者熟悉整个数据分析的生命期；
- 通过众多案例分析，使参与者熟悉大数据在不同相关行业中的角色和用途；
- 提供基本及高级大数据分析以及可视化的完整技术诀窍，帮助他们分析数据、创建统计模型和提供业务洞察力；
- 最后包含一个完整的项目，使参与者能够实施分析生命周期。

学习者的必备条件

要阅读这个系列图书，读者必须具备以下基础知识：

- 统计学基础知识，包括主要趋势和平均值计量、分散度计量、概率；
- 基本图表、直方图和散点图的创建；
- 基本熟悉数据库、表和字段，包括电子表格与计算。

建议的学习时间

《大数据分析权威教程》由 7 个学习模块组成（第 1 卷包括 4 个模块，第 2 卷包括 3 个模块）。

根据参与者的技能水平，可以选择任何数量的模块以积累特定领域的技能，每个模块的学习目标会在后面列出。

对于入门级的参与者来说，建议选择全部 7 个模块，以便为成为合格大数据分析师做好充足的准备。专业人士或者已经拥有某些必备技能的参与者则可以选择能够帮助自己加强特定领域技能的模块。

每个模块可能占用大约 10 小时的学习时间，因此完整的学习时间大约是 70 小时。

模块清单

《大数据分析权威教程》第 1 卷的 4 个模块的具体名称和学习目标如表 1 所示。

表 1

模块编号	模块名称	模块目标
模块 1	大数据入门	<ul style="list-style-type: none">● 了解大数据的角色和重要性● 讨论大数据在各行各业中的使用和应用● 讨论大数据相关的主要技术● 解释 Hadoop 生态系统中各种组件的角色● 解释 MapReduce 的基础概念和它在 Hadoop 生态系统中的作用

续表

模块编号	模块名称	模块目标
模块 2	分析和 R 编程入门	<ul style="list-style-type: none"> ● 讨论高级分析的重要性 ● 介绍分析方法和工具的发展 ● 讨论各种分析工具的特性 ● 用 R 语言开发脚本 ● 用 R 语言中的各种附加编辑器执行脚本 ● 用 R 语言执行读写操作 ● 用 R 语言操纵数据
模块 3	使用 R 语言进行数据分析	<ul style="list-style-type: none"> ● 使用 R 脚本和函数 ● 使用 R 函数环境和方法 ● 执行数据样本总结步骤 ● 使用积累的统计数据和汇总表 ● 用 R 创建列表、矩阵和数据帧 ● 使用 R 中的循环和条件执行 ● 安装 R Hadoop 和创建用户定义函数 ● 用 R 实现图表分析 ● 用 R 进行假设检验
模块 4	用 R 语言进行高级分析	<ul style="list-style-type: none"> ● 描述线性回归分析及其应用 ● 在 R 语言中应用线性回归分析的知识 ● 从应用角度理解非线性回归 ● 在 R 语言中应用非线性回归分析 ● 解释聚类分析技术 ● 用 R 实现聚类分析 ● 探索用于构建决策树的基本概念 ● 用 R 构建决策树 ● 将 R 与 Hadoop 集成, 以进行统计分析

《大数据分析师权威教程》第 2 卷的 3 个模块的具体名称和学习目标如表 2 所示。

表 2

模块编号	模块名称	模块目标
模块 1	机器学习的概念	<ul style="list-style-type: none"> ● 讨论机器学习在技术上和商业上的应用 ● 理解图模型的用途 ● 用 R 实现图模型 ● 理解贝叶斯网络表示法及其解读 ● 用贝叶斯网络解决预测问题 ● 探索神经网络及其结构和学习规则 ● 阐述神经网络的训练 ● 用 R 实现神经网络 ● 用因子分析和主成分分析实现降维 ● 从给定的预测因素列表识别最大影响因子/维度 ● 解释支持向量机 ● 用 R 语言实现支持向量机

模块编号	模块名称	模块目标
模块 2	社交媒体、移动分析和可视化	<ul style="list-style-type: none"> ● 应用可用于大数据实现的解决方案设计过程 ● 分析业务环境中社交媒体所承担的角色 ● 实施社交媒体分析 ● 执行基本移动分析 ● 讨论数据可视化及其重要性 ● 使用表格进行数据可视化 ● 有效地准备求职面试
模块 3	大数据分析的行业应用	<ul style="list-style-type: none"> ● 理解保险业中的数据分析应用 ● 理解金融机构中数据分析的实施 ● 理解电信行业中的分析工具 ● 实施在线客户细分中的分析

学习方法和特色

本书开发了一套独特的学习方法，这种专门设计的方法不仅以最大限度地学习大数据概念为目标，还注重对真实专业环境下应用这些概念的全面理解。

本书的独特方法和丰富特性简单介绍如下。

- 涵盖了大数据分析师必备的所有**大数据和 Hadoop 基础组件及相关组件的基本知识**，使学习者有可能在一个系列书中获得对所有相关知识、新兴技术和平台的了解。
- 在与大数据分析师关系最为密切的**描述性和预测性分析技术**上培养全面、结构化的技能，逐步理解**各种技术在 R 语言上的实施**（R 语言是最通用、使用最广泛的统计软件之一）。
- **基于场景的学习方法**，通过多种有代表性的现实场景的使用和案例研究，将 IT 基础知识融入现实环境，鼓励学习者积极、全面地学习和研究，实现体验式教学。
- 强调目标明确、**基于成果的学习**。每一讲都以“本讲目标”开始，该目标会进一步关联整个教程的更广泛的目标。
- **简明、循序渐进的编程和编码指导**，清晰地解释每行代码的基本原理。
- 强调**高效、实用的过程和技术**，帮助学习者深入理解巧妙、合乎道德的专业方法及其对业务的影响。

学习工具

下列学习工具将确保学习者高效地使用本教程。

- **模块目标**：列出某一讲所属模块的目标。
- **本讲目标**：列出与模块目标对应的本讲目标。
- **预备知识**：说明对特定部分或者整体概念的理解有特定作用的预备知识点。

- **交叉参考**：将整个模块中学到的相关概念联系起来，启发参与者理解分析中的不同功能、职责和挑战，确保任何概念都不是孤立地学习的。
- **总体情况**：不断提醒参与者，某个主题为什么是相关的，在行业中如何应用，从而为学习提供实践维度。
- **快速提示**：提供明智、高效地运用概念的简便技巧。
- **与现实生活的联系**：提供简短的案例分析和剪报，阐述概念在现实世界中的适用性。
- **技术材料**：提供加强技术诀窍理解的技巧和信息。
- **定义**：定义重要概念或者术语。
- **附加知识**：提供相关的附加信息。
- **知识检测点**：提出互动式课堂讨论的问题，强化每一讲之后的学习。
- **练习**：在每一讲结束时提出以知识为基础的实践问题，评估理解情况。
- **测试你的能力**：提供基于应用的实践问题。
- **备忘单**：提供本讲涵盖的重要步骤及过程的快速参考。

关键的大数据技术术语

大数据是一个非常年轻的行业，新的技术和术语每周都会出现。这种快节奏的环境是由开源社区、新兴技术公司以及 IBM、Oracle、SAP、SAS 和 Teradata 这样的业界巨人推动的。不用说，建立一个持久的权威术语表是很难的。鉴于这样的风险，我们在这里只提供一个小型的大数据词汇表，如表 3 所示。

表 3

术 语	定 义
算法	用来分析数据的数学方法。一般情况下，是一段计算过程；计算一个功能的指令列表；在软件中，这样一个过程以编程语言来实际实现
分析	一组用于查询和梳理平台数据的分析工具和计算能力
装置	专为特定活动集建立的一组优化硬件和软件
Avro	一个可编码 Hadoop 文件模式的数据序列化系统，特别擅长于数据解析，是 Apache Hadoop 项目的一部分
批处理	在后台运行、不与人发生交互的作业或进程
大数据	大数据事实上的标准定义是超越了传统的 3 个维度（数据量、多样性、速度）限制的数据。这 3 个维度的结合使得数据的摄取、处理和呈现更加复杂
Big Insights	IBM 的具有企业级增值组件的 Hadoop 商业发行版
Cassandra	由 Apache 软件基金会管理的开源列式数据库
Clojure	基于 LISP(从 20 世纪 50 年代起的人工智能编程语言事实标准)的动态编程语言，读作“closure”。通常用于并行数据处理
云	用以指代任何计算机运作的软件、硬件或服务资源的通用术语。它作为一种服务通过网络传送

术 语	定 义
Cloudera	Hadoop 的第一个商业分销商。Cloudera 提供了 Hadoop 发行版的企业级增值组件
列式数据库	按列进行的数据存储与优化。使用基于列的数据，对于一些分析处理特别有用
复杂事件处理(CEP)	对实时发生事件进行分析并采取措施的过程
数据挖掘	利用机器学习，从数据中发现模式、趋势和关系的过程
分布式处理	在多个 CPU 上的程序执行
Dremel	一个可扩展、交互式、点对点分析查询系统，有能力在数秒内对数万亿行的表进行聚合查询
Flume	一种从 Web 服务器、应用服务器、移动设备等目标抓取数据填充 Hadoop 的框架
网格	松散耦合的服务器通过网络连接起来，并行处理工作负载
Hadapt	一家提供 Hadoop 相关插件的商业供应商，这个插件可以通过高速连接器在 HDFS 和关系型表中移动数据
Hadoop	一个开源项目框架，可以在计算机集群（网格）中存储大量的非结构化数据（HDFS）并在其中对其进行处理（MapReduce）
HANA	来自 SAP 的内存处理计算平台，为大容量事务和实时分析而设计
HBase	一种分布式、列式存储的 NoSQL 数据库
HDFS	Hadoop 文件系统，是 Hadoop 的存储机制
Hive	一种 Hadoop 的类 SQL 查询语言
Norton	具有企业级增值工作组件的 Hadoop 商业发行版
HPC	高性能计算。通俗地说，就是为高速浮点处理、内存磁盘并行化而设计的设备
HASstreaming	为 Hadoop 提供实时 CEP（复杂事件处理）的 Hadoop 商业插件
机器学习	从经验数据中学习，然后利用这些经验教训去预测未来新数据的结果的算法技术
Mahout	为 Hadoop 创建可伸缩机器学习算法库的 Apache 项目，主要用 MapReduce 实现
MapR	具有企业级增值组件的 Hadoop 商业发行版
MapReduce	一种 Hadoop 计算批处理框架，其中的作业大部分用 Java 编写。作业将较大的问题分解为较小的部分，并将工作负载分布到网格中，使多个作业能够同时进行（mapper）。主作业（reducer）收集所有中间结果并将其组合起来
大规模并行处理(MPP)	能协调并行程序执行的系统（操作系统、处理器和内存）
MPP 装置	带有处理器、内存、磁盘和软件，能够并行处理工作负载的集成平台
MPP 数据库	一种已为 MPP 环境优化的数据库
MongoDB	一种用 C++编写的可扩展、高性能的开源 NoSQL 数据库
NoSQL 数据库	一个用以描述数据库的术语。这种数据库不使用 SQL 作为数据库的数据主检索，且可以是任意类型。NoSQL 拥有有限的传统功能，并为可扩展性和高性能检索及添加而设计。通常情况下，NoSQL 数据库利用键值对存储数据，能够很好地处理在本质上不相关的数据
Oozie	一个工作流处理系统，允许用户定义一系列用各种语言（如 MapReduce、Pig 和 Hive）编写的作业
Pig	一种使用查询语言（Pig Latin）的分布式处理框架，用以执行数据转换。目前，Pig Latin 程序被转换为 MapReduce 作业，在 Hadoop 上运行

续表

术 语	定 义
R	一种开源的语言和环境，用以统计计算和图形化
实时	今天，通俗地说，它被定义为即时处理。实时处理起源于 20 世纪 50 年代，当时多任务处理机提供了为更高优先级任务的执行而“中断”一个任务的能力。这些类型的机器为空间计划、军事应用和多种商业控制系统提供了动力
关系型数据库	按照行和列存储和优化数据
Scoring	使用预测模型，预测新数据的未来结果
半结构化数据	依靠可用的格式描述符，把非结构化的数据放入结构中
Spark	内存分析计算处理的高性能处理框架，通常被用来做实时查询
SQL（结构化查询语言）	关系型数据库中，存储、访问和操作数据的语言
Sqoop	一种命令行工具，具有把单个表或整个数据库导入 Hadoop 文件中能力
Sorm	分布式、容错、实时分析处理的开源框架
结构化数据	有预设数据格式的数据
非结构化数据	无预先设定结构的数据
Whirr	一套用于运行云服务的库
YARN	Apache Hadoop 的下一代计算框架，除了 MapReduce 之外还支持编程范式

提示

本书提供配套的网上下载资源，包括预备知识内容、PowerPoint 幻灯片、模拟试题和其他附加资源（包括额外的面试题）。以上所有资源均为英文资料。^①

“知识检测点”和“测试你的能力”环节中的问题可能需要使用特定数据集。读者可以使用本书配套的网上下载资源中提供的数据集，也可以使用从网上找到的合适的数据或者自己生成数据。

^① 本书配套的网上下载资源请登录异步社区（www.epubit.com.cn），访问本书对应页面下载。——编者注

目 录

模块 1 大数据入门

第 1 讲 大数据简介	3	3.1.2 虚拟化及其对大数据的重要性	47
1.1 什么是大数据	4	3.2 Hadoop 简介	47
1.1.1 大数据的优势	5	3.3 云计算和大数据	50
1.1.2 挖掘各种大数据源	6	3.3.1 大数据计算的特性	50
1.2 数据管理的历史——大数据的演化	7	3.3.2 云部署模型	51
1.3 大数据的结构化	9	3.3.3 云交付模型	52
1.4 大数据要素	13	3.3.4 大数据云	52
1.4.1 数据量	13	3.3.5 大数据云市场中的供应商	53
1.4.2 速度	14	3.3.6 使用云服务所存在的问题	54
1.4.3 多样性	14	3.4 大数据内存计算技术	54
1.5 大数据在商务环境中的应用	14	练习	56
1.6 大数据行业中的职业机会	16	备忘单	58
1.6.1 职业机会	17	第 4 讲 了解 Hadoop 生态系统	59
1.6.2 所需技能	17	4.1 Hadoop 生态系统	60
1.6.3 大数据的未来	19	4.2 用 HDFS 存储数据	61
练习	20	4.2.1 HDFS 架构	62
备忘单	22	4.2.2 HDFS 的一些特殊功能	65
第 2 讲 大数据在商业上的应用	23	4.3 利用 Hadoop MapReduce 处理数据	65
2.1 社交网络数据的重要性	24	4.3.1 MapReduce 是如何工作的	66
2.2 金融欺诈和大数据	30	4.3.2 MapReduce 的优点和缺点	66
2.3 保险业的欺诈检测	32	4.3.3 利用 Hadoop YARN 管理资源和应用	67
2.4 在零售业中应用大数据	36	4.4 利用 HBase 存储数据	68
练习	40	4.5 使用 Hive 查询大型数据库	69
备忘单	42	4.6 与 Hadoop 生态系统的交互	70
第 3 讲 处理大数据的技术	43	4.6.1 Pig 和 Pig Latin	70
3.1 大数据的分布式和并行计算	44	4.6.2 Sqoop	71
3.1.1 并行计算技术	46		

4.6.3 Zookeeper	72	5.3.1 硬件/网络拓扑	85
4.6.4 Flume	72	5.3.2 同步	86
4.6.5 Oozie	73	5.3.3 文件系统	86
练习	74	5.4 MapReduce 的应用	86
备忘单	76	5.5 HBase 在大数据处理中的角色	87
第5讲 MapReduce 基础	77	5.6 利用 Hive 挖掘大数据	89
5.1 MapReduce 的起源	78	练习	91
5.2 MapReduce 是如何工作的	79	备忘单	94
5.3 MapReduce 作业的优化技术	85		

模块2 分析和R编程入门

第1讲 理解分析	97	2.3.2 点解决方案的大爆发	123
1.1 分析与报告的对比	98	2.3.3 数据可视化工具	125
1.1.1 报告	99	2.4 一些流行的分析工具	127
1.1.2 分析	100	2.4.1 用于统计计算的R项目	127
1.2 基本和高级分析	102	2.4.2 IBM SPSS	128
1.3 进行分析——需要考虑的事项	105	2.4.3 SAS	130
1.3.1 正确限定问题的范围	105	2.5 分析工具之间的对比	131
1.3.2 统计显著性还是业务重要性	105	练习	133
1.3.3 样本与总体	107	备忘单	135
1.3.4 推理与计算统计数字的对比	109	第3讲 探索R	136
1.4 构建分析团队	110	3.1 安装R	137
1.4.1 成为分析师的必备技能	110	3.2 使用脚本工作	138
1.4.2 IT与分析的融合	111	3.2.1 RGui	138
练习	113	3.2.2 RStudio	140
备忘单	115	3.2.3 “Hello world!”	141
第2讲 分析方法与工具	116	3.2.4 简单数学运算	141
2.1 分析方法的演变	117	3.2.5 R中的数学运算	142
2.1.1 集成方法	117	3.2.6 使用向量	143
2.1.2 商品化模型	118	3.2.7 保存和计算数值	144
2.1.3 文本分析	120	3.2.8 回应用户	146
2.1.4 文本分析的挑战	121	3.3 浏览工作区	149
2.2 分析工具的演变	122	3.3.1 操纵工作区内容	149
2.3 分析工具分类	123	3.3.2 保存工作	150
2.3.1 图形用户界面的兴起	123	3.3.3 检索工作	150
		练习	151

备忘单	153	5.1 确定最合适的数据结构	178
第4讲 将数据集读入R, 从R导出数据		5.2 创建数据的子集	179
数据	154	5.2.1 指定子集	179
4.1 使用 <code>c()</code> 命令创建数据	155	5.2.2 构造数据帧的子集	180
4.1.1 输入数值项作为数据	155	5.2.3 从数据中取得样本	180
4.1.2 输入文本项作为数据	156	5.2.4 数据子集的应用	182
4.2 在R中使用 <code>scan()</code> 命令获取数据	157	5.3 在数据中添加计算得到的字段	184
4.2.1 输入文本作为数据	158	5.3.1 在数据帧列上执行算术运算	184
4.2.2 使用剪贴板制作数据	158	5.3.2 创建数据子组或者 bin	184
4.2.3 从磁盘读取数据文件	160	5.4 在R中组合和合并数据集	186
4.3 读取更大的数据文件	162	5.4.1 创建样本数据以说明合并的方法	187
4.3.1 <code>read.csv()</code> 命令	163	5.4.2 使用 <code>merge()</code> 函数	188
4.3.2 在R中读取数据的其他命令	164	5.4.3 合并类型	189
4.3.3 数据文件中的缺失值	167	5.4.4 使用查找表	190
4.4 从R导出数据	169	5.5 分类和排序数据	190
4.5 在R中保存你的工作	169	5.5.1 向量的排序	191
4.5.1 将数据文件保存到磁盘	170	5.5.2 数据帧的排序	191
4.5.2 保存命名对象	170	5.5.3 用 <code>apply()</code> 函数遍历数据	193
4.5.3 保存所有操作	170	5.6 公式接口简介	196
4.5.4 以文本文件形式保存数据到磁盘	171	5.7 数据整形	196
4.5.5 将向量对象写入磁盘	171	5.7.1 理解长格式和宽格式数据	197
4.5.6 将矩阵和数据帧对象写入磁盘	172	5.7.2 从 <code>reshape2</code> 程序包入手	198
4.5.7 将列表对象写入磁盘	172	5.7.3 将数据“融化”为长格式	199
练习	174	练习	202
备忘单	176	备忘单	204
第5讲 在R中操纵和处理数据	177		

模块3 使用R进行数据分析

第1讲 使用R中的函数和包	207	1.2 巧妙地使用参数	214
1.1 从脚本到函数	209	1.2.1 增加更多参数	214
1.1.1 创建脚本	209	1.2.2 使用点参数	216
1.1.2 将脚本转变为函数	210	1.2.3 使用函数作为参数	218
1.1.3 使用函数	211	1.3 函数作用域	219
1.1.4 减少行数	212	1.3.1 外部函数	219
		1.3.2 使用内部函数	221

1.4 指派方法	222	3.1.1 矩阵	271
1.4.1 寻找函数背后的方法	223	3.1.2 列表	272
1.4.2 以 UseMethod()函数使用方法	223	3.1.3 数据帧——数据集	273
1.5 程序包	225	3.2 向量、矩阵和列表的索引	273
1.5.1 为 Windows 安装程序包	225	3.2.1 向量的索引	273
1.5.2 为 Linux 安装程序包	225	3.2.2 矩阵的索引	274
1.6 程序包的使用	227	3.2.3 列表的索引	275
1.6.1 加载程序包	227	3.3 R 编程	276
1.6.2 卸载程序包	227	3.3.1 表达式、赋值和算术运算符	276
练习	228	3.3.2 成组的表达式	277
备忘单	230	3.3.3 条件执行——if 和 ifelse	278
第 2 讲 R 中的描述性统计	231	3.3.4 重复执行——循环	278
2.1 汇总命令	232	3.4 RHadoop	280
2.2 名称命令	234	3.4.1 安装 RHadoop	281
2.3 汇总样本	235	3.4.2 创建用户定义函数	281
2.4 累积统计信息	239	练习	283
2.4.1 简单累计命令	239	备忘单	285
2.4.2 复杂累积命令	241	第 4 讲 R 中的图形分析	286
2.5 数据帧的汇总统计	242	4.1 为单变量绘图	287
2.5.1 数据帧的通用汇总命令	242	4.1.1 直方图	288
2.5.2 专用的行和列汇总命令	243	4.1.2 索引图	292
2.5.3 用于行/列汇总的 apply()命令	243	4.1.3 时间序列图	293
2.6 矩阵对象的汇总统计	244	4.1.4 饼图	294
2.7 列表的汇总统计	246	4.1.5 stripchart 函数	294
2.8 列联表	247	4.2 绘制双变量图表	295
2.8.1 建立列联表	247	4.2.1 根据两个连续解释变量绘制 图表：散点图	296
2.8.2 选择表对象的各个部分	253	4.2.2 使用分类解释变量绘图	309
2.8.3 测试表对象	255	4.3 多重比较图表	312
2.8.4 复杂（扁平）表	256	4.4 绘制多变量图表	315
2.8.5 测试“扁平”表对象	260	4.4.1 pairs 函数	315
2.8.6 表的汇总命令	260	4.4.2 coplot 函数	316
2.9 交叉表	262	4.4.3 相互作用图表	316
练习	267	4.5 特殊图表	317
备忘单	269	4.5.1 设计图	318
第 3 讲 用函数、循环和数据帧分析 数据	270	4.5.2 气泡图	318
3.1 矩阵、列表和数据帧	271	4.5.3 有许多相同值的图表	319

4.6 将图形保存到外部文件	320	5.3 u 检验	333
练习	322	5.3.1 双样本 u 检验	333
备忘单	324	5.3.2 单样本 u 检验	334
第5讲 R 中的假设检验	325	5.3.3 u 检验中的公式语法和样本子集 构建	335
5.1 统计假设简介	326	5.4 配对 t 检验和 u 检验	338
5.1.1 假设检验	327	5.4.1 相关和协方差	340
5.1.2 决策错误	327	5.4.2 协方差	342
5.2 使用学生 t 检验	327	5.4.3 相关检验中的显著性检验	343
5.2.1 使用不相等方差的双样本 t 检验	328	5.4.4 公式语法	343
5.2.2 使用相等方差的双样本 t 检验 ..	328	5.5 关联分析检验	346
5.2.3 单样本 t 检验	328	5.6 拟合优度检验	348
5.2.4 t 检验中的公式语法和样本子集 构建	329	练习	352
		备忘单	354
 模块4 使用 R 进行高级分析 			
第1讲 R 中的线性回归	357	1.3.2 校正的平方和及乘积和	372
1.1 线性回归分析基础知识	358	1.3.3 分散度	372
1.1.1 简单线性回归	358	1.3.4 回归中的方差分析	373
1.1.2 多重线性回归	359	1.3.5 AIC	373
1.1.3 最小二乘估计	360	1.3.6 参数不可靠性的估算	373
1.1.4 检查模型适当性	361	1.3.7 用拟合模型预测	374
1.1.5 回归输出的解读	363	1.3.8 检查模型	374
1.1.6 回归假设	364	1.4 线性模型结果对象	375
1.1.7 多重共线性	365	1.4.1 系数	377
1.1.8 检测多重共线性	365	1.4.2 拟合值	377
1.2 使用线性回归进行工作	367	1.4.3 残差	378
1.2.1 确定 x 和 y 变量	367	1.4.4 公式	378
1.2.2 检查条件	368	1.4.5 最佳拟合线	378
1.2.3 回归线的计算	368	1.5 模型的构建	379
1.2.4 求取斜率	369	1.5.1 用前向逐步回归增加项	380
1.2.5 求取 y 截距	369	1.5.2 用后向删除方法删除项	382
1.2.6 回归线的解读	369	1.5.3 模型的比较	383
1.2.7 做出正确的预测	371	1.6 曲线回归	384
1.3 R 中的简单线性回归	371	练习	386
1.3.1 R 的 5 个著名函数	371	备忘单	389

第2讲 非线性回归	390	3.2 凝聚层次聚类	425
2.1 非线性回归分析简介	391	3.2.1 主要距离	426
2.2 非线性回归和广义线性模型	391	3.2.2 密度估算方法	427
2.3 逻辑回归	392	3.3 相似性聚合聚类	428
2.3.1 解读逻辑回归中的 β 系数	394	3.3.1 相似性聚合的原理	428
2.3.2 计算 β 系数	395	3.3.2 相似性聚合聚类的实施	428
2.3.3 具有交互变量的逻辑回归	395	3.4 R amap包的用法	429
2.3.4 具有指示变量的逻辑回归	396	3.5 k 均值聚类	431
2.3.5 逻辑回归模型适当性检查	396	3.6 R 聚类示例：欧洲人的蛋白质 摄入	431
2.3.6 使用逻辑回归线进行预测	397	3.7 R 聚类示例：美国月度失业率	434
2.4 用MLE进行线估算	400	3.8 在R中实施层次聚类	435
2.5 将非线性模型转化为线性模型	401	3.8.1 例1：重温欧洲人蛋白质摄入	435
2.6 其他非线性回归模型	402	3.8.2 例2：重温美国月度失业率	436
2.7 广义加性模型	406	练习	437
2.8 自启动函数	407	备忘单	439
2.8.1 自启动Michaelis-Menten模型	407	第4讲 决策树	440
2.8.2 自启动渐近指数模型	408	4.1 决策树的应用	441
2.8.3 轮廓似然	409	4.2 决策树原理	444
2.8.4 自启动逻辑	409	4.2.1 选择变量——创建树的 第1步	444
2.8.5 自启动四参数逻辑	409	4.2.2 拆分标准	445
2.8.6 自启动Weibull增长函数	410	4.2.3 为节点分配数据——创建树的 第2步	447
2.8.7 自启动一阶房室函数	411	4.2.4 修剪——创建树的第3步	447
2.9 用拔靴法建立一个非线性回归 家族	411	4.3 构建决策树	448
2.10 逻辑回归的应用	413	4.3.1 决策树如何确定纯度?	449
2.10.1 贷款接纳	414	4.3.2 使用决策树时的实际考虑 因素	450
2.10.2 德国信用数据	414	4.3.3 决策树选项	451
2.10.3 延误的航班	415	4.4 CART、C5.0和CHAID树	451
练习	416	4.4.1 CART	452
备忘单	418	4.4.2 C5.0	454
第3讲 聚类分析	419	4.4.3 CHAID	455
3.1 聚类简介	421	4.4.4 决策树对比	456
3.1.1 聚类的应用	421	4.5 用决策树预测	457
3.1.2 聚类的复杂性	422	4.6 决策树的优缺点	458
3.1.3 距离计量	422		
3.1.4 簇内和簇间平方和	423		
3.1.5 高效聚类的属性	424		