

HZ BOOKS
华章教育

WILEY

统计学精品译丛

(原书第5版)

例解回归分析

Regression Analysis by Example

(Fifth Edition)



(美) Samprit Chatterjee 著
Ali S. Hadi

郑忠国 许静 译



机械工业出版社
China Machine Press

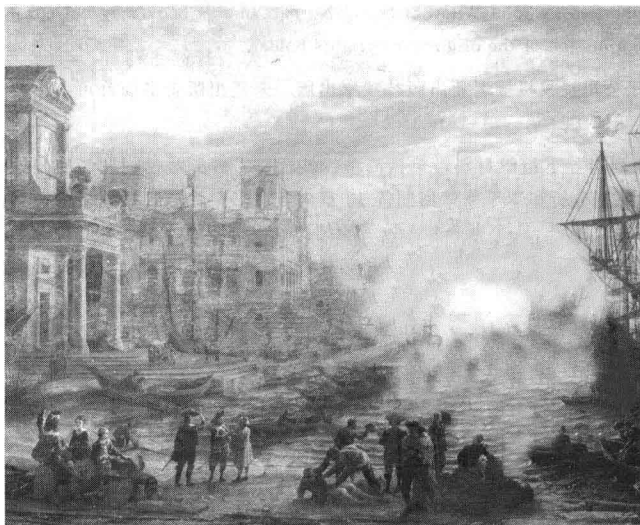
统计

(原书第5版)

例解回归分析

Regression Analysis by Example

(Fifth Edition)



(美) Samprit Chatterjee 著
Ali S.Hadi

郑忠国 许静 译



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

例解回归分析 (原书第 5 版) / (美) 查特吉 (Chatterjee, S.), (美) 哈迪 (Hadi, A. S.) 著; 郑忠国, 许静译. —北京: 机械工业出版社, 2013. 8

(统计学精品译丛)

书名原文: Regression Analysis by Example, Fifth Edition

ISBN 978-7-111-43156-5

I. 例… II. ①查… ②哈… ③郑… ④许… III. 回归分析 IV. O212.1

中国版本图书馆 CIP 数据核字 (2013) 第 145908 号

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问 北京市展达律师事务所

本书版权登记号: 图字: 01-2013-2602

All Rights Reserved. This translation published under license. Authorized translation from the English language edition, entitled *Regression Analysis by Example, Fifth Edition*, ISBN 978-0-470-90584-5, by Samprit Chatterjee and Ali S. Hadi, Published by John Wiley & Sons. No part of this book may be reproduced in any form without the written permission of the original copyrights holder.

本书中文简体字版由约翰-威利父子公司授权机械工业出版社独家出版。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。

本书在探索性数据分析的思想和原则指导下组织材料, 包括简单线性回归、多元线性回归、回归诊断、定性预测变量、变量变换、共线性数据分析和逻辑斯谛回归等 13 章内容。书中强调数据分析的技巧而不是统计理论的发展, 几乎是手把手地教读者如何去分析数据、检验结论、改进分析。作者精心挑选了丰富的实例, 形象生动而又系统详尽地阐述了回归分析的基本理论和具体的应用技术, 还辅以启发式的推理和直观的图形方法。

本书既可以作为非统计学专业回归分析的入门教材, 又可以作为统计学专业理论回归分析的补充教材, 对于从事数据分析的人员来说, 本书更是必备的参考书。

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: 明永玲

三河市杨庄长鸣印刷装订厂印刷

2013 年 8 月第 1 版第 1 次印刷

186mm×240mm·19.25 印张

标准书号: ISBN 978-7-111-43156-5

定 价: 69.00 元

凡购本书, 如有缺页、倒页、脱页, 由本社发行部调换

客服热线: (010) 88378991 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzsj@hzbook.com

中文版序

听闻《例解回归分析（第5版）》的中文版即将出版，我们非常高兴。我们相信，对中国的学生来说，中文版更容易理解，价格也更优惠。中国的学生是知名的高智商群体，你们正努力掌握先进的数据分析方法，解决科学、技术和社会科学中的实际问题，为你们的祖国的不断进步和繁荣做出贡献。希望这本书能对你们有所启发，对你们的应用工作有所帮助。

在此，我们特别感谢郑忠国和许静两位学者，他们承担了本书的翻译工作。在翻译期间，他们仔细地阅读了全文，发现了原书的一些错误，在中文版中得以改正。我们也希望获得中国读者关于本书的反馈意见。

A Note to Our Chinese Readers

We are delighted that the book is being published in Chinese, because it will become more accessible to a group of Chinese students who are highly intelligent, and want to contribute to their countries progress by devising better methods to analyze their data and contribute to the countries growth, prosperity, and process improvements. We hope they find the book stimulating, and helpful in their applied work.

We are thankful to Dr. Zhongguo Zheng and Jing Xu for translating the book. During the translation process they read the book very carefully and has discovered a few errors and corrected them in the Chinese version. We would be happy to hear feedback from our Chinese readers.

Samprit Chatterjee

Ali S. Hadi

译者序

回归分析是统计学的一个重大分支学科。无论从理论研究方面还是统计应用领域来说，回归分析都是国人熟悉的课题。国内已经有不少介绍回归分析的书籍和教材，但是，本书有它的独特之处。

1. 它没有把系统叙述回归分析的定义、模型和理论作为起始点和主要目的，也不强调读者的数学基础和逻辑推断能力。这就使得本书具有广泛的读者范围，即无论是具有较深理论基础的专业统计工作者还是需要利用回归分析作为数据分析工具的实际工作者，阅读本书都不会产生困难，也不会由于过多数据或逻辑推理而心里烦躁，而是会受到探索数据内在规律的启发。

2. 作者以数据实例分析贯穿始终，读者往往被数据中隐藏的关于事物本质的谜底所吸引而不感觉枯燥。推理往往是启发式的，有时候用直观的图形方法，这反映了探索性数据分析的特点。它不是对假设模型的理论进行推断，而是不设前提地对隐藏在事物背后的规律进行探索。这种方法与传统上用例子说明理论结果的目的是不同的。

3. 本书不将读者的计算机能力作为阅读本书的必备条件，但是具备这个条件将如虎添翼。对于现代的学者特别是年轻学生来说，具备计算机能力不是一个苛刻的要求。

4. 本书作者精心安排数据例子，使得读者在读完本书以后就可以系统地掌握回归分析的技巧和方法。对于那些热衷于回归分析方法的理论根源的读者，作者也提供了相关的参考文献，以便深化对回归分析的认识。此外，在某些章节还增加了附录，扩展介绍了一些方法，例如第 10 章介绍代理岭回归的概念，这是近年研究的内容。当然，作为教材，和大部分教科书一样，本书也提供了丰富的习题以供巩固所学。

基于以上特点，我们乐于向读者推荐本书，并建议大学教师将本书作为“回归分析”课程的教材，尝试一种新的教学方法。我们翻译本书的过程也是一个学习的过程、享受的过程。在翻译的过程中我们得到了作者的帮助，受益匪浅，在此特向作者表示感谢。

译者

2013 年 5 月 4 日

前 言

我们很高兴在此把《例解回归分析》第5版介绍给大家，本书初版于1977年。统计界一直对此书十分关心和支持，我们也从他们对本书的诸多改进意见中获益良多。

对于分析多因素数据资料，回归分析已经成为应用最为广泛的统计分析工具之一。它之所以广受欢迎，是因为它对分析变量之间的函数关系提供了概念上简单明了的方法。回归分析的标准方法是：对数据拟合一个模型，然后利用诸如 t 、 F 和 R^2 等统计量对拟合的方程进行评估。本书的方法比这些传统的方法更加广泛。我们将回归分析看成考察各个变量之间关系的一种数据分析的工具。本书并不强调形式化的统计检验和概率计算。我们的目标是挖掘数据内在的结构。

我们在这些数据直观表现的基础上，进行大量传统的和一些不那么传统的统计分析。我们主要依靠这些数据的图形表示，经常利用许多种类的回归残差图进行分析。我们不强调精确的概率计算利用残差图的图形方法可以展现模型的缺陷，找出某些病态的观测值。进一步追溯这些病态观测值，通常会发现它们有时候比正常的观测值更具信息价值。我们发现，快速一瞥残差图比形式化地进行某个限定的原假设的显著性检验能获取更多的信息。可以这么说，本书是在探索性数据分析的思想和原则指导下写成的。

我们通过精心设计的例子来解释和展现回归分析的各种基本概念和方法。每个例子中，我们总是集中介绍一两种回归分析技术。因此在选择数据的时候，我们仔细琢磨，精心挑选，以便突出我们所介绍的技术。在实际工作中，对于一个数据集合，通常要涉及许多不同的分析技术。但是本书例子的安排，使得分析数据时各种分析技巧有序出场，不需要在不同的例子中重复地介绍和解释同一个分析技巧。我们希望读者在学完本书以后，能够系统地掌握回归分析的各种技巧，并且能够融会贯通地处理所遇到的数据分析问题。

本书强调的是分析数据的技术，而不是统计公式、假设检验和置信区间。因此，我们的重点不在于这些分析技术的推导。当然，我们在分析数据时会介绍这些分析工具，并且给出它们的使用条件，最后，在具体的例子中给出使用效果的评价。虽然我们没有给出这些分析技术的推导，但是我们会给出这些技术的来源，有兴趣的读者可以参考并进一步钻研其理论。

我们假定读者能够接触到计算机和统计软件。现在，线性模型分析领域有了质的飞跃，从模型拟合到建模、从一般的检验到临床数据的检测、从宏观分析到微观分析，所有这些都需要计算机，因此我们假定大家手头具备这一工具。几乎所有我们用的分析工具，现有软件包里都能找到。特别是，在互联网上可以找到软件包 R，这个软件包具有很强的计算能力和图形功能。同时，它是免费的！

本书的读者对象是涉及分析数据的各层次人员。这本书对于具有统计基本知识的人员是颇有帮助的。在大学中，它可以作为“回归分析”课程的教材，课程的授课对象是非统计专业的学生，但是这些学生在他们的专业领域内又特别需要回归分析这个数据处理工具。对于统计专业的学生，如果他们修过“回归分析”这门课程，而课程的水平如 Rao

(1973)、Seber(1977) 或 Sen and Srivastava(1990) 那样, 那么本书的内容是他们所学的理论回归分析的补充, 从实际应用的角度去深化他们对回归分析的认识. 在大学以外, 对于那些应用标准统计方法 (如 t 、 F 、 R^2 、标准误等) 进行回归分析解决实际问题的, 如果他们要对多因素数据进行更加深入的分析, 那么这本书是非常有用的.

本书的配套网站是: <http://www.aucegypt.edu/faculty/hadi/RABE5>. 该网站包含本书的所有数据, 当然还有一些其他数据和内容.

本书的第 5 版语言更加流畅, 去掉一些模棱两可的说法, 纠正了一些错误, 这些错误是由读者指出的或由作者自己纠正的. 在第 1 章中加入了新的数据集的例子. 将第 4 版第 9 章中关于数据的中心化和规范化的材料移到 3.6 节. 第 9 章和第 10 章的材料经过重新组织, 使得概念上循序渐进, 学习起来更加通俗易懂. 第 10 章的附录简单描述了代理岭回归, 这是近年提出的新的研究内容. 第 5 版中还增加了新的参考文献. 在每一章的最后, 我们增加了习题, 对某些习题还加以改写. 我们认为, 做习题能够巩固和加强对前面所学内容的理解.

我们努力让更多人获益, 因此本书的读者对象是来自各种不同领域的数据处理的工作者. 本书强调的重点是数据分析的技巧, 而不是统计理论的发展.

我们很幸运地得到多位朋友、同事及合作者的鼓励和帮助. 在纽约大学和康奈尔大学的几个同行将本书的部分材料作为他们课程的教材, 并且将他们的评论以及学生的意见提供给我们. 特别要提到的是我们的朋友兼前同事 Jeffrey Simonoff (纽约大学), 他给出很好的审稿意见, 提出建议并给予很多其他帮助. 我们的“回归分析”课程上, 许多学生也对本书做出了贡献, 他们提出了许多深刻的问题, 还要求有意义且可以理解的答案. 我们也要特别感谢 Nedret Billor (Cukurova 大学, Turkey) 和 Sahar El-Sheneity (康奈尔大学), 他们仔细阅读了本书的早期版本. 同样, Amy Hendrickson 为本书准备了 Latex 文件并回答了有关 Latex 的问题, Dean Gonzalez 协助制作某些图形, 在此一并表示感谢.

Samprit Chatterjee
Ali S. Hadi
Brooksville, Maine
Cairo, Egypt

目 录

中文版序		
译者序		
前言		
第 1 章 概述	1	
1.1 什么是回归分析	1	
1.2 公用数据集	1	
1.3 回归分析应用实例选讲	2	
1.3.1 农业科学	2	
1.3.2 劳资关系	3	
1.3.3 政府	5	
1.3.4 历史	8	
1.3.5 环境科学	8	
1.3.6 工业生产	9	
1.3.7 挑战者号航天飞机	11	
1.3.8 医疗费用	12	
1.4 回归分析的步骤	14	
1.4.1 问题陈述	14	
1.4.2 选择相关变量	15	
1.4.3 收集数据	15	
1.4.4 模型设定	16	
1.4.5 拟合方法	17	
1.4.6 模型拟合	18	
1.4.7 模型评价和选择	18	
1.4.8 回归分析的目标	19	
1.5 本书的内容和结构	20	
习题	21	
第 2 章 简单线性回归	22	
2.1 引言	22	
2.2 协方差与相关系数	22	
2.3 实例：计算机维修数据	26	
2.4 简单线性回归模型	27	
2.5 参数估计	28	
2.6 假设检验	30	
2.7 置信区间	34	
2.8 预测	34	
2.9 拟合效果度量	35	
2.10 过原点的回归直线	38	
2.11 平凡的回归模型	39	
2.12 文献	40	
习题	40	
第 3 章 多元线性回归	45	
3.1 引言	45	
3.2 数据和模型的描述	45	
3.3 实例：主管人员业绩数据	46	
3.4 参数估计	47	
3.5 回归系数的解释	48	
3.6 中心化和规范化	50	
3.6.1 含截距模型的中心化和规范化	50	
3.6.2 无截距模型的规范化	51	
3.7 最小二乘估计的性质	52	
3.8 复相关系数	53	
3.9 单个回归系数的推断	54	
3.10 线性模型中的假设检验	55	
3.10.1 检验所有预测变量的回归系数为 0	56	
3.10.2 检验某些回归系数为 0	58	
3.10.3 检验某些回归系数相等	60	
3.10.4 带约束的回归参数的估计和检验	61	
3.11 预测	62	
3.12 小结	63	
习题	63	
附录 多元回归的矩阵表示	69	

第 4 章 回归诊断：违背模型假定的 检测	71	5.4.2 斜率相同但截距不同的模型 ...	107
4.1 引言	71	5.4.3 截距相同但斜率不同的模型 ...	108
4.2 标准回归假定	71	5.5 示性变量的其他应用	109
4.3 各种残差	72	5.6 季节性	109
4.4 图形方法	74	5.7 回归参数随时间的稳定性	111
4.5 拟合模型前的图形	76	习题	115
4.5.1 一维图	76	第 6 章 变量变换	121
4.5.2 二维图	77	6.1 引言	121
4.5.3 旋转图	78	6.2 线性化变换	122
4.5.4 动态图	78	6.3 X 射线灭菌	124
4.6 拟合模型后的图形	79	6.3.1 线性模型的不适用性	125
4.7 检查线性和正态性假定的图形	79	6.3.2 对数变换实现线性化	125
4.8 杠杆、强影响点和异常值	80	6.4 稳定方差的变换	126
4.8.1 响应变量的异常值	81	6.5 异方差误差的检测	130
4.8.2 预测变量中的异常值	81	6.6 消除异方差性	131
4.8.3 伪装和淹没问题	82	6.7 加权最小二乘	132
4.9 观测影响的度量	83	6.8 数据的对数变换	132
4.9.1 Cook 距离	84	6.9 幂变换	134
4.9.2 Welsch-Kuh 度量	84	6.10 总结	137
4.9.3 Hadi 影响度量	85	习题	137
4.10 位势-残差图	86	第 7 章 加权最小二乘法	141
4.11 如何处理异常点	87	7.1 引言	141
4.12 回归方程中变量的作用	88	7.2 异方差模型	142
4.12.1 添加变量图	88	7.2.1 主管人员数据	142
4.12.2 残差加分量图	88	7.2.2 大学教育花费数据	143
4.13 添加一个预测变量的效应	92	7.3 两阶段估计	144
4.14 稳健回归	92	7.4 教育费用数据	145
习题	93	7.5 拟合剂量-反应关系曲线	151
第 5 章 定性预测变量	97	习题	152
5.1 引言	97	第 8 章 相关误差问题	153
5.2 薪水调查数据	97	8.1 引言：自相关	153
5.3 交互变量	100	8.2 消费支出和货币存量	153
5.4 回归方程组：两个组的比较	102	8.3 Durbin-Watson 统计量	155
5.4.1 斜率和截距都不同的模型	103	8.4 利用变换消除自相关性	157
		8.5 当回归模型具有自相关误差时的	

迭代估计法	158	附录 10. B 岭回归	216
8.6 变量的缺失和模型的自相关性	159	附录 10. C 代理岭回归	218
8.7 住房开工规模的分析	160	第 11 章 变量选择	219
8.8 Durbin-Watson 统计量的局限性	162	11.1 引言	219
8.9 用示性变量消除季节效应	164	11.2 问题的陈述	219
8.10 两个时间序列之间的回归	166	11.3 删除变量的后果	220
习题	167	11.4 回归方程的用途	221
第 9 章 共线性数据分析	171	11.4.1 描述和建模	221
9.1 引言	171	11.4.2 估计和预测	221
9.2 共线性对推断的影响	172	11.4.3 控制	221
9.3 共线性对预测的影响	176	11.5 评价回归方程的准则	222
9.4 共线性的检测	178	11.5.1 残差均方	222
9.4.1 共线性的简单征兆	179	11.5.2 Mallows 的 C_p 准则	223
9.4.2 方差膨胀因子	182	11.5.3 信息准则	223
9.4.3 条件指数	184	11.6 共线性和变量选择	224
习题	186	11.7 评价所有可能的回归模型	225
第 10 章 共线性数据的处理	189	11.8 变量选择方法	225
10.1 引言	189	11.8.1 前向选择方法	226
10.2 主成分	189	11.8.2 后向剔除方法	226
10.3 利用主成分的计算	192	11.8.3 逐步回归法	226
10.4 施加约束条件	194	11.9 变量选择的一般注意事项	227
10.5 搜索模型中回归系数的		11.10 对主管人员业绩的研究	227
线性函数	195	11.11 共线性数据的变量选择	231
10.6 回归系数的有偏估计	198	11.12 凶杀数据	231
10.7 主成分回归	199	11.13 利用岭回归进行变量选择	234
10.8 消除数据中的共线性	200	11.14 空气污染研究中的变量选择	234
10.9 回归系数的约束条件	202	11.15 拟合回归模型的可能策略	243
10.10 主成分回归中的注意事项	203	11.16 文献	244
10.11 岭回归	205	习题	244
10.12 岭估计法	206	附录 误设模型的影响	247
10.13 岭回归: 几点注解	209	第 12 章 逻辑斯谛回归	249
10.14 小结	210	12.1 引言	249
10.15 文献	210	12.2 定性数据的建模	249
习题	211	12.3 Logit 模型	250
附录 10. A 主成分	214		

12.4	例子：破产概率的估计	251	第 13 章	进一步的论题	268
12.5	逻辑斯谛回归模型诊断	254	13.1	引言	268
12.6	决定变量的去留	255	13.2	广义线性模型	268
12.7	逻辑斯谛回归的拟合度	257	13.3	泊松回归模型	269
12.8	多项 Logit 模型	258	13.4	引进新药	269
12.8.1	多项逻辑斯谛回归	259	13.5	稳健回归	270
12.8.2	例子：确定化学糖尿病	259	13.6	拟合一个二次式模型	271
12.8.3	顺序值逻辑斯谛回归	263	13.7	美国海湾中 PCB 的分布	272
12.8.4	例子：重新考察化学糖尿病的 确定问题	264	习题		275
12.9	分类问题：另一种方法	264	附录 A	统计表	276
习题		266	参考文献		283
			索引		291

第1章 概述

1.1 什么是回归分析

回归分析的基本思想很简单,是研究变量间函数关系的一种方法.一个房地产评估师可能会将一栋房屋的销售价格和与之有关的因素联系起来,例如该建筑物的主要结构特征以及须支付的税费(地方税、教育税、国家税)等.我们也许想知道香烟的消费量是否与各种社会经济变量和人口统计学变量有关,例如年龄、教育程度、收入和香烟的价格等.变量之间的这种关系可以表示为方程或模型的形式,该方程或模型将响应变量或因变量与一个或多个解释变量或预测变量联系起来.在香烟消费的例子中,响应变量是香烟的消费量(用某州一年内人均购买的香烟包数计算),解释变量或预测变量是一些社会经济变量和人口统计学变量.在房地产评估的例子中,响应变量是房屋的价格,解释变量或预测变量是房屋的结构特征和拥有该房屋所承担的税费.

我们用 Y 表示响应变量,用 X_1, X_2, \dots, X_p 表示预测变量,其中 p 是预测变量的个数, Y 和 X_1, X_2, \dots, X_p 之间的真实关系可近似地用下述回归模型刻画

$$Y = f(X_1, X_2, \dots, X_p) + \epsilon, \quad (1.1)$$

其中 ϵ 是随机误差,它代表在近似过程中产生的偏差,也就是模型不能精确拟合数据的原因.函数 $f(X_1, X_2, \dots, X_p)$ 刻画了 Y 和 X_1, X_2, \dots, X_p 之间的关系,最简单的情形是线性回归模型

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon, \quad (1.2)$$

其中 $\beta_1, \beta_2, \dots, \beta_p$ 称为回归系数, β_0 称为截距,它们都是未知常数,称为模型的回归参数,这些未知参数可由数据确定(估计).通常,我们用希腊字母表示未知参数.

解释变量或预测变量也可以称为独立变量、协变量、回归变量或因素.虽然我们经常使用独立变量这种名称,但由于实际中的预测变量之间很少是相互独立的,所以这个名称并不贴切.

1.2 公用数据集

回归分析有广泛的应用领域,包括经济、金融、商业、法律、气象、医学、生物、化学、工程、物理、教育、体育、历史、社会学和心理学等.在1.3节将给出一些应用例子.回归分析是帮助读者分析数据最有效的方法.各位读者可以考虑一下在你们工作、研究或感兴趣的问题中有哪些可以用回归分析方法加以解决.当然,在进行回归分析之前,首先要收集相关的数据,然后将本书介绍的回归分析方法应用于这些数据.为了便于读者查找

⊖ 此页码为英文原书页码,与索引页码一致.

实际数据，本节给出大量公用数据集的一些资源链接。

一部分数据集可以从书籍和互联网获得。Hand et al. (1994) 所著的书中包含了许多领域的数据集，这些数据集的容量不大，适合练习使用。Chatterjee, Handcock and Simonoff (1995) 的书中提供了来自不同领域的大量数据集，这些数据集保存在随书所附的光盘中，也可以从相关网站上获得^①。

数据集也可在互联网的其他很多网站上获得，下面给出的一些网站允许直接复制和粘贴数据到统计软件包中，而其他的网站则需要下载数据文件，然后导入统计软件包中。部分网站还进一步提供了与其他数据集或统计相关网站的链接。

数据和故事图书馆 (Data and Story Library, DASL, 读作 “dazzle”) 是最有意思的一个网站，不但有很多数据集，而且还介绍了每一个数据集的“故事”或背景。DASL 是一个介绍基本统计方法应用的数据文件和背景的在线资料库^②，其中的数据集涵盖了广泛的研究领域。DASL 是一个查找感兴趣的数据和背景资料的强大的搜索引擎。

另一个提供数据集的网站是电子数据集服务 (Electronic Dataset Service)^③，其中的数据集按所用的分析方法组织安排。该网站同样提供了许多其他数据资源的网络链接。

最后，本书有一个网站^④，除了包含本书所有的数据集之外，还有很多其他的资料。书中的数据集和其他数据集都可在该网站获得。

1.3 回归分析应用实例选讲

回归分析提供了建立变量间函数关系的简便方法，是应用最广泛的统计工具之一，已经广泛应用于很多学科领域。前面提到的香烟消费问题和房地产评估问题就是其中两例。本节将给出一些其他的例子，以说明回归分析在现实生活中的广泛应用。这里用到的一些数据集在以后还会用来介绍回归方法或出现在各章末尾的习题中。

1.3.1 农业科学

在纽约州北部地区的奶牛改进合作组织 (Dairy Herd Improvement Cooperative, DHI) 收集并分析牛奶产量数据。我们感兴趣的问题是如何建立一个合适的模型，通过一些可测量的变量值预测牛奶产量。表 1-1 给出了响应变量 (以磅计量的本月的牛奶产量) 和预测变量。每个月抽取一次产奶的样本。母牛产奶的时期称为产奶期，产奶期数是采样时该母牛经历的产犊 (或产奶时期) 的次数。一种推荐的管理方法是，让奶牛产奶约 305 天，然后休息 60 天，再开始下一个产奶期。这个数据集有 199 个观测值，来自于 DHI 的牛奶产量记录，该牛奶产量数据可在本书的网站上获得。

① <http://www.stern.nyu.edu/~jsirnono/Casebook>

② <http://lib.stat.cmu.edu/DASL>

③ <http://www-unix.oit.umass.edu/~statdata>

④ <http://www.aucegypt.edu/faculty/hadi/RABE5>

表 1-1 牛奶产量数据中的变量

变 量	定 义
Current	本月牛奶产量（单位：磅）
Previous	前一个月牛奶产量（单位：磅）
Fat	牛奶中的脂肪百分比
Protein	牛奶中的蛋白质百分比
Days	自本次产奶期开始至今的总天数
Lactation	产奶期数
179	示性变量（Days≤79 时，值为 0；Days>79 时，值为 1）

1.3.2 劳资关系

1974 年，美国国会通过了针对瓦格纳法案的塔夫脱-哈特利修正案。最初的瓦格纳法案允许工会使用一种闭门合同[⊖]，除非州法律禁止这种做法。塔夫脱-哈特利修正案宣布闭门合同是不合法的，并且赋予各州进一步禁止以加入工会作为雇佣条件[⊕]的权利。这些劳动权利法已经在劳工运动中引起了不小的关注。我们感兴趣的问题是：这些法律对于美国一个中等收入的四口之家的生活支出有什么影响。要回答这个问题，研究者从各种渠道收集了由 38 个地区的信息构成的数据集。表 1-2 给出了数据集中的各个变量。表 1-3 给出了劳动权利法数据，该数据也可在本书的网站上获得。

4

表 1-2 工作权利法数据中的变量

变 量	定 义
COL	一个四口之家的生活支出
PD	人口密度（每平方英里的人数）
URate	1978 年州工会入会率
Pop	1975 年州人口数
Taxes	1972 年的物业税
Income	1974 年的人均收入
RTWL	示性变量（该州执行工作权利法，值为 1；否则，值为 0）

表 1-3 工作权利法数据

城市	COL	PD	URate	Pop	Taxes	Income	RTWL
Atlanta	169	414	13.6	1 790 128	5 128	2 961	1
Austin	143	239	11	396 891	4 303	1 711	1
Bakersfield	339	43	23.7	349 874	4 166	2 122	0

⊖ 闭门合同规定：所有雇员在被雇佣期间必须是工会会员，并且受雇的前提条件是必须保持会员资格。

⊕ 加入工会作为雇佣条件是指：雇员在被雇佣期间可以不是工会会员，但必须在两个月内成为会员，才能使雇主在雇佣决定中有完整的自由裁量权。

(续)

城市	COL	PD	URate	Pop	Taxes	Income	RTWL
Baltimore	173	951	21	2 147 850	5 001	4 654	0
Baton Rouge	99	255	16	411 725	3 965	1 620	1
Boston	363	1 257	24.4	3 914 071	4 928	5 634	0
Buffalo	253	834	39.2	1 326 848	4 471	7 213	0
Champaign-Urbana	117	162	31.5	162 304	4 813	5 535	0
Cedar Rapids	294	229	18.2	164 145	4 839	7 224	1
Chicago	291	1 886	31.5	7 015 251	5 408	6 113	0
Cincinnati	170	643	29.5	1 381 196	4 637	4 806	0
Cleveland	239	1 295	29.5	1 966 725	5 138	6 432	0
Dallas	174	302	11	2 527 224	4 923	2 363	1
Dayton	183	489	29.5	835 708	4 787	5 606	0
Denver	227	304	15.2	1 413 318	5 386	5 982	0
Detriot	255	1 130	34.6	4 424 382	5 246	6 275	0
Green Bay	249	323	27.8	169 467	4 289	8 214	0
Hartford	326	696	21.9	1 062 565	5 134	6 235	0
Houston	194	337	11	2 286 247	5 084	1 278	1
Indianapolis	251	371	29.3	1 138 753	4 837	5 699	0
Kansas City	201	386	30	1 290 110	5 052	4 868	0
Lancaster, PA	124	362	34.2	342 797	4 377	5 205	0
Los Angeles	340	1 717	23.7	6 986 898	5 281	1 349	0
Milwaukee	328	968	27.8	1 409 363	5 176	7 635	0
Minneapolis, St. Paul	265	433	24.4	2 010 841	5 206	8 392	0
Nashville	120	183	17.7	748 493	4 454	3 578	1
New York	323	6 908	39.2	9 561 089	5 260	4 862	0
Orlando	117	230	11.7	582 664	4 613	782	1
Philadelphia	182	1 353	34.2	4 807 001	4 877	5 144	0
Pittsburgh	169	762	34.2	2 322 224	4 677	5 987	0
Portland	267	201	23.1	228 417	4 123	7 511	0
St. Louis	184	480	30	2 366 542	4 721	4 809	0
San Diego	256	372	23.7	1 584 583	4 837	1 458	0
San Francisco	381	1 266	23.7	3 140 306	5 940	3 015	0
Seattle	195	333	33.1	1 406 746	5 416	4 424	0
Washington	205	1 073	21	3 021 801	6 404	4 224	0
Wichita	206	157	12.8	384 920	4 796	4 620	1
Raleigh-Durham	126	302	6.5	468 512	4 614	3 393	1

1.3.3 政府

一个国家的人从一个州或地区迁移至另一个州或地区称为国内移民，国内移民信息对于州和地区政府来说是很重要的。我们希望建立模型预测国内移民的情况，并且研究为什么人们会离开一个地方去另一个地方。许多因素会影响国内移民，例如天气情况、犯罪、税收和失业率等。我们建立了包含美国本土 48 个州的数据集，而阿拉斯加和夏威夷不在分析之列，因为这两州的环境明显不同于其他 48 个州，而且其地理位置也阻碍了外来移民。这里的响应变量是国内净移民率，指的是 1990—1994 年间，移入和移出某州的人数之差除以该州的总人数（以百分数表示）。表 1-4 给出了影响国内移民率的 11 个预测变量。表 1-5 和表 1-6 给出了国内移民数据，该数据也可在本书的网站上获得。

表 1-4 国内移民数据中的变量

变 量	定 义
State	州名
NDIR	1990—1994 年的国内净移民率
Unemp	1994 年民用劳动力的失业率
Wage	1994 年制造业工人的平均时薪
Crime	1993 年每十万人中的暴力犯罪率
Income	1994 年家庭收入的中位数
Metrop	1992 年生活在大都市地区的州人口百分比
Poor	1994 年生活在贫困线以下人口的百分比
Taxes	1993 年人均州税和地方税总额
Educ	1990 年 25 岁及以上人口中受高中及以上教育的人口百分比
BusFail	1993 年破产企业的数量除以该州人口数
Temp	1993 年该州 12 个月平均温度的平均值（华氏度）
Region	该州的地理区位（东北部、南部、中西部和西部）

6

表 1-5 国内移民数据中的前 6 个变量

州	NDIR	Unemp	Wage	Crime	Income	Metrop
Alabama	17.47	6.0	10.75	780	27 196	67.4
Arizona	49.60	6.4	11.17	715	31 293	84.7
Arkansas	23.62	5.3	9.65	593	25 565	44.7
California	-37.21	8.6	12.44	1 078	35 331	96.7
Colorado	53.17	4.2	12.27	567	37 833	81.8
Connecticut	-38.41	5.6	13.53	456	41 097	95.7
Delaware	22.43	4.9	13.90	686	35 873	82.7
Florida	39.73	6.6	9.97	1 206	29 294	93.0
Georgia	39.24	5.2	10.35	723	31 467	67.7
Idaho	71.41	5.6	11.88	282	31 536	30.0

(续)

州	NDIR	Unemp	Wage	Crime	Income	Metrop
Illinois	-20.87	5.7	12.26	960	35 081	84.0
Indiana	9.04	4.9	13.56	489	27 858	71.6
Iowa	0.00	3.7	12.47	326	33 079	43.8
Kansas	-1.25	5.3	12.14	469	28 322	54.6
Kentucky	13.44	5.4	11.82	463	26 595	48.5
Louisiana	-13.94	8.0	13.13	1 062	25 676	75.0
Maine	-9.770	7.4	11.68	126	30 316	35.7
Maryland	-1.55	5.1	13.15	998	39 198	92.8
Massachusetts	-30.46	6.0	12.59	805	40 500	96.2
Michigan	-13.19	5.9	16.13	792	35 284	82.7
Minnesota	9.46	4.0	12.60	327	33 644	69.3
Mississippi	5.33	6.6	9.40	434	25 400	34.6
Missouri	6.97	4.9	11.78	744	30 190	68.3
Montana	41.50	5.1	12.50	178	27 631	24.0
Nebraska	-0.62	2.9	10.94	339	31 794	50.6
Nevada	128.52	6.2	11.83	875	35 871	84.8
New Hampshire	-8.72	4.6	11.73	138	35 245	59.4
New Jersey	-24.90	6.8	13.38	627	42 280	100.0
New Mexico	29.05	6.3	10.14	930	26 905	56.0
New York	-45.46	6.9	12.19	1 074	31 899	91.7
North Carolina	29.46	4.4	10.19	679	30 114	66.3
North Dakota	-26.47	3.9	10.19	82	28 278	41.6
Ohio	-3.27	5.5	14.38	504	31 855	81.3
Oklahoma	7.37	5.8	11.41	635	26 991	60.1
Oregon	49.63	5.4	12.31	503	31 456	70.0
Pennsylvania	-4.30	6.2	12.49	418	32 066	84.8
Rhode Island	-35.32	7.1	10.35	402	31 928	93.6
South Carolina	11.88	6.3	9.99	1 023	29 846	69.8
South Dakota	13.71	3.3	9.19	208	29 733	32.6
Tennessee	32.11	4.8	10.51	766	28 639	67.7
Texas	13.00	6.4	11.14	762	30 775	83.9
Utah	31.25	3.7	11.26	301	35 716	77.5
Vermont	3.94	4.7	11.54	114	35 802	27.0
Virginia	6.94	4.9	11.25	372	37 647	77.5
Washington	44.66	6.4	14.42	515	33 533	83.0
West Virginia	10.75	8.9	12.60	208	23 564	41.8
Wisconsin	11.73	4.7	12.41	264	35 388	68.1
Wyoming	11.95	5.3	11.81	286	33 140	29.7