

统计应用与实战系列

Yingyong Tongji Fenxi yu
R Yuyan Shizhan

应用统计分析 与R语言实战

吕书龙 梁飞豹 / 主 编

刘文丽 薛美玉 / 副主编



北京大学出版社
PEKING UNIVERSITY PRESS

统计应用与实战系列

Yingyong Tongji Fenxi yu
R Yuyan Shizhan

应用统计分析 与R语言实战



吕书龙 梁飞豹 / 主 编

刘文丽 薛美玉 / 副主编



北京大学出版社
PEKING UNIVERSITY PRESS

图书在版编目(CIP)数据

应用统计分析与 R 语言实战/吕书龙, 梁飞豹主编. —北京: 北京大学出版社, 2017. 11
ISBN 978-7-301-28590-9

I. ①应… II. ①吕… ②梁… III. ①统计分析—统计程序—高等学校—教材
IV. ①C819

中国版本图书馆 CIP 数据核字 (2017) 第 195757 号

书 名: 应用统计分析与 R 语言实战
著作责任者: 吕书龙 梁飞豹 主编
责任编辑: 潘丽娜
标准书号: ISBN 978-7-301-28590-9
出版发行: 北京大学出版社
地 址: 北京市海淀区成府路 205 号 100871
网 址: <http://www.pup.cn> 新浪官方微博: @北京大学出版社
电子信箱: zpup@pup.cn
电 话: 邮购部 62752015 发行部 62750672 编辑部 62752021
印 刷 者: 三河市北燕印装有限公司
经 销 者: 新华书店
787 毫米×960 毫米 16 开本 31.25 印张 760 千字
2017 年 11 月第 1 版 2017 年 11 月第 1 次印刷
定 价: 69.00 元

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

举报电话: 010-62752024 电子信箱: fd@pup.pku.edu.cn

图书如有印装质量问题, 请与出版部联系, 电话: 010-62756370

内 容 简 介

本书以应用概率统计方法为主线,引入自由开放的统计软件 R,并按照数理统计的逻辑顺序,逐步展开理论分析、方法介绍和软件实现.主要内容有:统计软件与概率计算、数理统计初步与模拟计算、参数估计、假设检验、案例的直观分析、方差分析与正交试验设计、回归分析、多元统计初步、R 统计软件、R 软件的图形环境、R 软件中的数学运算.本书着重介绍统计方法的基本思想、背景、应用条件、实际意义及软件计算,使学生能够对数理统计方法的应用有一个较为系统、全面的了解,同时确保本书提到的所有的统计方法都能在 R 软件中得到快速实现,使学生对统计方法的实现和应用有一个标准参照.书中各章具有内容小结、知识点网络结构图和适量的习题,书末配有大部分习题的参考答案、附录中给出经典问题的索引页、部分源程序、常见软件包等.本书编写得到福州大学研究生“应用概率统计”重点课程建设项目和福州大学教务处重点教材建设项目的资助.

本书可作为高等院校工学、经济、管理、农学、医学等非数学类硕士、博士研究生以及高年级本科生学习应用概率统计(或统计分析方法或应用统计分析)课程的教材,也可作为相关学科和工程技术人员的工具参考书.为了方便教师教学和学生学习,我们将提供基本的课件、数据、程序脚本等,请通过电子邮件 wujispace@fzu.edu.cn 与作者联系.

作者简介

吕书龙，福建闽侯人，理学硕士。现为福州大学数学与计算机科学学院副教授。近年参编《应用统计方法》等3部教材；在《实验技术与管理》等核心刊物发表论文十余篇。主持或参与多项省级、校级课题以及横向课题。目前主要从事概率统计、统计应用与软件研发等教学与研究工作。

梁飞豹，福建莆田人，理学学士。现为福州大学数学与计算机科学学院副教授。近年来主编教材《概率论与数理统计》《应用统计方法》等4部；在SCI等核心刊物上发表学术论文二十余篇。目前主要从事概率统计的教学与研究工作。

刘文丽，湖北枝江人，理学硕士。现为福州大学数学与计算机科学学院副教授。近年参编《概率论与数理统计》等4部教材；在 *Expert Systems With Applications* 等核心刊物发表论文十余篇。主持或参与多项省级、校级课题及横向课题。目前主要从事概率统计、统计应用与系统工程等教学与研究工作。

薛美玉，福建福清人，理学硕士。现为福州大学数学与计算机科学学院讲师。近年参编《应用统计方法》等多部教材；在《福州大学学报》等核心刊物发表论文若干篇，主持或参与多项省级、校级教改课题。目前主要从事概率统计、时间序列等教学与研究工作。

前 言

本书是根据编者多年来积累的研究生教学经验和经历,结合各专业研究生的专业特点和科研需求,并参照教育部“工学硕士研究生应用概率统计课程教学基本要求”编写。本书编写融合概率统计理论及其应用、统计计算及软件模拟,内容编排由浅入深,直观易懂。它可以作为高等院校工学、经济、管理、农学、医学、心理学等非数学类硕士、博士研究生以及高年级本科生学习应用概率统计(或统计分析方法或应用统计分析)课程的教材,也可作为相关学科和工程技术人员的参考书。

应用统计是关于收集、整理、分析和解释受随机因素影响的数据的一门科学,是一种认识世界的方法论,是现代应用数学的一个重要组成部分,也是统计学一级学科的重要分支。它以概率论为基础,以统计软件为实现途径,旨在探索数据内在的数量规律性,以达到对客观事物的科学认识。随着计算机和互联网的发展,大数据时代的悄然而至以及人们对数据处理的需求和要求的不断提升,各种统计软件也应运而生,发展迅猛。不同学科与应用统计的交叉融合,使得应用统计在各个领域得到了广泛的应用,也推动了应用统计的纵深发展。正因为如此,应用统计课程成为高等院校非数学类硕士研究生最重要的公共基础课之一,也是非数学类硕士研究生观测世界所必备的文化修养,统计软件也成为硕士研究生进行学科深造和科学研究的必备工具之一。由于这门课程自身的特点,以及非数学类硕士研究生在本科阶段所学课程和统计软件应用水平的差异,我们在教材的编写过程中,力求体现以下几个特点:

(1) 以统计方法为主,着重介绍各种统计方法的应用背景、统计思想、适用条件、实际意义和软件实现,使学生能够对统计方法及其应用有一个系统、全面、直观的了解,以实用为主,尽量避免纯数学化的论证。

(2) 在内容编排上,考虑到不同学科和不同程度学生的需要,以概率论基础及直观模拟作为引入形成基础篇;着重介绍各种基本的统计方法,同时加入一些应用广泛的统计方法,如优化方法、非参数方法、回归专题、多元统计方法等形成方法篇;最后专门介绍 R 软件形成软件篇。不同的安排适合不同的教学模式。建议第一章和第二章作为数理统计的基础与第九章和第十章作为 R 软件的基础同步授课,并且早于其他章讲授。

(3) 考虑到专业的差异性,为了便于对照学习和应用,各章配备了一定数量的专业性统计案例分析以及完整的 R 代码实现,方便学生操作实践。

(4) 本书各章都配备适量针对性强的习题,并给出大部分习题的解答和 R 程序实现,便于教师教学和学生自学。

(5) 本书附录给出书中各种统计问题的索引页,便于读者快速检索;附录还给出丰富的 R 命令集函数、网络资源及数据集、常见软件包的介绍等。

本书主要由吕书龙、梁飞豹、刘文丽、薛美玉编写,并在集体讨论的基础上,由吕书龙和梁飞豹对全书进行统编、定稿。在此对所有参与本书编写、校对的人员表示衷心的感谢!

本书在编写、出版的过程中,得到福州大学研究生院、福州大学教务处以及北京大学出

出版社的大力支持；本书列入福州大学高水平建设规划教材，得到福州大学研究生处优质课程建设项目和福州大学教务处重点教材建设项目的资助。

限于编者的水平，书中难免存在不足之处甚至错误，恳请同行和广大读者批评指正。

编者

2016 年 12 月于福州大学

目 录

基础篇 R 软件与概率统计

第一章 统计软件与概率计算	3
1.1 R 软件简介	3
一、R 软件的获取与安装	3
二、R 软件的基本使用	4
三、R 软件的脚本	8
四、R 软件包的下载、安装与使用	9
五、R 软件的 IDE 工具 RStudio	
IDE	11
1.2 R 软件在概率论中的应用	11
一、R 软件中的集合运算、排列与组合	11
二、R 软件中的随机变量及概率计算	13
三、数值积分在 R 软件中的实现	15
四、数字特征	16
五、极限理论的模拟与计算	18
内容小结	20
习题一	21
第二章 数理统计初步与模拟计算	23
2.1 数理统计的基本概念	23
一、总体与样本	23
二、自助样本	24
三、统计量	25
四、R 软件的实现	29
2.2 经验分布函数、直方图与核密度	30
一、经验分布函数	30
二、直方图	32
三、核密度估计	35
四、R 软件的模拟与计算	38
2.3 常用的概率分布及分位点	41
一、分布及性质	41

二、概率分布的分位点	44
三、R 软件的模拟与计算	46
2.4 常用的抽样分布	47
一、正态总体的抽样分布	47
二、非正态总体的一些抽样分布	50
三、R 软件的模拟与计算	51
2.5 Monte Carlo 方法	52
一、Monte Carlo 方法求圆周率 π	53
二、Monte Carlo 方法求定积分	55
三、Monte Carlo 方法的精度分析	58
四、拟蒙特卡罗方法 (Quasi Monte Carlo) 法	61
五、系统模拟	63
2.6 Bootstrap 方法	64
一、Bootstrap 样本	64
二、参数型 Bootstrap 方法	65
三、非参数型 Bootstrap 方法	65
内容小结	65
习题二	66

方法篇 应用统计分析

第三章 参数估计	71
3.1 点估计	71
一、矩法	71
二、极大似然估计法	76
3.2 估计量的评价标准	82
一、无偏性	83
二、有效性	84
三、均方误差	85
四、一致性	86
3.3 区间估计	88
一、求未知参数 θ 的双侧置信区间	89
二、求未知参数 θ 的单侧置信区间	89

3.4 正态总体参数的区间估计	89
一、单总体的情形	89
二、双总体的情形	96
3.5 非正态总体参数的区间估计	105
一、指数分布参数的区间估计	105
二、0-1 分布参数的区间估计	106
3.6 Bootstrap 区间估计	110
内容小结	111
习题三	112
第四章 假设检验	114
4.1 假设检验的基本概念	114
一、问题的提出	114
二、原假设的讨论	117
三、 p 值检验法	118
4.2 参数型假设检验	118
一、单正态总体参数的假设检验	118
二、双正态总体参数的假设检验	122
三、非正态总体参数的假设检验	127
四、假设检验的 R 软件实现	130
4.3 非参数型假设检验	137
一、分布函数的拟合优度检验	137
二、正态性检验	143
三、基于列联表的检验	146
四、单总体秩检验	155
五、多总体秩检验	164
内容小结	167
习题四	168
第五章 案例的直观分析	171
5.1 实验对照数据的直观分析	171
一、案例及研究问题	171
二、分析过程	172
5.2 考试成绩的直观分析	175
一、案例及研究问题	175
二、分析过程	176
5.3 时间-空间数据的直观分析	186
内容小结	191
习题五	191
第六章 方差分析与正交试验设计	192
6.1 单因素方差分析	192
一、单因素试验	193
二、提出假设	193
三、统计分析	193
四、例题分析	195
6.2 双因素方差分析	198
一、有交互作用的双因素方差分析	199
二、无交互作用的双因素方差分析	202
6.3 方差齐性和均值差异的检验	206
一、方差齐性检验	206
二、均值差异性多重检验	208
6.4 正交试验设计	209
一、正交表	210
二、无交互作用的正交试验	210
三、有交互作用的正交试验	217
内容小结	222
习题六	223
第七章 回归分析	226
7.1 相关分析	226
一、相关系数	226
二、相关系数的检验	228
7.2 回归模型简介	230
一、回归的由来	230
二、回归分析的基本概念	230
7.3 线性回归模型	231
7.4 最小二乘估计及其性质	232
一、最小二乘估计	232
二、一元线性回归	235
三、最小二乘估计的性质	236
7.5 回归方程和回归系数的检验及区间估计	238
一、复相关系数	238
二、回归方程的 F 检验	240
三、回归系数的显著性检验	240
四、回归系数的区间估计	241
五、R 软件的实现	241
7.6 自变量选择	243
一、自变量选择的准则	244
二、最优回归方程	245

三、挑选变量的 R 软件实现	245	三、贝叶斯判别 (Bayes)	304
四、逐步回归	249	四、Fisher 判别	310
五、AIC 准则下逐步回归的 R 软件实现	251	五、线性判别的 R 软件实现	312
7.7 预测与控制	252	8.4 聚类分析	312
一、预测	252	一、系统聚类法的基本思想	312
二、控制	253	二、系统聚类法的步骤	313
7.8 非线性回归	256	三、系统聚类的 R 软件实现	316
一、可线性化的非线性模型	256	四、K-means 聚类的 R 软件实现	318
二、一般的非线性回归	261	8.5 主成分分析	318
7.9 非参数回归	264	一、基本原理	319
一、Nadaraya-Watson 核估计	264	二、计算步骤	319
二、近邻核估计	267	三、主成分分析的 R 软件实现	322
三、局部线性估计	268	8.6 典型相关分析	324
四、局部 p 阶多项式估计	269	一、总体典型相关	325
五、核估计的 R 软件实现	269	二、样本典型相关	325
7.10 分位数回归	272	三、典型相关变量和典型相关系数求解	326
一、回归原理及模型	272	四、典型相关系数的显著性检验	326
二、回归系数的求解	273	五、样本典型变量的得分	327
三、系数和拟合曲线的比较	273	六、例子分析和程序实现	327
四、结果解释	274	内容小结	330
7.11 关于定性变量的回归	275	习题八	331
一、虚拟变量模型	276		
二、Logistic 回归模型	278		
内容小结	280		
习题七	281		
第八章 多元统计初步	285		
8.1 多维随机变量	285		
一、多维随机变量	285		
二、多元正态分布	286		
三、抽样与统计量	287		
四、参数估计	291		
五、数据直观描述	292		
8.2 距离与相似性	294		
一、距离	294		
二、相似性	297		
8.3 判别分析	299		
一、问题描述	299		
二、距离判别及其实现	299		
		软件篇 R 软件	
		第九章 R 统计软件	335
		9.1 基本操作与控制	335
		一、脚本文件	335
		二、R 软件的帮助系统	336
		三、工作空间的保存与加载	338
		四、管理变量列表	339
		五、执行外部程序	339
		六、计算程序/代码执行的时间	340
		9.2 语法与数据类型	341
		一、赋值语句	341
		二、基本数据类型	341
		三、复杂数据类型	343
		四、以 SQL 方式操作数据框	347
		9.3 输入与输出	347

一、读取剪贴板数据	347	四、常规几何平面图	376
二、读取文本文件数据	348	五、为图形添加网格线	377
三、读取数据库数据 —— 基于 Windows 的 RODBC 包	350	六、数学标注	377
四、保存数据	352	七、指定图形窗口尺寸	378
五、sink 文本定向输出	352	八、打开新的图形窗口	378
六、利用 foreign 包读取外部数据 ..	352	九、输出图形文件	378
9.4 流程控制	352	10.2 高级绘图	379
一、一行多条语句	353	一、常用绘图函数	379
二、if / else 分支语句	353	二、条形图	382
三、switch 多分支语句	354	三、箱线图	383
四、循环结构	354	四、三维图形显示	384
五、括号的作用	355	五、lattice 软件包	388
9.5 函数与数据集	355	10.3 多图及特殊图形	390
一、基本函数介绍	355	一、打开多个图形设备	390
二、数据集	361	二、绘图区域分割	390
9.6 自定义函数	362	三、交互式图形环境	392
一、定义二元运算	362	四、绘制特殊图形	393
二、一般形式	363	五、快速矩阵绘图	395
三、缺省值和命名参数	363	10.4 表格式分组统计	396
四、省略号参数 (...)	364	一、一维数据	396
五、嵌套函数	364	二、二维列联表	396
六、递归函数	365	三、三维以上列联表	397
9.7 软件包	365	四、分组统计	397
9.8 R 软件的可视化工具与接口 ..	366	10.5 动画展示	398
一、R Commander	366	内容小结	399
二、RStudio IDE	366	习题十	399
三、Windows 平台下的可视化接口 函数	367	第十一章 R 软件中的数学运算	401
四、R 与 Matlab 的文件接口	369	11.1 矩阵运算	401
9.9 R 软件的相关网站	369	一、矩阵定义及基本性质	401
内容小结	369	二、矩阵运算的实现	402
习题九	370	11.2 数值方法	405
第十章 R 软件的图形环境	372	一、数值积分在 R 软件中的实现 ..	405
10.1 自定义绘图	372	二、函数求导	407
一、初级绘图	372	三、求极限	409
二、旋转文本输出	375	四、方程求根	409
三、在作图区域外输出文本	376	五、线性方程组求解	410
		六、非线性方程及方程组求解	410
		七、求极值	412

11.3 最优化·····	414	习题答案·····	429
一、混合整数线性规划·····	414	附录·····	474
二、运输问题·····	417	附录 A 常见正交表生成程序及表 头设计·····	474
三、指派问题·····	419	附录 B 实现常见分布的分布函数 和分位点表的 R 程序·····	477
四、线性目标规划问题·····	419	附录 C 部分问题集与索引·····	479
五、非线性规划·····	422	附录 D 部分软件包简介·····	481
六、图与网络规划·····	426	参考文献·····	483
内容小结·····	427		
习题十一·····	428		

基础篇 R 软件与概率统计

第一章 统计软件与概率计算

本章重点阐述统计软件及其在概率论中的应用。数理统计是统计学中一个应用极其广泛的学科分支,它以随机现象中的问题为研究对象并确立为总体,然后通过
对总体进行试验或观察得到数据,再采用合适的方法对数据进行分析,最后对总体
的客观规律性作出合理的估计或推断。数理统计是以概率论为理论基础,以统计软
件为实现手段。统计软件集成了各种统计方法和运算,可实现数据的预处理与分析、
提供报表和图形输出、集成编程和计算环境等,是进行统计分析不可或缺的工具。

本章第 1 节介绍统计软件的基本使用,第 2 节通过统计软件把抽象的理论和知
识直观化,以便读者对概率论的内容有个全面且直观的回顾。

1.1 R 软件简介

R 软件是由 Auckland(奥克兰)大学的 Robert Gentleman 和 Ross Ihaka 及广大志愿者
开发的一套包含数学运算、统计计算与数据分析、图形制作、程序设计等功能的自由软件系
统,可在 Unix/Linux, Windows 和 Macintosh 操作系统上运行。R 软件提供一个开放的、可
互动的、可编程的、解释型的统计计算环境。它支持命令执行、脚本执行和远程执行,并内嵌
一个非常实用的帮助系统和用户手册,还可外连互联网帮助系统。R 软件完全免费使用,更
新迅速,使用方便,其网络资源丰富,特别是软件包数量庞大、功能丰富。

一、R 软件的获取与安装

R 软件的官方网站提供最新的 R 软件和各种软件包的下载,官方网站位于:

<http://www.r-project.org>

建议到最近的镜像站点下载 R 软件

<http://mirrors.xmu.edu.cn/CRAN/>

对于 Windows 操作系统上运行的 R 软件,目前最新的版本是 R-3.3.2(32/64 bit),于 2016
年 10 月 31 日发布,大概 70MB 大小。它支持 Windows 所有的 32 位和 64 位操作系统。

R 软件的安装非常简单,双击下载的 R-3.3.2-win.exe,按照提示(或每次弹出对话框都单
击下一步)即可顺利安装。安装成功后,会在 Windows 系统的开始菜单和桌面上出现 R 的图
标,双击该图标即可运行 R 软件。

启动 R 软件后,出现的是 R 软件的图形界面,包括主菜单、工具栏和一个窗口,其标题
为“R Console”,如图 1.1.1 所示。该窗口可接受用户的命令输入,同时也输出执行命令后的
结果,但有些结果会显示在新建的窗口中(特别是图形)。界面中“>”表示命令提示符,闪烁

的“|”符号表示光标,出现它们意味着可以输入 R 命令了,或者说上一条命令执行完成,可继续下一条命令了.该窗口是用户通过命令交互来使用 R 软件的主要场所.用户输入命令时如果命令太长,可按下回车键(Enter 键),在下一行续写命令,R 程序会自动识别命令的完整性,直到命令结束,用户再按下回车键(Enter 键)确认.

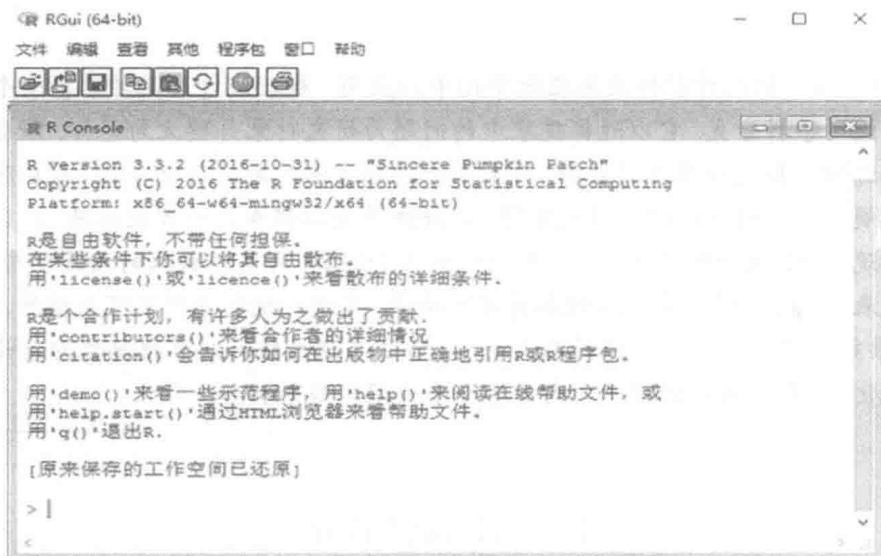


图 1.1.1 R 软件主界面

二、R 软件的基本使用

下面通过几个例子来阐述 R 软件的基本使用.

例 1.1.1 对于收集到的 (x, y) 数据,见表 1.1.1, 绘制散点-折线图.

表 1.1.1 (x, y) 数据

x	1	2	3	4	5	6	7	8	9	10
y	5	7	6	9	8	10	11	4	6	8

解 R 软件执行过程如下:

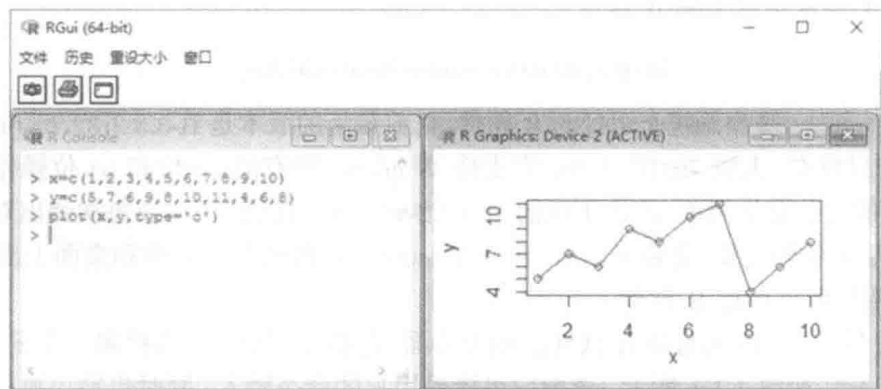
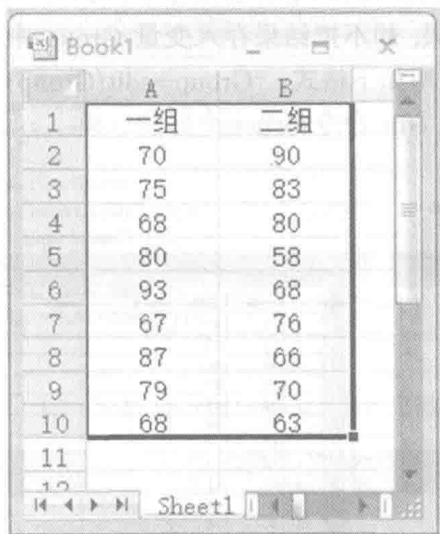


图 1.1.2 例 1.1.1 的命令和输出

对于第一条和第二条命令, 等号表示赋值语句, 其左侧是变量名, 右侧是一维数组, 由“c()”表示; 而第三条命令采用“plot()”将 x, y 当作坐标绘制在二维平面上, 并将图形显示在新建窗口中。

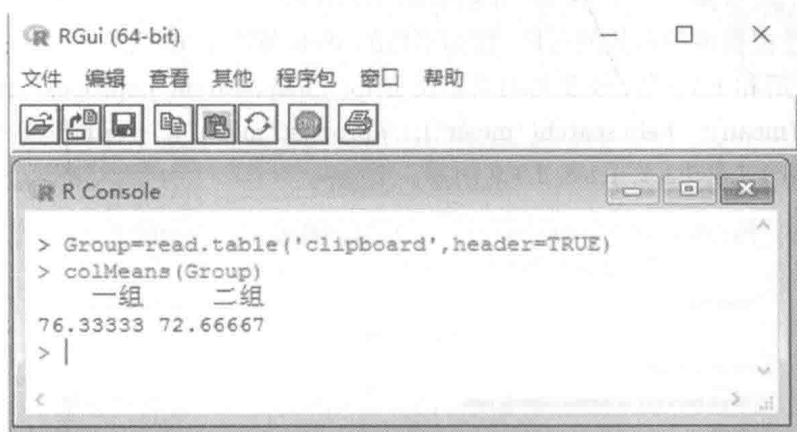
例 1.1.2 在 Excel 表格中有一批数据, 如图 1.1.3 所示, 分别计算两组数据的平均值。



	A	B
1	一组	二组
2	70	90
3	75	83
4	68	80
5	80	58
6	93	68
7	67	76
8	87	66
9	79	70
10	68	63

图 1.1.3 Excel 表格数据

解 在计算平均值之前, 必须将数据载入 R 程序中. 此处通过剪贴板实现数据的转移, 先在 Excel 中选中数据并复制到剪贴板 (如图 1.1.3 所示), 然后在 R 的命令窗口中输入如下命令:



```
> Group=read.table('clipboard',header=TRUE)
> colMeans(Group)
  一组  二组
76.33333 72.66667
> |
```

图 1.1.4 Excel 数据求均值

如图 1.1.4 所示, 第一条命令通过“read.table()”函数, 从剪贴板对象“clipboard”中获取数据, 存入“Group”变量中, 其中“header=TRUE”表示该数据的各列变量有名字. 第二条命令“colMeans(Group)”将数据按列计算平均值, 并将结果输出在命令窗口中。

例 1.1.3 将例 1.1.2 中第一组和第二组数据分别增加 60, 80 和 70, 75。