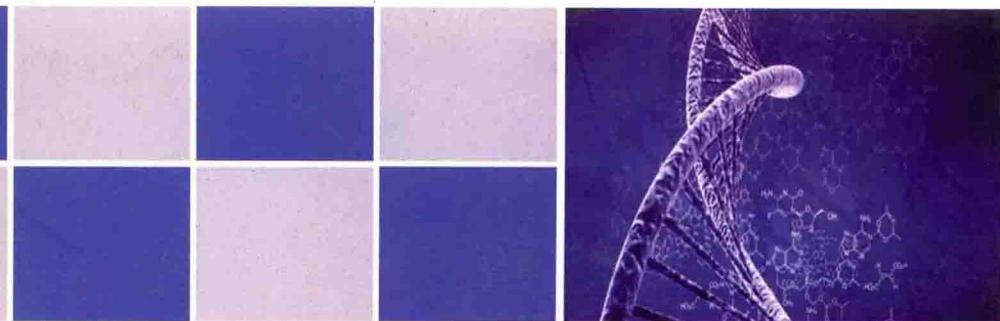


R and Its Applications
in Environmental Epidemiology



R 软件 及其环境流行病学应用

彭晓武 王家春 余松林 著

中国环境出版社

环保公益性行业科研专项经费项目系列丛书

课题：珠江三角洲地区灰霾天气对人群健康影响评估影响（200709004）

R 软件及其环境流行病学应用

彭晓武 王家春 余松林 著

中国环境出版社·北京

图书在版编目 (CIP) 数据

R 软件及其环境流行病学应用/彭晓武, 王家春, 余松林著. —北京: 中国环境出版社, 2013.9

ISBN 978-7-5111-1470-9

I. ①R… II. ①彭… ②王… ③余… III. ①统计分析—应用软件 IV. ①C819

中国版本图书馆 CIP 数据核字 (2013) 第 107103 号

出版人 王新程

责任编辑 李卫民

责任校对 唐丽虹

封面设计 刘丹妮

出版发行 中国环境出版社

(100062 北京市东城区广渠门内大街 16 号)

网 址: <http://www.cesp.com.cn>

电子邮箱: bjgl@cesp.com.cn

联系电话: 010-67112765 (编辑管理部)

发行热线: 010-67125803, 010-67113405 (传真)

印 刷 北京中科印刷有限公司

经 销 各地新华书店

版 次 2013 年 9 月第 1 版

印 次 2013 年 9 月第 1 次印刷

开 本 787×1092 1/16

印 张 34.5

字 数 832 千字

定 价 80.00 元 (附光盘)

【版权所有。未经许可, 请勿翻印、转载, 违者必究】

如有缺页、破损、倒装等印装质量问题, 请寄回本社更换

《环保公益性行业科研专项经费项目系列丛书》

编著委员会

顾 问：吴晓青

组 长：赵英民

副组长：刘志全

成 员：禹 军 陈 胜 刘海波

“十一五”环保公益性行业科研专项经费项目系列丛书

序 言

我国作为一个发展中的人口大国，资源环境问题是长期制约经济社会可持续发展的重大问题。党中央、国务院高度重视环境保护工作，提出了建设生态文明、建设资源节约型与环境友好型社会、推进环境保护历史性转变、让江河湖泊休养生息、节能减排是转方式调结构的重要抓手、环境保护是重大民生问题、探索中国环保新道路等一系列新理念新举措。在科学发展观的指导下，“十一五”环境保护工作成效显著，在经济增长超过预期的情况下，主要污染物减排任务超额完成，环境质量持续改善。

随着当前经济的高速增长，资源环境约束进一步强化，环境保护正处于负重爬坡的艰难阶段。治污减排的压力有增无减，环境质量改善的压力不断加大，防范环境风险的压力持续增加，确保核与辐射安全的压力继续加大，应对全球环境问题的压力急剧加大。要破解发展经济与保护环境的难点，解决影响可持续发展和群众健康的突出环境问题，确保环保工作不断上台阶出亮点，必须充分依靠科技创新和科技进步，构建强大坚实的科技支撑体系。

2006年，我国发布了《国家中长期科学和技术发展规划纲要（2006—2020年）》（以下简称《规划纲要》），提出了建设创新型国家战略，科技事业进入了发展的快车道，环保科技也迎来了蓬勃发展的春天。为适应环境保护历史性转变和创新型国家建设的要求，原国家环境保护总局于2006年召开了第一次全国环保科技大会，出台了《关于增强环境科技创新能力的若干意见》，确立了科技兴环保战略，建设了环境科技创新体系、环境标准体系、环境技术管理体系三大工程。五年来，在广大环境科技工作者的努力下，水体污染控制与治理科技重大专项启动实施，科技投入持续增加，科技创新能力显著增强；发布了502项新标准，现行国家标准达1263项，环境标准体系建设实现了跨越式发展；完成了100余项环保技术文件的制修订工作，初步建成以重点行业污染防治技术政策、技术指南和工程技术规范为主要内容的国家环境技术管理体系。环境

科技为全面完成“十一五”环保规划的各项任务起到了重要的引领和支撑作用。

为优化中央财政科技投入结构，支持市场机制不能有效配置资源的社会公益研究活动，“十一五”期间国家设立了公益性行业科研专项经费。根据财政部、科技部的总体部署，环保公益性行业科研专项紧密围绕《规划纲要》和《国家环境保护“十一五”科技发展规划》确定的重点领域和优先主题，立足环境管理中的科技需求，积极开展应急性、培育性、基础性科学的研究。“十一五”期间，环境保护部组织实施了公益性行业科研专项项目 234 项，涉及大气、水、生态、土壤、固废、核与辐射等领域，共有包括中央级科研院所、高等院校、地方环保科研单位和企业等几百家单位参与，逐步形成了优势互补、团结协作、良性竞争、共同发展的环保科技“统一战线”。目前，专项取得了重要研究成果，提出了一系列控制污染和改善环境质量技术方案，形成一批环境监测预警和监督管理技术体系，研发出一批与生态环境保护、国际履约、核与辐射安全相关的关键技术，提出了一系列环境标准、指南和技术规范建议，为解决我国环境保护和环境管理中急需的成套技术和政策制定提供了重要的科技支撑。

为广泛共享“十一五”期间环保公益性行业科研专项项目研究成果，及时总结项目组织管理经验，环境保护部科技标准司组织出版“十一五”环保公益性行业科研专项经费项目系列丛书。该丛书汇集了一批专项研究的代表性成果，具有较强的学术性和实用性，可以说是环境领域不可多得的资料文献。丛书的组织出版，在科技管理上也是一次很好的尝试，我们希望通过这一尝试，能够进一步活跃环保科技的学术氛围，促进科技成果的转化与应用，为探索中国环保新道路提供有力的科技支撑。

中华人民共和国环境保护部副部长

吴晓青

2011 年 10 月

序

目前国际上有许多很好的统计软件，如 SAS (Statistical Analysis System), SPSS (Statistical Package for the Social Sciences, 现更名为 Statistical Product and Service System) 等。虽然它们都具有强大的统计计算功能，并得到国际公认，但都属于商业软件，对于经济欠发达国家的用户来说难以承受，而且这些软件中模块的更新较慢。

R 软件是一种免费软件，具有强大的统计计算和绘图功能，最初由新西兰奥克兰大学统计学系的 Robert Gentleman 和 Ross Ihaka 编写。其源代码于 1995 年公布于众。由于该软件具有自编程序的功能，许多统计学者不断添加自编的新程序，使其功能日益完善，受到全世界统计工作者的青睐。本书的目的在于介绍该统计软件，为使用者提供一种新的软件选择。

随着经济的快速发展，环境保护越来越受到社会的重视，环境与健康关系的流行病学研究日益深入。本书向广大的环境流行病学工作者推荐应用 R 软件作为计算工具，书中内容包括处理独立观察资料的常用统计方法和处理非独立观察资料的时间序列方法。

本书的内容按由浅入深、循序渐进的顺序编排。全书内容包括三部分：第一部分为 R 基础，第二部分为常用统计方法，第三部分为时间序列分析方法。各部分的内容安排为：R 基础包含 R 软件的安装、窗口的介绍和 R 的基本工作原理，R 的数据运算和数据集操作，R 的绘图和编程等。常用统计方法包括定量变量的描述性统计，分类变量的描述性统计，区间估计与假设检验 (t 检验， χ^2 检验，单向和双向方差分析等)，二项分布与泊松分布，生存率的计算和比较方法，线性回归与相关，Logistic 回归，Cox 比例风险模型等。时间序列分析方法包括时间序列的经典分析方法，平稳时间序列分析方法，非平稳时间序列分析方法，季节非平稳时间序列模型，带输入变量的时间序列模型和泊松分布广义加性模型，以及经济效应分析等。

初学者在开始接触 R 软件时往往会感到不知所措，为了给大家的学习提供方便，我们编写了各部分的示例程序 (R 代码)，大家可以在其导引下，更容易地学习和理解书中的内容。由于对一个新软件的学习与掌握是一个从不熟悉到熟悉的过程，建议读者仔细阅读 R 基础，并在计算机上进行大量的操作练习。为了顺利地学习和掌握第三部分的时间序列分析方法，读者应具有一定的统计学基础。虽然从学习统计学的角度来说，已具备坚实的基本统计知识的读者可以跳过第二部分而直接开始学第三部分。但是从学习 R 软件的角度来说，我们建议从长远考虑，最好还是全书完整地学习一遍，以便于系统地掌握 R 语言的基础知识。

本书为环境流行病学工作者介绍 R 软件，使之应用于环境污染与健康关系的研究。鉴于本书内容的系统性和叙述的详尽性，并有配套的案例分析，亦可作为一般统计工作者的参考书以及高等院校的统计学软件教材。

在本书的编写过程中，得到华中科技大学同济医学院公共卫生学院领导和环境保护部华南环境科学研究所领导的关心和支持，得到了我们的很多同事、朋友以及家人的关心、鼓励和帮助，谨此致以谢意。崔伊薇教授还对本书的时间序列分析部分做了仔细阅读和认真修改，特此一并致谢。还要特别感谢李卫民编辑，她为本书的编审工作付出了辛勤劳动。她不仅对书稿进行了认真细致的逐字逐句审核，还帮助作者发现了许多原书稿中的错误和不足之处。本书的出版与她的努力是分不开的。

当今统计学和相应软件的发展一日千里。限于作者的学术水平，本书难免存在许多错误和不足之处，诚望广大读者给予批评指正，以便再版时予以修正。

本书中大部分例子引自所附参考文献，特别是引自余松林主编的《医学统计学》（人民卫生出版社，2002年）。谨向相关编委王洪源教授、王增珍教授、宇传华教授、张菊英教授、周燕荣教授、骆福添教授和曹素华教授致以谢意。

本书得到环境保护部环保公益性行业科研专项经费项目“珠江三角洲地区灰霾天气对人群健康影响评估研究”（200709004）的资助。

彭晓武 王家春 余松林

2013年3月2日

目 录

第一篇 R 基础

第1章 绪论	3
1.1 R 的起源和发展	3
1.2 R 的功能和特点	3
1.3 R 软件的获取与安装	4
1.4 R 工作基本原理	5
1.5 R 在线帮助	7
1.6 获取关于 R 和系统的信息	9
第2章 R 的数据操作	11
2.1 数的简单运算	11
2.2 数学函数	14
2.3 向量	16
2.4 矩阵	18
2.5 数组	24
2.6 因子向量	29
2.7 随机序列	30
第3章 对象和数据框	32
3.1 对象的种类与属性	32
3.2 改变对象的属性	35
3.3 对象的使用	37
第4章 R 数据的生成、导入和导出	49
4.1 创建 R 数据集	49
4.2 从文件读取数据	52
4.3 从其他应用软件所产生的数据文件导入数据	62
4.4 存储数据	65
4.5 在 R 中显示数据	67
第5章 数据集的整理	73
5.1 数据集的检查	73
5.2 数据集的修改	75
5.3 变量值的替换或取出数据子集	79
5.4 向量和矩阵的合并与删除	81

第 6 章 R 程序包	85
6.1 R 程序包的种类	85
6.2 程序包的安装	87
6.3 关于程序包操作的函数	89
6.4 程序包及其帮助	90
第 7 章 R 函数	100
7.1 函数的调用与查询	100
7.2 用户自定义函数	102
7.3 几种特殊的函数	108
7.4 泛型函数	111
第 8 章 R 绘图	114
8.1 管理绘图	115
8.2 绘制图形	120
8.3 绘图参数与绘图符号	138
8.4 几种复杂图形的绘制	143
第 9 章 控制流	153
9.1 if 条件语句	153
9.2 ifelse() 函数	155
9.3 switch() 函数	156
9.4 for() 语句	158
9.5 while() 语句	159
9.6 repeat 语句	160
第 10 章 R 编程实践	163
10.1 一个非线性模型的编程	163
10.2 编写一个两独立样本 t 检验的 R 程序	165
10.3 独立样本 2×2 差异性检验的自定义函数	165
10.4 计算线性回归参数估计值的程序	167
10.5 对三个不同种属的鸟绘图	168
10.6 编写用 Newton-Raphson 迭代法求解非线性方程组的根的程序	169
10.7 用递归函数计算积分的程序	171
10.8 正态分布概率密度函数动画程序	172
10.9 一个猜数字的小游戏	173
10.10 程序的运行方式	174

第二篇 常用统计方法

第 11 章 定量变量的描述性统计	177
11.1 频数分布	177
11.2 集中趋势	180
11.3 离散趋势	185

11.4 正态分布	188
11.5 医学参考值的估计	191
第 12 章 分类变量的描述性统计	193
12.1 常用的比例指标及其意义	193
12.2 相对危险度与优势比	195
12.3 率的标准化法	197
12.4 动态数列	201
12.5 比例指标应用时的注意事项	203
第 13 章 抽样误差、区间估计与假设检验	204
13.1 均数的抽样误差	204
13.2 均数的抽样误差的分布—— <i>t</i> 分布	207
13.3 总体均数的可信区间估计	208
13.4 方差的抽样误差与可信区间估计	209
13.5 率的抽样误差与可信区间估计	210
13.6 假设检验	212
第 14 章 χ^2 检验	229
14.1 χ^2 分布	229
14.2 拟合优度检验	230
14.3 独立性检验	233
14.4 趋势检验	242
14.5 多个四格表的联合分析	243
14.6 四格表的费歇尔精确概率检验	244
第 15 章 方差分析	246
15.1 单向方差分析	246
15.2 双向方差分析	256
第 16 章 二项分布与泊松分布	271
16.1 二项分布的概念	271
16.2 二项分布的性质	273
16.3 二项分布的应用	275
16.4 泊松分布的概念	279
16.5 Poisson 分布的性质	280
16.6 Poisson 分布的应用	283
第 17 章 生存时间资料的非参数分析方法	288
17.1 生存时间资料的特点	288
17.2 小样本生存率的 Kaplan-Meier 估计	293
17.3 大样本生存率的寿命表法估计	296
17.4 生存曲线比较的假设检验	298
第 18 章 回归与相关	304
18.1 直线回归与相关	304

18.2 多元线性回归与相关.....	321
第 19 章 Logistic 回归.....	338
19.1 Logistic 回归的模型结构.....	338
19.2 回归参数的估计及其假设检验.....	341
19.3 回归参数的解释.....	344
19.4 回归模型拟合情况的分析.....	345
19.5 应用 Logistic 回归时值得注意的几个问题.....	349
19.6 匹配设计资料的 Logistic 回归.....	354
第 20 章 Cox 比例风险模型.....	359
20.1 模型结构与参数估计.....	359
20.2 回归模型的应用.....	361
20.3 风险函数和生存函数的估计.....	365
20.4 比例风险假设的检验.....	372
20.5 时依协变量.....	377

第三篇 时间序列分析方法

第 21 章 时间序列的特点.....	381
21.1 时间序列资料的组分.....	382
21.2 时间序列的自相关性.....	382
21.3 时间序列的平稳性概念.....	386
21.4 几种基本的平稳时间序列模型.....	386
21.5 时间序列平稳性检验.....	388
第 22 章 时间序列的经典分析方法.....	391
22.1 经典组分分解法.....	391
22.2 线性回归分析法.....	398
22.3 调和季节模型 (harmonic seasonal models)	403
22.4 指数匀滑与 Holt-Winters 指数匀滑法.....	410
第 23 章 平稳时间序列分析.....	415
23.1 差分算子和后向移位算子.....	415
23.2 自回归模型.....	416
23.3 移动平均模型.....	420
23.4 自回归移动平均模型.....	424
23.5 平稳时间序列模型的配合.....	428
第 24 章 非平稳时间序列分析.....	441
24.1 非平稳时间序列的平稳化.....	441
24.2 ARIMA 模型.....	445
24.3 ARIMA(p,d,q)模型的预报.....	452
第 25 章 季节非平稳时间序列模型.....	455
25.1 单纯季节自回归求和移动平均模型.....	455

25.2	复合性季节自回归求和移动平均模型	461
第 26 章	带输入变量的时间序列模型	478
26.1	具有自相关残差的回归模型	479
26.2	干预模型	485
26.3	传递函数模型	489
第 27 章	广义加性模型	503
27.1	广义加性模型的结构	503
27.2	广义加性模型配合的例子	504
第 28 章	不良健康效应的经济损失分析	513
28.1	健康效应模型	513
28.2	经济损失的估计	519
附表 1	标准正态分布曲线下的面积	523
附表 2	t 界值表	524
附表 3	卡方界值表	526
附表 4	F 分布的上侧临界值表（供方差分析用）	527
附表 5	q 界值表	530
附表 6-1	百分率的可信区间	531
附表 6-2	百分率的可信区间	532
附表 6-3	百分率的可信区间	533
附表 7	Poisson 分布的可信区间	535
附表 8	r 界值表（Pearson 相关系数检验用）	535

第一篇 R 基础

第1章 绪论

1.1 R 的起源和发展

R 是一种统计计算和制图的语言或环境，同时又是一款统计软件。R 最初是 20 世纪 80 年代问世的 S 统计计算语言系统的一个分支，因此可以把 R 看成 S 语言的一种补充。所以很多用 S 语言编写的程序代码可以直接在 R 环境中运行。R 语言最先由新西兰奥克兰大学统计系的 Robert Gentleman 和 Ross Ihaka 合作创建，其源代码于 1995 年公布于众。R 语言的基本参考书是由 Richard A. Becker, John M. Chamber 和 Allan R. Wilks 所著的 *The New S Language: A Programming Environment for Data Analysis and Graphic*。与 S 语言一样，R 语言在设计理念上与其他统计软件系统存在着许多不同。虽然 R 出自 S，但 R 与 S 也还是有一些不同之处，读者不要混淆：

①在 R 中，统计分析通常按一系列步骤来完成，并将部分统计分析结果直接显示在屏幕上，某些中间结果（如 P 值、回归系数、残差等）保存在所谓的“对象”中。②做模型拟合时，SAS 软件和 SPSS 软件会直接给出可以复制的输出结果，而 R 则将其计算结果保存在一个“拟合对象”中，需要用另外的 R 语句（称作“函数”）来访问，直接输出的结果很少。③SAS 软件和 SPSS 软件对英文字母的大小写是不区分的，而 R 则是区分字母的大小写的，即如果指定了大写字母，而你用了小写字母的话，程序将停止运行并提示错误。

R 是在 GNU 协议（General Public Licence）下免费发行的。它的开发及维护工作由 R 开发核心小组（R Development Core Team）具体负责。由科学家维护，为科学家服务。R 自公布以来，很快赢得了很多用户，并不断有贡献者加入新的程序包，以改善和扩展 R 的功能。软件及其程序包的版本更新也是很快的，几个月内就有新的版本出现。

1.2 R 的功能和特点

R 内部包含了许多实用的统计分析函数和作图函数。作图函数能够在一个独立的窗口中展示它所产生的图形。用户还可以将这些图形保存为（甚至直接输出进）各种格式的图形文件（如 bmp、emf、jpg、png、ps、pdf、tif 等）。统计分析结果也能够直接显示出来，一些中间结果（如 P 值，回归系数和残差等）既可以用作进一步分析的对象，也可以保存到指定的文件。

在 R 语言中，用户既可以使用循环语句来连续分析多个数据集，也可以将多个不同的统计函数放在同一个语句中，执行更为复杂的分析。用户还可以借鉴网上提供的用 S 语言编写的大量的程序，而且大多数都能被 R 直接调用。

用户在初学阶段可能会觉得 R 比较复杂多变，令人琢磨不定。其实，R 的一个突出的优点正是它的灵活性。一般的统计软件通常会直接展示其分析结果，而 R 则是把这些结果

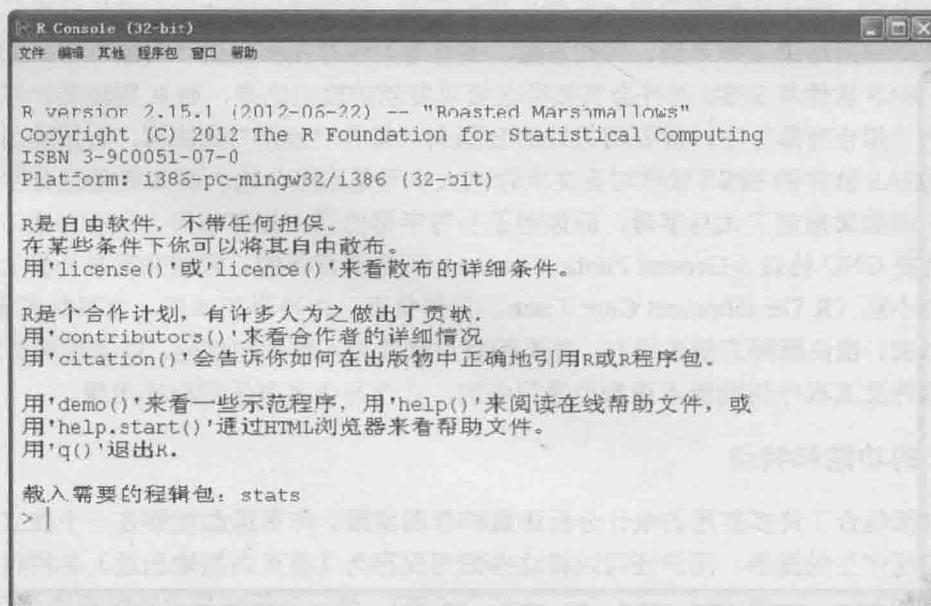
保存在对象(object)中，所以常常在分析执行结束后并不或很少显示任何结果，初学者可能会对此感到困惑。其实这个特点是很有用的，因为用户随时可以从大量的结果中选择性地提取出感兴趣的部分。例如，我们要运行 20 个回归分析而只想比较其回归系数，在 R 中就可以选择只显示所有分析所得出的回归系数，这样，结果就只占了一排，而用别的软件就可能要一下打开 20 个窗口。在下面的章节中，我们会看到更多 R 相对于传统软件更为灵活和优越的例子。

1.3 R 软件的获取与安装

R 由“综合 R 档案馆网站”(Comprehensive R Archive Network, CRAN)发布。下载网址为 <http://cran.r-project.org/mirrors.html>。到本书完成初稿时为止，R 的 Windows 版本为 R 2.15.1，大小约为 47MB。

R 的安装过程很简单，登录网站 <http://cran.r-project.org> → download R → 选择一个靠近你的地区(如 China 下的 <http://mirrors.ustc.edu.cn/CRAN/>) → Download R for Windows → base → Download R 2.15.1 for Windows → 指定一个文件夹以保存所下载的 R 安装文件 → 双击该文件以开始安装 R 软件 → 遵循自动安装中的提示指引 → 完成 R 的安装。

安装完成后，双击桌面上 R 软件图标即可启动 R。显示控制台(R Console)窗口如下。



在以上窗口中，最底下的一行为红色的“>”，后边还有一个红色的“|”在闪动，这个“>”就是 R 提示符，它的右边如果没有任何显示内容，则表明 R 正在等待你的命令(这时候可以键入命令，然后按回车键以执行此命令)。如果窗口中最底下的一行不是提示符，则表明 R 正在执行命令(这时候需要你等待 R，直到窗口中最底下的一行出现提示符为止)。

从以上过程不难看出，我们使用 R 的过程就是我们与 R 对话的过程(就像使用 QQ 那样)，这就是 R 的交互式操作。

当控制台窗口中的内容不再需要了，我们可以按 Ctrl-L 键，或用右键单击窗口内部的此为试读，需要完整PDF请访问：www.ertongbook.com