

SAS

高级统计分析教程 (第2版)

◎ 胡良平 主编
◎ 高 辉 审校



中国工信出版集团

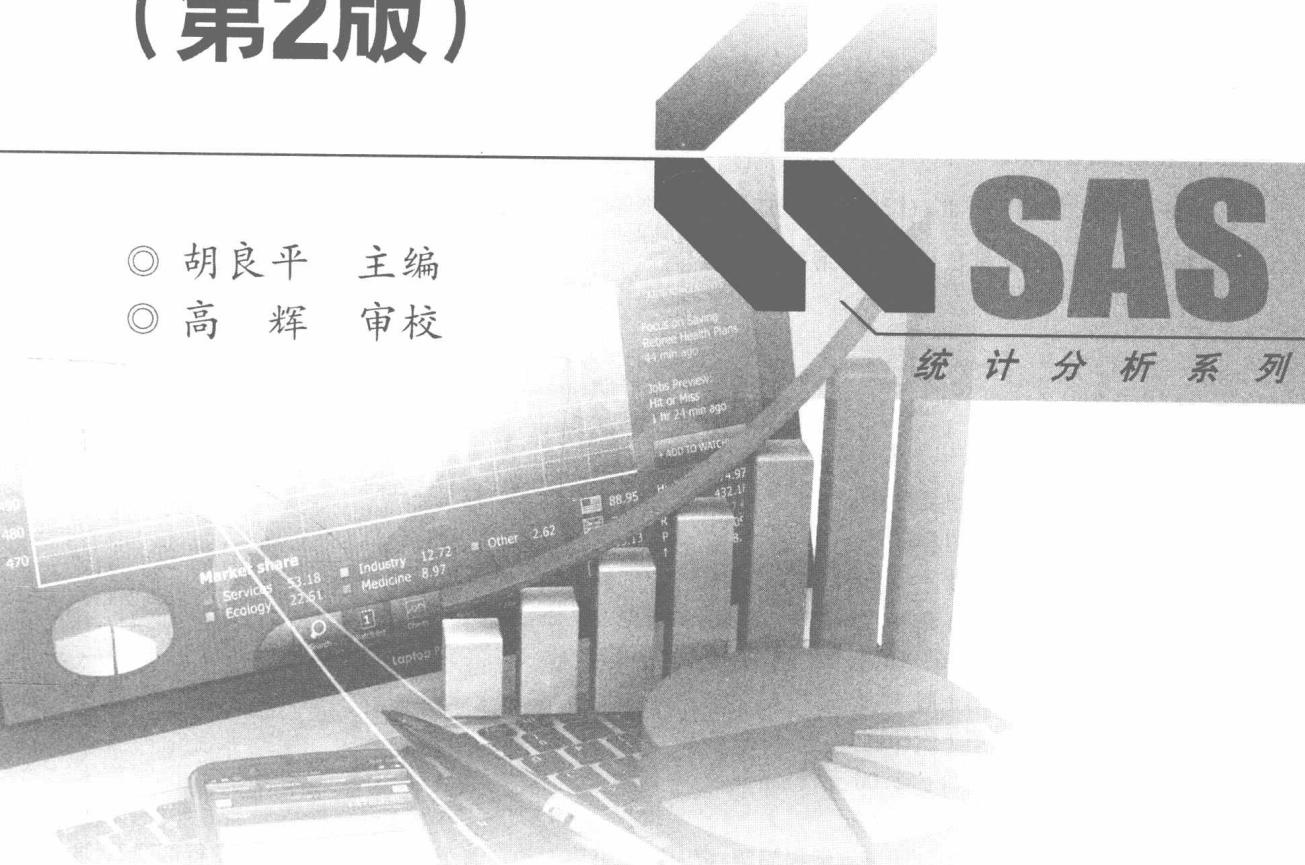


电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

SAS

高级统计分析教程 (第2版)

◎ 胡良平 主编
◎ 高 辉 审校



电子工业出版社

Publishing House of Electronics Industry

北京 · BEIJING

内 容 简 介

本书基于 SAS 9.3 版本, 内容丰富且新颖, 适用面宽且可操作性强, 涉及统计学基础和现代多元统计分析。这些内容高质量、高效率地解决了各种多元统计分析、数据挖掘、遗传资料统计分析和 SAS 实现及结果解释等人们迫切需要解决却又十分棘手的问题。

本书内容共 6 篇, 第 1 篇包括第 1~4 章, 回答了 4 个基础性问题, 即“如何确保数据是值得分析的”、“如何选择统计图并用 SAS 绘制”、“如何给统计分析方法分类与合理选用统计分析方法”和“如何基于偏好数据确定多因素的最佳水平组合”; 第 2 篇包括第 5~12 章, 介绍了研究变量之间相互和依赖关系的 8 种多元统计分析方法; 第 3 篇包括第 13~16 章, 介绍了评价样品间亲疏、优劣或相对位置的 4 种多元统计分析方法; 第 4 篇包括第 17~19 章, 介绍了评价变量与样品之间关联性的 3 种多元统计分析方法; 第 5 篇包括第 20~24 章, 第 6 篇包括第 25~26 章, 介绍了数据挖掘、生物信息学和遗传资料分析 3 大领域方面的知识和技术。另有配套的辅助资料, 可在华信教育资源网 www.hxedu.com.cn 查询。

本书适合需要运用现代多元统计分析以及数据挖掘和遗传资料分析等相关领域知识解决实际问题的本科生、研究生、博士生、科研和管理工作者、临床医生和杂志编辑学习和使用。

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

SAS 高级统计分析教程/胡良平主编.—2 版.—北京: 电子工业出版社, 2016.1

(统计分析系列)

ISBN 978-7-121-27640-8

I. ①S… II. ①胡… III. ①统计分析—应用软件—高等学校—教材 IV. ①C819

中国版本图书馆 CIP 数据核字(2015)第 281679 号

策划编辑: 秦淑灵

责任编辑: 苏颖杰

印 刷: 三河市鑫金马印装有限公司

装 订: 三河市鑫金马印装有限公司

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: 787×1092 1/16 印张: 24.25 字数: 621 千字

版 次: 2016 年 1 月第 1 版

印 次: 2016 年 1 月第 1 次印刷

印 数: 3000 册 定价: 55.00 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 zlts@phei.com.cn, 盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线:(010)88258888。

编 委 会

主 编 胡良平

审 校 高 辉

副主编 葛 毅 李长平 柳伟伟 胡纯严 郭 晋

编 委(以单位及姓氏笔画为序)

广东医学院 徐秀娟

山西医科大学 张岩波 罗艳虹

中日友好医院 周诗国

天津中医药大学 赵铁牛

天津医科大学 李长平

北京邮电大学 张晓航

军事医学科学院 李伍举 柳伟伟 胡良平 胡纯严

解放军卫生信息中心 葛 毅

解放军 95969 部队 高 辉

哈尔滨医科大学 李 霞 张瑞杰

济南军区疾控中心 李子建

首都医科大学 刘惠刚 罗艳侠 郭秀花 郭 晋

第三军医大学 伍亚舟 易 东

PPD 医药公司 毛 玮

第2版前言

本书第1版于2010年6月由电子工业出版社正式出版后的2~3年间，经受住了市场的考验，得到了广大读者的褒奖、支持和鼓励。大约在2013年下半年，因市场上此书脱销，故出版社征求本人意见后，重印了数千册。最近我在主讲全国统计学培训班期间，听见部分学员说此书在全国各大网店缺货。对于可能出现脱销现象，在2014年年底，此书编辑早就预见到了，并敦促我抓紧时间修订。因为SAS软件内容更新很快，第1版主要基于SAS 9.1.3版本，而现在，SAS 9.2和SAS 9.3版本较之前增添了许多新功能。另外，本书第1版出版后，在使用过程中也发现了个别不妥和疏漏之处，有必要对其中的差错之处进行纠正；同时，有必要结合SAS软件的特点和实际使用统计学的需求，对全书内容和布局进行必要的修改、补充和调整。

由于增加的新内容比较多，修订后的篇幅大大超过了原书，本书编辑提出了一个大胆的方案——将第2版改为两册出版，书名分别为《SAS常用统计分析教程(第2版)》与《SAS高级统计分析教程(第2版)》，本书是后者。

本书的核心内容是现代多元统计分析(不包括多元方差和协方差分析、判别分析，这些内容已被收录入《SAS常用统计分析教程(第2版)》中)，其内容散布在第2~4篇共15章之中，分别冠以篇名“变量间相互与依赖关系分析”、“样品间亲疏优劣或相对位置分析”和“样品与变量或原因与结果之间的关联性分析”，以使条理更加清楚，便于读者学习和使用。

为了适应科学技术的发展对统计学提出的需求，书中保留了数据挖掘、生物信息学和遗传资料统计分析这三大领域方面的知识和技术，有利于开拓读者思路和激发其对统计学更广泛的兴趣。

在本书出版之际，我情不自禁地要提及，本室2013和2014级硕士研究生刘一松、胡完和郭春雪，他们对全书做了认真的校对。最后，请允许我感谢直接和间接为本书第1版和第2版付出过辛勤劳动的所有同志和朋友！

由于编者水平有限，书中难免会出现这样或那样的不妥，甚至错误之处，恳请广大读者不吝赐教，以便再版时修正。为便于与读者沟通和交流，特将我的电子邮箱地址和有关网址呈现在此：lphu812@sina.com；www.statwd.com；www.huasitai.com。

主编 胡良平
于北京军事医学科学院
2015年8月10日

第1版前言

众所周知，统计学内容非常丰富，学习和正确使用它难度很大；SAS 软件功能非常强大，实用面极宽，SAS 语言又十分繁杂，学习和全面掌握其用法并非易事；显然，实际工作者要想在较短的时间内，学会用 SAS 软件方便快捷且正确地解决各种实验设计、统计表达与描述、常见和多元统计分析、现代回归分析、数据挖掘和基因表达谱分析等问题，几乎是天方夜谭，然而，本书却使其成为现实！

笔者有何灵丹妙药呢？“面向问题”的思维模式和写作手法是解决复杂问题、并使其化繁为简、实用方便的“锦囊妙计”。本书各章针对拟解决的每个具体问题，首先给出“问题、数据和统计分析方法的选择”，接着，用编制好的 SAS 程序分析给定的资料，并给出程序修改指导、主要输出结果及其解释。为了降低书的出版费用，笔者千方百计对书稿进行精简，使以纸质印刷的篇幅约为原稿的三分之一，而原稿中三分之二的内容以附录的形式放在与本书配套的光盘上。本书正文内容共分 8 篇 47 章。第 1 篇 对定量结果进行差异性分析；第 2 篇 对定性结果进行差异性分析；第 3 篇 对定量结果进行预测性分析；第 4 篇 对定性结果进行预测性分析；第 5 篇 多变量间相互与依赖关系分析；第 6 篇 变量或样品间亲疏关系或近似程度分析；第 7 篇 数据挖掘技术与基因表达谱分析简介；第 8 篇 用编程法绘制统计图与实现实验设计。各章末尾均注明参编者详细名单。值得一提的是：本书中的所有计算基于 SAS 9.2 版本，少量 SAS 过程（特别是涉及 SAS/Genetics 模块）在 SAS 9.2 以前的环境下不能正常运行。

与本书配套光盘上的内容有：附录 1 与 SAS 语言有关的内容简介，包括第 48 章 SAS 语句简介、第 49 章 SAS 过程简介、第 50 章 SAS 函数简介、第 51 章 SAS 宏简介、第 52 章 SAS ODS 简介、第 53 章 SAS SQL 简介、第 54 章 SAS 数组简介和第 55 章 SAS/IML 简介；附录 2 四个非编程模块简介，包括第 56 章 SAS/ASSIST 模块用法简介、第 57 章 SAS/ANALYST 模块用法简介、第 58 章 SAS/INSIGHT 模块用法简介和第 59 章 SAS/ADX 模块用法简介；附录 3 数据挖掘技术与基因表达谱分析，包括第 60 章 数据挖掘的概念及常用统计分析技术、第 61 章 基因表达谱的概念及数据分析技术和第 62 章 生物信息学；附录 4 各章实例和数据；附录 5 各章 SAS 程序；附录 6 各章 SAS 输出结果；附录 7 各章计算原理和计算公式；附录 8 各章参考文献和附录 9 胡良平统计学专著与配套软件简介。

在本书即将出版之际，笔者真挚地感谢为本书作出过突出贡献的来自全国十多所高等院校的教授、副教授和青年学者，如山西医科大学张岩波、余红梅和郭东星教授，中山医科大学张晋昕教授，武汉大学毛宗福教授，哈尔滨医科大学李霞和张瑞杰教授，第三军医大学易东教授，首都医科大学郭秀花教授，北京大学医学部曹波副教授，北京邮电大学张晓航副教授等；还要感谢所有为本书付出过辛勤劳动的人们，他们的名字被写在本书的编委名单之中。特别是高辉、柳伟伟、周诗国、郭晋为本书的审阅工作付出了大量细致而又卓有成效

的劳动，而本室 2009 级四名硕士研究生(鲍晓蕾、贾元杰、关雪、梦冰)也为本书的校对做出了很多贡献。正是由于他们的积极参与、不懈努力和真心奉献，才使这部历时四年的专著能够问世！

由于笔者水平有限，书中难免会出现这样或那样的不妥，甚至错误之处，恳请广大读者不吝赐教，以便再版时修正。

主编 胡良平
于北京军事医学科学院
生物医学统计学咨询中心
2010 年 1 月 14 日

目 录

第1篇 统计分析基础

第1章 应确保数据是值得分析的	(1)
1.1 什么是数据和/或统计资料	(1)
1.1.1 数据不等于统计资料	(1)
1.1.2 统计资料的要素	(2)
1.2 确保数据值得分析的第一道关——制订科学完善的课题设计方案 ...	(3)
1.2.1 什么叫科学研究	(3)
1.2.2 科学研究与课题之间是什么关系	(3)
1.2.3 做课题之前为什么要制订课题设计方案	(3)
1.2.4 课题设计方案有哪些种类 ...	(3)
1.2.5 科学完善的科研设计方案的标志	(5)
1.3 确保数据值得分析的第二道关——实时进行严格的过程质量控制 ...	(6)
1.3.1 必须严格控制课题实施过程中的质量	(6)
1.3.2 进行质量控制的必要性	(6)
1.3.3 进行质量控制的环节与措施	(7)
1.4 确保数据值得分析的第三道关——确保数据的原始性没有被破坏 ...	(7)
1.4.1 应有切实可行的措施确保收集的数据具有原始性	(7)
1.4.2 与常见试验设计类型对应的规范化统计表	(7)
1.5 常见不值得分析的数据种类	(17)
1.5.1 人为编造的数据是不值得分析的	(17)
1.5.2 产生于质量控制不严的数据是不值得分析的	(20)
1.5.3 经过错误的方法加工整理后的数据是不值得分析的 ...	(20)
1.5.4 不符合特定统计分析方法要求的数据是不值得分析的 ...	(21)
1.5.5 盲目解释基于误用统计分析方法所得到的分析结果是不可取的	(23)
1.5.6 缺失值过多的数据是不值得分析的	(24)
1.6 本章小结	(24)
第2章 绘制统计图	(25)
2.1 问题、数据及统计描述方法的选择	(25)
2.1.1 问题与数据	(25)
2.1.2 对数据结构的分析	(27)
2.1.3 分析目的与统计描述方法的选择	(27)
2.1.4 统计图概述	(28)
2.2 绘制单式条图	(28)
2.2.1 程序及说明	(28)
2.2.2 输出单式条图	(29)
2.3 绘制复式条图	(29)
2.3.1 程序及说明	(29)
2.3.2 输出复式条图	(30)
2.4 绘制百分条图	(30)
2.4.1 程序及说明	(30)
2.4.2 输出百分条图	(31)
2.5 绘制圆图	(32)
2.5.1 程序及说明	(32)
2.5.2 输出圆图	(32)
2.6 绘制箱式图	(33)
2.6.1 程序及说明	(33)

2.6.2 输出箱式图	(33)	3.2.1 合理选用统计分析方法的四要素	(47)
2.7 绘制直方图	(34)	3.2.2 合理选用统计分析方法的实例演示	(48)
2.7.1 程序及说明	(34)	3.3 面对实际问题合理选用统计分析方法的要领	(50)
2.7.2 输出直方图	(34)	3.3.1 描述性统计分析	(50)
2.8 绘制散点图	(35)	3.3.2 探索性统计分析	(51)
2.8.1 程序及说明	(35)	3.3.3 传统差异性统计分析	(58)
2.8.2 输出散点图	(35)	3.3.4 相关分析	(59)
2.9 绘制普通线图	(35)	3.3.5 回归分析	(59)
2.9.1 程序及说明	(35)	3.3.6 广义综合评价	(60)
2.9.2 输出普通线图	(36)	3.4 本章小结	(60)
2.10 绘制半对数线图	(37)	第4章 结合分析	(61)
2.10.1 程序及说明	(37)	4.1 问题与数据结构	(61)
2.10.2 输出半对数线图	(37)	4.1.1 实例	(61)
2.11 绘制P-P图和Q-Q图	(38)	4.1.2 对数据结构的分析	(63)
2.11.1 程序及说明	(38)	4.1.3 统计分析目的与分析方法的选择	(63)
2.11.2 输出P-P图	(38)	4.2 结合分析内容简介	(63)
2.12 本章小结	(39)	4.2.1 基本概念	(63)
第3章 统计分析方法的分类与合理选用的关键技术	(40)	4.2.2 基本原理	(64)
3.1 统计分析方法的分类	(40)	4.3 结合分析的应用	(65)
3.1.1 概述	(40)	4.3.1 用SAS分析例4-1中的资料	(65)
3.1.2 描述性统计分析	(40)	4.3.2 用SAS分析例4-2中的资料	(67)
3.1.3 探索性统计分析	(40)	4.4 本章小结	(70)
3.1.4 广义差异性统计分析	(41)		
3.1.5 广义相关与回归分析	(42)		
3.1.6 广义综合评价	(45)		
3.2 合理选用统计分析方法的关键技术	(47)		

第2篇 变量间相互与依赖关系分析

第5章 路径分析	(71)	5.2.2 适合进行路径分析的数据结构	(73)
5.1 问题与数据结构	(71)	5.2.3 路径分析的基本概念	(74)
5.1.1 实例	(71)	5.2.4 路径分析的基本原理	(74)
5.1.2 对数据结构的分析	(72)	5.2.5 路径分析的步骤	(78)
5.1.3 分析目的与统计分析方法的选择	(72)	5.3 路径分析的应用	(79)
5.2 路径分析内容简介	(73)	5.3.1 用REG过程实现路径分析	(79)
5.2.1 路径分析概述	(73)	5.3.2 用CALIS过程实现路径分析	(82)

5.3.3	如何处理非同质资料的思考	(85)	第8章	典型相关分析	(112)
5.3.4	用逐步多重线性回归分析方法分析例5-2的资料	(87)	8.1	问题与数据结构	(112)
5.4	本章小结	(88)	8.1.1	实例	(112)
第6章	主成分分析	(89)	8.1.2	对数据结构的分析	(112)
6.1	问题与数据结构	(89)	8.1.3	分析目的与统计分析方法的选择	(112)
6.1.1	实例	(89)	8.2	典型相关分析内容简介	(113)
6.1.2	对数据结构的分析	(89)	8.2.1	典型相关分析概述	(113)
6.1.3	分析目的与统计分析方法的选择	(90)	8.2.2	适合进行典型相关分析的数据结构	(113)
6.2	主成分分析内容简介	(90)	8.2.3	典型相关变量和典型相关系数的定义及解法	(113)
6.2.1	主成分分析概述	(90)	8.2.4	典型相关系数的假设检验	(115)
6.2.2	主成分分析的基本原理	(90)	8.2.5	典型冗余分析	(116)
6.2.3	主成分的计算步骤及性质	(91)	8.2.6	CANCORR过程简介	(117)
6.2.4	与主成分分析有关的其他内容	(94)	8.3	典型相关分析的应用	(118)
6.2.5	PRINCOMP过程简介	(94)	8.3.1	SAS程序	(118)
6.3	主成分分析的应用	(96)	8.3.2	主要分析结果及解释	(119)
6.3.1	SAS程序	(96)	8.4	本章小结	(125)
6.3.2	主要分析结果及解释	(98)	第9章	多元多重线性回归分析	(126)
6.4	本章小结	(101)	9.1	问题与数据结构	(126)
第7章	变量聚类分析	(102)	9.1.1	实例	(126)
7.1	问题与数据结构	(102)	9.1.2	对数据结构的分析	(126)
7.1.1	实例	(102)	9.1.3	统计分析目的与统计分析方法的选择	(126)
7.1.2	对数据结构的分析	(102)	9.2	多元多重线性回归分析内容简介	(127)
7.1.3	分析目的与统计分析方法的选择	(102)	9.2.1	基于普通最小二乘法筛选自变量的思路	(127)
7.2	变量聚类分析内容简介	(103)	9.2.2	何为偏最小二乘回归分析	(127)
7.2.1	变量聚类分析的概念	(103)	9.2.3	偏最小二乘回归分析的基本原理与步骤	(127)
7.2.2	变量聚类分析的聚类统计量	(103)	9.3	偏最小二乘回归分析的应用	(128)
7.2.3	适合进行变量聚类分析的数据结构	(103)	9.3.1	问题与数据结构	(128)
7.2.4	VARCLUS过程简介	(103)	9.3.2	用两种检验方法来决定抽取几对主成分变量	(128)
7.3	变量聚类分析的应用	(106)	9.4	如何获得较多统计量的计算结果	(133)
7.3.1	SAS程序	(106)	9.5	本章小结	(136)
7.3.2	主要分析结果及解释	(107)			
7.4	本章小结	(111)			

第 10 章	探索性因子分析	(137)
10.1	问题与数据结构	(137)
10.1.1	实例	(137)
10.1.2	对数据结构的分析	… (137)
10.1.3	分析目的与统计分析方法 的选择	(137)
10.2	探索性因子分析内容简介	… (138)
10.2.1	概述	(138)
10.2.2	探索性因子分析的数学 模型	(138)
10.2.3	探索性因子分析中载荷 矩阵 A 的统计意义	… (139)
10.2.4	因子载荷矩阵 A 的估计 方法	… (140)
10.2.5	公因子个数的确定 方法	(141)
10.2.6	因子旋转	(142)
10.2.7	因子得分	(142)
10.2.8	FACTOR 过程简介	… (143)
10.3	探索性因子分析的应用	… (145)
10.3.1	SAS 程序	… (145)
10.3.2	主要分析结果及解释	… (146)
10.4	本章小结	… (152)
第 11 章	证实性因子分析	(154)
11.1	问题与数据结构	… (154)
11.1.1	实例	(154)
11.1.2	对数据结构的分析	… (154)
11.1.3	分析目的与统计分析方法 的选择	(155)
11.2	证实性因子分析简介	… (155)
11.2.1	概述	… (155)
11.2.2	CALIS 过程简介	… (155)
11.3	证实性因子分析的应用	… (156)
11.3.1	SAS 程序	… (156)
11.3.2	主要分析结果及解释	… (158)
11.4	本章小结	… (160)
第 12 章	结构方程模型分析	(161)
12.1	问题与数据结构	… (161)
12.1.1	实例	… (161)
12.1.2	对数据结构的分析	… (161)
12.1.3	分析目的与统计分析方法 的选择	(162)
12.2	结构方程模型简介	… (162)
12.2.1	概述	… (162)
12.2.2	基本原理	… (163)
12.3	结构方程模型分析的应用	… (164)
12.3.1	SAS 程序	… (164)
12.3.2	主要分析结果及解释	… (165)
12.4	本章小结	… (168)

第3篇 样品间亲疏、优劣或相对位置分析

第 13 章	传统综合评价	(169)
13.1	问题与数据结构	(169)
13.1.1	实例	(169)
13.1.2	对数据结构的分析	...	(170)
13.1.3	分析目的与统计分析方法 的选择	(171)
13.2	传统综合评价方法内容 介绍	(172)
13.2.1	综合评分法	(172)
13.2.2	Topsis 法	(173)
13.2.3	层次分析法	(174)
13.2.4	RSR 综合评价法	(176)
13.3	传统综合评价方法的应用	(177)
13.3.1	用综合评分法对例 13-1 的 资料进行综合评价	...	(177)
13.3.2	用 Topsis 法对例 13-2 的 资料进行综合评价	...	(181)
13.3.3	用层次分析法对例 13-3 的 资料进行综合评价	...	(183)
13.3.4	用 RSR 综合评价法对 例 13-4 的资料进行 综合评价	(186)
13.4	本章小结	(188)
第 14 章	无序样品聚类分析	(189)
14.1	问题与数据结构	(189)

14.1.1 实例	(189)	15.2.1 概述	(214)
14.1.2 对数据结构的分析	… (189)	15.2.2 有序样品聚类分析的基本概念	(214)
14.1.3 分析目的与统计分析方法的选择	(189)	15.2.3 有序样品聚类分析的计算原理	(215)
14.2 无序样品聚类分析简介	… (190)	15.3 有序样品聚类分析的应用	(217)
14.2.1 概述	… (190)	15.3.1 SAS 程序	… (217)
14.2.2 无序样品聚类分析方法分类	(190)	15.3.2 主要分析结果及解释	… (219)
14.2.3 类的特征与个数的确定	(191)	15.4 本章小结	… (222)
14.2.4 无序样品聚类分析的计算原理	(193)	第 16 章 多维尺度分析	… (223)
14.2.5 CLUSTER 过程等简介	(200)	16.1 问题与数据结构	… (223)
14.3 无序样品聚类分析的应用	… (204)	16.1.1 实例	… (223)
14.3.1 SAS 程序	… (204)	16.1.2 对数据结构的分析	… (224)
14.3.2 主要分析结果及解释	… (206)	16.1.3 分析目的与统计分析方法的选择	… (224)
14.4 本章小结	… (212)	16.2 多维尺度分析内容简介	… (224)
第 15 章 有序样品聚类分析	… (213)	16.2.1 概述	… (224)
15.1 问题与数据结构	… (213)	16.2.2 度量型多维尺度分析的计算原理	… (224)
15.1.1 实例	… (213)	16.2.3 非度量型多维尺度分析的计算原理	… (227)
15.1.2 对数据结构的分析	… (214)	16.3 多维尺度分析的应用	… (228)
15.1.3 分析目的与统计分析方法的选择	(214)	16.3.1 SAS 程序	… (228)
15.2 有序样品聚类分析内容简介	… (214)	16.3.2 主要分析结果及解释	… (229)
		16.4 MDS 过程简介	… (231)
		16.5 本章小结	… (233)

第 4 篇 样品与变量或原因与结果之间的关联性分析

第 17 章 定量资料对应分析	… (234)	17.2.3 定量资料对应分析的实施步骤	… (236)
17.1 问题与数据结构	… (234)	17.3 定量资料对应分析的应用	… (238)
17.1.1 实例	… (234)	17.3.1 SAS 程序	… (238)
17.1.2 对数据结构的分析	… (234)	17.3.2 主要分析结果及解释	… (238)
17.1.3 分析目的与统计分析方法的选择	(235)	17.4 本章小结	… (240)
17.2 定量资料对应分析简介	… (235)	第 18 章 定性资料对应分析	… (241)
17.2.1 概述	… (235)	18.1 问题与数据结构	… (241)
17.2.2 定量资料对应分析的基本原理	(235)	18.1.1 实例	… (241)
		18.1.2 对数据结构的分析	… (241)

18.1.3	分析目的与统计分析方法 的选择	(242)	19.1.3	统计分析目的与分析方法 的选择	(248)
18.2	定性资料对应分析内容 简介	(242)	19.2	Shannon 信息量分析内容 简介	(248)
18.3	定性资料对应分析的 应用	(242)	19.2.1	概述	(248)
	18.3.1 SAS 程序	(242)	19.2.2	Shannon 信息量分析的基本 原理	(248)
	18.3.2 主要分析结果及解释	(243)	19.3	Shannon 信息量分析的 应用	(250)
18.4	本章小结	(246)	19.3.1	对例 19-1 的资料进行 Shannon 信息量分析	(250)
第 19 章	Shannon 信息量分析	(247)	19.3.2	对例 19-2 的资料进行 Shannon 信息量分析	(251)
19.1	问题与数据结构	(247)	19.4	本章小结	(252)
	19.1.1 实例	(247)			
	19.1.2 对数据结构的分析	(248)			

第 5 篇 数据挖掘与分析

第 20 章	决策树分析	(253)	22.1.1	数据挖掘的背景	(295)
20.1	决策树简介	(253)	22.1.2	数据挖掘的基本概念	(295)
20.2	决策树的基本原理	(253)	22.1.3	数据挖掘任务的分类	(295)
20.3	决策树种类及决策树构造 思路	(254)	22.1.4	数据挖掘的应用	(296)
20.4	递归分割的分裂准则	(255)	22.2	SAS 企业数据挖掘器介绍	(296)
20.5	变量重要性检测	(259)	22.3	关联规则与序列规则	(296)
20.6	实际应用与结果解释	(259)	22.3.1	关联规则分析	(296)
20.7	用数据挖掘模块近似实现各种 决策树算法	(272)	22.3.2	关联规则挖掘实例 分析	(297)
20.8	本章小结	(273)	22.3.3	序列规则分析	(301)
第 21 章	神经网络分析	(274)	22.3.4	序列规则挖掘实例 分析	(301)
21.1	前馈型神经网络简介	(274)	22.4	分类预测	(305)
21.2	多层感知器的学习	(276)	22.4.1	数据准备	(306)
21.3	模型过拟合	(279)	22.4.2	数据探索与数据转换	(306)
21.4	模型复杂性的评价	(279)	22.4.3	构造预测模型	(307)
	21.4.1 模型泛化能力(Generalization) 的评价	(279)	22.4.4	模型评估与数据预测	(308)
	21.4.2 模型选择的标准	(281)	22.5	本章小结	(308)
21.5	实际应用与结果解释	(281)	第 23 章	基因表达谱分析	(309)
21.6	本章小结	(294)	23.1	基因表达谱的概念	(309)
第 22 章	数据挖掘与分析	(295)	23.2	基因表达谱的数据获取及 标准化	(309)
22.1	数据挖掘的基本概念	(295)	23.2.1	基因表达谱的数据 获取	(309)

23.2.2	基因表达数据的 标准化	(310)	24.1.1	生物学问题	(323)
23.3	基因表达数据分析技术	(311)	24.1.2	生物数据	(323)
23.3.1	差异表达基因的筛选	… (311)	24.1.3	计算工具	(323)
23.3.2	基因表达的聚类分析 方法	(311)	24.2	统计学在生物信息学中的 应用	(324)
23.4	基因调控网络分析	(320)	24.2.1	基于基因表达谱的样本 分型研究	(324)
23.5	本章小结	(322)	24.2.2	基于基因表达谱的样本 分类研究	(330)
第24章	生物信息分析	(323)	24.3	本章小结	(334)
24.1	生物信息学定义	(323)			

第6篇 遗传资料统计分析

第25章	用SAS实现遗传资料统计 分析	(335)	第26章	遗传流行病学资料的统计 分析	(352)
25.1	SAS/Genetics简介	(335)	26.1	基因、基因型频率测定与 哈代-温伯格(Hardy-Weinberg)平衡 定律的验证	(352)
25.2	ALLEL、HAPLOTYPE和 HTSNP过程简介	(336)	26.1.1	问题与数据	(352)
25.2.1	数据格式	(336)	26.1.2	SAS程序中重要内容的 说明	(352)
25.2.2	ALLEL过程的语法 结构	(338)	26.1.3	主要分析结果及解释	… (353)
25.2.3	HAPLOTYPE过程的 语法结构	(341)	26.2	连锁不平衡与单体型分析	(353)
25.2.4	HTSNP过程的语法结构 及其应用	(343)	26.2.1	问题与数据	(354)
25.3	利用CASECONTROL和FAMILY 进行关联分析	(344)	26.2.2	SAS程序中重要内容的 说明	(354)
25.3.1	CASECONTROL过程的 语法结构	(344)	26.2.3	主要分析结果及解释	… (354)
25.3.2	FAMILY过程的语法结构 及其应用	(345)	26.3	多位点基因型与疾病关联 分析	(355)
25.4	亲缘系数和近交系数	(347)	26.3.1	问题与数据	(355)
25.5	结果校正和图形输出	(349)	26.3.2	SAS程序中重要内容的 说明	(356)
25.5.1	平滑处理和多重检验 校正	(349)	26.3.3	主要分析结果及解释	… (356)
25.5.2	PSMOOTH过程的语法结构 及其应用	(349)	26.4	标签SNP的确认与SAS 程序	(357)
25.5.3	%TPLOT宏及其应用	… (350)	26.4.1	问题与数据	(357)
25.6	本章小结	(351)	26.4.2	SAS程序中重要内容的 说明	(357)
			26.4.3	主要分析结果及解释	… (358)

26.5 一般人群病例对照遗传资料的 关联分析 (358)	26.6 家系数据的关联分析 (360)
26.5.1 问题与数据 (359)	26.6.1 问题与数据 (360)
26.5.2 SAS 程序中重要内容的 说明 (359)	26.6.2 SAS 程序中重要内容的 说明 (361)
26.5.3 主要分析结果及解释 ... (360)	26.6.3 主要分析结果及解释 ... (362)
	26.7 本章小结 (362)

附录

附录 A 胡良平统计学专著及配套软件 简介 (364)

第1篇 统计分析基础

第1章 应确保数据是值得分析的

一提起统计分析，很多人脑海中闪现出来的第一个问题是：应该选择什么统计分析方法分析面前的资料？这是对统计学一知半解的一个明显标志！统计分析应该被使用在“值得分析的资料”上，而不是被使用在任何资料上，更不是被使用在不知任何专业背景的所谓“大数据”上！

1.1 什么是数据和/或统计资料

1.1.1 数据不等于统计资料

【问题】下面列出了300个数据，请问：利用这些数据可以做哪些统计分析？

0.6	0.8	0.9	1.1	1.2	1.4	1.5	1.6	1.8	1.9	2.1	2.2	2.4	2.5
2.6	2.8	2.9	3.1	3.2	3.4	3.5	3.6	3.8	3.9	10.1	10.2	10.4	10.1
4.2	4.2	4.3	4.4	4.5	4.5	4.6	4.7	4.8	4.8	4.9	5.0	5.1	5.2
5.2	5.3	5.4	5.5	5.5	5.6	5.7	5.8	5.8	5.9	6.0	6.1	6.2	6.2
6.3	6.4	6.5	6.5	6.6	6.7	6.8	6.8	6.9	7.0	7.1	7.2	7.2	7.3
7.4	7.5	7.5	7.6	7.7	7.8	8.1	8.1	8.2	8.3	8.3	8.4	8.5	8.5
8.6	8.7	8.7	8.8	8.9	8.9	9.0	9.1	9.1	9.2	9.3	9.3	9.4	9.5
9.5	9.6	9.7	9.7	9.8	9.9	9.9	10.0	10.1	10.1	10.2	10.3	10.3	10.4
10.5	10.5	10.6	10.7	10.7	10.8	10.9	10.9	11.0	11.1	11.1	11.2	11.3	11.3
11.4	11.5	11.5	11.6	11.7	11.7	11.8	11.9	12.1	12.2	12.2	12.3	12.4	12.5
12.5	12.6	12.7	12.8	12.8	12.9	13.0	13.1	13.2	13.2	13.3	13.4	13.5	13.5
13.6	13.7	13.8	13.8	13.9	14.0	14.1	14.2	14.2	14.3	14.4	14.5	14.5	14.6
14.7	14.8	14.8	14.9	15.0	15.1	15.2	15.2	15.3	15.4	15.5	15.5	15.6	15.7
15.8	15.8	16.1	16.2	16.3	16.3	16.4	16.5	16.6	16.7	16.8	16.8	16.9	17.0
17.1	17.2	17.3	17.3	17.4	17.5	17.6	17.7	17.8	17.8	17.9	18.0	18.1	18.2
18.3	18.3	18.4	18.5	18.6	18.7	18.8	18.8	18.9	19.0	19.1	19.2	19.3	19.3

19.4	19.5	19.6	19.7	19.8	20.2	20.3	20.5	20.7	20.8	21.0	21.2	21.3	21.5
21.7	21.8	22.0	22.2	22.3	22.5	22.7	22.8	23.0	23.2	23.3	23.5	23.7	24.2
24.4	24.7	24.9	25.1	25.3	25.6	25.8	26.0	26.2	26.4	26.7	26.9	27.1	27.3
28.3	28.6	28.9	29.1	29.4	29.7	30.0	30.3	30.6	30.9	31.1	31.4	32.4	32.8
33.2	37.6	34.0	34.4	34.8	35.2	35.6	27.8	27.7	37.7	38.2	38.8	39.3	9.9
40.6	9.9	41.9	42.5	45.3	46.5								

【解答】 研究奥林匹克数学的人们会试图找出任何几个相邻数据之间的内在关系，最理想的结果是“推导”出一个“递推公式”，给定数据所在的序号 i ，就能推算出 x_i 等于多少。而纯数理统计学家们会不假思索地说，可以绘制频数直方图或箱式图，形象化地呈现这些定量数据的频数分布规律；还可计算三种区间（即置信区间、容许区间和预测区间），进而做出某些统计推断；若给定了标准值，还可对样本所抽自的总体均值与给定标准值之间差别是否具有统计学意义进行差异性分析，简称为单样本均值的假设检验。

其实，上面两种人所做的“统计工作”都有一个重要的前提条件，那就是此数据是值得分析的。然而，并非所有数据都是值得分析的！换一句话说，对数据进行分析之前，必须证明前提条件是成立的，否则，就是数字游戏！

事实上，值得分析的数据最起码要明确如下几条前提条件：收集这些数据的目的是什么？数据取自的总体是什么？这些数据分别在哪些条件下收集到的？数据的专业含义是什么？数据的名称是什么？数据的性质（定量或定性）是什么？若属于定量数据其单位是什么？测量数据的精确度有多高？由此可知，不明确这些前提条件的“一串数”可以称为“数据”；而明确这些前提条件的“一串数”才可以称为“统计资料”。

现在，全球有很多人在炒作“大数据”。大多数人或研究课题所呼喊的“大数据”，其实对应的正是前述提及的“数据”，而不是值得分析的“统计资料”！

1.1.2 统计资料的要素

1. 统计资料的基本要素

统计资料的基本要素有三条，即变量名称、专业含义和度量单位。例如，前面提及的那300个数据。若用变量 x 代表每个数据，且已知其专业含义为“尿汞值”，又知其单位是“ $\mu\text{g}/\text{L}$ ”，此时就可以说，这300个数据提供的信息已基本清晰了，它们基本上可以被称为一组统计资料了。

2. 统计资料的全部要素

严格地说，仅仅知道一组统计资料的基本要素还是不能确定它是否值得分析。例如，前面提及的那300个数据，即使知道了它们是 x （尿汞值， $\mu\text{g}/\text{L}$ ）的具体取值，但不知道收集这些数据的目的是什么、这些数据所测自的受试对象是什么、这些受试对象是否来自同一个总体。如果研究目的是估计某地区正常成年人尿汞值的80%参考值范围，这个地区全部正常成年人构成了一个特定的总体，而这300个数据所测自的受试对象正是从这个总体中完全随机或按重要非试验因素分层随机抽取的（目的是为了提高样本对于总体的代表性），那么便有理由认为，对这300个数据组成的统计资料进行各种统计表达与描述和各种统计分析，是值得的。

由此可知，统计资料的全部要素除了它应具备基本要素外，还应具备：有明确定义的与特