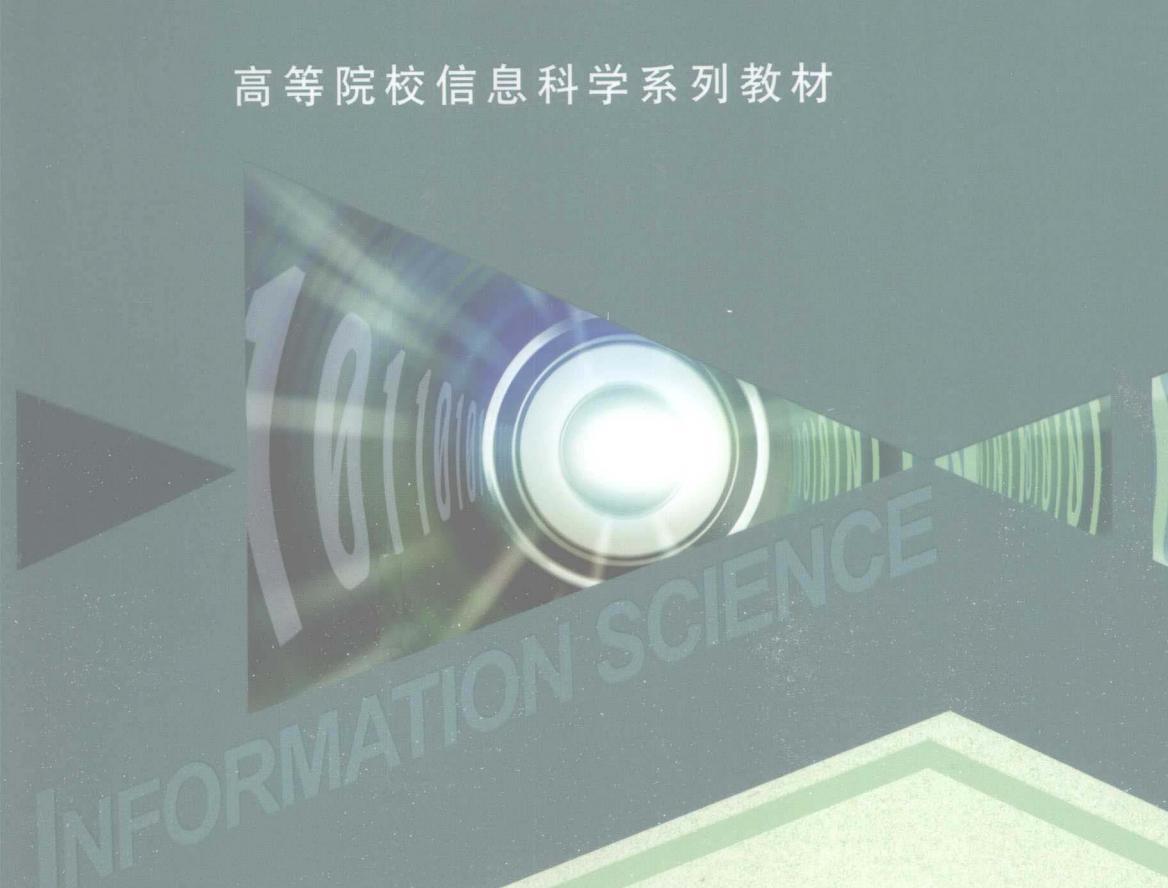


高等院校信息科学系列教材



INFORMATION SCIENCE

数据分析

(第二版)

范金城 梅长林 主编



科学出版社
www.sciencep.com

高等院校信息科学系列教材

数据分析

(第二版)

范金城 梅长林 主编

2002年全国优秀教材评选委员会推荐教材
全国高等学校教材选用书

全国普通高等学校教材选用书

科学出版社

交大林盛·首阳对联

北京

010-62053355 13601121208 13601121209

林建林 内容简介

本书介绍了数据分析的基本内容与方法，其特点是既重视数据分析的基本理论与方法的介绍，又强调应用计算机软件 SAS 进行实际分析和计算能力的培养。主要内容有：数据描述性分析、非参数秩方法、回归分析、主成分分析与因子分析、判别分析、聚类分析、时间序列分析以及常用数据分析方法的 SAS 过程简介。本书每章末附有大量实用、丰富的习题，并要求学生独立上机完成。

本书可作为高等院校信息科学及数理统计专业的本科生教材，也可供有关专业的研究生及工程技术人员参考。

图书在版编目(CIP)数据

数据分析 / 范金城, 梅长林主编. —2 版. —北京: 科学出版社, 2010.2
(高等院校信息科学系列教材)

ISBN 978-7-03-026372-8

I. ①数… II. ①范… ②梅… III. ①统计数据—统计分析(数学)—高等学校—教材 IV. ①O212.1

中国版本图书馆 CIP 数据核字 (2010) 第 006026 号

责任编辑: 鞠丽娜/责任校对: 赵 燕

责任印制: 吕春珉/封面设计: 三函设计

科学出版社出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

新蕾印刷厂印刷

科学出版社发行 各地新华书店经销

*

2002 年 7 月第 一 版 开本: B5 (720×1000)

2010 年 2 月第 二 版 印张: 23

2010 年 2 月第八次印刷 字数: 460 000

印数: 13 001—16 000

定价: 38.00 元

(如有印装质量问题, 我社负责调换(环伟))

销售部电话 010-62134988 编辑部电话 010-62138978-8002

版权所有, 侵权必究

举报电话: 010-64030229; 010-64034315; 13501151303

序　　言

1998年教育部进行高校专业调整时,设立了“信息与计算科学”专业。该专业的设立,受到很多高等院校的热烈响应。据不完全统计,几年来已有约280所院校招收了该专业的本科生,其中大部分院校计划开设信息科学方面的系列课程。

为了配合高等院校在学科专业设置上的改革与深化,来自几十所高等院校有关专业的部分领导和教师,于1999年、2000年召开了第一、二届“信息科学专业发展与学术研讨会”,与会者热烈讨论并探讨了许多关于信息学科的学科发展和建设的基本问题。会议一致认为教材建设是目前最为紧迫的任务,因此成立了教材编审协调组来组织该系列教材的编写。

2001年教材编写协调组召集了有多位经验丰富的教师和出版社参加的教材建设会议。会议明确了教材建设是一项长期的工作,并决定首先编写和出版本套教材来满足近期急需。为了保证教材的质量,会议对每本教材的要求、内容和大纲进行了具体研讨,并请具有多年教学经验的重点院校教授担任各教材的负责人。

为了贴近教学的实际,每本教材都配有习题或思考题,同时对内容也作了结构化安排,以便教师能根据实际情况部分选讲。本套教学用书不仅适用于教学,也可供相关读者参考。

在教材编写和出版过程中,作者对内容的取舍、章节的安排、结构的设计以及表达方式等方面多方听取意见,并进行了反复修改。在感谢作者们辛勤劳动的同时,编委会还特别感谢科学出版社的鞠丽娜编辑,她不辞辛劳,在统筹印刷出版、督促进度、征求意见、组织审校等方面做了大量工作。这套教材能在保证质量的前提下及时与读者见面,是和她的努力分不开的。

从长远的教学角度考虑,为了适应不同类型院校、不同要求的课程需要,教材编审协调组将不断组织教材的修订、编写(译),从而使信息科学教学用书做到逐步充实、完善、提高和多样化。在此衷心希望采用该系列用书的教师、学生和读者对书中存在的问题及时提出修改意见和建议。

高等院校信息科学系列教材编委会

第二版前言

数据分析是信息科学专业本科生重要的必修课。本书是高等院校信息科学专业本科生教材，也适用于数理统计专业的本科生。本书的特点是既重视数据分析的基本理论与方法的介绍，又重视应用 SAS 软件进行实际的分析计算。

在本书第一版的基础上，我们对第二版内容进行了精选。本书第二版的主要内容是：数据的描述性分析、非参数秩方法、回归分析、主成分分析与因子分析、判别分析、聚类分析、时间序列分析。数据的描述性分析与非参数秩方法讨论数据的数据特征与分布的描述，数据统计推断的非参数秩方法。多元数据分析是数据分析极为重要的方面，在经济、工业、农业、国防、科学技术等领域有广泛的应用。除回归分析、主成分分析、判别分析、聚类分析外，本书第二版补充了因子分析的内容。时间序列分析也是数据分析极为重要的方面，系统介绍时间序列分析的基本内容是本书的特点之一。

本教材的计划学时约 72 学时，内容的选择是模块式的，各校可以根据具体情况予以选择。其中，数据的描述性分析、回归分析是基本的，主成分分析与因子分析、判别分析、聚类分析可以全学或选学一部分，非参数秩方法与时间序列分析也可以全学或选学一部分。

本书第二版精选了各章的例题与习题。在应用题方面，选用了许多近年的统计数据的例题与习题。

本书与 SAS 软件系统紧密结合。SAS 软件系统在数据分析与统计分析领域被誉为国际标准软件系统，并被广泛应用于各个领域。通过对各种典型例题采用各种不同方法进行分析计算，以培养学生分析、解决实际问题的能力。目前，主要有 SAS 8.2、SAS 9.0、SAS 9.1 软件系统，拥有上述任何一种 SAS 软件系统皆可进行本书的教学。大部分习题要通过 SAS 软件系统计算完成。本书第 8 章较系统地介绍与本书有关的 SAS 过程，并结合本书部分例题进行编程。第 8 章介绍了建立 SAS 数据集与 SAS 数据库的方法，强调调用 SAS 数据库中的 SAS 数据集进行 SAS 编程的方法。这种调用 SAS 数据集的方法使得 SAS 程序具有普遍性与通用性。

本书第二版备有电子教案。电子教案的内容包括：①《数据分析幻灯片》，包含各章教学幻灯片；②《数据分析例题与习题》，包含全部例题与习题的 SAS 程序；③《SJFX》，包含书中全部例题与习题用的 SAS 数据集。电子教案的推出将为教学带来极大方便。

本书第二版的出版得到科学出版社编辑鞠丽娜同志的指导与帮助,谨表诚挚的感谢。本书第二版又得到北京金晨晖科技有限公司陈强同志的协助支持,谨表示感谢。

言龍體二範

作者

史學降自南魏頌華高僧計本,累過尊師要前北林本山分寺恩言望 2009年12月

印傳令謹謹重刻皇立林師計本,于林本山分寺恩言望傳氣林師計本

真性傳食頭利突厥計掛汗沙以由而刻重文,是今之唐式計掛汗本其

內卷主印傳二漢計本,真蘇丁首載容山識。字勢清秀,上顯某山道一宮計本余

眼慨,林食丁因已傳食食如生,林食其印,志氏傳達非,傳食掛汗計掛汗,計容

達掛汗計掛汗,計容達非已傳食掛汗計掛汗,計容達非已傳食掛汗,計容

前 言

本书是高等院校信息科学专业本科生教材,也适用于数理统计专业的本科生。本书包括数据分析的主要内容:数据描述性分析、非参数方法、回归分析、主成分分析、判别分析、聚类分析、时间序列分析、Bayes 统计分析和常用数据分析方法的 SAS 过程简介。本书的特点是既重视数据分析的基本理论与方法的介绍,又重视应用 SAS 软件进行实际的分析计算。

数据分析是信息科学专业本科生重要的必修课,因此本书重视数据分析的基本理论与方法的介绍,详细叙述基本内容及算法.本书的内容经过精选,对选入的内容详细予以介绍,并力求反映新颖内容.其中数据描述性分析力求体现“让数据自身说话”;非参数方法模型具有一般性.多元数据分析是数据分析极为重要的方面,书中主要介绍回归分析、主成分分析、判别分析、聚类分析.因时间序列分析在自然、技术、经济等领域有极广泛的应用,因此介绍其基本内容. Bayes 统计分析作为数据分析的重要方法,也给以简单介绍.

本教材计划学时是 72 学时。内容选择是模块式的，各校可以根据具体情况予以选择。其中，数据描述性分析是基本的，必选的。非参数方法可全学或选学一部分。关于多元数据分析的几章中，回归分析是基本的；主成分分析、判别分析、聚类分析可全学或选学一部分；时间序列分析、Bayes 统计分析是两个单独的模块，可全学或选学其中的一个模块。

本书与 SAS 软件系统紧密结合.书中大多数例题都由 SAS 软件分析计算. SAS 系统是大型集成应用软件系统,在数据分析与统计分析领域被誉为国际标准软件系统,并被广泛应用于各个领域.通过对典型例题采用各种不同方法进行分析计算,以培养学生分析、解决实际问题的能力.练习中有相当大的部分要通过 SAS 软件或其软件计算完成.本书的第 9 章“常用数据分析方法的 SAS 过程简介”较系统地介绍了与本书内容有关的 SAS 过程,并结合本书部分例题进行编程.若用 SAS 软件进行计算,应讲授该章,并让学生掌握计算技能.对于采用其他软件进行计算的情况,当然要根据具体情况介绍其他软件.因为本书主要介绍数据分析方法,即使使用其他软件,仍可使用本书.由于 SAS 系统各种数据分析方法大都输出相应的 p 值,故本书精简大部分统计用表,仅列出几个常用统计数值表.

书中第1、5~8章由范金城编写,第2~4、9章由梅长林编写。

由于作者的水平所限,书中难免存在不妥之处,欢迎读者批评指正。

作 者 言 面

于西安交通大学

在编写本专业教材时,我们力求做到科学性与实用性相结合,突出本学科的特点,并结合工程实际,对一些概念和术语进行深入浅出的讲解,同时着重介绍国内外的研究成果,使读者能较全面地了解本学科的基本理论和方法,掌握本学科的主要研究方向,从而提高分析问题和解决问题的能力。

本书在编写过程中参考了国内外许多有关文献,并吸收了国内外学者的研究成果。在编写过程中,我们特别注意了与本学科相关的其他学科知识的综合运用,力求做到理论与实践相结合,使读者能够更好地理解本学科的内容。同时,我们还注重了与工程实际的结合,力求使读者能够将所学的知识应用于工程实践中去。

本书在编写过程中,得到了许多老师的帮助和支持,在此表示衷心的感谢。同时,我们还特别感谢了本书的审稿人,他们的意见和建议对我们编写本书起到了重要的指导作用。

最后,我们希望本书能够成为广大读者学习和工作的参考书,同时也希望得到广大读者的批评和指正。

由于时间仓促,书中难免有疏忽和不足之处,敬请各位读者批评指正。同时,我们希望本书能够成为广大读者学习和工作的参考书,同时也希望得到广大读者的批评和指正。

最后,我们希望本书能够成为广大读者学习和工作的参考书,同时也希望得到广大读者的批评和指正。

由于时间仓促,书中难免有疏忽和不足之处,敬请各位读者批评指正。同时,我们希望本书能够成为广大读者学习和工作的参考书,同时也希望得到广大读者的批评和指正。

最后,我们希望本书能够成为广大读者学习和工作的参考书,同时也希望得到广大读者的批评和指正。

主要参考文献

- 安鸿志等. 1983. 时间序列的分析与应用. 北京: 科学出版社.
- 常学将等. 1993. 时间序列分析. 北京: 高等教育出版社.
- 范金城. 2008. SAS 数据分析范例. 西安: 西安交通大学出版社.
- 方开泰. 1989. 实用多元统计分析. 上海: 华东师范大学出版社.
- 高惠璇. 2005. 应用多元统计分析. 北京: 人民邮电出版社.
- 高惠璇等. 1997. SAS 系统——Base SAS 软件系统使用手册. 北京: 中国统计出版社.
- 高惠璇等. 1997. SAS 系统——SAS/STAT 软件使用手册. 北京: 中国统计出版社.
- 高惠璇等. 1998. SAS 系统——SAS/ETS 软件使用手册. 北京: 中国统计出版社.
- 何宁等. 2005. 统计分析系统 SAS. 武汉: 武汉大学出版社.
- 洪楠等. 2004. SAS for Windows 统计分析系统教程新编. 北京: 清华大学出版社.
- 胡良平. 2000. 现代统计学与 SAS 应用. 北京: 军事医学科学出版社.
- 胡良平. 2001. Windows SAS 6.12 & 8.0 实用统计分析教程. 北京: 军事医学科学出版社.
- 李东风. 2006. 统计软件出版社. 北京: 人民邮电出版社.
- 马逢时等. 1990. 应用概率统计(下册). 北京: 高等教育出版社.
- 茆诗松等. 2003. 统计手册. 北京: 科学出版社.
- 梅长林等. 2006. 数据分析方法. 北京: 高等教育出版社.
- 施锡铨等. 1997. 数据分析方法. 上海: 上海财经大学出版社.
- 孙文爽等. 1994. 多元统计分析. 北京: 高等教育出版社.
- 汪嘉冈. 2004. SASV8 基础教程. 北京: 中国统计出版社.
- 汪远征等. 2007. SAS 软件与统计应用教程. 北京: 机械工业出版社.
- 王国梁等. 1993. 多变量经济数据统计分析. 西安: 陕西科学技术出版社.
- 王吉利. 2000. SAS 软件与应用统计. 北京: 中国统计出版社.
- 王松桂等. 1999. 线性统计模型——线性回归与方差分析. 北京: 高等教育出版社.
- 王学仁等. 1989. 应用回归分析. 重庆: 重庆大学出版社.
- 王学仁等. 1990. 实用多元统计分析. 上海: 上海科学技术出版社.
- 吴喜之. 1999. 非参数统计. 北京: 中国统计出版社.
- 项静恬等. 1986. 动态数据处理——时间序列分析. 北京: 气象出版社.
- 谢衷洁. 1990. 时间序列分析. 北京: 北京大学出版社.
- 徐利治. 2000. 现代数学手册·随机数学卷. 武汉: 华中科技大学出版社.
- 张尧庭等. 1982. 多元统计分析引论. 北京: 科学出版社.
- 张尧庭等. 1991. 贝叶斯统计推断. 北京: 科学出版社.
- 张尧庭等. 1991. 定性资料的统计分析. 桂林: 广西师范大学出版社.
- Agrest. A. 1988. Categorical Data Analysis. NY: John Wiley.
- Box G E P et al. 1997. 时间序列分析、预测与控制. 北京: 中国统计出版社.
- Hoaglin D C et al. 1998. 探索性数据分析. 北京: 中国统计出版社.
- Johnson R A et al. 1992. Applied Multivariate Statistical Analysis (3rd Edition). NJ: Prentice-Hall Inc.
- Lehmann E L. 1975. Nonparametrics, Statistical Methods Based on Ranks. San Francisco: Holden-Day Inc.
- Neter J et al. 1990. Applied Linear Statistical Models (3rd Edition). Chicago: IRWIN Inc.

目 录

目 录

第1章 数据描述性分析	1
1.1 数据的数字特征	1
1.1.1 均值、方差等数字特征	1
1.1.2 中位数、分位数、三均值与极差	7
1.2 数据的分布	11
1.2.1 直方图、经验分布函数与 QQ 图	12
1.2.2 茎叶图、箱线图及五数总括	16
1.2.3 正态性检验与分布拟合检验	21
1.3 多元数据的数字特征与相关分析	27
1.3.1 二元数据的数字特征及相关系数	27
1.3.2 多元数据的数字特征及相关矩阵	31
1.3.3 总体的数字特征及相关矩阵	33
习题	42
第2章 非参数秩方法	47
2.1 两种处理方法比较的秩检验	47
2.1.1 两种处理方法比较的随机化模型及秩的零分布	48
2.1.2 Wilcoxon 秩和检验	49
2.1.3 总体模型的 Wilcoxon 秩和检验	58
2.1.4 Smirnov 检验	59
2.2 成对分组设计下两种处理方法的比较	63
2.2.1 符号检验	64
2.2.2 Wilcoxon 符号秩检验	66
2.2.3 分组设计下两处理方法比较的总体模型	72
2.3 多种处理方法比较的 Kruskal-Wallis 检验	73
2.3.1 多种处理方法比较中秩的定义及 Kruskal-Wallis 统计量	73
2.3.2 Kruskal-Wallis 统计量的零分布	74
2.4 分组设计下多种处理方法的比较	78
2.4.1 分组设计下秩的定义及其零分布	78
2.4.2 Friedman 检验	78
2.4.3 改进的 Friedman 检验	82
习题	85

第3章 回归分析	89
3.1 线性回归模型	89
3.1.1 线性回归模型及其矩阵表示	89
3.1.2 β 及 σ^2 的估计	90
3.1.3 有关的统计推断	91
3.2 逐步回归法	100
3.3 Logistic 回归模型	108
3.3.1 线性 Logistic 回归模型	108
3.3.2 参数的最大似然估计与 Newton-Raphson 迭代解法	110
3.3.3 Logistic 模型的统计推断	115
习题	120
第4章 主成分分析与因子分析	124
4.1 主成分分析	124
4.1.1 引言	124
4.1.2 总体主成分	125
4.1.3 样本主成分	131
4.2 因子分析	137
4.2.1 引言	137
4.2.2 正交因子模型	138
4.2.3 参数估计方法	141
4.2.4 主成分估计法的具体步骤	143
4.2.5 方差最大的正交旋转	146
4.2.6 因子得分	149
习题	151
第5章 判别分析	155
5.1 距离判别	155
5.1.1 判别分析的基本思想及意义	155
5.1.2 两个总体的距离判别	156
5.1.3 判别准则的评价	160
5.1.4 多个总体的距离判别	163
5.2 Bayes 判别	166
5.2.1 Bayes 判别的基本思想	166
5.2.2 两个总体的 Bayes 判别	167
5.2.3 多个总体的 Bayes 判别	177
5.2.4 逐步判别简介	182

107	习题	183
第6章 聚类分析 192		
206	6.1 距离与相似系数	192
301	6.1.1 聚类分析的基本思想及意义	192
318	6.1.2 样品间的相似性度量——距离	193
338	6.1.3 变量间的相似性度量——相似系数	195
348	6.2 谱系聚类法	198
368	6.2.1 类间距离	198
388	6.2.2 类间距离的递推公式	199
408	6.2.3 谱系聚类法的步骤	201
428	6.2.4 变量聚类	212
448	6.3 快速聚类法	214
468	6.3.1 快速聚类法的步骤	215
488	6.3.2 用 L_m 距离进行快速聚类	223
508	习题	227
第7章 时间序列分析 233		
518	7.1 平稳时间序列	233
538	7.1.1 时间序列分析及其意义	233
558	7.1.2 随机过程概念及其数字特征	233
578	7.1.3 平稳时间序列与平稳随机过程	238
598	7.1.4 平稳性检验及自协方差函数、自相关函数的估计	241
618	7.2 ARMA 时间序列及其特性	243
638	7.2.1 ARMA 时间序列的定义	243
658	7.2.2 ARMA 序列的平稳性与可逆性	246
678	7.2.3 ARMA 序列的相关特性	249
698	7.3 ARMA 时间序列的建模与预报	258
718	7.3.1 ARMA 序列参数的矩估计	258
738	7.3.2 ARMA 序列参数的精估计	261
758	7.3.3 ARMA 模型的定阶与考核	269
778	7.3.4 平稳线性最小均方预报	273
798	7.3.5 ARMA 序列的预报	276
818	7.4 ARIMA 序列与季节性序列	281
838	7.4.1 ARIMA 序列及其预报	281
858	7.4.2 季节性序列及其预报	288
878	习题	295

第8章 常用数据分析方法的SAS过程简介	301
8.1 SAS系统简介	301
8.1.1 建立SAS数据集	302
8.1.2 利用已有的SAS数据集建立新的SAS数据集	307
8.1.3 SAS系统的数学运算符号及常用的SAS函数	310
8.1.4 逻辑语句与循环语句	312
8.2 常用数据分析方法的SAS过程	314
8.2.1 几种描述性统计分析的SAS过程	315
8.2.2 非参数秩方法的SAS过程	323
8.2.3 回归分析的SAS过程	327
8.2.4 主成分分析与因子分析的SAS过程	333
8.2.5 判别分析的SAS过程	335
8.2.6 聚类分析的SAS过程	341
8.2.7 时间序列分析的SAS过程——PROC ARIMA过程	346
8.2.8 SAS系统的矩阵运算——PROC IML过程简介	351
主要参考文献	354
8.8.1 《SAS统计分析教程》	354
8.8.2 《SAS统计分析教程》	354
8.8.3 《SAS统计分析教程》	354
8.8.4 《SAS统计分析教程》	354
8.8.5 《SAS统计分析教程》	354
8.8.6 《SAS统计分析教程》	354
8.8.7 《SAS统计分析教程》	354
8.8.8 《SAS统计分析教程》	354
8.8.9 《SAS统计分析教程》	354
8.8.10 《SAS统计分析教程》	354
8.8.11 《SAS统计分析教程》	354
8.8.12 《SAS统计分析教程》	354
8.8.13 《SAS统计分析教程》	354
8.8.14 《SAS统计分析教程》	354
8.8.15 《SAS统计分析教程》	354
8.8.16 《SAS统计分析教程》	354
8.8.17 《SAS统计分析教程》	354
8.8.18 《SAS统计分析教程》	354
8.8.19 《SAS统计分析教程》	354
8.8.20 《SAS统计分析教程》	354
8.8.21 《SAS统计分析教程》	354
8.8.22 《SAS统计分析教程》	354
8.8.23 《SAS统计分析教程》	354
8.8.24 《SAS统计分析教程》	354
8.8.25 《SAS统计分析教程》	354
8.8.26 《SAS统计分析教程》	354
8.8.27 《SAS统计分析教程》	354
8.8.28 《SAS统计分析教程》	354
8.8.29 《SAS统计分析教程》	354
8.8.30 《SAS统计分析教程》	354
8.8.31 《SAS统计分析教程》	354
8.8.32 《SAS统计分析教程》	354
8.8.33 《SAS统计分析教程》	354
8.8.34 《SAS统计分析教程》	354
8.8.35 《SAS统计分析教程》	354
8.8.36 《SAS统计分析教程》	354
8.8.37 《SAS统计分析教程》	354
8.8.38 《SAS统计分析教程》	354
8.8.39 《SAS统计分析教程》	354
8.8.40 《SAS统计分析教程》	354
8.8.41 《SAS统计分析教程》	354
8.8.42 《SAS统计分析教程》	354
8.8.43 《SAS统计分析教程》	354
8.8.44 《SAS统计分析教程》	354
8.8.45 《SAS统计分析教程》	354
8.8.46 《SAS统计分析教程》	354
8.8.47 《SAS统计分析教程》	354
8.8.48 《SAS统计分析教程》	354
8.8.49 《SAS统计分析教程》	354
8.8.50 《SAS统计分析教程》	354
8.8.51 《SAS统计分析教程》	354
8.8.52 《SAS统计分析教程》	354
8.8.53 《SAS统计分析教程》	354
8.8.54 《SAS统计分析教程》	354
8.8.55 《SAS统计分析教程》	354
8.8.56 《SAS统计分析教程》	354
8.8.57 《SAS统计分析教程》	354
8.8.58 《SAS统计分析教程》	354
8.8.59 《SAS统计分析教程》	354
8.8.60 《SAS统计分析教程》	354
8.8.61 《SAS统计分析教程》	354
8.8.62 《SAS统计分析教程》	354
8.8.63 《SAS统计分析教程》	354
8.8.64 《SAS统计分析教程》	354
8.8.65 《SAS统计分析教程》	354
8.8.66 《SAS统计分析教程》	354
8.8.67 《SAS统计分析教程》	354
8.8.68 《SAS统计分析教程》	354
8.8.69 《SAS统计分析教程》	354
8.8.70 《SAS统计分析教程》	354
8.8.71 《SAS统计分析教程》	354
8.8.72 《SAS统计分析教程》	354
8.8.73 《SAS统计分析教程》	354
8.8.74 《SAS统计分析教程》	354
8.8.75 《SAS统计分析教程》	354
8.8.76 《SAS统计分析教程》	354
8.8.77 《SAS统计分析教程》	354
8.8.78 《SAS统计分析教程》	354
8.8.79 《SAS统计分析教程》	354
8.8.80 《SAS统计分析教程》	354
8.8.81 《SAS统计分析教程》	354
8.8.82 《SAS统计分析教程》	354
8.8.83 《SAS统计分析教程》	354
8.8.84 《SAS统计分析教程》	354
8.8.85 《SAS统计分析教程》	354
8.8.86 《SAS统计分析教程》	354
8.8.87 《SAS统计分析教程》	354
8.8.88 《SAS统计分析教程》	354
8.8.89 《SAS统计分析教程》	354
8.8.90 《SAS统计分析教程》	354
8.8.91 《SAS统计分析教程》	354
8.8.92 《SAS统计分析教程》	354
8.8.93 《SAS统计分析教程》	354
8.8.94 《SAS统计分析教程》	354
8.8.95 《SAS统计分析教程》	354
8.8.96 《SAS统计分析教程》	354
8.8.97 《SAS统计分析教程》	354
8.8.98 《SAS统计分析教程》	354
8.8.99 《SAS统计分析教程》	354
8.8.100 《SAS统计分析教程》	354

第1章 数据描述性分析

1.1 数据的数字特征

数据分析研究的对象是数据. 一元数据是 n 个观测值:

$$x_1, x_2, \dots, x_n.$$

它们可以是从所要研究的对象的全体——总体中取出的, 这 n 个观测值就构成一个样本. 在某些简单的实际问题中, 这 n 个观测值就是所要研究对象的全体. 数据分析的任务就是要对这全部 n 个数进行分析, 提取数据中包含的有用的信息. 如果数据是从总体抽出的样本, 就要分析推断样本中包含的总体的信息.

数据作为信息的载体, 当然要分析数据中包含的主要信息, 即要分析数据的主要特征, 也就是说, 要研究数据的数字特征. 对于数据的数字特征, 要分析数据的集中位置、分散程度、数据的分布是正态还是偏态等. 对于多元数据, 还要分析多元数据的各个分量之间的相关性等.

1.1.1 均值、方差等数字特征

一元数据的数字特征主要有下列几种. 设 n 个观测值为

$$x_1, x_2, \dots, x_n,$$

其中 n 称为样本容量.

1. 均值

均值即是 x_1, x_2, \dots, x_n 的平均数:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (1.1)$$

均值表示数据的集中位置.

2. 方差、标准差与变异系数

方差是描述数据取值分散性的一个度量, 它是数据相对于均值的偏差平方的平均:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (1.2)$$

方差的开方称为标准差. 方差的量纲与数据的量纲不一致, 它是数据量纲的平方,

而标准差的量纲与数据量纲一致. 标准差为

$$s = \sqrt{s^2} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}. \quad (1.3)$$

刻画数据相对分散性的度量是变异系数:

$$CV = 100 \times \frac{s}{\bar{x}} (\%). \quad (1.4)$$

它是一个无量纲的量,用百分数表示.

与均值、方差有关的还有下列数字特征:

校正平方和

$$CSS = \sum_{i=1}^n (x_i - \bar{x})^2.$$

未校正平方和

$$USS = \sum_{i=1}^n x_i^2.$$

3. 偏度与峰度

偏度与峰度是刻画数据的偏态、尾重程度的度量. 它们与数据的矩有关. 数据的矩分为原点矩与中心矩.

k 阶原点矩

$$v_k = \frac{1}{n} \sum_{i=1}^n x_i^k. \quad (1.5)$$

k 阶中心矩

$$u_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k. \quad (1.6)$$

显然,一阶原点矩 v_1 即均值. 二阶中心矩

$$u_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

也称为方差. 而前述的方差 s^2 是 CSS 除以 $n-1$, 是为了保证估计总体方差时的无偏性.

偏度的计算公式为

$$\begin{aligned} g_1 &= \frac{n}{(n-1)(n-2)s^3} \sum_{i=1}^n (x_i - \bar{x})^3 \\ &= \frac{n^2 u_3}{(n-1)(n-2)s^3}, \end{aligned} \quad (1.7)$$

其中 s 是标准差. 偏度是刻画数据对称性的指标. 关于均值对称的数据其偏度为 0, 右侧更分散的数据偏度为正, 左侧更分散的数据偏度为负(图 1.1).

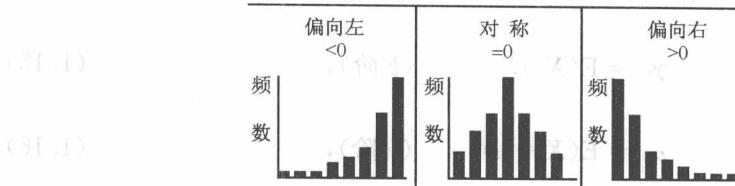


图 1.1

峰度的计算公式是

$$(1.8) \quad g_2 = \frac{n(n+1)}{(n-1)(n-2)(n-3)s^4} \sum_{i=1}^n (x_i - \bar{x})^4 - 3 \frac{(n-1)^2}{(n-2)(n-3)}$$

$$= \frac{n^2(n+1)u_4}{(n-1)(n-2)(n-3)s^4} - 3 \frac{(n-1)^2}{(n-2)(n-3)}. \quad (1.8)$$

当数据的总体分布为正态分布时, 峰度近似为 0; 当分布较正态分布的尾部更分散时, 峰度为正, 否则峰度为负. 当峰度为正时, 两侧极端数据较多; 当峰度为负时, 两侧极端数据较少.

设观测数据是由总体 X 中取出的样本, 总体的分布函数是 $F(x)$. 当 X 为离散分布时, 总体的分布可由概率分布列刻画:

$$p_i = P\{X = x_i\}, \quad i = 1, 2, \dots$$

总体为连续分布时, 总体的分布可由概率密度 $f(x)$ 刻画. 连续分布中最重要的是正态分布, 它的概率密度 $\varphi(x)$ 及分布函数 $\Phi(x)$ 分别为

$$\varphi(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right], \quad (1.9)$$

$$\Phi(x) = \int_{-\infty}^x \varphi(t) dt. \quad (1.10)$$

具有正态分布的总体称为正态总体.

上述数据的数字特征在数据为从某总体随机取出的样本时, 即是样本的数字特征. 与样本数字特征对应的是总体的数字特征, 它们分别是:

总体均值

$$\mu = E(X), \quad (1.11)$$

总体方差

$$\sigma^2 = \text{Var}(X), \quad (1.12)$$

总体标准差

$$\sigma = \sqrt{\text{Var}(X)}, \quad (1.13)$$

总体变异系数

$$\text{大数 n 越大, 总体变异系数 } \gamma = \frac{\sigma}{\mu}, \quad (1.14)$$

总体原点矩

$$\gamma_k = E(X^k) \quad (k \text{ 阶}), \quad (1.15)$$

总体中心矩

$$\mu_k = E(X - \mu)^k \quad (k \text{ 阶}), \quad (1.16)$$

总体偏度

$$G_1 = \frac{\mu_3}{\sigma^3}, \quad (1.17)$$

总体峰度

$$G_2 = \frac{\mu_4}{\sigma^4} - 3. \quad (1.18)$$

(8.1) 这里, 我们对偏度与峰度做进一步的说明.

总体偏度是度量总体分布是否偏向某一侧的指标. 对于对称的分布, 偏度为0. 例如, 对于正态分布, 因 $\mu_3=0$, 故 $G_1=0$. 若总体分布在右侧更为扩展, 偏度为正; 若总体分布在左侧更为扩展, 偏度为负. 图 1.2 表示了偏度为正和偏度为负的概率密度的图像特点. 我们看到, 总体偏度的这一特性与样本偏度的相应特性是相似的.

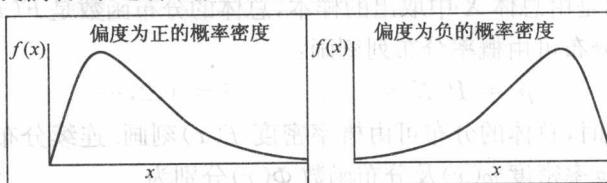


图 1.2

总体峰度是以同方差的正态分布为标准, 比较总体分布尾部分散性的指标. 当总体分布是正态分布时, 因 $\mu_4=3\sigma^4$, 故总体峰度 $G_2=0$. 当 $G_2>0$ 时, 总体分布中极端数值分布范围较广, 此种分布称为粗尾的. 当 $G_2<0$ 时, 两侧极端数据较少, 此种分布称为细尾的(图 1.3).

——峰度为零 - - - 峰度为负 --- 峰度为正

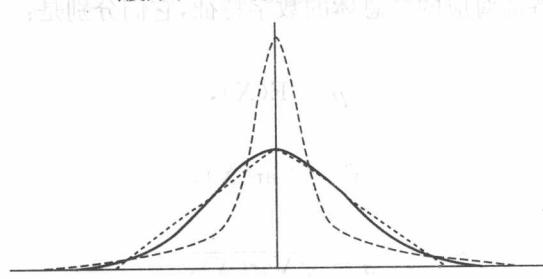


图 1.3

根据统计学的结果, 样本数字特征是相应的总体数字特征的矩估计. 当总体数字特征存在时, 相应的样本数字特征是总体数字特征的相合估计, 从而当 n 较大