

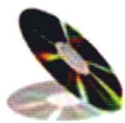
工程设计与分析系列

# SAS

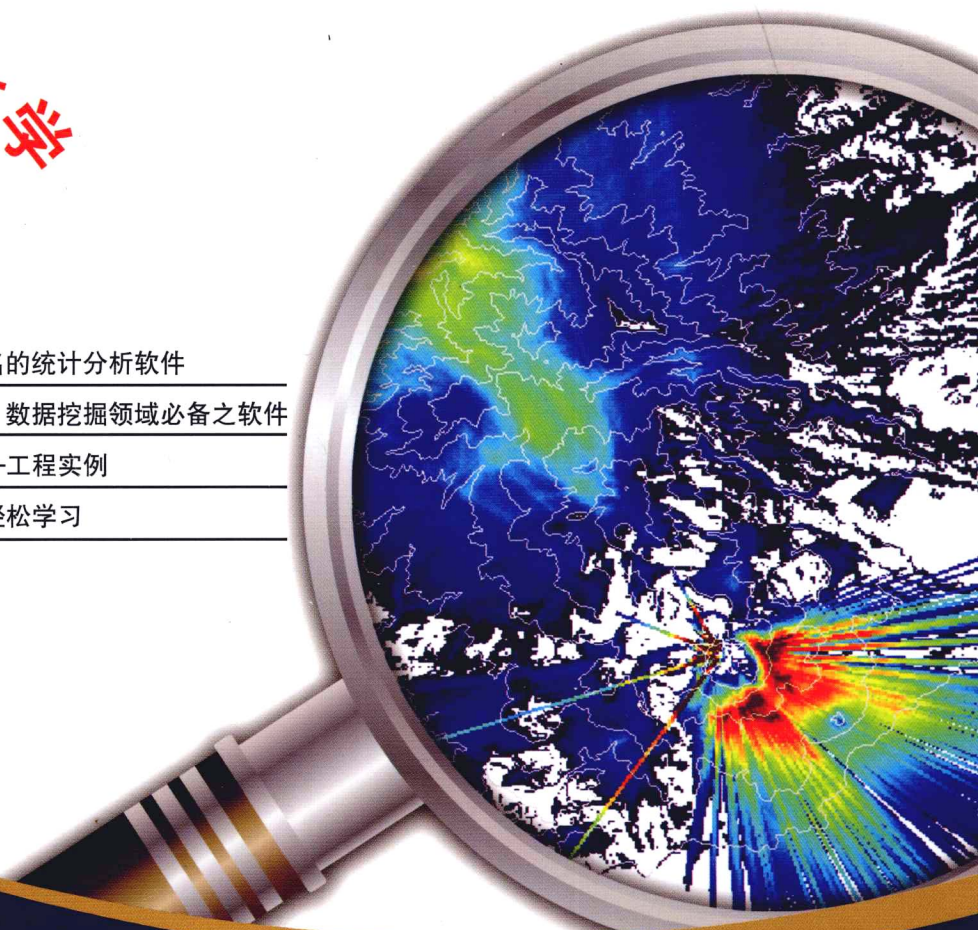
## 统计分析与数据挖掘

谢龙汉 尚涛 编著

视频教学



- ★ SAS——全球最为著名的统计分析软件
- ★ SAS——金融、科研、数据挖掘领域必备之软件
- ★ 基础知识—实训实例—工程实例
- ★ 实例操作视频教学，轻松学习



电子工业出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

<http://www.phei.com.cn>

工程设计与分析系列

# SAS 统计分析与数据挖掘

谢龙汉 尚 涛 编著

电子工业出版社

## 内 容 简 介

本书基于 SAS 9.2 版本编写, 从 SAS 编程出发, 用案例形式介绍 SAS 数据挖掘在各领域的广泛应用。全书分为 SAS 基础篇、提高篇、应用篇, 每章均给出大量分析案例, 具体内容为 SAS 软件与数据挖掘简介, SAS 编程基础, 图形与报表制作, 描述性分析, 假设检验, 回归分析, 方差分析与因子分析, 相关分析与对应分析, 判别分析, 聚类分析, 生存分析, 时间序列分析, 以及 SAS 在具体数据挖掘项目中的应用等。

本书最大特点是抛弃了其他同类书籍中只说理论、缺少案例分析的弊病, 全书给出大量数据挖掘分析案例, 为读者展示 SAS 在数据整合、数据挖掘、商业智能、金融数据分析、金融风险管理等项目中的强大应用技术。

本书读者对象为高等院校统计、信息等相关专业的本科生、研究生, 科研单位的科技人员和企事业单位计算机工作者, 以及数据分析、商业咨询、金融工程、商业智能工作者。

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 统计分析与数据挖掘/谢龙汉, 尚涛编著. —北京: 电子工业出版社, 2012.1  
(工程设计与分析系列)

ISBN 978-7-121-14888-0

I. ①S… II. ①谢… ②尚… III. ①统计分析—应用软件, SAS IV. ①C819

中国版本图书馆 CIP 数据核字 (2011) 第 217233 号

策划编辑: 许存权

责任编辑: 陈韦凯 特约编辑: 刘丽丽

印 刷: 涿州市京南印刷厂

装 订: 涿州市桃园装订有限公司

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: 29.25 字数: 749 千字

印 次: 2012 年 1 月第 1 次印刷

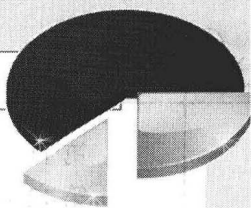
册 数: 4 000 册 定价: 59.00 元 (含 DVD 光盘 1 张)



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 [zltz@phei.com.cn](mailto:zltz@phei.com.cn), 盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

服务热线: (010) 88258888。



## 前 言

SAS 是当今国际最著名的数据分析软件。本书从 SAS 编程出发,用案例的形式介绍 SAS 数据分析在各个领域的广泛应用。


本书是在结合作者多年的使用经验和工作经验的基础上编写的,在编写过程中,突出了以下特点:

(1) 直观易懂性。全书以图解实例的形式介绍基础知识和实例操作,所有的知识块和案例分析都尽可能详细,直观易懂,使读者能够在最短的时间内获取最多的知识。

(2) 先进性。以最新的 SAS 9.2 版为蓝本进行讲解,广泛吸收国内外优秀教材的成果进行内容设计,在系统介绍基本理论和基本方法的同时,注意介绍新的成熟内容,以及统计学在实际问题中的应用。

(3) 实用性。全书采用了基础知识介绍和实例操作相结合的方法,互相补充,本书的实例大多来源于经济生活之中,使读者在学完本书后能够快速地将知识应用于实践。

(4) 结构清晰,讲解详尽。全书采用基础知识—程序实现—综合实例分析的讲解方法,循序渐进地提高读者的 SAS 编程知识,而且每个知识点和实例都做了尽可能详细的讲解,使读者学习起来轻松自如。

(5) 配有全部的案例数据、程序与多媒体示范。本书的配套光盘中提供了所有实例的数据、SAS 程序、视频操作  动画演示,读者可以在观看录像中增强对知识点的理解。

本书共分为 21 章,依次介绍了 SAS 9.2 的基本编程知识、基本统计分析、高级统计分析、SAS 数据挖掘、SAS 在金融中的应用,以及各章节中的案例分析等内容。

第 1 章 数据挖掘概述。介绍数据挖掘的涵义,为什么要进行数据挖掘,数据挖掘的用途、过程等,最后介绍 SAS 软件系统在数据挖掘中的地位。

第 2 章 SAS 模块概述。详细介绍 SAS 软件系统和 SAS 各个常用的模块。

第 3 章 SAS 程序设计基础。详细介绍 SAS 编程中的各个步骤,包括 SAS 安装过程, SAS 数据步和过程步,程序调试及 SAS 函数,掌握 SAS 编程基础。

第 4 章 数据预处理。内容包括数据输入、数据整理、数据步变量控制、数据修改与选择,最后介绍 SAS 与 SPSS 软件之间的数据转换。

第 5 章 数据汇总与报表制作。介绍最基本的数据报表生成过程, PROC PRINT 和 PROC TABULATE 过程。

第 6 章 SAS 绘图。统计图是统计描述的重要工具,本章介绍 SAS 绘制图形的 GPLOT、GCHART、G3D 过程。

第 7 章 数据描述。介绍利用统计图、统计量和数据分布进行数据描述的过程。

第 8 章 描述性统计分析。叙述描述性统计分析中的平均数、中位数、众数等度量集中趋势,用极差、标准差、变异系数等度量。

第 9 章 ANALYST 模块。主要介绍利用 ANALYST 模块进行数据管理。

第 10 章 参数估计和假设检验。介绍对于定量资料的统计描述和简单推断,主要有 UNIVARIATE、MEANS、TTEST 过程。



第 11 章 方差分析与协方差分析。介绍方差分析与协方差分析的基本原理及其 SAS 过程。

第 12 章 回归分析。介绍线性回归、REG 过程、多项式回归、逐步回归 LOGISTIC 过程及非线性回归。

第 13 章 主成分分析与因子分析。介绍主成分分析与因子分析的原理、数学模型及 SAS 过程，然后应用于案例分析。

第 14 章 相关分析和对应分析。介绍相关分析和对应分析的原理、数学模型及案例分析。

第 15 章 判别分析。介绍判别分析的最大似然法、Fisher 判别分析法、Bayes 判别分析法、逐步判别分析法等，并应用于案例分析。

第 16 章 聚类分析。介绍聚类分析的基本原理，以及 CLUSTER、FASTCLUS、VARCLUS、TREE 过程。

第 17 章 生存分析。介绍 SAS 对生存分析的应用，并进行研究探讨。

第 18 章 时间序列分析。介绍时间序列的数学模型，SAS 的 ARIMA 过程及案例分析。

第 19 章 SAS 数据挖掘应用。主要介绍 SAS 数据挖掘方法论——SEMMA，SAS 的企业数据挖掘套件 SAS/EM 及其案例分析。

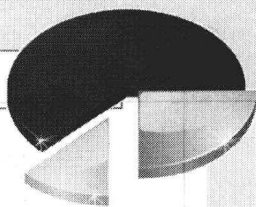
第 20 章 SAS 在数据预测中的应用。介绍利用 SAS 软件系统对有关数据预测的一些案例进行分析研究，为读者展示数据预测的相关技术和方法。

第 21 章 SAS 在金融数据分析中的应用。介绍利用 SAS 软件对部分金融学中数学模型进行分析求解，展示 SAS 软件在金融数据挖掘中的应用。

本书主要由尚涛完成，香港中文大学谢龙汉博士提供技术支持和指导，并对本书进行了校对和完善。参加本书编写和光盘开发的还有林伟、魏艳光、林木议、王悦阳、林伟洁、林树财、郑晓、吴苗、李翔、莫衍、朱小远、唐培培、耿煜、邓奕、张桂东、鲁力、刘文超、刘新东等，同时也非常感谢腾龙工作室其他成员的帮助和支持。

由于时间仓促，书中难免有疏漏之处，请读者谅解。读者可通过电子邮件 reader.toptech@gmail.com 或者 reader.toptech@163.com 与我们交流。

编著者



## 目 录

<b>第 1 章 数据挖掘概述</b> ..... 1	2.4.4 SAS/INSIGHT 模块.....34
1.1 数据挖掘简介 .....1	2.4.5 SAS/EM 模块.....36
1.1.1 数据挖掘的含义.....1	<b>第 3 章 SAS 程序设计基础</b> ..... 38
1.1.2 数据挖掘的起源.....2	3.1 SAS 编程基础 .....38
1.1.3 统计学与数据挖掘.....2	3.1.1 SAS 语言基础.....39
1.1.4 数据挖掘相关的一些问题.....5	3.1.2 SAS 语言构成.....43
1.2 数据挖掘用途 .....10	3.1.3 SAS 结构化编程语句.....46
1.3 数据挖掘过程 .....11	3.1.4 SAS 程序编写规则.....48
1.3.1 数据挖掘用户.....11	3.2 SAS 程序的数据步 .....49
1.3.2 数据挖掘工具.....14	3.2.1 DATA 语句.....49
1.3.3 数据挖掘步骤.....14	3.2.2 INPUT 语句.....50
1.4 SAS——数据挖掘领域的领导者.....15	3.2.3 CARDS 与 CARDS4 语句.....50
1.5 SAS 在各种商业解决方案中的应用.....16	3.2.4 INFILE 语句.....51
1.5.1 SAS 数据挖掘技术的实现.....17	3.2.5 SET 语句.....52
1.5.2 SAS 在商业领域中的应用.....18	3.2.6 MERGE 语句.....53
<b>第 2 章 SAS 模块概述</b> ..... 20	3.3 SAS 数据步循环与转移控制 .....54
2.1 SAS 简介.....20	3.3.1 IF 语句.....54
2.1.1 SAS 的设计思想.....21	3.3.2 SELECT 语句.....55
2.1.2 SAS 的功能.....21	3.3.3 DO 语句.....56
2.1.3 SAS 的特点.....22	3.3.4 GO TO 语句.....58
2.2 SAS 软件安装、启动与退出 .....22	3.3.5 RETURN 语句.....59
2.2.1 SAS 软件的安装.....22	3.3.6 CONTINUE 语句与 LEAVE
2.2.2 SAS 软件的启动.....22	语句.....59
2.2.3 SAS 软件的退出.....23	3.3.7 如何跳出选择结构和循环体 .....59
2.3 SAS 界面.....24	3.4 SAS 程序的过程步 .....60
2.3.1 Explorer 窗口.....25	3.4.1 SAS 过程步用法.....60
2.3.2 Editor 窗口.....25	3.4.2 VAR 与 MODLE 语句.....60
2.3.3 Results 窗口.....26	3.4.3 ID 与 WHERE 语句.....61
2.3.4 Log 窗口.....27	3.4.4 BY 与 CLASS 语句.....61
2.3.5 Output 窗口.....27	3.4.5 OUTPUT 语句.....62
2.4 SAS 模块介绍.....28	3.4.6 FERQ 与 WEIGHT 语句.....62
2.4.1 SAS/BASE 模块.....30	3.4.7 LABEL 与 FORMAT 语句.....62
2.4.2 SAS/ANALYSIS 模块.....31	3.5 SAS 函数.....63
2.4.3 SAS/ASSIST 模块.....32	3.5.1 数学函数.....63
	3.5.2 数组函数.....64

3.5.3	日期时间函数	64	5.1.2	使用中文列标题	106
3.5.4	概率分布函数	65	实例 5-2	修改标题实例	107
3.5.5	分位数函数	66	5.1.3	标题和脚注	107
3.5.6	样本统计函数	66	实例 5-3	修改标题实例	107
3.5.7	随机函数	67	5.1.4	用 BY 语句分组处理	108
<b>第 4 章</b>	<b>数据预处理</b>	<b>69</b>	5.2	使用过程 TABULATE 制作汇 总报表	109
4.1	数据输入	69	实例 5-4	汇总报表实例	110
4.1.1	原始数据的读取	70	实例 5-5	绘制统计量表	112
4.1.2	数据导入	71	<b>第 6 章</b>	<b>SAS 绘图</b>	<b>114</b>
4.2	数据整理	73	6.1	GPLOT 过程	114
4.2.1	数据集选项	73	实例 6-1	GPLOT 过程绘制图形 编程操作	115
4.2.2	整理数据集	74	6.2	GCHART 过程	115
4.2.3	缺失值处理	84	实例 6-2	GCHART 过程绘制 条形图	116
4.2.4	UPDATE 语句更新数据集	86	实例 6-3	GCHART 过程绘制 GDP 数据的 BLOCK 图形	117
4.2.5	数据清洗	87	6.3	G3D 过程	118
4.3	数据步变量控制	92	实例 6-4	绘制二维正态分布曲面 图形	118
4.3.1	ARRAY 语句	92	实例 6-5	绘制 $z = \sin \sqrt{x^2 + y^2}$ 函数的三维图形	120
4.3.2	INFORMAT 语句与 FORMAT 语句	93	<b>第 7 章</b>	<b>数据描述</b>	<b>123</b>
4.3.3	LABEL 语句	94	7.1	统计图	124
4.3.4	ATTRIB 语句	96	7.1.1	直方图	124
4.3.5	DROP 语句与 KEEP 语句	97	实例 7-1	GCHART 过程绘制 直方图	124
4.3.6	RENAME 语句与 RETAIN 语句	97	7.1.2	条形图	126
4.4	数据修改与选择	98	实例 7-2	GCHART 过程绘制 条形图	126
4.4.1	赋值语句	98	7.1.3	散点图	127
4.4.2	累加语句	98	实例 7-3	GPLOT 过程绘制散点图	128
4.4.3	DELETE 语句与 LOSTCARD 语句	99	7.1.4	饼图	129
4.4.4	STOP 语句与 ABORT 语句	100	实例 7-4	GCHART 过程绘制饼图	129
4.4.5	WHERE 语句	101	7.1.5	盒形图	130
4.4.6	REMOVE 语句与 REPLACE 语句	101	实例 7-5	BOXPLOT 过程绘制 盒形图	131
4.4.7	MISSING 语句	102	7.1.6	茎叶图	132
<b>第 5 章</b>	<b>数据汇总与报表制作</b>	<b>103</b>			
5.1	使用过程 PRINT 制作报表	103			
5.1.1	基本用法	104			
实例 5-1	PROC PRINT 操作 实例	104			

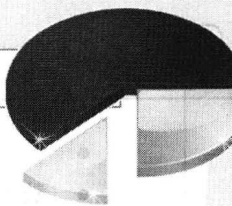
实例 7-6 UNIVARIATE 过程绘制	实例 8-8 SUMMARY 语句实例.....165
茎叶图 .....132	8.2.2 统计值的图形表示: PROC
7.1.7 时间序列图 .....133	CHART .....166
实例 7-7 TIMEPLOT 过程绘制	实例 8-9 绘制数据分布图形 .....168
时间序列图 .....133	实例 8-10 利用 CHART 过程的
7.2 统计量 .....135	VBAR 及 HBAR 命令
7.2.1 集中趋势 .....135	绘制条形图 .....169
实例 7-8 利用 MEAN 函数求	8.2.3 一般制图: PROC PLOT .....171
平均数 .....136	实例 8-11 PLOT 过程绘制图形 .....172
7.2.2 离散程度 .....137	<b>第 9 章 ANALYST 模块</b> .....173
实例 7-9 利用函数 VAR 和 STD	9.1 ANALYST 模块概述 .....173
求方差和标准差 .....139	9.1.1 ANALYST 模块简介 .....173
7.2.3 分布状态 .....141	9.1.2 ANALYST 菜单介绍 .....177
实例 7-10 利用 SKEWNESS 和	9.2 数据集的窗口操作 .....177
KURTOSIS 函数求偏度	9.2.1 数据集输入 .....177
和峰度 .....142	9.2.2 数据表修改 .....178
7.3 数据分布 .....143	9.2.3 数据保存 .....180
实例 7-11 SAS 中的部分概率分布	9.3 绘制统计图 .....180
函数的应用 .....144	9.3.1 条形图 .....180
<b>第 8 章 描述性统计分析</b> .....146	9.3.2 饼图 .....181
8.1 SAS 编程进行统计分析 .....146	9.3.3 散点图 .....183
8.1.1 基本概念 .....147	9.4 统计分析 .....184
8.1.2 FREQ 过程 .....149	<b>第 10 章 参数估计与假设检验</b> .....187
实例 8-1 频数表的生成实例 .....151	10.1 参数估计和假设检验概述 .....187
实例 8-2 绘制实验数据表格 .....153	10.1.1 参数估计 .....187
8.1.3 MEANS 过程 .....154	10.1.2 假设检验 .....189
实例 8-3 求平均增长率 .....156	10.2 假设检验的 SAS 过程 .....190
实例 8-4 利用 MEANS 过程求各种	10.2.1 UNIVARIATE 过程 .....190
统计量 .....156	10.2.2 MEANS 过程 .....191
8.1.4 UNIVARIATE 过程 .....159	10.2.3 TTEST 过程 .....192
实例 8-5 利用 UNIVARIATE 过程	10.3 不同类型的均值和方差的检验 .....192
求各种统计量 .....160	10.3.1 单变量均值 t 检验 .....192
实例 8-6 求样本的极差、上四分位	实例 10-1 TTEST 过程的实例数据
数和下四分位数 .....161	分析 .....193
8.1.5 TABULATE 过程 .....162	实例 10-2 总体均值检验 .....194
实例 8-7 制作数据表格 .....162	10.3.2 样本均数与总体均数差异的
8.2 其他描述性统计过程 .....165	t 检验 .....194
8.2.1 产生描述性统计值的输出	实例 10-3 均值的显著性差别
文件: PROC SUMMARY .....165	检验 .....195



- 10.3.3 配对资料的 t 检验 ..... 195
- 实例 10-4 乳酸饮料实验数据的  
配对 t 检验 ..... 195
- 实例 10-5 均值有无差异的检验 ..... 197
- 10.3.4 两样本均数比较的 t 检验 ..... 198
- 实例 10-6 均数差别的显著性  
检验 ..... 198
- 实例 10-7 数据比例的显著性  
检验 ..... 198
- 10.4 正态性检验 ..... 200
- 实例 10-8 样本数据的正态性检验  
实例 1 ..... 200
- 实例 10-9 样本数据的正态性检验  
实例 2 ..... 201
- 第 11 章 方差分析与协方差分析** ..... 204
- 11.1 方差分析的基本原理 ..... 204
- 11.1.1 自由度与平方和分解 ..... 206
- 11.1.2 F 检验 ..... 207
- 11.2 单因素方差分析 ..... 208
- 11.2.1 单因素方差分析步骤 ..... 208
- 11.2.2 判断与结论 ..... 210
- 11.2.3 ANOVA 过程 ..... 210
- 实例 11-1 分析饲料营养效果是否  
有明显差异 ..... 211
- 实例 11-2 分析不同实验室试制的  
纸张光滑度有无差异 ..... 212
- 实例 11-3 研究 6 种棉花种子包衣剂  
对棉花生长的影响 ..... 214
- 11.3 双因素方差分析 ..... 216
- 11.3.1 只考虑主效应的多因素  
方差分析 ..... 217
- 11.3.2 存在交互效应的多因素  
方差分析 ..... 219
- 实例 11-4 某药物对某癌细胞株增殖  
影响的研究 ..... 221
- 11.4 协方差分析 ..... 222
- 实例 11-5 分析三种饲料的营养价值  
之间有无显著性差别 ..... 225
- 第 12 章 回归分析** ..... 230
- 12.1 线性回归 ..... 230
- 12.1.1 线性回归模型 ..... 231
- 12.1.2 回归方程的显著性检验 ..... 231
- 12.1.3 预测问题 ..... 233
- 12.2 REG 过程 ..... 234
- 实例 12-1 分析我国内地可支配  
收入和消费性支出之间  
的关系 ..... 237
- 实例 12-2 利用多元线性回归分析  
学生肺活量及有关变量  
的关系 ..... 240
- 12.3 多项式回归 ..... 243
- 12.3.1 曲线回归的基本原理 ..... 243
- 12.3.2 RSREG 过程 ..... 243
- 实例 12-3 确定最佳经济用肥量的  
多项式回归模型 ..... 244
- 12.4 逐步回归 ..... 246
- 实例 12-4 人体血糖、胰岛素及生  
长素的多元线性回归  
关系 ..... 246
- 12.5 LOGISTIC 回归 ..... 248
- 12.5.1 逻辑回归模型概述 ..... 249
- 12.5.2 LOGISTIC 过程 ..... 250
- 实例 12-5 对照研究单因素两暴露  
水平及多暴露水平资料  
的统计分析 ..... 251
- 12.6 非线性回归 ..... 255
- 12.6.1 非线性回归分析的基本  
原理 ..... 255
- 12.6.2 NLIN 过程 ..... 256
- 实例 12-6 酵母种群增长的拟合  
生长模型 ..... 257
- 实例 12-7 最佳生长模型的 LOGISTIC  
拟合 ..... 259
- 第 13 章 主成分分析与因子分析** ..... 262
- 13.1 主成分分析 ..... 262
- 13.1.1 主成分分析的数学原理 ..... 263
- 13.1.2 用 PRINCOMP 过程进行  
主成分分析 ..... 264
- 实例 13-1 我国 2006 年经济发展  
情况的主成分分析 ..... 265

13.2 因子分析 .....	270	15.3 综合实例 .....	321
13.2.1 因子分析的基本原理 .....	271	实例 15-1 国内各省市农民家庭	
13.2.2 因子分析的基本步骤和		收支情况的研究 .....	321
过程 .....	273	实例 15-2 基于判别分析法的上市	
13.2.3 利用 FACTOR 过程进行		公司财务分析研究 .....	328
因子分析 .....	274	<b>第 16 章 聚类分析</b> .....	337
实例 13-2 中国房地产经济区的		16.1 聚类分析的基本原理 .....	337
研究分析 .....	276	16.1.1 聚类的数学原理 .....	338
13.3 主成分分析和因子分析的区别 .....	282	16.1.2 SAS 中的聚类过程 .....	344
<b>第 14 章 相关分析和对应分析</b> .....	284	16.2 聚类分析的步骤和过程 .....	345
14.1 相关分析 .....	284	16.2.1 CLUSTER 过程 (系统聚类	
14.1.1 相关关系 .....	285	过程) .....	345
14.1.2 相关图形和相关系数 .....	286	实例 16-1 中国城镇居民消费结构的	
14.1.3 简单相关分析的 CORR		聚类分析 .....	346
过程 .....	287	16.2.2 FASTCLUS 过程 (快速聚类	
实例 14-1 简单相关系数的计算 .....	288	过程) .....	351
14.2 典型相关分析 .....	290	实例 16-2 聚类分析在客户定位中	
14.2.1 典型相关分析的基本原理 .....	290	的应用研究 .....	352
14.2.2 典型相关分析的 CANCORR		16.2.3 VARCLUS 过程 (变量聚类	
过程 .....	291	过程) .....	355
实例 14-2 城市竞争力与基础设施的		实例 16-3 变量聚类在多指标系统	
典型相关分析 .....	292	评价中的应用 .....	357
实例 14-3 城镇居民收入和支出的		16.2.4 TREE 过程 (画树状图	
典型相关分析 .....	298	过程) .....	360
14.3 对应分析 .....	305	实例 16-4 对全球各国信息设施的	
14.3.1 对应分析的基本原理 .....	306	发展情况进行聚类分析	
14.3.2 对应分析的 CORRESP		研究 .....	362
过程 .....	307	<b>第 17 章 生存分析</b> .....	365
实例 14-4 对应分析在市场细分中		17.1 生存分析基本概述 .....	365
的应用 .....	308	17.1.1 生存分析的基本概念 .....	365
<b>第 15 章 判别分析</b> .....	313	17.1.2 生存资料的特点 .....	367
15.1 判别分析的基本原理 .....	313	17.1.3 生存分析方法 .....	368
15.1.1 判别分析的含义 .....	314	17.2 生存分析的 LIFETEST 过程 .....	369
15.1.2 判别分析的数学模型与判别		实例 17-1 生存分析在医学课题研	
方法 .....	315	究中的应用 .....	370
15.2 判别分析的 SAS 过程 .....	317	17.3 COX 模型回归分析 .....	373
15.2.1 DISCRIM 过程 .....	317	17.3.1 COX 回归模型 .....	373
15.2.2 CANDISC 过程 .....	319	17.3.2 PHREG 过程 .....	375
15.2.3 STEPDISC 过程 .....	319	实例 17-2 COX 模型的分析应用 .....	376

<b>第 18 章 时间序列分析</b> .....	380	实例 19-2 购物篮问题分析	423
18.1 时间序列概述 .....	380	<b>第 20 章 SAS 在数据预测中的应用</b> .....	427
18.1.1 时间序列的组成部分 .....	381	20.1 数据预测简介 .....	427
18.1.2 时间序列的数学模型 .....	381	20.1.1 数据预测 .....	427
18.1.3 时间序列的因素分析 .....	382	20.1.2 SAS 中的预测分析模块 .....	430
18.1.4 随机时间序列分析 .....	386	20.2 数据预测案例分析 .....	430
18.1.5 时间序列的分析步骤 .....	388	实例 20-1 国民生产总值的预测 .....	430
18.2 SAS 的 ARIMA 过程 .....	388	实例 20-2 SAS/Time Series	
18.3 综合实例 .....	389	Forecasting System	
实例 18-1 化工生产数据的时间		模块应用 .....	435
序列分析 .....	389	<b>第 21 章 SAS 在金融数据分析中的</b>	
实例 18-2 国内金融及保险业每人		<b>应用</b> .....	439
每月平均薪资趋势		21.1 现金流贴现分析 .....	439
分析 .....	394	实例 21-1 现金流贴现的计算 .....	440
实例 18-3 运用 ARIMA 过程对上证		实例 21-2 企业现金流的贴现	
指数日线数据进行拟合		计算 .....	441
分析 .....	406	实例 21-3 利用金融函数 compound	
<b>第 19 章 SAS 数据挖掘应用</b> .....	410	计算复利率 .....	442
19.1 SAS 数据挖掘 .....	410	21.2 股票分类 .....	442
19.2 SAS 数据挖掘方法论——		实例 21-4 利用 CLUSTER 过程对	
SEMMA .....	414	股票进行聚类分析 .....	443
19.2.1 数据取样 .....	414	21.3 资本资产定价模型 (CAPM	
19.2.2 数据探索 .....	414	模型) .....	448
19.2.3 问题明确化、数据调整和		实例 21-5 CAPM 模型实例研究 .....	449
技术选择 .....	415	21.4 B-S 模型期权定价 .....	454
19.2.4 模型研发 .....	416	实例 21-6 B-S 期权定价的 SAS	
19.2.5 模型评估 .....	416	程序实现 .....	457
19.3 数据挖掘套件 SAS/EM .....	417		
实例 19-1 SAS/EM 聚类分析 .....	418		



# 第1章 数据挖掘概述

每个企业每天都会产生大量的数据，这些数据来自于不同的数据源。数据挖掘是对大量的原始数据进行选择、分析和建模，从中发现以前没有发现的趋势和模式，通过数据和文本挖掘得到的信息对企业战略决策有很大帮助。因此，数据挖掘使得人们从原始数据到更加明智的业务决策。本章将详细叙述数据挖掘的过程、工具、用法，以及 SAS 软件系统在数据挖掘中的位置等，为利用 SAS 软件系统进行数据挖掘打下基础。



## 本章内容

- 数据挖掘简介
- 数据挖掘用途
- 数据挖掘过程
- SAS——数据挖掘领域的领导者
- SAS 在各种商业解决方案中应用

## 1.1 数据挖掘简介

数据挖掘，简言之就是从大量繁杂的数据中获取隐含其中的信息，如对顾客分类、聚类、欺诈甄别、潜在顾客识别等，现在应用领域很广，如设计、零售、金融、银行、医疗、政府决策、企业财务、商业决策。下面将从数据挖掘的概念、起源等说起。

### 1.1.1 数据挖掘的含义

数据挖掘 (Data Mining) 就是从大量数据中获取有效的、新颖的、潜在有用的、最终可理解的模式的非平凡过程。美国 SAS 软件研究所将数据挖掘定义为“数据挖掘是按照既定的业务目标，对大量的企业数据进行探索、揭示隐藏其中的规律性并进一步模型化的先进、有效的方法。”数据挖掘的广义观点：数据挖掘就是从存放在数据库，数据仓库或其他信息库中的大量的数据中“挖掘”有趣知识的过程。数据挖掘又称数据库中知识发现 (Knowledge Discovery in Database, KDD)，也有人把数据挖掘视为数据库中知识发现过程的一个基本步骤。知识发现过程由以下步骤组成：



- 数据清理;
- 数据集成;
- 数据选择;
- 数据变换;
- 数据挖掘;
- 模式评估;
- 知识表示。

数据挖掘可以与用户或知识库交互。

并非所有的信息发现任务都被视为数据挖掘。例如，使用数据库管理系统查找个别的记录，或通过互联网的搜索引擎查找特定的 Web 页面，则是信息检索（Information Retrieval）领域的任务。虽然这些任务是重要的，可能涉及使用复杂的算法和数据结构，但是它们主要依赖传统的计算机科学技术和数据的明显特征来创建索引结构，从而有效地组织和检索信息。尽管如此，数据挖掘技术也已用来增强信息检索系统的能力。

## 1.1.2 数据挖掘的起源

需要是发明之母。近年来，数据挖掘引起了信息产业界的极大关注，其主要原因是存在大量数据，可以广泛使用，并且迫切需要将这此数据转换成有用的信息和知识。获取的信息和知识可以广泛用于各种应用，包括商务管理、生产控制、市场分析、工程设计和科学探索等。

在国内，数据挖掘应用最多的行业也是数据积累最快的行业，如电信、金融、保险、互联网、零售等行业。

数据挖掘利用了来自如下一些领域的思想：

- 来自统计学的抽样、估计和假设检验；
- 人工智能、模式识别和机器学习的搜索算法、建模技术和学习理论。

数据挖掘也迅速地接纳了来自其他领域的思想，这些领域包括最优化、进化计算、信息论、信号处理、可视化和信息检索，一些其他领域也起到重要的支撑作用。特别的是，需要数据库系统提供有效的存储、索引和查询处理支持。源于高性能（并行）计算的技术在处理海量数据集方面常常是重要的。分布式技术也能帮助处理海量数据，并且当数据不能集中到一起处理时更是至关重要。

## 1.1.3 统计学与数据挖掘

统计学和数据挖掘有着共同的目标：发现数据中的结构。事实上，由于它们的目标相似，有些人（尤其是统计学家）认为数据挖掘是统计学的分支。这是一个不切合实际的想法。因为数据挖掘还应用了其他领域的思想、工具和方法，尤其是计算机学科，例如，数据库技术和机器学习，而且它所关注的某些领域和统计学家所关注的有很大不同。

### 1.1.3.1 统计学的性质

试图为统计学下一个太宽泛的定义是没有意义的。尽管可能做到，但会引来很多异

议。相反，更要关注统计学不同于数据挖掘的特性。

差异之一同 1.1.2 节中最后一段提到的相关，即统计学是一门比较保守的学科，目前有一种趋势是越来越精确。当然，这本身并不是坏事，只有越精确，才能避免错误，发现真理。但是，如果过度的话，则是有害的。这个保守的观点源于统计学是数学的分支这样一个看法，尽管统计学确实以数学为基础（正如物理和工程也以数学为基础，但没有被认为是数学的分支），但它同其他学科还有紧密的联系。

数学背景和追求的确加强了这样一个趋势：在采用一个方法之前先要证明，而不是像计算机科学和机器学习那样注重经验。这就意味着有时候和统计学家关注同一问题的其他领域的研究者提出一个很明显有用的方法，但它却不能被证明（或还不能被证明）。统计杂志倾向于发表经过数学证明的方法而不是一些特殊方法。数据挖掘作为几门学科的综合，已经从机器学习那里继承了实验的态度。这并不意味着数据挖掘工作者不注重精确，而只是说明如果方法不能产生结果的话就会被放弃。

正是统计文献显示了统计的数学精确性。同时还显示了其对推理的侧重。尽管统计学的一些分支也侧重于描述，但是浏览一下统计论文的话，就会发现这些文献的核心问题就是在观察了样本的情况下如何去推断总体。当然这常常也是数据挖掘所关注的。下面会提到数据挖掘的一个特定属性，就是要处理的是一个大数据集。这就意味着，由于可行性的原因，常常得到的只是一个样本，但是需要描述样本取自的那个大数据集。然而，数据挖掘问题常常可以得到数据总体，例如，关于一个公司的所有职工数据，数据库中的所有客户资料，去年的所有业务。在这种情形下，推断就没有价值了（如年度业务的平均值），因为观测到的值也就是估计参数。这就意味着，建立的统计模型可能会利用一系列概率表述（如一些参数接近于 0，则会从模型中剔除掉），但当总体数据可以获得的话，在数据挖掘中则变得毫无意义。在这里，可以很方便地应用评估函数：针对数据的足够的表述。事实是，常常所关注的是，模型是否合适而不是它的可行性，在很多情形下，使得模型的发现很容易。例如，在寻找规则时常常会利用吻合度的单纯特性（如应用分支定理），但当应用概率陈述时则不会得到这些特性。

统计学和数据挖掘部分交迭的第三个特性是在现代统计学中起核心作用的“模型”。“模型”这个术语更多的含义是变化的。一方面，统计学模型是基于分析变量间的联系，但另一方面这些模型关于数据的总体描述确实没有道理。关于信用卡业务的回归模型可能会把收入作为一个独立的变量，因为一般认为高收入会导致大的业务。这可能是一个理论模型（尽管基于一个不牢靠的理论）。与此相反，只需在一些可能具有解释意义的变量基础上进行逐步的搜索，从而获得一个有很大预测价值的模型，尽管不能做出合理的解释（通过数据挖掘去发现一个模型的时候，常常关注的就是后者）。

还有其他方法可以区分统计模型，但在这里将不做探讨。这里关注的是，现代统计学是以模型为主的。而计算、模型选择条件是次要的，只是如何建立一个好的模型。但在数据挖掘中，却不完全如此。在数据挖掘中，准则起了核心的作用。在很多情形下，模型的选择并不都是显而易见的，选择一个合适的模型是不可能的，最合适的计算方法也是不可行的。在这种情形下，从另外一个角度出发，应用设计的一系列技术来回答 MVA 问题，暂不考虑模型和最优判别的选择。

相对于统计学而言，准则在数据挖掘中起着更为核心的作用并不奇怪，数据挖掘所继承的学科如计算机科学及相关学科。数据集的规模常常意味着传统的统计学准则不适合数据挖掘问题，不得不重新设计。尽管一些统计学的准则已经得到发展，但更多的应用是机器学习。

### 1.1.3.2 数据挖掘的性质

由于统计学基础的建立在计算机的发明和发展之前，所以，常用的统计学工具包含很多可以手工实现的方法。因此，对于很多统计学家来说，1000 个数据就已经是很大的了。但这个“大”对于英国大的信用卡公司每年 3.5 亿笔业务或 AT&T 每天 2 亿个长途呼叫来说相差太远了。很明显，面对这么多的数据，则需要设计不同于那些“原则上可以用手工实现”的方法。这意味着计算机（正是计算机使得大数据可能实现）对于数据的分析和处理是关键，分析者直接处理数据将变得不可行。相反，计算机在分析者和数据之间起到了必要的过滤的作用，这也是数据挖掘特别注重准则的另一原因。尽管有必要，把分析者和数据分离开很明显导致了一些关联任务。这里就有一个真正的危险：非预期的模式可能会误导分析者，这一点下面会讨论。

在现代统计中计算机是一个重要的工具，它们确实是，并不是因为数据的规模。对数据的精确分析方法（如 Bootstrap 方法、随机测试、迭代估计方法）及比较适合的复杂的模型，正是有了计算机才是可能的。计算机已经使得传统统计模型的视野大大的扩展了，还促进了新工具的飞速发展。

统计学很少会关注实时分析，然而数据挖掘问题常常需要这些。例如，银行事务每天都会发生，没有人能等 3 个月得到一个可能的欺诈的分析。类似的问题发生在总体随时间变化的情形。我的研究组有明确的例子，显示银行债务的申请随时间、竞争环境、经济波动而变化。

### 1.1.3.3 讨论

主要有几个问题：

① 数据挖掘有时候是一次性的实验，这是一个误解。它更应该被看作是一个不断的过程（尽管数据集是确定的）。从一个角度检查数据可以解释结果，以相关的观点检查可能会更接近等。关键是，除了极少的情形下，很少知道哪一类模式是有意义的。数据挖掘的本质是发现非预期的模式，同样非预期的模式要以非预期的方法来发现。

② 与把数据挖掘作为一个过程的观点相关联的是认识到结果的新颖性，许多数据挖掘的结果是所期望的。然而，可以解释这个事实并不能否定挖掘出它们的价值。没有这些实验，可能根本不会想到这些。实际上，只有那些可以依据过去经验形成的合理的解释的结构才会有价值的。

③ 在大数据集中发现模式的可能性当然存在，然而，也不应就此掩盖危险：所有真实的数据集（即使那些是以完全自动方式搜集的数据）都有产生错误的可能。关于人的数据集（如事务和行为数据）尤其有这种可能。

## 1.1.4 数据挖掘相关的一些问题

### 1.1.4.1 数据挖掘和统计分析

非要去区分数据挖掘和统计的差异其实是没有太大意义的。一般将其定义为数据挖掘技术的 CART (Classification and Regression Trees)、CHAID (Chi-Square Automatic Interaction Detector) 或模糊计算等理论方法,也都是由统计学者根据统计理论所发展衍生。换另一个角度看,数据挖掘有相当大的比重是由高等统计学中的多变量分析所支撑。但是为什么数据挖掘的出现会引发各领域的广泛关注呢?主要原因在于相对传统统计分析而言,数据挖掘有下列几项特性:

① 处理大量实际数据更强势,且无须专业的统计背景去使用数据挖掘的工具;

② 数据分析趋势为从大型数据库获取所需数据并使用专属计算机分析软件,数据挖掘的工具更符合企业需求;

③ 就纯理论的基础点来看,数据挖掘和统计分析有应用上的差别,毕竟数据挖掘目的是方便企业终端用户使用而非给统计学家检测用的。

### 1.1.4.2 数据挖掘和数据仓库

若将数据仓库 (Data Warehousing) 比作“矿坑”,数据挖掘就是深入矿坑采矿的工作。毕竟数据挖掘不是一种无中生有的魔术,也不是点石成金的炼金术。若没有足够丰富完整的数据,是很难期待数据挖掘能挖掘出什么有意义的信息的。

要将庞大的数据转换成为有用的信息,必须先有效率地收集信息。随着科技的进步,功能完善的数据库系统就成了最好的收集数据的工具。数据仓库,简单地说,就是收集来自其他系统的有用数据,存放在一个整合的储存区内。所以,其实就是一个经过处理整合,且容量特别大的关系型数据库,用于储存决策支持系统 (Design Support System, DSS) 所需的数据,供决策支持或数据分析使用。从信息技术的角度来看,数据仓库的目标是在组织中,在正确的时间,将正确的数据交给正确的人。

许多人对于数据仓库和数据挖掘时常混淆,不知如何分辨。其实,数据仓库是数据库技术的一个新主题,利用计算机系统帮助操作、计算和思考,让作业方式改变,决策方式也跟着改变。

数据仓库本身是一个非常大的数据库,它储存着由组织作业数据库中整合而来的数据,特别是指 OLTP (On-Line Transactional Processing, 事务处理系统) 所得来的数据。将这些整合过的数据置放于数据仓库中,而公司的决策者则利用这些数据作决策;但是,这个转换及整合数据的过程,是建立一个数据仓库最大的挑战。因为将作业中的数据转换成有用的策略性信息是整个数据仓库的重点。

综上所述,数据仓库应该具有这些数据:整合性数据 (Integrated Data)、详细和汇总性的数据 (Detailed and Summarized Data)、历史数据、解释数据的数据。从数据仓库挖掘出对决策有用的信息与知识,是建立数据仓库与使用数据挖掘的最大目的,两者的本质与过程是两回事。换句话说,数据仓库应先行建立完成,数据发掘才能有效率的进行,因为数据仓库本身所含数据是干净 (不会有错误的数据参杂其中)、完备,且经过整合的。因此,两者之间的关系或许可解读为数据挖掘是从巨大数据仓库中找出有用信息的一种过程与技术。



### 1.1.4.3 数据挖掘和 OLAP

OLAP (Online Analytical Process) 意指由数据库所连接出来的在线分析处理程序。有些人会说：已经有 OLAP 的工具了，所以，不需要数据挖掘。事实上两者间是截然不同的，主要差异在于数据挖掘用在产生假设，OLAP 则用于查证假设。简单来说，OLAP 是由使用者所主导，使用者先有一些假设，然后利用 OLAP 来查证假设是否成立；而数据挖掘则是用来帮助使用者产生假设。所以，在使用 OLAP 或其他 Query 的工具时，使用者是自己在做探索 (Exploration)，但数据挖掘是用工具在帮助做探索。

例如，市场分析师在为超市规划货品架柜摆设时，可能会先假设婴儿尿布和婴儿奶粉会是常被一起购买的产品，接着便可利用 OLAP 的工具去验证此假设是否为真，又成立的证据有多明显；但数据挖掘则不然，执行数据挖掘的人将庞大的结账数据整理后，并不需要假设或期待可能的结果。透过数据挖掘技术可找出存在于数据中的潜在规则，于是可能得到（如尿布和啤酒常被同时购买的）意料之外的发现，这是 OLAP 所做不到的。

数据挖掘常能挖掘出超越归纳范围的关系，但 OLAP 仅能利用人工查询及可视化的报表来确认某些关系，是以数据挖掘此种自动找出甚至不会被怀疑过的数据模型与关系的特性。事实上已超越了经验、教育、想象力的限制，OLAP 可以和数据挖掘互补，但这项特性是数据挖掘无法被 OLAP 取代的。

### 1.1.4.4 数据挖掘步骤

以下提供一个数据挖掘的进行步骤作为参考：

- 理解业务与理解数据；
- 获取相关技术与知识；
- 整合与查询数据；
- 去除错误或不一致及不完整的数据；
- 由数据选取样本先行试验；
- 建立数据模型；
- 实际数据挖掘的分析工作；
- 测试与检验；
- 找出假设并提出解释；
- 持续应用于企业流程中。

由上述步骤看出，数据挖掘牵涉了大量的准备工作与规划过程，事实上许多专家皆认为整套数据挖掘的进行有 80% 的时间精力是花费在数据前置作业阶段，其中包含数据的净化与格式转换甚至表格的连接。由此可知，数据挖掘只是信息挖掘过程中的一个步骤而已，在进行此步骤前还有许多的工作要先完成。

### 1.1.4.5 数据挖掘理论

数据挖掘是近年来数据库应用技术中相当热门的议题，看似神奇、听来时髦，实际上却也不是什么新东西。因其所用如预测模型、数据分割、连结分析 (Link Analysis)、偏差侦测 (Deviation Detection) 等，美国早在第二次世界大战前就已应用，运用在人口普查及军事等方面。