

高等院校信息科学系列教材

数据分 析

范金城 梅长林 编著



科学出版社

高等院校信息科学系列教材

数 据 分 析

范金城 梅长林 编著

科学出版社

2002

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,梅长林编著。—北京:科学出版社,2002
(高等院校信息科学系列教材)
ISBN 7-03-010458-7

I . 数… II . ①范… ②梅… III . 统计分析(数学) IV . 0212.1

中国版本图书馆 CIP 数据核字(2002)第 035339 号

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 蕉 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:720×1000

2002年7月第一次印刷 印张:26 1/2

印数:1—5 000 字数:520 000

定 价:32.00 元

(如有印装质量问题,我社负责调换〈路通〉)

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

2002年3月

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, 书中尚存在一些不妥之处, 欢迎读者批评指正.

作者
于西安交通大学
2002 年 1 月

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、方差等数字特性	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	15
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
习题一	41
第2章 非参数方法	45
2.1 两种处理方法比较的秩检验.....	45
2.1.1 两种处理方法比较的随机化模型及秩的零分布	45
2.1.2 Wilcoxon 秩和检验	47
2.1.3 总体模型的 Wilcoxon 秩和检验	55
2.1.4 Smirnov 检验	56
2.2 成对分组设计下两种处理方法的比较.....	60
2.2.1 符号检验	61
2.2.2 Wilcoxon 符号秩检验	63
2.2.3 分组设计下两处理方法比较的总体模型	68
2.3 多种处理方法比较的 Kruskal-Wallis 检验	69
2.3.1 多种处理方法比较中秩的定义及 Kruskal-Wallis 统计量	69
2.3.2 Kruskal-Wallis 统计量的零分布	70
2.4 分组设计下多种处理方法的比较.....	73
2.4.1 分组设计下秩的定义及其零分布	73
2.4.2 Friedman 检验	74

2.4.3 改进的 Friedman 检验	77
2.5 列联表的独立性检验.....	80
2.5.1 定性变量与列联表	80
2.5.2 二维 $r \times s$ 列联表的独立性检验	83
2.5.3 三维 $r \times s \times t$ 列联表的独立性检验	84
习题二	90
第3章 回归分析	94
3.1 线性回归模型.....	94
3.1.1 线性回归模型及其矩阵表示	94
3.1.2 β 及 σ^2 的估计	95
3.1.3 有关的统计推断.....	96
3.2 残差分析	103
3.2.1 误差项的正态性检验	104
3.2.2 残差图分析	106
3.3 回归方程的选取与系统建模概述	109
3.3.1 穷举法	109
3.3.2 逐步回归法	118
3.3.3 系统建模过程概述	122
3.4 Logistic 回归模型	124
3.4.1 线性 Logistic 回归模型	124
3.4.2 参数的最大似然估计与 Newton-Raphson 迭代解法	126
3.4.3 Logistic 模型的统计推断	131
习题三	135
第4章 主成分分析.....	141
4.1 引言	141
4.2 总体主成分	142
4.2.1 总体主成分的定义	142
4.2.2 总体主成分的求法	143
4.2.3 总体主成分的性质	144
4.2.4 标准化变量的主成分	146
4.3 样本主成分	148
习题四	154
第5章 判别分析.....	159
5.1 距离判别	159

5.1.1 判别分析的基本思想及意义	159
5.1.2 两个总体的距离判别	160
5.1.3 判别准则的评价	163
5.1.4 多个总体的距离判别	168
5.2 Bayes 判别	171
5.2.1 Bayes 判别的基本思想	171
5.2.2 两个总体的 Bayes 判别	172
5.2.3 多个总体的 Bayes 判别	182
5.3 逐步判别	187
5.3.1 判别效果的检验	188
5.3.2 逐步判别的步骤	193
习题五	198
第6章 聚类分析	205
6.1 距离与相似系数	205
6.1.1 聚类分析的基本思想及意义	205
6.1.2 样品间的相似性度量——距离	206
6.1.3 变量间的相似性度量——相似系数	208
6.2 谱系聚类法	210
6.2.1 类间距离	211
6.2.2 类间距离的递推公式	212
6.2.3 谱系聚类法的步骤	214
6.2.4 谱系聚类法的统计量	221
6.2.5 变量聚类	226
6.3 快速聚类法	228
6.3.1 快速聚类法的步骤	228
6.3.2 用 L_m 距离进行快速聚类	237
习题六	241
第7章 时间序列分析	245
7.1 平稳时间序列	245
7.1.1 时间序列分析及其意义	245
7.1.2 随机过程概念及其数字特征	245
7.1.3 平稳时间序列与平稳随机过程	250
7.1.4 平稳定性检验及自协方差函数、自相关函数的估计	253
7.2 ARMA 时间序列及其特性	255

7.2.1 ARMA 时间序列的定义	255
7.2.2 ARMA 序列的平稳性与可逆性	258
7.2.3 ARMA 序列的相关特性	261
7.3 ARMA 时间序列的建模与预报	269
7.3.1 ARMA 序列参数的矩估计	270
7.3.2 ARMA 序列参数的精估计	272
7.3.3 ARMA 模型的定阶与考核	280
7.3.4 平稳线性最小均方预报	284
7.3.5 ARMA 序列的预报	288
7.4 ARIMA 序列与季节性序列	292
7.4.1 ARIMA 序列及其预报	292
7.4.2 季节性序列及其预报	299
习题七	305
第8章 Bayes 统计分析	310
8.1 Bayes 统计模型	310
8.1.1 Bayes 统计分析的基本思想及意义	310
8.1.2 Bayes 统计模型	310
8.1.3 Bayes 统计推断原则	315
8.1.4 先验分布的 Bayes 假设与不变先验分布	317
8.1.5 共轭先验分布	320
8.1.6 先验分布中超参数的确定	324
8.1.7 后验分布的计算	329
8.2 Bayes 统计推断	331
8.2.1 Bayes 参数点估计	331
8.2.2 Bayes 区间估计	338
8.2.3 Bayes 假设检验	343
习题八	349
第9章 常用数据分析方法的 SAS 过程简介	352
9.1 SAS 系统简介	352
9.1.1 数据的输入与输出	353
9.1.2 利用已有的 SAS 数据集建立新的 SAS 数据集	355
9.1.3 SAS 系统的数学运算符号及常用的 SAS 函数	357
9.1.4 逻辑语句与循环语句	360
9.2 常用数据分析方法的 SAS 过程	362

9.2.1 几种描述性统计分析的 SAS 过程	362
9.2.2 非参数方法的 SAS 过程	372
9.2.3 回归分析的 SAS 过程	376
9.2.4 主成分分析的 SAS 过程——PROC PRINCOMP 过程	383
9.2.5 判别分析的 SAS 过程	385
9.2.6 聚类分析的 SAS 过程	389
9.2.7 时间序列分析的 SAS 过程——PROC ARIMA 过程	394
9.2.8 SAS 系统的矩阵运算——PROC IML 过程简介	400
9.2.9 Bayes 统计分析计算实例	402
常用统计数值表	405
主要参考文献	411

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,它们是 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 USS 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).

峰度的计算公式是

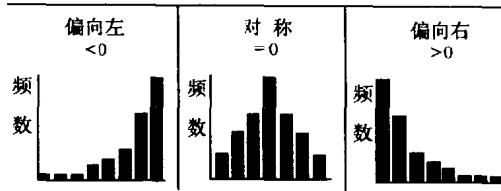


图 1.1

$$\begin{aligned} g_2 &= \frac{n(n+1)}{(n-1)(n-2)(n-3)s^4} \sum_{i=1}^n (x_i - \bar{x})^4 - 3 \frac{(n-1)^2}{(n-2)(n-3)} \\ &= \frac{n^2(n+1)u_4}{(n-1)(n-2)(n-3)s^4} - 3 \frac{(n-1)^2}{(n-2)(n-3)}. \end{aligned} \quad (1.8)$$

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:

$$\text{总体均值} \quad \mu = E(X), \quad (1.11)$$

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

$$\text{总体标准差} \quad \sigma = \sqrt{\text{Var}(X)}, \quad (1.13)$$

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

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

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

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

$$\text{总体峰度} \quad G_2 = \frac{\mu_4}{\sigma^4} - 3, \quad (1.18)$$

这里,我们对偏度与峰度作进一步的说明.

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.

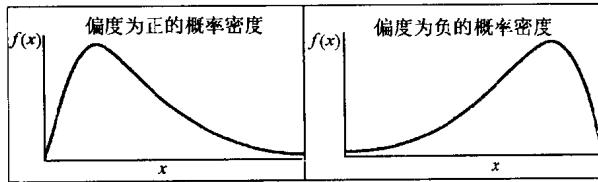


图 1.2

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).

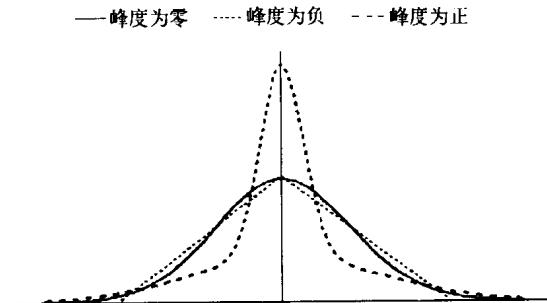


图 1.3

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 较大时, 有

$$\begin{aligned}\mu &\approx \bar{x}, \\ \sigma^2 &\approx s^2, \\ \sigma &\approx s, \\ \gamma &\approx CV, \\ \gamma_k &\approx v_k, \\ \mu_k &\approx u_k, \\ G_1 &\approx g_1, \\ G_2 &\approx g_2.\end{aligned}$$

这里, 特别要强调下列情况: 当观测数据 x_1, x_2, \dots, x_n 是所要研究对象的全体时,

数据的分布即总体分布. 我们认为取得每一个观测数据 x_i 是等可能性的, 即为 $\frac{1}{n}$;
总体分布是离散均匀分布:

$$P\{X=x_i\}=\frac{1}{n}, \quad i=1, 2, \dots, n.$$

对这种情况, 数据数字特征即总体数字特征. 许多实际数据属于这种情况, 它更能体现数据分析的特点——让数据本身说话. 实际上, 我们也可以把这种情况看作取自确定性模型的数据, 而上述数字特征仍有相应的统计意义.

计算数据的上述数字特征可以通过 SAS 系统 proc means 过程或 proc univariate 过程来实现.

例 1.1 从 19 个杆塔上的普通盘形绝缘子测得该层电导率(μs)的数据如下:

9.89	8.00	6.40	6.17	5.39	7.27	9.08
10.40	11.20	8.75	6.45	11.90	10.30	9.58
9.24	7.75	6.20	8.95	8.33		

计算均值、方差、标准差、变异系数、偏度、峰度.

解 通过计算, 得

$$\begin{aligned}\bar{x} &= 8.487, & s^2 &= 3.406, \\ s &= 1.845, & CV &= 21.745, \\ g_1 &= 0.035, & g_2 &= -0.852.\end{aligned}$$

这里需要注意:

$$CV = 100 \times \frac{s}{\bar{x}}$$

取的是百分数. $g_1=0.035$, 向右微偏. g_1, g_2 的绝对值较小, 可以认为是来自正态总体的数据. 关于正态性的检验, 以后还将介绍.

例 1.2 某单位对 100 名女学生测定血清总蛋白含量(g/L), 数据如下:

74.3	78.8	68.8	78.0	70.4	80.5	80.5	69.7	71.2	73.5
79.5	75.6	75.0	78.8	72.0	72.0	72.0	74.3	71.2	72.0
75.0	73.5	78.8	74.3	75.8	65.0	74.3	71.2	69.7	68.0
73.5	75.0	72.0	64.3	75.8	80.3	69.7	74.3	73.5	73.5
75.8	75.8	68.8	76.5	70.4	71.2	81.2	75.0	70.4	68.0
70.4	72.0	76.5	74.3	76.5	77.6	67.3	72.0	75.0	74.3
73.5	79.5	73.5	74.7	65.0	76.5	81.6	75.4	72.7	72.7
67.2	76.5	72.7	70.4	77.2	68.8	67.3	67.3	67.3	72.7
75.8	73.5	75.0	73.5	73.5	73.5	72.7	81.6	70.3	74.3
73.5	79.5	70.4	76.5	72.7	77.2	84.3	75.0	76.5	70.4

计算均值、方差、标准差、变异系数、偏度、峰度.

解 通过计算,得

$$\begin{aligned}\bar{x} &= 73.660, & s^2 &= 15.524, \\ s &= 3.940, & CV &= 5.349, \\ g_1 &= 0.061, & g_2 &= 0.034.\end{aligned}$$

从计算结果看,偏度、峰度的绝对值皆较小,可以认为数据是取自正态总体的样本,即数据的总体分布是正态分布.

对于容量 n 较大的数据,有时可以采用分组的办法进行统计.这也可以通过 SAS 程序进行计算,可参看下例.

例 1.3 某电瓷厂的某种悬式绝缘子机电破坏负荷试验数据(单位:吨)分组表示如表 1.1. 计算这批分组数据的均值、方差、标准差、变异系数、偏度、峰度.

表 1.1 绝缘子机电负荷破坏试验数据

组 段	组 中 值	组 频 数
5.5~6.0	5.75	4
6.0~6.5	6.25	3
6.5~7.0	6.75	15
7.0~7.5	7.25	42
7.5~8.0	7.75	49
8.0~8.5	8.25	78
8.5~9.0	8.75	50
9.0~9.5	9.25	31
9.5~10.0	9.75	5

解 通过计算,得

$$\begin{aligned}\bar{x} &= 8.100, & s^2 &= 0.629, \\ s &= 0.793, & CV &= 9.789, \\ g_1 &= -0.382, & g_2 &= 0.057.\end{aligned}$$

这批数据的容量 $n=277$. 需要注意,对于这批分组数据,实际上是以组中值当成各组段中实际数据的代表(它不一定是实际数据),因此,以上算得的各数字特征是原始数据的数字特征的近似.

下面一个例子是确定型数据的类型,我们来进行分析.

例 1.4 1952~1997 年我国人均国内生产总值数据如表 1.2(单位:元). 计算这批数据的数字特征.