

# 数理统计理论、方法、应用 和软件计算

刘顺忠 编著



华中科技大学出版社  
E-mail: hustpp@wuhan.cngb.com

- 强调理论、方法、实例和软件计算相统一
- 注重基本方法的使用技巧，有助于提高实际应用能力
- 通俗易懂，特别适合非数学专业人员学习
- 依据例题，逐步操作，轻松学习SPSS13.0

# 数理统计理论、方法、应用

## 和软件计算

刘顺忠 编著



华中科技大学出版社

## 图书在版编目(CIP)数据

数理统计理论、方法、应用和软件计算/刘顺忠 编著  
武汉:华中科技大学出版社,2005年9月  
ISBN 7-5609-3445-5

I. 数…  
II. 刘…  
III. 数理统计-理论-应用  
IV. O212

## 数理统计理论、方法、应用和软件计算

刘顺忠 编著

责任编辑:沈旭日

封面设计:刘卉

责任校对:刘飞

责任监印:熊庆玉

出版发行:华中科技大学出版社

武昌喻家山 邮编:430074 电话:(027)87557437

录 排:武汉佳年华科技有限公司

印 刷:湖北新华印务有限公司

开本:787×960 1/16

印张:13.25

字数:225 000

版次:2005年9月第1版

印次:2005年9月第1次印刷

定价:18.80元

ISBN 7-5609-3445-5/O·361

(本书若有印装质量问题,请向出版社发行部调换)

## 内 容 提 要

本书将数理统计的原理、方法、计算、实例和SPSS软件计算相结合，以数理统计原理和方法为出发点，将手工计算过程和SPSS软件计算过程相结合，使读者能够从数学角度深刻理解SPSS软件的数据处理方法、过程和处理结果。本书介绍了大量的非参数统计方法，这些方法是当前社会科学实证研究的锐利武器，这也是本书的独特之处。

全书共有11章和两个附录。前2章介绍数理统计基本原理，从第3章到第11章介绍了实际工作中经常使用的数理统计方法及其SPSS软件计算方法。分别为：数理统计的基本概念、参数估计和假设检验、相关分析、回归分析、定性数据统计分析、非参数检验、方差分析、聚类分析、典型相关分析、主成分分析、因子分析。

本书可作为各大专院校非数学专业如经济学、统计学、市场营销学、医学、心理学、人文地理学、社会学、管理学等专业学生学习数理统计及SPSS软件的教材，也可供统计人员、科研人员以及广大自然科学工作者做科研定量分析时参考。同时，还可作为高等院校财经类专业研究生和本科生掌握SPSS统计分析方法和软件使用方法的工具书。

# 序 言

---

数据量庞大和数据关系越来越复杂是目前统计应用中遇到的主要问题,只有快速准确地处理庞大的数据量和复杂的数据关系,才能够有效解决实际工作和科研中遇到的问题。电子计算机的广泛普及和软件的不断完善,为统计分析中的数据处理和计算提供了强有力的工具,促进了统计方法在各个领域中的广泛应用。目前有关统计分析的软件主要有SPSS(Statistics Package for Social Science,社会科学统计软件包)和SAS(Statistical Analysis System,统计分析软件系统),这些软件大大方便了统计数据的处理和分析,特别是SPSS软件由于使用Windows的窗口操作方式进行数据的处理和分析,避免了繁琐的编程工作,是深受非数学专业人员欢迎的统计分析软件。统计软件是专用软件,不同于不需要专业知识的普通通用软件(如Word)。只有掌握了统计分析理论,才能够选择相应的软件模块对数据进行处理分析,并对软件输出结果进行科学解释,从而为科学决策提供有力的数据支持。

具有数理统计专业基础的读者在学习SPSS软件过程中经常遇到这样的问题:首先,如何组织和加工数据,才能利用SPSS软件进行数据处理和分析?其次,怎样操作SPSS软件对数据进行处理和分析以及如何理解和合理解释分析结果?……随着计算机的普及和数据处理的需要,目前越来越多的科研工作者、高校学生、市场营销人员和管理工作者使用统计软件对统计数据进行处理和分析,他们不具备数理统计专业基础,在用SPSS软件对数据处理过程中不但会面对统计专业人员面临的SPSS软件使用问题,更为重要的是需要学习和理解各种数理统计方法的原理及其应用。

我在市场调研公司和广告公司工作期间,以及在攻读博士期间都深感统计分析在实际工作和科研中的重要作用。从2000年开始,我就开始了本书的编写工作,为了适合非数学专业人员的特点,本书在讲述数理统计理论方法的同时,侧重从这些方法的计算过程出发讲解数理统计分析方法的应用。本书对各种数理统计方法的计算过程都给出了手工计算实例,不但便于读者理解数理统计方法,也方便读者理解SPSS软件的计算过程。经过给研究生和高年级本科生授课,使初稿不断完善,形成了现在的这本书。

本书将数理统计的原理、方法、计算、实例和SPSS软件计算相结合,以数理统计原理和方法为出发点,将手工计算过程和SPSS软件计算过程相结合。本书各章结构基本相同,都是数理统计原理介绍,计算方法,手工实例计算和统计软件实例计算。

本书的前 2 章介绍数理统计的基本原理,第 3 章到第 11 章介绍了实际工作中经常使用的数理统计方法及其 SPSS 软件计算方法,书籍的附录介绍了数理统计的基本概念、线性代数基本概念和 SPSS 软件的基本操作。

本书可作为各大专院校非数学专业如经济学、统计学、市场营销学、医学、心理学、人文地理学、社会学、管理学等专业学生学习数理统计及 SPSS 软件教科书,也可供统计人员、科研人员以及广大自然科学工作者做科研定量分析时参考。同时,还可作为高等院校财经类专业研究生和本科生掌握 SPSS 统计分析方法和软件使用的工具书。由于水平所限,书中的问题与错误在所难免,恳请读者批评指正。

刘顺忠

2005 年 5 月于武昌南湖

# 目 录

---

---

<b>第1章 数理统计的基本概念</b>	.....	(1)
1.1 抽样方法和测量水平	.....	(1)
1.1.1 总体和样本	.....	(1)
1.1.2 抽样方法	.....	(2)
1.1.3 测量水平	.....	(3)
1.2 统计量及其计算	.....	(4)
1.3 统计量分布	.....	(7)
1.3.1 正态分布	.....	(7)
1.3.2 $\chi^2$ 分布	.....	(8)
1.3.3 $t$ 分布	.....	(9)
1.3.4 $F$ 分布	.....	(10)
1.4 分位数的基本概念	.....	(10)
1.5 用 SPSS 软件计算本章例题	.....	(11)
1.5.1 计算方法与步骤	.....	(11)
1.5.2 输出结果解释	.....	(12)
<b>第2章 参数估计和假设检验</b>	.....	(13)
2.1 点估计和区间估计的基本原理	.....	(13)
2.1.1 点估计的基本原理	.....	(13)
2.1.2 区间估计的基本原理	.....	(15)
2.2 假设检验的基本原理	.....	(16)
2.2.1 假设检验的基本思想	.....	(16)
2.2.2 假设检验的两类错误	.....	(18)
2.3 正态分布的假设检验	.....	(18)
2.3.1 单样本 $t$ 检验	.....	(18)
2.3.2 两个独立样本 $t$ 检验	.....	(19)
2.3.3 两个配对样本 $t$ 检验	.....	(21)
2.3.4 正态分布假设检验的 SPSS 求解	.....	(23)
<b>第3章 相关分析</b>	.....	(28)
3.1 简单相关分析	.....	(29)
3.1.1 简单相关系数的计算	.....	(29)

---

3.1.2 简单相关系数的显著性检验	(31)
3.1.3 相关系数及其显著性的计算实例	(32)
3.1.4 简单相关系数及其显著性检验的SPSS求解	(33)
3.2 秩相关分析	(34)
3.2.1 Spearman秩相关系数	(34)
3.2.2 Kendall秩相关系数	(35)
3.2.3 秩相关及其显著性检验的SPSS求解	(37)
3.3 偏相关分析	(39)
3.3.1 偏相关系数的计算	(39)
3.3.2 偏相关系数计算实例	(40)
3.3.3 偏相关系数及其显著性检验的SPSS求解	(41)
<b>第4章 回归分析</b>	(43)
4.1 回归分析概述	(43)
4.2 一元线性回归	(44)
4.2.1 一元线性回归模型的数学形式	(44)
4.2.2 一元线性模型的参数估计及其性质	(46)
4.2.3 一元线性回归方程的显著性检验	(47)
4.2.4 一元线性回归残差分析	(49)
4.2.5 一元线性回归模型计算实例	(50)
4.2.6 一元线性回归模型的SPSS计算	(52)
4.3 多元线性回归	(55)
4.3.1 多元线性回归模型的数学形式	(55)
4.3.2 多元线性回归模型的基本假定	(56)
4.3.3 多元回归线性模型的参数估计及其性质	(57)
4.3.4 多元线性回归方程的显著性检验和残差分析	(58)
4.3.5 标准化回归系数	(60)
4.3.6 多元线性回归模型计算实例	(61)
4.3.7 多元线性回归模型的SPSS计算	(64)
4.4 回归模型违反假设的处理	(66)
4.4.1 自相关问题及其解决	(66)
4.4.2 异方差问题及其解决	(70)
4.4.3 多重共线性问题及其解决	(74)
4.5 其他回归模型	(79)
4.5.1 虚拟变量回归模型	(79)
4.5.2 可以化为线性回归的曲线回归	(80)

---

4.5.3 虚拟变量回归模型和可以化为线性回归的曲线回归模型的 SPSS 求解 .....	(82)
<b>第5章 定性数据统计分析 .....</b>	(84)
5.1 列联表分析 .....	(84)
5.1.1 列联表的概念 .....	(84)
5.1.2 列联表的分析方法 .....	(85)
5.1.3 列联表分析的计算实例 .....	(85)
5.1.4 列联表求解的 SPSS 计算步骤 .....	(86)
5.2 Logistic 回归分析 .....	(88)
5.2.1 Logistic 回归分析的模型 .....	(88)
5.2.2 Logistic 模型的参数估计 .....	(89)
5.2.3 Logistic 回归系数的含义及其显著性检验 .....	(89)
5.2.4 SPSS 求解 Logistic 回归实例和步骤 .....	(90)
5.3 多维尺度法 .....	(92)
5.3.1 古典 MDS .....	(92)
5.3.2 古典 MDS 求解实例 .....	(93)
5.3.3 非度量 MDS .....	(94)
5.3.4 SPSS 求解非度量 MDS 计算 .....	(95)
<b>第6章 非参数检验 .....</b>	(98)
6.1 单样本非参数检验 .....	(98)
6.1.1 $\chi^2$ 检验 .....	(98)
6.1.2 二项式分布检验 .....	(101)
6.1.3 游程检验 .....	(104)
6.1.4 柯尔莫哥洛夫-斯米诺夫单样本检验 .....	(107)
6.2 两个独立样本的非参数检验 .....	(109)
6.2.1 Mann-Whitney 检验 .....	(109)
6.2.2 K-S 双样本检验 .....	(111)
6.2.3 两个独立样本的非参数检验 SPSS 求解步骤 .....	(113)
6.3 多个独立样本的非参数检验 .....	(114)
6.3.1 中位数检验 .....	(114)
6.3.2 Kruskal-Wallis 检验 .....	(116)
6.3.3 多个独立样本的非参数检验 SPSS 求解 .....	(118)
6.4 两个相关样本的非参数检验 .....	(119)
6.4.1 Mcnemar 变化显著性检验 .....	(119)
6.4.2 Mcnemar 变化显著性检验的 SPSS 求解 .....	(121)

6.4.3 符号检验 .....	(121)
6.4.4 两相关样本 Wilcoxon 符号平均秩检验 .....	(123)
6.4.5 两相关样本符号和 Wilcoxon 平均秩检验的 SPSS 求解 .....	(124)
<b>6.5 多个相关样本的非参数检验 .....</b>	<b>(125)</b>
6.5.1 Cochran Q 检验 .....	(126)
6.5.2 Cochran Q 检验的 SPSS 求解 .....	(128)
6.5.3 Friedman 双向评秩方差检验 .....	(128)
6.5.4 Kendall 协和系数检验 .....	(130)
6.5.5 Friedman 双向评秩方差检验和 Kendall 协和系数检验的 SPSS 求解 .....	(131)
<b>第7章 方差分析 .....</b>	<b>(132)</b>
7.1 单因素方差分析 .....	(132)
7.1.1 单因素方差分析的原理 .....	(132)
7.1.2 单因素方差分析的计算过程 .....	(134)
7.1.3 单因素方差分析的求解实例 .....	(135)
7.1.4 单因素方差分析的 SPSS 求解 .....	(135)
7.2 无交互作用双因素方差分析 .....	(136)
7.2.1 无交互作用双因素方差分析原理 .....	(136)
7.2.2 无交互作用双因素方差分析的计算过程 .....	(139)
7.2.3 无交互作用双因素方差分析的求解实例 .....	(139)
7.2.4 无交互作用双因素方差分析的 SPSS 求解 .....	(140)
7.3 有交互作用双因素方差分析 .....	(142)
7.3.1 有交互作用双因素方差分析原理 .....	(142)
7.3.2 有交互作用双因素方差分析的计算过程 .....	(145)
7.3.3 有交互作用双因素方差分析的求解实例 .....	(146)
7.3.4 有交互作用双因素方差分析的 SPSS 求解 .....	(147)
<b>第8章 聚类分析 .....</b>	<b>(149)</b>
8.1 相似度的测量 .....	(149)
8.1.1 样本点间距离的测量 .....	(149)
8.1.2 类间距离测量 .....	(150)
8.2 系统聚类法 .....	(150)
8.2.1 系统聚类法的思想和计算过程 .....	(150)
8.2.2 系统聚类法的计算实例 .....	(150)
8.2.3 SPSS 系统聚类法求解步骤 .....	(152)
8.3 快速聚类法 .....	(153)

---

8.3.1 K-Means 法思想 .....	(153)
8.3.2 SPSS 的 K-Means 法求解 .....	(154)
<b>第9章 典型相关分析 .....</b>	<b>(156)</b>
9.1 典型相关的基本思想 .....	(156)
9.2 典型相关分析的原理 .....	(157)
9.2.1 典型相关变量的确立原则 .....	(157)
9.2.2 典型相关变量的求解方法 .....	(157)
9.2.3 典型相关的检验 .....	(159)
9.2.4 典型相关变量的求解步骤 .....	(160)
9.3 典型相关分析的求解实例 .....	(160)
9.4 SPSS 求解典型相关 .....	(163)
<b>第10章 主成分分析 .....</b>	<b>(165)</b>
10.1 主成分分析的基本思想 .....	(165)
10.2 主成分分析的基本原理 .....	(166)
10.2.1 求解主成分的相关定理 .....	(166)
10.2.2 主成分求解步骤 .....	(166)
10.2.3 主成分的辅助分析技术 .....	(167)
10.2.4 主成分求解的数据预处理 .....	(168)
10.3 主成分分析求解实例 .....	(169)
10.4 SPSS 求解主成分 .....	(170)
<b>第11章 因子分析 .....</b>	<b>(173)</b>
11.1 因子分析的基本思想 .....	(173)
11.2 因子分析的原理 .....	(174)
11.2.1 因子负荷的统计意义 .....	(174)
11.2.2 因子负荷矩阵的估计 .....	(175)
11.2.3 因子旋转 .....	(176)
11.2.4 因子分数 .....	(176)
11.2.5 因子分析的求解步骤 .....	(177)
11.3 因子分析的计算实例 .....	(177)
11.4 SPSS 求解因子分析 .....	(179)
<b>附录一 SPSS 软件及其基本操作 .....</b>	<b>(184)</b>
<b>附录二 本书例题录入示例及文件名 .....</b>	<b>(191)</b>
<b>附录三 线性代数和概率论基本知识 .....</b>	<b>(195)</b>
<b>参考文献 .....</b>	<b>(197)</b>



## 数理统计的基本概念

由于偶然性和不确定性,数据的随机性常常是无法避免的。数理统计是对带有随机性影响的数据进行有效搜集、整理、分析和推断,将杂乱无章的各种数据转化为简洁有序和容易解释的信息,并对观察到的现象作出推断或预测,为决策提供依据的数学方法。掌握数理统计的基本概念是理解和使用各种数理统计分析方法的前提,下面首先介绍数理统计分析中常用的基本概念。

### 1.1 抽样方法和测量水平

#### 1.1.1 总体和样本

在数理统计中,把研究对象的全体称为总体,而把组成总体的每一个单元称为个体。例如,若研究某城市人口的年龄结构,则该城市全部人口就构成了总体,而其中每个城市公民就是个体;若研究某纺织厂一天生产的 20000 件服装的合格率,则这 20000 件服装就组成一个总体,而每一件服装就是一个个体。

在实际研究中,我们关心的常常是研究对象的某个数量指标  $X$ (如人口寿命,产品是否合格),它是一个随机变量。因此,总体通常是指某个随机变量取值的全体,每个个体就是其中的一个实数。因此,可以把总体和数量指标  $X$  可能取值的全体组成的集合等同起来,随机变量  $X$  的分布就是总体的分布。

从一个总体  $X$  中,随机抽取  $n$  个个体  $X_1, X_2, \dots, X_n$ ,其中  $x_i (i=1, 2, \dots, n)$  是第  $i$  个个体抽样观察结果,称  $x_1, x_2, \dots, x_n$  为总体  $X$  的一组样本观察值。对于某一次抽样结果来说,它是完全确定的一组数,但是,由于抽样的随机性,所以它又是随每次抽样而改变的,这样每个  $x_i$  都可以看作某一个随机变量  $X_i (i=1, 2, \dots, n)$  所取的观察值。将  $X_1, X_2, \dots, X_n$  称为容量为  $n$  的样本,  $x_1, x_2, \dots, x_n$  就是样本  $X_1, X_2, \dots, X_n$  的一组观察值,称为样本值。

比如,在研究某城市人口年龄结构时,从全部城市公民  $X$ (总体)中随机抽取

1000人,记这1000人为 $X_1, X_2, \dots, X_{1000}$ ,则 $X_1, X_2, \dots, X_{1000}$ 构成样本。若样本 $X_1, X_2, \dots, X_{1000}$ 的取值记为 $x_1, x_2, \dots, x_{1000}$ ,则 $x_1, x_2, \dots, x_{1000}$ 就是样本 $X_1, X_2, \dots, X_{1000}$ 的一组观察值。

## 1.1.2 抽样方法

若要获得总体中的样本,就必须采取一定抽取样本的方法。抽样方法主要有两种:概率抽样和非概率抽样。

### 1. 概率抽样

所谓概率抽样就是根据已知的概率来选取样本,无须在抽样过程中对样本进行判断和抽选。概率抽样主要有以下4种。

简单随机抽样是最基本的抽样形式,它是完全随机地选取样本。该方法要求有一个完美的抽样框,或每一个个体的详尽名单。如在调查某高校学生月消费的研究中,按照随机数表选取给定学号的学生进行抽样调查,就是简单随机抽样。

分层抽样是先将总体分成不同的层,然后在每一层内进行简单随机抽样。分层抽样可以防止简单随机抽样造成的样本构成同总体构成不成比例的现象。如在研究城市人口年收入状况时,先将一个城市按照不同城区进行分组,然后在各个组中采取简单随机抽样的方法进行抽样,这就是分层随机抽样。

整群抽样是将一组被调查者视作一个抽样单位而不是个体的抽样方法。如在市场研究中,选取的是某个小区的指定楼号的所有居民进行调查,这就是整群抽样。

等距抽样又称系统抽样,是在抽样框中每隔一定距离抽取一个样本。如在研究某大学新生家庭收入情况时,按照学号每隔100号抽取一个学生进行调研,就是等距抽样。

### 2. 非概率抽样

非概率抽样不是完全按照随机原则,而是按照研究人员主观判断或研究方便的原则选取样本。非概率抽样主要有以下4种。

方便抽样又称偶遇抽样,是根据调查者的方便来抽取样本的抽样方法。如在市场调查中经常使用的街头拦截访问就是方便抽样。

判断抽样又称目的抽样,它是凭研究人员的主观意愿、经验和知识,从总体中抽取典型代表性样本的方法,主要是根据研究目的选取具有代表性的样本或极端样本。如企业质量管理人员根据以往的经验,从所有会出现次品的产品中抽取样本,这就是判断抽样。

配额抽样是非概率抽样中最流行的一种。它首先将总体按照一定标志分为若

干组,然后在每组中用方便抽样或判断抽样来选取给定数量的样本。如在研究某种新型移动通信服务时,将目前手机用户按照月话费额进行分组,然后在每组中根据方便抽样原则抽取给定数量的样本,这就是配额抽样。

雪球抽样是以滚雪球的方式来抽取样本,即在少量样本的基础上不断扩充以获得更多样本的信息。这种抽样方式运用的前提是总体各样本之间存在一定的联系,可以在不甚了解总体的情况下对总体或总体部分样本情况进行把握。如在研究知识密集型服务业创新机理过程中,由于知识密集型服务业显著特征就是相互之间联系密切,因此可以以少量知识密集型服务业企业作为原始“雪球”,采取“滚雪球”的抽样方式不断扩充对知识密集型服务业企业的抽样,这就是雪球抽样。

随机抽样是以概率论为基础,容易进行估计、检验和控制抽样误差,因此以概率论为基础的数理统计分析适合处理随机抽样得到的数据。从原则上讲,数理统计分析的方法不适合处理非随机抽样得到的数据,但是,由于非随机抽样简单易行,适合对数据进行探索性分析。

### 1. 1. 3 测量水平

当我们对样本进行测量时,也就是赋予给定变量以具体数值。变量测量水平主要分为以下4类,不同类别的测量水平所容许的运算是不同的。

#### 1. 分类量表

分类量表又称为称名量表。它是使用一些数字或其他符号来识别不同事务所属类别的一种方法。这些数字或符号就构成一个称名量表,它是测量水平最低的量表。

比如,在对性别进行分类时,男性用1表示,女性用2表示,这时,数字1和2就构成一个称名量表。又比如,在对学历进行分类时,学士用B表示,硕士用S表示,博士用D表示,这时,符号B,S,D也构成了一个分类量表。

分类量表运算是把一个给定的集合分成一个个互不相交的子集合,它们所包含的关系是等价关系,即任何一个子集成员在被标度的特征方面是相同的。比如在使用数字1和2对性别进行标度时,所有被标度为数字1的都是男性,所有被标度为数字2的都是女性。

#### 2. 顺序量表

顺序量表又称评秩量表。如果在分类量表中,所有子集之间都保持一个顺序关系,则该量表构成一个顺序量表。

比如,按照年龄对人进行划分,则人可以划分为儿童(A)、少年(B)、青年(C)、中年(D)和老年(E)5个子集。每个子集如果按照年龄排序,则存在一个偏序关系,

即  $A < B < C < D < E$ 。这样,不同年龄的人之间就存在一个顺序关系,因此也就构成了一个顺序量表。

顺序量表包含偏序关系和等价关系,但是,还不能进行算术运算。

### 3. 间隔量表

当一个量表具有顺序量表的特征,并且量表中任意两个数之间的距离为已知时,其测量水平就高于顺序量表,称为间隔量表。间隔量表的特征是有公共的、不变的测量单位,它将一个实数赋给样本中的每个个体。在间隔量表中,零点和单位是任意的。

比如温度的测量,通常采用两种不同的量表——摄氏量表和华氏量表。在测量温度时,测量单位是任意的;这两种量表具有不同的零点和单位,却包含了相同的信息量,一个量表中温度差的比值和另一个量表中温度差的比值是相同的。

在摄氏温度量表和华氏温度量表中,表 1-1 所示的温度是相同的。

表 1-1 摄氏和华氏温度对照表

量表	相同温度			
	摄氏 /°C	0	10	30
华氏 /°F	32	50	86	212

在摄氏温度量表中  $\frac{30-10}{10-0}=2$ , 在华氏温度量表中  $\frac{86-50}{50-32}=2$ , 在使用不同间隔量表时其温度差的比值是相同的,与零点和单位没有关系,由于间隔量表是将实数赋给样本观测值,所以该量表可以进行算术运算。

### 4. 比例量表

当间隔量表以真正的零点作为原点时,就称为比例量表。比例量表是测量水平最高的量表。

比如,在测量物体重量时,可以使用公斤和市斤,两个量表都具有零点。任意两个重量间的比值与测量单位无关。

比如,在测量两个物体重量时,公斤量表分别为 10 和 50,市斤量表分别为 20 和 100,则  $50/10=5$ ,  $100/20=5$ ,两个重量间比值与测量单位无关,比例量表可以进行算术运算。

## 1.2 统计量及其计算

通过试验或观察取得的样本值往往呈现为一堆杂乱无章的数据。因此,在利用

样本数据推断总体特征时,往往不能够直接利用样本数据,而需要对它们进行一定的加工和整理,通过计算出来的一些量来对总体进行推断,否则无法有效利用样本数据中的信息。

**【例 1-1】** 某系随机抽取 30 名学生的英语考试成绩,以调查该系学生英语考试情况。现得到下列数据,试对该系学生英语学习成绩和相互之间的差距状况进行大致分析。

73	61	67	71	90	88
66	91	75	80	71	61
80	61	68	66	87	60
64	80	83	80	86	75
44	75	60	98	65	56

**【解】** 显然,面对这些参差不齐的数据,很难得出什么印象。但是,只要对这些数据稍事加工,便能进行分析,并得出结论。设这 30 名学生的英语成绩分别为  $x_1, x_2, \dots, x_{30}$ , 则

$$\bar{x} = \frac{1}{30} \sum_{i=1}^{30} x_i = 72.73 \quad (1-1)$$

$$S = \sqrt{\frac{1}{30} \sum_{i=1}^{30} (x_i - \bar{x})^2} = 12.24 \quad (1-2)$$

由  $\bar{x}$  得出该系学生英语平均成绩为中等,由  $S$  得出该系学生英语成绩差距不大的结论。因此,对样本数据进行加工和处理是十分重要的。

对样本数据进行加工和处理主要就是构造统计量。如果用数学的语言来描述,则统计量就是一个不含未知参数的样本已知函数。设样本为  $X_1, X_2, \dots, X_n$ , 则统计量通常记为

$$T = T(X_1, X_2, \dots, X_n) \quad (1-3)$$

设  $X_1, X_2, \dots, X_n$  为总体  $X$  的样本,其观测值为  $x_1, x_2, \dots, x_n$ , 则常用统计量如下:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1-4)$$

$\bar{x}$  称为样本均值。

$$S = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (1-5)$$

$S$  称为样本标准差,  $S^2$  称为样本方差。

众数就是一组样本中出现次数最多的变量值,记为  $M_0$ 。对统计量的描述还有最大值

$$\text{MAX} = \max(x_1, x_2, \dots, x_n) \quad (1-6)$$

最小值

$$\text{MIN} = \min(x_1, x_2, \dots, x_n) \quad (1-7)$$

样本极差

$$D^* = \text{MAX} - \text{MIN} \quad (1-8)$$

已知  $X_1, X_2, \dots, X_n$  的取值  $x_1, x_2, \dots, x_n$ , 假设  $x_1 < x_2 < \dots < x_n$ , 则

$$M^* = \begin{cases} x_{\frac{n+1}{2}}, & n \text{ 为奇数} \\ \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}), & n \text{ 为偶数} \end{cases} \quad (1-9)$$

$M^*$  称为样本中位数, 中位数是将样本从中间等分为两部分。

与中位数类似的还有四分位数, 四分位数也称为四分位点, 它是通过 3 个点将全部样本数据分为四部分, 其中每部分包含 25% 的样本, 处于分位点上的样本值就是四分位数。很显然, 中间的四分位数就是中位数, 通常所说的四分位数是指第一个四分位数(下四分位数)和第 3 个分位数(上四分位数)。

设下四分位数位为  $Q_L$ , 中间四分位数为  $Q_M$ , 上四分位数为  $Q_U$ , 则  $Q_L$  的位置是  $\frac{N+1}{4}$ ,  $Q_M$  的位置是  $\frac{2(N+1)}{4}$ ,  $Q_U$  的位置是  $\frac{3(N+1)}{4}$ , 显然  $Q_M = M^*$ 。

同中位数计算相同, 当四分位数不在某一个数值上且不为整数时, 可以根据四分位数的位置, 按照比例分摊四分位数两侧数值的差值。

根据例 1-1 的数据计算上述各统计量, 首先对例 1-1 中 30 个成绩数据从小到大排序, 结果如表 1-2 所示。

表 1-2 成绩数据升序排序表

序号	分数	序号	分数	序号	分数
1	44	11	66	21	80
2	56	12	67	22	80
3	60	13	68	23	80
4	60	14	71	24	83
5	61	15	71	25	86
6	61	16	73	26	87
7	61	17	75	27	88
8	64	18	75	28	90
9	65	19	75	29	91
10	66	20	80	30	98

根据表 1-2 中的数据, 可计算出该班分数样本数据的各个统计量的值: