



全国工程专业学位研究生教育国家级规划教材

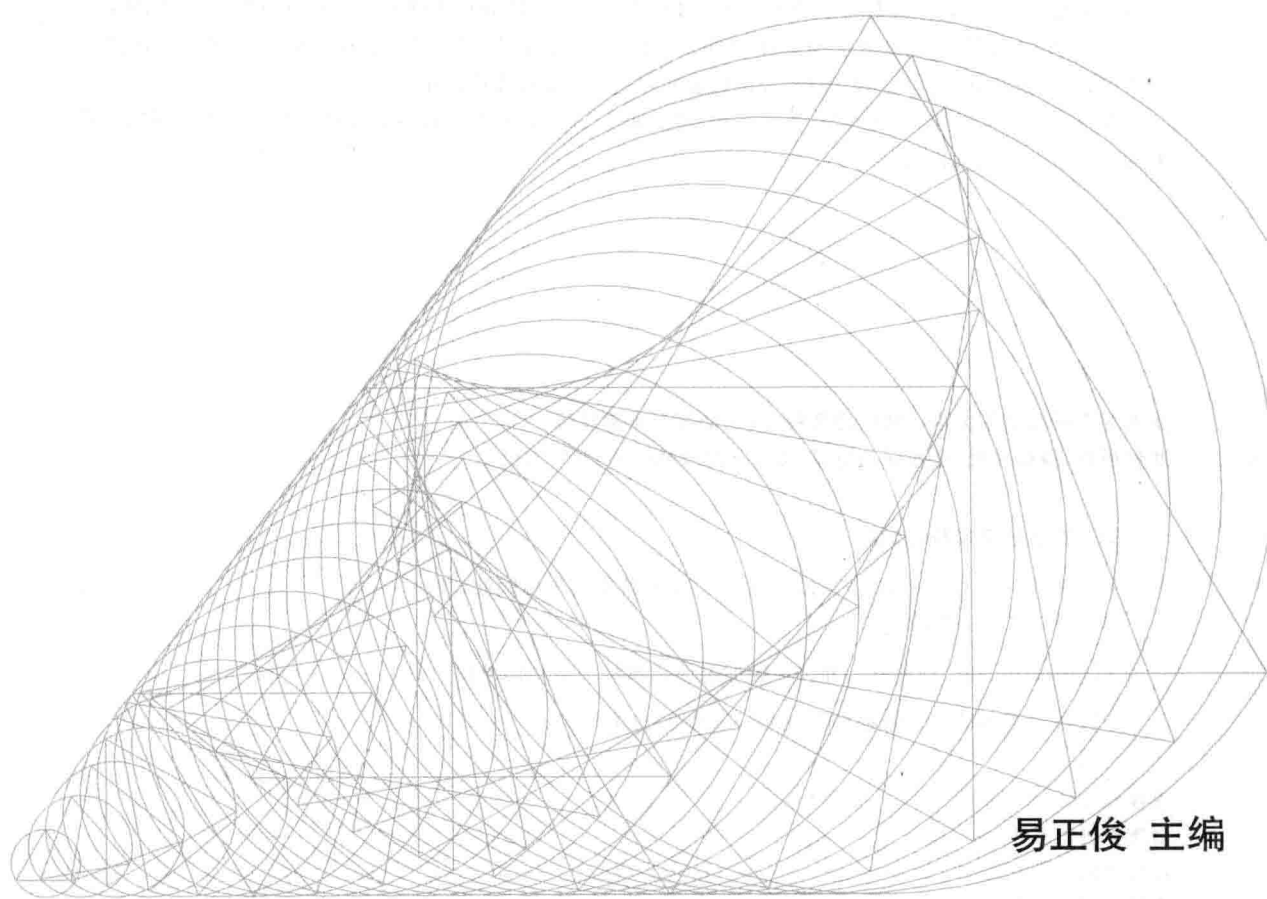
易正俊 主编

数理统计及其工程应用



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

清华大学出版社



易正俊 主编

数理统计及其工程应用

清华大学出版社

内 容 简 介

本书是专为工程硕士和专业硕士学习数理统计及其工程应用而编写的教材. 全书共 8 章, 主要内容有: 统计的基本概念及抽样分布、参数估计、假设检验、方差分析、正交试验设计、回归分析、系统聚类分析和主成分分析; 每章配有 R 软件、SPSS 软件或 Excel 软件等统计分析软件及相应的训练案例; 习题的设置依据培养学生不同能力的要求分为 A、B 两组, A 组主要是训练学生的应用能力, B 组是提升学生的理论基础水平, 书后附有概率基础知识回顾、分位数表和习题的部分答案或提示.

本书讲解简明扼要, 图文并茂, 注重应用, 覆盖面广, 也可以作为统计专业本科学生的教材及实际工作者的应用参考书和工具书.

本书封面贴有清华大学出版社防伪标签, 无标签者不得销售.

版权所有, 侵权必究. 侵权举报电话: 010-62782989 13701121933

图书在版编目(CIP)数据

数理统计及其工程应用/易正俊主编. --北京: 清华大学出版社, 2014

ISBN 978-7-302-36438-2

I. ①数… II. ①易… III. ①数理统计—教材 IV. ①O212

中国版本图书馆 CIP 数据核字(2014)第 095923 号

责任编辑: 刘 颖

封面设计: 常雪影

责任校对: 王淑云

责任印制: 刘海龙

出版发行: 清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址: 北京清华大学学研大厦 A 座 邮 编: 100084

社 总 机: 010-62770175 邮 购: 010-62786544

投稿与读者服务: 010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈: 010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者: 北京嘉实印刷有限公司

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 14.75 字 数: 355 千字

版 次: 2014 年 7 月第 1 版 印 次: 2014 年 7 月第 1 次印刷

印 数: 1~3000

定 价: 33.00 元

产品编号: 057930-01

编者从事数理统计课程教学二十几年,主要是讲授数理统计的原理和方法,侧重于理论方面的教学;学生学习这门课程也主要是应付期终考试,虽然学生取得了相当满意的成绩,但在实际应用中无法用所学的数理统计知识来解决实际问题,学生在工作中遇到数理统计的应用问题只有回到母校找老师请教解决的方法.出现这种情况主要是由于我们的教材和对学生的考核标准有问题,在实际中有用的内容老师没有讲,教材没有写(即便是写了也是打了*号),不作为考核内容.这给学生一种错误的感觉:数理统计课程只是考试有用,在实际中没有用.

全国工程专业学位研究生教育指导委员会根据社会实际工程领域对人才的需求,提出了工程硕士课程教学改革设想和指导性意见,旨在服务于行业创新发展的需求,提升职业能力,注重解决实际问题,提高在实践中发现问题、分析问题和解决问题的能力,通过整理和提炼实践工作中的问题,综合运用所学知识分析并解决问题,培养在实际工作中发现问题的敏感性、分析问题的科学性、处理问题的有效性.

“数理统计及其工程应用”是工程硕士培养的一门重要公共基础课程,在工程领域有广泛的应用.全国工程专业学位研究生教育指导委员会指派重庆大学易正俊教授组织编写《数理统计及其工程应用》教材,紧扣解决工程实际问题能力这一核心目标,为工程硕士、专业硕士的现有职业或未来职业提供专业支持,培养具有竞争优势的应用型创新创业人才.教材初稿完成后,由全国工程专业学位研究生教育指导委员会副主任陈子辰教授(浙江大学)和秘书长沈岩(清华大学)组织专家组成员齐欢教授(华中科技大学)、周杰教授(清华大学)、李大美教授(武汉大学)、韩中庚教授(解放军信息工程大学)对教材初稿进行讨论与修改.

本教材是专为工程硕士和专业硕士而编写的,具有以下特色:

1. 注重质量提升,突出职业需求导向,加强案例教学.案例的选取参考了国内外优秀教材和学术论文,博采众家之长,体现案例的实用性和趣味性,激发学生学习的积极性.

2. 训练学生借助统计分析软件(如 Excel 软件、R 软件和 SPSS 软件等)解决工程领域的实际问题,培养学生工程应用的意识和素质,提高解决工程实际问题的能力.

3. 基本概念和基本理论尽可能从学生熟悉的背景知识引入,采用几何图形等方法加强学生对基本理论和基本方法的理解,淡化比较复杂的理论推导,增强教材的可读性和可接受性.

4. 重视反例在学生理解、掌握基本概念和基本理论中的重要作用.

5. 习题的设置依据培养学生不同能力的要求分为 A, B 两组, A 组主要是训练学生的应用能力, B 组是提升学生的理论基础水平.

本教材由易正俊教授担任主编, 参加编写的作者还有颜军、刘朝林、荣腾中、彭智军、曹术存、袁玉兴和罗秀娟.

教材的编写得到教育部学位管理与研究生教育司、重庆大学研究生院、重庆大学数学与统计学院的资助; 重庆大学数学与统计学院穆春来和西南大学原校长宋乃庆对教材的编写提出了宝贵的意见. 编者在此深表感谢!

写好这样一本具有实际应用价值的教材, 编者深感难度很大. 由于编者学识有限, 书中不妥之处真诚地欢迎读者批评指正.

编者

2014 年 4 月

第 1 章 统计的基本概念及抽样分布	1
1.1 统计的基本概念	1
1.1.1 总体、样本与统计量	1
1.1.2 样本的联合分布函数和联合分布密度函数	2
1.1.3 统计量	2
1.2 顺序统计量、经验分布函数和直方图	4
1.2.1 顺序统计量	4
1.2.2 最大最小顺序统计量的分布	4
1.2.3 经验分布函数与直方图	5
1.3 抽样分布及分位数	7
1.3.1 正态分布的导出分布	7
1.3.2 抽样分布定理	10
1.3.3 下分位数	13
1.4 案例及统计分析软件的训练	14
训练项目 1 统计量的数字特征求解	14
训练项目 2 常用统计量的分布	21
训练项目 3 直方图、经验分布函数图	24
习题 1	28
第 2 章 参数估计	29
2.1 参数的点估计	29
2.1.1 矩估计法	29
2.1.2 极大似然估计	31
2.1.3 点估计的优良评价准则	33
2.2 参数的区间估计	38
2.2.1 置信区间的定义	38
2.2.2 单个正态总体参数的区间估计	38
2.2.3 双正态总体参数的区间估计	41
2.3 案例及统计分析软件的训练	44
训练项目 1 单个正态总体均值的区间估计	44
训练项目 2 两个正态总体均值差的区间估计	46

训练项目 3 两个正态总体方差比的区间估计	49
习题 2	49
第 3 章 假设检验	52
3.1 假设检验的基本概念	52
3.1.1 统计假设的设置	52
3.1.2 假设检验的基本思想	54
3.1.3 假设检验的步骤	56
3.2 参数假设检验	57
3.2.1 正态总体的参数假设检验	57
3.2.2 非正态总体的参数假设检验	69
3.3 非参数假设检验	70
3.3.1 总体分布函数的假设检验	71
3.3.2 独立性的假设检验	73
3.3.3 两总体分布比较的假设检验	77
3.4 案例及统计分析软件的训练	79
训练项目 1 单个正态总体均值的假设检验	79
训练项目 2 两个正态总体均值的假设检验	81
训练项目 3 两个正态总体方差的假设检验	84
训练项目 4 单总体分布的假设检验	84
训练项目 5 两个总体独立性假设检验	88
习题 3	90
第 4 章 方差分析	93
4.1 单因素方差分析	93
4.1.1 方差分析的基本原理	93
4.1.2 单因素方差分析	94
4.2 双因素方差分析	99
4.2.1 无交互作用的双因素方差分析	99
4.2.2 有交互作用的双因素方差分析	101
4.3 案例及统计分析软件训练	104
训练项目 1 单因素方差分析	104
训练项目 2 双因素方差分析	107
习题 4	110
第 5 章 正交试验设计	113
5.1 正交表与正交试验设计	113
5.1.1 正交表	113

5.1.2	正交试验设计	115
5.2	正交试验的结果分析	116
5.2.1	直观分析法	116
5.2.2	方差分析法	117
习题 5	119
第 6 章	回归分析	122
6.1	一元线性回归分析	122
6.1.1	一元线性回归模型	122
6.1.2	一元线性回归方程	123
6.1.3	回归参数的最小二乘估计	123
6.1.4	最小二乘估计的性质	124
6.1.5	显著性检验	128
6.1.6	预测与控制	129
6.2	非线性回归	132
6.3	多元线性回归	136
6.3.1	多元线性回归的数学模型	136
6.3.2	参数 β 的最小二乘估计	136
6.3.3	最小二乘估计的性质	140
6.3.4	显著性检验	141
6.4	案例及统计分析软件的训练	148
训练项目 1	一元线性回归分析	148
训练项目 2	多元线性回归分析	152
习题 6	154
第 7 章	系统聚类分析	158
7.1	系统聚类分析的原理	158
7.1.1	相似性度量	158
7.1.2	系统聚类法	160
7.2	案例及统计分析软件的训练	165
训练项目	系统聚类	165
习题 7	172
第 8 章	主成分分析	174
8.1	主成分分析的原理	174
8.1.1	主成分的几何解释	174
8.1.2	主成分的导出	175
8.1.3	特征值因子的筛选	177

8.1.4 主成分分析法	178
8.2 案例及统计分析软件的训练	181
训练项目 主成分分析	181
习题 8	189
附录 概率基础知识回顾	193
附表 常用数理统计表	200
附表 1 标准正态分布表	200
附表 2 t 分布分位数表	201
附表 3 卡方分布分位数表	203
附表 4 F 分布分位数表	204
附表 5 常用正交表	212
附表 6 符号检验临界值表	216
附表 7 秩和临界值表	216
部分习题参考答案	218
参考文献	226

统计的基本概念及抽样分布

1.1 统计的基本概念

1.1.1 总体、样本与统计量

总体(X)是指研究对象的全体; **个体**是组成总体的每个单元; **样本**是从总体中随机抽取的 n 个个体 X_1, X_2, \dots, X_n , n 称为**样本容量**. 一旦抽取了一个样本 X_1, X_2, \dots, X_n 进行试验, 抽取样本后得到一组数据 (x_1, x_2, \dots, x_n) , 这组数据称为样本 X_1, X_2, \dots, X_n 的一组观察值, 样本观察值是随着抽样的变化而变化的.

抽样的目的是用样本的特征去代表总体的特征. 在有些情况下, 我们不可能对总体的每个个体进行逐一试验, 特别是具有破坏性的试验. 如检验一批炮弹是否合格, 检验办法是将抽取的炮弹进行试爆, 这个检验就是破坏性的; 对总体的量很大, 个体比较小的检验, 虽然检验可能不是破坏性的, 但对总体的检验也不可能逐一进行. 如对豌豆种子的检验就是总体大、个体小. 抽样包括简单随机抽样、分群抽样和分层抽样. 本教材的抽样是指简单随机抽样, 简单随机抽样所获得的样本 X_1, X_2, \dots, X_n 称为简单样本 (simple sample). 它满足以下两点:

- (1) 独立性: 要求样本 X_1, X_2, \dots, X_n 为相互独立的随机变量;
- (2) 代表性: 要求每个样本 $X_i (i=1, 2, \dots, n)$ 与总体 X 具有相同的分布.

从简单样本的定义可以看出: 简单样本是有放回地抽取得到的样本. 在实际工作中, 我们的抽样都是无放回地抽样, 从理论上说就不再是简单样本. 但总体中个体的数目很大, 从中抽取一些个体对总体成分没有太大的影响, 可近似地看成有放回的抽样, 其样本仍可看成是独立同分布的.

简单随机抽样有很多的益处: 抽样单元的随机选取排除了调查者的偏见, 这种偏见可能调查者并没有意识到; 与完全枚举相比, 小样本减少很多成本, 调查更省时; 小样本的结论实际上可能比完全枚举更精确. 小样本的数据质量更容易监控, 完全枚举需要更多的员工去实施; 随机抽样技术使得抽样误差的估计变得可能; 在抽样设计时, 通常可以确定出满足预设误差水平的样本尺寸.

例 1.1 某灯泡厂进行技改并扩大了生产规模, 要求生产的灯泡寿命在 1000h 以上才算合格品. 现从技改后生产的第一批灯泡中随机抽取 4 个, 测得其使用寿命分别为 1200, 1120, 980, 1350 (单位: h), 试叙述总体、样本、样本容量、样本观察值.

解 总体是灯泡的使用寿命 X 的取值全体; 抽取 4 个个体 X_1, X_2, X_3, X_4 就是一个样本, 样本容量是 4; 1200, 1120, 980, 1350 是一组样本观察值.

1.1.2 样本的联合分布函数和联合分布密度函数

设总体 X 的分布函数为 $F(x)$, X_1, X_2, \dots, X_n 是来自总体 X 的样本, 则该样本的联合分布函数为

$$F(x_1, x_2, \dots, x_n) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) = \prod_{i=1}^n F(x_i).$$

当总体 X 是连续型随机变量且具有密度函数 $f(x)$ 时, 则样本的联合密度函数为

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i).$$

当总体 X 是离散型随机变量且具有分布律 $P(X=x_i)=p_i$ 时, 则样本 X_1, X_2, \dots, X_n 的联合分布律为

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \prod_{i=1}^n P(X = x_i).$$

例 1.2 假设总体 X 服从指数分布, 密度函数为

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0, \\ 0, & x < 0, \end{cases}$$

从总体中抽取容量为 n 的样本 X_1, X_2, \dots, X_n , 写出样本的联合分布密度函数.

解 X_i 的密度函数为

$$f(x_i) = \begin{cases} \lambda e^{-\lambda x_i}, & x_i \geq 0, \\ 0, & x_i < 0, \end{cases}$$

X_1, X_2, \dots, X_n 是简单随机样本, 且相互独立, 所以 X_1, X_2, \dots, X_n 的联合密度函数等于边际密度函数的乘积, 即

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i) = \lambda^n e^{-\sum_{i=1}^n x_i}, \quad x_i \geq 0, i = 1, 2, \dots, n.$$

例 1.3 设总体 $X \sim B(1, p)$, X_1, X_2, \dots, X_n 是来自总体的一个样本, 写出样本的联合密度函数.

解 总体 X 的概率函数可以写成

$$f(x) = p^x (1-p)^{1-x}, \quad x = 0, 1.$$

X_i 的概率函数为

$$f(x_i) = p^{x_i} (1-p)^{1-x_i}, \quad x_i = 0, 1.$$

样本 X_1, X_2, \dots, X_n 是相互独立的, 所以联合概率函数等于边际密度函数的乘积, 为

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i) = p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}, \quad x_i = 0, 1.$$

1.1.3 统计量

1. 统计量的定义

统计量是不含参数的样本 X_1, X_2, \dots, X_n 的函数或即便含参数, 但参数是已知的, 记为 $T=T(X_1, X_2, \dots, X_n)$, 一旦获得样本的观察值, 代入统计量得到一个数值. 例如 X_1, X_2, \dots, X_n 为来自总体 X 的一个样本, $T = \sum_{i=1}^n X_i$ 是一个统计量.

若 $X \sim N(\mu, \sigma^2)$, X_1, X_2, \dots, X_n 为 X 的一个样本, $\sum_{i=1}^n \frac{X_i - \mu}{\sigma}$ 就不是统计量, 因为含有未知的参数 μ 和 σ .

2. 常见统计量

(1) 样本均值: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.

(2) 样本方差: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)$.

(3) 样本标准差: $S = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}$.

(4) 样本 k 阶原点矩: $M_k = \frac{1}{n} \sum_{i=1}^n X_i^k, k = 1, 2, \dots$

(5) 样本 k 阶中心矩: $M_k^* = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^k, k = 1, 2, \dots$

3. 样本均值 \bar{X} 的性质

(1) $\sum_{i=1}^n (X_i - \bar{X}) = 0$.

(2) 若总体 X 的均值、方差存在, 且 $EX = \mu, DX = \sigma^2$, 则

$$E\bar{X} = \mu, \quad D\bar{X} = \frac{\sigma^2}{n}.$$

(3) 当 $n \rightarrow \infty$ 时, $\bar{X} \xrightarrow{P} \mu$.

证明 (1) $\sum_{i=1}^n (X_i - \bar{X}) = \sum_{i=1}^n X_i - n\bar{X} = n \frac{\sum_{i=1}^n X_i}{n} - n\bar{X} = n\bar{X} - n\bar{X} = 0$.

(2) $E\bar{X} = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} \sum_{i=1}^n EX_i = \frac{1}{n} \sum_{i=1}^n EX = \mu$,

$$D\bar{X} = D\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n DX_i = \frac{1}{n^2} \sum_{i=1}^n DX = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}.$$

(3) 由概率论中的大数定律知, 当 $n \rightarrow \infty$ 时, $\bar{X} \xrightarrow{P} \mu$.

4. 样本方差 S^2 的性质

(1) 如果 DX 存在, 则 $ES^2 = DX, EM_2^* = \frac{n-1}{n}DX$;

(2) 对任意实数 μ , 有 $\sum_{i=1}^n (X_i - \bar{X})^2 \leq \sum_{i=1}^n (X_i - \mu)^2$.

证明 (1) $ES^2 = E\left[\frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2\right)\right] = \frac{1}{n-1} \left(\sum_{i=1}^n EX_i^2 - nE\bar{X}^2\right)$

$$= \frac{n}{n-1} (EX^2 - E\bar{X}^2) = \frac{n}{n-1} (DX + (EX)^2 - D\bar{X} - (E\bar{X})^2)$$

$$= \frac{n}{n-1} \left(DX + (EX)^2 - \frac{DX}{n} - (EX)^2 \right) = DX,$$

$$M_2^* = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{(n-1)S^2}{n},$$

$$EM_2^* = E \frac{(n-1)S^2}{n} = \frac{(n-1)ES^2}{n} = \frac{(n-1)DX}{n};$$

$$\begin{aligned} (2) \sum_{i=1}^n (X_i - \bar{X})^2 &= \sum_{i=1}^n ((X_i - \mu) + (\mu - \bar{X}))^2 \\ &= \sum_{i=1}^n (X_i - \mu)^2 + n(\mu - \bar{X})^2 + 2(\mu - \bar{X}) \sum_{i=1}^n (X_i - \mu) \\ &= \sum_{i=1}^n (X_i - \mu)^2 + n(\mu - \bar{X})^2 - 2(\mu - \bar{X})(n\mu - \sum_{i=1}^n X_i) \\ &= \sum_{i=1}^n (X_i - \mu)^2 - n(\mu - \bar{X})^2 \leq \sum_{i=1}^n (X_i - \mu)^2. \end{aligned}$$

例 1.4 设总体 $X \sim U[0, \theta]$, $\theta > 0$, X_1, X_2, \dots, X_n 为样本, 求 $E\bar{X}, D\bar{X}, EM_2^*$.

解 $E\bar{X} = EX = \frac{\theta}{2}, D\bar{X} = \frac{1}{n}DX = \frac{1}{n} \frac{(\theta-0)^2}{12} = \frac{\theta^2}{12n}.$

$$EM_2^* = \frac{n-1}{n}DX = \frac{(n-1)\theta^2}{12n}.$$

1.2 顺序统计量、经验分布函数和直方图

1.2.1 顺序统计量

X_1, X_2, \dots, X_n 为总体 X 的样本, x_1, x_2, \dots, x_n 为样本观察值, 将样本观察值 x_1, x_2, \dots, x_n 按从小到大的递增顺序进行排列: $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$, $X_{(1)}$ 取最小观察值, $X_{(2)}$ 取次小值, \dots , $X_{(n)}$ 取最大观察值, $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ 称为顺序统计量. $X_{(1)}$ 称为最小顺序统计量, $X_{(n)}$ 称为最大顺序统计量. $R = X_{(n)} - X_{(1)}$ 称为极差, 极差在实际中用来衡量方差的大小, 反映了随机变量 X 取值的分散程度.

$$\tilde{X} = \begin{cases} X_{(\frac{n+1}{2})}, & n \text{ 为奇数,} \\ \frac{1}{2} \left(X_{(\frac{n}{2})} + X_{(\frac{n}{2}+1)} \right), & n \text{ 为偶数} \end{cases}$$

为样本中位数. 样本中位数反映了随机变量 X 在实轴上分布的位置特征.

1.2.2 最大最小顺序统计量的分布

设 $F(x), \varphi(x)$ 分别为总体 X 的分布函数和分布密度函数, X_1, X_2, \dots, X_n 为 X 的样本, $F_{X_{(n)}}(x), \varphi_{X_{(n)}}(x)$ 分别为 $X_{(n)}$ 的分布函数和分布密度函数, $F_{X_{(1)}}(x), \varphi_{X_{(1)}}(x)$ 分别为 $X_{(1)}$ 的分布函数和分布密度函数, 则对任意的实数 x , 有

$$\begin{aligned} F_{X_{(n)}}(x) &= P(X_{(n)} \leq x) = P(\max\{X_1, X_2, \dots, X_n\} \leq x) \\ &= P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) = \prod_{i=1}^n P(X_i \leq x) = F^n(x), \end{aligned}$$

$$\varphi_{X_{(n)}}(x) = [F^n(x)]' = nF^{n-1}(x)\varphi(x),$$

$$\begin{aligned} F_{X_{(1)}}(x) &= P(X_{(1)} \leq x) = P(\min\{X_1, X_2, \dots, X_n\} \leq x) \\ &= 1 - P(\min\{X_1, X_2, \dots, X_n\} > x) \\ &= 1 - P(X_1 > x, X_2 > x, \dots, X_n > x) \\ &= 1 - P(X_1 > x)P(X_2 > x) \cdots P(X_n > x) \\ &= 1 - [1 - F(x)]^n, \end{aligned}$$

$$\varphi_{X_{(1)}}(x) = -n[1 - F(x)]^{n-1}(-\varphi(x)) = n\varphi(x)[1 - F(x)]^{n-1}.$$

例 1.5 设总体 $X \sim U[0, \theta]$, $\theta > 0$, X_1, X_2, \dots, X_5 为 X 的样本, 分别求 $X_{(1)}, X_{(5)}$ 的密度函数 $\varphi_{X_{(1)}}(x), \varphi_{X_{(5)}}(x)$.

解 因为 $X \sim U[0, \theta]$, $\theta > 0$, 所以 X 的密度函数与分布函数分别为

$$\varphi(x) = \begin{cases} \frac{1}{\theta}, & x \in [0, \theta], \\ 0, & x \notin [0, \theta], \end{cases} \quad F(x) = \begin{cases} 0, & x \leq 0, \\ \frac{x}{\theta}, & 0 < x \leq \theta, \\ 1, & x > \theta, \end{cases}$$

$$\varphi_{X_{(5)}}(x) = 5F^4(x)\varphi(x) = \begin{cases} \frac{5x^4}{\theta^5}, & x \in [0, \theta], \\ 0, & x \notin [0, \theta], \end{cases}$$

$$\varphi_{X_{(1)}}(x) = 5\varphi(x)[1 - F(x)]^4 = \begin{cases} \frac{5}{\theta} \left(1 - \frac{x}{\theta}\right)^4, & x \in [0, \theta], \\ 0, & x \notin [0, \theta]. \end{cases}$$

1.2.3 经验分布函数与直方图

样本是总体的代表和反映, 总体 X 的分布函数 $F(x)$ 称为理论分布, 往往是未知的, 在实际工作中一般可用经验分布函数去推断总体的分布, 用直方图去描述(推断)总体 X (连续)的密度函数.

1. 经验分布函数

设 x_1, x_2, \dots, x_n 为来自总体 X 的样本观察值, 将这些值从小到大排序:

$$x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)},$$

对任意实数 x , 有

$$F_n(x) = \begin{cases} 0, & x < x_{(1)}, \\ \frac{k}{n}, & x_{(k)} \leq x < x_{(k+1)}, \quad k = 1, 2, \dots, n-1, \\ 1, & x \geq x_{(n)}, \end{cases}$$

则称 $F_n(x)$ 为总体 X 的经验分布函数. 经验分布函数具有下列性质:

- (1) $0 \leq F_n(x) \leq 1$;
- (2) $F_n(-\infty) = 0, F_n(+\infty) = 1$;
- (3) $F_n(x+0) = F_n(x)$ (右连续性).

经验分布函数 $F_n(x)$ 同样满足总体分布函数 $F(x)$ 的三个基本性质, 值得注意的是: 对于样本的不同观察值 x_1, x_2, \dots, x_n 得到的经验分布函数 $F_n(x)$ 是不同的, 在试验之前, 对固

定的 x 值, $F_n(x)$ 是一个随机变量, 当然也是一个统计量. 样本容量越大, 用经验分布函数 $F_n(x)$ 作为分布函数 $F(x)$ 的估计将会越精准. 图 1.1 给出了某正态总体 X 的理论分布函数 $F(x)$ 的曲线和经验分布函数 $F_n(x)$ 的曲线拟合情况, 从此图形可以看出, 经验分布函数可以近似代替总体的分布函数.

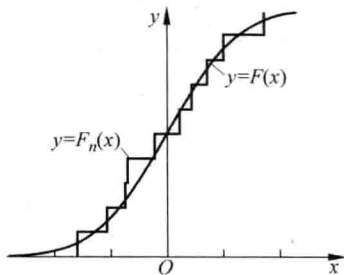


图 1.1 正态分布与经验分布拟合曲线

2. 直方图

直方图是用于近似连续型总体密度函数的曲线, 当样本容量 n 越大, 且分组比较细时, 近似程度也就越好.

假设 x_1, x_2, \dots, x_n 为连续型总体 X 的样本观察值. 构造直方图的步骤:

步骤 1 求出样本观察值 x_1, x_2, \dots, x_n 的极差 $x_{(n)} - x_{(1)}$.

步骤 2 确定组数与组距, 将包含 $x_{(1)}, x_{(n)}$ 的区间 $[a, b]$ 分成 m 个小区间: $[t_{i-1}, t_i)$ ($i=1, 2, \dots, m$), 其中 a 略小于 $x_{(1)}$, b 略大于 $x_{(n)}$. 一般组数由经验公式 $m \approx 1.87(n-1)^{0.4}$ 确定, 组距 $= \frac{b-a}{m}$. $t_i = t_{i-1} + \frac{b-a}{m}$, ($i=1, 2, \dots, m, t_0 = a, t_m = b$).

步骤 3 计算落入各区间样品个数, 记落入区间 $[t_{i-1}, t_i)$ 内的样品个数为 v_i , 称它为样本落入第 i 个区间的频数, 称 $f_i = \frac{v_i}{n}$ 为样本落入区间 $[t_{i-1}, t_i)$ 内的频率.

步骤 4 作图. 在 xOy 平面上, 以 x 轴上第 i 个小区间 $[t_{i-1}, t_i)$ 为底, 以 $y_i = \frac{f_i}{t_i - t_{i-1}}$ 为高作第 i 个长方形, 这样一排竖着的长方形所构成的图形就叫做直方图. 第 i 个长方形的面积为 f_i , 所有长方形面积之和为 1. 沿直方图边缘的曲线就是连续型总体的密度函数曲线的近似曲线.

例 1.6 某轧钢厂生产一批同型号的钢材, 为研究这批钢材的抗张力, 从中随机抽取了 76 个样品做张力实验, 测出数据见表 1.1.

表 1.1 钢材抗张力数据表

kg/cm²

41.0	37.0	33.0	44.2	30.5	27.0	45.0	28.5	31.2	33.5	38.5	41.5
42.0	45.5	42.5	39.0	38.8	35.5	32.5	29.6	32.6	34.5	37.5	39.5
42.8	45.1	42.8	45.8	39.8	37.2	33.8	31.2	29.0	35.2	37.8	41.2
43.8	48.0	43.6	41.8	36.6	34.8	31.0	32.0	33.5	37.4	40.8	44.7
40.2	41.3	38.8	34.1	31.8	34.6	38.3	41.3	30.0	35.2	37.5	40.5
38.1	37.3	37.1	41.5	29.5	29.1	27.5	34.8	36.5	44.2	40.0	44.5
40.6	36.2	35.8	31.5								

根据表中的数据, 作出直方图.

解 根据作直方图的步骤, 计算结果如表 1.2, 其图形如图 1.2 所示.

表 1.2 直方图计算表

分组区间	频数 ν_i	频率 f_i	纵坐标值 y_i
[27,30)	8	0.105	0.035
[30,33)	10	0.132	0.044
[33,36)	12	0.158	0.053
[36,39)	17	0.224	0.074
[39,42)	14	0.184	0.061
[42,45)	11	0.145	0.048
[45,48)	4	0.053	0.018

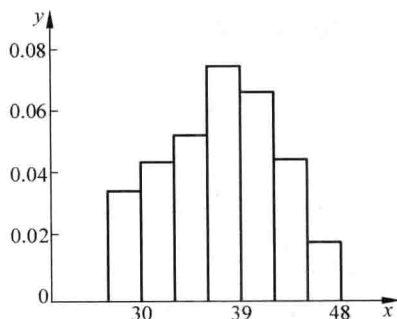


图 1.2 直方图

1.3 抽样分布及分位数

统计量是对总体分布和总体所含参数进行推断的基础,由于统计量是一个随机变量,称统计量的分布为抽样分布(sampling distribution).一般确定一个统计量的分布是十分复杂的,要用到许多概率知识.本节将讨论正态总体下一些常用的抽样分布.

1.3.1 正态分布的导出分布

1. χ^2 (卡方)分布

(1) χ^2 分布的定义

设 X_1, X_2, \dots, X_n 相互独立,且 $X_i \sim N(0,1) (i=1,2,\dots,n)$,则称 $\sum_{i=1}^n X_i^2$ 为服从自由度为

n 的卡方分布.记为 $\chi^2 = \sum_{i=1}^n X_i^2 \sim \chi^2(n)$.

(2) χ^2 分布的密度函数及其图像

$\chi^2(n)$ 分布的密度函数为

$$f(x, n) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, & x > 0, \\ 0, & x \leq 0. \end{cases}$$

密度函数图像在第一象限内是非负的,如图 1.3 所示.

(3) χ^2 分布的性质

① 若 $X \sim \chi^2(n)$, 则 $EX = n, DX = 2n$;

② (可加性) 若 $X \sim \chi^2(n_1), Y \sim \chi^2(n_2)$, 且 X, Y 相互独立, 则 $X+Y \sim \chi^2(n_1+n_2)$.

证明 ① 因为 $X = \sum_{i=1}^n X_i^2, X_i \sim N(0, 1) (i=1, 2, \dots, n)$, 故

$$EX_i = 0, DX_i = 1.$$

又因为

$$DX_i = EX_i^2 - (EX_i)^2,$$

故

$$EX_i^2 = DX_i + (EX_i)^2 = 1,$$

所以

$$EX = \sum_{i=1}^n EX_i^2 = n.$$

又因

$$\begin{aligned} EX_i^4 &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^4 e^{-\frac{x^2}{2}} dx = \frac{-1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^3 (-x e^{-\frac{x^2}{2}}) dx \\ &= \frac{-1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^3 dx e^{-\frac{x^2}{2}} = \frac{-1}{\sqrt{2\pi}} \left(x^3 e^{-\frac{x^2}{2}} \Big|_{-\infty}^{+\infty} - 3 \int_{-\infty}^{+\infty} x^2 e^{-\frac{x^2}{2}} dx \right) \\ &= 3 \int_{-\infty}^{+\infty} x^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 3EX_i^2 = 3, \end{aligned}$$

所以

$$DX_i^2 = EX_i^4 - (EX_i^2)^2 = 3 - 1 = 2, \quad \text{故} \quad DX = \sum_{i=1}^n DX_i^2 = \sum_{i=1}^n 2 = 2n.$$

② 由于 X 是 n_1 个独立的标准正态分布的平方和, Y 是 n_2 个独立的标准正态分布的平方和, 又因 X, Y 是相互独立的, 所以 $X+Y$ 是 n_1+n_2 个独立的标准正态分布的平方和, 根据定义有 $X+Y \sim \chi^2(n_1+n_2)$.

例 1.7 设 X_1, X_2, \dots, X_n 独立同分布于 $N(\mu, \sigma^2)$, 求 $E \sum_{i=1}^n (X_i - \mu)^2, D \sum_{i=1}^n (X_i - \mu)^2$.

解 因为 $\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 = \frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} \sim \chi^2(n)$, 故

$$E \frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} = n, D \frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2} = 2n,$$

于是有

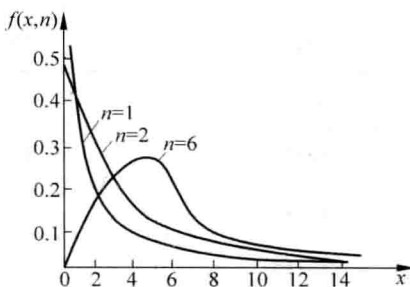


图 1.3 χ^2 分布的密度函数图像