

R

R

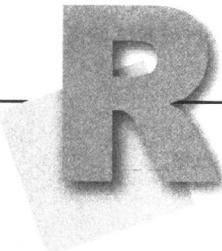
统计建模与R软件

薛 毅 陈立萍 编著

http://www.tup.com.cn

清华大学出版社





统计建模与R软件

薛 毅 陈立萍 编著

清华大学出版社
北京

内 容 简 介

本书以统计理论为基础,按照数理统计教材的章节顺序,在讲明统计的基本概念的同时,以 R 软件为辅助计算手段,介绍统计计算的方法,从而有效地解决统计中的计算问题.

书中结合数理统计问题对 R 软件进行科学、准确和全面的介绍,以便使读者能深刻理解该软件的精髓和灵活、高效的使用技巧.此外,还介绍了在工程技术、经济管理、社会生活等各方面的丰富的统计问题及其统计建模方法,通过该软件将所建模型进行求解,使读者获得从实际问题建模入手,到利用软件进行求解,以及对计算结果进行分析的全面训练.

本书可作为理工、经济、管理、生物等专业学生数理统计课程的辅导教材或教学参考书,也可作为统计计算课程的教材和数学建模竞赛的辅导教材.

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 13501256678 13801310933

图书在版编目 (CIP) 数据

统计建模与 R 软件/薛毅, 陈立萍编著. —北京: 清华大学出版社, 2007. 4

ISBN 978-7-302-14366-6

I. 统… II. ① 薛… ② 陈… III. 统计分析—应用软件 IV. C819

中国版本图书馆 CIP 数据核字(2006)第 158986 号

责任编辑: 刘 颖 赵从棉

责任校对: 赵丽敏

责任印制: 孟凡玉

出版发行: 清华大学出版社 地址: 北京清华大学学研大厦 A 座

<http://www.tup.com.cn> 邮 编: 100084

c-service@tup.tsinghua.edu.cn

社 总 机: 010-62770175 邮购热线: 010-62786544

投稿咨询: 010-62772015 客户服务: 010-62776969

印 刷 者: 北京市清华园胶印厂

装 订 者: 三河市源深装订厂

经 销: 全国新华书店

开 本: 185×230 印 张: 33.75 字 数: 716 千字

版 次: 2007 年 4 月第 1 版 印 次: 2007 年 4 月第 1 次印刷

印 数: 1~4000

定 价: 49.00 元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系
调换。联系电话: (010)62770177 转 3103 产品编号: 020971 -01

前　　言

本书力求将实用统计方法的介绍与在计算机上 R 软件如何实现这些方法紧密地联系起来, 不仅介绍了各种数理统计方法的统计思想、实际背景、统计模型和计算方法, 并且结合 R 软件, 给出相应解决问题的步骤并对计算结果进行分析.

关于数理统计的教材或教科书已非常多. 这类教材主要以数理统计的理论为基础, 讲述其理论、方法与应用背景, 但对于计算, 讲的较少, 基本是以手工计算为主, 目的是为了帮助读者理解相应的统计方法, 可操作性不强.

关于统计计算的书也有不少. 目前, 统计计算的教材一般是讲算法 (这一点与数值分析或计算方法差不多), 而没有相应的软件做支撑, 有些内容是数值分析内容的重复, 统计味不足.

结合软件讲统计的书, 目前最多的是结合 SAS 软件、SPSS 软件. 这类书籍基本上相当于软件使用说明书, 虽然谈到一些统计概念, 但讲得很少.

本书既不是一本单纯的关于数理统计或统计计算的教科书, 也不是一本关于 R 软件的使用手册, 而是一本将两者相结合的教科书. 其特点是结合 R 软件讲述数理统计的基本概念与计算方法.

R 软件是一种共享的统计软件, 也是一种数学计算环境. 它提供了有弹性的、互动的环境来分析和处理数据; 它提供了若干统计程序包, 以及一些集成的统计工具和各种数学计算、统计计算的函数. 用户只需根据统计模型, 指定相应的数据库及相关的参数, 便可灵活机动地进行数据分析等工作, 甚至创造出符合需要的新的统计计算方法. 使用 R 软件可以简化数据分析过程, 从数据的存取, 到计算结果的分享, R 软件提供了非常方便的计算工具, 帮助用户更好地分析和解决问题. 通过 R 软件的许多内嵌统计函数, 可以很容易学习和掌握 R 软件的语法, 也可以编制自己的函数来扩展现有的 R 语言, 完成科研工作.

本教材的编写风格是: (1) 以目前常见的数理统计教材的内容为基准, 首先对数理统计的基本概念、基本方法作一个简单、清晰的介绍, 在注重基础的同时, 侧重统计思想和统计方法的介绍. (2) 以 R 语言为主, 编写相应的计算程序. 这部分内容的目的有两个: 第一是学习 R 软件的编程方法, 掌握 R 软件的基本技巧; 第二是通过编程加深对统计方法的了解与掌握, 同时, 还可以通过编程, 加深对 R 软件中相关函数的了解. (3) 介绍相关的计算函数. 针对许多统计方法, R 软件提供了大量的内嵌计算函数, 使用者只需输入数据并且调用相应的内嵌函数, 就可得到相应的结果. 本书这一部分的写作重点放在对计算结果的统计解释上, 即如何通过结果来分析已有数据中所包含的统计信息, 着重介绍相应的统计建模方法. 这些是本教材最主要的特色, 也是不同于其他与软件有关的教材之处. 本

书着重强调统计建模, 以及如何使用 R 软件得到其计算结果和相应的结果解释.

本书的主要内容: 第 1 章, 概率统计的基本知识. 主要是复习统计的基本知识, 便于对后面各章内容的理解. 第 2 章, R 软件的使用. 主要介绍 R 软件的基本使用方法. 第 3 章, 数据描述性分析. 从数据描述开始分析数据, 主要介绍数据的基本特征, 如均值、方差, 还有与数据有关的各种图形, 如直方图、散点图等. 第 4 章, 参数估计. 介绍参数估计的基本方法, 如点估计和区间估计. 着重介绍 R 软件中与估计有关的函数. 第 5 章, 假设检验. 介绍假设检验的基本方法, 一类是参数检验, 另一类是非参数检验. 重点是介绍 R 软件中与假设检验有关的 R 函数及相关的使用方法. 第 6 章, 回归分析. 介绍回归分析的基本方法, 着重介绍回归分析的过程与方法和如何使用 R 软件作回归分析. 除一般的回归方法外, 还谈到逐步回归、非线性回归等内容. 第 7 章, 方差分析. 介绍单因素方差分析、双因素方差分析, 以及正交试验设计与方差分析之间的关系. 第 8 章, 应用多元分析 (I). 介绍判别分析和聚类分析, 这些内容与判别和分类有关. 第 9 章, 应用多元分析 (II). 介绍主成分分析、主因子分析和典型相关分析, 它是应用多元分析中降维计算的内容. 第 10 章, 计算机模拟. 介绍计算机模拟的 Monte Carlo 方法, 以及系统模拟方法, 最后介绍模拟方法在排队论中的应用. 此外, 还包括两个附录, 内容分别是作者自编函数的索引和 R 软件中函数的索引.

在学习本书的内容之后, 可以发现, 尽管有些统计内容其计算相当复杂, 但在使用 R 软件之后, 这些问题可以很轻松地得到解决.

本书所编写的 R 函数, 以及所介绍的 R 函数均以 R-2.1.1 版为基础 (目前的版本是 R-2.3.1, 大约每 3 至 4 个月版本会更新一次), 而且全部程序均运行通过, 读者如果需要作者自编的 R 程序, 可以通过电子邮件向作者索取, 邮件地址: xueyi@bjut.edu.cn.

本书是为理工、经济、管理、生物等专业学生或专业人员解决统计计算问题编写的, 可以作为上述专业学生数理统计课程的辅导教材或教学参考书, 也可作为统计计算课程的教材和数学建模竞赛的辅导教材.

由于编者水平有限, 书中一定存在不足甚至错误之处, 欢迎读者不吝指正, 我们的电子邮件地址是: xueyi@bjut.edu.cn (薛毅); chenliping@bjut.edu.cn (陈立萍).

编　者

2006 年 7 月

于北京工业大学

目 录

第 1 章 概率统计的基本知识	1
1.1 随机事件与概率	1
1.1.1 随机事件	1
1.1.2 概率	3
1.1.3 古典概型	4
1.1.4 几何概型	5
1.1.5 条件概率	6
1.1.6 概率的乘法公式、全概率公式、Bayes 公式	7
1.1.7 独立事件	8
1.1.8 n 重 Bernoulli 试验及其概率计算	9
1.2 随机变量及其分布	9
1.2.1 随机变量的定义	9
1.2.2 随机变量的分布函数	10
1.2.3 离散型随机变量	10
1.2.4 连续型随机变量	12
1.2.5 随机向量	16
1.3 随机变量的数字特征	21
1.3.1 数学期望	21
1.3.2 方差	22
1.3.3 几种常用随机变量分布的期望与方差	22
1.3.4 协方差与相关系数	23
1.3.5 矩与协方差矩阵	24
1.4 极限定理	26
1.4.1 大数定律	26
1.4.2 中心极限定理	28
1.5 数理统计的基本概念	29
1.5.1 总体、个体、简单随机样本	29
1.5.2 参数空间与分布族	31

1.5.3 统计量和抽样分布.....	31
1.5.4 正态总体样本均值与样本方差的分布.....	37
习题.....	38
第2章 R 软件的使用.....	41
2.1 R 软件简介	41
2.1.1 R 软件的下载与安装	42
2.1.2 初识 R 软件.....	42
2.1.3 R 软件主窗口命令与快捷方式	48
2.2 数字、字符与向量	58
2.2.1 向量	58
2.2.2 产生有规律的序列.....	60
2.2.3 逻辑向量	61
2.2.4 缺失数据	62
2.2.5 字符型向量	63
2.2.6 复数向量	63
2.2.7 向量下标运算	64
2.3 对象和它的模式与属性	67
2.3.1 固有属性: mode 和 length	67
2.3.2 修改对象的长度	68
2.3.3 attributes()和 attr()函数	68
2.3.4 对象的 class 属性	69
2.4 因子	69
2.4.1 factor()函数	70
2.4.2 tapply()函数	70
2.4.3 gl()函数	71
2.5 多维数组和矩阵	71
2.5.1 生成数组或矩阵	71
2.5.2 数组下标	73
2.5.3 数组的四则运算	74
2.5.4 矩阵的运算	76
2.5.5 与矩阵(数组)运算有关的函数	82

2.6 列表与数据框.....	84
2.6.1 列表	84
2.6.2 数据框.....	86
2.6.3 列表与数据框的编辑	89
2.7 读、写数据文件.....	89
2.7.1 读纯文本文件.....	89
2.7.2 读其他格式的数据文件.....	92
2.7.3 链接嵌入的数据库.....	94
2.7.4 写数据文件	95
2.8 控制流	96
2.8.1 分支语句	96
2.8.2 中止语句与空语句.....	97
2.8.3 循环语句	97
2.9 编写自己的函数.....	99
2.9.1 简单的例子	99
2.9.2 定义新的二元运算.....	102
2.9.3 有名参数与默认参数	102
2.9.4 递归函数	104
习题	105
第 3 章 数据描述性分析.....	107
3.1 描述统计量	107
3.1.1 位置的度量	107
3.1.2 分散程度的度量	112
3.1.3 分布形状的度量	114
3.2 数据的分布	116
3.2.1 分布函数	116
3.2.2 直方图、经验分布图与 QQ 图	118
3.2.3 茎叶图、箱线图及五数总括	123
3.2.4 正态性检验与分布拟合检验	128
3.3 R 软件中的绘图命令	130
3.3.1 高水平作图函数	130

3.3.2 高水平绘图中的命令	137
3.3.3 低水平作图函数	138
3.4 多元数据的数据特征与相关分析	140
3.4.1 二元数据的数字特征及相关系数	140
3.4.2 二元数据的相关性检验	142
3.4.3 多元数据的数字特征及相关矩阵	145
3.4.4 基于相关系数的变量分类	148
3.5 多元数据的图形表示方法	154
3.5.1 轮廓图	154
3.5.2 星图	155
3.5.3 调和曲线图	158
习题	160
第 4 章 参数估计	162
4.1 点估计	162
4.1.1 矩法	163
4.1.2 极大似然法	166
4.2 估计量的优良性准则	174
4.2.1 无偏估计	174
4.2.2 有效性	176
4.2.3 相合性(一致性)	177
4.3 区间估计	177
4.3.1 一个正态总体的情况	178
4.3.2 两个正态总体的情况	182
4.3.3 非正态总体的区间估计	190
4.3.4 单侧置信区间估计	191
习题	200
第 5 章 假设检验	203
5.1 假设检验的基本概念	203
5.1.1 基本概念	203
5.1.2 假设检验的基本思想与步骤	205
5.1.3 假设检验的两类错误	205

5.2 重要的参数检验	206
5.2.1 正态总体均值的假设检验	206
5.2.2 正态总体方差的假设检验	215
5.2.3 二项分布总体的假设检验	220
5.3 若干重要的非参数检验	222
5.3.1 Pearson 拟合优度 χ^2 检验	222
5.3.2 Kolmogorov-Smirnov 检验	227
5.3.3 列联表数据的独立性检验	229
5.3.4 符号检验	235
5.3.5 秩统计量	239
5.3.6 秩相关检验	240
5.3.7 Wilcoxon 秩检验	243
习题	249
第 6 章 回归分析	253
6.1 一元线性回归	253
6.1.1 数学模型	253
6.1.2 回归参数的估计	255
6.1.3 回归方程的显著性检验	256
6.1.4 参数 β_0 与 β_1 的区间估计	258
6.1.5 预测	259
6.1.6 控制	260
6.1.7 计算实例	260
6.2 R 软件中与线性模型有关的函数	265
6.2.1 基本函数	265
6.2.2 提取模型信息的通用函数	266
6.3 多元线性回归分析	267
6.3.1 数学模型	267
6.3.2 回归系数的估计	268
6.3.3 显著性检验	269
6.3.4 参数 β 的区间估计	271
6.3.5 预测	272

6.3.6 修正拟合模型	272
6.3.7 计算实例	273
6.4 逐步回归	279
6.4.1 “最优”回归方程的选择	279
6.4.2 逐步回归的计算	279
6.5 回归诊断	284
6.5.1 什么是回归诊断	284
6.5.2 残差	288
6.5.3 残差图	291
6.5.4 影响分析	296
6.5.5 多重共线性	304
6.6 广义线性回归模型	307
6.6.1 与广义线性模型有关的 R 函数	307
6.6.2 正态分布族	308
6.6.3 二项分布族	309
6.6.4 其他分布族	316
6.7 非线性回归模型	318
6.7.1 多项式回归模型	319
6.7.2 (内在)非线性回归模型	323
习题	331
第 7 章 方差分析	336
7.1 单因素方差分析	336
7.1.1 数学模型	337
7.1.2 方差分析	338
7.1.3 方差分析表的计算	339
7.1.4 均值的多重比较	342
7.1.5 方差的齐次性检验	345
7.1.6 Kruskal-Wallis 秩和检验	348
7.1.7 Friedman 秩和检验	351
7.2 双因素方差分析	353
7.2.1 不考虑交互作用	353

7.2.2 考虑交互作用	356
7.2.3 方差齐性检验	359
7.3 正交试验设计与方差分析	361
7.3.1 用正交表安排试验	361
7.3.2 正交试验的方差分析	364
7.3.3 有交互作用的试验	366
7.3.4 有重复试验的方差分析	369
习题	371
第 8 章 应用多元分析(I)	375
8.1 判别分析	375
8.1.1 距离判别	375
8.1.2 Bayes 判别	385
8.1.3 Fisher 判别	393
8.2 聚类分析	397
8.2.1 距离和相似系数	397
8.2.2 系统聚类法	403
8.2.3 动态聚类法	418
习题	420
第 9 章 应用多元分析(II)	423
9.1 主成分分析	423
9.1.1 总体主成分	423
9.1.2 样本主成分	427
9.1.3 相关的 R 函数以及实例	429
9.1.4 主成分分析的应用	434
9.2 因子分析	441
9.2.1 引例	442
9.2.2 因子模型	443
9.2.3 参数估计	445
9.2.4 方差最大的正交旋转	455
9.2.5 因子分析的计算函数	458
9.2.6 因子得分	461

9.3 典型相关分析.....	463
9.3.1 总体典型相关.....	464
9.3.2 样本典型相关.....	466
9.3.3 典型相关分析的计算.....	467
9.3.4 典型相关系数的显著性检验	471
习题	473
第 10 章 计算机模拟	476
10.1 概率分析与 Monte Carlo 方法	476
10.1.1 概率分析	476
10.1.2 Monte Carlo 方法	477
10.1.3 Monte Carlo 方法的精度分析.....	480
10.2 随机数的产生.....	485
10.2.1 均匀分布随机数的产生.....	485
10.2.2 均匀随机数的检验.....	486
10.2.3 任意分布随机数的产生.....	488
10.2.4 正态分布随机数的产生.....	489
10.2.5 用 R 软件生成随机数.....	490
10.3 系统模拟	490
10.3.1 连续系统模拟	491
10.3.2 离散系统模拟	492
10.4 模拟方法在排队论中的应用	497
10.4.1 排队服务系统的基本概念	497
10.4.2 排队模型模拟的关键	500
10.4.3 等待制排队模型的模拟	501
10.4.4 损失制与混合制排队模型	507
习题	513
附录 索引	515
附录 A 自编写的函数(程序)	515
附录 B R 软件中的函数(程序)	517
参考文献	526

第1章 概率统计的基本知识

本书是一本统计建模与软件应用相结合的教科书, 讲述重点是数理统计的基本方法和用 R 软件进行相应的计算. 众所周知, 数理统计是以概率论为基础, 应用非常广泛的数学学科分支, 是通过对试验或观察数据进行分析, 来研究随机现象以达到对研究对象的客观规律性做出合理的估计和推断的目的, 因此在介绍统计建模和 R 软件知识之前, 有必要先回顾一下相关的概率与数理统计的基本概念, 以及数理统计的各个应用分支.

本章用 4 节内容简单回顾概率论的基础知识, 用 1 节内容简单介绍数理统计的基本概念. 这样做的目的是使读者对已有概率论的知识有一个全面的了解与回顾, 对数理统计的概念有一个基本的认识.

1.1 随机事件与概率

1.1.1 随机事件

1. 随机事件

在一定条件下, 所得的结果不能预先完全确定, 而只能确定是多种可能结果中的一种, 称这种现象为随机现象. 例如, 抛掷一枚硬币, 其结果有可能是出现正面, 也有可能是出现反面; 电话交换台在 1 min 内接到的呼叫次数, 可能是 0 次、1 次、2 次、……; 在同一工艺条件下生产出的灯泡, 其使用寿命有长有短; 测量同一物体的长度时, 由于仪器及观察受到环境的影响, 多次测量的结果往往有差异, 等等. 这些现象都是随机现象.

使随机现象得以实现和对它观察的全过程称为随机试验 (random experiment), 记为 E . 随机试验满足以下条件:

- (1) 可以在相同条件下重复进行;
- (2) 结果有多种可能性, 并且所有可能结果事先已知;
- (3) 做一次试验究竟哪个结果出现, 事先不能确定.

称随机试验的所有可能结果组成的集合为样本空间 (sample space), 记为 Ω . 试验的每一个可能结果称为样本点 (sample point), 记为 ω .

称 Ω 中满足一定条件的子集为随机事件 (random event), 用大写字母 A, B, C, \dots 表示.

若一个随机事件只含一个不可再分的试验结果称为一个基本事件 (即一个样本点所组成的集合 $\{\omega\}$).

在试验中, 称一个事件发生是指构成该事件的一个样本点出现. 由于样本空间 Ω 包含了所有的样本点, 所以在每次试验中, 它总是发生, 因此称 Ω 为必然事件 (certain event). 空集 \emptyset 不包含任何样本点, 且在每次试验中总不发生, 所以称 \emptyset 为不可能事件 (impossible event).

2. 随机事件之间的关系

若事件 A 的发生必然导致事件 B 的发生, 则称事件 A 包含于事件 B , 或事件 B 包含事件 A , 记为 $A \subset B$, 亦称为事件的包含 (contain) 关系.

若 $A \subset B$, 且 $B \subset A$, 则称事件 A 与事件 B 等价 (equivalent), 记为 $A = B$.

若事件 A 与事件 B 至少有一个发生, 则称为事件的和 (union), 记为 $A \cup B$. 若 n 个事件 A_1, A_2, \dots, A_n 中至少有一个发生, 则称为 n 个事件的和, 记为 $A_1 \cup A_2 \cup \dots \cup A_n$ 或 $\bigcup_{i=1}^n A_i$.

同样, 可以定义可列无穷个事件的和 $A_1 \cup A_2 \cup \dots \cup A_n \cup \dots$ 或 $\bigcup_{i=1}^{\infty} A_i$, 表示无穷个事件中至少有一个发生.

若事件 A 发生而事件 B 不发生, 则称为事件 A 与事件 B 的差, 记为 $A - B$.

若事件 A 与 B 同时发生, 则称事件 A 与事件 B 的积 (intersection), 记为 $A \cap B$ 或 AB . 若 n 个事件 A_1, A_2, \dots, A_n 同时发生, 则称为 n 个事件的积, 记为 $A_1 \cap A_2 \cap \dots \cap A_n$ 或 $\bigcap_{i=1}^n A_i$.

同样, 可以定义可列无穷个事件的积 $A_1 \cap A_2 \cap \dots \cap A_n \cap \dots$ 或 $\bigcap_{i=1}^{\infty} A_i$, 表示无穷个事件同时发生.

若事件 A 与 B 不能同时发生, 则称事件 A 与事件 B 为互斥事件 (mutually exclusive event) 或不相容事件 (incompatible event), 记为 $AB = \emptyset$.

在一次试验中, 基本事件之间是两两互斥的.

若 A 为随机事件, 称“事件 A 不发生”的事件为事件 A 的对立事件 (opposite event) 或逆事件 (complementary event), 记为 \bar{A} . 事件与对其立事件有如下关系:

$$A \cup \bar{A} = \Omega, \quad A\bar{A} = \emptyset.$$

由定义可知: 对立事件一定是互斥事件, 但互斥事件不一定是对立事件.

3. 随机事件的运算律

(1) 交换律

$$A \cup B = B \cup A, \quad AB = BA. \tag{1.1}$$

(2) 结合律

$$(A \cup B) \cup C = A \cup (B \cup C), \quad (A \cap B) \cap C = A \cap (B \cap C). \tag{1.2}$$

(3) 分配律

$$(A \cup B)C = (AC) \cup (BC), \quad A \cup (BC) = (A \cup B)(A \cup C). \quad (1.3)$$

(4) De Morgan(德·摩根) 律

$$\overline{A_1 \cup A_2} = \overline{A_1} \cap \overline{A_2}, \quad \overline{A_1 \cap A_2} = \overline{A_1} \cup \overline{A_2}. \quad (1.4)$$

对于 n 个或可列无穷个事件有

$$\overline{\bigcup_{k=1}^n A_k} = \bigcap_{k=1}^n \overline{A_k}, \quad \overline{\bigcap_{k=1}^n A_k} = \bigcup_{k=1}^n \overline{A_k}, \quad \overline{\bigcup_{k=1}^{\infty} A_k} = \bigcap_{k=1}^{\infty} \overline{A_k}, \quad \overline{\bigcap_{k=1}^{\infty} A_k} = \bigcup_{k=1}^{\infty} \overline{A_k}. \quad (1.5)$$

(5) 减法满足

$$A - B = A\overline{B} \quad \text{或} \quad A - B = A \cap \overline{B}. \quad (1.6)$$

1.1.2 概率

1. 概率的公理化定义

在概率论中并非样本空间 Ω 的任何子集均可以看作事件, 所定义的事件之间应满足一定的代数结构.

定义 1.1 设随机试验 E 的样本空间为 Ω , \mathcal{F} 是 Ω 的子集组成的集族, 满足

- (1) $\Omega \in \mathcal{F}$;
- (2) 若 $A \in \mathcal{F}$, 则 $\overline{A} \in \mathcal{F}$; (对逆运算封闭)
- (3) 若 $A_i \in \mathcal{F}$, $i = 1, 2, \dots$, 则 $\bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$. (对可列并运算封闭)

则称 \mathcal{F} 为 Ω 的一个 σ 代数 (事件体), \mathcal{F} 中的集合称为事件. 样本空间 Ω 和 σ 代数的二元体 (Ω, \mathcal{F}) 称为可测空间.

定义 1.2 随机试验 E 的样本空间为 Ω , (Ω, \mathcal{F}) 是可测空间, 对于每个事件 $A \in \mathcal{F}$, 定义一个实数 $P(A)$ 与之对应, 若函数 $P(\cdot)$ 满足条件:

- (1) 对每个事件 A , 均有 $0 \leq P(A) \leq 1$;
- (2) $P(\Omega) = 1$;
- (3) 若事件 A_1, A_2, \dots 两两互斥, 即对于 $i, j = 1, 2, \dots, i \neq j$, $A_i A_j = \emptyset$ 均有

$$P(A_1 \cup A_2 \cup \dots) = P(A_1) + P(A_2) + \dots.$$

则称 $P(A)$ 为事件 A 的概率 (probability), 称 (Ω, \mathcal{F}, P) 为概率空间.

2. 概率的性质

性质 1 $P(\emptyset) = 0$, 即不可能事件的概率为零.

但性质 1 反过来不成立, 即 $P(A) = 0 \not\Rightarrow A = \emptyset$.

性质 2 若事件 A_1, A_2, \dots, A_n 两两互斥, 则有

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = P(A_1) + P(A_2) + \dots + P(A_n), \quad (1.7)$$

即互斥事件和的概率等于它们各自概率的和.

性质 3 对任一事件 A , 均有 $P(\bar{A}) = 1 - P(A)$.

性质 4 对两个事件 A 和 B , 若 $A \subset B$, 则有

$$P(B - A) = P(B) - P(A), \quad P(B) \geq P(A). \quad (1.8)$$

性质 5(加法公式) 对任意两个事件 A 和 B , 有

$$P(A \cup B) = P(A) + P(B) - P(AB). \quad (1.9)$$

性质 5 可以推广为

$$\begin{aligned} P(A_1 \cup A_2 \cup A_3) &= P(A_1) + P(A_2) + P(A_3) - P(A_1 A_2) - P(A_1 A_3) \\ &\quad - P(A_2 A_3) + P(A_1 A_2 A_3), \end{aligned} \quad (1.10)$$

$$P(A_1 \cup A_2 \cup \dots \cup A_n) = S_1 - S_2 + S_3 - S_4 + \dots + (-1)^{n-1} S_n, \quad (1.11)$$

其中 $S_1 = \sum_{i=1}^n P(A_i)$, $S_2 = \sum_{1 \leq i < j \leq n} P(A_i A_j)$, $S_3 = \sum_{1 \leq i < j < k \leq n} P(A_i A_j A_k)$, \dots , $S_n = P(A_1 A_2 \dots A_n)$.

1.1.3 古典概型

设随机事件 E 的样本空间中只有有限个样本点, 即 $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$, 其中 n 为样本点总数. 每个样本点 $\omega_i (i = 1, 2, \dots, n)$ 出现是等可能的, 并且每次试验有且仅有一个样本点发生, 则称这类现象为古典概型 (classical probability model). 若事件 A 包含 m 个样本点, 则事件 A 的概率定义为

$$P(A) = \frac{m}{n} = \frac{\text{事件 } A \text{ 包含的基本事件数}}{\text{基本事件总数}}. \quad (1.12)$$

例 1.1 设有 k 个不同的 (可分辨) 球, 每个球都能以同样的概率 $1/l$ 落到 l 个格子 ($l \geq k$) 的每一个中, 且每个格子可容纳任意多个球, 试分别求如下两事件 A 与 B 的概率.

A : 指定的 k 个格子中各有一个球;

B : 存在 k 个格子, 其中各有一个球.

解 由于每个球可以落入 l 个格子中的任一个, 并且每一个格子中可落入任意多个球, 所以 k 个球落入 l 个格子中的分布情况相当于从 l 个格子中选取 k 个的可重复排列, 故样本空间共有 l^k 种等可能的基本结果.