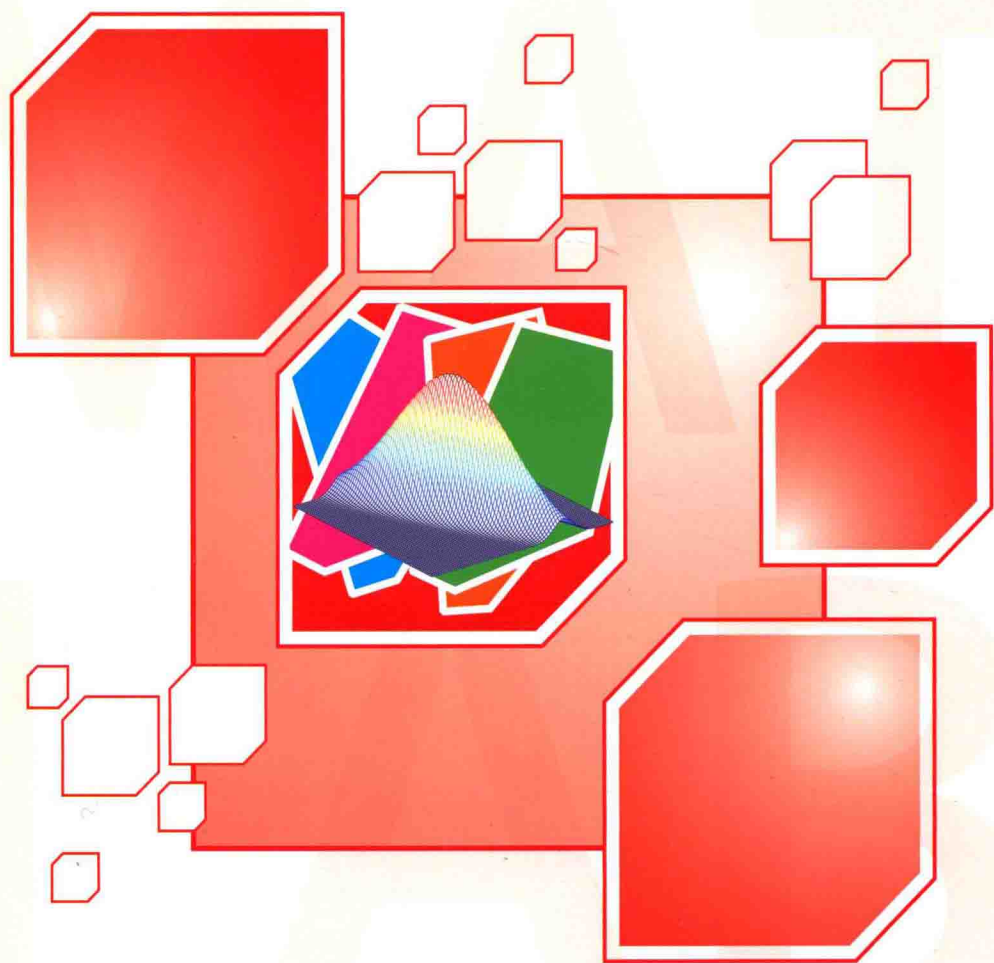


实用数据分析 与MATLAB软件

薛毅 陈立萍◎编著



北京工业大学出版社

本书以MATLAB软件为工具，结合大量的工程实例，深入浅出地介绍了MATLAB软件在工程中的应用。全书共分10章，主要内容包括：MATLAB的基本操作、MATLAB的图形用户界面设计、MATLAB的符号运算、MATLAB的数值计算、MATLAB的字符串处理、MATLAB的数据库接口、MATLAB的并行计算、MATLAB的Web应用、MATLAB的部署与发布、MATLAB的国际化。

实用数据分析与 MATLAB 软件

薛毅 陈立萍 编著

北京工业大学出版社

内 容 简 介

本书以数据处理的过程为基准,介绍相关的基本概念、基本方法,在注重统计思想的同时,侧重于将统计方法与计算机软件相结合。淡化计算、注重方法是本书的宗旨。因此,书中数据处理的计算工作将由 MATLAB 软件来完成。本书着重介绍了 MATLAB 统计工具箱中与数据分析相关的函数。有了这些函数,读者就可以轻松地完成过去认为不可能完成的工作,体会到使用 MATLAB 软件的乐趣。

本书是为理工、经济、管理、生物等专业学生或专业人员编写的,目的是为了解决数据分析中可能遇到的问题。本书可作为上述专业学生相关课程的辅导教材或教学参考书,也可作为信息与计算科学专业学生数据分析或统计计算课程的教材和概率论、数理统计、多元分析、回归分析与实验设计等课程的辅助教材,还可作为数学建模竞赛的辅导教材,以及作为科技工作者和工程技术人员学习和参考用书。

图书在版编目(CIP)数据

实用数据分析与 MATLAB 软件/薛毅,陈立萍编著. —北京:北京工业大学出版社, 2015.7

ISBN 978 - 7 - 5639 - 4375 - 3

I. ①实… II. ①薛… ②陈… III. ①Matlab 软件 - 应用 - 统计分析 IV. ①C819

中国版本图书馆 CIP 数据核字(2015)第 144064 号

实用数据分析与 MATLAB 软件

编 著:薛 毅 陈立萍

责任编辑:李周辉 贺 帆

封面设计:何 强

出版发行:北京工业大学出版社

(北京市朝阳区平乐园 100 号 邮编:100124)

010 - 67391722(传真)bgdcbs@sina.com

出 版 人:郝 勇

经销单位:全国各地新华书店

承印单位:徐水宏远印刷有限公司

开 本:787 毫米×1092 毫米 1/16

印 张:32.5

字 数:687 千字

版 次:2015 年 8 月第 1 版

印 次:2015 年 8 月第 1 次印刷

标准书号:ISBN 978 - 7 - 5639 - 4375 - 3

定 价:58.00 元

版权所有 翻印必究

(如发现印装质量问题,请寄本社发行部调换 010 - 67391106)

前 言

当今的时代是信息爆炸的时代，数据无处不在，从数据中得到有用的信息是分析和解决问题的重要环节。因此，如何从大量的、看似杂乱无章的数据中揭示隐藏在内部的规律、挖掘出有用的信息，以指导人们进行科学的推断和决策，这就需要用到数据分析。

顾名思义，数据分析就是分析和处理数据的方法，从这个层面上讲，数据分析不存在固定的模式与解决方法，一切与分析 and 处理数据有关的方法都包含在其中。但由于数据带有一定的随机性，因此狭义地讲，数据分析就是用数理统计的方法来分析和处理数据。本书正是在这个层面上开展工作的。

本书之所以命名为《实用数据分析与 MATLAB 软件》，其目的是使用概率或数理统计的方法来处理数据，借用 MATLAB 软件强大的计算功能，完成各种统计方法的计算，从而将研究工作从繁杂的计算中解放出来，将研究问题的重点放在对数据分析和处理的方法上。

MATLAB 是美国 MathWork 公司于 1984 年开发的主要用于工程计算的计算机高级语言，是当今国际上优秀的数学软件之一。MATLAB 以矩阵形式处理数据，具有强大的科学计算、图形处理、可视化、开放式和可扩展的功能，MATLAB 及其附带的几十种面对不同领域的工具箱(ToolBox)，能够广泛应用于数值分析、信号与图像处理、控制系统设计、通信仿真、工程优化、数学建模和统计分析等各个领域。

从教学层面来讲，概率论、数理统计、多元分析等统计课程，都有大量实例可供学生完成，但复杂的计算过程，会削弱学生的学习兴趣。使用 MATLAB 之后，计算工作交给计算机完成，学生学习的重点可放在问题的分析与理解方面。本书着重介绍了 MATLAB 数据统计工具箱中与数据分析相关的函数。有了这些函数，读者就可以轻松地完成过去认为不可能完成的工作，体会到学习的乐趣，通过对实例的计算，加深对所学理论知识的理解。

本书的主要内容共分 8 章来介绍。第 1 章，概率论与数理统计基础，简单回顾一下相关的概率论与数理统计的基本概念。第 2 章，数据描述性分析，介绍数据位置和分散程度的度量，数据的分布，以及数据的图形表示。第 3 章，参数估计，介绍分布参数的估计方法，第 4 章，假设检验，介绍常用的参数检验和非参数检验方法。第 5 章，回归分析。介绍回归分析和回归诊断的内容。第 6 章，方差分析，介绍单、双和

多因素分析方法。第 7 章,应用多元分析(I),介绍与分类有关的判别分析和聚类分析。第 8 章,应用多元分析(II),介绍主成分分析、因子分析、典型相关分析和非负矩阵分解的知识。本书的基本结构是:在介绍统计知识之后,重点介绍如何使用 MATLAB 数据统计工具箱中的函数完成统计方法的计算,同时介绍如何利用 MATLAB 提供的函数,编写一些有效地便于求解问题的程序。为帮助读者对文中 MATLAB 函数或程序的理解,本书给出了两个附录:附录 A, MATLAB 软件简介;附录 B, MATLAB 数理统计工具箱。

本书既不是纯粹的概率统计教材,也不是 MATLAB 软件的使用说明书,而是以数据处理的过程为基准,介绍相关的基本概念、基本方法,在注重统计思想的同时,侧重于将统计方法与计算机软件相结合,淡化计算、注重方法是本书的宗旨。因此,数据处理的计算工作将由 MATLAB 软件来完成。具体地说,绝大多数的计算都由 MATLAB 数据统计工具箱中的函数完成。在学习这些函数之后,可以达到以下三个目的:第一,学会使用 MATLAB 的相关函数求解问题,在科研与实际工作中会选择合适的函数进行计算。第二,读懂 MATLAB 的计算结果,会运用计算结果来解释所要处理的问题。第三,将 MATLAB 的函数作为基础函数,编写自己需要的函数程序,能够“站在巨人的肩膀上”工作。

本书的全部程序均通过计算检验,书中的程序已在 MATLAB 7.9.0(R2009b)环境下运行通过。读者所持的软件的版本可能与编者的不一致,这基本不会影响到书中的程序的运行,但书中的某些函数不能在较低版本环境下运行,请读者见谅。

本书是为理工、经济、管理、生物等专业学生或专业人员编写的,目的是为了解决数据分析中可能遇到的问题。本书可作为上述专业学生相关课程的辅导教材或教学参考书,也可作为信息与计算科学专业学生数据分析或统计计算课程的教材和概率论、数理统计、多元分析、回归分析与实验设计等课程的辅助教材,还可作为数学建模竞赛的辅导教材,以及作为科技工作者和工程技术人员学习和参考用书。

由于受编者水平限制,可能在内容的取舍、结构的编排及课程的讲法上存在着不妥之处,希望使用本书的教师、学生和同行专家及其他读者提出宝贵的批评和建议。

在本书出版之际,谨向对本书提供过帮助的各位教师和专家表示感谢,感谢北京工业大学研究生院对数学建模课程的支持,同时感谢北京工业大学出版社为本书的出版做的大量工作。

目 录

第 1 章 概率论与数理统计基础	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(贝叶斯)公式	7
1.1.7 独立事件	7
1.1.8 n 重伯努利试验及其概率计算	8
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	9
1.2.3 离散型随机变量	9
1.2.4 连续型随机变量	11
1.2.5 分位数	13
1.2.6 MATLAB 软件中的分布函数	13
1.2.7 随机向量	16
1.3 随机变量的数字特征	23
1.3.1 数学期望	23
1.3.2 方差	23
1.3.3 几种常用随机变量的期望与方差	24
1.3.4 协方差与相关系数	24
1.3.5 矩与协方差矩阵	25
1.4 极限定理	27
1.4.1 大数定律	28
1.4.2 中心极限定理	29

1.5 数理统计的基本概念	32
1.5.1 总体、个体、简单随机样本	32
1.5.2 参数空间与分布族	34
1.5.3 统计量	34
1.5.4 常用的分布	36
1.5.5 MATLAB 统计工具箱中的分布函数	41
习题 1	44
第 2 章 数据描述性分析	47
2.1 描述统计量	47
2.1.1 数据的分类	47
2.1.2 位置的度量	48
2.1.3 分散程度的度量	58
2.1.4 分布形状的度量	60
2.2 数据的分布	63
2.2.1 直方图与核密度估计	63
2.2.2 经验分布图与 QQ 图	67
2.2.3 箱线图	73
2.3 多元数据的数字特征与相关分析	77
2.3.1 协方差矩阵	78
2.3.2 相关系数与相关性检验	80
2.4 多元数据的图表示方法	84
2.4.1 轮廓图	85
2.4.2 星图	87
2.4.3 调和曲线图	89
习题 2	91
第 3 章 参数估计	93
3.1 点估计	93
3.1.1 总体矩、样本矩、矩法	94
3.1.2 用 MATLAB 作矩估计	96
3.1.3 极大似然法	98
3.1.4 MATLAB 中作极大似然估计的函数	102
3.2 区间估计	106
3.2.1 一个正态总体的情况	107
3.2.2 一个正态总体情况的 MATLAB 计算	109

3.2.3	两个正态总体的情况	113
3.2.4	两个正态总体情况的 MATLAB 计算	115
3.2.5	非正态总体的区间估计	118
	习题 3	120
第 4 章	假设检验	123
4.1	假设检验的基本概念	123
4.1.1	基本概念	123
4.1.2	假设检验的基本思想与步骤	125
4.1.3	假设检验的两类错误	125
4.2	重要的参数检验	126
4.2.1	单个正态总体均值的检验	126
4.2.2	两个正态总体均值差的检验	130
4.2.3	单个正态总体方差的假设检验	135
4.2.4	两个正态总体方差的假设检验	137
4.2.5	非正态总体参数的假设检验	139
4.3	分布检验	141
4.3.1	Pearson(皮尔森)拟合优度 χ^2 检验	142
4.3.2	Kolmogorov-Smirnov(科尔莫戈罗夫-斯米尔诺夫)检验	147
4.3.3	Jarque-Bera(雅克-贝拉)正态性检验	151
4.3.4	Lilliefors 检验	152
4.4	随机性检验和独立性检验	154
4.4.1	游程检验	154
4.4.2	列联表数据的独立性检验	158
4.4.3	Fisher(费希尔)精确独立检验	162
4.5	符号检验	165
4.5.1	单个总体样本中位数检验	165
4.5.2	成对数据的符号检验	167
4.6	秩检验	169
4.6.1	秩统计量	169
4.6.2	单个总体样本的符号秩检验	171
4.6.3	成对数据的符号秩检验	173
4.6.4	秩和检验	174
4.7	相关性检验	178
4.7.1	Pearson(皮尔森)相关检验	178

4.7.2	Spearman(斯皮尔曼)相关检验	179
4.7.3	Kendall(肯达尔)相关检验	180
4.7.4	MATLAB 函数作相关检验的计算	180
	习题 4	183
第 5 章 回归分析		189
5.1	一元线性回归	189
5.1.1	数学模型	190
5.1.2	回归参数的估计	191
5.1.3	参数 β_0 和 β_1 的区间估计	193
5.1.4	回归方程的显著性检验	194
5.1.5	预测	198
5.2	多元线性回归分析	200
5.2.1	数学模型	200
5.2.2	回归系数的估计	201
5.2.3	参数 β 的区间估计	201
5.2.4	显著性检验	203
5.2.5	预测	205
5.2.6	计算实例	206
5.3	逐步回归	211
5.3.1	“最优”回归方程的选择	211
5.3.2	逐步回归的计算	211
5.4	回归诊断	216
5.4.1	什么是回归诊断	216
5.4.2	残差	220
5.4.3	残差图	225
5.4.4	影响分析	231
5.4.5	多重共线性	236
5.4.6	岭估计	238
5.5	稳健回归	242
5.5.1	稳健回归的基本概念	242
5.5.2	极大似然型稳健回归——M 估计	244
5.5.3	用 MATLAB 内置函数计算稳健回归	245
5.6	非线性回归模型	247
5.6.1	多项式回归模型	248

5.6.2 (内在)非线性回归模型	250
5.7 广义线性回归模型	256
5.7.1 相关的 MATLAB 函数	257
5.7.2 Logistic 回归模型	259
5.7.3 其他分布族	264
习题 5	268
第 6 章 方差分析	275
6.1 单因素方差分析	275
6.1.1 数学模型	276
6.1.2 方差分析	277
6.1.3 方差分析表的计算	278
6.1.4 均值的多重比较	281
6.1.5 方差的齐次性检验	284
6.1.6 Kruskal-Wallis(克鲁斯卡尔-沃利斯)秩和检验	287
6.2 双因素方差分析	290
6.2.1 不考虑交互作用	290
6.2.2 考虑交互作用	292
6.2.3 方差分析表的计算	294
6.2.4 交互效应图	297
6.2.5 Friedman(弗里德曼)秩和检验	298
6.3 正交试验设计与方差分析	301
6.3.1 用正交表安排试验	301
6.3.2 正交试验的方差分析	304
6.3.3 有交互作用的试验	307
6.3.4 有重复试验的方差分析	310
习题 6	312
第 7 章 应用多元分析(I)	316
7.1 判别分析	316
7.1.1 距离判别	317
7.1.2 Bayes(贝叶斯)判别	327
7.1.3 Fisher 判别	334
7.1.4 用 MATLAB 软件中的函数作判别分析	338
7.2 聚类分析	345
7.2.1 距离和相似系数	345

7.2.2	样本间距离的 MATLAB 计算	349
7.2.3	系统聚类法	351
7.2.4	系统聚类法的 MATLAB 实现	354
7.2.5	类个数的确定	361
7.2.6	实例	365
7.2.7	动态聚类法	370
	习题 7	372
第 8 章 应用多元分析(II)		377
8.1	主成分分析	377
8.1.1	总体主成分	377
8.1.2	样本主成分	380
8.1.3	相关的 MATLAB 函数	383
8.1.4	实例	386
8.1.5	主成分分析的应用	390
8.2	因子分析	395
8.2.1	引例	395
8.2.2	因子模型	396
8.2.3	参数估计	398
8.2.4	因子旋转	409
8.2.5	因子得分	414
8.2.6	因子分析的计算函数	416
8.3	典型相关分析	421
8.3.1	总体典型相关	421
8.3.2	样本典型相关	423
8.3.3	典型相关系数的显著性检验	424
8.3.4	典型相关分析的计算	425
8.4	非负矩阵分解	430
8.4.1	非负矩阵分解的理论与方法	431
8.4.2	非负矩阵分解的 MATLAB 函数	431
	习题 8	433
附录 A MATLAB 软件简介		437
A.1	MATLAB 的工作界面	437
A.1.1	MATLAB 系统的安装	437
A.1.2	MATLAB 的工作界面	437

A. 1. 3	MATLAB 的帮助系统	439
A. 2	矩阵与数组的运算	441
A. 2. 1	向量与矩阵的表示	441
A. 2. 2	矩阵运算	442
A. 2. 3	数组运算	445
A. 2. 4	关系运算	449
A. 2. 5	逻辑运算	449
A. 2. 6	矩阵运算函数	450
A. 2. 7	基本函数	453
A. 3	程序设计	455
A. 3. 1	控制流	455
A. 3. 2	M 文件	457
A. 4	数据的导入和导出	460
A. 4. 1	低级函数	460
A. 4. 2	高级函数	463
A. 4. 3	读写 Excel 表	467
A. 5	绘图	469
A. 5. 1	二维绘图	469
A. 5. 2	三维绘图	474
A. 5. 3	与图形有关的函数	478
A. 5. 4	图形的保存	481
附录 B	MATLAB 数理统计工具箱	483
附录 C	答案	497
参考文献	505

第 1 章 概率论与数理统计基础

数据分析就是分析和处理数据的理论与方法。就广义而言，目前还没有固定的数据分析方法；但就狭义而言，数据分析需要用数理统计、多元分析的方法进行分析和处理。

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，在介绍用统计方法分析数据之前，有必要先回顾一下相关的概率论与数理统计的基本概念，以及数理统计的各个应用分支。

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 分钟内接到的呼叫次数，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，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。

随机试验的所有可能结果组成的集合称为样本空间，记为 Ω 。试验的每一个可能结果称为样本点，记为 ω 。

称 Ω 中满足一定条件的子集为随机事件，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。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。

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 $A \subset B$ ，亦称为事件的包含关系。

若 $A \subset B$ 且 $B \subset A$ ，则称事件 A 与事件 B 等价，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，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，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 为互斥事件或互不相容事件，记为 $AB = \emptyset$ 。

在一次试验中，基本事件之间是两两互斥的。

若事件 A 与事件 B 必有一个发生，且仅有一个发生，则称事件 A 与事件 B 互为对立事件或互逆事件，记为 \bar{A} 。事件 A 与事件 B 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

$$A \cup B = B \cup A, \quad AB = BA. \quad (1.1)$$

(2) 结合律

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

(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) 德·摩根律

$$\overline{A_1 \cup A_2} = \bar{A}_1 \cap \bar{A}_2, \quad \overline{A_1 \cap A_2} = \bar{A}_1 \cup \bar{A}_2. \quad (1.4)$$

对于 n 个或可列无穷个事件有

$$\begin{aligned} \overline{\bigcup_{k=1}^n A_k} &= \bigcap_{k=1}^n \bar{A}_k, & \overline{\bigcap_{k=1}^n A_k} &= \bigcup_{k=1}^n \bar{A}_k, \\ \overline{\bigcup_{k=1}^{\infty} A_k} &= \bigcap_{k=1}^{\infty} \bar{A}_k, & \overline{\bigcap_{k=1}^{\infty} A_k} &= \bigcup_{k=1}^{\infty} \bar{A}_k. \end{aligned} \quad (1.5)$$

(5) 减法满足

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

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关系。

定义 1.1 设随机试验 E 的样本空间 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 $\bar{A} \in \mathcal{F}$; (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 设随机试验 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, 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

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

即互斥事件之和的概率等于它们各自的概率之和。

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

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

性质 5 对任意两个事件 A 和 B ,

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

性质 5 可以推广为

$$P(A_1 \cup A_2 \cup A_3) = P(A_1) + P(A_2) + P(A_3) - P(A_1 A_2) - P(A_1 A_3) - P(A_2 A_3) + P(A_1 A_2 A_3), \quad (1.10)$$

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

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

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 $\Omega = \{\omega_1, \omega_2, \cdots, \omega_n\}$, 其中 n 为样本点总数。每个样本点 $\omega_i (i=1, 2, \cdots, n)$ 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。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, 试分别求以下两事件 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。

事件 A 所含基本结果数应是 k 个球在指定的 k 个格子中的全排列数, 即 $k!$, 所以

$$P(A) = \frac{k!}{l^k}.$$

为了算出事件 B 所含的基本事件数, 可设想分两步进行: 因为 k 个格子可以是任意选取的, 故可先从 l 个格子中任意选出 k 个来, 选法共有 $\binom{l}{k}$ 种; 对于每种选定的 k 个格子, 依上述各有一个球的推理, 则有 $k!$ 个基本结果, 故 B 含有 $\binom{l}{k} k!$ 个基本结果。所以

$$P(B) = \binom{l}{k} \frac{k!}{l^k} = \frac{l!}{(l-k)! l^k}.$$

概率论的历史上有一个颇为著名的问题——生日问题：求 k 个同班同学没有两人生日相同的概率。

若把这 k 个同学看作例 1.1 中的 k 个球，而把一年 365 天看作格子，即 $l=365$ ，则上述的 $P(B)$ 就是所要求的概率。例如， $k=40$ 时， $P(B)=0.109$ 。或者换句话说，40 个同学中至少两个人同一天过生日的概率是： $P(\bar{B})=1-0.109=0.891$ ，其概率大得出乎意料。

1.1.4 几何概型

当随机试验的样本空间是某一可度量(长度、面积或体积)的区域，并且任意一点落在度量相同的子区域内是等可能的，则事件 A 的概率定义为

$$P(A) = \frac{S_A}{S} = \frac{\text{构成事件 } A \text{ 的子区域的度量}}{\text{样本空间的度量}}. \quad (1.13)$$

这种概率模型称为几何概型。

例 1.2 蒲丰投针问题。设平面上画有等距为 a 的一簇平行线。取一枚长为 l ($l < a$) 的针随意扔到平面上，求针与平行线相交的概率。

解 设 x 表示针的中心到最近一条平行线的距离， θ 表示针与此直线间的交角，如图 1.1(a) 所示，则 (θ, x) 完全决定针所落的位置。针的所有可能的位置为

$$\Omega = \left\{ (\theta, x) : 0 \leq \theta \leq \pi, 0 \leq x \leq \frac{a}{2} \right\},$$

它可用 $\theta-x$ 坐标平面上的一个矩形来表示，如图 1.1(b) 所示。

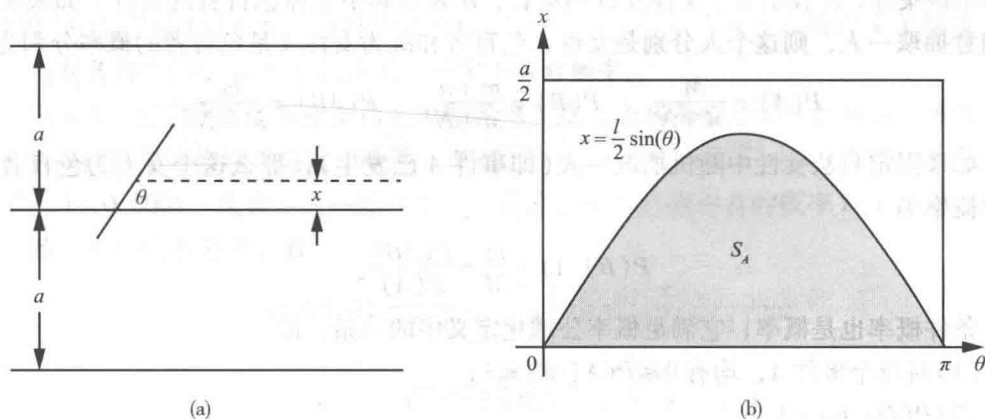


图 1.1 蒲丰投针的几何概率

针与平行线相交的充分必要条件是 $x \leq \frac{l}{2} \sin \theta$ ，即图 1.1(b) 中阴影部分，它的面积为