



北京大学数学教学系列丛书

本科生
数学基础课教材

概率与统计

(第二版)
(统计学分册)

郑忠国 陈家鼎 编著

馆外借



北京大学出版社
PEKING UNIVERSITY PRESS

北京大学数学教学系列丛书

概率与统计

(第二版)

(统计学分册)



北京大学出版社
PEKING UNIVERSITY PRESS

图书在版编目(CIP)数据

概率与统计. 统计学分册/郑忠国, 陈家鼎编著.—2 版.—北京: 北京大学出版社, 2017. 7

(北京大学数学教学系列丛书)

ISBN 978-7-301-28006-5

I. ①概… II. ①郑… ②陈… III. ①概率论 ②数理统计 IV. ①O21

中国版本图书馆 CIP 数据核字(2017)第 021154 号

书名 概率与统计(第二版)(统计学分册)

GAILÜ YU TONGJI

著作责任者 郑忠国 陈家鼎 编著

责任编辑 曾琬婷

标准书号 ISBN 978-7-301-28006-5

出版发行 北京大学出版社

地址 北京市海淀区成府路 205 号 100871

网址 <http://www.pup.cn>

电子信箱 zpup@pup.cn

新浪微博 @北京大学出版社

电话 邮购部 62752015 发行部 62750672 编辑部 62767347

印刷者 北京大学印刷厂

经销商 新华书店

890 毫米×1240 毫米 A5 9.625 印张 314 千字

2007 年 8 月第 1 版

2017 年 7 月第 2 版 2017 年 7 月第 1 次印刷

定价 35.00 元

未经许可, 不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有, 侵权必究

举报电话: 010-62752024 电子信箱: fd@pup.pku.edu.cn

图书如有印装质量问题, 请与出版部联系, 电话: 010-62756370

内 容 简 介

本书系统论述概率论和统计学的概念、方法、理论及其应用,是一部为高等院校本科生学习概率论和数理统计而编写的教材或教学参考书。本书不仅提供了这个学科领域的基本内容,而且叙述了在日常生活、自然科学、技术科学、人文社会科学及经济管理等各方面的应用例子。全书分为两册:概率论分册和统计学分册。概率论分册共五章,内容包括:随机事件与概率,随机变量与概率分布,随机向量,概率极限定理,随机过程。统计学分册共五章,内容包括:统计学中的基本概念,估计,假设检验,回归分析,统计决策和贝叶斯分析简介。本书恰当处理逻辑严谨性与生动直觉的辩证关系,使学生既有严谨的抽象思维能力,又对随机现象具有直觉想象力;认真贯彻理论联系实际,应用举例贴近时代生活;概率论部分强调了随机现象在社会生活和科学技术中的广泛性及所具有的内在规律,统计学部分则强调了其数据处理的功能,二者都以认识随机性、恰当处理随机性(包括决策和行动)为目标;内容选取上注意对难点进行化解,叙述通俗易懂,结构层次分明,使学生易于理解与掌握。

本书可作为高等学校理工类本科学生的教材或教学参考书,也可供经济管理和财经类等有关专业的研究生和从事统计计算的科技人员阅读。

“北京大学数学教学系列丛书”编委会

名誉主编：姜伯驹

主编：张继平

副主编：李忠

编委：（按姓氏笔画为序）

王长平 刘张炬 何书元 张平文

陈大岳 郑志明 柳彬

编委会秘书：方新贵

责任编辑：刘勇

作者简介

郑忠国 北京大学数学科学学院教授、博士生导师，1962年毕业于北京大学数学力学系，1965年北京大学研究生毕业。长期从事数理统计的教学和科研工作，研究方向是非参数统计、可靠性统计以及统计计算，发表论文近百篇。主持完成国家自然科学基金项目“不完全数据统计理论及其应用（1999—2001）”，教育部博士点基金项目“应用统计方法研究”和“工业与医学中的应用统计研究”等。研究项目“随机加权法”获国家教委科技进步二等奖。出版的教材有《高等统计学》（北京大学出版社，1995）。

陈家鼎 北京大学数学科学学院教授、博士生导师，1959年毕业于北京大学数学力学系。长期从事数理统计的教学和科研工作，研究方向是不完全数据的统计推断、序贯统计及其在可靠性工程上的应用，发表论文50多篇。曾任北京大学概率统计系主任、北京大学数学科学学院副院长、中国概率统计学会理事长、中国统计学会副会长。主持完成“序贯分析”“生存分析与可靠性的若干前沿问题”等多项国家自然科学基金和教育部博士点基金项目。主编的教材《数理统计学讲义》获国家教委优秀教材一等奖（高等教育出版社，1995）。与郑忠国等合作的项目“可靠性评定的数学理论与应用”获北京市科技进步二等奖（2002）。

序　　言

自 1995 年以来,在姜伯驹院士的主持下,北京大学数学科学学院根据国际数学发展的要求和北京大学数学教育的实际,创造性地贯彻教育部“加强基础,淡化专业,因材施教,分流培养”的办学方针,全面发挥我院学科门类齐全和师资力量雄厚的综合优势,在培养模式的转变、教学计划的修订、教学内容与方法的革新,以及教材建设等方面进行了全方位、大力度的改革,取得了显著的成效。2001 年,北京大学数学科学学院的这项改革成果荣获全国教学成果特等奖,在国内外产生很大反响。

在本科教育改革方面,我们按照加强基础、淡化专业的要求,对教学各主要环节进行了调整,使数学科学学院的全体学生在数学分析、高等代数、几何学、计算机等主干基础课程上,接受学时充分、强度足够的严格训练;在对学生分流培养阶段,我们在课程内容上坚决贯彻“少而精”的原则,大力压缩后续课程中多年逐步形成的过窄、过深和过繁的教学内容,为新的培养方向、实践性教学环节,以及为培养学生的创新能力所进行的基础科研训练争取到了必要的学时和空间。这样既使学生打下宽广、坚实的基础,又充分照顾到每个人的不同特长、爱好和发展取向。与上述改革相适应,积极而慎重地进行教学计划的修订,适当压缩常微、复变、偏微、实变、微分几何、抽象代数、泛函分析等后续课程的周学时。并增加了数学模型和计算机的相关课程,使学生有更大的选课余地。

在研究生教育中,在注重专题课程的同时,我们制定了 30 多门研究生普选基础课程(其中数学系 18 门),重点拓宽学生的专业基础和加强学生对数学整体发展及最新进展的了解。

教材建设是教学成果的一个重要体现。与修订的教学计划相配合,我们进行了有组织的教材建设。计划自 1999 年起用 8 年的时间

修订、编写和出版 40 余种教材。这就是将陆续呈现在大家面前的“北京大学数学教学系列丛书”。这套丛书凝聚了我们近十年在人才培养方面的思考，记录了我们教学实践的足迹，体现了我们教学改革的成果，反映了我们对新世纪人才培养的理念，代表了我们新时期数学教学水平。

经过 20 世纪的空前发展，数学的基本理论更加深入和完善，而计算机技术的发展使得数学的应用更加直接和广泛，而且活跃于生产第一线，促进着技术和经济的发展，所有这些都正在改变着人们对数学的传统认识。同时也促使数学研究的方式发生巨大变化。作为整个科学技术基础的数学，正突破传统的范围而向人类一切知识领域渗透。作为一种文化，数学科学已成为推动人类文明进化、知识创新的重要因素，将更深刻地改变着客观现实的面貌和人们对世界的认识。数学素质已成为今天培养高层次创新人才的重要基础。数学的理论和应用的巨大发展必然引起数学教育的深刻变革。我们现在的改革还是初步的。教学改革无禁区，但要十分稳重和积极；人才培养无止境，既要遵循基本规律，更要不断创新。我们现在推出这套丛书，目的是向大家学习。让我们大家携起手来，为提高中国数学教育水平和建设世界一流数学强国而共同努力。

张继平

2002 年 5 月 18 日

于北京大学蓝旗营

第二版前言

本书第二版对第一版进行了少量修改和补充,改动不大。概率论部分主要修订内容是:修改了个别不妥的文字和不正确的数字;举例说明强大数律与大数律的差别;对于初学者来说过于困难的几道习题,有的予以删除,有的予以改换。统计学部分主要是对原书中的笔误做了改正,对某些地方的表达方式做了一些修正,使得表达更精确和通顺。另外,由于统计学是面向实际应用的学科,近年来出现许多新方法,十分热门和流行。在第六章中,我们介绍了“大数据”这一方向,阐明它与统计学的关系,希望引起读者对这一当今热门对象的关注。在回归分析变量选择部分,我们还介绍了近年出现的 Lasso 方法,希望引起关注。

另外,考虑到现今许多高等院校理工类本科“概率论与数理统计”课程已改为“概率论”和“统计学”两门课程,本次修订我们将全书分为概率论分册和统计学分册,以满足课程改革的需要。

我们要特别强调的是,第二版和第一版一样,是为高等学校各专业本科学生学习“概率论与数理统计”而编写的教材,只要求学生预先学过“微积分”和“线性代数”的基础知识,不要求较深的数学知识(如实变函数、测度论)。但有一些内容打上 * 号或小字排印,这些内容或者难度较大,或者涉及较深的数学知识,均不属于教学大纲的范围,只供有余力的学生进一步学习时参考。

陈家鼎 郑忠国

2015 年 11 月

第一版前言

概率论是研究自然界、人类社会及技术过程中随机现象的数量规律的一门数学。数理统计学则是以概率论为指导，研究如何有效地收集和分析数据，以对所考查的问题进行推断或预测，直至为采取一定的决策和行动提供依据和建议。随着现代科学技术的迅速发展和人类生活条件的不断改进，概率论和数理统计学得到了蓬勃的发展。二者不仅形成了系统的理论，而且在自然科学、人文社会科学、工程技术及经济管理等方面有越来越广泛的应用。很多院校都开设“概率论”课、“数理统计”课或“概率统计”课。

最近几年我们二人一直担任北京大学数学科学学院为全校本科生开设的基础课程——“概率论”和“数理统计”的教学工作。这两门课各有 60 学时，学生来自文科、理科和医科的多个不同院系。本书正是在我们讲稿的基础上经过修改、扩充而成的，其中第一章至第五章由陈家鼎编写，第六章至第十章由郑忠国编写。

我们在编写过程中参考了国内外已有的特别是近十年出版的多部优秀教材（见本书的参考文献），注意吸收这些教材中好的讲法和具体例子。我们在编写中注意了下面三点：

(1) 恰当处理逻辑严谨性与生动直觉的关系，使学生既有严谨的抽象思维能力又有概率统计的直觉与对随机性的想象力。通过各方面的例子介绍有关的概念、方法和定理的实际含义，注意引导学生的思维从直觉和想象上升到科学的抽象。例如，既介绍了概率的“频率定义”和“主观定义”，又介绍了“公理化定义”，说明后者是在前者基础上的科学抽象。先介绍随机变量的直观含义和直观描述，然后介绍随机变量的严格定义。在介绍数学期望时先用加权平均的思想介绍离散型随机变量的期望，然后对一般的随机变量用离散型随机变量逼近的办法定义期望。对每个定理都给出确切的论述，能不用测度论证明的尽量写出证明，但由于教学时数的限制，许多证明打上 * 号或用小字排印，不要求

学生掌握。例如,对“两个随机变量之和的期望等于两个随机变量的期望之和”这一重要定理,我们在正文中只叙述了结论,但其详细证明则放在附录里小字排印。对“中心极限定理”和有关充分统计量的“因子分解定理”则不叙述证明。

(2) 认真贯彻理论联系实际的原则。既要使学生掌握概率和统计的基本理论,又要使学生认识这些理论如何灵活运用于实际,从而培养学生解决实际问题的能力。要做到这一点,必须要用心地列举贴近时代生活的,使学生感兴趣的多方面的应用例子。本书努力朝这个方向做。除了叙述日常生活、工业、商业、医学及管理等方面的应用例子(包括一些著名例子)外,还介绍一些较复杂的灵活应用例子。例如,第一章中作为独立试验序列的应用,介绍了乒乓球赛制的概率分析;第二章讲述随机变量取值的分散性时,除了“方差”外还介绍了经济学中常用的“基尼系数”;在讲述正态分布的性质之后,介绍了当今工业质量管理工作广泛关注的“ 6σ ”;第九章中作为回归分析的应用介绍了高考作文评分的监控方法,等等。

本书特别注重对理论联系实际的难点进行化解。例如,对“假设检验”,避免单纯从逻辑推理进行论述,着重从多方面的应用实例说明假设检验问题的提法、零假设的设置及两类错误的概率。把实际中的检验问题分成两大类:决策性检验问题和显著性检验问题。有些检验问题强调控制第一类错误的概率(例如第八章例 1.4),有些检验问题则重点在控制第二类错误的概率(例如第八章例 1.5)。本书还用一定篇幅介绍 p 值的概念和用法。又如,介绍“回归分析”的应用时把自变量分为两类:可控制的和不可控制的,把自变量和反应变量之间的关系分为两类:因果关系和非因果性的相关关系。本书还特别关注数据的来源和变量的性质。

(3) 在叙述方法与内容编排上注意基本内容与进一步内容、重点与非重点的界限,力求做到层次分明,便于教和学。我们认为,大学教材应比教学大纲规定的多一些,更应比课堂实际讲授的多一些。这样做有利于教师根据实际情况灵活掌握,有利于学生课外阅读,使有余力的学生可以选学更多的东西。本书中凡打 * 号和小字排印的部分均不是基本内容,不要求学生掌握。有些内容虽未标上 * 号也非小字排印,教师此为试读,需要完整PDF请访问: www.ertongbook.com

也可根据实际情况确定为非基本内容.

本教材虽是按两学期的教学安排(“概率论”课一学期,“数理统计”课一学期)编写的,但是也可作为一学期的“概率统计”课的教材.作为后者使用时,应选定书中最基本的部分.笔者建议选择下列内容:

第一章(不含 § 1.7),第二章(不含 § 2.8),第三章(不含 § 3.7, § 3.8),第四章 § 4.2 和 § 4.3 的部分内容,第五章的 § 5.1,第六章,第七章的 § 7.1, § 7.5,第八章的 § 8.1, § 8.2, § 8.4 中关于正态总体参数的检验方法,§ 8.6 中的 χ^2 检验,第九章的 § 9.1, § 9.2 及 § 9.3 至 § 9.5 中方法的应用部分,第十章的 § 10.1.

北京大学出版社刘勇和曾琬婷同志对本书的出版付出了辛勤的劳动,我们在此向他们表示感谢.

由于我们水平有限,本书一定有不少缺点和谬误,欢迎读者和专家批评指正.

陈家鼎 郑忠国

2007 年 6 月于北京大学数学科学学院

目 录

第六章 统计学中的基本概念	1
§ 6.1 引言	1
§ 6.2 若干基本概念	3
§ 6.3 若干统计问题	9
习题六	17
第七章 估计	19
§ 7.1 最大似然估计	19
§ 7.2 矩估计	30
§ 7.3 估计的无偏性	34
§ 7.4 无偏估计的优良性	39
§ 7.5 估计的相合性	52
§ 7.6 估计的渐近分布	57
§ 7.7 置信区间和置信限	62
习题七	81
第八章 假设检验	86
§ 8.1 问题的提法	86
§ 8.2 N-P 引理和似然比检验	95
§ 8.3 单参数模型中的检验	101
§ 8.4 广义似然比检验和关于正态总体参数的检验	114
§ 8.5 关于比率的检验	135
§ 8.6 拟合优度检验	142
习题八	156
第九章 回归分析	160
§ 9.1 引言	160

§ 9.2 一元线性回归	166
§ 9.3 多元线性回归	175
§ 9.4 多元线性回归中的参数检验	185
§ 9.5 预测和控制	197
*§ 9.6 模型检验	206
*§ 9.7 变量选择	214
§ 9.8 方差分析	224
*§ 9.9 逻辑斯谛回归	237
习题九	240
 第十章 统计决策和贝叶斯分析简介	246
§ 10.1 统计决策问题概述	246
§ 10.2 贝叶斯统计	253
§ 10.3 先验分布的确定	260
习题十	265
 习题答案与提示	267
附表 1 标准正态分布数值表	275
附表 2 t 分布临界值表	276
附表 3 χ^2 分布临界值表	277
附表 4 F 分布临界值表	278
附表 5 柯氏检验临界值表	284
参考文献	286
名词索引	288

第六章 统计学中的基本概念

§ 6.1 引言

在学习数理统计学之前,我们必须弄明白什么是数理统计学,数理统计学的研究对象是什么.为回答这些问题,我们要引用《中国大百科全书·数学》(中国大百科全书出版社,1992)中关于数理统计学的定义:数理统计学研究怎样去有效地收集、整理和分析带有随机性的数据,以对所考查的问题做出推断或预测,直至为采取一定的决策和行动提供依据和建议.这句话规定了数理统计学的研究内容.这句话或许太抽象,有些学术化,不易吸引初学者的兴趣.我们再引一段话,它来自David Freedman的名著《统计学》(见文献[33]):“统计学是什么?统计学是对令人困惑费解的问题做出数学设想的艺术.应该怎样设计实验来测定新药的疗效?什么东西引起父母与孩子之间的相像,并且那种力量有多强?通货膨胀率如何测定?失业率呢?它们怎样联系起来?赌场为什么在轮盘赌上得益?盖洛普民意测验怎么能够使用仅仅几千人的样本预测美国大选结果?”David Freedman用生动的事例描述如何应用数理统计这个工具去解决这些人们困惑的问题.由这些生动的叙述看出,统计学试图解决人们困惑的问题,并且涉及的范围非常广泛,从天文地理、尖端科技、社会经济直至日常生活.在日常生活中,也有使人感兴趣的问题.例如,某学校希望从两名考生中录取一名学生,他们的考试成绩如表6.1.1所示.

表 6.1.1 学生成绩表

(单位:分)

学生	语文	数学	英文	物化
甲	80	70	60	85
乙	70	90	65	80

学校希望了解这两名学生在学习能力上有没有差别。单由成绩的高低，并不能说明两者之差异。因为即使是同一个人，在两次水平相同的考试中也可能具有不同的成绩；即使同一份考卷，不同阅卷教师也可能给出不同的分数。要解决这个问题，必须求助于统计学。在上述问题中，学校的目的是希望选择一个学习能力较强的学生继续培养。为了达到这个目的，学校必须对这两位学生进行考查。对此问题，单凭一双慧眼就能成为识别千里马的伯乐是几乎不可能的。有经验的考查者必须收集一些资料，对这些资料进行分析，从中得到所需要的结论。由于学生的能力只能从他吸收知识，对外界信息的反应上表达出来，因此考查者必须设计一些方法来收集这方面的资料。比较传统的方法是上述的考试，看他们之间的差异，从中挑选能力强的学生。上面提到的挑选学生的过程就反映了数理统计学的任务。首先考查者必须收集资料，这些资料就是数理统计学中的数据。为了达到这个目的，考查者必须提出合适的问题，或者出几份合适的考卷。这就是一个如何有效地收集数据的问题。如果在试卷中，出一些要求学生死记硬背的题目，这些题目就不能反映学生的学习能力，这样的方法就不是有效的方法，从死记硬背得到的答案看不出其能力。现在即使考卷都合适，而且两位学生的成绩已经列出，仍然还存在判断两学生能力差异的问题。在对数据进行分析的基础上，还需做出一些判断，比如乙在某些方面的能力比甲的强。这些判断将为录取学生提供根据。

上述考查学生的例子能够完整地说明数理统计学的任务。其中，出题等收集数据问题虽然也涉及其他专业的有关知识，但也是数理统计学的任务。不过在本书中，我们不将收集数据作为重点。本书的重点是介绍分析数据的技术。数据是现实生活中大量存在的。统计工作者将这些数据看成从某个信息源以某种方式释放出来的信号。统计工作的任务是分析数据——将这些信号进行加工处理，提取其中的信息。显然，分析数据是统计工作的核心任务。

附带说明，“统计学”与“数理统计学”实质是同一个学科。通常，在统计研究者强调数学方法时，将“统计学”前面加上一个“数理”的形容词。

§ 6.2 若干基本概念

统计学的研究对象是数据. 什么是数据呢? 广义地讲, 数据就是我们在实际工作中的记录. 例如, 某工厂为了考查某些电子产品的使用寿命, 随机地抽取了 18 台产品做试验, 测得寿命数据(单位: h)如下:

17, 29, 50, 68, 100, 130, 140, 270, 280, 340,
410, 450, 520, 620, 190, 210, 800, 1100.

这 18 个寿命值就是数据, 就是我们的研究对象. 又例如, 某社会工作者调查某城市中成年吸烟者占成年人口的比例, 共调查了 339 人, 其中 205 人为吸烟者, 134 人为不吸烟者. 同样, 数据 339, 205, 134 是我们的研究对象. 若不对这些数据进行合理的抽象, 就不可能对这些数据进行深层次的分析, 从中获得更多的信息. 在数据处理时, 我们通常用 x 表示数据, 此处 x 既可以是一个数, 也可以是一个向量或其他的量. 当明确表示向量或向量与它的分量同时出现时, 我们用黑体 x 表示之. 这时数据的主要形式是 $x = (x_1, \dots, x_n)$. 在实际问题中, 有时候单一个字母是不够用于表达数据的. 例如, 在连续 10 天的气象记录中, 得到 $m_1, \dots, m_{10}, M_1, \dots, M_{10}$, 其中 m_i ($i=1, \dots, 10$) 是每天的最低气温, M_i ($i=1, \dots, 10$) 是每天的最高气温, 此时的数据为 $x = \{(m_i, M_i), i=1, \dots, 10\}$. 但是, 在学习统计学的时候, 用 $x = (x_1, \dots, x_n)$ 表示数据是最方便并且能够抓住数据本质的一种方法. 本节开头引入的寿命数据可表达成 $x = (x_1, \dots, x_{18})$ 或 $x = (x_1, \dots, x_n), n=18$.

引入数据的概念以后, 我们要记住统计工作的核心任务是对数据进行分析, 进而对所考查的问题做出推断. 在寿命数据的问题中, 我们的任务是考查该厂生产的电子产品的使用寿命. 我们收集到的 18 台电子产品的寿命数据是该厂生产的一部分产品的数据. 此处特别强调, 我们的目的是要了解该厂生产的电子产品的使用寿命, 而不是这 18 台产品的使用寿命. 这 18 台产品的使用寿命是已经明摆着的数据, 不必再进行细究. 为了研究产品的使用寿命, 我们必须弄明白, 什么是工厂生产的电子产品的使用寿命, 而且还要弄清楚这 18 台产品的寿命与该厂此为试读, 需要完整PDF请访问: www.ertongbook.com