

精·品·课·程·立·体·化·教·材·系·列

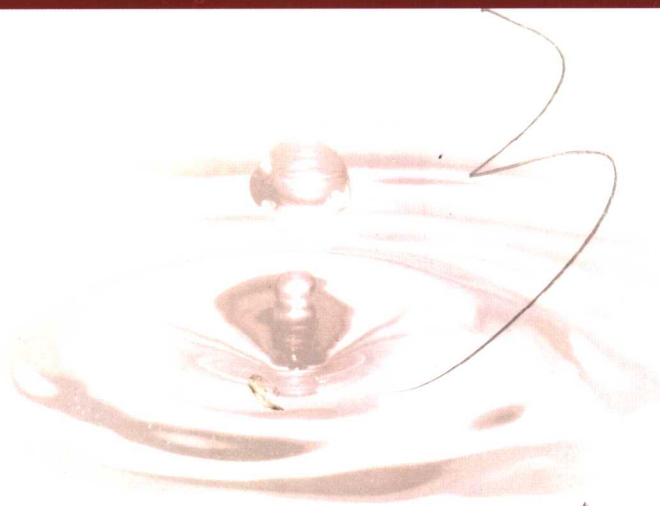


应用统计学：

数理统计方法、数据获取与SPSS应用

(精要版)

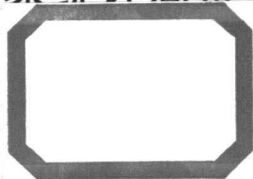
马庆国 编著



 科学出版社
www.sciencep.com

精·品·课·程·立·体·化·教·材·系·列

2004年国家统计局推荐教材

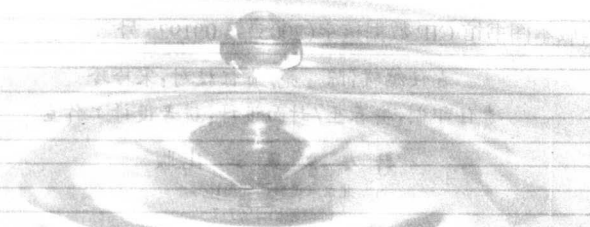


应用统计学：

数理统计方法、数据获取与SPSS应用

(精要版)

马庆国 编著



科学出版社

北京

内 容 简 介

本书是《管理统计》(马庆国著,科学出版社2002版)的精要版。它依托社会科学领域使用最为广泛的SPSS统计软件为基本工具,结合实际案例讲述了社会科学特别是经济管理学科中常用的统计学原理、方法与技术。本书突破了注重数学推理与公式证明的传统统计学教材的框架,真正把讲述的重点放在如何应用统计学原理、方法来解决科学研究与社会生活中的实际问题上,在强调透彻理解统计学原理和方法的基础上,着重论述了如何从解决实际问题的需要出发,进行数据收集,设计统计调查方案,并利用SPSS软件的强大功能,分析数据并解释分析结果,从而构建起一个贯通统计学原理、SPSS软件使用技巧与统计学应用研究方法的应用型统计学知识体系,可迅速、有效地提高学生分析和解决问题的综合能力。

本书配备多媒体教学课件、教学录像以及习题案例集等多种立体化教学支持,可作为经济管理类专业以及其他社会科学领域的统计学教材,也可供其他对应用统计学和SPSS软件有兴趣的实际工作者参考。

图书在版编目(CIP)数据

应用统计学:数理统计方法、数据获取与SPSS应用(精要版)/马庆国编著. —北京:科学出版社,2005

精品课程立体化系列=中国科学院规划教材

ISBN 7-03-016143-2

I. 应… II. 马… III. 应用统计学-教材 IV. C8

中国版本图书馆CIP数据核字(2005)第091925号

责任编辑:陈 亮/责任校对:宋玲玲

责任印制:安春生/封面设计:耕者设计工作室

科学出版社 出版

北京东黄城根北街16号

邮政编码:100717

<http://www.sciencep.com>

双青印刷厂 印刷

科学出版社发行 各地新华书店经销

*

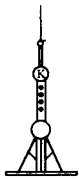
2005年11月第 一 版 开本: B5 (720×1000)

2005年11月第一次印刷 印张: 21 1/4

印数: 1—4 000 字数: 402 000

定价: 28.00元 (含光盘)

(如有印装质量问题, 我社负责调换〈环伟〉)



精要版序言

长期以来，我国非数学专业的《应用统计》课程都是以数学证明和手工计算为主要内容的。本来，在大学非数学专业设置数理统计类课程的初衷，是为非数学专业的学生提供研究本专业问题、处理本专业的数据、挖掘深层信息的工具，但是以数学证明为主要内容的课程安排，显然与这个初衷不协调。非数学专业的学生在学习时，用了大量的时间和精力，学会了统计问题的数学证明，但在毕业后需要用数理统计处理本专业数据时，却发现根本不需要证明什么数理统计的结论，需要的是对数理统计原理和方法的正确理解，以及正确地处理数据（计算）。而解决本专业实际问题所需要的知识和能力，恰恰是以往数理统计类课程没有关注的。手工计算根本不能面对实际中涌现的大量数据。因此，这类课程的内容设置，必须根本变革，以适应非数学专业的研究和处理实际问题的需要。

如果说，在 20 世纪 80 年代中期，受限于计算工具的不发达（计算机不普及，几乎没有专门的应用统计软件），我们为非数学专业设置的数理统计类课程，以数学证明和手工计算为主要内容，多少有一些无奈的话，那么，到了今天，PC 机是如此的普及，任何一所大学的本科生上机练习，都不是一件困难的事，而且相当多的学生都有自己的手提电脑，同时，专用的各类应用统计软件是如此发达，在这种情形下，我们还是固守 20 世纪 80 年代中期的构思，为非数学专业的学生讲授统计学，就太保守了。这类课程的设置思路和内容的变革，不仅是必须的，而且有了广泛的物质手段和计算工具的基础。

基于现代信息技术和工具，非数学专业的数理统计类的课程变革，是必然的、不可逆转的历史发展。

本教材，就是依据这样的判断设计的。

本教材删去了数理统计学中的几乎所有的数学证明，给出了得出相应定理和方法的思路，以及相应统计方法的实质和特点，以利于学生“从总体上把握”和“从应用角度把握”相应的统计方法和原理，避免“只见树木，不见森林”的弊端。

在从总体上把握某一统计原理和方法之后，紧接着给出运用这一统计方法处理数据的相应软件的应用，以便学生能够针对本专业的问题，应用统计方法，处理数据，挖掘信息。

考虑到本教材主要是针对社会科学专业的，所以选用了更适用于社会科学的SPSS软件，作为处理数据的软件工具。

虽然SPSS软件是从解决社会科学问题的统计分析起家的(Statistics Package for Social Science)，但同样可以解决非社会科学的统计分析问题。一方面，从数学角度看，统计分析方法本来就是针对问题结构的，只要具有同样的问题结构，不论是社会科学的问题，还是非社会科学的问题，都可以用同样的方法解决，所以，也就可以用执行同样统计分析任务的程序来计算，而不论这段程序包含在哪个软件包中。另一方面，SPSS公司从市场角度出发，也已经注意到了该软件的更宽的适用性问题，所以在功能延伸后，SPSS已经被解释为：统计产品和服务解决方案(Statistical Product and Service Solutions)。

当然，学习者可以用其他任何恰当的软件来处理数据，做统计分析。一般来说，有了使用一个统计软件的基础，再学习其他统计分析软件是不困难的。

在实际应用中，不论什么专业，数据都不是从天上掉下来的。因此，在以应用为中心的统计课程中，增加数据收集的基本概念、理论和方法，是有意义的。本教材的第2章，就是为此目的设置的。

第5~9章的统计分析方法，都是基于概率论的。为了正确理解这些统计处理方法，把握概率论的基础知识是必要的，因此，本教材的第1章对概率论的基础知识做了一个高度浓缩和深入浅出的介绍。对于学过概率论的读者而言，这有利于他们从更高的角度，概括性地理解概率论，通过对基础知识的透彻理解，产生对概率论认识的升华。对于没有学过概率论的读者而言，依照如此扼要的版本，为着应用的目的，来学习概率论的主脉络，也不失为一个较好的选择。

本教材是在科学出版社出版的专著《管理统计：数据获取、统计原理、SPSS工具与应用研究》的基础上缩编、增补而成的。

本教材的各章结构如下：

第1章“概率论基础知识”。对概率论基础知识作概要回顾和复习。

第2章“数据与数据的获取”。重点介绍了获取数据的问卷方法和实验方法，介绍信度和效度问题，讨论了这一领域的常见错误和必须注意的关键问题。

第3章“样本数据特征的初步分析”，是对样本数据集合有关特征的初步认识，基本上不用概率论的基础知识。

第4章“SPSS的简单应用”，主要目的是初步介绍SPSS软件，为第5~9章（在介绍有关统计知识后）应用这个软件来处理数据奠定基础。

第5章“总体分布、样本分布与参数估计”。从本章以后，所有统计学的知识都建立在概率论的基础上了。本章的主要目的是介绍统计学的基本定理、常用的统计量和有关性质，同时对参数估计和估计效果的优劣判别准则，做了一个简单介绍。

第6章“参数假设检验”，介绍了用样本数据检验母体分布的参数假设和检验两个母体参数的比较问题。重点纠正了不做参数检验、直接用样本估计值代替母体参数的错误做法。

第7章“方差分析”，介绍了单因素和双因素方差分析问题。这是处理数据，特别是实验数据，寻找关键因素的、应用面十分广泛的方法。

第8章“相关分析”，介绍了常用的 Pearson 积矩相关、基于曲线变换的非线性相关、基于顺序级数据的 Spearman 等级相关、偏相关、基于一个二值名义级变量的点双列相关和两个二值名义级变量的 ϕ 相关系数。

第9章“线性回归分析”，介绍了多元线性回归分析的基本概念、有关参数的估计方法、几何解释，逐步回归和名义级解释变量的处理方法。

第5~9章中，在介绍了相应的统计原理和方法之后，都介绍了SPSS中相应模块的应用，以便学生在学习相应数理统计方法后，能够直接借助软件处理数据，解决本专业的实际问题。

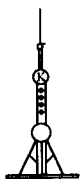
总之，本书的主要特点可以简单地概括为十六个字：略去证明，讲清原理，依托软件，突出实用。

目的只有一个：学以致用。使学生掌握研究群体现象的基本方法，获得处理实际问题的本领。

马庆国

2005年7月

于浙江大学求是园



目 录

精要版序言

第 1 章

概率论基础知识..... 1

1.1 随机实验、样本空间、概率与条件概率..... 1

1.2 随机变量与概率分布的基本概念..... 12

1.3 几个典型的概率分布..... 23

思考与练习题 26

第 2 章

数据与数据的获得..... 27

2.1 总体、个体、特征与数据..... 27

2.2 数据类型..... 28

2.3 获得数据的抽样调查方法..... 32

2.4 问卷与问卷设计..... 38

2.5 获得数据的信度与效度的基本概念..... 50

2.6 获得数据的实验方法..... 57

思考与练习题 63

第3章

样本数据特征的初步分析..... 66

3.1 样本数据结构的基本特征: 频次与频率..... 66

3.2 观察刻度级样本数据结构的茎叶图与直方图方法..... 74

3.3 样本数据的位置特征: 对数据中心的描述..... 82

3.4 样本数据的离散特征..... 86

3.5 样本数据特征的综合表达: 箱形图..... 92

思考与练习题..... 96

第4章

SPSS 的简单应用..... 100

4.1 使用 SPSS 的基础知识..... 100

4.2 SPSS 的简单应用..... 119

思考与练习题..... 143

第5章

总体分布、样本分布与参数估计..... 145

5.1 总体分布与样本分布..... 145

5.2 统计量与统计量的分布..... 150

5.3 点估计..... 156

5.4 判断点估计的优劣标准..... 159

5.5 区间估计..... 161

5.6 SPSS 在参数估计中的应用..... 166

思考与练习题..... 173

第6章

参数假设检验..... 175

6.1 假设检验的基本概念..... 175

6.2 一个正态总体下的参数假设检验..... 176

6.3	一个 0-1 总体分布下的参数假设检验	186
6.4	两个正态总体下的参数假设检验	191
6.5	大样本下两个任意总体的均值检验	195
6.6	用 SPSS 作假设检验	198
	思考与练习题	212

第 7 章**方差分析**

7.1	单因素方差分析	215
7.2	用 SPSS 作单因素方差分析	219
7.3	无重复实验的双因素方差分析	225
7.4	重复实验的双因素方差分析	229
7.5	用 SPSS 作双因素方差分析	235
	思考与练习题	244

第 8 章**相关分析**

8.1	两个随机变量的总体相关与样本相关	247
8.2	非线性相关	249
8.3	Spearman 等级相关	250
8.4	偏相关	251
8.5	相关系数异于零的显著性检验	253
8.6	SPSS 对普通相关分析的处理	255
8.7	至少有一个变量是二值名义级的相关系数	270
	思考与练习题	275

第 9 章**线性回归分析**

9.1	一元线性回归	278
-----	--------------	-----

9.2 多元线性回归	285
9.3 逐步回归	292
9.4 虚拟解释变量问题	295
9.5 用 SPSS 处理经典回归问题	299
9.6 曲线回归与 SPSS 的应用	313
思考与练习题	317
主要参考文献	320
附录	
常用数理统计表	321
附表 1 标准正态分布表	321
附表 2 χ^2 分布表	322
附表 3 t 分布表	324
附表 4 F 分布表	325

第 1 章 概率论基础知识

1.1 随机实验、样本空间、概率与条件概率

1.1.1 一些基本概念

1. 随机实验(Random Trial, or Random Experiment)

“抛硬币”就是一个简单的随机实验,可用此例子来理解、记忆如下概念:①在同一条件下可无限次重复的实验;②实验结果有多个,且不确定;③事前不知实验结果(Outcome)。

2. 基本事件(Elementary Event)

一次随机实验的可能结果,称为基本事件或基本随机事件。

若随机实验 E 是“抛 2 次(或先后 2 枚)硬币”,其基本事件就是“正、反”,“正、正”,“反、正”,“反、反”。

3. 样本空间(Sample Space)

所有基本事件所组成的集合,称为样本空间或基本空间。

例如,随机实验 E 是“抛 2 次(或先后 2 枚)硬币”,其样本空间就是集合: {“正、反”,“正、正”,“反、正”,“反、反”}。

4. 随机事件(Random Event)

随机事件简称事件,指一些基本事件所组成的集合。

例如,随机实验 E 是“抛 2 次(或先后 2 枚)硬币”,事件“两枚出现相同面”,就由两个基本事件组成:“正、正”、“反、反”。事件“至少出现 1 个正面”,就由 3 个基本事件“正、正”、“正、反”、“反、正”组成。

5. 相容事件 (Mutually Inclusive Events) 与不相容事件 (Mutually Exclusive Events)

在随机实验中,不能同时发生或其交集为空集的几个事件,称为不相容事件,

反之就称为相容事件。

例如,随机实验 E 为“先后投 2 枚硬币”, A 事件“两枚出现相同面”(由事件“正、正”和事件“反、反”构成)与 B 事件“两枚出现不同的面”(由事件“正、反”和“反、正”构成),就是两个不相容的事件。 C 事件“至少出现 1 枚正面”(由“正、正”、“正、反”、“反、正”组成)与 D 事件“至少出现 1 枚反面”(由“反、反”、“正、反”、“反、正”组成),就是两个相容的事件。

若形象地把事件看成平面上的点集,那么,若 A 与 B 没有共同的点,则 A 与 B 就是不相容事件;若 A 与 B 有共同的点,则 A 与 B 就是相容事件。

6. 概率(Probability)

用通俗的语言说,概率指在随机实验中,对事件出现的可能性大小的一种严格的度量。所谓严格,是指从无限次重复的角度看,度量结果具有惟一性。

随机实验 E “抛 1 枚均匀的骰子”,做 60 次实验,1 点(俗称大点)朝上的次数,可能是 9, 10 或 11, 那么 1 点朝上的频率分别是 $9/60, 1/6, 11/60$ 。所以,频率不是随机实验中事件出现的可能性的严格度量,不是概率。但是,随着实验的次数的增加,频率约等于概率。所以,可以通俗地将频率理解为“概率的模糊的影子”。

概率的定义:设 E 是随机实验, S 是其样本空间,给 E 的每一个事件 A 赋一个实数 $P(A)$,若 $P(A)$ 满足如下条件(Postulate),就称为 A 的概率:①对每一个事件 A , 都有 $0 \leq P(A) \leq 1$;② $P(S) = 1$;③对于两两互不相容的事件 $A_k (k = 1, 2, \dots)$, 有: $A_k (k = 1, 2, \dots)$ 的并集的概率等于各个 A_k 的概率之和:

$$P\left\{\bigcup_k A_k\right\} = P\{A_1 \cup A_2 \cup \dots\} = P(A_1) + P(A_2) + \dots = \sum_k P(A_k)$$

7. 概率运算的主要性质(Properties of Probability)

(1) 设 \bar{A} 是 A 的对立事件(若把事件看成平面上的点集,则 \bar{A} 是样本空间 S 中除 A 以外的所有的点),则 $P(A) = 1 - P(\bar{A})$ 。当 \bar{A} 为样本空间 S 时, A 为空集 \emptyset , 从而可以得出空集 \emptyset 的概率为 0:

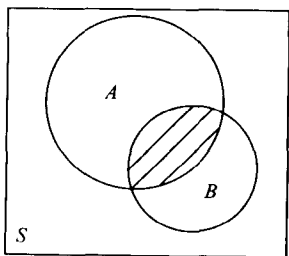


图 1.1.1 事件并集的概率计算

$$P(A) = 1 - P(\bar{A}) = 1 - P(S) = 0$$

(2) 对任意两个事件 A 与 B , 有

$$P\{A \cup B\} = P(A) + P(B) - P(AB)$$

其中, AB 是 A 与 B 的交集 $A \cap B$ 的缩写。为了便于理解上式,可参见图 1.1.1。

关于交集概率 $P\{A \cap B\}$ 的计算公式,我们将在下一节的“概率的乘法定理”中给出。

若 A 与 B 的交集 $A \cap B$ 为空集(即 A 与 B 不

相容或不相交),则 A 与 B 的并集 $A \cup B$ 的概率,就是 A 、 B 两个事件的概率之和:

$$P\{A \cup B\} = P(A) + P(B)$$

(3)若事件 $A \subset B$,则 $P(A) \leq P(B)$ 。

8. 等概率随机实验(Equally Likely Outcomes)

若一个随机实验的基本事件的个数有限,且基本事件出现的概率相等,则该随机实验称为等概率随机实验或等可能概型。抛均匀的硬币、抛均匀的骰子,都是等概率随机实验的例子。

在等概率随机实验中,事件 A 的概率计算公式是:

$$P(A) = \frac{A \text{ 包含的基本事件个数}}{\text{该实验中基本事件的总个数}}$$

例如,随机实验 E “先后抛 2 枚均匀的硬币”,共有 4 个基本事件:“正、正”、“正、反”、“反、正”、“反、反”。若我们要考察的事件 A 是“至少有 1 枚正面朝上”,那么,事件 A 包括 3 个基本事件“正、正”、“正、反”、“反、正”,所以 $P(A) = 3/4$ 。

1.1.2 条件概率与概率乘法定理

1. 条件概率(Conditional Probability)

我们用以下实例来介绍条件概率。

例 1.1.1 一个包装箱里有 6 个产品。假设其中 4 个为一级品,2 个为二级品。若随机实验 E 是“从该包装箱中抽取 1 个产品”,那么,显然每次抽取到二级品的概率是 $1/3$ 。

设事件 A 是“第一次抽取并抽到二级品”,事件 B 是“第二次抽取并抽到二级品”,那么,在事件 A 发生的条件下,再从剩余的 5 个产品中抽出 1 个产品,事件 B “第二次抽到二级品”发生的概率就是 $1/5$ 。我们称这样的概率为“事件 A 发生条件下,事件 B 发生的概率”,简称为“事件 B 的条件概率”,记为 $P\{B|A\}$ 。

在例 1.1.1 中, $P\{B|A\} = 1/5$ 。

为对比条件概率与非条件概率的区别,我们看 $P(B)$ 等于多少。首先,事件 B 可以表达为不相容的 2 个事件 AB 与 \overline{AB} 之和,即 $B = AB + \overline{AB}$ (若将 B 看成 S 平面上的点集,相当于用 A 与 \overline{A} 把 B 分成不相交的两块),于是,由不相交事件的概率的性质,有

$$P(B) = P(AB) + P(\overline{AB}) = \frac{2}{6} \times \frac{1}{5} + \frac{4}{6} \times \frac{2}{5} = \frac{1}{3}$$

其实,对于原来的随机实验(也是原来的样本空间)而言,很显然,无论是事件 A (第一次抽取且抽到二级品),还是事件 B (第二次抽取且抽到二级品),出现的概率

都是 $1/3$, 即 $P(A) = P(B) = 1/3$ 。

由上例, 不难理解如下公式:

$$P\{B|A\} = \frac{P(AB)}{P(A)} = \frac{1/15}{1/3} = \frac{1}{5}, \quad P(A) > 0$$

一般都把这个式子, 作为条件概率的定义式。

条件概率的定义

对样本空间 S 中的两个事件 A, B , 若 $P(A) > 0$, 则条件概率 $P\{B|A\}$ 由如下定义式给出:

$$P\{B|A\} = \frac{P(AB)}{P(A)}$$

2. 概率的乘法公式(定理)(Multiplication Theorem)

由条件概率定义式, 易得

$$P(AB) = P\{B|A\} \times P(A) = P\{A|B\} \times P(B)$$

这就是所谓概率的乘法公式或乘法定理。

概率的乘法公式

对样本空间中任意两个事件 A, B , 有

$$P(AB) = P\{B|A\} \times P(A) = P\{A|B\} \times P(B)$$

条件概率的本质是, 事件 A 的出现, 改变了产生事件 B 的范围条件, 即改变了样本空间, 也就是改变了随机实验(注意: 产品的个数不是样本空间)。事件 A, B 之间也不一定要有特定的时间先后的关系。

3. 条件概率计算表

我们可以把上述数据列成表 1.1.1:

表 1.1.1 条件概率计算表

	B : 第二次抽且二级品	\bar{B} : 第二次抽且一级品	总计
A : 首次抽且二级品	$\frac{1}{3} \times \frac{1}{5}$	$\frac{1}{3} \times \frac{4}{5}$	$\frac{1}{3}$
\bar{A} : 首次抽且一级品	$\frac{2}{3} \times \frac{2}{5}$	$\frac{2}{3} \times \frac{3}{5}$	$\frac{2}{3}$
总 计	$\frac{1}{3}$	$\frac{2}{3}$	1

中间4个方格中的数据,是相应行、列两个事件同时发生的概率(注意:完全可以是其他的数量描述;而不是两两事件联合发生的概率,见下例)。

这个表格对于理解和计算条件概率非常有用,我们称之为“条件概率计算表”,简称“条件概率表”。

在该表中,若以最初的随机实验为基础,考虑 A 与 B 发生的概率问题,应当分别看“总计”列与“总计”行。

对 $P(A)$ 而言,应看“总计”列。“总计”列的最后一个数据,即整个表格的右下角数据(在本例中是“1”),应当作随机实验的样本空间概率1,所以应当用它去除其他相应的数据,从而求得相应事件的概率。因此,有

$$P(A) = \frac{1/3}{1} = \frac{1}{3}$$

对 $P(B)$ 而言,应看“总计”行,有

$$P(B) = \frac{1/3}{1} = \frac{1}{3}$$

在发生事件 A 的条件下,考虑事件 B 的发生概率相当于考虑一个新的随机实验,应看事件 A 所在的行,以该行的“总计”数 $1/3$ 为子(新)样本空间的概率1,用 $1/3$ 去除该行的每个数据,以求得相应的概率。于是,有

$$P\{B | A\} = \frac{P(AB)}{P(A)} = \frac{1/15}{1/3} = \frac{1}{5},$$

$$P\{\bar{B} | A\} = \frac{P(A\bar{B})}{P(A)} = \frac{4/15}{1/3} = \frac{4}{5}$$

同样,可以用类似的方法考虑并求出事件 \bar{A} 条件下事件 B 与事件 \bar{B} 概率:

$$P\{B | \bar{A}\} = \frac{P(\bar{A}B)}{P(\bar{A})} = \frac{4/15}{2/3} = \frac{2}{5},$$

$$P\{\bar{B} | \bar{A}\} = \frac{P(\bar{A}\bar{B})}{P(\bar{A})} = \frac{6/15}{2/3} = \frac{3}{5}$$

如前所述,特定时间并不是条件概率中的必要概念。由概率的乘法公式,我们完全可以有

$$P\{A | B\} = \frac{P(AB)}{P(B)}, \text{ 其中 } P(B) > 0$$

由表 1.1.1 中的数据,也可以直接求出上式的条件概率:

$$P\{A | B\} = \frac{1/15}{1/3} = \frac{1}{5}$$

在例 1.1.1 中,由于事件 A 、 B 分别被冠以“第一次”、“第二次”的字样,使得我们对于条件概率的实质难以有比较深刻的理解,为此,我们再来看例 1.1.2。

例 1.1.2 某城市市民的肝炎患病率只有 0.01%,但是在某验血指标为阳性的人群中,得肝炎的概率是 90%。那么,在该城市中任意碰到一个人,该人是肝炎

患者的概率就是 0.0001。若碰到的这个人的验血指标是阳性,则该人为肝炎患者的概率是 0.9。若我们规定事件 A 为“碰到的是验血指标为阳性的人”,事件 B 为“碰到的是肝炎患者”,那么, $P(B) = 0.0001$, $P\{B|A\} = 0.9$ 。事件 A 改变了产生事件 B 的范围条件,本质上是改变了样本空间,也就改变了随机实验。

设 R 为验血指标为阳性的人口占城市人口的比例,根据以上数据,由“总人口 $\times R \times 0.9 = 0.0001 \times$ 总人口”,我们容易得出: $R = 0.0001/0.9 = 1/9000$ 。我们可以把上述数据列成条件概率计算表 1.1.2:

表 1.1.2 条件概率计算表

	B :肝炎	\bar{B} :非肝炎	总 计
A :阳性	0.9 ^①	0.1 ^①	1 ^①
\bar{A} :非阳性	0 ^②	8999 ^⑤	8999 ^③
总计	0.9 ^⑥	8999.1 ^⑦	9000 ^②

注:表中的顺序号①、②、…,表示表中数据填入的顺序。例如, A 行的数据是第①步填入的。它表示:在该城市中碰到 1 个单位的验血为阳性的人,其中有 0.9 个单位的人为肝炎患者,0.1 个单位的人为非肝炎患者, A 行的总计是 1 个单位的人。第②步是依据 $R = 1/9000$ 来填写的。如果城市中验血为阳性的人是 1 个单位,那么该城市的总人口就是 9000 个单位。第③步……。

由表 1.1.2 的“总计”行,易验证: $P(B) = 0.9/9000 = 0.0001$ 。为计算发生事件 A 条件下的 B 事件的概率,应当看事件 A 所在行,易得 $P\{B|A\} = 0.9/1 = 0.9$ 。进而,可以计算 $P\{A|B\}$ (碰到该城市的一个肝炎患者,其验血指标是阳性的概率)。看条件事件 B 所在列,有 $P\{A|B\} = 0.9/0.9 = 1$,就是说,在该城市碰到一个肝炎患者,其验血指标必定是阳性。

这是两个事件没有特定的时间先后,可以分别互为条件的例子。

4. 全概率公式(Law of Total Probability)

若随机实验 E 中的一组事件 A_1, A_2, \dots, A_n , 满足:① $A_i \cap A_k = \emptyset, i \neq k$, 即这些事件之间互不相容(互不相交),② $A_1 \cup A_2 \cup \dots \cup A_n = S$, 即这些事件覆盖了整个样本空间 S ;那么,就称这组事件是样本空间 S 的一个划分。

这样,任意事件 B ,也就可以被 A_1, A_2, \dots, A_n 划分成了互不相交(互不相容)的子事件: BA_1, BA_2, \dots, BA_n 。于是,由概率的性质,有

$$\begin{aligned} P(B) &= P(BA_1) + P(BA_2) + \dots + P(BA_n) \\ &= P\{B|A_1\}P(A_1) + P\{B|A_2\}P(A_2) + \dots + P\{B|A_n\}P(A_n) \end{aligned}$$

$$= \sum_{k=1}^n P(B | A_k)P(A_k)$$

这就是所谓全概率公式。

为了便于理解全概率公式,可参看图 1.1.2。

显然, A 与 \bar{A} 是 S 的一个最简单的划分,于是,对 S 中的任意事件 B ,有

$$\begin{aligned} P(B) &= P(AB) + P(\bar{A}B) \\ &= P\{B | A\}P(A) + P\{B | \bar{A}\}P(\bar{A}) \end{aligned}$$

其实,在例 1.1.1 中我们已经用过这个公式。

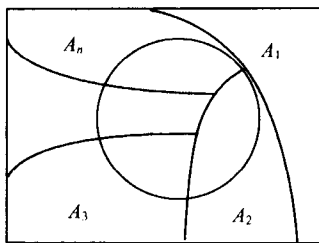


图 1.1.2 样本空间 S 的一个划分与全概率公式示意

全概率公式

若 A_1, A_2, \dots, A_n 是对样本空间 S 的一个划分,则对 S 中的任意事件 B ,有全概率公式如下:

$$P(B) = \sum_{k=1}^n P(BA_k) = \sum_{k=1}^n P(B | A_k)P(A_k)$$

公式的第一个等号的两端,在条件概率计算表中,就是 B 所在的列(或行)所有非“总计”数据之和(也就是 B 所在列或行的“总计”数据),除以表格的右下角数据。

5. 条件概率表的一般表达

设一组事件 A_1, A_2, \dots, A_n 是样本空间 S 的一个划分,另一组事件 B_1, B_2, \dots, B_m 也是样本空间 S 的一个划分,那么条件概率计算表的一般格式如表 1.1.3 所示。

表 1.1.3 条件概率计算表的一般格式

	B_1	B_2	...	B_m	总计
A_1	$F(A_1B_1)$	$F(A_1B_2)$...	$F(A_1B_m)$	ΣA_1 行
A_2	$F(A_2B_1)$	$F(A_2B_2)$...	$F(A_2B_m)$	ΣA_2 行
\vdots	\vdots	\vdots		\vdots	\vdots
A_n	$F(A_nB_1)$	$F(A_nB_2)$...	$F(A_nB_m)$	ΣA_n 行
总计	ΣB_1 列	ΣB_2 列	...	ΣB_m 列	$\Sigma \Sigma$

要点 1:表头的事件(无论在上表头还是左表头)都必须是样本空间 S 的一个划分。

要点 2:表中的 $F(A_iB_j)$ 表示对事件 A_i 与事件 B_j 同时发生的数量描述。该数量描述既可以是概率描述,也可以不是概率描述,还可以部分是概率描述,但必