

Bayesian Statistics with R

贝叶斯统计 及其R实现

黄长全◎编著



清华大学出版社



经济学系列
经济管理精品教材

21
世纪

Bayesian Statistics with R

贝叶斯统计及其R实现

黄长全◎编著



清华大学出版社
北京

内 容 简 介

贝叶斯统计学是现代统计学中非常有特色的内容,应用范围极其广泛。本书系统地介绍了贝叶斯统计的基本思想及其来龙去脉、先验分布和后验分布的概念以及寻求方法,贝叶斯统计推断, MCMC 计算方法以及统计决策理论等。为使初学者更好地理解贝叶斯统计并培养起对贝叶斯统计的兴趣,本书引入了丰富的案例,涉及经济、管理、天文、医药、生物、体育等领域,本书专门制作了一个专用 R 软件包,把书中所有案例数据和主要程序都放入了此压缩包中,增强了师生教学与互动的效果,以便激发初学者对贝叶斯统计的兴趣,掌握贝叶斯统计的精髓,为贝叶斯统计的应用打好基础。

本书可作为高等院校统计、经济、金融、管理、医药、生物等专业高年级本科生和研究生贝叶斯统计课程的教材,也可作为对贝叶斯统计感兴趣人士的参考用书。

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

贝叶斯统计及其 R 实现 / 黄长全编著. —北京:清华大学出版社,2017
(21世纪经济管理精品教材·经济学系列)
ISBN 978-7-302-46785-4

I. ①贝… II. ①黄… III. ①贝叶斯统计量 IV. ①O212.8

中国版本图书馆 CIP 数据核字(2017)第 052773 号

责任编辑: 吴雷

封面设计: 李召霞

责任校对: 王凤芝

责任印制: 王静怡

出版发行: 清华大学出版社

<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: 三河市金元印装有限公司

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 10.75 字 数: 239 千字

版 次: 2017 年 5 月第 1 版 印 次: 2017 年 5 月第 1 次印刷

印 数: 1~3000

定 价: 35.00 元

产品编号: 073181-01

前 言 R

贝叶斯统计学是现代统计学中重要而独特的部分,不仅在统计学本身而且在众多其他学科中也有重要应用。近二十多年来,有关贝叶斯统计本身和贝叶斯统计应用的论文频频出现在各类统计以及非统计刊物上,贝叶斯统计解决了大量经典统计难以解决的复杂问题。可以这么说,没有学习过贝叶斯统计,就不能说了解过现代统计学。因此,贝叶斯统计理应成为大学统计类专业的一门必修课。

厦门大学经济学院统计系(原计划统计系)于2003年第一次正式开设了贝叶斯统计学课程,从那时起,我就一直担任该课程的主讲教师。光阴荏苒、白驹过隙,十多年的时间一晃就过去了。这十多年来,如何教好这门在统计学中独一无二的课程一直是萦绕在我脑海中挥之不去的一个问题,在此期间我既有教训也积累了不少教学经验。因此,在几年前我就萌发了用自己的教学经验和教学观点撰写一本有些许自己风格的贝叶斯统计教科书的念头。

有了撰写教材的想法后,自然而然地就会考虑:如何写出一本有特色的好教材呢?一本好教材的标准又是什么呢?我想就统计教学而言,一本好教材绝不仅仅是教给学生一些统计知识,更重要的是要培养和激发学生对统计学的兴趣和热爱,因为兴趣是最好的老师。那么怎样培养和激发学生对统计学的兴趣呢?多年的统计学科的教学经历使我认识到,要培养和激发学生对统计学的兴趣,一定要首先培养学生的“数据感”。众所周知,球类运动员要培养“球感”,语言学习者要培养“语感”,这些对他们而言都是极为重要的练习过程。对于统计专业以及任何学习统计的学生来说,在学习过程中培养自身的数据感同样极为重要。有了良好的数据感,才会对统计产生亲切感,从而才能激发起自身对统计的兴趣,这实际上也是专业素质的培养。如果大学本科四年不能培养起学生良好的数据感,就不能说是成功的本科统计教育。基于这种教学认识,本书以培养学生的数据感和激发学生的学习兴趣为写作方向。为了使本教材充满统计意味,我们从一开始就介绍贝叶斯统计学的最新有趣应用,同时,全书的案例丰富多彩,涉及经济、管理、天文、医药、生物、体育等领域,也有和日常生活息息相关的例子,使学生觉得贝叶斯统计不再是枯燥无味的,而是既有用又富有生活气息的。本书也专门制作了一个专用R软件包,把书中所有案例数据和主要程序都放入了此压缩包中,增强了师生之间的互动效果。此外,R软件的使用贯穿全书,目的就是通过数据和实际案例分析,加深学生对理论的理解并培养学生良好的

数据感，强化学生的动手操作能力。

本书共七章内容：第1章从一个贝叶斯统计学的真实应用开始，介绍贝叶斯统计的基本概念和公式，概述贝叶斯统计学的历史和发展趋势以及与经典统计学的比较；第2章引入共轭先验和充分统计量等概念，初步讨论后验分布的寻求以及共轭先验下的后验分布特性；第3章介绍先验分布的重要性和一系列先验分布的寻求方法，包括杰弗里斯先验等；第4章研究贝叶斯统计推断理论并介绍了贝叶斯统计在一系列不同领域的应用案例；第5章讨论贝叶斯统计决策理论，引入决策函数等一系列概念；第6章从实用的角度介绍了马尔可夫链蒙特卡罗(MCMC)方法的思想和简史以及马氏链样本的收敛检验问题；第7章则简要讨论统计决策理论，包括贝叶斯风险准则与后验风险准则的等价性等问题。另外，本书附带有R软件包、课件、部分习题参考答案，读者可通过扫描书中的二维码，联系出版社进行下载学习。

本书可作为高等院校统计、经济、金融、管理、医药、生物等专业高年级本科生和研究生的贝叶斯统计课程的教材或参考书。关于教学内容建议：对本科生而言，讲授前五章的全部内容，可加选讲第6、7章；对于研究生则应讲授全部七章的内容。

本书得以出版要感谢清华大学出版社；感谢吴雷编辑，他在组织出版的过程中做了大量的工作。此外，本书的初稿在厦门大学经济学院统计系和王亚南经济研究院双学位课程班讲授过，所以也要感谢各位学习这门课程的同学，是他们的认真学习，触动了我去思考如何教好这门课程。

坦率地说，撰写教材是一件吃力不讨好的工作。但我认为撰写教材是教师的职责之一，当一名教师在某门课程上认真教学了多年，有了教学上的经验与教训，那么就应该把它写出来。最后，本书若能激发读者对贝叶斯统计的兴趣，有助于读者学习贝叶斯统计，那将是对笔者最大的慰藉。当然，由于自身学识所限，本书一定存在许多不足和错误之处，恳望读者朋友指正。

黄长全
2017年1月于厦门大学
Email:cqhuang@xmu.edu.cn



前言 i

第①章 贝叶斯统计基本概念 001

1.1 引言	001
1.1.1 一个美国书呆子的故事	001
1.1.2 贝叶斯统计简史	001
1.1.3 经典统计方法	002
1.1.4 贝叶斯统计方法	003
1.2 概率空间与随机事件贝叶斯公式	003
1.2.1 概率空间与随机事件贝叶斯公式	003
1.2.2 两例：她怀孕了吗？“非典”时期病人为何要测量体温？	004
1.2.3 案例：自动语音识别——神奇的语音输入法	006
1.3 三种信息与先验分布	007
1.3.1 总体与总体信息	007
1.3.2 样本信息	007
1.3.3 先验信息与先验分布	008
1.4 一般形式的贝叶斯公式与后验分布	009
1.4.1 知识准备	009
1.4.2 R 语言与 R 软件包	010
1.4.3 一般形式的贝叶斯公式	011
1.4.4 计算后验分布示例	012
本章要点小结	014
思考与练习	014

第②章 共轭先验分布与充分统计量 016

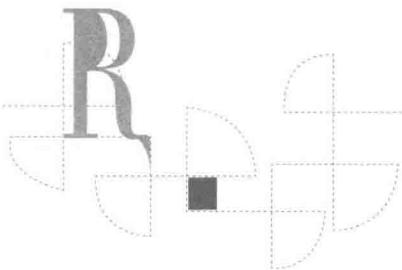
2.1 共轭先验分布	016
------------------	-----

2.1.1	后验分布的核	016
2.1.2	共轭先验分布	017
2.2	多参数先验与后验分布	021
2.2.1	联合先(后)验密度函数	021
2.2.2	多参数共轭先验示例	021
2.3	充分统计量与应用	023
2.3.1	充分统计量概念	024
2.3.2	充分统计量示例	024
本章要点小结		026
思考与练习		027
第3章	先验分布寻求方法	028
3.1	先验分布类型已知时超参数估计	028
3.2	由边际分布确定先验分布	031
3.2.1	混合分布与混合样本	032
3.2.2	寻求先验密度的II型最大似然法	033
3.2.3	寻求先验密度的边际矩法	034
3.3	用主观概率作为先验概率	036
3.3.1	为什么需要主观概率	036
3.3.2	确定主观概率的方法	037
3.4	无信息先验分布	038
3.4.1	非正常先验与贝叶斯假设	039
3.4.2	位置参数的无信息先验	041
3.4.3	尺度参数的无信息先验	043
3.4.4	杰弗里斯先验	044
本章要点小结		048
思考与练习		049
第4章	贝叶斯统计推断	051
4.1	贝叶斯估计	051
4.1.1	点估计	051
4.1.2	贝叶斯估计优良性准则	052
4.1.3	区间估计	053

4.2 泊松分布参数的估计	055
4.2.1 后验分布	055
4.2.2 参数估计	055
4.2.3 案例：受教育程度不同的妇女生育率相同吗？	055
4.3 指数分布参数的估计	057
4.3.1 参数估计	057
4.3.2 案例：国产彩电的寿命有多长？	057
4.4 正态分布参数的估计	059
4.4.1 方差已知时均值的估计	059
4.4.2 均值已知时方差的估计	059
4.4.3 均值和方差的同时估计	060
4.4.4 案例：无先验信息如何估计马拉松成绩分布的参数	062
4.5 贝叶斯假设检验	063
4.5.1 贝叶斯假设检验与贝叶斯因子	064
4.5.2 简单假设对简单假设	065
4.5.3 复杂假设对复杂假设	066
4.5.4 简单假设对复杂假设	068
4.5.5 案例：哪个疗效更好？	069
4.6 模型的比较与选择	071
4.6.1 模型比较与选择	071
4.6.2 案例：足球队进球数量的分布是什么？	074
4.7 统计预测	075
4.7.1 预测原理	075
4.7.2 统计预测示例	076
本章要点小结	078
思考与练习	078
第 5 章 决策概念与贝叶斯决策	080
5.1 决策基本概念	080
5.1.1 决策问题三要素	080
5.1.2 行动的容许性与先验期望准则	083
5.1.3 先验期望准则两性质	085
5.2 损失函数	086
5.2.1 什么是损失函数	086

5.2.2 损失函数下的先验期望准则	087
5.2.3 二行动线性决策问题的损失函数	089
5.3 贝叶斯决策	090
5.3.1 什么是贝叶斯决策	090
5.3.2 决策函数	091
5.3.3 后验风险与后验风险准则	094
5.3.4 常用损失函数下的贝叶斯估计	097
5.3.5 贝叶斯决策下的假设检验	101
5.4 抽样的价值	103
5.4.1 完全信息期望值	103
5.4.2 抽样信息期望值	105
5.4.3 最佳样本量的确定	107
本章要点小结	110
思考与练习	111
第 6 章 贝叶斯统计计算方法	114
6.1 什么是 MCMC 方法	114
6.1.1 蒙特卡罗法	114
6.1.2 马尔可夫链	117
6.1.3 马氏链蒙特卡罗法	120
6.2 吉布斯抽样	121
6.2.1 二阶段吉布斯抽样	121
6.2.2 多阶段吉布斯抽样	125
6.3 梅切波利斯-哈斯廷斯算法	127
6.4 MCMC 的收敛性问题	130
本章要点小结	137
思考与练习	137
第 7 章 统计决策概要	139
7.1 风险函数	139
7.1.1 风险函数与一致最优决策函数	139
7.1.2 统计决策框架中的经典推断	140

7.2 决策函数的容许性与最小最大准则	142
7.2.1 容许性	142
7.2.2 最小最大准则	144
7.3 贝叶斯风险准则与贝叶斯解	146
7.3.1 贝叶斯风险准则	146
7.3.2 贝叶斯解的性质	149
本章要点小结	153
思考与练习	154
拓展资源:专用 R 软件包	155
附录 常用概率分布表	156
参考文献	159



第 1 章

贝叶斯统计基本概念

俗话说，万事开头难。为了提高读者的学习兴趣，本章从一个贝叶斯统计的真实应用开始，介绍贝叶斯统计的基本概念和公式，概述贝叶斯统计学的历史和发展趋势以及与经典统计学的比较。

1.1 引言

1.1.1 一个美国书呆子的故事

在 2012 年美国总统大选期间，一个一直都被称作“书呆子”的美国人纳特·西尔弗 (Nate Silver, 生于 1978 年 1 月 13 日) 用以统计为主要工具的模型准确预测了美国全部 50 个州的选举结果。在大选日当天早晨，他的模型最新预测到时任总统巴拉克·奥巴马 (Barack Obama) 将有 90.9% 的可能获得多数选举人票从而连任，而选举结果确确实实就是奥巴马总统赢得了这次美国总统大选。于是，他凭借自己的模型及其准确的预测打败了所有时事政治记者、政党媒体顾问和政治评论员。“你们知道谁是今晚(大选日当夜)的赢家吗？”美国全国广播公司新闻节目主播自问自答，“是纳特·西尔弗”。其实，早在 2008 年的美国总统大选期间，西尔弗就准确预测了整个美国 50 个州中 49 个州的选举结果。两次极为准确的预测，让这个“书呆子”扬眉吐气、名声大震，各种荣誉接踵而来，甚至于被四所大学授予了四个荣誉博士学位，当然这也让我们从事统计领域的人士大感骄傲。西尔弗的预测模型有什么神秘之处呢？答案就是其利用了大数据和我们将要学习的贝叶斯统计理论和方法。

1.1.2 贝叶斯统计简史

贝叶斯统计学是以英国人托马斯·贝叶斯 (Thomas Bayes, 1702—1761) 的名字命名

的。贝叶斯是一位英国牧师,但他却热衷于概率统计等科学研究,还是英国皇家学会会员。遗憾的是,现在人们对他的生平却知之甚少,甚至没有人知道贝叶斯的相貌如何,现存所有他的画像都是传说,并不能证实是他的真容。贝叶斯统计学起源于贝叶斯逝世后公开发表的一篇论文——《论一个概率理论问题的求解》(*An Essay Towards Solving a Problem in the Doctrine of Chances*)。在贝叶斯去世两年之后,这篇论文由他的朋友理查德·普莱斯(Richard Price)介绍到英国皇家学会,引起了该学会的注意和讨论,并于1763年发表在《皇家学会哲学会刊》上。在该篇论文中,贝叶斯首次提出了贝叶斯统计的基本思想和归纳推理方法。

五十一年后,法国数学、统计学、天文学和物理学家拉普拉斯(P. S. Laplace, 1749—1827)在1814年出版了著作《关于概率的哲学评述》(*A Philosophical Essay on Probabilities*),在该著作中他将贝叶斯提出的公式进行了推广并导出了一些很有意义的新结果。然而,之后相当长的一段时间里虽然有一些理论和应用研究,但由于其理论与经典统计学相比显得另类,而且人们对它的理解还不够深刻,在应用上其计算复杂且计算量巨大,因此贝叶斯统计理论和方法长期未被普遍接受,甚至被一些学者看作一种旁门左道。直到20世纪中叶开始,有一批统计学家,例如杰弗里斯(H. Jeffreys, 1939)、萨维奇(L. J. Savage, 1954)、雷法和施莱弗(H. Raiffa and R. Schlaifer, 1961)以及伯杰(J. O. Berger, 1985)等,才对贝叶斯统计做了更加深入的研究,特别是罗马尼亚(匈牙利)裔美国统计学家阿布拉汉·瓦尔德(Abraham Wald, 1939, 1950)通过将损失函数引入统计学并利用决策概念和思想把经典统计推断纳入决策理论框架中而形成了统计决策理论,这样经典统计学和贝叶斯统计学通过决策理论有机地联系到了一起,才得到了很有意义的理论结果。从20世纪中叶开始,在一批学者的努力下,人们对贝叶斯统计在观点、方法和理论上的认识不断加深。从20世纪90年代以来,伴随着计算机科学技术的发展和有效的贝叶斯统计计算方法的发现和应用,贝叶斯统计解决了相当一批经典统计难以解决的实际问题,从而得到了人们极大的重视。现在,贝叶斯理论和方法获得了人们的普遍接受,贝叶斯统计不仅在统计学本身而且在众多学科中都得到了广泛的应用,解决了各个不同学科中大量的复杂统计问题。贝叶斯统计表现出了勃勃生机和欣欣向荣的景象,在统计学领域牢牢地站稳了一席之地,也成为现代统计学的重要分支,可以这么说,没有学习过贝叶斯统计,就不能说了解过现代统计学。

1.1.3 经典统计方法

我们先来回顾一下经典统计学的思想方法,以便与下一小节的贝叶斯统计思想方法进行比较。回顾一下概率统计课程中概率的定义,便容易明白经典统计学思想方法也就是“频率方法”,它把概率定义为频率的极限,也就是说如果随着随机试验重复次数的增多,随机事件发生的频率会稳定在一个常数附近,这个常数就是该随机事件发生的概率。同时,它认为总体的数字特征(如均值、方差)和别的参数仅仅是未知的常数,可以用样本统计量来估计。而且,它又认为样本是随机变量,从而样本统计量也是随机变量,因此具有概率分布,即它的抽样分布。如果统计量的分布可以求出,利用该分布,就可以进行区

间估计和假设检验等统计推断。然而,我们知道寻求统计量的概率分布和进行区间估计以及假设检验等都不是容易的事,而且参数的区间估计既不容易理解也不容易解释。

1.1.4 贝叶斯统计方法

贝叶斯统计学虽然也认可经典统计学的概率定义,但它同时把概率理解为人们对随机事件发生可能性的一种信念(有时被称为“可信度”),当然,这种信念不是信口开河,而是基于学识和经验之上的审慎度量。其次,贝叶斯统计把任意一个未知量(参数)都看作一个随机变量,可用一个概率分布去描述它。我们说这种观点是合理的,因为即使是一个确定性的未知量,也可以把它看成随机变量的特殊情形,即服从0—1分布的随机变量。所以说,任一个未知量都可用一个适当的概率分布去描述它。这个概率分布利用历史数据或其他历史信息或研究人员的经验和学识而确定,称为该未知量(参数)的先验分布。而后利用新样本信息(即抽样信息)对先验分布进行更新,更新之后的这个新概率分布称为该未知量的后验分布。由此,未知参数的点估计、区间估计和假设检验等统计推断都是基于后验分布来进行的,而且参数的区间估计既容易理解也容易解释,假设检验则简单明了。

经典统计学把概率定义为频率的极限,初看起来似乎客观、严谨,但是在现实世界中要进行重复试验需要花费大量的人力、物力,而且有时根本无法重复,例如,我们无法重复昨天的天气和去年的经济活动。因此,用频率的极限来定义概率在实际应用中受到了极大的限制。相反,贝叶斯统计把概率理解为人们对随机事件发生可能性的信念,则在实际应用中没有任何限制,因为它不需要重复,事件甚至可以一次都没有发生。而且,在贝叶斯统计中一旦后验分布建立起来了,所有的统计推断都是基于后验分布来进行的,因此,至少从理论上而言,贝叶斯统计推断比经典统计推断要简单明了得多。当然,现代统计学的发展趋势是,根据实际问题的条件和需要挑选经典统计方法或贝叶斯统计方法,有时甚至是综合利用这两种统计理论和方法进行统计推断。所以,不管是经典统计还是贝叶斯统计,能够解决问题的就是“好统计”!

对于经典统计学与贝叶斯统计学的比较,有待学完本书的内容后才能有更深刻的体会,因此希望读者在研读完本书后,再好好对它们做一个详细的比较分析。

1.2 概率空间与随机事件贝叶斯公式

1.2.1 概率空间与随机事件贝叶斯公式

我们从概率论知道概率空间是三位一体的一个研究对象(Ω, F, P),其中 Ω 是样本点全体,也称为样本空间; F 是事件域(简单说就是所要研究的随机事件全体,包含必然事件 Ω 和不可能事件 Φ); P 是定义在事件域 F 上的概率(测度),满足以下三条公理:

- (1) 非负性:对于任意事件 A ,其概率 $P(A) \geq 0$;
- (2) 规范性:必然事件 Ω 的概率等于1,即 $P(\Omega) = 1$;

(3) 可列可加性: 如 $\{A_i\}_{i=1}^{\infty}$ 是一列事件, 满足 $A_i A_j = \Phi (i \neq j)$ (称为两两互不相容), 则

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = P\left(\sum_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$$

这一公理体系称为柯尔莫哥洛夫概率论公理体系, 是苏联著名数学家柯尔莫哥洛夫于 1933 年建立的, 得到了概率统计学者们的广泛认可, 从而为概率论建立了坚实的理论基础。

另外, 对于任意两个事件 A, B 且 $P(A) > 0$, 定义在 A 发生的条件下, B 发生的条件概率为

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

从而, $P(AB) = P(A)P(B|A)$, 这就是乘法公式。推而广之, 设 $\{A_k\}_{k=1}^n$ 是任意 n 个随机事件, 则有更一般的乘法公式

$$P(A_1 A_2 \cdots A_n) = P(A_1)P(A_2 | A_1)P(A_3 | A_1 A_2) \cdots P(A_n | A_1 A_2 \cdots A_{n-1})$$

现设 $\{A_i\}_{i=1}^{\infty}$ 是事件域 F 中的一列事件, 若 $\bigcup_{i=1}^{\infty} A_i = \Omega$, 且 $A_i A_j = \Phi (i \neq j)$, 则称 $\{A_i\}_{i=1}^{\infty}$ 为 Ω 的一个划分(也称为 Ω 的完全事件组, 这里事件的个数也可以是有限多个, 比如说 n 个, 这相当于 $k > n$ 时都有 $A_k = \Phi$)。显然, 任一个事件 A 与其补 \bar{A} 就是 Ω 的一个划分。现在设 $\{A_i\}_{i=1}^{\infty}$ 为 Ω 的一个划分且 $P(A_i) > 0$, 则对任一个事件 $B \in F$ 有全概率公式

$$P(B) = \sum_{i=1}^{\infty} P(A_i)P(B|A_i)$$

事实上, 由

$$B = B\left(\bigcup_{i=1}^{\infty} A_i\right) = \bigcup_{i=1}^{\infty} (A_i B) \text{ 且 } (A_i B) \cap (A_j B) = (A_i A_j)B = \Phi, i \neq j$$

利用可列可加性及乘法公式就得

$$P(B) = P\left(\bigcup_{i=1}^{\infty} A_i B\right) = \sum_{i=1}^{\infty} P(A_i B) = \sum_{i=1}^{\infty} P(A_i)P(B|A_i)$$

现在将全概率公式以及乘法公式应用到条件概率 $P(A_j | B)$ 的公式上就有

$$P(A_j | B) = \frac{P(A_j B)}{P(B)} = \frac{P(A_j)P(B|A_j)}{\sum_{i=1}^{\infty} P(A_i)P(B|A_i)} \quad j = 1, 2, \dots, n, \dots$$

这就是著名的随机事件形式的贝叶斯公式(定理或法则), 也称为逆概率公式, 这里 $\{A_i\}$ 可以认为是事件 B 发生的所有可能的原因, 而贝叶斯公式就是计算在已知事件 B 发生的条件下每个原因的可能性大小(概率), 也就是说由结果去推测原因, 因此叫逆概率公式。在贝叶斯公式中, $P(A_j)$ 称为 A_j 的先验概率, 因为这是事先已知的, 而 $P(A_j | B)$ 自然称为 A_j 的后验概率。

1.2.2 两例: 她怀孕了吗? “非典”时期病人为何要测量体温?

贝叶斯公式与全概率公式都是概率论中的著名公式, 在许多学科中都有重要应用, 下

下面我们来看两个例子。

例 1.1 (她怀孕了吗?)根据历史资料知道:女性一次性交后怀孕的概率为 15%。假如一个女性某次性交后怀疑自己怀孕了,但又不能确定。于是,她做了个准确率为 90% 的验孕测试,即 90% 的怀孕案例会给出阳性反应的检验结果,同时知道该测试当未怀孕时阳性反应占 10%。她当然想知道在检验结果为阳性的条件下的怀孕概率。然而,她不懂贝叶斯统计,所以请你帮助她算出该概率。

解 已知

$$P(\text{怀孕}) = 0.15, P(\text{检测阳性} | \text{怀孕}) = 0.90, P(\text{检测阳性} | \text{未怀孕}) = 0.10$$

由已知得, $P(\text{未怀孕}) = 0.85$ 。由贝叶斯公式知在检验结果为阳性的条件下的怀孕概率:

$$\begin{aligned} P(\text{怀孕} | \text{检验阳性}) &= \frac{P(\text{检验阳性} | \text{怀孕})P(\text{怀孕})}{P(\text{检验阳性} | \text{怀孕})P(\text{怀孕}) + P(\text{检验阳性} | \text{未怀孕})P(\text{未怀孕})} \\ &= \frac{0.90 \times 0.15}{0.90 \times 0.15 + 0.10 \times 0.85} = \frac{0.135}{0.135 + 0.085} = 0.614 \end{aligned}$$

这里 $P(\text{怀孕}) = 0.15$ 就是怀孕的先验概率, $P(\text{怀孕} | \text{检验阳性}) = 0.614$ 就是怀孕的后验概率, 它是在观察数据(阳性测试)后怀孕概率的更新, 表明如果测验呈阳性, 则怀孕的可能性大大提高。

例 1.2 (非典时期病人为何要测量体温?)“非典(SARS)”患者的主要病症表现为发热、干咳。根据某地区历史资料,已知人群中既发热又干咳的病人患“非典”的概率为 5%;仅发热的病人患“非典”的概率为 3%;仅干咳的病人患“非典”的概率为 1%;无上述病症而患“非典”的概率为 0.01%;现对该区 25 000 人进行检查,发现其中既发热又干咳的病人有 250 人,仅发热的病人为 500 人,仅干咳的病人为 1 000 人,试求:

- (1) 该地区中某人患“非典”的概率;
- (2) “非典”患者是仅发热的病人的概率。

解 引入记号

$$A = \{\text{既发热又干咳的病人}\}, B = \{\text{仅发热的病人}\},$$

$$C = \{\text{仅干咳的病人}\}, D = \{\text{无明显症状的人}\},$$

$$E = \{\text{“非典”患者}\}$$

易知 A, B, C, D 构成了一个划分。根据对该区 25 000 人进行检查的结果, 有

$$P(A) = \frac{250}{25\,000}, P(B) = \frac{500}{25\,000}, P(C) = \frac{1\,000}{25\,000},$$

$$P(D) = \frac{25\,000 - (250 + 500 + 1\,000)}{25\,000} = \frac{23\,250}{25\,000}$$

由全概率公式得患“非典”的概率:

$$\begin{aligned} P(E) &= P(A)P(E|A) + P(B)P(E|B) + P(C)P(E|C) + P(D)P(E|D) \\ &= \frac{250}{25\,000} \times 5\% + \frac{500}{25\,000} \times 3\% + \frac{1\,000}{25\,000} \times 1\% + \frac{23\,250}{25\,000} \times 0.01\% = 0.001\,593 \end{aligned}$$

由贝叶斯公式知,“非典”患者是仅发热的病人的概率:

$$P(B|E) = \frac{P(B)P(E|B)}{P(E)} = \frac{\frac{500}{25\,000} \times 3\%}{0.001\,593} = 0.376\,647\,8$$

同理,可以算出“非典”患者是既发热又干咳、仅干咳、无明显症状的病人的概率分别为

$$P(A|E) = \frac{P(A)P(E|A)}{P(E)} = \frac{\frac{250}{25\,000} \times 5\%}{0.001\,593} = 0.313\,873\,2$$

$$P(C|E) = \frac{P(C)P(E|C)}{P(E)} = \frac{\frac{1\,000}{25\,000} \times 1\%}{0.001\,593} = 0.251\,098\,6$$

$$P(D|E) = \frac{P(D)P(E|D)}{P(E)} = \frac{\frac{23\,250}{25\,000} \times 0.01\%}{0.001\,593} = 0.058\,380\,41$$

不难看出

$$P(A|E) + P(B|E) + P(C|E) + P(D|E) = 1$$

而一个人患“非典”时最可能的症状是发热。这就是为什么在非典时期要测量病人体温的原因。

1.2.3 案例:自动语音识别——神奇的语音输入法

你的手机里安装了讯飞语音输入法或其他语音输入法了吗?是不是觉得它很神奇呢?想不想知道它为什么能够把你的话转换为文字呢?这个转换过程其实就是自动语音识别。简单地说,自动语音识别是指由机器自动将语音信号转换为文字的方法和过程。人类的语言可以说是各种信息里最复杂和最动态的一种,著名语言学家乔姆斯基(A. N. Chomsky)和信息论的祖师爷香农(C. Shannon)等学者都关注过自动语音识别问题,然而那时自动语音识别并没有获得很大进展。在这个领域率先取得突破的是捷克裔美国语音和语言处理大师贾里尼克(F. Jelinek)。从20世纪60年代开始,贾里尼克开创性地将语音识别问题看成一个通信问题,认为语音识别就是根据接收到的信号序列推測说话人实际发出的信号序列(即说的话)和要表达的意思,并且用贝叶斯公式和两个隐含马尔可夫模型建立起统计语音识别系统,把对应的一套模型称为声学模型和语言模型,从而极大地改变了这一领域的研究方向。此外,他还与其他合作者提出了数字通信领域最重要的算法之一——BCJR(L. R. Bahl, J. Cocke, F. Jelinek, J. Raviv, 1974)算法。难能可贵的是,这种统计语音识别系统不但能够识别静态的词库里的语音,而且对动态变化的词库语音具有很好的适应性,即对新出现的词汇,只要这个词已经被高频使用,可用于训练的数据量足够多,系统就能通过训练而正确地识别之。这实际上表明贝叶斯公式对新词汇语音信息有非常好的适应能力。由于本书的性质,这里我们不可能对问题展开详细的讨论,有兴趣者可以去研读有关文献资料。但我们从已经开发出来的语音输入法知道这种统计语音识别系统是非常成功的!

1.3 三种信息与先验分布

在1.1节中,我们初步了解到统计学中有两个主要学派:经典统计学派与贝叶斯统计学派。在本节我们将从这两个学派使用的信息种类来讨论它们之间的异同。首先我们来了解统计推断问题中存在的三种信息。

1.3.1 总体与总体信息

我们从已学课程知道统计学中总体就是根据一定的目的和要求所确定的研究对象的全体。例如,如果要统计调查全国大学男生的身高,那么,我们就可以把全国大学男生的集合作为总体,而大学男生身高这个指标就是关于该总体的一个数量,可以用一个符号 X 来标记它。由于在对随机抽出的一个大学男生具体测量之前,并不知道该大学男生的确切身高,而且人的身高是受遗传、营养等随机因素影响而确定的,所以 X 是一个随机变量并且服从某种概率分布。再如,我们要考察一个经济指标 Y (可以把它设想为某一只股票的收益率或一个国家的GDP),由于受各种各样的随机因素的影响, Y 是一个随机变量,它的所有可能取值就构成了一个总体并且也服从某一种概率分布。由于一个随机变量的概率分布完全刻画了该随机变量的统计规律性,因此,我们实际上甚至可以抽象地把这个随机变量的概率分布看作总体。总体信息就是我们对总体概率分布的了解或知识,一般而言,对总体信息最大的了解是知道总体概率分布所属的分布族,例如,若我们知道总体服从正态分布族 $N(\mu, \sigma^2)$,虽然这时两个参数还是未知的,我们也知道它的密度函数是一条关于总体均值对称的钟形曲线并且它的各阶矩都存在,同时也知道第一个参数 μ 是分布的均值,第二个参数 σ^2 是分布的方差。当然,总体到底服从怎样的概率分布族对一个新研究问题而言通常不得而知,这正是统计学的一个分支——非参数统计所要研究的。显而易见,要获得总体信息往往必须投入大量的人力、物力,例如,美国军队为了获得某种新的电子元件的寿命分布,购买了上万个此种电子元件,做了大量的寿命实验,在获得大量数据后才确认其寿命概率分布是什么。简言之,总体信息非常重要,要获得它虽然不容易但又是必须要做的,因为它是统计推断的基础。

1.3.2 样本信息

为了对所研究的总体有更多的了解,我们必须从总体抽取(观察或收集)一定的样本 $x=(x_1, x_2, \dots, x_n)$,这些样本给我们提供的信息就是样本信息,也称为抽样信息。样本信息两种最重要的表现形式是样本的联合分布与样本统计量的抽样分布,其次是样本对总体特征的各种估计,例如,样本均值、样本方差(标准差)等。样本是统计学(无论是频率学派还是贝叶斯学派)的粮食,没有样本就如同巧妇难为无米之炊一样,做不成统计学上的任何事情,也就没有统计学了。

仅仅基于总体信息和样本信息进行统计推断的统计学理论和方法称为经典统计学。它的历史悠久,但大发展却是从19世纪末到20世纪上半叶。由于统计学家皮尔逊