

研究生教学用书

教育部研究生工作办公室推荐

信息论与编码

Information Theory & Coding

(第二版)

姜 丹 编著

中国科学技术大学出版社

研究生教学用书

教育部研究生工作办公室推荐

信息论与编码

Information Theory & Coding

SECOND EDITION

(第二版)

姜 澍 编著

中国科学技术大学出版社

内 容 简 介

本书系统论述香农信息论的基本理论,介绍编码的基本方法.全书共分九章.内容包括:信息的定义、信息论的基本思路;单符号离散信源与信道、信息熵、互信息、信道容量、数据处理定理、加权熵、效用信息熵;多符号离散信源与信道、极限熵、独立并列信道的信道容量;连续信源与信道、相对熵、高斯白噪声加性信道的信道容量;无失真信源编码定理、霍夫曼编码方法;抗干扰信道编码定理、线性分组码、汉明码与扩展汉明码;限失真信源编码定理、信息率失真函数、数据压缩原理、信息价值、广义信息率失真函数;网络信息理论等.

本书可作为高等院校、科研院所相关专业的研究生、高年级本科生的教材或教学参考书,也可供从事信息理论、信息技术和信息科学的教学、科研和工程技术人员参考.

图书在版编目(CIP)数据

信息论与编码/姜丹 编著. —2版. —合肥:中国科学技术大学出版社,2004.7

ISBN 7-312-01693-6

I. 信… II. 姜… III. ①信息论 ②信源编码-编码理论 ③信道编码-编码理论 IV. TN911.2

中国版本图书馆 CIP 数据核字(2004)第 048238 号

中国科学技术大学出版社 出版发行

(安徽省合肥市金寨路 96 号,邮编:230026)

合肥学苑印务有限公司印刷

全国新华书店经销

开本:787×960/16 印张:44.75 字数:1110千

2004年8月第2版 2004年8月第3次印刷

印数:10001—15000册

ISBN 7-312-01693-6/TN·57 定价:55.00元

前 言

随着科学技术,特别是信息技术的发展,信息理论在通信领域中发挥越来越重要的作用,显示出解决通信领域中有关问题的有力工具的本色.同时,由于信息理论解决问题的思路和方法的独特、新颖和有效,在当今信息时代,信息理论已渗透到其它相关的自然科学,甚至社会科学领域,与电子技术、自动控制、计算机网络以及管理科学、生物医学工程、遗传工程、人工智能、心理学等学科密切结合,显示出它的勃勃生机和不可估量的发展前景.信息论是信息科学中最成熟、最完整、最系统的重要组成部分,它是信息科学的发展起源与基石.

本书以香农(Claude E. Shannon)信息论为基础,论述近代信息理论的基本概念和主要结论.

作者鉴于 20 余年的教学经验,为了便于读者正确认识通信领域中信息的定义和本质,理解信息论解决问题的思路和方法,在“引言”中归纳、提炼出香农信息论的三大理论支柱.为了便于读者建立信息流通的完整系统概念,把信息论的基础理论部分由传统的“信源一条线”、“信道一条线”的“纵向结构”,改变成由“单符号离散通信系统”(第一章、第二章)、“多符号离散通信系统”(第三章)、“单维连续通信系统”(第四章)、“多维连续通信系统”(第五章)等四个“横向教学板块”组成的“横向结构”,由简单到复杂、由浅入深、循序渐进地安排教学内容.信息论是一门具有严密的数学演绎系统和高度抽象性、概括性的科学理论.为了帮助读者排除学习信息论过程中经常遇到的数学分析方面的困难,结合有关内容,系统而简明地介绍必要的数学基础知识,给出导致重要结论的数学推演过程,提供不同的证明方法和途径.为了帮助读者正确理解有关结论的物理含意,提供通俗易懂、富有哲理的诠释.

本书按照理论联系实际的原则,在全面系统地论述信息论的基础理论的基础上,进而严密论证了“无失真信源编码定理”、“抗干扰信道编码定理”、“限失真信源编码定理”等信息论中的三大定理,介绍“霍夫曼(Huffman)码”、“线性分组码”等无失真信源编码、抗干扰信道编码的实际编码方法,阐明信息率失真理论在限失真信源编码、实施“数据压缩”方面的应用.使读者既能看到实现有效而可靠的通信系统的光明前景,又能掌握某些实现通信系统“最优化”的实际编码方法.

本书对如何构建“加权熵”、“效用信息熵”;如何运用信息率失真理论定义“信息价值”;如何凝炼“信息率失真函数”的数学精髓,构建“广义信息率失真函数”,估算通信系统的有关指标界限等问题,进行了探索性的讨论.以“多用户信道”的容量界限为重点,对网络信息传输的有关特性,作了初步探讨,给读者提供探究当今正在蓬勃兴起的互联网通信理论的初步基础知识.

本书既适用于信息论初学者,也有助于已具信息论初步知识的读者在更高层次上对信息理论的研究和应用.它可作为高等院校、科研院所的研究生、高年级本科生的教材或教学参考书,也可供从事信息理论、信息技术和信息科学的教学、科研以及工程技术人员参考.

作者在撰写本书过程中,得到中国科学院电子学研究所陈宗骥教授、中国科学院研究生院钱玉美教授的热情指导和帮助,在此一并表示衷心感谢.

热忱希望广大读者对书中的错误和不当之处予以批评指正.

姜 丹

2000年11月于北京

再版前言

作者编著的《信息论与编码》一书,于2001年由中国科学技术大学出版社首次出版发行,并于2003年入选教育部研究生工作办公室推荐的“研究生教学用书”。

作者鉴于在研究生教学实践中的体会和经验,对首次出版的《信息论与编码》部分章、节作了修订,进一步充实、完善了教材的内容和结构。现以《信息论与编码》(第二版),由中国科学技术大学出版社再版发行。

真诚欢迎广大读者对书中的错误和不当之处予以批评指正。

姜 丹

2004年2月于北京

引 言

信息论是人们在长期通信实践活动中,由通信技术与概率论、随机过程、数理统计等学科相结合而逐步发展起来的一门新兴交叉学科.美国科学家香农(C. E. Shannon)于1948年发表的著名论文《通信的数学理论》,奠定了信息论的理论基础.

随着科学技术,特别是信息技术的迅猛发展,在当今信息时代中,“信息”这个词逐渐被人们所接受,并得到越来越广泛的应用.当人们收到由电话、传真、电子邮件、广播、电视等媒体传播的消息后,往往说成获得了“信息”.一般以数字、数据、图表、曲线等形式出现的计算机通信、运算和处理所需要的条件、内容和结果,人们也习惯称之为“输入信息”、“输出信息”、“补充信息”、“反馈信息”等等.人们通常也把由眼、耳、鼻、口等感觉器官直接感觉到的颜色、声音、气味、味道、气温、湿度等说成是感知到了某种外界环境“信息”.确实,在人类社会中,信息的传递与交换,是时时处处都发生着的事情.无数看不见的信息溪流,昼夜川流不息,聚成汹涌澎湃的信息洪涛,汇成浩瀚无垠、波澜壮阔的信息海洋.人类就生活在这样一个信息的海洋之中.

通信是人类活动中最为普遍的现象之一.在频繁的信息传递和交换中,人们总是希望有效、可靠地传递信息.那么,自然而然地会提出这样一个问题,用什么标准来衡量通信的有效程度和可靠程度?怎样判断通信方法的优劣?显然,解决这些问题的关键在于解决信息的度量问题.例如,某人收到一封来信,谈的是同学最近的工作、学习情况.同时又收到一封家信,谈的是家人的健康情况.显然他从这两封信中都获得了信息.但若问:他从哪一封信中获得了更多的信息?也许,按某种想当然的感觉,他会给出某种模糊的回答,如“家信中得到了更多的信息”.这个结论可靠吗?就算这个结论不错,如进一步问:“家信中含有的信息比同学来信中含有的信息多了多少?”一般来说,很难回答这个问题.

为什么很难回答以上的问题呢?其主要原因在于对信息的本质缺乏明确的认识,经验性地把“信息”与“消息”混为一谈.

众所周知,消息是用文字、符号、数据、语言、音符、图片、图像等能被人们的感觉器官所感知的形式,对客观物质运动和主观思维活动状态的一种表述.不同的消息,不仅有不同的形式,而且含有不同的语义和不同的语用.例如,“中国男子体操队获二十七届奥运会团体冠军”这条消息.在形式上,可把它看作是从汉字表中挑选19个字的一种选择.在语义上,这条消息含有多个语义层次:是“中国”,而不是“美国”、“俄罗斯”等其它国家;是“男子”,而不是“女子”;是“体操队”,而不是“足球队”、“篮球队”、“排球队”等其它代表队;是“二十七届”,而不是“二十六届”、“二十五届”等其它届;是“奥运会”,而不是“全运会”、“亚运会”等其它运动会;是“团体”,而不是“个人”、“双人”等其它项目;是“冠军”,而不是“亚军”、“季军”等等,所以从语义上来分析就显得相当复杂.从语用上来说,这条消息对中国人和俄罗斯等外国人所引起的反响程度和语用效果显然是大不相同的.

要解决信息的度量问题,必然要应用数学工具,进行数量的运算.我们知道,数学是刻画物质运

动形式的工具,用数学对消息的形式进行刻画,不存在法则上的困难.但如何运用数学工具刻画消息含有的语义乃至语用,至今仍然是一个巨大的难题.所以,如把信息与消息混为一谈,把形式、语义、语用三个因素交织在一起,综合地解决信息度量问题,必然面临头绪纷繁、无从下手的僵局.

美国科学家香农针对人类通信活动的特点,精辟地提出了“形式化假说”、“非决定论”、“不确定性”等三个论点,以新颖的思想和方法,打破了这个僵局,跨出了用数学方法定量描述信息的关键一步,开创了通信领域信息理论新局面.

(一)形式化假说

通过对通信活动的基本功能的观察分析,香农指出,“通信的基本问题,是在消息的接收端精确地或近似地复制发送端所挑选的消息.通常消息是有意义的,即是说,它按某种关系与某些物质或概念的实体联系着.通信的语义方面的问题与工程问题是没有关系的.”这就是说,通信的任务只是在接收端把发送端发出的消息从形式上复制出来,通信工程并不须要对复制出来的消息的语义作任何处理和判断.对消息的语义内容的处理和判断,是接收者自己的事,不是通信工程本身的任务,与通信工程无关.至于消息的效用问题,更应该是接收者自己的感受问题,与传送消息的通信系统无关.例如电视屏幕上出现一则消息,有的观众看了兴高采烈;有的观众看了满腔愤怒,甚至把电视机从楼上掷到楼下;有的观众看了漠不关心,毫无反应.不论不同的观众有什么不同的效用反应,对电视通信工程来说,已经完成了它本身的任务.这就是香农对通信活动的“形式化”假说.

这种通信工程的“形式化”假说,大胆地去掉了消息的语义、语用因素,巧妙地保留了能用数学描述的形式因素,这使应用数学工具定量度量信息成为可能,打开了信息理论进入科学殿堂的大门.

(二)非决定论

经过对通信活动的对象的分析研究,香农指出,“重要的是,一个实际的消息,总是从可能发生的消息集合中选择出来的.因此,系统必须设计得对每一种选择都能工作,而不是只适合工作于某一种选择.因为,各种消息的选择是随机的,设计者事先无法知道什么时候会选择什么消息来传送.”这就是说,一切有通信意义的消息的发生都是随机的,是事先无法预料的.消息传递过程中遇到的噪声干扰也是随机的,通信系统的工程设计者也是无法事先预料的.面对公众的通信系统,不是针对某一特定的通信对象设计的,什么样的用户,什么时候使用,传递什么样的消息都是无法始料的.显然,根据通信工程的这些特点,必须采用概率论、随机过程、数理统计等数学工具,从大量不可预料的随机消息(包括噪声)中,寻求其统计规律,作为通信工程设计的依据,用非决定论观点揭示信息的本质.这就是香农看待通信活动的“非决定论”观点.

这种“非决定论”观点,是对通信活动的总的认识观,它从原则上回答了应采用什么样的数学工具来解决信息度量的问题.

(三)不确定性

通过对通信活动的机制和作用的剖析研究,香农一针见血地指出“人们只在两种情况下有通信的需要.其一,是自己有某种形式的消息要告知对方,而估计对方“不知道”这个消息;其二,是自己有某种“疑问”要询问对方,而估计对方能作出一定的解答”.这里的所谓“不知道”、“疑问”,就是通

信前对某事件可能发生的若干种结果不能作出明确的判断,存在某种知识上的“不确定性”。通信后,通过消息的传递,由原先的“不知道”到“知道”,或由“知之不多”到“知之甚多”;原先的“疑问”得到了解答,或部分解答,由原先的“疑问”到“明白”,或部分“明白”。这就是说,通信后,消除或部分消除了通信前存在的“不确定性”。所以,通信的作用就是通过消息的传递,使接收者从收到的消息中获取了一样“东西”,因而消除了通信前存在的“不确定性”。这种“东西”,就是“信息”。这样,我们就有理由给“信息”下一个明确的定义:“信息就是用来消除不确定性的东西”,进而,可合理地推断:通信后接收者获取的“信息”,在数量上等于通信前后“不确定性”的消除量。这就是香农从“不确定性”观点出发,给“信息”下的明确的定义。

我们知道,“可能性”的大小在数学上可以用概率的大小来表示:概率大即表示出现的“可能性”大;概率小即表示出现的“可能性”小。我们同样知道,“不确定性”与“可能性”是有联系的:“可能性”大就意味着“不确定性”小;“可能性”小就意味着“不确定性”大。这样,“不确定性”就可与消息发生的概率联系起来。例如,“中国女子乒乓球队夺取亚运会冠军”这条消息,根据中国女子乒乓球队历来的表现,夺取亚运会冠军的概率很大,即可能性很大,也就意味着“不确定性”很小。这个消息一旦发生,消除的不确定性也很小,收信者从这条消息中获取的信息量也很小。相反,“中国男子足球队夺取世界杯赛冠军”这条消息,根据中国男子足球队历来的表现,夺取世界杯赛冠军的概率很小,即“可能性”很小,也就意味着“不确定性”很大。若有朝一日这个消息真的发生了,消除的“不确定性”很大,收信者从这条消息中获取的信息量也很大,甚至惊喜万分、欢呼跳跃。由此可见,“不确定性”与消息发生的概率有内在联系,它应该是消息发生概率的某一函数。

根据香农关于信息的定义,通信后收信者从消息中获取的“信息”,从数量上等于通信前后“不确定性”的消除。既然“不确定性”一定是消息发生概率的某一函数,那么,“不确定性”的“消除量”也一定是消息发生概率的某一函数。当然,通信后获取的信息量也应该是消息发生概率的某一函数。对于随机消息来说,我们虽然不能精确预料它能否发生,但表示随机消息的可能性大小的“概率”一定是一个精确的数量。所以,香农对“信息”的定义,从理论原则上完全解决了信息的度量问题。

香农从“不确定性”观点出发对“信息”的明确定义告诉我们,“信息”与“消息”两者之间既有联系,又有区别,两者不应混为一谈。“消息”是表达“信息”的形式,是载荷“信息”的客体;“信息”是“消息”统计特性的函数,是“消息”的抽象本质。不同形式的“消息”,可能有相同数量的“信息”;相同形式的“消息”,可能有不同数量的“信息”。信息论的研究对象不是具体的消息,而是抽象于各种不同形式的“消息”的“信息”。所以,专门讨论“信息”的产生、传输和处理规律的信息论,是一门具有高度抽象性和概括性的学科。

信息论起源于通信领域,它所讨论的问题局限于通信领域的范围。经过半个多世纪的充实、完善、发展和提高,它已成为信息科学中最完善、最系统、最成熟的重要组成部分,是信息科学发展的起点与基石。随着信息科学和信息技术的迅猛发展,必将显示出它的勃勃生机和光明的前景。

目 录

引 言	(i)
第一章 单符号离散信源	(1)
第一节 信源的数学模型	(1)
第二节 信源符号的自信息量	(3)
第三节 信源的信息熵	(7)
第四节 信息熵的代数性质	(12)
第五节 信息熵的解析性质	(20)
第六节 信息熵的最大值	(26)
第七节 熵函数的公理构成	(31)
第八节 加权熵及其数学特性	(36)
第九节 加权熵的公理构成	(46)
第十节 效用信息熵	(62)
习 题	(70)
第二章 单符号离散信道	(72)
第一节 信道的数学模型	(72)
第二节 信道的交互信息量	(76)
第三节 条件交互信息量	(82)
第四节 平均交互信息量	(89)
第五节 平均交互信息量的非负性	(95)
第六节 平均交互信息量的极值性	(98)
第七节 平均交互信息量的不增性	(105)
第八节 平均交互信息量的上凸性	(117)
第九节 信道容量及其一般算法	(121)
第十节 几种无噪信道的信道容量	(136)
第十一节 几种对称信道的信道容量	(140)
第十二节 可逆矩阵信道的信道容量	(152)
第十三节 信道容量的迭代计算	(156)
习 题	(166)
第三章 多符号离散信源与信道	(171)
第一节 离散平稳信源的数学模型	(171)

第二节	离散平稳无记忆信源的信息熵	(174)
第三节	离散平稳有记忆信源的信息熵	(178)
第四节	离散平稳有记忆信源的极限熵	(189)
第五节	马尔柯夫(Markov)信源的极限熵	(192)
第六节	信源的剩余度与结构信息	(215)
第七节	离散无记忆信道的数学模型	(217)
第八节	离散无记忆信道的信道容量	(224)
第九节	独立并列信道的信道容量	(230)
	习 题	(234)
第四章	单维连续信源与信道	(237)
第一节	相对熵与平均交互信息量	(237)
第二节	几种单维连续信源的相对熵	(246)
第三节	相对熵的极值性	(249)
第四节	相对熵的上凸性	(253)
第五节	最大相对熵定理	(255)
第六节	信息变差与熵功率	(261)
第七节	连续熵的变换	(263)
第八节	平均交互信息量的不变性	(267)
第九节	数据处理定理	(269)
第十节	连续信源的信息测量	(275)
第十一节	连续信道的信道容量	(281)
第十二节	高斯加性信道的容量	(286)
	习 题	(291)
第五章	多维连续信源与信道	(295)
第一节	随机过程的离散化	(295)
第二节	多维连续信源的熵	(314)
第三节	多维熵的最大值	(324)
第四节	多维熵的变换	(329)
第五节	多维连续信道的传输特性	(334)
第六节	高斯白噪声	(339)
第七节	高斯白噪声加性信道的容量	(342)
第八节	独立并列信道的最大容量	(350)
	习 题	(358)
第六章	无失真信源编码	(360)
第一节	单义可译码	(361)

第二节	非延长码及其构成	(363)
第三节	单义可译定理	(366)
第四节	平均码长与有效性	(371)
第五节	平均码长的界限定理	(375)
第六节	信源扩展与数据压缩	(382)
第七节	无失真信源编码定理	(387)
第八节	霍夫曼(Huffman)有效码	(390)
习 题		(408)
第七章	抗干扰信道编码	(411)
第一节	译码规则	(411)
第二节	译码规则的选择准则	(415)
第三节	信道编码的编码原则	(420)
第四节	抗干扰信道编码定理	(429)
第五节	分组码及其检纠能力	(439)
第六节	线性分组码的代数结构	(451)
第七节	线性分组码及其生成矩阵	(472)
第八节	一致校验矩阵与伴随式	(481)
第九节	标准阵列与译码表	(493)
第十节	检纠能力与一致校验矩阵的关系	(509)
第十一节	完备码	(515)
第十二节	汉明码与扩展汉明码	(522)
习 题		(531)
第八章	限失真信源编码	(537)
第一节	平均交互信息量的下凸性	(537)
第二节	平均失真度	(543)
第三节	信息率失真函数 $R(D)$ 与数据压缩	(548)
第四节	$R(D)$ 函数的数学特性	(564)
第五节	离散信源的 $R(D)$ 函数	(569)
第六节	离散信源 $R(D)$ 函数的参量表述	(580)
第七节	二元离散信源 $R(D)$ 函数的参量计算	(587)
第八节	正向与反向试验信道的转换	(592)
第九节	$R(D)$ 函数的迭代计算	(597)
第十节	高斯连续信源的 $R(D)$ 函数	(601)
第十一节	连续信源 $R(D)$ 函数的参量表述	(610)
第十二节	高斯连续信源 $R(D)$ 函数的参量计算	(613)

第十三节 正向与反向高斯加性试验信道的转换	(619)
第十四节 限失真信源编码定理	(627)
第十五节 $R(D)$ 函数与信息价值	(642)
第十六节 广义信息率失真函数	(652)
习题	(661)
第九章 网络信息理论	(665)
第一节 双输入单输出信道的信道容量	(665)
第二节 离散二址接入信道的容量界限	(669)
第三节 高斯加性二址接入信道的容量界限	(678)
第四节 单输入双输出信道的信道容量	(685)
第五节 高斯链式接续信道的容量界限	(688)
第六节 相关信源的边信息与公信息	(694)
习 题	(697)
附 录 《供熵函数计算用的几种函数表》	(699)
参考文献	(702)

第一章 单符号离散信源

通信系统一般由信源、信道和信宿三部分组成(如图 1.1)。“信源”就是信息的源泉.信息不是消息本身,但它又包含在消息之中.信源是由含有信息的消息组成的集合.若信源是由有限或无限

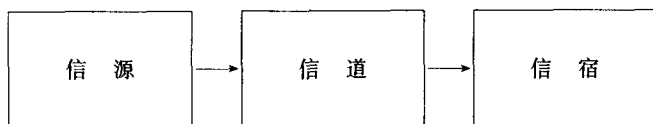


图 1.1

可列个取值离散的符号(如文字、字母、数字等)组成的离散集合,则这种信源称为离散信源.又若一个符号就代表一个完整的消息,则这种离散信源又称为单符号离散信源.单符号离散信源是最简单的离散信源.

第一节 信源的数学模型

信源要含有一定的信息,必须具有随机性,以一定的概率发出各种不同的符号.单符号离散信源是具有一定概率分布的离散符号的集合.基于对信源的这种认识,我们可用一个离散随机变量的可能取值,表示信源可能发出的不同符号;用离散随机变量的概率分布,表示信源发出不同符号可能性的大小.总之,我们可用一个离散随机变量来代表一个单符号离散信源.

例如,掷一个六面质地均匀的骰子,每次出现朝上一面的点数是随机的.如把出现朝上一面的点数作为这个随机试验的结果,并把试验的结果看作信源的输出消息,无疑,这个随机试验可看作是一个信源.这个信源输出有限种离散数字,其组成的集合为 $A: \{1, 2, 3, 4, 5, 6\}$,而且每一个数字代表一个完整的消息.所以,这个信源是单符号离散信源.我们可用离散随机变量 X 来表示这个单符号离散信源; X 的可能取值就是信源可能发出的各种不同符号,其状态空间就是信源可能发出的各种不同符号组成的集合 $A: \{1, 2, 3, 4, 5, 6\}$; X 的概率分布,就是信源发出各种不同符号的先验概率,其概率空间就是信源发出各种不同符号的先验概率组成的概率空间 $P: \left\{ P(X=1) = \frac{1}{6}, P(X=2) = \frac{1}{6}, \dots, P(X=6) = \frac{1}{6} \right\}$.所以,这个单符号离散信源的数学模型可完整地表示为

$$[X \cdot P]: \begin{cases} X: & 1 & 2 & 3 & 4 & 5 & 6 \\ P(X): & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} & \frac{1}{6} \end{cases}$$

我们把 $[X \cdot P]$ 称为信源 X 的“信源空间”。信源 X 输出的符号只可能是集合 $A: \{1, 2, 3, 4, 5, 6\}$ 中的任何一种,不可能是集合 A 以外的其它任何符号。信源 X 的概率空间是一个完备集,即有

$$P(X=1) + P(X=2) + \cdots + P(X=6) = 1$$

在这个典型实例的启发下,我们可构建一般单符号离散信源的数学模型。若某信源可能发出 r 种不同的符号 a_1, a_2, \dots, a_r ,相应的先验概率分别是 $p(a_1), p(a_2), \dots, p(a_r)$ 。我们用随机变量 X 表示这个信源,其信源空间可表示为

$$[X \cdot P]: \begin{cases} X: & a_1 & a_2 & \cdots & a_r \\ P(X): & p(a_1) & p(a_2) & \cdots & p(a_r) \end{cases} \quad (1.1)$$

其中

$$\begin{aligned} 0 \leq p(a_i) \leq 1 \quad (i=1, 2, \dots, r) \\ \sum_{i=1}^r p(a_i) = 1 \end{aligned} \quad (1.2)$$

不同信源对应不同的信源空间。如信源给定,这就意味着相应的信源空间已经确定。反之,如信源空间已经确定,这就意味着相应的信源已经给定。用信源空间表示信源的数学模型的必要前提,就是信源可能发出的各种不同符号的概率先验可知,或事先可测定。测定信源的概率空间是构建信源空间的关键。例如,在一个箱子中,有红、黄、蓝、白四种不同颜色的彩球,它们的大小、质量和重量完全一样。若从这个箱子中任意摸取出一个球,并把球的颜色当作试验的结果。显然,这个随机试验就可看作是一个单符号离散信源,信源的输出符号集就是四种不同的颜色 $A: \{\text{红, 黄, 蓝, 白}\}$ 。构建这个信源的信源空间的关键在于测定出现各种不同颜色的概率。在这个问题中,我们可把各种不同颜色的彩球的出现频率,近似地看作其出现的概率。假如,箱子中共有32个球,其中:红球16个;黄球8个;蓝球和白球各4个。则可得各种彩球出现频率,即各种彩球出现的先验概率分别为

$$\text{出现红球的概率} \quad P(\text{红}) = \frac{16}{32} = \frac{1}{2};$$

$$\text{出现黄球的概率} \quad P(\text{黄}) = \frac{8}{32} = \frac{1}{4};$$

$$\text{出现蓝球的概率} \quad P(\text{蓝}) = \frac{4}{32} = \frac{1}{8};$$

$$\text{出现白球的概率} \quad P(\text{白}) = \frac{4}{32} = \frac{1}{8}.$$

若用随机变量 X 表示这个信源,其信源空间为

$$[X \cdot P]: \begin{cases} X: & \text{红} & \text{黄} & \text{蓝} & \text{白} \\ P(X): & \frac{1}{2} & \frac{1}{4} & \frac{1}{8} & \frac{1}{8} \end{cases}$$

用一个离散随机变量 X 代表一个单符号离散信源,这就是我们用数学描述信源的基本原则。随机变量 X 的状态空间和概率空间,是构建信源空间 $[X \cdot P]$ 的两个基本要素,而概率空间是决定性要素。概率可测是香农信息论的基本前提。

第二节 信源符号的自信息量

由信息的定义可知,在通信过程中,收信者所获取的信息量,在数量上等于通信前后不确定性的消除.例如,一位教授作学术报告,由于扩音设备不好,在报告的声音中夹杂悦耳的音乐声.学生听到的是报告声音的一种变型.听了这场报告后,虽然知道了报告的大致内容,但某些具体细节没有听清,仍然存在一定的不确定性.学生从这场报告中所获取的信息量,应该是听报告前对报告内容的不确定性,减去听报告后仍然存在的不确定性所得之差,即不确定性的消除.

若信源发某符号 a_i ,由于信道中噪声的随机干扰,收信者收到的是 a_i 的某种变型 b_j .收信者收到 b_j 后,从 b_j 中获取关于 a_i 的信息量用 $I(a_i; b_j)$ 表示,则有

$$\begin{aligned} I(a_i; b_j) &= [\text{收到 } b_j \text{ 前,收信者对信源发 } a_i \text{ 的不确定性}] \\ &\quad - [\text{收到 } b_j \text{ 后,收信者对信源发 } a_i \text{ 仍然存在的不确定性}] \\ &= \text{收信者收到 } b_j \text{ 前、后,对信源发 } a_i \text{ 的不确定性的消除} \end{aligned} \quad (1.3)$$

当信道中没有噪声的随机干扰(无噪信道)时,信源发出的符号 a_i 可以不受任何干扰传递给收信者,收信者收到的 b_j 就是 a_i 本身.由于收信者确切无误地收到了信源发出的符号 a_i ,当然就完全消除了对信源发符号 a_i 的不确定性,即

$$\begin{aligned} &[\text{收到 } b_j \text{ 后,收信者对信源发 } a_i \text{ 仍然存在的不确定性}] \\ &= [\text{收到 } a_i \text{ 后,收信者对信源发 } a_i \text{ 仍然存在的不确定性}] \\ &= 0 \end{aligned} \quad (1.4)$$

这时,(1.3)式就可改写为

$$I(a_i; a_i) = [\text{收到 } a_i \text{ 前,收信者对信源发 } a_i \text{ 的不确定性}] \quad (1.5)$$

式(1.5)中的 $I(a_i; a_i)$ 表示收到 a_i 后,收信者从 a_i 中获取关于 a_i 的信息量.当然,这就是信源符号 a_i 所含有的全部信息量.我们把 $I(a_i; a_i)$ 称为信源符号 a_i 的自信息量,并用 $I(a_i)$ 表示.由(1.5)式即可有

$$I(a_i) = [\text{收到 } a_i \text{ 前,收信者对信源发 } a_i \text{ 的不确定性}] \quad (1.6)$$

(1.6)式表明,信源符号 a_i 的自信息量 $I(a_i)$ 的度量问题,已转变为信源发符号 a_i 的不确定性的度量问题.我们知道,不确定性是与可能性相联系的,而可能性又可由概率的大小来表示.可以断言,自信息量 $I(a_i)$ 一定是信源发符号 a_i 的先验概率 $p(a_i)$ 的某一函数,即

$$I(a_i) = f[p(a_i)] \quad (i = 1, 2, \dots, r) \quad (1.7)$$

从客观事实和人们的习惯概念出发,函数 $I(a_i) = f[p(a_i)]$ ($i = 1, 2, \dots, r$) 必须满足以下四个公理性条件:

(1) 若有两条消息:一条是“中国男子足球队获取世界杯冠军”.根据中国男子足球队的历来表现,获取世界杯赛的冠军的概率很小,即可能性很小.从习惯概念出发,认为“中国男子足球队获取世界杯冠军”这条消息的不确定性很大.这一事件一旦发生,人们就会获取很大的信息量,甚至会出现震动全国、万众欢呼的动人场面;另一条消息是“中国女子乒乓球队获取亚运会冠军”.根据中国

女子乒乓球队的历来表现,中国女子乒乓球队获取亚运会冠军的概率很大,即可能性很大.从习惯概念出发,认为“中国女子乒乓球队获取亚运会冠军”的不确定性很小,这一事件一旦发生,从这条消息中获取的信息量要小得多.

这个大家都承认的公理,可以这样来表述.如信源符号 a_i 和 a_j 的先验概率分别为 $p(a_i)$ 和 $p(a_j)$,且 $0 < p(a_i), p(a_j) < 1$. 若 $p(a_i) > p(a_j)$, 则

$$I(a_i) = f[p(a_i)] < I(a_j) = f[p(a_j)] \quad (1.8)$$

即函数 $I(a_i) = f[p(a_i)]$ 是先验概率 $p(a_i)$ 的单调递减函数.

(2)众所周知,“太阳从西边升起”的概率等于零,是不可能事件.从习惯概念出发,“太阳从西边升起”这条消息的不确定性应是无穷大.这个事件一旦发生,将会天翻地覆,人们获取无穷大的信息量.这就是说,如信源符号 a_i 的先验概率 $p(a_i) = 0$, 则

$$I(a_i) = f[p(a_i)] \rightarrow \infty \quad (1.9)$$

(3)人们同样公认,“太阳从东边升起”的概率等于1,是确定事件.从习惯概念出发,“太阳从东边升起”这条消息不存在任何不确定性.如果你把这条消息告诉别人,凡听到这条消息的人都会认为你讲的是废话,不会得到任何信息量.这就是说,如信源符号 a_i 的先验概率 $p(a_i) = 1$, 则

$$I(a_i) = f[p(a_i)] = 0 \quad (1.10)$$

(4)人们的习惯概念认为,两个统计独立事件的联合信息量,应等于它们各自信息量之和.比如,一般可把“天安门广场有人在照像”与“中国科学院研究生院上信息论课”看作为相互统计独立、互不相关的两件事.若我们同时得知“天安门广场有人照像”和“中国科学院研究生院上信息论课”这两条消息,从这两条消息中得到的联合信息量,应等于这两条消息各自信息量之和.

这个公理条件可这样来表述.设有两个统计独立的信源 X 和 Y . 信源 X 的符号 a_i 的先验概率为 $p(a_i)$; 信源 Y 的符号 b_j 的先验概率为 $p(b_j)$. 符号 a_i 和 b_j 组成的联合消息 (a_i, b_j) 的先验概率是联合概率 $p(a_i, b_j)$. 则

$$\begin{aligned} I(a_i, b_j) &= f[p(a_i, b_j)] \\ &= I(a_i) + I(b_j) \end{aligned} \quad (1.11)$$

从数学上可以证明,满足(1.8)、(1.9)、(1.10)、(1.11)四个公理条件的函数 $I(a_i)$, 是符号 a_i 的先验概念 $p(a_i)$ 的倒数的对数,即

$$\begin{aligned} I(a_i) &= \log \frac{1}{p(a_i)} \\ &= -\log p(a_i) \quad (i = 1, 2, \dots, r) \end{aligned} \quad (1.12)$$

很容易验证,(1.12)式满足公理条件(1.8)、(1.9)、(1.10). 同样,因为信源 X 和 Y 统计独立,所以

$$p(a_i, b_j) = p(a_i)p(b_j)$$

则由(1.12)式可得

$$I(a_i, b_j) = \log \frac{1}{p(a_i, b_j)} = \log \frac{1}{p(a_i)p(b_j)}$$