JOURNAL OF ELECTRONICS AND INFORMATION TECHNOLOGY

信源信道联合编码的一种方法 1

周延蕾 梁 钊* - ılı 张有为*

(北京航空航天大学电子工程系 北京 100083) *(五邑大学信息科学研究所 广东省江门市 529020)

摘 要 对于单个符号组成的信源,该文提出了一种在 Shannon 编码中嵌入信道编码的编码方法,在 [0,1) 区间内选择最小汉明距离为 d_{\min} 的一系列码字。在有信道误差的情况下,译码器运用码字间的最小汉明距离和 码字与符号的对应关系以及信源的统计特性,恢复出原始的符号序列,而不用采用已知的信道编码技术。 d_{\min} 提供了译码器所需的纠错能力。一般适用于带宽有限、误码率较高的二进制对称信道。

信源信道联合编码, Shannon 编码

TN911.2 中图号

1 引 言

分离定理是 Shannon 信息理论的一个基本定理。它表明: 只要满足信源的熵小于信道容 量,信源编码和信道编码的分离进行可实现无差错传输 [1]。为了达到这个理论极限、往往要求 信源编码有长的码字以及信道编码产生码字无限长的分组码。这在实际中是难以实现的。而且 实际的信源编码、例如霍夫曼编码和算术编码易受差错的影响、压缩数据中的一个错误往往会 一直延续下去、导致整个序列译码错误 [2] , 因此, 在带宽有限的噪声信道中传输压缩数据时, 需要用到可靠的信道编码技术。用信道编码附加冗余码来降低传输错误概率。传统的信源编码 最大限度地压缩信源以提高通信系统的有效性,但是在实际的通信系统中,信源编码器并不能 对信源进行完全压缩, 压缩数据中仍包含了冗余信息; 而且信道误差的引入使译码器的输入并 不完全等同于编码器的输出,因此,人们重新考虑这样一个问题: 信源是否值得压缩呢 [3] ?

近些年人们开始尝试利用信源中的冗余以及在信源编码内嵌入信道编码来降低系统的误码 率 [4,5],提高信息传送的可靠性。G.F.Elmasry 将信道编码嵌入了传统的算术编码 (ACECC)。 研究表明、信源概率不等时、 ACECC 比分别编码可获得更低的误码率 [6] , 本文介绍了一种将 信道编码嵌入 Shannon 编码的方法, 具有与 ACECC 相同的纠错能力, 并具体给出了编码和解 码的步骤,软件实现和试验结论。

2编码过程

对一个有 N 个不同符号的信源, $A\{a_1,a_2,\cdots,a_N\}$, 这里 $N=2^m(m$ 为整数) . 如果已知 信源以概率 p_i 发射符号 a_i , 那么每个符号 a_i 的累积概率 $P(a_i)$ 定义为

$$P(a_i) = \sum_{k=1}^{i-1} p_k \tag{1}$$

这里 $i=1,2,\cdots,N,$ $P(a_1)=0$ 。为了运算简便,将信源的概率舍入为 $2^{-k}(k)$ 为整数)。编码器 令每个符号的码字之间有最小汉明距离 d_{\min} , 也就是说, 对于一系列的码字 $C = c_1, c_2, \cdots, c_N$ 任何两个码字之间的汉明距离 $d(c_i, c_i) > d_{\min}(i \neq j)$.

^{1 2000-01-12} 收到, 2000-06-12 定稿

2.1 编码步骤

对于如上所述的 N 个符号的信源, 当 $d_{min} = 2$ 时, 编码步骤如下:

- (1) 将信源发出的 N 个消息符号按其概率的递减次序依次排列,并计算每一符号的累积概率。
- (2) 划分区间 [0,1) 为 N 个相等的小区间,每个符号分配一个小区间。汉明距离差别大的区间分配给具有相同或相近概率的符号。
 - (3) 将每个符号的累积概率映射到对应的小区间上、每个符号对应于其小区间内的一点。
 - (4) 去除小数点, 为每个输入符号分配码字, 顺序输出,

2.2 例子

设排列后的无记忆信源的符号集为 $A:\{a,b,c,d\}$,其相应概率为: p(0.5,0.25,0.125,0.125),累积概率用二进制表示为 P(0.0 , 0.10 , 0.110 , 0.111 。

当 $d_{\min}=2$ 时,根据步骤 2 ,区间 [0,1) 被分为四个相等的小区间 (S_0,S_1,S_2,S_3) ,如图 1(a) 中所示,00 ,01 ,10 和 11 ,分别对应于区间 S_0,S_1,S_2,S_3 。值得注意的是,区间 S_0 中的任一码字与 S_1 和 S_2 区间内的码字之间满足 $d_{\min} \geq 1$,与 S_3 的码字之间满足 $d_{\min} \geq 2$ 。由于符号与码字是一一对应的,以及各符号的码字与符号概率有关,那些具有相同或相近概率的符号被分配给汉明距离差别大的区间。因此,区间 S_0 和 S_3 分别分配给符号 c ,d ,区间 S_1 , S_2 分别分配给符号 a ,b 。步骤 3 将每个符号的累积概率映射到符号的对应区间上,得到区间内的一点 C ,例如将 a 的累积概率 0.0 映射到区间 S_1 上得到 C(a)=0.010 ,去除小数点,a 的码字即为 010 。同理,为 b ,c ,d 分配的码字分别为 1010 ,00110 ,11111 ,此码组满足两码元间的最小汉明距离为 2 。

当 $d_{\min}=3$ 的情况下,编码器把区间 [0,1) 划分为 2N 个而不是 N 个相等的小区间,其它步骤同上。因此编码器要在八个小区间中选择差别最大的四个区间。为简单起见,考虑下列 $d_{\min}=1$ 的二进制码组: 00 , 01 , 10 , 11 ,如果每个码字后面都加上一位奇偶校验位,就得到 $d_{\min}=2$ 的一系列码字 000 , 011 , 101 和 110 。如图 1(b) 中所示,这一组码字对应着区间 S_0 , S_3 , S_5 和 S_6 ,任何两区间的汉明距离都为 2 。将 a ,b ,c ,d 的累积概率分别映射到区间 S_0 , S_3 , S_5 和 S_6 上,结果分别对应于点 0.0000 , 0.01110 , 0.101110 , 0.110111 。因此各符号码字为: 0000 , 01110 , 101110 , 110111 。

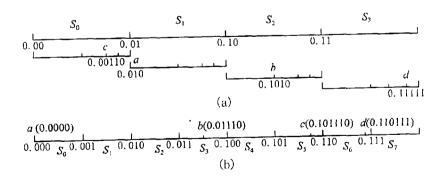


图 1 码区间划分示意

2.3 构造编码器

假设在符号集 $A\{a_1,a_2,\cdots,a_N\}$ 中,为符号 a_i 分配的区间的起始位置记为 S_{is} ,编码后 a_i 对应其区间内的一点 C_i 。根据上述的编码过程,当 $d_{\min}=2$ 时, a_i 的码字可由下式得出

$$C_i = S_{is} + P(a_i)/N \tag{2}$$

而当 $d_{\min} = 3$ 时, (2) 式变为

$$C_i = S_{is} + P(a_i)/(2N) \tag{3}$$

对于上例的情况,当 $d_{\min}=2$ 时, S_{is} 分别对应于 0.01 , 0.10 , 0.00 , 0.11 ; 当 $d_{\min}=3$ 时, S_{is} 分别对应于 0.000 , 0.011 , 0.101 , 0.110 。我们观察加入冗余码之前的情况: 0 , 10 , 111 ,这个结果与对单个符号进行 Shannon 编码的结果是一致的,只是由于符号概率的舍入,可以看作 Shannon 编码在信源概率为 1/2 的倍数时的特殊例子。最后得到的 N 个码字是具有最小汉明距离的 d_{\min} 的异前置码。这类异前置码可以通过在一系列传统的异前置码(例如霍夫曼码)的基础上附加一系列($d_{\min}-1$)的等长码简单地得到。

3译码过程

译码器同时进行信道和信源译码,对收到的二进制数据流 r 进行最大后验概率 (MAP) 译码 (即译码器选择使 p(y'/r) 有最大值的符号序列, y' 是估计符号序列)。 d_{\min} 提供了在错误概率已知的二进制对称信道 (BSC) 传输所需的检错能力。

3.1 译码步骤

- (1) 根据编码后各符号的码字长度,确定译码时输入数据的比特数,译码概率最大和最小的符号所需的输入比特数分别记为 l_{\min} 和 l_{\max} .
- (2) 从输入流的第一个比特开始译码,每次至少输入 l_{\min} 比特数据,最多输入 l_{\max} 比特数据。
- (3) 检查输入码与各符号码字之间的汉明距离,若找到一个符号与输入码之间的汉明距离小于 d_{\min} ,则作为译出符号;如果找不到与输入码的汉明距离小于 d_{\min} 的符号,则拒绝译码;如果有多个符号与输入码之间的汉明距离小于 d_{\min} ,则每一符号代表一个序列分别进行译码。
- (4) 如果同一段输入数据可译码为多个不同的符号序列, 即序列重合的情况, 则采用最大后 验概率规则选择一个序列。
 - (5) 若译出的符号序列没有用尽所有输入数据,则拒绝接受该序列,请求重发。

3.2 例题分析

当 $d_{\min} = 2$ 时,假设待解码序列是" adaabacb",编码器将产生下列二进制数据流:

010111110100101010010001101010

假设译码器接收到的数据如下, 带下划线的是错误比特:

<u>10100</u>11101001<u>1</u>10<u>0</u>00100011<u>1</u>1010

每次译码需要输入的最小比特数 $l_{\min} = 3(010$ 译码为 a) ,最大比特数 $l_{\max} = 5(00110$ 或 11111 译码后为 c ,d 。 因此,对任何输入码序列进行译码时,设置节点 3n ,3n+1 和 3n+2 ,其中 n 为整数.译码器的输出如下面的篱笆图所示,横轴 n(近似于卷积译码器中的时间 t),纵 轴为节点 3n ,3n+1 和 3n+2 。译码器从 $3n|_{n=0}$ 开始译码,每次至少输入 3 bit 数据,由于编码器共输出 30 比特数据,应选择在节点 $3n|_{n=10}$ 处结束的序列。图中符号的下标表示正确译码该符号需要纠正的比特数(例如 a_1 表示若要解出 a 需纠正一位错误)。在译出符号需要纠正的比特数和符号概率的基础上进行最大后验概率判决。

从节点 $3n|_{n=0}$ 开始的输入数据与符号 b 之间的汉明距离是 0 ,与其它符号的汉明距离均大于或等于 2 ,因此译码器将其译码为 b ,输入比特数为 4 .然后译码器再从节点 $3n+1|_{n=1}$ 处开始译码,可以看到输入数据与符号 a 和 c 的汉明距离都为 1 ,同其它符号的汉明距离均大

于等于 d_{\min} . 解码器同时建立两个序列,如图 2 中所示, b_0a_1 表示第一个序列, b_0c_1 表示第二个序列。但序列 b_0c_1 在节点 $3n|_{n=3}$ 处被拒绝译码,因为接下来输入码与任一符号的汉明距离均不小于 d_{\min} ,因此选择序列 b_0a_1 。在节点 $3n|_{n=7}$,序列 $a_1b_1a_0$ 与 c_1c_1 重合了。译码器运用 MAP 判决,选择序列 aba 。在节点 $3n+1|_{n=9}$ 处结束的序列因没有用尽所有的输入码而被拒绝接受。序列 " bababacb " 是唯一幸存的译出序列。

由于信道传输过程中引入一突发性错误,输入数据一开始就出现 5bit 的连续误码,与原始的符号序列"adaabacb"相比,译码器译出的符号序列前三个符号被错误译码。由图 2 中可以看出,虽然部分序列译码错误,但序列仍可以恢复过来,输入数据的剩余部分被正确译码,这样就避免了一部分译码错误引发出剩余序列的译码错误。译码过程中对于序列重合的情况采用MAP 判决,避免选择错误序列,图 2 中可以看出,译码即使选择序列出错,也不会导致误差扩散,剩余的输入数据仍可以正确译码。但信源与信道编码分离进行就会出现这种情况,一个信道译码时的错误会导致信源译码器丧失符号与码字的对应关系,开始错误的译码,导致整个序列译码错误。

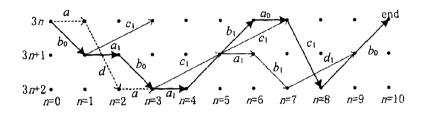


图 2 $d_{\min} = 2$ 时的译码过程

3.3 译码器的性能分析

发生符号译码错误的概率 PE 可以作为译码器的性能指标。对一个有 N 个符号的信源,对一个符号进行正确译码的概率 PC=1-PE 可从下面的表达式中得出:

$$PC = \sum_{i=0}^{d_{\min}-1} \sum_{i=1}^{N} p_i \begin{pmatrix} l_i \\ j \end{pmatrix} p_e^j (1 - p_e)^{l_i - j}$$
(4)

其中 p_e 是 BSC 信道的传输错误概率。 $d_{\min}=2$ 时,应用上面的公式我们得到:

$$PC = \sum_{i=1}^{N} p_i \begin{pmatrix} l_i \\ 0 \end{pmatrix} (1 - p_e)^{l_i} + \sum_{i=1}^{N} p_i \begin{pmatrix} l_i \\ 1 \end{pmatrix} (1 - p_e)^{l_i - 1}$$
$$= \sum_{i=1}^{N} p_i (1 - p_e)^{l_i - 1} (1 - p_e + l_i p_e)$$
(5)

将 $(1-p_e)^{l_i-1}$ 用幂级数代替,并舍去 p_e 的幂次大于等于 3 的项,考虑到 $\sum_{i=1}^N p_i=1$,并应用于前例的情况我们得到如下的表达式:

$$PE \approx \frac{p_e^2}{2} \sum_{i=1}^{N} p_i l_i (l_i - 1) \approx \frac{11 p_e^2}{2}$$
 (6)

同样当 $d_{min} = 3$ 时

$$PE \approx \frac{p_e^2}{6} \sum_{i=1}^{N} p_i l_i (l_i - 1) (l_i - 2) \approx \frac{19 p_e^3}{2}$$
 (7)

未编码时、 $d_{min} = 1$, 直接由 (3) 式得

$$PE = \sum_{i=1}^{N} p_i l_i = \frac{7p_e}{4}$$
 (8)

在给定的传输概率下, (6) 和 (7) 式对比了不同编码效率下符号的错误概率,均低于未编码时的情况,由此可见,此种译码器降低了符号的错误概率,提高了检错纠错的能力。

4 软件实现与结果分析

在实现仿真的过程中,我们采用的信源模型是例如手写脚本、程序一类的可读文本文件 (包括 256 个不同的符号)。编码过程采用先统计后编码的自适应编码方法。为了计算简便,统计后将文件中符号的个数近似为 2 的整数次幂,符号概率用 $2^{-k}(k$ 为整数) 去近似。输入待编码文件的文件名和希望达到的最小汉明距离 d_{\min} ,编码后输出编码结果和对应的码本两个数据文件。编码过程如前所述。

信道采用高斯噪声 (AWGN) 信道,由于会出现译码器拒绝译码的情况,差错控制方式采用自动请求重发方式。

译码过程首先从要译码的数据文件中读入一定比特的数据,将其与码本中对应长度的码字一比较,选汉明距离小于 d_{min} 码字作为译出码。开辟缓冲区暂存满足条件的所有译出符号序列,采用 MAP 判决,输出最终的符号序列,否则拒绝译码,反馈请求重发,清除缓冲区。试验结果分析:

- (1) 随着文件中不同符号个数的增加、输出文件的长度增加。
- (2) 采用这种信源信道的联合编码大大降低了符号译码错误概率。码字之间的汉明距离越大,检错与纠错能力越强。对于信道中的突发性错误、符号错误等都有很好的纠错能力,更不会使误差扩散,导致整个序列的译码错误。图 3 画出了不同的编码效率 R 下,符号误码率 PE 与信道传输错误概率 p_e 之间的关系曲线。
- (3) 在信道误码率过高的情况下,译码器找不到与输入码满足最小汉明距离的码字,使译码过程停止,要求发端重发,也就是说,译码器能够识别不可恢复性错误而停止译码。
- (4)由于编码结果得到的是一组变长码,实际应用时常需要有大容量的存储设备作为缓冲, 这也便于重发。同时也需要克服存储器的溢出和取空问题,如采用信息分段发送。

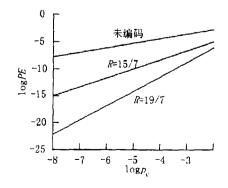


图 3 信道误码率与符号错误概率之间的关系曲线

5 结束语

对信源概率近似为 1/2 的倍数的信源,本文介绍了一种在 Shannon 编码中嵌入信道编码方法,得到一组具有最小汉明距离的码字。译码器运用码字间的最小汉明距离和码字与符号的对应关系以及信源的统计特性,恢复出原始的符号序列,码字间的汉明距离提供了译码器所需的纠错能力。与单纯的信源编码相比,这种联合编码的方法以降低编码效率为代价,获得了较低的符号误码率,并可以有效地防止误差的扩散。当误码率高,译码器找不到满足 d_{\min} 的符号时,译码器会停止译码,要求发端重发,适合误码率较高时的二进制对称信道。

参考文献

- [1] C. E. Shannon, The mathematical theory of communication., BSTJ, 1949, 27, 379-423.
- [2] 周炯磐, 丁晓明, 信源编码原理, 北京, 人民邮电出版社, 1994, 50-60.
- [3] G. Buch, F. Burkert, J. Hagenauer, B. Kukla, To compress or not to compress?, Proc. of IEEE Globecom Communication Theory Mini Conference, London, U. K., 1996, 198-203.
- [4] K. Ssyood, J. C. Borkenhagen, Use of residual redundancy in the design of joint source and channel coders., IEEE Trans. on Commun., 1991, 39(6), 838-846.
- [5] K. Sayoof, F. Liu, J. D. Gibson, A constrained joint source/ channel coder design., IEEE J. on SAC, 1994, 12(9), 1584-1593.
- [6] G. F. Elmasry, Embedding channel coding in arithmetic coding., IEE Proc.-I, 1999, 146(2), 73-78.

A JOINT SOURCE/CHANNEL CODING DESIGN

Zhou Yanlei Liang Zhao* Meng Shan Zhang Youwei*

(Dept. of EE, Beijing University of Aeronautics and Astronautics, Beijing 100083, China)

*(Institute of Information and Science, Wuyi University, Jiangmen 529020, China)

Abstract In this paper a method of joint source/channel coding is presented. For a source with separate symbols, codes having desired minimum Hamming distance d_{\min} are selected within range [0,1). Then, under the case having channel error, the decoder utilizes the d_{\min} , the corresponding relationship between code and symbol, and the statistical characteristic of source to recover the original sequence of symbols without the use of channel coding. The introduced scheme is suitable for bandwidth-limited binary symmetric channel (BSC) with fairly high error rates.

Key words Joint source/channel coding, Shannon coding

周延韶: 女, 1976年生, 硕士生, 研究方向为联合编码、分布式网络系统研究.

梁 钊: 男, 1947年生, 副教授, 主要研究领域为通信与编码.

蒙 山: 男,1973年生,博士生,主要研究领域为语音识别、图像处理和多媒体人机交互系统。

张有为: 男, 1937 年生, 教授, 博士生导师, 中国电子学会高级会员, IEEE 会员, 主要研究领域是信号与信息处理, 检测与估计理论, 多目标跟踪.