

第五章离散信道编码

5.1 模型与概念

陈兴同

中国矿业大学 数学学院

2021 年 8 月

内容提要

1 信道编码模型

内容提要

① 信道编码模型

② 译码准则

内容提要

① 信道编码模型

② 译码准则

③ 平均差错

内容提要

- ① 信道编码模型
- ② 译码准则
- ③ 平均差错
- ④ 最小汉明距离译码规则

引言

人们总是希望信息在信道中传输时既“快”又“准”，“快”体现在传输的信息量大，又称为效率高；“准”体现在准确率高，又称为可靠性高。但是由于信道容量的限制及传输过程中噪声的干扰使得传输的效率与可靠性可能达不到理想的状况，传输错误不可避免。但是仙农的信道编码定理说明可以使用合适的编码使信道传输的效率与错误概率控制在合理的范围内。本章将学习信道编码的基本原理，包括仙农的信道编码定理、信道编码的方法等。

传输模型图 5-1

下面将根据这个模型图定义信道编码概念。

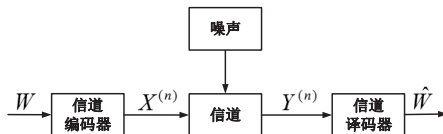


Figure: 图 5-1

定义: 5.1.1

(离散无记忆信道 (M,n) 编码) 离散无记忆信道 $\{\mathcal{X}, Q(y|x), \mathcal{Y}\}$ 的一个 (M,n) 信道码包括:

- (1) 信源编码器输出的 M 个消息, 它们正好是信道编码器的输入字符, 用随机变量 W 表示, 取值空间为 $\mathcal{W} = \{1, 2, \dots, M\}$ 。
- (2) 一个编码映射 $f: \mathcal{W} \rightarrow \mathcal{X}^n$, 每个消息都被编成一个 n 长码字, 共有 M 个码字: $x^{(n)}(i) = f(i), i \in \mathcal{W}$ 。全体码字集合称为码表。
- (3) 一个译码映射 $g: \mathcal{Y}^n \rightarrow \mathcal{W}$, 信道输出的每个 n 长序列 $y^{(n)}$ 都被译成一个消息, 这是一个从输出得到输入的规则。
- (4) 编码速率或码率 (单位按对数底)

$$R = \frac{1}{n} H(W),$$

当输入的消息等概率分布时

信道译码有多种

信道的编码方法有多种，后面将学习二元线性分组码的编码方法。译码规则的选择很有讲究，同一个信道的输出可以制定出多种译码规则。图 5-2 是一个二进对称信道，它可以有二种译码方法，比如：

$$g_1 : \quad b_1 \rightarrow a_1 \quad , \quad b_2 \rightarrow a_2,$$

$$g_2 : \quad b_1 \rightarrow a_2 \quad , \quad b_2 \rightarrow a_1.$$

如果信道输入是消息 i 的码字 a_i ，从输出 b_j 无法译成消息 i ，则译码错误，否则译码正确。不同译码方法有不同的译码错误。如果输入是 a_1 输出是 b_2 ，按译码规则 g_1 就会译错，而按译码规则 g_2 就会译对。

二进对称信道

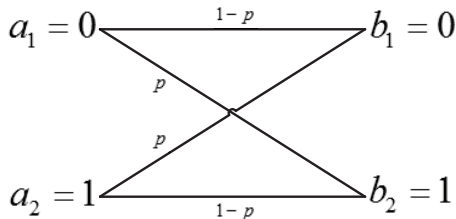


Figure: 图 5-2

译码错误定义:

关于定义 5.1.1 中的 (M, n) 码的译码错误定义。(1) 当信道输入为第 i 个消息的码字 $x^{(n)}(i)$ 时, 相应的**译码错误**为

$$e_i = \sum_{y^{(n)}} P \left\{ g(y^{(n)}) \neq i | X^{(n)} = x^{(n)}(i) \right\}. \quad (5.1)$$

(2) 最大译码错误为

$$e^{(n)} = \max_{1 \leq i \leq M} e_i. \quad (5.2)$$

(3) 在输入消息为等概率分布时平均译码错误为

$$P_e = \frac{1}{M} \sum_{i=1}^M e_i, \quad (5.3)$$

平均译码错误:

如果输入消息概率分布是 p_1, p_2, \dots, p_M , 则平均译码错误为

$$P_e = \sum_{i=1}^M p_i e_i. \quad (5.4)$$

注意到信道输入就是消息 i 的码字即一个 n 长字符串 $x^{(n)}(i)$, 记为随机变量 $X^{(n)}$; 信道输出也是一个 n 长字符串 $y^{(n)}$, 记为随机变量 $Y^{(n)}$, 则平均译码错误 (5.4) 也可以表为

$$P_e = \sum_{i=1}^M \sum_{g(y^{(n)}) \neq i} p(x^{(n)}(i), y^{(n)}) = P\{X^{(n)} \neq g(Y^{(n)})\}. \quad (5.5)$$

平均译码错误的大小与编码或译码准则的纠错性能力有关, 应当构造纠错能力强的编码与译码算法。

两种译码规则：

对同一个信道编码，可以有不同的译码规则，一般都是根据输入与输出之间的条件概率来建立，目标是使平均译码错误最小。最常用的有最大后验概率译码方法与最大似然概率译码方法。

(1) **最大后验概率译码方法**：这种译码方法是在输出为 $y^{(n)}$ 时检查哪个输入 $x^{(n)} \in \mathcal{C}$ 会使**后验概率** $p(x^{(n)}|y^{(n)})$ 最大，就将 $y^{(n)}$ 译成那个码字 $x^{(n)}$ 对应的消息。译码映射是

$$g(y^{(n)}) = \arg \max_{i \in \mathcal{W}} p(x^{(n)}|y^{(n)}).$$

(2) **最大似然概率译码方法**：这种译码方法是在输出为 $y^{(n)}$ 时检查哪个输入 $x^{(n)} \in \mathcal{C}$ 会使**似然概率** $p(y^{(n)}|x^{(n)})$ 最大，就将 $y^{(n)}$ 译成那个码字 $x^{(n)}$ 对应的消息。译码映射是

$$g(y^{(n)}) = \arg \max_{i \in \mathcal{W}} p(y^{(n)}|x^{(n)}).$$

两种译码方法比较：

通常最大后验概率译码方法的建立要已知每个输入 $x^{(n)}(i)$ 的概率，才能计算出后验概率 $p(x^{(n)}|y^{(n)})$ ，然后才能建立译码规则；事实上：按公式

$$p(x^{(n)}|y^{(n)}) = \frac{p(x^{(n)}, y^{(n)})}{p(x^{(n)})},$$

可以求出当输入分布已知时后验概率矩阵，在它的每一行中找最大元素所在的列号 i ，此时就可以将 $y^{(n)}$ 译成码字 $x^{(n)}(i)$ 。//但最大似然概率译码方法不需要信道的输入分布，可根据信道矩阵直接建立，只需要在信道矩阵每列中找最大元素所在的行号 i ，此时就可以将 $y^{(n)}$ 译成码字 $x^{(n)}(i)$ ，所以最大似然概率译码方法实现更方便。

两种译码的等价:

命题 5.1

当信道输入是等概分布时，最大后验概率译码方法等同于最大似然概率译码方法。

证明：因为信道输入为等概分布，故可设这个概率为常数 p 即 $p(x^{(n)}) = p$ ，则由 Bayes 公式得：

$$\begin{aligned} p(x^{(n)}|y^{(n)}) &= \frac{p(x^{(n)}, y^{(n)})}{p(y^{(n)})} \\ &= \frac{p(y^{(n)}|x^{(n)})p(x^{(n)})}{\sum_{\tilde{x}^{(n)} \in \mathcal{C}} p(y^{(n)}|\tilde{x}^{(n)})p(\tilde{x}^{(n)})} \\ &= \frac{p(y^{(n)}|x^{(n)})}{\sum_{\tilde{x}^{(n)} \in \mathcal{C}} p(y^{(n)}|\tilde{x}^{(n)})} \end{aligned}$$

由于当输出 $y^{(n)}$ 固定时，和 $\sum_{\tilde{x}^{(n)} \in \mathcal{C}} p(y^{(n)}|\tilde{x}^{(n)})$ 是常量。故当后验概率 $P\{x^{(n)}|y^{(n)}\}$ 最大时，似然概率 $p(y^{(n)}|x^{(n)})$ 也最大；若似然概率 $p(y^{(n)}|x^{(n)})$ 最大，则后验概率 $p(x^{(n)}|y^{(n)})$ 也最大。

例题 5.1.1

设信道矩阵如下，试确定最大似然概率译码规则与最大后验概率译码规则。

$$\begin{pmatrix} 0.5 & 0.3 & 0.2 \\ 0.2 & 0.3 & 0.5 \\ 0.3 & 0.3 & 0.4 \end{pmatrix}.$$

解：设信道字符集为 $\mathcal{X} = \{a_1, a_2, a_3\}$, $\mathcal{Y} = \{b_1, b_2, b_3\}$ 。则有条件概率

$$\begin{aligned} p(b_1|a_1) &= 0.5 & p(b_2|a_1) &= 0.3 & p(b_3|a_1) &= 0.2, \\ p(b_1|a_2) &= 0.2 & p(b_2|a_2) &= 0.3 & p(b_3|a_2) &= 0.5, \\ p(b_1|a_3) &= 0.3 & p(b_2|a_3) &= 0.3 & p(b_3|a_3) &= 0.4, \end{aligned}$$

最大似然概率译码规则：

$$g : b_1 \rightarrow a_1, \quad \textcircled{b_2 \rightarrow a_1 \text{ 或 } a_2 \text{ 或 } a_3}, \quad b_3 \rightarrow a_2。$$

(Handwritten red note: $b_2 \rightarrow a_3$)

如果信道有输入分布 $X \sim p_X(x) = (0.25, 0.25, 0.5)$ ，则可求出输出分布

$$Y \sim p_Y(y) = (0.325, 0.3, 0.375),$$

以及后验概率

$$\begin{aligned} p(a_1|b_1) &= 5/13 & p(a_2|b_1) &= 2/13 & p(a_3|b_1) &= 6/13, \\ p(a_1|b_2) &= 1/4 & p(a_2|b_2) &= 1/4 & p(a_3|b_2) &= 1/2, \\ p(a_1|b_3) &= 2/15 & p(a_2|b_3) &= 5/15 & p(a_3|b_3) &= 8/15, \end{aligned}$$

由此可得最大后验概率译码规

则： $g : b_1 \rightarrow a_3$ ， $b_2 \rightarrow a_3$ ， $b_3 \rightarrow a_3$ ，即不论输出什么全译成 a_3 。

例题 5.1.2

平均错误 P_e 不仅与译码规则有关，也与编码规则有关。这里以二进对称信道为例来说明不同编码规则时最大似然译码规则的译码错误。

→ (无重复编码) 图 5.2 是一个二进对称信道，当 $p = 0.01$ 时其容量为 $C = 0.9192\text{bits}$ 。来自信源编码器的输出字符仅有两个 $\mathcal{W} = \{0, 1\}$ 。现在进行无重复编码，编码方案为

$$f: w_1 = 0 \rightarrow a_1 = 0, w_2 = 1 \rightarrow a_2 = 0^1; \quad g: b_1 \rightarrow w_1, b_2 \rightarrow w_2,$$

于是使用该信道传输时译码错误概率为

$$e_1 = \sum_{b_j} P \{g(b_j) \neq w_1 | a_1\} = p(b_2 | a_1) = p = 10^{-2},$$

$$e_2 = \sum_{b_j} P \{g(b_j) \neq w_2 | a_2\} = p(b_1 | a_2) = p = 10^{-2},$$

因为 $e_1 = e_2$ ，故对任何输入分布平均错误都是 $P_e = 10^{-2}$ 。

例题 5.1.3

(三重复编码) 将来自信源编码器的输出 $\mathcal{W} = \{0, 1\}$ 重复三次进行编码

$$f: w_1 = 0 \rightarrow a_1 = 000, w_2 = 1 \rightarrow a_2 = 111.$$

使用图 5-2 所示二进对称信道传输 ($p = 0.01$)，在输出端 3 长
输出可能有 8 个 传错

$$b_1 = 000, b_2 = 001, b_3 = 010, b_4 = 011, b_5 = 100, b_6 = 101, b_7 = 110, b_8 = 111$$

它们与输入之间的条件概率矩阵为

$$Q(b_j|a_i) = \begin{pmatrix} q^3 & q^2p & q^2p & qp^2 & q^2p & qp^2 & qp^2 & p^3 \\ p^3 & p^2q & p^2q & pq^2 & p^2q & pq^2 & pq^2 & q^3 \end{pmatrix}, q = 1-p,$$

最大似然译码规则为

$$\left. \begin{matrix} b_1 = 000 \\ b_2 = 001 \\ b_3 = 010 \\ b_5 = 100 \end{matrix} \right\} \rightarrow w_1, \quad \left. \begin{matrix} b_4 = 011 \\ b_6 = 101 \\ b_7 = 110 \\ b_8 = 111 \end{matrix} \right\} \rightarrow w_2,$$

续：例题 5.1.1

于是译码错误概率

$$e_1 = \sum_{b_j} P \{g(b_j) \neq w_1 | a_1\} = p(b_4 | a_1) + p(b_6 | a_1) + p(b_7 | a_1) + p(b_8 | a_1) = 3$$

$$e_2 = \sum_{b_j} P \{g(b_j) \neq w_2 | a_2\} = p(b_1 | a_2) + p(b_2 | a_2) + p(b_3 | a_2) + p(b_5 | a_2) = 3$$

因为 $e_1 = e_2$ ，故对任何输入分布平均错误都为： $P_e = 3 \times 10^{-4}$ 。

重复编码好处：

这说明重复编码有助于降低平均译码错误，但是每个信道码字符所能携带的信息量即码率或编码速率 $R = 1/3\text{bits}$ 较小。

可以继续重复编码，错误概率继续下降，码率也不断减小。

比如：当 $n = 5$ 时 $P_e = 10^{-5}$, $R = 1/5\text{bits}$ ；当 $n = 7$ 时 $P_e = 10^{-7}$, $R = 1/7\text{bits}$ ；当 $n = 9$ 时 $P_e = 10^{-8}$, $R = 1/9\text{bits}$ ；当 $n = 11$ 时 $P_e = 10^{-10}$, $R = 1/11\text{bits}$ 。

重复编码实际上是增加被传输信息的冗余，那些重复的码符可以抵消噪声的影响，减少译码错误，提高传输的可靠性，但却降低了传输的效率，因为每信道字符所携带的信息量减小。是否可以在提高可靠性的同时又提高传输的有效性，让可靠性与有效性都能满足我们的要求？那就是对长信息进行编码。

消

例题 5.1.4

（消息串编码） 设信源输出字符为 $\mathcal{W} = \{0, 1\}$ ，现在对 2 长消息串进行信道编码，此时 $\mathcal{W} = \{00, 01, 10, 11\}$ 。还是使用 3 长的信道码，方案如下：

$$f: w_1 = 00 \rightarrow a_1 = 000, w_2 = 01 \rightarrow a_2 = 010, w_3 = 10 \rightarrow a_3 = 100, w_4 = 11 \rightarrow a_4 = 110$$

使用图 5-2 所示信道 ($p = 0.01$) 传输，在输出端 3 长输出可能有 8 个

$$b_1 = 000, b_2 = 001, b_3 = 010, b_4 = 011, b_5 = 100, b_6 = 101, b_7 = 110, b_8 = 111$$

它们与输入之间的条件概率为

$$Q(b_j|a_i) = \begin{pmatrix} q^3 & q^2p & q^2p & qp^2 & q^2p & qp^2 & qp^2 & p^3 \\ pq^2 & p^2q & q^3 & q^2p & p^2q & p^3 & pq^2 & qp^2 \\ pq^2 & p^2q & qp^2 & p^3 & q^3 & pq^2 & q^2p & qp^2 \\ p^2q & p^3 & q^2p & p^2q & pq^2 & qp^2 & q^3 & pq^2 \end{pmatrix}, q = 1-p,$$

续：例题 5.1.4

最大似然译码规则为

$$\left. \begin{array}{l} b_1 = 000 \\ b_2 = 001 \\ b_5 = 100 \\ b_6 = 101 \end{array} \right\} \rightarrow w_1, \quad \left. \begin{array}{l} b_3 = 010 \\ b_4 = 011 \\ b_7 = 110 \\ b_8 = 111 \end{array} \right\} \rightarrow w_2, \\ \left. \begin{array}{l} b_5 = 100 \\ b_6 = 101 \end{array} \right\} \rightarrow w_3, \quad \left. \begin{array}{l} b_7 = 110 \\ b_8 = 111 \end{array} \right\} \rightarrow w_4,$$

译码错误

$$\begin{aligned} e_1 &= \sum_{b_j} P \{g(b_j) \neq w_1 | a_1\} = 1 - \sum_{b_j} P \{g(b_j) = w_1 | a_1\} \\ &= 1 - p(b_1 | a_1) - p(b_2 | a_1) = 1 - q^3 - q^2 p = e_2 = e_3 = e_4. \end{aligned}$$

平均错误为 $P_e = 1.99 \times 10^{-2}$ ，码率为 $R = 2/3 \text{bits}$ 。和上例相比，本例说明增加要编码的消息串长度可以提高效率，但同时又降低了可靠性。实际上效率与可靠性之间总是一对矛盾。

例题 5.1.5

对图 5-2 中二进对称信道 ($p = 0.01$) 输入四个 2 长消息

$$w_1 = 00, w_2 = 01, w_3 = 10, w_4 = 11$$

进行 5 长编码，方案如下：

$$x^{(5)} = (x_1, x_2, x_3, x_4, x_5) = (u_1, u_2)A \quad \text{其中 } A = \begin{pmatrix} 1 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \end{pmatrix},$$

生成的码字是

$$f : w_1 \rightarrow c_1 = 00000, w_2 \rightarrow c_2 = 01101,$$

$$w_3 \rightarrow c_3 = 10111, w_4 \rightarrow c_4 = 11010,$$

续：例题 5.1.5

在输出端可能会有 32 个 5 长序列输出，采用最大似然译码规则

$$\begin{array}{l}
 00000 \\
 00001 \\
 00010 \\
 00100 \\
 01000 \\
 10000 \\
 00011 \\
 10001
 \end{array} \left. \vphantom{\begin{array}{l} 00000 \\ 00001 \\ 00010 \\ 00100 \\ 01000 \\ 10000 \\ 00011 \\ 10001 \end{array}} \right\} \rightarrow w_1, \quad
 \begin{array}{l}
 01101 \\
 01100 \\
 01111 \\
 01001 \\
 00101 \\
 11101 \\
 11100 \\
 01110
 \end{array} \left. \vphantom{\begin{array}{l} 01101 \\ 01100 \\ 01111 \\ 01001 \\ 00101 \\ 11101 \\ 11100 \\ 01110 \end{array}} \right\} \rightarrow w_2, \quad
 \begin{array}{l}
 10111 \\
 10110 \\
 10101 \\
 10011 \\
 11111 \\
 00111 \\
 00110 \\
 10100
 \end{array} \left. \vphantom{\begin{array}{l} 10111 \\ 10110 \\ 10101 \\ 10011 \\ 11111 \\ 00111 \\ 00110 \\ 10100 \end{array}} \right\} \rightarrow w_3, \quad
 \begin{array}{l}
 11010 \\
 11011 \\
 11000 \\
 11110 \\
 10010 \\
 01010 \\
 01011 \\
 11001
 \end{array} \left. \vphantom{\begin{array}{l} 11010 \\ 11011 \\ 11000 \\ 11110 \\ 10010 \\ 01010 \\ 01011 \\ 11001 \end{array}} \right\} \rightarrow w_4$$

续：例题 5.1.4

可以计算出条件概率 $P\{y^{(5)}|c_i\}$ 为

$$e_i = \sum_{y^{(5)}} P\{g(y^{(5)}) \neq w_i | X^{(5)} = c_i\} = 1 - q^5 - 5q^4p - 2q^3p^2 \approx 7.86 \times 10^{-4}$$

其中 $q = 1 - p$ ，从而平均错误为 $P_e \approx 7.86 \times 10^{-4}$ ，码率为

$$R = \frac{\log 4}{5} = \frac{2}{5} \text{ bits} < C.$$

和例题 5.1.4 相比，这个例子说明适当增大消息长度及编码长度，可以达到同时提高传输效率与传输可靠性。

更一般地仙农的信道编码定理将会指出：只要给定的码率 $R < C$ ，则对任意小的正数 ε ，存在码率为 R 的 n 长信道编码，它的平均错误满足 $P_e \leq \varepsilon$ 。

汉明距离定义 5.1.3

也可以将输出序列译成与它最接近的输入序列，但如何度量两个序列的接近程度？可以利用汉明距离。

设 $\mathcal{F} = \{0, 1, 2, \dots, D-1\}$ 是 D 进字符集， \mathcal{F}^n 表示所有 n 长字符串的集合，两 n 长字符串 $x^{(n)} = (x_1, x_2, \dots, x_n), y^{(n)} = (y_1, y_2, \dots, y_n)$ 之间汉明距离（记为 $d(x^{(n)}, y^{(n)})$ ）是指这两个字符串对应分量中不相等元素个数，即集合

$$\{i | x_i \neq y_i, i = 1, 2, \dots, n\}$$

中元素个数。

汉明距离只在两个等长序列之间定义，不等长序列没有汉明距离；可以用汉明距离来衡量两个序列之间的接近程度。若汉明距离为 0，则两个序列完全相同。

命题 5.1.2: 汉明距离性质

- (1) 非负性或正定性: $d(x^{(n)}, y^{(n)}) \geq 0$, 并且
 $d(x^{(n)}, y^{(n)}) = 0 \Leftrightarrow x^{(n)} = y^{(n)}$ 。
- (2) 对称性: $d(x^{(n)}, y^{(n)}) = d(y^{(n)}, x^{(n)})$ 。
- (3) 三角不等式:

$$d(x^{(n)}, y^{(n)}) \leq d(x^{(n)}, z^{(n)}) + d(z^{(n)}, y^{(n)})$$

练习：

证明第（3）条结论即三角不等式。

练习:

已知字符串 $u \in \mathcal{F}^n$, 分别求在 D 进 n 长字符串空间 \mathcal{F}^n 中有多少个字符串 v 使 $d(u, v) = i$, 其中 $i = 0, 1, 2, \dots, n$ 。

命题和定义:

命题 5.1.3

对于两个等长的二进字符串 $x^{(n)}, y^{(n)}$, 若它们按位异或和为字符串 $z^{(n)} = x^{(n)} \oplus y^{(n)}$, 则它们汉明距离恰好是 $z^{(n)}$ 中非 0 元的个数。

比如两个字符串 $x^{(6)}=100111$ 与 $z^{(6)} = 111111$ 的异或和为 $x^{(6)} \oplus z^{(6)} = 011000$, 因此汉明距离为 2。

定义 5.1.7

最小汉明距离译码规则 如果对给定的 n 长输出字符串 $y^{(n)}$, 用汉明距离来衡量与它最接近的 n 长输入字符串为 $x^{(n)}$, 则将 $y^{(n)}$ 译成码字 $x^{(n)}$ 对应的消息 i 。

命题 5.1.4:

设二元对称信道（如图 5-2）中字符传错概率 $0 < p < 0.5$ ，则最大似然译码规则与最小汉明距离译码规则等同。

证明:

对已知输出字符串 $y^{(n)}$, 它与输入字符串 $x^{(n)}$ 之间似然概率可以用它们之间的汉明距离 $d = d(x^{(n)}, y^{(n)})$ 来表示

$$p\left(y^{(n)}|x^{(n)}\right) = p^d(1-p)^{n-d} = \left[\frac{p}{1-p}\right]^d (1-p)^n.$$

注意到当 $0 < p < 0.5$ 时, $0 < \frac{p}{1-p} < 1$, 所以当 d 最小时似然概率最大, 反之也对, 从而这两种译码规则是等价的。

最小汉明距离译码规则也可以使用在非二进对称信道中, 只是它未必与最大然概率译码规则等价。