

信源与信源无失真编码

2020年4月3日 21:59

0. 概述:

我们可以从通过信息论来讨论无失真压缩的极限。这也是香农信息论的第一定理。信源无失真压缩定理。基本思路则是在信息量不变的前提下，提高每个码字的信息量，从而降低码字长度。那么我们就一定是从非等概率分布(均匀分布)等概率分布。我们可以通过:渐进等同分割性质和典型序列是成同一目标。下面来讨论一下。

1. 渐进等同分割性质

首先回顾一下概率论中的大数定律。伯努利在1713年首先给出了频率稳定性定理表述。该结论是大多数定律的第一个。

在一个“成功”概率为 p 的伯努利试验序列中

$$X_n = \begin{cases} 1, & \text{第一次试验A发生.} \\ 0, & \text{第一次试验A不发生.} \end{cases} \quad S_n = \sum_{i=1}^n X_i, \quad n=1, 2, \dots$$

对于任意正数 ϵ 和 δ , 存在正整数 N . 当 $n > N$ 时, 有

$$P\left(\left|\frac{S_n}{n} - p\right| \geq \epsilon\right) < \delta \quad \text{或} \quad P\left(\left|\frac{S_n}{n} - p\right| < \epsilon\right) \geq 1 - \delta$$

通俗点讲就是频率等于概率。经典分布趋于统计分布。

那么渐进等同分割 (Asymptotic Equipartition Property) 就是对一段随机变量序列进行等同分割。让长度 n 进行到无穷大。

由弱大数定律可知, 当分割长度无限大时, 任何事件出现的频率都趋于概率。那么差不多所有的事件都将等概率出现。对于数据压缩的极限也认识清晰。(即每一个事件出现的个数 n_k 有 $\left|\frac{n_k}{n} - P(k)\right| < \epsilon$)

0 等同分割 \Rightarrow 等概率 \Rightarrow 最大熵原理 \Rightarrow 定长编码定理。

0 暗示了一种定长编码的方法。

下面我们对其进行严格的数学推导:

如令 X_1, X_2, \dots 是独立同分布的高维随机变量, 分布服从 $P(x)$, 则 $-\frac{1}{n} \log p(X_1, X_2, \dots, X_n) \xrightarrow{P} H(X)$

“ P ”代表以概率 P 收敛于 $H(X)$ 。

对于这个数学表述可以这样理解。首先我们考虑正都是无记忆的稳恒信源, 即

$$P(w_i, w_{i+1}, \dots, w_{i+N} = A) = P(w_j, w_{j+1}, \dots, w_{j+N} = A)$$

其中 A 为固定的字母序列。(这种是离散平稳随机过程。那么对于某一随机变量的信息量即为 $-\log p(x_1, x_2, \dots, x_n)$)

当 n 趋于无穷大时, 则有弱大数定律。成为等概率情况。 $-\frac{1}{n} \log p(x_1, x_2, \dots, x_n)$ 即为平均一个码元携带的信息量。那么即为信息熵。偏离这个等概率分布的概率很小。

因此我们就可以给出一个典型序列。我们给出数学上的严格定义。

相对于分布 $P(x)$ 和序列 $(x_1, x_2, \dots, x_n) \in X_n$. 典型序列集合 $A_\epsilon^{(n)}$ 定义为满足下列不等式约束的所有序列集合

$$2^{-n(H(X)+\epsilon)} \leq P(x) = p(x_1, x_2, \dots, x_n) \leq 2^{-n(H(X)-\epsilon)}$$

这个其中也是由弱大数定律推导而来。

$$P\left(\left|\frac{S_n}{n} - p\right| \geq \epsilon\right) < \delta \rightarrow \log P(w) = \sum_{k=1}^K n_k \log p(a_k) \Rightarrow n \cdot \sum_{k=1}^K \frac{n_k}{n} \log p(a_k)$$

我们再用正数定理, 即有

$$2^{-n(H(X)+\epsilon)} \leq P(x) = p(x_1, x_2, \dots, x_n) \leq 2^{-n(H(X)-\epsilon)}$$

由上述的证明过程和定义, 我们发现典型序列满足以下性质。

1. 若 $x \in A_\epsilon^{(n)}$, 则 $H(X) - \epsilon \leq -\frac{1}{n} \log p(x) \leq H(X) + \epsilon$.

2. 若 n 趋向无穷大, $P(A_\epsilon^{(n)}) \rightarrow 1 - \epsilon$. (这里我们给出一个证明方法, 见附录)

$$2^{-n(H(X)-\epsilon)} \leq |A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$$

2. 若 \$n\$ 足够大, \$Pr\{A_\epsilon^{(n)}\} > 1 - \epsilon\$. (这里我们给出一个证明方法, 见附录)

3. $(1 - \epsilon) 2^{n(H(x) - \epsilon)} \leq |A_\epsilon^{(n)}| \leq 2^{n(H(x) + \epsilon)}$

性质1和性质2可以直接由定义得到. 简单证明一下性质3

$$1 = \sum_{x^n \in X^n} P(x^n) \geq \sum_{x^n \in A_\epsilon^{(n)}} P(x^n) \geq \sum_{x^n \in A_\epsilon^{(n)}} 2^{-n(H(x) + \epsilon)} = 2^{-n(H(x) + \epsilon)} |A_\epsilon^{(n)}|$$

$$1 - \epsilon \leq Pr\{A_\epsilon^{(n)}\} \leq \sum_{x^n \in A_\epsilon^{(n)}} 2^{-n(H(x) - \epsilon)} = 2^{-n(H(x) - \epsilon)} |A_\epsilon^{(n)}|. \Rightarrow |A_\epsilon^{(n)}| \geq (1 - \epsilon) 2^{n(H(x) - \epsilon)}$$

由此性质3即证明完毕.

典型序列的个数大约为 $2^{n(H(x) + \epsilon)}$ 个, 而 $|X|^n$ 有很多. 典型序列的集合只是所有可能序列的一个子集, 但好处是他几乎全是典型序列. 致此, 我们完成了几乎没有失真. 从非等长码到等长码的过程, 我们可以举于这个达到无失真的变长编码定理.

2. 变长编码定理.

变长编码定理就是选重把 \$N\$ 取得很大, 从而使序列几乎都是典型序列. 从而在出错率较小的情况下, 尽可能的压缩码长. 下面我们给出定理

设 X^n 是由独立同分布离散随机变量 $x \sim p(x)$ 构成的序列. 对于任意正数 ϵ , 总有足够大的 n , 可以找到一个一一映射, 将 X^n 映射成二进制序列, 且满足:

$$E\left[\frac{1}{n} L(x^n)\right] \leq H(x) + \epsilon.$$

下面我们简单证明并给出一个例子. 最后总结一下.

我们的编码的平均码长为 $E[L(x^n)]$

$$E[L(x^n)] = \sum_{x^n \in X^n} P(x^n) L(x^n) = \sum_{x^n \in A_\epsilon^{(n)}} P(x^n) L(x^n) + \sum_{x^n \in \bar{A}_\epsilon^{(n)}} P(x^n) L(x^n)$$

$$\leq \sum_{x^n \in A_\epsilon^{(n)}} P(x^n) \cdot [n(H(x) + \epsilon) + 1 + 1] + \sum_{x^n \in \bar{A}_\epsilon^{(n)}} P(x^n) [n \log |X| + 2]$$

$$= Pr\{A_\epsilon^{(n)}\} [n(H(x) + \epsilon) + 2] + Pr\{\bar{A}_\epsilon^{(n)}\} [n \log |X| + 2]$$

$$\leq n(H(x) + \epsilon) + 2 + n\epsilon \log |X| + 2\epsilon$$

$$= n(H(x) + \epsilon'). \quad \text{其中 } \epsilon' = \epsilon + \epsilon \log |X| + \frac{2\epsilon}{n}$$

解释说明:

平均码长可分为典型序列和非典型序列之和

套用典型序列的定义, 进行缩放. 非典型序列

直接缩放到最大值.

$+1, +1$ 是由二进制取整保持误差为典型序列

代入定义, 继续缩放.

可见 $n \rightarrow \infty$ 时, $\epsilon' \rightarrow 0$.

于证得证.

!

取一个实际的例子.

设信源传输符号分布为

$$\begin{pmatrix} X \\ P(x) \end{pmatrix} = \begin{pmatrix} a_1 & a_2 & a_3 & a_4 & a_5 & a_6 & a_7 & a_8 \\ 0.40 & 0.18 & 0.10 & 0.10 & 0.07 & 0.08 & 0.05 & 0.04 \end{pmatrix}$$

要求编码错误概率小于 10^{-8} , 且编码效率大于 90%. 求最小可行分组长度.

首先验证一下错误概率, 也就是非典型序列的概率.

$$Pr\{\bar{A}_\epsilon^{(n)}\} \leq \frac{D(I(x))}{N\epsilon^2}$$

这个公式目前不知道怎么计算.

接着, 我们定义编码效率. 编码效率的极限自然是完全去掉冗余度, 即使条件时复杂, 信息冗余也尽可能不会减小.

那么极限自然就是 $H(x)$ 即 $H(x)$. 用 P_n 我们所能达到的则是 $H(x) + \epsilon'$. 这在上面的证明中也给出了. 编码效率平均为

$$\frac{H(x)}{H(x) + \epsilon'}, \quad \frac{H(x)}{H(x) + \epsilon'} > 0.9 \Rightarrow \epsilon < 0.28 \Rightarrow Pr\{\bar{A}_\epsilon^{(n)}\} \leq 10^{-8} \Rightarrow N > 10^8.$$

上可证其 < 属于熵的熵理, $\epsilon < \epsilon'$.

$$H(X) + \epsilon, H(X) + \epsilon > 0. \Rightarrow \ll 0.28 \implies N > 10^0$$

关于定长编码定理的证明过程:

首先我们需要知道为什么会有定长编码. 为此我们需要有一个码率的极限, 这个极限描述的是随着码组长度 \$N\$ 的增长, 平均每个字母的码长 \$H(N)\$ 的极限值. 数学上即 $H(U) = \lim_{N \rightarrow \infty} \frac{1}{N} H(U_1, U_2, \dots, U_N)$. 而对于离散平稳信源, 有一条定理 $H(U) = \lim_{N \rightarrow \infty} \frac{1}{N} H(U_1, U_2, \dots, U_{N-1})$, 证明过程可以参考教材. 较为核心几点如下:

根据离散平稳信源的定义

$$H(U_N | U_2, U_3, \dots, U_{N-1}) = H(U_{N-1} | U_1, \dots, U_{N-2})$$

这体现了时间的无记忆性, 而根据条件熵的性质:

我们可以得到每一字母的码长随着 \$N\$ 的增大而减小.

$$NH(N) = H(U_1) + H(U_2 | U_1) + \dots + H(U_N | U_1, U_2, \dots, U_{N-1})$$

$$\geq NH(U_N | U_1, U_2, \dots, U_{N-1})$$

我们可以知道, 信源输出中前后字母之间的统计关系是使信息熵减小的重要原因.

另一方面: $NH(N) = H(U_N | U_1, U_2, \dots, U_{N-1}) + (N-1)H_{N-1}(U)$, 故 $NH(N) \geq H(U_N | U_1, U_2, \dots, U_{N-1})$

可知得到 $H(N) \leq H_{N-1}(U)$, 下证有界.

我们可知 \$N\$ 减少时, 码长增大. 由熵的极值性可知 $H(U) \leq \log k$, \$k\$ 为信源字母表进数.

因此, 当信源的一维分布不均匀时, 即又不存在冗余关系, 也不能无限增大码长以携带信息.

那么反映信源有效程度的物理量就是冗余度和相对冗余度.

冗余度: $\log k - H(U)$, 相对冗余度: $1 - \frac{H(U)}{\log k}$. 我们进行定长编码就是使非独立变为独立, 非等概率变为等概率的过程.

关于定长编码定理我们还有一个形式的定理:

设信源为 \$H(U)\$ 的离散无记忆信源被分成 \$N\$ 个字母组, 并用长为 \$M\$ 的码字母组进行编码, 字母字母表大小为 \$J\$. 则对于任意 \$\epsilon > 0\$ 和 \$\delta > 0\$, 只要 \$N\$ 足够大, 且满足不等式

$$\frac{M}{N} \log J > H(U) + \delta$$

则信源字母组没有自己特有的码字, 其码平均长度 \$p_0\$ 可小于 \$\epsilon\$.

这个定理证明比较简单.

$$N \epsilon = 2^{N(H(U) + \delta)}$$

我们由编码定理可知, 那么冗余度的进数大于典型序列数.

$$J^M \geq 2^{N(H(U) + \delta)} \implies M \log J \geq N(H(U) + \delta) \implies \frac{M}{N} \log J \geq H(U) + \delta$$

但是定长编码定理只是证明了理论上的可行性. 实际上, 正如阿诺斯所, 一个较好的定长编码需要 \$10^8\$ 分组长度. 这无论何时都远, 这是不现实的. 都是无法接受的. 我们也无法找到一个满意的折衷方案. 这也证明了应用范围不广的弊端.

3. 变长码

3.1 变长码的类型

研究无失真编码一定要保证是一一映射, 那么我们称为非奇异码.

若一个码 \$C\$ 可将不同的 \$X\$ 映射为不同的 \$D^*\$ 的序列, 即

$$X = X' \implies C(X) \neq C(X')$$

则称该码是非奇异的.

但是即使是——映射也会有麻烦. 例如, 对于不等长的编码

$$C(x_1) = 0, C(x_2) = 1, C(x_3) = 01$$

01 码的解码出现了歧义. 那么我们为了改进见有3位-可译码, 如下:

称码 \$C^+\$ 是码 \$C\$ 的扩展. 当 \$C^+\$ 是有限长 \$X\$ 序列到有限长 \$D^+\$ 序列的映射, 且满足

$$C^+(x_1 x_2 \dots x_n) = C^+(x_1) C^+(x_2) \dots C^+(x_n)$$

那么其扩展是非奇异的, 我们则称为唯一可译码.

这其实是一个定义上的解答. 并没有可构造性的方案. 换言之, 没有码字是码字的组合

例如 $C(x_1) = 0, C(x_2) = 01, C(x_3) = 1$. 则序列 010110010 唯一可译为 $x_1 x_1 x_3 x_1 x_2 x_1$

从冗余性看 \$P\$ 的计算我们可以台视. 但是有一个逐字译码直接译码. 按码字译码.

这其实是一个定义上的解法，并没有可建设性的元素、换言之，没有码字是码字的组合

例如 $C(x_1) = 0$, $C(x_2) = 01$, $C(x_3) = 11$. 树序列 01110010 唯一可译为 $x_1 x_1 x_3 x_1 x_2 x_1$

但是经过自己的计算我们可以发现，这是一个译字并重建译的过程，需要回退。

虽然唯一可译码可以解决译码的问题，但对于解码效率的要求，于是我们需要一种速度更快的码，称为即时码，也称为前缀码。

当没有码字是其他码字的前缀，那么则称为前缀码或即时码。

例如: $C(1) = 0$, $C(2) = 10$, $C(3) = 110$, $C(4) = 111$. 树序列 01101111010 直接可译为 12432.

这是一种很好性质的码。

简单提一下各码之间的关系:

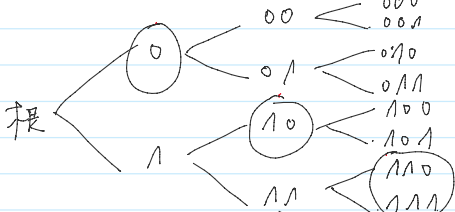
即时码 \subset 唯一可译码 \subset 非奇异码 \subset 所有码。

某些非前缀码不是唯一可译码，需要发报者的手译。

汉语: 没有鸡鸭鱼也可没有鱼肉也可青菜豆腐不可少。这是歧义句。没有汉语没有标点也不是唯一可译码。

3.2 Kraft不等式与香农第一定理

我们已经知道了前缀码具有较好的性能，我们从定义中也可以直接得到，前缀码不能使用最短的码字，否则前缀的条件无法确定，我们可以用二进数来表示。



我们很容易发现，无论是定长还是变长都对应树的叶子。例如，我们可以这样选取前缀码。

那么我们就需要回答如何构造前缀码？前缀码存在什么条件？

实际上通过树图，我们已经可以很直观的发现其中的规律，下面我们每次使用数学语言来描述。

由于前缀码是在树上，D进制树为D叉树。若被选为码字的叶子不会被删除，设其码长度为 l_i ，最长为 l_{max} 。所以有

$$D^{-l_i} = D^{-l_{max}} \cdot D^{l_i + l_{max}}$$

于是，按含有 k 个不同字母的信源要用 J 个字母的码字母表进行变长编码，则当且仅当各码字长度 l_1, l_2, \dots, l_k 满足 Kraft 不等式

$$\sum_{i=1}^k J^{-l_i} \leq 1.$$

我们已经证明了必要性，关于充分性的证明差不多，我们只是要砍掉其他的叶子，所以这个过程一定是可行的。所以这一页是个必要条件，但是仅仅从编码的构造考虑，与码是否存在无关，所以不用再考虑最优性。

关于码的长度，与变长编码一致，有如下定理。

对于随机变量 X 进行 D -进制前缀码编码，得到码长函数 $L \geq H_D(X)$ ，等号当且仅当 $D^{-l_i} = p_i$ 时成立。

证明如下:

$$L - H_D(X) = \sum p_i l_i - \sum p_i \log_D \frac{1}{p_i}$$

$$\begin{aligned} &= \sum p_i \cdot \log_D D^{l_i} - \sum p_i \log_D \frac{1}{p_i} = \sum p_i \log_D p_i^{-l_i} - \sum p_i \log_D D^{-l_i} \\ &= \sum p_i \log_D p_i^{-l_i} - \sum p_i \log_D \frac{D^{l_i}}{D^{l_i}} = \sum p_i \log_D p_i^{-l_i} - \sum p_i \log_D \frac{D^{l_i}}{D^{l_i}} - \log_D D^0 \end{aligned}$$

其中 $\Delta = \sum D^{-l_i}$ ，那么上式变为

$$L - H_D(X) = \sum p_i \log_D \frac{p_i}{D^{-l_i}} + \log_D \frac{1}{\Delta}$$

根据熵的信道对数求和不等式，第一项大于0，由 Kraft 不等式，第二项也大于0，用 p_m 证明

与是我们找到了一个下界，我们发现当 $D^{-l_i} = p_i$ 时，可以直接取等。

这里稍微解释一下这个定理，对于一个有 k 个字母的信源，每个字母出现的概率为 p_i ，那么这个信源的熵可以理解为他的熵带的信息量。而我们的信道长度为 L 是一个大于等于 $H_D(X)$ 的数，所以 $L - H_D(X)$ 又叫做冗余度。一个离散信源的熵带分布 (2个字母)

与是我们找到了一个下界。我们发现当 $D^{-1} = p_i$ 时，可以直接取等。
 这里稍微解释一下这个原理，对于一个有 k 个符号的信源，每个符号出现的概率为 p_i ，那么这个信源的熵可以为他所携带的信息量。而我们的二进制编码本质上是一个 D 叉树，于是我们利用了这个 D 叉树实际上直接反了一个新的概率分布。这个新的等概率性一定会比之前的更复杂，至少也是相同。所以我们就有了这样一个结论。
 实际上这个结论也告诉了我们一个道理，要让这个平均长度也等于这个新的分布信源的熵，我们更应该让他的分布更可能的靠近原始分布。这样可以让信源的熵也随着下降。
 那么最佳的码长是多少呢？有以下定理。这个定理也被称为香农第一定理。

对随机变量 X 进行 D -进制前缀编码，得到的最佳码长不等式满足下列不等式

$H_D(X) \leq L^* \leq H_D(X) + 1$
 $H_D(X)$ 对于非独立的情况应严格记为 $H_D(W) / \log D$
 这个定理的证明也是非常巧妙的，香农也构造了一个香农码。其码长为 $\lceil \log \frac{1}{p_i} \rceil$ 。那么

$$\sum D^{-\lceil \log \frac{1}{p_i} \rceil} \leq \sum D^{-\log \frac{1}{p_i}} = \sum p_i = 1$$

所以满足 Kraft 不等式。一定是一个唯一可译码。

$$\log_D \frac{1}{p_i} \leq \log_D \lceil \frac{1}{p_i} \rceil \leq \log_D \frac{1}{p_i} + 1$$

$$\Rightarrow H_D(X) \leq \sum p_i l_i \leq H_D(X) + 1$$

$$\Rightarrow H_D(X) \leq L^* \leq H_D(X) + 1$$

最佳码一定不能比这个码更差。所以巧妙证明的结论。

那么定长编码告诉我们我们可以做到平均码长和之前相比只有 ϵ 的误差。而变长编码则有 1 比特。
 我们还可以通过两者结合做的更好。

对于信源 X 进行分组前缀编码，得到每消息符号数 L_n^* 满足不等式

$$\frac{H(X_1, X_2, \dots, X_n)}{n} \leq L_n^* \leq \frac{H(X_1, \dots, X_n)}{n} + \frac{1}{n}$$

这个过程就是独立分组。然后进行前缀编码。若信源独立，那么 $L_n^* \sim H(X)$ 。

如果我们对信源分布的估计出了偏差，那么性能一定不会变好。因为鉴别信息的非负性。愿意的分析有如下定理

对于服从 $p(x)$ 信源 X 进行前缀编码，若码字长度取 $l(x) = \lceil \log \frac{1}{q(x)} \rceil$ ，则平均码长满足

$$H(p) + D(p||q) \leq E_p[l(x)] < H(p) + D(p||q) + 1$$

因为 $H(q) = H(p) + D(p||q)$ ，所以这个定理实际上并不需要证明。

3.3. 几种前缀变长码分析

3.3.1 香农码

香农码在一般情况下并非最优情况。

我们举一个简单例子。

$$k=2, p(k_1) = 0.9999, p(k_2) = 0.0001, \text{ 那么 } l_1 = 1, l_2 = 14$$

我们可以计算一下。

$$H(k) = p(k_1) \log \frac{1}{p(k_1)} + p(k_2) \log \frac{1}{p(k_2)} = 1.47 \times 10^{-3}$$

$E(l_{\text{Shannon}}) = 0.9999 \times 1 + 0.0001 \times 14 = 1.0013$ 。但是我们可以用 1 比特精度，且他离最佳的解差很远。

主要原因也是由局尺。当对 $\log \frac{1}{p_i}$ 向上取整时，误差可能会非常大。这时我们还有办法减小码字很大的码长。

3.3.2 Huffman 编码

Huffman 作为于 1952 年提出，可以证明 Huffman 编码具有最优性。即在给定信源符号出现的概率分布和字母因表的前提下，没有其他码长序列比 Huffman 编码更短的平均码长。

先抛开 Huffman 编码不谈，我们对上面例子进行更深一点的分析。

$$H(k) = p(k_1) \log \frac{1}{p(k_1)} + p(k_2) \log \frac{1}{p(k_2)} = p(k_1) \cdot 1.44 \times 10^{-3} + p(k_2) \cdot 13.2$$

$$E(l) = p(k_1) l_1 + p(k_2) l_2$$

通过二叉树我们可以发现，让 k_2 的长度可以为 1，最后得到长度为 1 的平均码长。

当然，这个过程违背了鉴别信息对于平均码长的理论指导。但是，我们经过人为的修正，可以更少的减少误差。Huffman 事被这个假说让概率大的信源码长更可能短。当然前提是不能有重复。

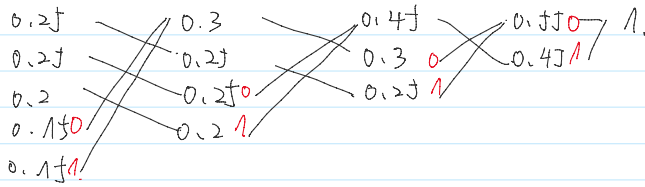
这里可以理解为：当信源的概率分布为 p_i 时，我们总可以找到一个 q_i 使得 $p_i \leq q_i$ 且 q_i 是 2^{-n} 的倍数。

当然,这个过程违背了熵的信息对于平均码长的理论指导,但,我们经过人为的修正,可以更少的减少误差, Huffman 就是通过这种方法让概率大的信源的码长更可能短,当然前提是有能短码长。

编写过程实际上是非常简单的,按概率由大到小分布,然后进行交叉对排列,我们举一例来说明。

考虑 $\begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 0.25 & 0.25 & 0.2 & 0.15 & 0.15 \end{pmatrix}$

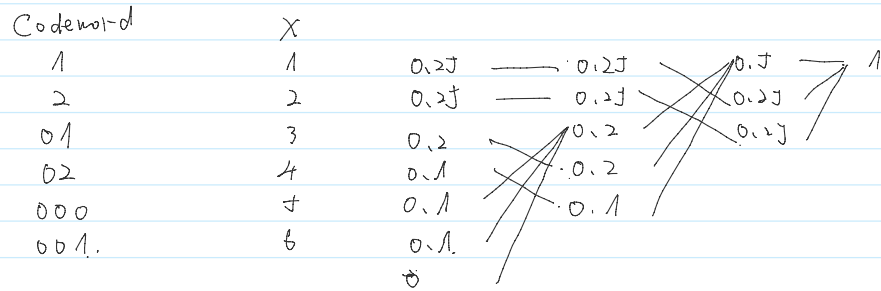
Huffman 编码过程:



最后得到

1 → 01
2 → 10
3 → 11
4 → 000
5 → 001

二进制有三个合并, 这里出现了问题, 可以从调整一些分支后, 如下所示4种情况, 每次选择缩减减少 $(D-1)$, 由于最大的一支缩减, 所以由 D 缩减为 $D-1$, 填充后的进位一变为 $D-1$ $k(D-1)$



Huffman 的精确在于他最优性, 我们接下来重新从最简单的二元情形证明他, 这个证明也比较容易推广到其他形式, 我们稍后会分析也发现了, 对于二元情形, 我们要让他尽可能短, 一般是依次排列, 但概率最小的码字尽可能长, 于是就有两条定理

1. 对任一给定的离散无记忆信源, 存在一个最优的二元前缀码, 这个码中最小发生的两个码字必定具有相同的长度, 且码中相同长度的码字有两个或两个以上时, 其中必有两个码字的差别只在最后一位。
2. 设 C' 是某信源经缩减后得到的缩减信源的最优前缀码, 将 C' 中由原信源的较小概率的两个字码缩减得到的字码所对应的码字后各加 0 和 1, 作为原信源的最小概率的两个码字, 而其余码字不变, 则这样得到的码 C 对原信源也是最优的。

第一个定理证明可以采用反证法, 第二个定理可以用第一个定理得到, 详细的证明过程可以参见《应用信息论卷二》。

3.3.3 算术码

在之前的证明提到, 当 $-\log_2 p(a_k)$ 不为整数时, 我们希望通过扩展信源的办法来改进它的效果, 即需要分组, 这自然会使得编码更复杂, 我们可以用树形来解决这个问题。

对于离散无记忆信源树形码之一, 一种长码编码方法是出自 P. Elias 最早在一篇未发表的研究报告中提到, 后经 J. Jelinek, J. Rissanen 等的改进和发表投入实用, 现在人们一般称之为算术码 (arithmetic coding)。

3.3.3.1 过程描述:

设信源符号集为 $\{a_1, a_2, \dots, a_k\}$, 符号 a_k 的概率为 $p(a_k)$ ($k=1, 2, \dots, k$), 将符号按降序排列并记作 $a_1 > a_2 > \dots > a_k$, 于是, 可以定义 X 个符号 a_k 的累积概率

$$F(a_k) = \sum_{a_i \geq a_k} P(a_i)$$

再经 X 修正概率

$$\bar{F}(a_k) = \sum_{a_i \geq a_k} P(a_i) + \frac{1}{2} p(a_k)$$

那么我们会得到一个概率值 $F(a_k)$ 和 $\bar{F}(a_k)$, 将这个值用相应的进制表示, 则可以得到码字, 举一个例子。

a_k	$p(a_k)$	$F(a_k)$	$\bar{F}(a_k)$	$\bar{F}(a_k)$ 二进制表示	码字	Huffman 编码
a_1	0.25	0.25	0.125	0.001	001	10
a_2	0.5	0.75	0.5	0.10	10	0
a_3	0.125	0.875	0.8125	0.1101	1101	110
a_4	0.125	1	0.9375	0.1111	1111	111

我们可知看到，对于信源序列直接进行编码效果不是很好，这不如 Huffman 编码。这是因为变长码是修正的 F(ak) 代表也是一个线段长度，所以是有冗余的，也有对信源做扩展。还是当线段长度缩短，当信源长度趋于无限时，P 所在在区间成为一点，开始冗余则随着长度增加而降低。

举一个两码分布的例子 $X \sim B(0.6)$,

设信源在输出第 $n-1$ 个符号后，P 所在区间为 (A_{n-1}, B_{n-1}) ，那么有

$$\begin{cases} A_n = A_{n-1} \\ B_n = A_{n-1} + 0.6(A_{n-1} - B_{n-1}) \end{cases} \quad \begin{matrix} m_n = a_0 \\ m_n = a_1 \end{matrix}$$

我们可以看出随着 n 增大区间越来越小，而收敛于一个点。

当信源输出第 N 个信源符号时，码字输出结束此时 A_N 和 B_N 中表示一致的部分。此时输出二进制数的位数。二元码的码长应为 M 。

$$\frac{1}{2^M} > \prod_{m=1}^N p(m_n) > \frac{1}{2^{M+1}}$$

3.3.3.2 极限情况下的压缩效果。

对于长的信源符号序列，同样可以加码排序。

例如若有 $V_N = v_1 v_2 \dots v_N$, $V'_N = v'_1 v'_2 \dots v'_N$
若存在 $t \leq N$, 使

$$v_n = v'_n \quad \text{当 } n < t \text{ 时}$$

$$v_t > v'_t \quad \text{当 } n = t \text{ 时}$$

则记 $V_N > V'_N$ ，我们对此信源符号序列进行排序。可定义 V_N 的累积概率为

$$F(V_N) = \sum_{V_N > V'_N} p(V'_N)$$

视该序列 V_N 后紧接着 v_{N+1} ，形成新序列 $V_{N+1} = V_N v_{N+1}$ ，则此序列的概率为

$$p(V_{N+1}) = p(V_N) p(v_{N+1})$$

此时的累积概率为

$$\begin{aligned} F(V_{N+1}) &= \sum_{V_{N+1} > V'_{N+1}} p(V'_{N+1}) \\ &= F(V_N) + p(V_N) \sum_{v_{N+1} > v'_{N+1}} p(v'_{N+1}) \\ &= F(V_N) + p(V_N) F(v_{N+1}) \end{aligned}$$

由于累积概率随着序列的平滑增加，所以只要将区间

$$[F(V_N), F(V_N) + p(V_N) F(v_{N+1})]$$

用有限精度的二进制数加以表示，并将此二进制数与 V_N 加以对应。就可以实现无损压缩。由于这种情况下存在冗余条件仍然是 Kraft 不等式，所以在极限情况下，算术码可以实现对信源的无损压缩。

3.4 离散马尔可夫信源的熵率

关于马尔可夫信源的基本内容这里不在描述，参阅随机信号处理。

求熵率首先要求条件熵。

设过某时刻的状态为 $S_1 = i$ ，现在要求 n 时刻信源输出符号的熵

$$H(U_n | U_1 U_2 \dots U_{n-1}, S_1 = i)$$

$$= - \sum_{U_1 U_2 \dots U_n} p(U_1 U_2 \dots U_n | S_1 = i) \log p(U_n | U_1 U_2 \dots U_{n-1}, S_1 = i)$$

由马尔可夫信源的性质可知。

1. 某时刻信源输出哪个符号与此时信源所处的状态有关，而与以前的状态以及以前的输出符号均无关。

2. 信源某时刻所处的状态，由当前输出符号和前一时刻信源的状态可唯一确定。

那么，我们可以得到

$$p(U_n | U_1 U_2 \dots U_{n-1}, S_1 = i) = p(U_n | S_n)$$

$$P(U_1, U_2, \dots, U_n | S_1 = i) = P(U_1, U_2, U_{n-1}, S_n | S_1 = i) \cdot P(U_n | S_n)$$

将这两个式子代入原式得：

$$H(U_1, U_2, \dots, U_{n-1} | S_1 = i)$$

$$= - \sum_{S_n=U_1, U_2, \dots, U_n} P(S_n, U_1, U_2, \dots, U_{n-1} | S_1 = i) \cdot P(U_n | S_n) \log(P(U_n | S_n))$$

$$= \sum_{j=1}^S P(S_n = j | S_1 = i) H(U_n | S_n = j)$$

那么我们将代入所有的 S_1 即可得到条件熵

$$H(U_n | U_1, U_2, \dots, U_{n-1}, S_1) = \sum_{j=1}^S P(S_n = j) H(U_n | S_n = j)$$

对于平稳马尔可夫链， $P(S_n = j)$ 与 S_1 的根号分布无关，为平稳分布 $P(j)$ ，且 $H(U_n | S_n = j)$ 与 n 无关，可以记做 $H(U | S = j)$ 。那么

$$H(U_n | U_1, U_2, \dots, U_{n-1}, S_1) = \sum_{j=1}^S P(j) H(U | S = j)$$

有了条件熵之后，我们就可以根据可加性求熵

$$\begin{aligned} \frac{1}{N} H(U_1, U_2, \dots, U_N | S_1) &= \frac{1}{N} \sum_{n=1}^N H(U_n | U_1, U_2, \dots, U_{n-1}, S_1) \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^S P(j) H(U | S = j) \\ &= \sum_{j=1}^S P(j) H(U | S = j) \end{aligned}$$

如果不是平稳变化的，可定义

$$P_{(1, N)}(j) = \frac{1}{N} \sum_{n=1}^N P(S_n = j)$$

代入可得

$$H(U_1, U_2, \dots, U_N | S_1) = \sum_{j=1}^S P_{(1, N)}(j) H(U | S = j)$$

最后我们可以求 $H_{\infty}(U)$ 。

$$H(U_1, U_2, \dots, U_N) = I(S_1; U_1, \dots, U_N) + H(U_1, U_2, \dots, U_N | S_1)$$

其中

$$I(S_1; U_1, \dots, U_N) = H(S_1) - H(S_1 | U_1, U_2, \dots, U_N) \leq H(S_1) \leq \log S$$

代入 $H_{\infty}(U)$ 的表达式中

$$\begin{aligned} H_{\infty}(U) &= \lim_{N \rightarrow \infty} \frac{1}{N} H(U_1, U_2, \dots, U_N) \\ &= \lim_{N \rightarrow \infty} \left[\frac{1}{N} I(S_1; U_1, \dots, U_N) + \frac{1}{N} H(U_1, U_2, \dots, U_N | S_1) \right] \end{aligned}$$

我们可知互信息 $I(S_1; U_1, \dots, U_N)$ 是有界的，所以 $\lim_{N \rightarrow \infty} \frac{1}{N} \log S = 0$ 。

$$H_{\infty}(U) = \lim_{N \rightarrow \infty} \frac{1}{N} H(U_1, U_2, \dots, U_N | S_1) = \sum_{j=1}^S P(j) H(U | S = j)$$

高维马尔可夫信源的熵率是由条件熵组成的，由条件熵求熵可知马尔可夫信源的熵率比之前提到的离散无记忆信源的熵率更小。

马尔可夫信源之所以这么重要，也是因为生活中很多情形下都可以用这个数学模型表达，例如天气预报、语言文字等。若马尔可夫信源输出字母根号与前面 m 个字母输出字母有关，与更前面的输出字母无关，则称之为 m 阶马尔可夫信源。若信源输出字母的总数为 k ，则一般情况下， m 阶马尔可夫信源的总数 $S = k^m$ 。显然 m 越大，对马尔可夫信源的数学模型就越复杂。

如何选种恰当阶数 m 的马尔可夫信源来做实际信源在工程上需要考虑的重要问题，例如英文字母通常认为有 26 个字母，如果 m 选太大或太小，一般取在 3 到 5 之间。

3.5 高维马尔可夫信源的编码原理与最优编码

我们之前提到一般平稳遍历的信源在编码原理中可以取现型在信源，这也就是成立的。

对于非平稳遍历的信源，其熵率 $H_{\infty}(U)$ 与 $H(U)$ 不同，且 $H(U) < H_{\infty}(U)$ 。对于非平稳遍历的信源，其熵率 $H_{\infty}(U)$ 与 $H(U)$ 不同，且 $H(U) < H_{\infty}(U)$ 。

3.3 熵与平均码长的关系

我们之前提到，一般平均码长的信源在编码处理时可以近似理想压缩，这也就是成立的。

对于变长编码，我们也可以简单地将离散无记忆信源的变长编码推广到离散马尔可夫信源。

设马尔可夫信源为离散，状态为 S_1, \dots, S_N ，信源输出序列为 $U = U_1 \dots U_N$ ，每一序列可用变长编码的方法得到一个对应的码字，其码字长度 $L(U)$ 满足。

$$J^{-L(U)} \leq P(U_1 U_2 \dots U_N | S_1 = i) \leq J^{-L(U)+1}$$

此时有

$$\sum_U J^{-L(U)} \leq \sum_U P(U_1 U_2 \dots U_N | S_1 = i) = 1$$

那么所有码字满足 Kraft 不等式，我们应用香农第一定理。

$$\frac{H(U_1 U_2 \dots U_N | S_1 = i)}{\log J} \leq \bar{L} \leq \frac{H(U_1 U_2 \dots U_N | S_1 = i)}{\log J} + 1$$

平均到每个字符则是

$$\frac{H(U_1 U_2 \dots U_N | S_1)}{N \log J} \leq \bar{L} \leq \frac{H(U_1 U_2 \dots U_N | S_1)}{N \log J} + \frac{1}{N}$$

让 N 趋于无穷，代入熵率。

$$\lim_{N \rightarrow \infty} \frac{1}{N} H(U_1 U_2 \dots U_N | S_1) = \sum_{j=1}^S P(j) H(U=j)$$

即

$$H_{\infty}(U) = \sum_{j=1}^S P(j) H(U=j)$$

于是有

$$\frac{H_{\infty}(U)}{\log J} \leq \bar{L} \leq \frac{H_{\infty}(U)}{\log J} + \frac{1}{N}$$

就可以得到马尔可夫信源的变长编码定理。

当用 J 个字母的码字表对以熵率为 $H_{\infty}(U)$ 的离散马尔可夫信源进行变长编码时，其平均码长满足

$$\frac{H_{\infty}(U)}{\log J} \leq \bar{L} \leq \frac{H_{\infty}(U)}{\log J} + \frac{1}{N}$$

其中 N 是信源字码分组的长度。

那么当 N 足够大时， \bar{L} 是可以接近 $H_{\infty}(U)/\log J$ ，从而达到理想的压缩。

附录

1.

$$P(X_n \in A_\epsilon) > 1 - \epsilon.$$

$$\text{设 } Z_n = -\frac{1}{n} \log P(X_n) = -\frac{1}{n} \log \left[\prod_{i=1}^n P(X_i) \right] = -\frac{1}{n} \sum_{i=1}^n \log P(X_i)$$

$$\Rightarrow E(Z_n) = -\frac{1}{n} \cdot n \cdot H(X_k) = H(X)$$

$$\text{var}(Z_n) = \frac{1}{n} \text{var}(-\log_2 P(X))$$

$$P(X_n \in A_\epsilon) = P\left(-\frac{1}{n} \log_2 P(X_i) - H(X) \leq \epsilon\right), \text{ 由弱大数定律可得，上面我们也就得到它。}$$

$$= P(|Z_n - E(Z_n)| \leq \epsilon)$$

二、由大数定律可得，上面我们也就得到它。

$$= P(|z_n - E(z_n)| \leq \varepsilon)$$

应用马尔可夫不等式或者切比雪夫不等式:

$$P(X_n \in A_\varepsilon) = 1 - P(|z_n - E(z_n)| \geq \varepsilon) \geq 1 - \frac{\text{Var}(z_n)}{\varepsilon^2} = 1 - \frac{1}{n} \cdot \frac{\text{Var}(-\log_2 P(X_1))}{\varepsilon^2}$$

我们想要 $P(X_n \in A_\varepsilon) > 1 - \varepsilon$, 则

$$\frac{1}{n} \cdot \frac{\text{Var}(-\log_2 P(X_1))}{\varepsilon^2} < \varepsilon \Rightarrow n > \text{Var}(-\log_2 P(X_1)) / \varepsilon^3.$$

则有 $P(X_n \in A_\varepsilon) > 1 - \varepsilon$.