

第4章 情報源符号化とその限界

4.1 情報源符号化に必要な条件

異なる情報源記号にに対して同じ符号語が割り当てられている符号: 特異符号 (singular code)

一意復号可能な符号 \longleftrightarrow 一意復号可能な符号 等長符号 \longleftrightarrow 非等長符号

瞬時符号: 符号語のパターンが現れたとき、それを直ちに復号できる符号 \longleftrightarrow 非瞬時符号

特異でない等長符号は瞬時符号である。

コンマ符号: コンマ (comma) の役割をする符号。符号系列中、これだけを見れば、符号語の区切りを覚えることができる。(A:0 B:10 C:110 D:1110, 0はコンマ符号)

情報源符号化に必要な条件

(i) 一意復号可能であること、瞬時符号であることが望ましい (iii) 装置化があまり複雑とはならない。

(ii) 1 情報源記号当たりの平均符号長ができるだけ短い。

4.2 瞬時符号と符号の木

x: 01 y: 011 xはyの語頭 (prefix) である。語頭条件: どの符号語も他の符号語の語頭となてはならない。

符号の木: 符号語に対応する節点を白丸で表し、それ以外の節点を黒丸で表している

瞬時符号の符号語はすべて葉に対応づけられている。

4.3 クラフトの不等式 (Kraft's inequality)

長さ l_1, l_2, \dots, l_m とする M 個の符号語を持つ q 元符号が瞬時符号となるための必要かつ条件は

$$q^{-l_1} + q^{-l_2} + \dots + q^{-l_m} \leq 1 \text{ が成立することである。 (一意復号可能な符号が存在するための必要かつ条件でもある。マウランの不等式と呼ばれる)}$$

4.2 平均符号長の限界

$P(a_i) = p_i \ (i=1, 2, \dots, M)$, 1 情報源記号当たりの平均符号長は $L = l_1 p_1 + l_2 p_2 + \dots + l_m p_m$

一意復号可能な q 元符号に符号化したとき、平均符号長は $L \geq H(S)$ 。また、 $L < H(S) + 1$ とする瞬時符号を作ることができる。 $H(S)$ は S の 1 次エントロピーと呼ばれる量である

$$H_1(S) = -\sum_{i=1}^M P(a_i) \log_2 P(a_i) = -\sum_{i=1}^M p_i \log_2 p_i$$

シャノンの補助定理: $q_1, q_2, \dots, q_m \in \mathbb{N}$ と $q_1 + q_2 + \dots + q_m \leq 1$ を満たす任意の非負の数 p_i と q_i と $p_i \neq 0$ かつ $q_i \neq 0$ のとき

$$\text{then } -\sum_{i=1}^m p_i \log_2 q_i \geq -\sum_{i=1}^m p_i \log_2 p_i = H_1(S) \text{ 等号は } q_i = p_i \ (i=1, 2, \dots, m) \text{ のとき}$$

4.3 ハフマン符号

2元ハフマン符号構成法: (i) 各情報源記号対葉作成 (ii) 最小の2つを一つ父結点連結、全体の確率相加、作為一新行として処理 (iii) 重複 (iv) 直至只剩一個子

例: A: 0.6 B: 0.25 C: 0.1 D: 0.05 ハフマン符号は一つには決まらない。

(i) 各節点から出る二つの枝のどちらに0を割り当て、どちらに1を割り当てるかは全く任意だからである 最高次の葉: 根から最も遠い

(ii) 確率の等しい葉が現れるような場合 位置にある葉

平均符号長を最小とする符号をコンパクト符号 (compact code) という。ハフマン符号はコンパクト符号である

一般の q 元ハフマン符号: 確率の最小な q 枚の葉をまとめて符号の木を作っていくという過程で符号を構成できる。(情報源の数が $(q-1)m+1$ (m : 正整数) という形でないときは、このような形になるまで、

確率0の情報源記号を付け加えてから符号を構成する必要がある)

4.4 情報源符号化定理

4.4.1 ブロック符号化

一定個数の情報源記号ごとにとめて符号化する方法。それによって構成される符号をブロック符号 (block code) と呼ぶ。

一般に、 M 元情報源 S に対し、それが発生する n 個の情報源記号をまとめて一つの情報源記号と見れば、それを発生する M^n 元情報源を S の n 次拡大情報源と見、 S^n で表す。

4.4.2 情報源符号化定理

$$H_1(S^n) \leq L_n < H_1(S^n) + 1 \quad H_1(S^n) = -\sum_{x_0, \dots, x_{n-1}} P(x_0, \dots, x_{n-1}) \log_2 P(x_0, \dots, x_{n-1}) \quad H_1(S) = \frac{H_1(S^n)}{n}$$

$$H_n(S) \leq L < H_n(S) + \frac{1}{n} \quad H(S) = \lim_{n \rightarrow \infty} H_n(S) \quad H(S) \leq H_1(S)$$

情報源符号化定理: 任意の正数 ϵ に対して、1 情報源記号当たりの平均符号長 L が $H(S) < L < H(S) + \epsilon$ とするような q 元瞬時符号に符号化できる。

$$q \text{ 元: } \frac{H(S)}{\log_2 q} \leq L < \frac{H(S)}{\log_2 q} + \epsilon$$

4.5 基本的な情報源のエントロピー

4.5.1 記憶のない情報源のエントロピー (平均符号長の限界)

$$H_1(S^n) = n H_1(S) \Rightarrow H(S) = H_1(S) = -\sum_{i=1}^M p_i \log_2 p_i$$

エントロピー関数 (entropy function): $\mathcal{H}(x) = -x \log_2 x - (1-x) \log_2 (1-x)$

4.5.2 マルコフ情報源のエントロピー

$$H(S) = \sum_{i=1}^M \omega_i \left[-\sum_{j=1}^M P(a_j | S_i) \log_2 P(a_j | S_i) \right]$$

总结:

1. 情報源符号化に必要な条件

瞬時符号 符号の木 クラフトの不等式

2. 平均符号長の限界 - S の 1 次エントロピー

$$H_1(S) = -\sum_{i=1}^M P(a_i) \log_2 P(a_i)$$

符号化方法: I. ハフマン符号 (これはコンパクト符号)

注意: $(q-1)m+1$ という形でないときは、

このような形になるまで、確率0の情報記号を付け加える

4.6 基本的な情報源符号化法

4.6.1 ハフマンブロック符号化法

ハフマン符号化を行う。1 情報源記号当たりの平均符号長は L である。その下限に近づけること。理論的には限界に L まで近い符号化が行えることはあるが、装置化の面から見ると、 L に近づけることはまだ大変である。

4.6.2 非等長情報源系列の符号化

符号化すべき情報源系列を非等長にしておく。

方法: \rightarrow 以下の

4.6.3 ランレングス符号化法

1, 01, 001, 000 は0のランの長さがそれぞれ、

0, 1, 2, 3 である場合に現れるからである。

例: 011000010001 \Rightarrow 103230

同じ記号が連続する長さ (run length) を符号化する。

例: 情報源から発生する B, AB, AAB, AAA, B を用いて、情報源記号列を区切り、2 符号化法。 ABB AAAAAB BAAAB \Rightarrow AB-B-AAAA-AAB-AAAA と符号化する

N 個の情報源系列に対し、平均長:

$$\bar{n} = \frac{1 - P^n}{P}$$

$$n = \log_2 N$$

4.7 算術符号

各節点を符号語に変換していくというものであり、

情報源系列全体を一つの符号語に符号化してしまふのである。

4.7.1 情報源系列の累積確率

$i \ a_i \ P(a_i) \ C(a_i) \ C(a_i)_2$

0 000 0.343 0 $C(a_i)$ は a_i までの累積確率と叫び

1 001 0.147 0.343 $C(a_i) = \begin{cases} 0 & (i=0) \\ \sum_{j=0}^{i-1} P(a_j) & (i=1, 2, \dots, 2^n-1) \end{cases}$

2 010 0.147 0.49 長さが0の系列は空系列と呼ばれ、 $C(0) = C(0) + P(0)$

3 011 0.063 0.637 $C(1) = C(1) + P(1)$

4 100 0.147 0.7 $C(2) = C(2) + P(2)$

5 101 0.063 0.847 $C(3) = C(3) + P(3)$

6 110 0.063 0.91 $C(4) = C(4) + P(4)$

7 111 0.027 0.973 $C(5) = C(5) + P(5)$

4.7.2 基本的算術符号化法

a_i と累積確率 $C(a_i)$ は

1 対 1 に対応する。 $a_i \in C(a_i)$ の 2 進数表を符号化する。

$C(a_i)$ の 2 進数表: $C(a_i)_2$

$L_n = \sum_{i=1}^n P(a_i) l_i \quad \lim_{n \rightarrow \infty} L_n = H(S)$

$n \rightarrow \infty$ とするとき、無限精度の乗算を要することになる。

4.7.3 乗算の不要な算術符号化

系列 x の確率 $P(x)$ の近似 $A(x)$

累積確率 $C(x)$ の近似 $\tilde{C}(x)$

I. $A(x) = 1, \tilde{C}(x) = 0$

II. $A(x) = A(x) \cdot 2^{-a}$

$A(x_0) = \langle A(x) - A(x_1) \rangle_m$

$\tilde{C}(x) = \begin{cases} \tilde{C}(x) & x=0 \\ \tilde{C}(x) + A(x_0) & x=1 \end{cases}$

算術符号の復号法

(i) $x_0 = 1$ とおく

(ii) $k=0, 1, \dots, n-1$ について

$\tilde{C}(x) < \tilde{C}(x_k) + A(x_{k+1})$

であれば $x_{k+1} = x_k 0$ とし、

そうでなければ $x_{k+1} = x_k 1$ とする

操作を繰り返す

(iii) $x = x_n$ とおく

II. ブロック符号化

M 元情報源 S , n 個

情報源 \rightarrow 1 つの記号

記号

M^n 元情報源を

S の n 次拡大情報源

と見、 S^n で表す。