

情報源符号化の定式化・続

「情報源 alphabet と頻度との組」

が符号化の対象

S : 情報源 alphabet (有限集合)

$P : S \rightarrow [0, 1] \subset \mathbf{R}$: 生起確率 $\left(\sum_{s \in S} P(s) = 1 \right)$

情報源 $\mathcal{S} := (S, P)$

: 文字 $s \in S$ を確率 $P(s)$ で次々と発生

$\rightarrow w \in S^+$ を発生

(ここでの) 仮定 : 情報源の無記憶性

各 $s \in S$ の生起確率は、 s のみで決まり、
先立って発生した文字に依らない。

情報源符号化の定式化・続

T : 符号 alphabet(有限集合)
(しばしば $T = \{0, 1\}$)

$C : S \longrightarrow T^+ : \text{符号}$
 \longrightarrow 文字列を並べて $C^* : S^* \longrightarrow T^*$ に延長

$L(C) := \sum_{s \in S} P(s) |C(s)| : C$ の平均符号長
(1 文字の符号語長の期待値)

符号への要請

- 一意符号: $C^* : S^* \longrightarrow T^* : \text{単射}$
 - 瞬時符号: $C(x) = C(s)w \implies x = sy$
(最初に届いた符号語 $C(s)$ で
最初の文字 s が復元できる)
- (以上は生起確率 P には依らない)
- 効率が良い... 平均符号長 $L(C)$ が小さい
(これは生起確率 P に依る)

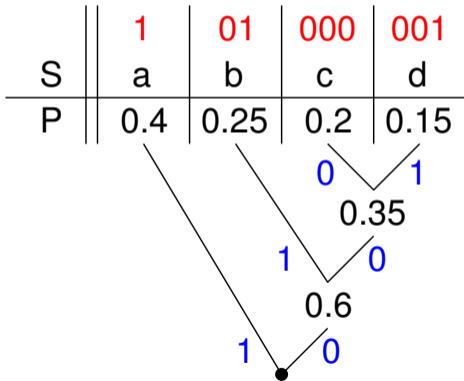
生起確率を考慮に入れて、

平均符号長の小さい符号の構成を考えよう

→ **Huffman 符号**

平均符号長の小さい符号の構成 (Huffman 符号)

例: $\#S = 4 = 2^2$, 平均符号長: $1.95 (< 2)$



Huffman 符号

- 瞬時符号の中での
平均符号長 $L(C)$ の最小値を実現
… 最適符号 (optimal code)
- 各文字の生起確率 $P(s)$ の
“ばらつき” が大きいほど効果的
- 弱点: 予め生起確率が判らないと
符号を構成できない

$S = 2$ だったら、

どうやっても

(生起確率に関わらず) $L(C) = 1$ か？

(何かうまい手はないか)

→ “**拡大情報源**” を考える

拡大情報源

情報源 $\mathcal{S} = (S, P)$ に対し、

“ n 文字づつまとめた情報源” \mathcal{S}^n を考える

$$\mathcal{S}^n := (S^n, P^{\otimes n})$$

$$S^n = S \times \cdots \times S = \{s_{i_1} \cdots s_{i_n} \mid s_{i_j} \in S\}$$

$$P^{\otimes n} : S^n \longrightarrow [0, 1] \subset \mathbf{R}$$

$$s_{i_1} \cdots s_{i_n} \longmapsto P(s_{i_1}) \cdots P(s_{i_n})$$

→ この情報源 \mathcal{S}^n を符号化せよ

問題: 次の生起確率を持つ情報源

$S = (S, P), S = \{a, b\}$ について、

S	a	b
P	0.8	0.2

- (1) 2 次の拡大情報源 $S^2 = (S^2, P^{\otimes 2})$
に対する **Huffman** 符号 C_2 を構成し、
“1 文字当たりの平均符号長”
 $L(C_2)/2$ を求めよ。
- (2) 3 次の拡大情報源 $S^3 = (S^3, P^{\otimes 3})$
に対しても同様の計算をせよ。

一般に、

n 次の拡大情報源 $S^n = (S^n, P^{\otimes n})$
の符号 C_n で、

n を大きくしてゆくと、

情報源 alphabet 1 文字当たりの平均符号長

$$\frac{L(C_n)}{n}$$

を下げる事が出来る

(符号は複雑になってゆくが)

n 次の拡大情報源 $S^n = (S^n, P^{\otimes n})$
の符号 C_n を利用すると、

情報源 alphabet 1 文字当たりの平均符号長

$$\frac{L(C_n)}{n}$$

はどこまで下げられるか？

→ 情報源 S が本来持っている “情報の量”
より小さくはならないだろう。

“エントロピー (entropy)”

“情報の量”

「或る事象 P が起こる」

という“情報の価値”は

どう評価したら良いか？

“情報の量”

基本的なアイデア:

確率 $\frac{1}{4}$ で起きる出来事を
教えてもらうことの価値は、

確率 $\frac{1}{2}$ で起きる出来事を 2 つ
教えてもらうのと同じ

→ “情報の量” が 2 倍

「事象 P が起こる」という“情報の価値” $I(P)$

要請:

(1) 生起確率 p のみに依る $\longrightarrow I(p) := I(P)$

(2) 独立な事象 P_1, P_2 に対して、

$$I(P_1 \wedge P_2) = I(P_1) + I(P_2)$$

$$\longrightarrow I(p_1 p_2) = I(p_1) + I(p_2)$$

(3) $I : (0, 1] \longrightarrow \mathbf{R}_{\geq 0}$: 連続関数

(const. 0 ではない)

$$\longrightarrow I(p) = C \log \frac{1}{p} = -C \log p \quad (C > 0)$$

「事象 P が起こる」という“情報の価値” $I(P)$

$$I(p) = C \log \frac{1}{p} = -C \log p \quad (C > 0)$$

定数 C の取り方

←→ \log の底の取り方

←→ “情報量” の単位の取り方

通常 2 を底に取って ($I(\frac{1}{2}) := 1$)、

“情報量” の単位とする: **bit (binary digit)**

情報源のエントロピー

情報源 $\mathcal{S} = (S, P)$ の

1 文字から得られる情報量の期待値

$$H(\mathcal{S}) := \sum_{s \in S} P(s) I(P(s))$$

: 情報源 \mathcal{S} の**エントロピー (entropy)**

$S = \{s_1, \dots, s_k\}$, $P(s_i) = p_i$ の時は、

$$H(\mathcal{S}) := \sum_{i=1}^k p_i \log \frac{1}{p_i} = - \sum_{i=1}^k p_i \log p_i$$

情報源のエントロピー

特に $\#S = 2$ ($S = \{ \text{起きる}, \text{起きない} \}$) の時、

- 起きる確率 p
- 起きない確率 $1 - p =: \bar{p}$

$$H(p) := H(S) = p \log \frac{1}{p} + \bar{p} \log \frac{1}{\bar{p}}$$

: 2 値エントロピー関数

問: $y = I(p), y = H(p)$ のグラフの概形を描け。
(最大値をとる p とその時の値は ?)

n 次の拡大情報源 $S^n = (S^n, P^{\otimes n})$
の符号 C_n を利用すると、

情報源 alphabet 1 文字当たりの平均符号長

$$\frac{L(C_n)}{n}$$

はどこまで下げられるか？

→ 情報源 S のエントロピー $H(S)$ より
小さくはないだろう

→ まず一般に、符号 C に対し、
平均符号長 $L(C)$ と $H(S)$ とを比べよう

定理

$\mathcal{S} = (S, P)$: 情報源

\mathcal{C} : \mathcal{S} の一意符号 (瞬時符号)

$$\implies \boxed{L(\mathcal{C}) \geq H(\mathcal{S})}$$

(但し、 \log の底は $r := \#T$ に取る)

$$\eta := \frac{H(\mathcal{S})}{L(\mathcal{C})} : \mathcal{C} \text{ の効率 (efficiency)}$$

$$\bar{\eta} = 1 - \eta : \mathcal{C} \text{ の冗長度 (redundancy)}$$

Kraft の不等式(再掲)

$$S = \{s_1, \dots, s_k\}, \quad \#T = r$$

自然数列 (ℓ_1, \dots, ℓ_k) に対し、

各符号語長 $|C(s_i)| = \ell_i$ なる

r 元瞬時符号が存在

$$\iff \sum_{i=1}^k \frac{1}{r^{\ell_i}} \leq 1$$

補題

$$x_i, y_i > 0 \quad (i = 1, \dots, k)$$

$$\sum_{i=1}^k x_i = \sum_{i=1}^k y_i = 1$$

$$\implies \sum_{i=1}^k x_i \log \frac{1}{x_i} \leq \sum_{i=1}^k x_i \log \frac{1}{y_i}$$

(等号は $\forall i : x_i = y_i$ のとき)

- $L(C) = H(S)$ は実現できるのか？
- $\inf_C L(C) = H(S)$ であるか？

Huffman 符号は確かに最適だが、
 $L(C)$ の上からの評価は難しい

→ **Shannon-Fano** 符号
(**Kraft-McMillan** の不等式の利用)