

# 平均符号長の限界と 情報源符号化定理

# 情報源の1次エントロピー

## ■ 情報源 $S$ :

- 情報源記号 $A = \{a_1, \dots, a_M\}$
- 各情報源記号の発生確率  $P(a_i) = p_i$

## ■ 情報源符号化

- 情報源記号 $a_i$ の符号化 $K(a_i) = c_i$
- 1情報源記号あたりの平均符号長 $L = p_1 l_1 + \dots + p_M l_M$   
ここで,  $l_i = |c_i|$

## ■ 情報源 $S$ の1次エントロピー

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

# 情報源の1次エントロピー

## 【補助定理3】

$p_1, \dots, p_M$ を  $p_1 + \dots + p_M = 1$ となる非負の数,  $q_1, \dots, q_M$ を  $q_1 + \dots + q_M \leq 1$ となる非負の数とする.  $p_i \neq 0$ のとき,  $q_i \neq 0$ とすれば,

$$-\sum_{i=1}^M p_i \log_2 q_i \geq -\sum_{i=1}^M p_i \log_2 p_i$$

が成立する. 等号はすべての*i*について $q_i = p_i$ のときに限り成立する.

## 【補助定理3】の適用例

$(p_1, p_2, p_3, p_4, p_5)$

(0.6, 0.2, 0.1, 0.07, 0.03)

$$-\sum_{i=1}^M p_i \log_2 q_i$$

$(q_1, q_2, q_3, q_4, q_5)$

(0.2, 0.2, 0.2, 0.2, 0.2)

$\approx 2.32193$

(0.9, 0.025, 0.025, 0.025, 0.025)

$\approx 2.21997$

(0.6, 0.2, 0.1, 0.07, 0.03)

$\approx 1.65908$

# 情報源の1次エントロピー

## 【補助定理3】

$p_1, \dots, p_M$  を  $p_1 + \dots + p_M = 1$  となる非負の数,  $q_1, \dots, q_M$   $q_1 + \dots + q_M \leq 1$  となる非負の数とする. .  $p_i \neq 0$  のとき,  $q_i \neq 0$  とすれば,

$$-\sum_{i=1}^M p_i \log_2 q_i \geq -\sum_{i=1}^M p_i \log_2 p_i$$

が成立する. 等号はすべての  $i$  について  $q_i = p_i$  のときに限り成立する.

## 【補助定理3】の証明の骨子

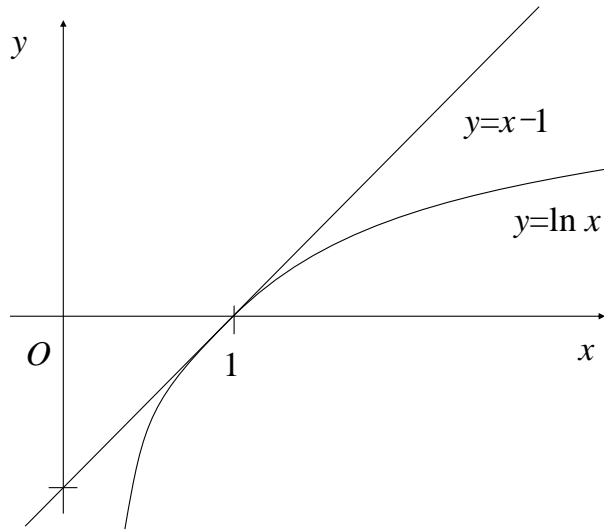
$$D = -\sum_{i=1}^M p_i \log_2 q_i + \sum_{i=1}^M p_i \log_2 p_i = -\sum_{i=1}^M p_i \log_2 \frac{q_i}{p_i} = -\sum_{i=1}^M \frac{p_i}{\ln 2} \ln \frac{q_i}{p_i}$$

と置く.  $\ln x \leq x - 1$  であることを利用すると,

$$D = -\sum_{i=1}^M \frac{p_i}{\ln 2} \ln \frac{q_i}{p_i} \geq -\sum_{i=1}^M \frac{p_i}{\ln 2} \left( \frac{q_i}{p_i} - 1 \right) = \frac{1}{\ln 2} \sum_{i=1}^M (p_i - q_i) = \frac{1}{\ln 2} \left( \sum_{i=1}^M p_i - \sum_{i=1}^M q_i \right) \geq 0$$

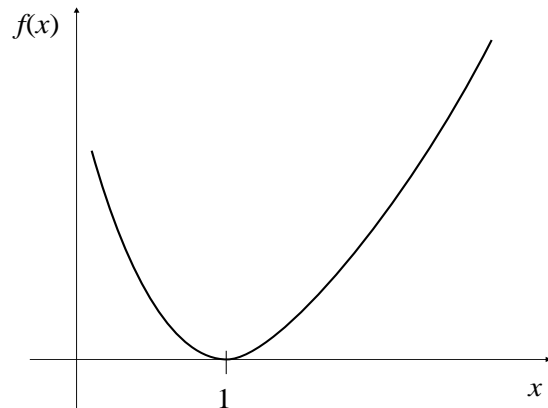


# 情報源の1次エントロピー



$$\ln x \leq x - 1$$

$$f(x) = x - 1 - \ln x$$



$$f(1) = 0, f'(x) = 1 - \frac{1}{x}$$

# 情報源の1次エントロピー

## 【定理】

情報源 $S$ の各情報源記号を一意復号可能な2元符号に符号化すると、平均符号長 $L$ は、 $H_1(S) \leq L$ を満足する.

また、平均符号長 $L$ が $L < H_1(S) + 1$ となる瞬時符号を作ることができる.

# 情報源の1次エントロピー

【定理】の証明の骨子 （前半:  $H_1(S) \leq L$  の証明）

- $q_i = 2^{-l_i}$  と置く ( $q_i$  を導入する).
- 一意復号可能な情報源符号は,  $l_1, l_2, \dots, l_M$  はマクミランの不等式を満足するので,  $2^{-l_1} + 2^{-l_2} + \dots + 2^{-l_M} \leq 1$ . つまり,  $q_1 + \dots + q_M \leq 1$  が満足される. また,  $q_i > 0$  も明らか.
- 従って補助定理から,  $-\sum_{i=1}^M p_i \log_2 q_i \geq -\sum_{i=1}^M p_i \log_2 p_i = H_1(S)$  となる.
- 左辺は,  $-\sum_{i=1}^M p_i \log_2 q_i = -\sum_{i=1}^M p_i \log_2 2^{-l_i} = \sum_{i=1}^M p_i l_i = L$  である. 等号が成立するのは, 全ての  $i$  について,  $p_i = 2^{-l_i}$  のときである. ■

# 情報源の1次エントロピー

【定理】の証明の骨子（後半:  $L < H_1(S) + 1$ を満足する瞬時符号が構成可能）

- $-\log_2 p_i \leq l_i < -\log_2 p_i + 1$ なる整数 $l_i$ が一意に決まる.
- 得られた $l_i$ に対して $2^{-l_i} \leq 2^{\log_2 p_i} = p_i$ が成り立つ.
- $\sum 2^{-l_i} \leq \sum p_i = 1$ となり, クラフトの不等式が満たされるので,  $l_i$ を符号長バッグとする瞬時符号をつくることができる.
- $L = \sum_{i=1}^M p_i l_i$ であるので,

この符号の平均符号長は,  $-p_i \log_2 p_i \leq p_i l_i < -p_i \log_2 p_i + p_i$ つまり,

$H_1(S) \leq L < H_1(S) + 1$ が満足される. ■



# 定理の意味すること

情報源S:  $\langle A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03 \rangle$

$$H_1(S) = -(0.6 \log_2 0.6 + 0.2 \log_2 0.2 + 0.1 \log_2 0.1 + 0.07 \log_2 0.07 + 0.03 \log_2 0.03) \\ \approx 1.65908$$

平均符号長が1.65908～2.65908となる瞬時符号を構成可能.

しかしこの定理は, そのような瞬時符号の構成法を教えてくれるわけではない.

⇒ ハフマン符号化などを用いてコンパクト符号を計算する.

$\langle A \rightarrow 0, B \rightarrow 10, C \rightarrow 110, D \rightarrow 1110, E \rightarrow 1111 \rangle$ はコンパクト符号で平均符号長1.7

情報源記号ごとに符号化するかぎり, 平均符号長を理論的下限に近づけることはできない

# 拡大情報源

- $q$ 元情報源 $S$ の $n$ 次の拡大情報源 $S^n$ :  $S$ の連続する $n$ 個の情報源記号列を情報源記号とする $q^n$ 元情報源.

- 無記憶情報源 $S$ の連続する $n$ 個の出力は互いに独立であるから, その結合確率分布は,

$$P(x_0, x_1, \dots, x_{n-1}) = P(x_0)P(x_1) \cdots P(x_{n-1})$$

となる.

- 無記憶情報源 $S$ の $n$ 次拡大情報源を $S^n$ とすれば,

$$H_1(S^n) = nH_1(S)$$

が成り立つ.

# 拡大情報源

$$\begin{aligned} H_1(S^n) = nH_1(S) \text{になるわけ} &\Rightarrow H_1(S^3) \\ &= -\sum_{x_0} \sum_{x_1} \sum_{x_2} P(x_0, x_1, x_2) \log_2 P(x_0, x_1, x_2) \\ &= -\sum_{x_0} \sum_{x_1} \sum_{x_2} P(x_0)P(x_1)P(x_2) \log_2 (P(x_0)P(x_1)P(x_2)) \\ &= -\sum_{x_0} \sum_{x_1} \sum_{x_2} P(x_0)P(x_1)P(x_2) \log_2 P(x_0) \\ &\quad - \sum_{x_0} \sum_{x_1} \sum_{x_2} P(x_0)P(x_1)P(x_2) \log_2 P(x_1) \\ &\quad - \sum_{x_0} \sum_{x_1} \sum_{x_2} P(x_0)P(x_1)P(x_2) \log_2 P(x_2) \\ &= -\sum_{x_0} P(x_0) \log_2 P(x_0) \\ &\quad - \sum_{x_1} P(x_1) \log_2 P(x_1) \\ &\quad - \sum_{x_2} P(x_2) \log_2 P(x_2) \\ &= 3H_1(S) \end{aligned}$$

## 計算の様子 (1/2)

$$\sum_{x_1 \in \{A,B,C\}} \sum_{x_2 \in \{D,E,F\}} \sum_{x_3 \in \{G,H,I\}} P(x_0)P(x_1)P(x_2) \log_2(P(x_0)P(x_1)P(x_2))$$

$P(A)P(D)P(G)\log_2(P(A)P(D)P(G))$   
 $P(A)P(D)P(H)\log_2(P(A)P(D)P(H))$   
 $P(A)P(D)P(I)\log_2(P(A)P(D)P(I))$   
 $P(A)P(E)P(G)\log_2(P(A)P(E)P(G))$   
 $P(A)P(E)P(H)\log_2(P(A)P(E)P(H))$   
 $P(A)P(E)P(I)\log_2(P(A)P(E)P(I))$   
 $P(A)P(F)P(G)\log_2(P(A)P(F)P(G))$   
 $P(A)P(F)P(H)\log_2(P(A)P(F)P(H))$   
 $P(A)P(F)P(I)\log_2(P(A)P(F)P(I))$   
 $P(B)P(D)P(G)\log_2(P(B)P(D)P(G))$   
 $P(B)P(D)P(H)\log_2(P(B)P(D)P(H))$   
 $P(B)P(D)P(I)\log_2(P(B)P(D)P(I))$   
 $P(B)P(E)P(G)\log_2(P(B)P(E)P(G))$   
 $P(B)P(E)P(H)\log_2(P(B)P(E)P(H))$   
 $P(B)P(E)P(I)\log_2(P(B)P(E)P(I))$   
 $P(B)P(F)P(G)\log_2(P(B)P(F)P(G))$   
 $P(B)P(F)P(H)\log_2(P(B)P(F)P(H))$   
 $P(B)P(F)P(I)\log_2(P(B)P(F)P(I))$   
 $P(C)P(D)P(G)\log_2(P(C)P(D)P(G))$   
 $P(C)P(D)P(H)\log_2(P(C)P(D)P(H))$   
 $P(C)P(D)P(I)\log_2(P(C)P(D)P(I))$   
 $P(C)P(E)P(G)\log_2(P(C)P(E)P(G))$   
 $P(C)P(E)P(H)\log_2(P(C)P(E)P(H))$   
 $P(C)P(E)P(I)\log_2(P(C)P(E)P(I))$   
 $P(C)P(F)P(G)\log_2(P(C)P(F)P(G))$   
 $P(C)P(F)P(H)\log_2(P(C)P(F)P(H))$   
 $P(C)P(F)P(I)\log_2(P(C)P(F)P(I))$

$P(A)P(D)P(G)(\log_2 P(A) + \log_2 P(D) + \log_2 P(G))$   
 $P(A)P(D)P(H)(\log_2 P(A) + \log_2 P(D) + \log_2 P(H))$   
 $P(A)P(D)P(I)(\log_2 P(A) + \log_2 P(D) + \log_2 P(I))$   
 $P(A)P(E)P(G)(\log_2 P(A) + \log_2 P(E) + \log_2 P(G))$   
 $P(A)P(E)P(H)(\log_2 P(A) + \log_2 P(E) + \log_2 P(H))$   
 $P(A)P(E)P(I)(\log_2 P(A) + \log_2 P(E) + \log_2 P(I))$   
 $P(A)P(F)P(G)(\log_2 P(A) + \log_2 P(F) + \log_2 P(G))$   
 $P(A)P(F)P(H)(\log_2 P(A) + \log_2 P(F) + \log_2 P(H))$   
 $P(A)P(F)P(I)(\log_2 P(A) + \log_2 P(F) + \log_2 P(I))$   
 $P(B)P(D)P(G)(\log_2 P(B) + \log_2 P(D) + \log_2 P(G))$   
 $P(B)P(D)P(H)(\log_2 P(B) + \log_2 P(D) + \log_2 P(H))$   
 $P(B)P(D)P(I)(\log_2 P(B) + \log_2 P(D) + \log_2 P(I))$   
 $P(B)P(E)P(G)(\log_2 P(B) + \log_2 P(E) + \log_2 P(G))$   
 $P(B)P(E)P(H)(\log_2 P(B) + \log_2 P(E) + \log_2 P(H))$   
 $P(B)P(E)P(I)(\log_2 P(B) + \log_2 P(E) + \log_2 P(I))$   
 $P(B)P(F)P(G)(\log_2 P(B) + \log_2 P(F) + \log_2 P(G))$   
 $P(B)P(F)P(H)(\log_2 P(B) + \log_2 P(F) + \log_2 P(H))$   
 $P(B)P(F)P(I)(\log_2 P(B) + \log_2 P(F) + \log_2 P(I))$   
 $P(C)P(D)P(G)(\log_2 P(C) + \log_2 P(D) + \log_2 P(G))$   
 $P(C)P(D)P(H)(\log_2 P(C) + \log_2 P(D) + \log_2 P(H))$   
 $P(C)P(D)P(I)(\log_2 P(C) + \log_2 P(D) + \log_2 P(I))$   
 $P(C)P(E)P(G)(\log_2 P(C) + \log_2 P(E) + \log_2 P(G))$   
 $P(C)P(E)P(H)(\log_2 P(C) + \log_2 P(E) + \log_2 P(H))$   
 $P(C)P(E)P(I)(\log_2 P(C) + \log_2 P(E) + \log_2 P(I))$   
 $P(C)P(F)P(G)(\log_2 P(C) + \log_2 P(F) + \log_2 P(G))$   
 $P(C)P(F)P(H)(\log_2 P(C) + \log_2 P(F) + \log_2 P(H))$   
 $P(C)P(F)P(I)(\log_2 P(C) + \log_2 P(F) + \log_2 P(I))$

$P(A)P(D)P(G)\log_2 P(A)$   
 $P(A)P(D)P(H)\log_2 P(A)$   
 $P(A)P(D)P(I)\log_2 P(A)$   
 $P(A)P(E)P(G)\log_2 P(A)$   
 $P(A)P(E)P(H)\log_2 P(A)$   
 $P(A)P(E)P(I)\log_2 P(A)$   
 $P(A)P(F)P(G)\log_2 P(A)$   
 $P(A)P(F)P(H)\log_2 P(A)$   
 $P(A)P(F)P(I)\log_2 P(A)$   
 $P(B)P(D)P(G)\log_2 P(B)$   
 $P(B)P(D)P(H)\log_2 P(B)$   
 $P(B)P(D)P(I)\log_2 P(B)$   
 $P(B)P(E)P(G)\log_2 P(B)$   
 $P(B)P(E)P(H)\log_2 P(B)$   
 $P(B)P(E)P(I)\log_2 P(B)$   
 $P(B)P(F)P(G)\log_2 P(B)$   
 $P(B)P(F)P(H)\log_2 P(B)$   
 $P(B)P(F)P(I)\log_2 P(B)$   
 $P(C)P(D)P(G)\log_2 P(C)$   
 $P(C)P(D)P(H)\log_2 P(C)$   
 $P(C)P(D)P(I)\log_2 P(C)$   
 $P(C)P(E)P(G)\log_2 P(C)$   
 $P(C)P(E)P(H)\log_2 P(C)$   
 $P(C)P(E)P(I)\log_2 P(C)$   
 $P(C)P(F)P(G)\log_2 P(C)$   
 $P(C)P(F)P(H)\log_2 P(C)$   
 $P(C)P(F)P(I)\log_2 P(C)$

$P(A)P(D)P(G)\log_2 P(D)$   
 $P(A)P(D)P(H)\log_2 P(D)$   
 $P(A)P(D)P(I)\log_2 P(D)$   
 $P(A)P(E)P(G)\log_2 P(E)$   
 $P(A)P(E)P(H)\log_2 P(E)$   
 $P(A)P(E)P(I)\log_2 P(E)$   
 $P(A)P(F)P(G)\log_2 P(F)$   
 $P(A)P(F)P(H)\log_2 P(F)$   
 $P(A)P(F)P(I)\log_2 P(F)$   
 $P(B)P(D)P(G)\log_2 P(D)$   
 $P(B)P(D)P(H)\log_2 P(D)$   
 $P(B)P(D)P(I)\log_2 P(D)$   
 $P(B)P(E)P(G)\log_2 P(E)$   
 $P(B)P(E)P(H)\log_2 P(E)$   
 $P(B)P(E)P(I)\log_2 P(E)$   
 $P(B)P(F)P(G)\log_2 P(F)$   
 $P(B)P(F)P(H)\log_2 P(F)$   
 $P(B)P(F)P(I)\log_2 P(F)$   
 $P(C)P(D)P(G)\log_2 P(D)$   
 $P(C)P(D)P(H)\log_2 P(D)$   
 $P(C)P(D)P(I)\log_2 P(D)$   
 $P(C)P(E)P(G)\log_2 P(E)$   
 $P(B)P(E)P(H)\log_2 P(E)$   
 $P(B)P(E)P(I)\log_2 P(E)$   
 $P(B)P(F)P(G)\log_2 P(F)$   
 $P(B)P(F)P(H)\log_2 P(F)$   
 $P(B)P(F)P(I)\log_2 P(F)$

$P(A)P(D)P(G)\log_2 P(G)$   
 $P(A)P(D)P(H)\log_2 P(H)$   
 $P(A)P(D)P(I)\log_2 P(I)$   
 $P(A)P(E)P(G)\log_2 P(G)$   
 $P(A)P(E)P(H)\log_2 P(H)$   
 $P(A)P(E)P(I)\log_2 P(I)$   
 $P(A)P(F)P(G)\log_2 P(G)$   
 $P(A)P(F)P(H)\log_2 P(H)$   
 $P(A)P(F)P(I)\log_2 P(I)$   
 $P(B)P(D)P(G)\log_2 P(G)$   
 $P(B)P(D)P(H)\log_2 P(H)$   
 $P(B)P(D)P(I)\log_2 P(I)$   
 $P(B)P(E)P(G)\log_2 P(G)$   
 $P(B)P(E)P(H)\log_2 P(H)$   
 $P(B)P(E)P(I)\log_2 P(I)$   
 $P(B)P(F)P(G)\log_2 P(G)$   
 $P(B)P(F)P(H)\log_2 P(H)$   
 $P(B)P(F)P(I)\log_2 P(I)$   
 $P(C)P(D)P(G)\log_2 P(G)$   
 $P(C)P(D)P(H)\log_2 P(H)$   
 $P(C)P(D)P(I)\log_2 P(I)$   
 $P(C)P(E)P(G)\log_2 P(G)$   
 $P(C)P(E)P(H)\log_2 P(H)$   
 $P(C)P(E)P(I)\log_2 P(I)$   
 $P(C)P(F)P(G)\log_2 P(G)$   
 $P(C)P(F)P(H)\log_2 P(H)$   
 $P(C)P(F)P(I)\log_2 P(I)$

# 計算の様子 (1/2)

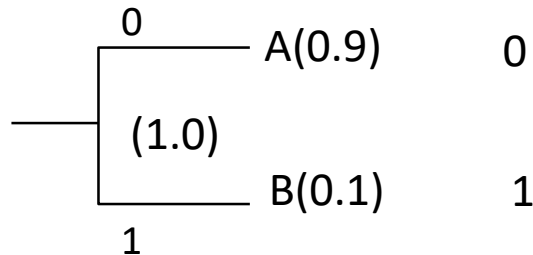
$$\begin{aligned}
 & \sum \begin{matrix} P(A)P(D)P(G)\log_2 P(A) \\ P(A)P(D)P(H)\log_2 P(A) \\ P(A)P(D)P(I)\log_2 P(A) \\ P(A)P(E)P(G)\log_2 P(A) \\ P(A)P(E)P(H)\log_2 P(A) \\ P(A)P(E)P(I)\log_2 P(A) \\ P(A)P(F)P(G)\log_2 P(A) \\ P(A)P(F)P(H)\log_2 P(A) \\ P(A)P(F)P(I)\log_2 P(A) \\ P(B)P(D)P(G)\log_2 P(B) \\ P(B)P(D)P(H)\log_2 P(B) \\ P(B)P(D)P(I)\log_2 P(B) \\ P(B)P(E)P(G)\log_2 P(B) \\ P(B)P(E)P(H)\log_2 P(B) \\ P(B)P(E)P(I)\log_2 P(B) \\ P(B)P(F)P(G)\log_2 P(B) \\ P(B)P(F)P(H)\log_2 P(B) \\ P(B)P(F)P(I)\log_2 P(B) \\ P(C)P(D)P(G)\log_2 P(C) \\ P(C)P(D)P(H)\log_2 P(C) \\ P(C)P(D)P(I)\log_2 P(C) \\ P(C)P(E)P(G)\log_2 P(C) \\ P(C)P(E)P(H)\log_2 P(C) \\ P(C)P(E)P(I)\log_2 P(C) \\ P(C)P(F)P(G)\log_2 P(C) \\ P(C)P(F)P(H)\log_2 P(C) \\ P(C)P(F)P(I)\log_2 P(C) \end{matrix} + \sum \begin{matrix} P(A)P(D)P(G)\log_2 P(D) \\ P(A)P(D)P(H)\log_2 P(D) \\ P(A)P(D)P(I)\log_2 P(D) \\ P(A)P(E)P(G)\log_2 P(E) \\ P(A)P(E)P(H)\log_2 P(E) \\ P(A)P(E)P(I)\log_2 P(E) \\ P(A)P(F)P(G)\log_2 P(F) \\ P(A)P(F)P(H)\log_2 P(F) \\ P(A)P(F)P(I)\log_2 P(F) \\ P(B)P(D)P(G)\log_2 P(D) \\ P(B)P(D)P(H)\log_2 P(D) \\ P(B)P(D)P(I)\log_2 P(D) \\ P(B)P(E)P(G)\log_2 P(E) \\ P(B)P(E)P(H)\log_2 P(E) \\ P(B)P(E)P(I)\log_2 P(E) \\ P(B)P(F)P(G)\log_2 P(F) \\ P(B)P(F)P(H)\log_2 P(F) \\ P(B)P(F)P(I)\log_2 P(F) \\ P(C)P(D)P(G)\log_2 P(D) \\ P(C)P(D)P(H)\log_2 P(D) \\ P(C)P(D)P(I)\log_2 P(D) \\ P(C)P(E)P(G)\log_2 P(E) \\ P(C)P(E)P(H)\log_2 P(E) \\ P(C)P(E)P(I)\log_2 P(E) \\ P(B)P(E)P(H)\log_2 P(E) \\ P(B)P(F)P(G)\log_2 P(F) \\ P(B)P(F)P(H)\log_2 P(F) \\ P(B)P(F)P(I)\log_2 P(F) \end{matrix} + \sum \begin{matrix} P(A)P(D)P(G)\log_2 P(G) \\ P(A)P(D)P(H)\log_2 P(H) \\ P(A)P(D)P(I)\log_2 P(I) \\ P(A)P(E)P(G)\log_2 P(G) \\ P(A)P(E)P(H)\log_2 P(H) \\ P(A)P(E)P(I)\log_2 P(I) \\ P(A)P(F)P(G)\log_2 P(G) \\ P(A)P(F)P(H)\log_2 P(H) \\ P(A)P(F)P(I)\log_2 P(I) \\ P(B)P(D)P(G)\log_2 P(G) \\ P(B)P(D)P(H)\log_2 P(H) \\ P(B)P(D)P(I)\log_2 P(I) \\ P(B)P(E)P(G)\log_2 P(G) \\ P(B)P(E)P(H)\log_2 P(H) \\ P(B)P(E)P(I)\log_2 P(I) \\ P(B)P(F)P(G)\log_2 P(G) \\ P(B)P(F)P(H)\log_2 P(H) \\ P(B)P(F)P(I)\log_2 P(I) \\ P(C)P(D)P(G)\log_2 P(G) \\ P(C)P(D)P(H)\log_2 P(H) \\ P(C)P(D)P(I)\log_2 P(I) \\ P(C)P(E)P(G)\log_2 P(G) \\ P(C)P(E)P(H)\log_2 P(H) \\ P(C)P(E)P(I)\log_2 P(I) \\ P(C)P(F)P(G)\log_2 P(G) \\ P(C)P(F)P(H)\log_2 P(H) \\ P(C)P(F)P(I)\log_2 P(I) \end{matrix} = \sum \begin{pmatrix} P(A)\log_2 P(A) \\ P(B)\log_2 P(B) \\ P(C)\log_2 P(C) \end{pmatrix} + \sum \begin{pmatrix} P(D)\log_2 P(D) \\ P(E)\log_2 P(E) \\ P(F)\log_2 P(F) \end{pmatrix} + \sum \begin{pmatrix} P(G)\log_2 P(G) \\ P(H)\log_2 P(H) \\ P(I)\log_2 P(I) \end{pmatrix} \\
 & \qquad \qquad \qquad \sum \begin{pmatrix} P(A)\log_2 P(A) \\ P(B)\log_2 P(B) \\ P(C)\log_2 P(C) \end{pmatrix} \times \sum \begin{pmatrix} P(D)\log_2 P(D) \\ P(E)\log_2 P(E) \\ P(F)\log_2 P(F) \end{pmatrix} + \sum \begin{pmatrix} P(D)\log_2 P(D) \\ P(E)\log_2 P(E) \\ P(F)\log_2 P(F) \end{pmatrix} \times \sum \begin{pmatrix} P(G)\log_2 P(G) \\ P(H)\log_2 P(H) \\ P(I)\log_2 P(I) \end{pmatrix} + \sum \begin{pmatrix} P(G)\log_2 P(G) \\ P(H)\log_2 P(H) \\ P(I)\log_2 P(I) \end{pmatrix} \times \sum \begin{pmatrix} P(A)\log_2 P(A) \\ P(B)\log_2 P(B) \\ P(C)\log_2 P(C) \end{pmatrix} \\
 & \qquad \qquad \qquad \sum \begin{pmatrix} P(D)P(G) \\ P(D)P(H) \\ P(D)P(I) \\ P(E)P(G) \\ P(E)P(H) \\ P(E)P(I) \\ P(F)P(G) \\ P(F)P(H) \\ P(F)P(I) \end{pmatrix} + \sum \begin{pmatrix} P(A)P(G) \\ P(A)P(H) \\ P(A)P(I) \\ P(B)P(G) \\ P(B)P(H) \\ P(B)P(I) \\ P(C)P(G) \\ P(C)P(H) \\ P(C)P(I) \end{pmatrix} + \sum \begin{pmatrix} P(A)P(D) \\ P(A)P(E) \\ P(A)P(F) \\ P(B)P(D) \\ P(B)P(E) \\ P(B)P(F) \\ P(C)P(D) \\ P(C)P(E) \\ P(C)P(F) \end{pmatrix} \\
 & \qquad \qquad \qquad \sum \begin{pmatrix} P(D)P(G) \\ P(D)P(H) \\ P(D)P(I) \\ P(E)P(G) \\ P(E)P(H) \\ P(E)P(I) \\ P(F)P(G) \\ P(F)P(H) \\ P(F)P(I) \end{pmatrix} = (P(D) + P(E) + P(F)) \times (P(G) + P(H) + P(I)) = 1
 \end{aligned}$$

# ブロック符号化

- ブロック符号化：情報源から発生する記号をまとめて符号化する方法.
- ブロック符号化をすることによって，平均符号長を情報源エントロピーに近づけることができる.

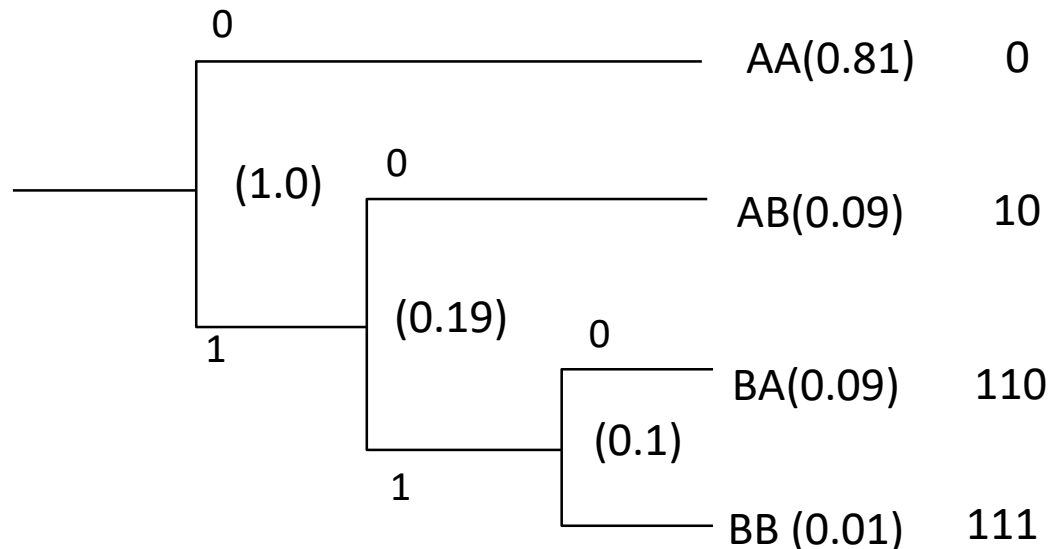
# ブロック符号化

〈A: 0.9, B: 0.1〉に対するハフマン符号化



⇒ 情報源記号あたりの平均符号長1

〈AA: 0.81, AB: 0.09, BA: 0.09, BB: 0.01〉に対するハフマン符号化



⇒ 情報源記号あたりの平均符号長0.645

# ブロック符号化

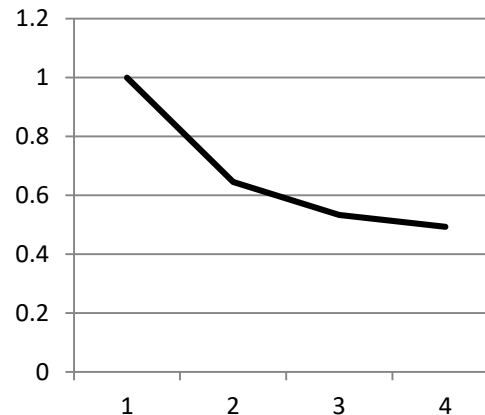
〈A: 0.9, B: 0.1〉の3次拡大情報源に対するハフマン符号化

AAA: 0.729 0  
AAB: 0.081 100  
ABA: 0.081 101  
BAA: 0.081 110  
ABB: 0.009 11100  
BAB: 0.009 11101  
BBA: 0.009 11110  
BBB: 0.001 11111

⇒ 情報源記号あたりの平均符号長約0.53267

## ブロック符号化の威力

情報源記号あたり平均符号長



ブロック長を大きくしたらどうなるか？

ブロック長



# 情報源符号化定理

情報源 $S$ について,

$$H_1(S^n) \leq L_n < H_1(S^n) + 1$$

が成立する.

情報源  $S$  の  $n$  次の拡張  $S^n$  についても,

$$H_1(S^n) \leq L_n < H_1(S^n) + 1$$

ここで,  $L_n$  は  $S$  の 1 記号あたりの平均符号長.  $S$  が無記憶であれば

$$H_1(S^n) = nH_1(S)$$

また,  $L_n = nL$  なので,

$$H_1(S) \leq L < H_1(S) + \frac{1}{n}$$

→ブロック符号化によって, 平均符号長を $H_1(S)$ に限りなく近づけられる

# 情報源符号化定理

## 【情報源符号化定理】

情報源 $S$ は、任意の正数 $\varepsilon$ に対して、1情報源記号あたりの平均符号長 $L$ が、

$$H(S) \leq L < H(S) + \varepsilon$$

となるような2元瞬時符号に符号化できる.

他方、どのような一意復号可能な2元符号を用いても、平均符号長がこの式の左辺より小さくなる符号化はできない.

# 情報源符号化定理

情報源 $S$ のエントロピー $H(S)$  :

$$H(S) \equiv \lim_{n \rightarrow \infty} H_n(S) = \lim_{n \rightarrow \infty} \frac{H_1(S^n)}{n}$$

ここで,  $H_1(S^n) = -\sum \cdots \sum P(x_0, \cdots, x_{n-1}) \log_2 P(x_0, \cdots, x_{n-1})$   
は $S$ の $n$ 次の拡大情報源  $S^n$  の1次エントロピー.

$H_n(S) \equiv \frac{H_1(S^n)}{n}$  は $S$ の1情報源記号あたりの $n$ 次エントロピー.

一般には,  $H(S) \leq H_n(S)$ であることが知られている.

# まとめ

- 情報源の1次エントロピー
- 補助定理3
- 情報源記号ごとに瞬時符号を構成した場合の平均符号長の下限
- 拡大情報源とその1次エントロピー
- ブロック符号化: 情報源から発生する記号をまとめて符号化する方法.
- ブロック符号化をすることによって, 平均符号長を情報源エントロピーに近づけることができる.
- 情報源符号化定理