

§ Shannonの第1定理

(source coding theorem)

Ω 種類の「文字」 ← 情報源

・文字 i が確率 p_i で出現 (ランダム情報源)

・符号化: $i \rightarrow l_i$ 個の 0,1 の列

$$\langle \hat{L} \rangle = \sum_{i=1}^{\Omega} l_i p_i$$

定理

・任意の接頭符号について

$\langle \hat{L} \rangle \geq H(P)$ が成立

・適切な接頭符号をすれば、必ず

$\langle \hat{L} \rangle \leq H(P) + 1$ となる

意味 $H(P) \gg 1$, $H(P) + 1 \approx H(P)$

最適な符号を用いると

N 文字の情報源 $\rightarrow N \langle \hat{L} \rangle \approx H(P)$

bit の
0,1 の列

(文字列) の「情報量」が $H(P)$ bit

$H(P) \gg 1$ ではないときは?

この場合

n 文字をくっつけてひとつの「文字」とみなせる

文字の種類 Ω^n

確率分布 $P(i_1, i_2, \dots, i_n) = p_{i_1} p_{i_2} \dots p_{i_n}$

$H(P_{new}) = nH(P_{old})$

P_{new} について Shannonの定理

N 文字の情報源 $\rightarrow \frac{N}{n}$ 文字の coding $\rightarrow \langle \hat{L}_{new} \rangle \frac{N}{n}$
 Ω 文字 $\rightarrow \Omega^n$ 文字
 $\approx H(P_{new}) \frac{N}{n}$
 $= H(P_{old}) N$

Shannonの定理

Ω 種類の文字 $i=1, \dots, \Omega$

各々の文字が確率 p_i で出現

(ランダム文章)

文章 \rightarrow 0,1 の文字列

符号: 文字 $i \rightarrow$ 長さ l_i の 0,1 の列

$$\langle \hat{L} \rangle = \sum_{i=1}^{\Omega} l_i p_i$$

文字列の
ビット数

定理: 任意の接頭符号について

$H(P) \leq \langle \hat{L} \rangle$

$H(P) \gg 1$

つまり、接頭符号をくっつけて

$\langle \hat{L} \rangle \leq H(P) + 1$

出てくる
文字
 $\langle \hat{L} \rangle \approx H(P)$

$H(P) \gg 1$ ではないときは?

n 文字をくっつけて「1文字」とみなす

N 文字

ABDA BADA BACA

n文字 n文字 n文字

$\frac{N}{n}$ 文字

$nH(P_{old}) = H(P_{new}) \gg 1$

N 文字の文章

文字列 $L \approx H(P_{new}) \cdot \frac{N}{n} = H(P_{old}) N$