

情報理論レポート1(情報量)解答例

提示：2008/10/8(水) 提出：2008/10/15(水)

1. 英文の情報量

(1) 単語 w が情報源 A から生成されたときの自己情報量を $I_A(w)$ 、情報源 B から生成されたときの自己情報量を $I_B(w)$ と表す。このとき、次の自己情報量をそれぞれ求めよ。

(解)

記号 α が情報源 A より生成されたときの自己情報量を $i_A(\alpha)$ と書き、情報源 B より生成されたときの自己情報量を $i_B(\alpha)$ とかく。

$$\begin{aligned} I_A(\text{electronics}) &= i_A(e) + i_A(l) + i_A(e) + i_A(c) + i_A(t) + i_A(r) + i_A(o) \\ &\quad + i_A(n) + i_A(i) + i_A(c) + i_A(s) \\ &= -2 \times \log(0.1183) - \log(0.0402) - 2 \times \log(0.0331) - \log(0.0849) \\ &\quad - \log(0.0645) - \log(0.0731) - \log(0.0702) - \log(0.0740) - \log(0.0723) \\ &\simeq 43.3 \quad [\text{bit}] \end{aligned}$$

$$\begin{aligned} I_A(\text{information}) &= i_A(i) + i_A(n) + i_A(f) + i_A(o) + i_A(r) + i_A(m) + i_A(a) \\ &\quad + i_A(t) + i_A(i) + i_A(o) + i_A(n) \\ &= -2 \times \log(0.0740) - 2 \times \log(0.0702) - \log(0.0206) \\ &\quad - 2 \times \log(0.0731) - \log(0.0645) - \log(0.0277) \\ &\quad - \log(0.0876) - \log(0.0849) \\ &\simeq 44.5 \quad [\text{bit}] \end{aligned}$$

$$\begin{aligned} I_A(\text{system}) &= i_A(s) + i_A(y) + i_A(s) + i_A(t) + i_A(e) + i_A(m) \\ &= -2 \times \log(0.0723) - \log(0.0186) - \log(0.0849) \\ &\quad - \log(0.1183) - \log(0.0277) \\ &\simeq 25.1 \quad [\text{bit}] \end{aligned}$$

情報源 B では、各記号が均等に現れるので、次式が成り立つ。

$$\begin{aligned}
i_B(a) &= i_B(b) = \dots = i_B(z) \\
&= -\log P(\alpha) \\
&= -\log \frac{1}{26} \\
&\simeq 4.70 \text{ [bit]}
\end{aligned}$$

よって、情報源 B から生成されたの単語の情報量は文字数だけで定まる。

$$I_B(\text{electronics}) = 11 \times i_B(\alpha) \simeq 51.7 \text{ [bit]}$$

$$I_B(\text{information}) = 11 \times i_B(\alpha) \simeq 51.7 \text{ [bit]}$$

$$I_B(\text{system}) = 6 \times i_B(\alpha) \simeq 28.2 \text{ [bit]}$$

これらの結果の注意点は、以下のとおり。

○記号の出現確率が不均等な場合、記号数が同じでも自己情報量は異なる。

○同じ記号列だとしても、どの情報源から生成されたかによって、自己情報量は異なる。

(2)情報源 A のエントロピー（平均情報量） $H(A)$ および情報源 B のエントロピー $H(B)$ をそれぞれ求めよ。

$$\begin{aligned}
H(A) &= -\sum_{\alpha \in A} P(\alpha) i_A(\alpha) \\
&= -\sum_{\beta \in B} P(\alpha) \log P(\alpha) \\
&= -\left(0.1183\right) \log \left(0.1183\right) - \left(0.0876\right) \log \left(0.0876\right) - \left(0.0849\right) \log \left(0.0849\right) - \left(0.0740\right) \log \left(0.0740\right) \\
&\quad - \left(0.0731\right) \log \left(0.0731\right) - \left(0.0723\right) \log \left(0.0723\right) - \left(0.0702\right) \log \left(0.0702\right) - \left(0.0645\right) \log \left(0.0645\right) \\
&\quad - \left(0.0441\right) \log \left(0.0441\right) - \left(0.0402\right) \log \left(0.0402\right) - \left(0.0402\right) \log \left(0.0402\right) - \left(0.0331\right) \log \left(0.0331\right) \\
&\quad - \left(0.0306\right) \log \left(0.0306\right) - \left(0.0277\right) \log \left(0.0277\right) - \left(0.0210\right) \log \left(0.0210\right) - \left(0.0206\right) \log \left(0.0206\right) \\
&\quad - \left(0.0189\right) \log \left(0.0189\right) - \left(0.0186\right) \log \left(0.0186\right) - \left(0.0176\right) \log \left(0.0176\right) - \left(0.0171\right) \log \left(0.0171\right) \\
&\quad - \left(0.0104\right) \log \left(0.0104\right) - \left(0.0093\right) \log \left(0.0093\right) - \left(0.0021\right) \log \left(0.0021\right) - \left(0.0018\right) \log \left(0.0018\right) \\
&\quad - \left(0.0009\right) \log \left(0.0009\right) - \left(0.0009\right) \log \left(0.0009\right) \\
&\simeq 4.19 \text{ [bit / 記号]}
\end{aligned}$$

$$\begin{aligned}
H(B) &= -\sum_{\beta \in B} P(\beta) i_B(\beta) \\
&= -\sum_{\beta \in B} P(\beta) \log P(\beta) \\
&= -26 \times \frac{1}{26} \log \frac{1}{26} \\
&= \log 26 \\
&\simeq 4.70 [\text{bit} / \text{記号}]
\end{aligned}$$

よって、 $H(A) \leq H(B)$ が成り立つ。通常の英文の（平均の）情報量は、情報源 B から得られる同じ長さの記号列の（平均の）情報量より、小さいことを意味する。

一般に、同じ記号を生成するような事象系同士 X, Y において、事象の実現確率が不均等な事象系 X のエントロピーより、確率が均等な事象系 Y のエントロピーのほうが大きい。すなわち、次式が成り立つ。

$$H(X) \leq H(Y)$$

なお、記号の出現確率が均等な場合、ある記号 1 個の自己情報量は、エントロピーと等しい。（平均の意味から考えて、当然である。）