

情報理論レポート1(情報量)解答例

提示：2010/10/6(水) 提出：2010/10/20(水)

1. 英文の情報量

26文字からなるアルファベット $\mathcal{A} = \{a, b, \dots, z\}$ の各記号 $\alpha \in \mathcal{A}$ の出現確率 $P(\alpha)$ は、次の表で与えられる。(<http://www7.plala.or.jp/dvorakjp/hinshutu.htm> 参照。)

α	e	a	t	i	o	s	n	r
$P(\alpha)$	0.1183	0.0876	0.0849	0.0740	0.0731	0.0723	0.0702	0.0645
α	h	l	d	c	u	m	p	f
$P(\alpha)$	0.0441	0.0402	0.0402	0.0331	0.0306	0.0277	0.0210	0.0206
α	g	y	w	b	v	k	j	x
$P(\alpha)$	0.0189	0.0186	0.0176	0.0171	0.0104	0.0093	0.0021	0.0018
α	q	z						
$P(\alpha)$	0.0009	0.0009						

この出現確率 $P(\alpha)$ 基づいた情報源を

$$A = \left\{ \begin{matrix} a & , & b & , & \cdots & , & z \\ P(a) & , & P(b) & , & \cdots & , & P(z) \end{matrix} \right\}$$

とし、

アルファベット中の各記号 $\alpha \in \mathcal{A}$ が均等に現れる情報源を

$$B = \left\{ \begin{matrix} a & , & b & , & \cdots & , & z \\ \frac{1}{26} & , & \frac{1}{26} & , & \cdots & , & \frac{1}{26} \end{matrix} \right\}$$

とする。

また、単語 w が情報源 A から生成される確率を $P_A(w)$ 、自己情報量を $I_A(w)$ と表し、情

報源 B から生成されたときの確率を $P_B(w)$ 、自己情報量を $I_B(w)$ と表す。

このとき次の問い合わせよ。

(1) $w_1 = \text{information}$ 、 $w_2 = \text{theory}$ 、 $w_3 = \text{shannon}$ とし、次の確率および自己情報量を有効数字 3 衔で求めよ。

$$\begin{aligned} & P_A(w_1), I_A(w_1), P_B(w_1), I_B(w_1), \\ & P_A(w_2), I_A(w_2), P_B(w_2), I_B(w_2), \\ & P_A(w_3), I_A(w_3), P_B(w_3), I_B(w_3) \end{aligned}$$

(解答)

次のように計算される。

$$\begin{aligned} P_A(w_1) &= \prod_{\alpha \in w_1} P_A(\alpha) \\ &= P_A(i) \times P_A(n) \times P_A(f) \times P_A(o) \times P_A(r) \\ &\quad \times P_A(m) \times P_A(a) \times P_A(t) \times P_A(i) \times P_A(o) \times P_A(n) \\ &\simeq 0.0740 \times 0.0702 \times 0.0206 \times 0.0731 \times 0.0645 \\ &\quad \times 0.0277 \times 0.0876 \times 0.0849 \times 0.0740 \times 0.0731 \times 0.0702 \\ &\simeq 3.95 \times 10^{-14} \end{aligned}$$

$$\begin{aligned} I_A(w_1) &= -\log_2 P_A(w_1) \\ &= \sum_{\alpha \in w_1} i_A(\alpha) \\ &= i_A(i) + i_A(n) + i_A(f) + i_A(o) + i_A(r) \\ &\quad + i_A(m) + i_A(a) + i_A(t) + i_A(i) + i_A(o) + i_A(n) \\ &\simeq -\log 0.0740 - \log 0.0702 - \log 0.0206 - \log 0.0731 - \log 0.0645 \\ &\quad - \log 0.0277 - \log 0.0876 - \log 0.0849 - \log 0.0740 - \log 0.0731 - \log 0.0702 \\ &\simeq 3.76 + 3.83 + 5.60 + 3.77 + 3.95 \\ &\quad + 5.17 + 3.51 + 3.56 + 3.76 + 3.77 + 3.83 \\ &\simeq 44.5 [\text{bit}] \end{aligned}$$

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。

$$P_B(a) = P_B(b) = \dots = P_B(z) = \frac{1}{26} \simeq 0.0038$$

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

よって、情報源 B からの文字列出現確率や自己情報量は文字数だけで定まる。

$$\begin{aligned} P_B(w_1) &= \prod_{\alpha \in w_1} P_B(\alpha) \\ &= P_B(i) \times P_B(n) \times P_B(f) \times P_B(o) \times P_B(r) \\ &\quad \times P_B(m) \times P_B(a) \times P_B(t) \times P_B(i) \times P_B(o) \times P_B(n) \\ &= \frac{1}{26^{11}} \\ &\simeq 2.72 \times 10^{-16} \end{aligned}$$

$$\begin{aligned} I_B(w_1) &= \sum_{\alpha \in w_1} i_B(\alpha) \\ &= 11 \times i_B(\alpha) \\ &= 11 \times \log 26 \\ &\simeq 51.7 \text{ [bit]} \end{aligned}$$

$$\begin{aligned} P_A(w_2) &= P_A(t) \times P_A(h) \times P_A(e) \times P_A(o) \times P_A(r) \times P_A(y) \\ &\simeq 3.88 \times 10^{-8} \end{aligned}$$

$$\begin{aligned} I_A(w_2) &= i_A(t) + i_A(h) + i_A(e) + i_A(o) + i_A(r) + i_A(y) \\ &\simeq 3.56 + 4.50 + 3.08 + 3.77 + 3.95 + 5.75 \\ &\simeq 24.6 \text{ [bit]} \end{aligned}$$

$$P_B(w_2) = \frac{1}{26^6} \simeq 3.88 \times 10^{-8}$$

$$I_B(w_2) = 6 \times i_B(\alpha) = 6 \times \log 26 \simeq 28.2 \text{ [bit]}$$

$$P_A(w_3) = P_A(s) \times P_A(h) \times P_A(a) \times P_A(n) \times P_A(n) \times P_A(o) \times P_A(n)$$

$$\simeq 7.06 \times 10^{-9}$$

$$I_A(w_3) = i_A(s) + i_A(h) + i_A(a) + i_A(n) + i_A(n) + i_A(o) + i_A(n)$$

$$\simeq 3.79 + 4.50 + 3.51 + 3.83 + 3.83 + 3.77 + 3.83$$

$$\simeq 27.1 [bit]$$

$$P_B(w_3) = \frac{1}{26^7} \simeq 7.06 \times 10^{-9}$$

$$I_B(w_3) = 7 \times i_B(\alpha) = 7 \times \log 26 \simeq 32.9 [bit]$$

これらの結果より、「記号の出現確率が不均等な場合、記号列（記号数および記号の並び方）が同じでも出現確率および自己情報量は異なる。」ことが分かる。

(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。

情報源 A に関しては次の通り。

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

情報源 B に関しては次の通り。

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

なお、情報源記号が同じ情報源に対して、エントロピーは、情報源記号の生成確率が一定の場合に最大となる。