



HOKKAIDO UNIVERSITY

Title	2005年度 情報理論講義ノート
Author(s)	井上, 純一; Inoue, Jun-ichi
Description	http://www005.upp.so-net.ne.jp/j_inoue/index.html http://chaosweb.complex.eng.hokudai.ac.jp/~j_inoue/
Issue Date	2005-11-18T09:19:52Z
Doc URL	https://hdl.handle.net/2115/772
Rights(URL)	https://creativecommons.org/licenses/by-nc-sa/2.1/jp/
Type	learning object
File Information	InfoTheory05_5.pdf, 第5回講義ノート



情報理論 配布資料 #5

担当：井上 純一 (情報科学研究科棟 8-13)

URL : http://chaosweb.complex.eng.hokudai.ac.jp/~j_inoue/

平成 17 年 5 月 23 日

目次

4.6	情報源符号の平均符号長	33
4.7	一意復号可能である場合の平均符号長の下限	33
4.8	クラフト不等式からの考察	34
4.9	n 元拡大情報源	35

演習問題 4 の解答例

演習問題 3 の結果と確率に関する積の公式：

$$P_{A,B}(x, y) = P_{A|B}(x|y)P_B(y) \quad (143)$$

から直ちに

$$P_{A,B}(H, T) = P_{A|B}(H|T)P_B(T) = \frac{q}{2} \quad (144)$$

$$P_{A,B}(H, H) = P_{A|B}(H|H)P_B(H) = \frac{1}{2}(1-q) \quad (145)$$

$$P_{A,B}(T, H) = P_{A|B}(T|H)P_B(H) = \frac{q}{2} \quad (146)$$

$$P_{A,B}(T, T) = P_{A|B}(T|T)P_B(T) = \frac{1}{2}(1-q) \quad (147)$$

と各々の同時分布が求まり、また、

$$P_A(H) = P_A(T) = \frac{1}{2}, \quad P_B(H) = P_B(T) = \frac{1}{2} \quad (148)$$

であったから、2 つの確率分布 $P_{A,B}(x, y), P_A(x) \cdot P_B(y)$ の間の距離である KL 情報量は

$$\begin{aligned} D(P_{A,B} || P_A \cdot P_B) &= \sum_{x=\{H,T\}} \sum_{y=\{H,T\}} P_{A,B}(x, y) \log \left\{ \frac{P_{A,B}(x, y)}{P_A(x) \cdot P_B(y)} \right\} \\ &= P_{A,B}(H, H) \log \left\{ \frac{P_{A,B}(H, H)}{P_A(H)P_B(H)} \right\} + P_{A,B}(H, T) \log \left\{ \frac{P_{A,B}(H, T)}{P_A(H)P_B(T)} \right\} \\ &+ P_{A,B}(T, H) \log \left\{ \frac{P_{A,B}(T, H)}{P_A(T)P_B(H)} \right\} + P_{A,B}(T, T) \log \left\{ \frac{P_{A,B}(T, T)}{P_A(T)P_B(T)} \right\} \\ &= \frac{1}{2}(1-q) \log \left\{ \frac{(1-q)/2}{1/4} \right\} + \frac{q}{2} \log \left\{ \frac{q/2}{1/4} \right\} \\ &+ \frac{q}{2} \log \left\{ \frac{q/2}{1/4} \right\} + \frac{1}{2}(1-q) \log \left\{ \frac{(1-q)/2}{1/4} \right\} \\ &= 1 + q \log q + (1-q) \log(1-q) \end{aligned} \quad (149)$$

となり、これは **演習問題 3** で求めた相互情報量 $I(A; B)$ と一致する。従って関係式：

$$I(A; B) = D(P_{A,B}|P_A \cdot P_B) \quad (150)$$

が成り立つ。

4.6 情報源符号の平均符号長

$\mathcal{A} = \{a_1, a_2, \dots, a_K\}$, $\phi: \mathcal{A} \mapsto \mathcal{B}^+$, $\mathcal{B}^+ = \{b_1, b_2, \dots, b_M\}$ とする。また、 l_i を符号 $\phi(a_i)$ の長さとし、 $p_i(a_i)$ を記号 a_i の出現確率とする。前に見た例で言うならば、 $\mathcal{A} = \{aa, ab, ba, bb\}$ を $\mathcal{B} = \{1, 0\}$ で符号化する場合を考えれば、 $K = 4, M = 2$ であり、 $a_1 = aa, a_2 = ab, a_3 = ba, a_4 = bb$ とすれば、それぞれの出現確率が $p_1 = p(aa), p_2 = p(ab), p_3 = p(ba), p_4 = p(bb)$ で与えられることになる。

このとき、この符号 ϕ の平均符号長 L は

$$L = \sum_{i=1}^K p_i l_i \quad (151)$$

で与えられる。

4.7 一意復号可能である場合の平均符号長の下限

$H_M(X)$ を対数の底を M とした場合の情報源のエントロピーであるとすれば、平均符号長 L からこのエントロピーを差し引いたものは

$$\begin{aligned} L - H_M(X) &= \sum_{i=1}^K p_i l_i - \left(- \sum_{i=1}^K p_i \log_M p_i \right) \\ &= - \sum_{i=1}^K p_i \log_M M^{-l_i} + \sum_{i=1}^K p_i \log_M p_i \end{aligned} \quad (152)$$

と書ける。そこで、定数 c を

$$M^{-l_i} \equiv cr_i \quad (153)$$

で定義するのであれば、(152) 式を次のように書き直すことができる。

$$\begin{aligned} L - H_M(X) &= - \sum_{i=1}^K p_i \log_M cr_i + \sum_{i=1}^K p_i \log_M p_i \\ &= - \sum_{i=1}^K p_i \{ \log_M c + \log_M r_i \} + \sum_{i=1}^K p_i \log_M p_i \\ &= \sum_{i=1}^K p_i \log_M \left\{ \frac{p_i}{r_i} \right\} - \log_M c \end{aligned} \quad (154)$$

ここで、もちろん $\sum_{i=1}^K p_i = 1$ が成り立っていることに注意する。そこで、(153) 式で定義した r_i が p_i と同様に「確率」の意味を持つのであれば、上式の右辺第 1 項は、確率 $\mathbf{p} = (p_1, \dots, p_i, \dots, p_K)$ と確率 $\mathbf{r} = (r_1, \dots, r_i, \dots, r_K)$ の間の KL 情報量を表していることになるので

$$L - H_M(X) = D_M(\mathbf{p}||\mathbf{r}) - \log_M c \quad (155)$$

と書けることになる。

さて、 $r_i = M^{-l_i}/c$ を「確率」とみなすためには、 r_i を全ての i に関して足しあげたものは必ず 1 にならなければならない(確率の規格化の条件)。式で書けば

$$\sum_{i=1}^K r_i = \frac{1}{c} \sum_{i=1}^K M^{-l_i} = 1 \quad (156)$$

すなわち、定数 c は

$$c = \sum_{i=1}^K M^{-l_i} \quad (157)$$

でなければならない。また、この定数 c のとりうる値の制限がクラフト不等式で与えられる。前回に見たクラフト不等式から、ここで考える符号が一意復号可能であるためには

$$c = \sum_{i=1}^K M^{-l_i} \leq 1 \quad (158)$$

が成り立たなければならない。また、この事実と KL 情報量は負にならないという条件から、結局

$$L - H_M(X) \geq 0 \quad (159)$$

つまり、

$$L \geq H_M(X) \quad (160)$$

が成り立つ。この不等式は、どのような符号化法により、情報源を圧縮しても、その符号が一意復号可能であるためには、その平均符号長を情報源のエントロピーより小さくすることができないことを意味している。どんなにがんばっても平均符号長が情報源のエントロピーより小さい一意復号可能な符号を構成することはできないのである。

4.8 クラフト不等式からの考察

一意復号可能であるための平均符号長の満たすべき長さに関しては、上で述べたようなプロセスを経ずに、クラフト不等式からダイレクトに導くこともできる。

記号 a_i を符号化したもの $\phi(a_i)$ の長さ l_i を

$$l_i = \left\lceil \log_M \frac{1}{p_i} \right\rceil \quad (161)$$

と置くと、この l_i ($i = 1, \dots, K$) はクラフトの不等式を満たす。ここで、記号 $\lceil w \rceil$ は w 以上の最小の整数を表すものと約束する。従って、(161) 式はこの記号を使わずに次のように書き直すことができる。

$$\log_M \frac{1}{p_i} \leq l_i < \log_M \frac{1}{p_i} + 1 \quad (162)$$

辺々に p_i をかけて $\sum_{i=1}^K (\dots)$ のように和をとると

$$\sum_{i=1}^K p_i \log_M \frac{1}{p_i} \leq \sum_{i=1}^K p_i l_i < \sum_{i=1}^K p_i \log_M \frac{1}{p_i} + 1 \quad (163)$$

すなわち

$$H_M(X) \leq L < H_M + 1 \quad (164)$$

が得られる。従って、一意復号可能な符号の平均符号長は上記の不等式を満たすことになる。

4.9 n 元拡大情報源

系列 $x_1 \cdots x_n$ に割り当てる符号語 $\phi(x_1 \cdots x_n)$ の長さを $l(x_1 \cdots x_n)$ とする. このとき, 情報源記号 1 つあたりの平均符号長を

$$L_n = \frac{1}{n} \sum_{x_1} \sum_{x_2} \cdots \sum_{x_n} p_n(x_1, \dots, x_n) l(x_1 \cdots x_n) \quad (165)$$

で定義する. ここで, x_1, \dots, x_n が全て独立に生成されるものとするれば, 結合分布 $p_n(x_1, \dots, x_n)$ はそれぞれの分布 $p(x_i)$ の積で書くことができ

$$p_n(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i) \quad (166)$$

である. こうすれば, n 源拡大情報源は同じ分布を持つ n 個の確率変数 X_1, \dots, X_n からなる情報源と考えることができる. 従って, 平均符号長に関する不等式は

$$H_M(X_1, \dots, X_n) \leq \sum_{x_1} \cdots \sum_{x_n} p_n(x_1, \dots, x_n) l(x_1 \cdots x_n) < H_M(X_1, \dots, X_n) + 1 \quad (167)$$

と書けるが, (166) 式から

$$H_M(X_1, \dots, X_n) = nH_M(X) \quad (168)$$

であるから, 不等式 (167) は

$$H_M(X) \leq L_n < H_M(X) + \frac{1}{n} \quad (169)$$

となり, 従って, $n \rightarrow \infty$ の極限で L_n は

$$L_n = H_M(X) \quad (170)$$

へと収束し, n を長くすることで情報源記号 1 つあたりの平均符号長は情報源のエントロピーに限りなく近づくことになる.

ここまで学んだ内容を確認するために, 次の例題 7, 例題 8 を見ていこう.

例題 7

次の表に与えた符号 ψ_2 (教科書 p.31 参照) について以下の問いに答えよ.

aa	00
ab	10
ba	11
bb	110

- (1) 情報源の生成確率が $p(aa) = p(ab) = p(ba) = p(bb) = 1/4$ のとき, 情報源のエントロピー H , 及び, 平均符号長 L を求めよ.
- (2) 情報源の生成確率を $p(aa) = p(ab) = p(ba) = p$, 及び, $p(bb) = 1 - 3p$ と選ぶとき

$$\Phi(p) \equiv L - H \quad (171)$$

を求め, Φ を最小にするような $p = p_*$ 及び, そのときの Φ の値 $\Phi(p_*)$ を求めよ.

- (3) (1)(2) の結果から, 符号 ψ_2 に関してわかることを簡潔に述べよ.

(解答例)

(1) 問題文中に与えた表に従って平均符号長 L , 及び, エントロピー H を計算すると

$$\begin{aligned} L &= \sum_{x=aa,ab,ba,bb} p(x)l(x) \\ &= p(aa)l(aa) + p(ab)l(ab) + p(ba)l(ba) + p(bb)l(bb) = \frac{9}{4} \end{aligned} \quad (172)$$

$$H = - \sum_{x=aa,ab,ba,bb} p(x) \log p(x) = \log 4 = 2 \quad (173)$$

となる.

(2) 情報源アルファベットの生成確率を $p(aa) = p(ab) = p(ba) = p, p(bb) = 1 - 3p$ のように p を用いて表現する場合, (1) と同様にして平均符号長, エントロピーとしてそれぞれ

$$L = 3 - 3p \quad (174)$$

$$H = -3p \log p - (1 - 3p) \log(1 - 3p) \quad (175)$$

が得られる. 従って, 両者の差である p の関数 Φ は

$$\begin{aligned} \Phi &= L - H \\ &= 3 - 3p + 3p \log p + (1 - 3p) \log(1 - 3p) \end{aligned} \quad (176)$$

であり, Φ を最小化する p の値は $(\partial\Phi/\partial p) = 0$ より

$$\log \left\{ \frac{p}{1 - 3p} \right\} = 1 \quad (177)$$

すなわち, $(p/1 - 3p) = 2$, つまり, $p_* = 2/7$ となり, このときの Φ の値として

$$\Phi(p_*) = 3 - \log 7 \simeq \underline{0.19} \quad (178)$$

が得られる. この値は, (1) で調べた全ての情報源アルファベットが等確率で現れる場合の $\Phi(1/4) = 0.25$ と比べて小さくなっていることがわかる.(3) $p(bb) = 1/7$ であるから, 生成確率 (出現確率) の低い情報源アルファベットには長い符号を割り振り, 逆に, 生成確率の高いアルファベットには短い符号を与えることにより, 平均符号長をエントロピーに近づけることができる. Φ は非負であるから (つまり, どんなに頑張っても平均符号長はエントロピーより小さくできない), この手続きにより平均符号長を短くすることができる¹

この例では4つの情報源アルファベットの生成確率を p で表し, それを Φ を最小化するという意味で最適化したが, 実際に個々の情報源アルファベットに生成確率が割りあてられた場合 (何らかの方法によりそれらの確率が事前に計測できた場合), いかにして符号を構成するかに関しては, 次回の講義で説明するハフマン符号が有効な方法の一つとして知られている (このハフマン符号は「平均符号長を最小にする」という意味で最適な符号化法 (圧縮法) である).

¹ 既に見たように, n 値エントロピー関数が最大となるのは n 個の事象が等確率で現れる場合であったから, 全ての情報源アルファベットが等確率で生じるという仮定の下でエントロピーは最大となっており, 平均符号長はこの最大エントロピーを下回らないわけだから, 非常に長い値を持つことになる. 情報源アルファベットの出現確率が偏り始め, ある文字が出現しやすい状況になるとエントロピーが減少し始める. そしてその分だけ平均符号長の下限も低くなる. この状況下で適切な戦略の下に符号化すればその「下限」に一致させる最適符号を構成することができる.

例題 8

N_k を長さ k である符号化系列の総数としよう.

情報源アルファベット	符号語
x_1	0
x_2	10
x_3	11

表 1 : この問題で考える符号.

表 1 の場合には

$$N_1 = 1 \quad (x_1)$$

$$N_2 = 3 \quad (x_1x_1, x_2, x_3)$$

$$N_3 = 5 \quad (x_1x_1x_1, x_1x_2, x_1x_3, x_2x_1, x_3x_1)$$

となる (括弧内は実際にその長さを与える情報源アルファベットの組み合わせ). このとき以下の問いに答えよ. ただし, この問題で考えるのは全て表 1 で与えられる符号であるとする.

- (1) w_n を表 1 における長さ n の符号語の個数とする. つまり, $w_1 = 1, w_2 = 2, w_n = 0 (n \geq 3)$ である. このとき, N_k を N_{k-1}, N_{k-2} , 及び, w_1, w_2 の中から必要なものを用いて表せ. ただし $k \geq 3$ とする.

- (2) (1) で得られた N_k に関する漸化式の解として

$$N_k = \lambda^k$$

を仮定する. このとき (1) で得られた漸化式を λ に関する方程式に書き直せ.

- (3) 初期条件: $N_0 = N_1 = 1$ のもとで (2) で得られた方程式を解くことにより, N_k を求めよ. また, 得られた N_k の正当性をチェックするために, 長さ $k = 4$ の符号化系列を与える情報源アルファベットの組み合わせを全て列挙せよ.

(解答例)

問題文の誘導に従えばよい.

- (1) 長さ k である符号化系列の総数は, 長さ $k-1$ である符号化系列と長さ 1 の符号化系列の並べ方, 長さ $k-2$ の符号化系列と長さ 2 の符号化系列の並べ方の和であるから

$$N_k = w_1 N_{k-1} + w_2 N_{k-2} \quad (k \geq 3) \tag{179}$$

と表すことができる.

(2) (1) で得られた N_k に関する 3 項間漸化式の解を

$$N_k = \lambda^k \quad (180)$$

と仮定しよう. すると, この (180) 式を (179) 式に代入することにより

$$\lambda^{k-2}(\lambda^2 - w_1\lambda - w_2) = 0 \quad (181)$$

が得られる. $\lambda \neq 0$ であり, 問題の表に与えられた符号語の個数より $w_1 = 1, w_2 = 2$ であるから, λ は 2 次方程式:

$$\lambda^2 - \lambda - 2 = 0 \quad (182)$$

の解である.

(3) この方程式 (182) は 2 つの異なる実根 $\lambda = 2, -1$ をもつ. 従って, 漸化式 (179) の解は 2^k と $(-1)^k$ の線形結合として与えられる. つまり, a_1, a_2 を定数として

$$N_k = a_1 2^k + a_2 (-1)^k \quad (183)$$

が問題の漸化式の解である. ただし, 定数 a_1, a_2 は初期条件 (ここではむしろ「境界条件」と言った方がよいかもしれない) から決まる. $N_0 = N_1 = 1$ より, この定数は直ちに $a_1 = 2/3, a_2 = 1/3$ と定まるので, 結局, 求める長さ k の符号化系列の総数 N_k は

$$N_k = \frac{2^{k+1}}{3} + \frac{1}{3}(-1)^k \quad (184)$$

となる.

さて, この解の正当性を確認しよう. まずは問題文にその値が与えられている N_1, N_2, N_3 は上 (184) 式に $k = 1, 2$, 及び, $k = 3$ を代入することにより, それぞれ

$$N_1 = \frac{2^2}{3} - \frac{1}{3} = 1 \quad (185)$$

$$N_2 = \frac{2^3}{3} + \frac{1}{3} = 3 \quad (186)$$

$$N_3 = \frac{2^4}{3} - \frac{1}{3} = 5 \quad (187)$$

となり, 一致することが確認できる. 問題の $k = 4$ のときには

$$N_4 = \frac{2^5}{3} + \frac{1}{3} = 11 \quad (188)$$

が得られるが, 実際にこの長さ 4 の符号化系列を与える情報源アルファベットの組み合わせを列挙してみると $\{x_1 x_1 x_1 x_1, x_1 x_1 x_2, x_1 x_2 x_1, x_2 x_1 x_1, x_1 x_1 x_3, x_1 x_3 x_1, x_3 x_1 x_1, x_2 x_3, x_3 x_2, x_2 x_2, x_3 x_3\}$ のように確かに 11 通りある.

演習問題 5

1. 本講義ノート中の式 (168) の成立を示せ.
2. 次の表で与えられる符号は一意復号可能であるか否かを判定せよ (そのように判定した理由も記すこと).

記号	符号語
x_1	a
x_2	c
x_3	ad
x_4	abb
x_5	bad
x_6	deb
x_7	bbcde