

語頭条件を満たさない単語集合をもつ Word-Valued Sourceの性質について

石田 崇[†] 後藤 正幸^{††} 松嶋 敏泰[†] 平澤 茂一[†]

[†] 早稲田大学工学部 〒169-8555 東京都新宿区大久保 3-4-1

^{††} 武蔵工業大学環境情報学部 〒224-0015 神奈川県横浜市都筑区牛久保西 3-3-1

E-mail: †{ishida,hirasawa}@hirasa.mgmt.waseda.ac.jp, †toshi@matsu.mgmt.waseda.ac.jp,
††goto@yc.musashi-tech.ac.jp

あらまし 情報源符号化における情報源モデルとして、“言語アルファベット情報源 (word-valued source)” が提案されている [1], [2]. 西新らは, i.i.d. 言語アルファベット情報源を, 可算アルファベット \mathcal{Y} 上の i.i.d. (定常無記憶) 情報源と, \mathcal{Y} から有限アルファベット \mathcal{X} の有限系列への写像 ϕ によって定義し, この情報源の漸近等分割性 (AEP) を示し, エントロピー・レートを与えた [1]. 後藤らはこれを定常エルゴード言語アルファベット情報源に対して一般化し, 同様の結果を示した [2]. これらの結果は, 写像 ϕ が prefix-free であるという条件もとで導かれている. 一方, 写像 ϕ が prefix-free でない場合については, エントロピー・レートの存在すら明らかではなく, 情報源のエントロピー密度レート [6] に対して西新ら [1] がその上界と石田ら [3] が下界を与えたにとどまっている. そこで本稿では, 数値計算によってエントロピー・レートや上界・下界の有効性について検証し, ϕ が prefix-free でない言語アルファベット情報源の性質について考察を行う.

キーワード 情報源符号化, 言語アルファベット情報源, エントロピー・レート, 漸近等分割性

Properties of a Word-valued Source with a Non-prefix-free Word Set

Takashi ISHIDA[†], Masayuki GOTO^{††}, Toshiyasu MATSUSHIMA[†], and Shigeichi HIRASAWA[†]

[†] School of Science and Engineering, Waseda University Okubo 3-4-1, Shinjuku-ku, Tokyo, 169-8555 Japan

^{††} Faculty of Environmental and Information Studies, Musashi Institute of Technology Ushikubo-nishi 3-3-1, Tuzuki-ku, Yokohama-shi, Kanagawa, 224-0015 Japan

E-mail: †{ishida,hirasawa}@hirasa.mgmt.waseda.ac.jp, †toshi@matsu.mgmt.waseda.ac.jp,
††goto@yc.musashi-tech.ac.jp

Abstract Recently, *word-valued source* is proposed as a new class of source models. A word-valued source is defined as a source which has a probability distribution over *word set*. When the word set is prefix-free, it has been shown that there exists entropy rate of the source with simple expression and the AEP holds. However, when the word set is not prefix-free, it has been shown only the upper bound and lower bound on the entropy rate of the source. In this paper, we verify the entropy rate of the source by numerical computations for some source models in order to clarify the properties of the word-valued source.

Key words source coding, word-valued source, entropy rate, asymptotic equipartition property (AEP)

1. まえがき

情報源符号化における情報源モデルとして、“言語アルファベット情報源 (word-valued source)” が提案されている [1], [2]. この情報源は, 情報源アルファベットの有限系列を単語と定義し, 単語単位で確率構造を有するモデルと解釈する事ができる.

西新らは, i.i.d. 言語アルファベット情報源 X を, 可算アルファベット \mathcal{Y} 上の i.i.d. (定常無記憶) 情報源 Y と, \mathcal{Y} から有限アルファベット \mathcal{X} の有限系列 (単語) への写像 ϕ によって定義し, この情報源の漸近等分割性 (AEP) を証明し, さらにエントロピー・レートを示した [1]. 後藤らはこれを一般化し, 定常エルゴード言語アルファベット情報源に対して, 同様の結果を導いた [2]. また, 後藤ら [2] や石田ら [4], [5] は言語アルファベット情報源に対するユニバーサル符号の漸近的な性質についても

議論をしている. これらの議論では, 主に写像 ϕ が prefix-free であることを仮定している. ここで, ϕ が prefix-free であるとは任意の単語が他の単語の語頭に一致していないことを意味し, このとき単語集合 \mathcal{W} は語頭条件を満たしているという. 一方, ϕ が prefix-free ではない言語アルファベット情報源は, エントロピー・レートが存在するのかについてもいまだ明らかにされておらず, 情報源のエントロピー密度レート [6] について西新ら [1] がその上界を, 石田ら [3] が下界を示しているのみである. したがって, これらの上界や下界の厳しさや, エントロピー・レートなどについてより詳細な解析が必要である.

そこで本稿では, ϕ が prefix-free ではない, すなわち単語集合 \mathcal{W} が語頭条件を満たさない言語アルファベット情報源の性質をさらに明らかにするために, i.i.d. 言語アルファベット情報源に対して, 様々な設定を与えて数値実験を行い, エントロ

ピー密度レートやその上界, 下界などに関する理論的結果の評価を行う。

2. 言語アルファベット情報源 [1], [2]

西新ら [1] によって i.i.d. 言語アルファベット情報源が提案され, 情報源のエントロピー・レートが示された。後に後藤ら [2] はそれを“定常エルゴード言語アルファベット情報源”に一般化し同様の結果を導いている。

2.1 定常エルゴード言語アルファベット情報源 [2]

確率変数 Y の無限系列 $\mathbf{Y} = (Y_1, Y_2, Y_3, \dots)$ を可算アルファベット \mathcal{Y} 上に値をとる定常エルゴード情報源とする。 \mathcal{X}^* は有限アルファベット \mathcal{X} 上の有限系列すべての集合を表すものとし, 写像 ϕ を $\phi: \mathcal{Y} \rightarrow \mathcal{X}^*$ で定める。

さらに w を $w = \phi(y)$ ($w \in \mathcal{X}^*, y \in \mathcal{Y}$) によって与えこれを単語とよぶ。ここで, ϕ の値域を \mathcal{W} と書き, これを単語集合とよぶ。さらに確率変数 W を $W = \phi(Y)$ によって定義し, $W = \phi(Y)$ の無限系列を $\mathbf{W} = (W_1, W_2, \dots)$ と表す。

ここで, 系列 $\mathbf{Y} = (Y_1, Y_2, Y_3, \dots)$ に対して $W_1 = \phi(Y_1), W_2 = \phi(Y_2), W_3 = \phi(Y_3), \dots$ の接続を $\phi(\mathbf{Y}) = \phi(Y_1)\phi(Y_2)\phi(Y_3)\dots = W_1W_2W_3\dots = \mathbf{W}$ と書く。このとき定常エルゴード言語アルファベット情報源 $\mathbf{X} = (X_1, X_2, X_3, \dots)$ は $\mathbf{X} \stackrel{\text{def}}{=} \phi(\mathbf{Y})$ で定義する。

確率変数列 \mathbf{X} の実現値 (情報源から出力される系列) を $\mathbf{x} = (x_1, x_2, x_3, \dots)$ で表す。また, $n = 1, 2, 3, \dots$ に対して長さ n の有限系列を $X^n = (X_1, X_2, \dots, X_n)$, $\mathbf{x}^n = (x_1, x_2, \dots, x_n)$ と書く。同様に $m = 1, 2, 3, \dots$ に対して長さ m の有限系列を $Y^m = (Y_1, Y_2, \dots, Y_m)$, $\mathbf{y}^m = (y_1, y_2, \dots, y_m)$, $W^m = (W_1, W_2, \dots, W_m)$, $\mathbf{w}^m = (w_1, w_2, \dots, w_m)$ とする。さらに写像 $\phi: \mathcal{Y}^m \rightarrow \mathcal{W}^m$ も同様に系列 $\phi(Y_1), \phi(Y_2), \dots, \phi(Y_m)$ の接続 $\phi(Y_1^m) = \phi(Y_1)\phi(Y_2)\dots\phi(Y_m) = W_1W_2\dots W_m = W^m$ によって定義する^(注1)。

言語アルファベット情報源からは, 系列は単語単位で出力されるが, それぞれの単語が接続された長さ n の \mathcal{X} 上の系列 x^n として観測される。すなわち $x^n = w^n$ であり, このとき任意の m に対して $n = |w_1| + |w_2| + \dots + |w_m|$ である。ここで $|w|$ は単語の長さを表す。本稿ではこれ以降, w^n を単語系列, x^n をシンボル系列とよぶことにする。

系列 \mathbf{Y}, \mathbf{W} の確率分布を次のように表す。

$$P_{Y^m}(\mathbf{y}^m) \stackrel{\text{def}}{=} \Pr\{Y^m = \mathbf{y}^m\}, \quad (1)$$

$$P_{W^m}(\mathbf{w}^m) \stackrel{\text{def}}{=} \Pr\{W^m = \mathbf{w}^m\}. \quad (2)$$

$m = 1$ の場合は $P_{W^1}(w) = P_W(w), P_{Y^1}(y) = P_Y(y)$ と書く。一方, 系列 \mathbf{X} の確率分布を同様に次のように表す。

$$P_{X^n}(\mathbf{x}^n) \stackrel{\text{def}}{=} \Pr\{X^n = \mathbf{x}^n\}. \quad (3)$$

$n = 1$ の場合は $P_{X^1}(x) = P_X(x)$ と書く。

2.2 エントロピー・レート

\mathbf{X} に対して

$$-\frac{1}{n} \log P_{X^n}(\mathbf{x}^n) \quad (4)$$

を情報源 \mathbf{X} のエントロピー密度レートという [6]。また, エントロピー密度レートの期待値について,

$$H(\mathbf{X}) = \lim_{n \rightarrow \infty} \left[-\frac{1}{n} \sum_{\mathbf{x}^n \in \mathcal{X}^n} P_{X^n}(\mathbf{x}^n) \log P_{X^n}(\mathbf{x}^n) \right], \quad (5)$$

(注1): W^m は確率変数列 (W_1, W_2, \dots, W_m) を表しているが, 特に問題がない場合には W_1, W_2, \dots, W_m の接続に対しても W^m という表記を用いることにする。 w^m についても同様。

が存在するとき $H(\mathbf{X})$ を言語アルファベット情報源 \mathbf{X} のエントロピー・レートとよぶ [7]。

単語系列 \mathbf{W} のエントロピー・レート $H(\mathbf{W})$, 平均単語長レート $E[|\mathbf{W}|]$ は以下で与えられる [2]。

$$H(\mathbf{W}) = \lim_{m \rightarrow \infty} \left[-\frac{1}{m} \sum_{\mathbf{w}^m \in \mathcal{W}^m} P_{W^m}(\mathbf{w}^m) \log P_{W^m}(\mathbf{w}^m) \right], \quad (6)$$

$$E[|\mathbf{W}|] \stackrel{\text{def}}{=} \lim_{m \rightarrow \infty} \frac{1}{m} E_{P_{W^m}} \left[\sum_{i=1}^m |W_i| \right]. \quad (7)$$

ここで, $E_P[\cdot]$ は確率分布 P による期待値を表す。

[定理 1] (定常エルゴード言語アルファベット情報源の AEP [2], [3]) \mathbf{Y} を定常エルゴード情報源とする。 $\mathbf{X} = \phi(\mathbf{Y})$, $H(\mathbf{Y}) < \infty, E[|\mathbf{W}|] < \infty$ のとき, 次式が成り立つ。

$$\limsup_{n \rightarrow \infty} \frac{1}{n} [-\log P_{X^n}(\mathbf{x}^n)] \leq \frac{H(\mathbf{Y})}{E[|\mathbf{W}|]}, \quad \text{a.s.} \quad (8)$$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} E_{P_{X^n}} [-\log P_{X^n}(\mathbf{x}^n)] \leq \frac{H(\mathbf{Y})}{E[|\mathbf{W}|]}. \quad (9)$$

が成り立つ。さらに, ϕ が prefix-free のとき,

$$\lim_{n \rightarrow \infty} \frac{1}{n} [-\log P_{X^n}(\mathbf{x}^n)] = \frac{H(\mathbf{W})}{E[|\mathbf{W}|]}, \quad \text{a.s.} \quad (10)$$

$$H(\mathbf{X}) = \frac{H(\mathbf{W})}{E[|\mathbf{W}|]}. \quad (11)$$

□

2.3 \mathcal{W} が語頭条件を満たさない言語アルファベット情報源

単語集合 \mathcal{W} が語頭条件を満たさない i.i.d. 言語アルファベット情報源を non-prefix-free i.i.d. 言語アルファベット情報源とよぶことにし, 次のように定義する [3]。

[定義 1] (non-prefix-free i.i.d. 言語アルファベット情報源 [3]) non-prefix-free i.i.d. 言語アルファベット情報源 $\mathbf{X} = \phi(\mathbf{Y})$ は, 有限アルファベット \mathcal{Y} 上に値をとる i.i.d. 情報源 \mathbf{Y} と, \mathcal{Y} から \mathcal{W} への 1 対 1 写像 $\phi: \mathcal{Y} \rightarrow \mathcal{W} = \cup_{i=1}^K \mathcal{X}^i$ によって与えられる。ただし, $\mathcal{X} = \{0, 1\}$ とする。 □

この情報源モデルでは, 単語集合 \mathcal{W} を深さ K の完全二分木で表現することができ, 各ノードが単語 w に対応する (図 1 参照)。このとき, 最大単語長は K であり, 単語集合 \mathcal{W} の要素数 $|\mathcal{W}| = 2(2^K - 1)$ である。単語集合 \mathcal{W} に対して確率分布の与え方によって様々な構造を持つ情報源モデルとなるが, 一般に \mathcal{W} は語頭条件を満たさない。石田らはこの情報源に対してエントロピー密度レートの下界式を導出した [3]。

[定理 2] (non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー密度レートの下界 [3]) non-prefix-free i.i.d. アルファベット情報源 \mathbf{X} について,

$$\liminf_{n \rightarrow \infty} \frac{1}{n} [-\log P_{X^n}(\mathbf{x}^n)] \geq \frac{H(\mathbf{W})}{E[|\mathbf{W}|]} - \frac{H(S)}{E[|\mathbf{W}|]}, \quad \text{a.s.} \quad (12)$$

$$\liminf_{n \rightarrow \infty} \frac{1}{n} E_{P_{X^n}} [-\log P_{X^n}(\mathbf{x}^n)] \geq \frac{H(\mathbf{W})}{E[|\mathbf{W}|]} - \frac{H(S)}{E[|\mathbf{W}|]}, \quad (13)$$

が成り立つ。ここで, $H(S)$ は単語の長さの確率分布 $P_S(s)$ のエントロピーを表し,

$$H(S) = - \sum_{s=1}^k P_S(s) \log P_S(s), \quad (14)$$

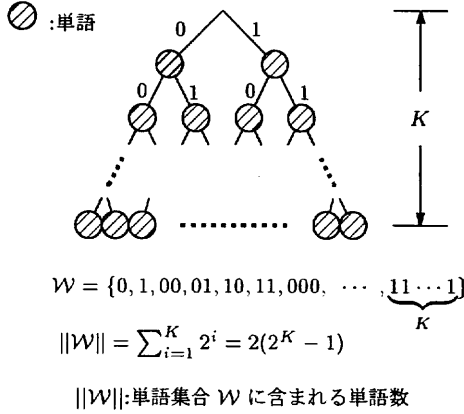


図1 単語集合 W の2分木表現

$$P_S(s) = \sum_{\{w:|w|=s\}} P_W(w), \quad (15)$$

で与えられる。 □

定理 1,2 より, non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レート $H(X)$ が存在するならば,

$$\frac{H(W)}{E[|W|]} - \frac{H(S)}{E[|W|]} \leq H(X) \leq \frac{H(W)}{E[|W|]}, \quad (16)$$

がいえる。

この情報源におけるエントロピー密度レートの上界式, 下界式は1本のシンボル系列 x^n に写像される単語系列 w^m の本数を評価することによって得られている。いま, 十分に長い単語系列 w^m は AEP により出現確率はほぼ $2^{-mH(W)}$ である。出現する x^n に対しては少なくとも1つの単語系列 w^m が写像されるので, $P_{X^n}(x^n) \geq 2^{-mH(W)}$ が成り立ち, これより上界式が得られる [1], [2]。

次に下界であるが, x^n を観測した場合, 複数の単語系列 w^m がこのシンボル系列に写像されており, 単語の切れ目がどこなのかを一意に決めることができない。しかし, 単語の切れ方を確定してしまえば, 同じような単語の切れ方をする異なった単語系列 w^m は必ず異なるシンボル系列 x^n に写像されることから, 単語の切れ目の入れ方の総数を評価すればよいことが分かり, これによって下界を導いている [3]。

以上のように, W が語頭条件を満たさない場合については, 定常エルゴード言語アルファベット情報源に対してエントロピー密度レートの上界 (式 (8),(9)), i.i.d. 言語アルファベット情報源に対してその下界 [3] が得られているのみで, エントロピー・レート $H(X)$ が存在するのか, いまだ明らかにされていない。また, 上界, 下界の有効性についても理論的に評価がされていない。

そこで本稿では数値実験によって, W が語頭条件を満たさない言語アルファベット情報源のエントロピー・レートやその性質に関する検証を行う。

3. 数値計算と考察

non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レートについてさらに議論を行うためには, シンボル系列 x^n の生起確率 $P_{X^n}(x^n)$ を陽な形で与える必要があると考えられるが, いまだその結果には至っていない。 $P_{X^n}(x^n)$ は x^n と w^m の写像によって決定されるが, この写像は単語集合 W の構造や, 単語系列 w^m における単語の並び方などに依存すると考えられ, 理論的な解析を行うことは容易ではない。

そこで本稿では, 数値計算によって以下の内容を検証する。

- non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レートは存在するのか。
- 上界, 下界の有効性はどうか。
- 情報源の確率構造によってエントロピー密度レートとその上界, 下界の性能はどのように変化するのか。

以上の内容を検証するために, いくつかの設定のもとで i.i.d. 言語アルファベット情報源に実際に単語の確率分布 $P_W(w)$ を数値で与え, 情報源から出力されたシンボル系列 x^n に対して次式によって生起確率 $P_{X^n}(x^n)$ を計算し, エントロピー密度レートを求める。

$$P_{X^n}(x^n) = \sum_{\{w^m: x^n=w^m\}} P_{W^m}(w^m). \quad (17)$$

なお, 計算ではすべて \log の底を 2 としている。

3.1 情報源のエントロピー密度レートの上界と下界の有効性について

まず始めに, 式 (16) について検証する。すなわち, 情報源のエントロピー・レートが存在するのか, また, 上界, 下界の性能はどれほどなのかを調べるために以下の4つの単語集合モデルについて数値計算を行う。

- [model 1] 全ての単語の長さが等しいモデル
- [model 2] W が語頭条件を満たすモデル
- [model 3] W が語頭条件を満たさないモデル
- [model 4] 単語の長さが全て異なるモデル

具体的なパラメータは以下のように与える。

[model 1] 全ての単語の長さが等しいモデル

$$W_1 = \{000, 001, 010, 011, 100, 101, 110, 111\}$$

$$\begin{aligned} P_W(000) &= 0.05, & P_W(001) &= 0.05, \\ P_W(010) &= 0.10, & P_W(011) &= 0.10, \\ P_W(100) &= 0.20, & P_W(101) &= 0.20, \\ P_W(110) &= 0.20, & P_W(111) &= 0.10. \end{aligned}$$

[model 2] W が語頭条件を満たすモデル

$$W_2 = \{0, 10, 110, 111\}$$

$$\begin{aligned} P_W(0) &= 0.50, & P_W(10) &= 0.20, \\ P_W(110) &= 0.20, & P_W(111) &= 0.10. \end{aligned}$$

[model 3] W が語頭条件を満たさないモデル

$$W_3 = \{0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, 101, 110, 111\}$$

$$\begin{aligned} P_W(0) &= 0.04, & P_W(1) &= 0.05, \\ P_W(00) &= 0.12, & P_W(01) &= 0.01, \\ P_W(10) &= 0.01, & P_W(11) &= 0.11, \\ P_W(000) &= 0.06, & P_W(001) &= 0.10, \\ P_W(010) &= 0.12, & P_W(011) &= 0.01, \\ P_W(100) &= 0.11, & P_W(101) &= 0.10, \\ P_W(110) &= 0.04, & P_W(111) &= 0.12. \end{aligned}$$

[model 4] 単語の長さが全て異なるモデル

$$W_4 = \{0, 10, 101\}$$

$$P_W(0) = 0.50, \quad P_W(10) = 0.30, \quad P_W(101) = 0.20.$$

単語集合モデル W_1, W_2, W_3, W_4 の木表現を図2に示す。ここで, 斜線の引かれたノードが出現する単語 ($P_W(w) > 0$) を表している。

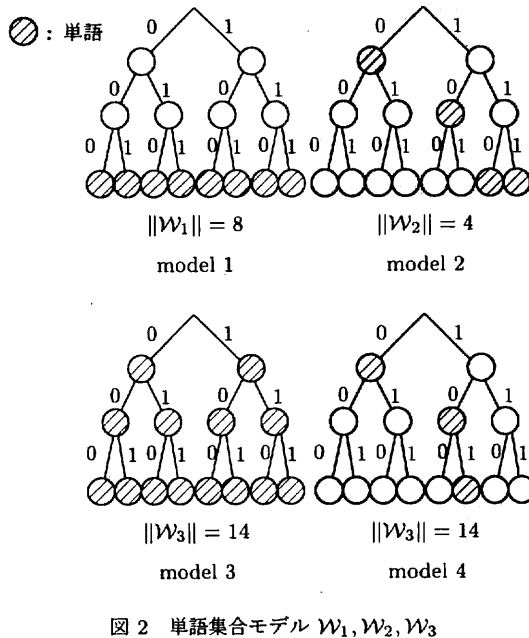


図2 単語集合モデル $\mathcal{W}_1, \mathcal{W}_2, \mathcal{W}_3$

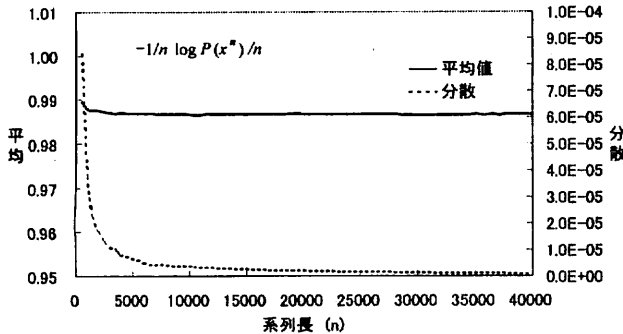


図3 系列長 n に対するエントロピー密度レートの推移

3.1.1 エントロピー・レートに関する結果と考察

まず、系列長 n に対して、エントロピー密度レートの値がどのように推移しているかを確認する。ここでは、model 3 についてのグラフを図3に示す。このグラフでは、 x^n の200系列それぞれのエントロピー密度レートの平均値（実線）と分散（破線）の推移が示されている。

この図よりエントロピー密度レートの平均値がほぼ収束していることから、non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レートが存在するのではないかと予想される。さらに、 n が大きくなるにつれて分散の値が限りなく0に近づいていることから、漸近等分割性も成り立っているのではないかと考えられる。さらに単語集合 \mathcal{W}_3 に対して model 3 とは異なる確率分布 $P_W(w)$ を割り振って計算を何通りか行ったが、いずれについても同様の結果が得られている。

3.1.2 エントロピー・レートの上界と下界の有効性に関する結果と考察

3.1.1の結果を踏まえて、以降の数値計算では $n = 40000$ のシンボル系列長に対して200系列のエントロピー密度レートの平均値を求め、これを non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レートであると考えことにする。

model 1 から model 4 までの計算結果を表1に示す。

model 1, model 2 は \mathcal{W} が語頭条件を満たしていることから、これらの情報源のエントロピー・レート $H(X)$ は存在し、式(11)から求められる。

(1) model 1 では上界と下界がエントロピー・レートと一致

表1 エントロピー密度レートの平均値と分散 ($n=40000$)

model	上界	平均	entropy rate	下界	分散
1	0.941	0.941	0.941	0.941	$3.92 \cdot 10^{-6}$
2	0.978	0.978	0.978	0.153	$1.60 \cdot 10^{-6}$
3	1.361	0.987	—	0.890	$8.50 \cdot 10^{-7}$
4	0.874	0.747	—	0.000	$3.21 \cdot 10^{-6}$

entropy rate : 式(11)より計算。

上界 : 式(16)の右辺より計算。

下界 : 式(16)の左辺より計算。

し、エントロピー密度レートの平均値もエントロピー・レートに一致していることが計算結果より分かる。式(16)より、上界と下界との差分は $H(S)/E\|\mathcal{W}\|$ で与えられるが、単語の長さが全て等しいモデルでは $H(S)$ が0となることから確かに下界が上界と一致するモデルであること分かる。単語の長さが全て等しいということはシンボル系列 x^n に対して単語の切れ目が一意に決定できることから、シンボル系列 x^n と単語系列 w^m が1対1に対応していることを意味している。

(2) model 2 は各単語の長さは一定値ではないものの \mathcal{W} が語頭条件を満たしていることから、上界とエントロピー・レートが一致するモデルである。しかし、このモデルでは $H(S)$ が必ずしも0とはならないことから下界は上界と一致しない。これは、下界を導出する過程で単語の長さのみに着目して x^n へ写像される w^m の本数を評価したために、本来1対1対応しているものを過大に評価しているためだといえる。

(3) model 3 は \mathcal{W} が語頭条件を満たさないため、陽な形でエントロピー・レートを求めることはできない。このモデルではエントロピー密度レートの平均値は上界にも下界にも一致しないが、計算結果より上界とは大きく離れて下界に近い値をとっていることが分かる。これは、下界の有効性を示しており、 x^n へ写像される w^m の本数は単語の長さの確率分布 $P_S(s)$ によってよい精度で評価できているのではないかと考えられる。単語集合 \mathcal{W}_3 に対して、異なる確率分布 $P_W(w)$ を与えたモデルについても、ほとんどの場合で同様の結果を得た。

(4) model 4 では下界値が0になっている。これは各単語の長さが全て異なり $H(S) = H(W)$ となるためで、式(16)からも明らかである。このとき下界は意味をなさない。これは語頭条件を満たしている場合にも起こる。

結果をまとめると以下のことが言える。

- non-prefix-free i.i.d. 言語アルファベット情報源のエントロピー・レートは存在すると考えられる。
- 一般的な non-prefix-free i.i.d. 言語アルファベット情報源に対して、上界はあまくなるが、下界はタイトであると考えられる。

3.2 情報源の確率構造とエントロピー・レートの振る舞いについて

前項の結果からも明らかなように、 \mathcal{W} が語頭条件を満たさないモデルでは確率分布 $P_W(w)$ の与え方によって様々な性質を有することがわかる。

そこで本項では \mathcal{W} が語頭条件を満たさないモデルの確率構造を明らかにすることを目的として、いくつかの場合を想定して単語の確率分布を与え、エントロピー密度レートの平均と上界、下界を計算し、それらがどのような性質をもつか数値計算によって検証する。

単語集合 \mathcal{W} が語頭条件を満たさない情報源モデルとして以下のような場合を考える。

- [case A] \mathcal{W} の語頭条件が徐々に崩れる場合
- [case B] $H(S)$ 一定のもとで $E\|\mathcal{W}\|$ が変化する場合
- [case C] 同じ長さの単語同士が等確率で出現する場合
- [case D] 上界と下界が一定の場合

表 2 model 5-1 - 5-4 の $P_W(w)$

	$P_W(0)$	$P_W(1)$	$P_W(10)$	$P_W(11)$
model 5-1	0.10	0.90	0.00	0.00
model 5-2	0.10	0.80	0.02	0.08
model 5-3	0.10	0.10	0.20	0.60
model 5-4	0.10	0.00	0.25	0.65

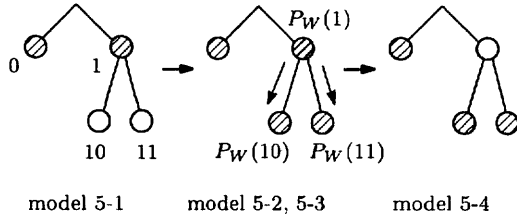


図 4 case A の単語集合

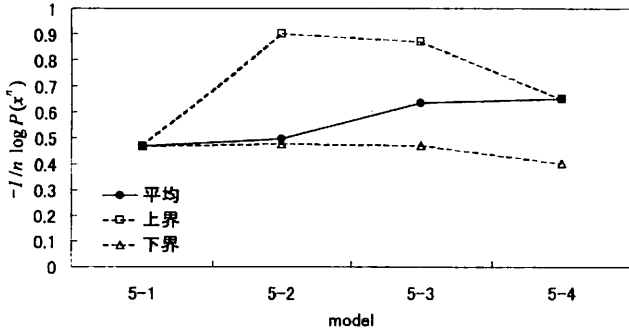


図 5 model 5-1 - 5-4 の数値計算結果

(1) case A : \mathcal{W} の語頭条件が徐々に崩れる場合

任意の s について単語長の確率分布 $P_S(s)$ が 0 か 1 の値をとるとき、このモデルは語頭条件を満たすモデルとなるが、わずかでも違う長さの単語に出現確率が割り当てられた瞬間、このモデルは語頭条件を満たさなくなる。

そこで、単純なモデルとして単語集合 $\mathcal{W}_5 = \{0, 1, 10, 11\}$ を考え、表 2 のように単語の出現確率 $P_W(w)$ を割り当てた model 5-1 - 5-4 を設定する。

これらのモデルは次のような構造になっている。model 5-1 は長さ 1 の単語しか出現しない、語頭条件を満たすモデルであるが、model 5-2, 5-3 になるにつれて次第に $P_W(1)$ を $P_W(10), P_W(11)$ に割り振っていき、model 5-4 で完全に $P_W(1) = 0$ となって再び語頭条件を満たすモデルとなる (図 4)。数値計算の結果を図 5 に示す。

model 5-1 は全ての単語長が等しいモデルなので上界、下界、エントロピー密度レートが情報源のエントロピー・レートと一致する特殊な場合である。model 5-2 は $P_W(1)$ を $P_W(10), P_W(11)$ に割り振ることによって model 5-1 の語頭条件がわずかに崩れたモデルとなっているが、上界が一気にエントロピー密度レートの平均値よりも大きくかけ離れてしまうことが分かる。一方、下界とエントロピー密度レートとの差はほとんどない。しかし、さらに $P_W(1)$ の割り振りを続けると徐々にエントロピー・密度レートの平均値と上界の値が近づいてゆき (model 5-3)、 $P_W(1) = 0$ となったとき (完全に語頭条件を満たすとき) 両者は完全に一致する。

(2) case B : $H(S)$ 一定のもとで $E[|W|]$ が変化する場合

case B では $H(S)$ が一定のもとで平均単語長 $E[|W|]$ が増加したときにエントロピー密度レートの値と上界、下界がどのような性能を示すのかを調べる。

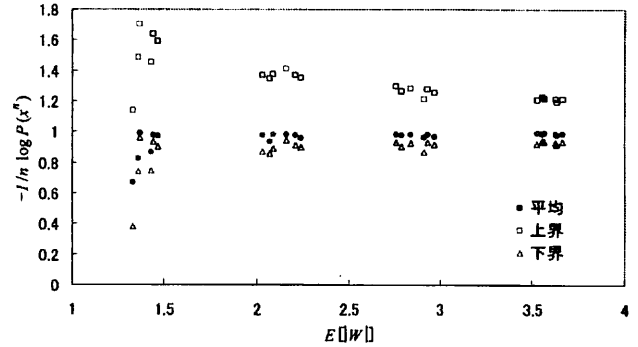


図 6 case B の数値計算結果

表 3 model 7-1 - 7-4 の $P_S(s)$

	$P_S(1)$	$P_S(2)$	$P_S(3)$	$P_S(4)$
model 7-1	1.00	0.00	0.00	0.00
model 7-2	0.40	0.60	0.00	0.00
model 7-3	0.20	0.30	0.40	0.00
model 7-4	0.10	0.20	0.30	0.40

このモデルでは、最大深さ $K = 4$ の完全 2 分木全てのノードに対応する単語集合 \mathcal{W}_6 を考える ($|\mathcal{W}_6| = 30$)。長さ s の出現確率 $P_S(s)$ ($s = 1, 2, 3, 4$) は 4 つの定数の組 $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ (ただし、 $\alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 = 1$) を割り当てることによって定めるものとする。このとき、 α の割り当て方の組み合わせは $4! = 24$ 通りあるので、それぞれを 1 つの情報源モデルとする (model 6-1 - model 6-24)。このとき、24 個全てのモデルの $H(S)$ は等しい。また、24 個の情報源モデルにおいて各単語の出現確率 $P_W(w)$ は割り当てられた $P_S(s)$ を満たすように、ランダムに割り振ることとする。

$(\alpha_1, \alpha_2, \alpha_3, \alpha_4) = (0, 80, 0.10, 0.07, 0.03)$ として、24 個のモデルに対する数値実験の結果を図 6 に示す。このグラフは 24 個それぞれのモデルについて求めたエントロピー密度レートの平均値、上界、下界の値を、横軸に平均単語長 $E[|W|]$ をとってプロットしたものである。

グラフより、 $H(S)$ 一定のもとで平均単語長 $E[|W|]$ の増加にともなって上界と下界の差が小さくなっていく様子が分かる。これは、上界と下界の差分が $H(P_S)/E[|W|]$ で与えられることから、明らかな結果である。

(3) case C : 同じ長さの単語同士が等確率で出現する場合

最大深さ $K = 4$ の完全 2 分木によって定まる単語集合 \mathcal{W}_7 について、長さが等しい単語同士が等確率で出現するような情報源モデルを考える。すなわち、単語長の確率分布 $P_S(s)$ が与えられたもとで、長さが s の単語の出現確率が $P_S(s)/2^s$ で与えられるモデルである。単語長の確率分布 $P_S(s)$ ($s = 1, 2, 3, 4$) を表 3 のように与える。

このモデルに対する数値計算結果を図 7 に示す。グラフから、エントロピー密度レート平均値と下界は常に一致し、1.0 となっていることが分かる。このような確率構造を持つモデルについては、長さ s の単語の個数が 2^s 個で、その出現確率が全て $P_S(s)/2^s$ であることから、式 (16) の左辺 (下界式) より

$$\begin{aligned} & \frac{H(W)}{E[|W|]} - \frac{H(S)}{E[|W|]} \\ &= \frac{-\sum_w P_W(w) \log P_W(w) + \sum_s P_S(s) \log P_S(s)}{E[|W|]} \\ &= \frac{-\sum_s P_S(s) \log(P_S(s)/2^s) + \sum_s P_S(s) \log P_S(s)}{E[|W|]} \end{aligned}$$

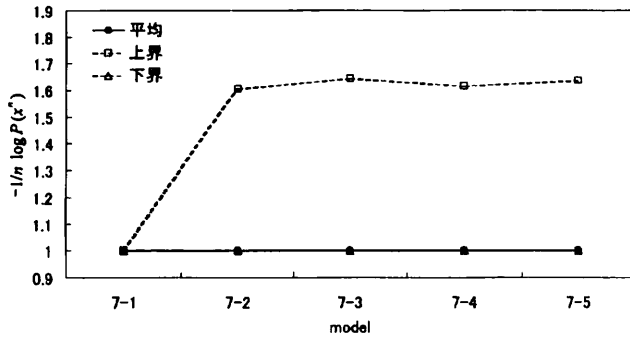


図7 case C の数値計算結果

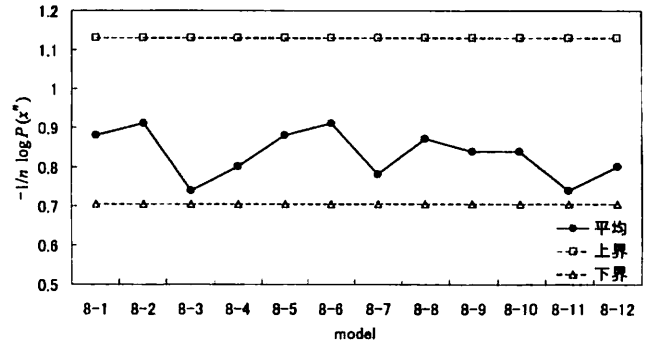


図8 model 8 の数値計算結果

表4 model 8-1 - 8-12 の $P_W(w)$

	$P_W(00)$	$P_W(01)$	$P_W(10)$	$P_W(11)$
model 8-1	0.03	0.07	0.00	0.00
model 8-2	0.07	0.03	0.00	0.00
model 8-3	0.00	0.00	0.03	0.07
model 8-4	0.00	0.00	0.07	0.03
model 8-5	0.03	0.00	0.07	0.00
model 8-6	0.07	0.00	0.03	0.00
model 8-7	0.03	0.00	0.00	0.07
model 8-8	0.07	0.00	0.00	0.03
model 8-9	0.00	0.03	0.07	0.00
model 8-10	0.00	0.07	0.03	0.00
model 8-11	0.00	0.03	0.00	0.07
model 8-12	0.00	0.07	0.00	0.03

全ての model について $P_W(0) = 0.2$, $P_W(1) = 0.7$

$$= \frac{\sum_s s \cdot P_S(s)}{E[|W|]} = 1, \quad (18)$$

より、下界値が1となることは容易に導かれる。また、エントロピー密度レートの値が下界と一致するという事は、[3]で下界を導出する際に評価した x^n へ写像される w^m の本数が等号で成り立っていることが分かる。

(4) case D: 上界と下界が一定の場合

最後に、上界と下界は一定であるが確率構造の異なる情報源モデルを考える。 $W_8 = \{0, 1, 00, 01, 10, 11\}$ とし、単語0と1の生起確率を $P_W(0) = 0.2$, $P_W(1) = 0.7$ とし、残りの0.1を表4に示すように割り当てる。すなわち、case Dでは長さ1の単語の生起確率を固定し、長さ2の単語について生起確率の値を入れ替える事によってできるモデルの集合を考えている。

各モデルで単語の生起確率 $P_W(w)$ がとる値の組み合わせはどのモデルについても共通であり、また、 $P_S(s)$ も共通であることから、どのモデルの $H(W)$ も $H(S)$ も等しい。さらに、平均単語長 $E[|W|]$ も共通なので、すべての情報源について式(16)によって計算される上界、下界は等しい。

case Dの数値計算結果を図8に示す。この結果から、上界、下界の値が同じでもエントロピー密度レートの値はモデルによって異なる事がわかる。すなわち、 W が語頭条件を満たさない言語アルファベット情報源のエントロピー・レートの性質を明らかにするにはこれまでの解析でなされてきた議論やパラメータだけでは不十分である可能性を示唆している。

しかし、model 8-1とmodel 8-5やmodel 8-2とmodel 8-6などはエントロピー密度レートが同じ値になっており、規則性が見られる。これらは $P_W(01)$ と $P_W(10)$ の確率が入れ替わっただけのモデル同士であるが、このような情報源は同じエントロピー・レートの値を持つと考えられる。

4. むすび

本稿では、non-prefix-free i.i.d. 言語アルファベット情報源について、エントロピー・レートや情報源の構造に関する性質を明らかにすることを目的として、様々な設定のもとで情報源モデルに対して数値計算を行い、エントロピー密度レートやその上界、下界の値について検証を行った。検証の結果として、以下のようなことが挙げられる。

non-prefix-free i.i.d. 言語アルファベット情報源について、

- エントロピー・レートが存在すると考えられる。
- AEP が成り立つと考えられる。
- 一般に式(16)の上界式はあく、下界式タイトである。
- 単語の確率分布 $P_W(w)$ の与え方によってエントロピー密度レートとその上界、下界は大きく異なる性質を有する。
- 情報源の確率構造をさらに解析するためには新たな視点が必要であると考えられる。

しかし、non-prefix-free i.i.d. 言語アルファベット情報源に関する理論的な解析はまだ十分になされていない。本稿における数値計算から得られた結果や、予想をもとにさらに語頭条件を満たさない情報源モデルの構造を明らかにし、エントロピー・レート $H(X)$ の収束性や、漸近等分割性についての理論的な解析を進めていく必要がある。また、 W が語頭条件を満たさないが、モデルに制約を加えることで情報源のエントロピー・レートが定義できるような情報源モデルがないか考えたい。

謝辞 著者の1人石田は、本研究を進めるにあたり熱心に議論をしてくださる湘南工科大学 小林学先生に感謝いたします。

文献

- [1] M.Nishiara and H.Morita, "On the AEP of word-valued sources," *IEEE Trans. Inform. Theory*, vol.IT-46, no.3, pp.1116-1120, May 2000.
- [2] Masayuki GOTO, Toshiyasu MATSUSHIMA, and Shigeichi HIRASAWA, "A source model with probability distribution over word set and recurrence time theorem," submitted to *IEICE Trans. Fundamentals. Special Section on Information Theory and Its Applications*, 2003.
- [3] 石田 崇, 後藤 正幸, 松嶋敏泰, 平澤 茂一, "単語単位で系列を出力する情報源の性質について," 第25回情報理論とその応用シンポジウム予稿集 (SITA2002), pp.695-698, 2002.
- [4] 石田 崇, 後藤 正幸, 平澤 茂一, "ブロック単位で系列を出力する情報源に対するベイズ符号と Ziv-Lempel 符号のユニバーサル性について," *電子情報通信学会論文誌*, vol.J84-A, no.9, pp.1167-1178, 2001.
- [5] 石田 崇, 後藤 正幸, 平澤 茂一, "単語単位で系列を出力する情報源に対する LZ78 符号のユニバーサル性について," 第24回情報理論とその応用シンポジウム予稿集 (SITA2001), pp.243-246, 2001.
- [6] 韓 太舜, 情報理論における情報スペクトル的方法, 培風館, 1998.
- [7] 韓 太舜, 小林 欣吾, 情報と符号化の数理, 培風館, 1999.