

6.符号化法(6章)

1

(情報源の)符号化定理と符号化法

情報源符号化定理(シャノンの第1定理)は、符号化の理論的限界を与えているだけであり、具体的な符号化の手法を与えてはいない。すなわち、ある情報源に対して、エントロピーは平均符号長の目標にすぎない。



ここでは、具体的な符号化法を与える。

2

代表的(情報源)符号化法

- シャノン・ファノ符号化(算術符号化)
  - 確率を2進数化して符号化する。
- ハフマン符号化(コンパクト符号化)
  - 最小の平均符号長を持つ符号(コンパクト符号)を実現する符号。符号の木の葉から符号を割り当てていく。

3

P進数と基数変換  
(シャノン・ファノ符号化法の準備)

4

p進数

p種類の記号  $\{0, 1, \dots, p-1\}$  を基に、数表現することができる。小数点以上n桁、小数点以下m桁のp進数  $(q_{n-1}q_{n-2} \dots q_0 \cdot q_{-1}q_{-2} \dots q_{-m})_p$  は以下の値を持つ。ただし、 $q_i \in \{0, 1, \dots, p-1\}$

基数を明示する表記

$$\sum_{i=-m}^n q_i \times p^i$$

整数部 (小数点の左の数、1以上)

$$= q_{n-1}p^{n-1} + q_{n-2}p^{n-2} + \dots + q_0p^0$$

小数部 (小数点の右の数、1未満)

$$+ q_{-1}p^{-1} + q_{-2}p^{-2} + \dots + q_{-m}p^{-m}$$

5

10進数と2進数

$$D = (d_{n-1}d_{n-2} \dots d_1d_0 \cdot d_{-1}d_{-2} \dots d_{-m})_{10}$$

$d_i \in \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$

各桁の数字は、その位の値が何個あるのかを示している。

小数点の位置が重要

基数を明示する表記(通常は10の省略)

$$B = (b_{s-1}b_{s-2} \dots b_1b_0 \cdot b_{-1}b_{-2} \dots b_{-t})_2$$

$b_i \in \{0, 1\}$

2<sup>s-1</sup> 2 2<sup>0</sup> = 1 2<sup>-1</sup> = 1/2 2<sup>-t</sup> = 1/2<sup>t</sup>

6

$$(19)_{10} = 1 \times 10^1 + 9 \times 10^0$$

$$= 10 + 9$$


$$16 + 2 + 1$$

$$= 1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

$$= (10011)_2$$

7

**練習**

次の値を求めよ。

(1)  $(11011.01)_2$       (3)  $(43.21)_5$

(2)  $(2102.2)_3$       (4)  $(72.3)_8$

8

**10進数から2進数への変換1**

入力  $D^+ = (d_{n-1}d_{n-2} \dots d_1d_0)_{10}$       出力  $B^+ = (b_{s-1}b_{s-1} \dots b_1b_0)_2$

**2進数変換アルゴリズム(整数部分)**

[step1]:  $D_0^+ := D^+, i = 0$ とする。      初期設定

[step2]:  $D_i^+ \neq 0$ の間以下を繰り返す;

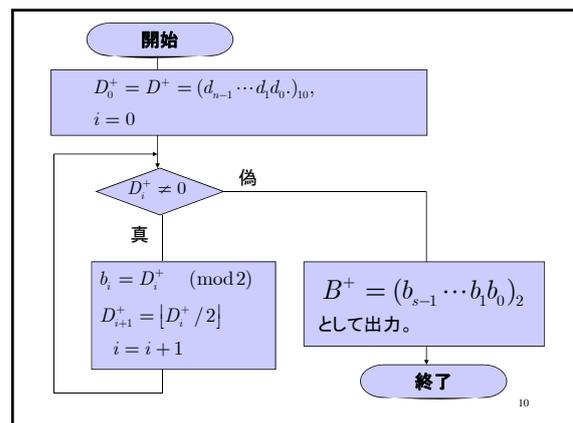
(2-1)  $b_i := D_i^+ \text{ mod } 2$       2で割った余り

(2-2)  $D_{i+1}^+ := \lfloor D_i^+ / 2 \rfloor$       2で割った商(切り捨て)

(2-3)  $i := i + 1$        $i$ を1増加させる。

[step3]:  $B^+ = (b_{s-1} \dots b_1b_0)_2$ を出力して終了する。

9



**例**

$(54)_{10}$  を2進数に変換せよ。

2) 54	2) $D_0^+$
2) 27 ... 0	2) $D_1^+$ ... $b_0$
2) 13 ... 1	2) $D_2^+$ ... $b_1$
2) 6 ... 1	2) $D_{s-1}^+$ ... $b_{s-2}$
2) 3 ... 0	0 ... $b_{s-1}$
2) 1 ... 1	
0 ... 1	

よって、

$$(54)_{10} = (110110)_2$$

11

**練習**

次の10進数を2進数に変換せよ。

(1)  $(35)_{10}$       (2)  $(63)_{10}$

(3)  $(48)_{10}$       (4)  $(41)_{10}$

12



この関係より、各  $i, 0 \leq i \leq s-1$  に対して  
次ように計算できる。

**ループ不変条件**

$$D_i^+ = B_i^+$$

$$= (b_{s-1} \cdots b_i)_2$$

$$= b_{s-1} \times 2^{s-i-1} + \cdots + b_{i+1} \times 2^1 + b_i \times 2^0$$

$$= (b_{s-1} \times 2^{s-i-2} + \cdots + b_{i+1} \times 2^0) \times 2 + b_i$$

$$= B_{i+1}^+ \times 2 + b_i$$

**アルゴリズムの動作**

$$\Leftrightarrow b_i = D_i^+ \pmod{2}$$

QED 19

### 10進数から2進数への変換2

**入力**  $D^- = (0.d_{-1}d_{-2} \cdots d_{-m})_{10}$  **出力**  $B^- = (0.b_{-1}b_{-2} \cdots b_{-t})_2$

**2進数変換アルゴリズム(小数部分)**

[step1]:  $D_{-1}^- := D^-, i = -1$  とする。 **初期設定**

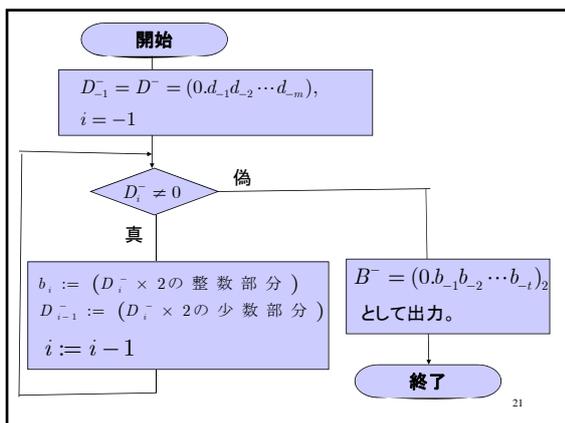
[step2]:  $D_i^- \neq 0$  の間以下を繰り返す;

(2-1)  $b_i := (D_i^- \times 2$  の整数部分)

(2-2)  $D_{i-1}^- := (D_i^- \times 2$  の少数部分)

(2-3)  $i := i - 1$

[step3]:  $B^- = (0.b_{-1}b_{-2} \cdots b_{-t})_2$  を出力して終了する。 20



**例**  $(0.5625)_{10}$  を2進数に変換せよ。

$2 \times 0.5625 = 1.125 \cdots 1 + 0.125$	$2 \times \overline{D_{-1}} = b_{-1} + \overline{D_{-2}}$
$2 \times 0.125 = 0.25 \cdots 0 + 0.25$	$2 \times \overline{D_{-2}} = b_{-2} + \overline{D_{-3}}$
$2 \times 0.25 = 0.5 \cdots 0 + 0.5$	⋮
$2 \times 1 = 1.0 \cdots 1 + 0.0$	$2 \times \overline{D_{-t}} = b_{-t} + 0.0$

よって、  
 $(0.5625)_{10} = (0.1001)_2$

**小数点に近い方から順に求まる。** 22

**練習**

次の10進数を2進数に変換せよ。

(1)  $(0.625)_{10}$                       (2)  $(0.53125)_{10}$

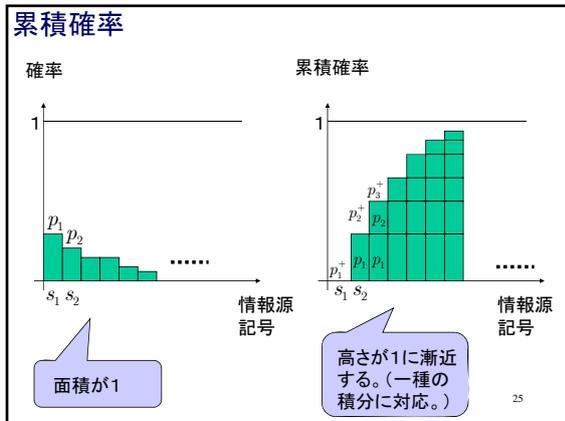
(3)  $(0.3)_{10}$                               (4)  $(0.59375)_{10}$

23

**練習**

小数部分の基数変換アルゴリズムにおけるループ不変条件を示し、正当性を証明せよ。

24



## シャノン・ファノ符号化法

26

### シャノン・ファノ符号化

入力: 情報源 (情報源記号の集合とその発生確率)  
 $S = \{s_1, \dots, s_n\}$   
 $\{p_1, \dots, p_n\}$

出力: 符号 (情報源記号に対応する符号語の集合)  
 $C = \{c_1, c_2, \dots, c_n\}$

ステップ1: 発生確率の大きい順に並べる。  
 (ここでは、添え字でこの順序が見たされるとする。)

ステップ2: 各符号語長を次式で求める。  
 $l_i = \lceil -\log p_i \rceil$

ステップ3: 次のような累積確率  $p_i^+$  を求める **切り上げ**  
 $p_i^+ = 0 (i=1), p_i^+ = \sum_{j=1}^{i-1} p_j (i \geq 2)$

ステップ4: 累積確率  $p_i^+$  を2進数  $(p_i^+)_2$  に変換する。

ステップ5: 2進数  $(p_i^+)_2$  の上位  $l_i$  桁を符号語  $c_i$  とする。

27

### 例

次の無記憶情報源  $S$  に対して、シャノン・ファノ符号を構成する。

$$S = \begin{Bmatrix} a & b & c & d \\ 0.2 & 0.3 & 0.1 & 0.4 \end{Bmatrix}$$

ステップ1 (降順に並び替え)

$$S = \begin{Bmatrix} d & b & a & c \\ 0.4 & 0.3 & 0.2 & 0.1 \end{Bmatrix}$$

$\therefore s_1 = d, s_2 = b, s_3 = a, s_4 = c$

28

### ステップ2 (符号語長の決定)

$d: -\log 0.4 \simeq 1.322$   
 $\therefore l_1 = l_d = \lceil -\log 0.4 \rceil = 2$

$b: -\log 0.3 \simeq 1.737$   
 $\therefore l_2 = l_b = \lceil -\log 0.3 \rceil = 2$

$a: -\log 0.2 \simeq 2.322$   
 $\therefore l_3 = l_a = \lceil -\log 0.2 \rceil = 3$

$c: -\log 0.1 \simeq 3.322$   
 $\therefore l_4 = l_c = \lceil -\log 0.1 \rceil = 4$

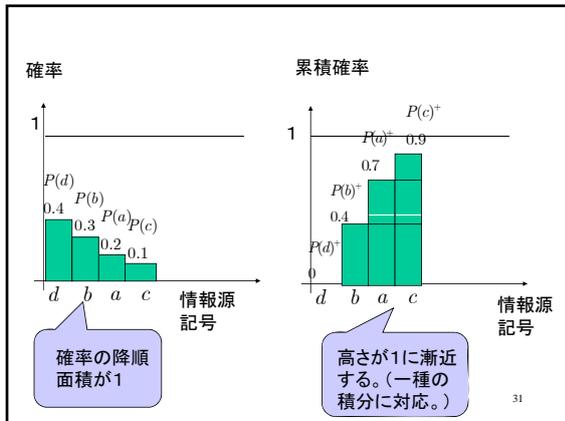
したがって、  
 $L = \{l_1, l_2, l_3, l_4\} = \{2, 2, 3, 4\}$   
 の符号語長を持つ符号を構成する。

29

### ステップ3 (累積確率の計算)

情報源記号	符号語長	確率	累積確率
$d$	$l_1 = 2$	$p_1 = P(d) = 0.4$	$p_1^+ = 0.0$
$b$	$l_2 = 2$	$p_2 = P(b) = 0.3$	$p_2^+ = 0.4$
$a$	$l_3 = 3$	$p_3 = P(a) = 0.2$	$p_3^+ = 0.7$
$c$	$l_4 = 4$	$p_4 = P(c) = 0.1$	$p_4^+ = 0.9$

30



ステップ4(累積確率の2進数化)

$$(p_1^+)_{10} = (0.0)_{10} \simeq (0.00000)_2$$

$$(p_2^+)_{10} = (0.4)_{10} \simeq (0.01100)_2$$

$$(p_3^+)_{10} = (0.7)_{10} \simeq (0.10110)_2$$

$$(p_4^+)_{10} = (0.9)_{10} \simeq (0.11100)_2$$

ステップ5(符号の割り当て)

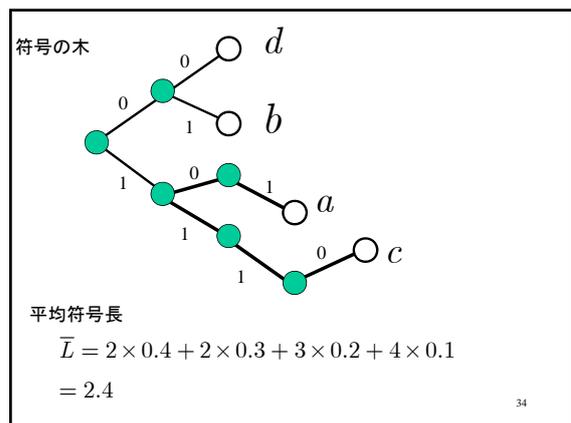
$$d \rightarrow 00$$

$$b \rightarrow 01$$

$$a \rightarrow 101$$

$$c \rightarrow 1110$$

$$\therefore C = \{c_1, c_2, c_3, c_4\}$$

$$= \{00, 01, 101, 1110\}$$


練習

次の情報源ジャンをシャノンファノの符号化法に従って符号化せよ。

$$\text{ジャン} = \left\{ \begin{array}{l} \text{グー} , \text{チョキ} , \text{パー} \\ 0.35 , \quad 0.25 , \quad 0.4 \end{array} \right\}$$

また、得られた符号に対する符号の木を示し、平均符号長、効率を求めよ。

シャノンファノ符号化法の平均符号長

$1 \leq i \leq n$  に対して、

$$l_i = \lceil -\log p_i \rceil$$

であるが、これは次式を満たす。

$$-\log p_i \leq l_i < -\log p_i + 1$$

したがって、平均符号長は次式で求められる。

$$-\sum_{i=1}^n p_i \log p_i \leq \sum_{i=1}^n p_i l_i < -\sum_{i=1}^n p_i \log p_i + \sum_{i=1}^n p_i$$

$$\therefore H(S) \leq \bar{L} < H(S) + 1$$

どこかで見た式

### シャノンファノ符号の非特異性

シャノンファノ符号化法で構成された符号は非特異符号である。

証明

符号語長の選び方より、次式が成り立つ。

$$-\log p_i \leq l_i < -\log p_i + 1$$

$$\therefore 2^{-l_i} \leq p_i < 2^{-l_i+1}$$

37

一方、

記号  $s_{i-1}$  の生成確率

$$p_i^+ - p_{i-1}^+ = p_{i-1}$$

に注意する。

累積確率の  $i$  番目

累積確率の  $i-1$  番目

よって、2進数表現して、次式が成り立つ。

$$(p_i^+)_2 - (p_{i-1}^+)_2 = (p_{i-1})_2$$

$$\geq 2^{\log p_{i-1}}$$

第  $-l_{i-1}$  桁で必ず異なることを意味する。

$$\geq 2^{-l_{i-1}}$$

Q.E.D 38

### 縮退情報源 (ハフマン符号化法の準備)

39

### 縮退情報源

情報源  $S$  に対して、 $S$  の2つ以上の情報源記号  $s_i, s_j \in S$  を一つにまとめた情報源を元の情報源の縮退情報源という。すなわち、

$$S = \left\{ \begin{matrix} s_1 & , & \dots & , & s_i & , & s_j & , & \dots & , & s_n \\ p_2 & , & \dots & , & p_i & , & p_j & , & \dots & , & p_n \end{matrix} \right\}$$

$$S^- = \left\{ \begin{matrix} s_1 & , & \dots & , & *s_k & , & \dots & , & s_n \\ p_1 & , & \dots & , & *p_k & , & \dots & , & p_n \end{matrix} \right\}$$

ここで、 $*s_k = \{s_i, s_j\}$      $*p_k = p_i + p_j$

$$|S^-| = |S| - 1$$

40

### 例

次の情報源の縮退情報源をいくつか示せ。

$$S = \left\{ \begin{matrix} a & , & b & , & c & , & d \\ 0.2 & , & 0.3 & , & 0.1 & , & 0.4 \end{matrix} \right\}$$

解)

$$S_1^- = \left\{ \begin{matrix} A & , & c & , & d \\ 0.5 & , & 0.1 & , & 0.4 \end{matrix} \right\}, A = \{a, b\}$$

$$S_2^- = \left\{ \begin{matrix} a & , & B & , & d \\ 0.2 & , & 0.4 & , & 0.4 \end{matrix} \right\}, B = \{b, c\}$$

$$S_3^- = \left\{ \begin{matrix} C & , & b & , & d \\ 0.3 & , & 0.3 & , & 0.4 \end{matrix} \right\}, C = \{a, c\}$$

41

### 練習

次の情報源に対して、2記号を1記号に縮退して得られる情報源をすべて示せ。

$$S = \left\{ \begin{matrix} a & , & b & , & c & , & d \\ 0.15 & , & 0.25 & , & 0.05 & , & 0.55 \end{matrix} \right\}$$

42

## ハフマン符号化法

43

### ハフマン符号化

入力: 情報源 (情報源記号の集合とその発生確率)

$$S = \begin{Bmatrix} s_1, \dots, s_n \\ p_1, \dots, p_n \end{Bmatrix}$$

出力: 符号 (情報源記号に対応する符号語の集合)

$$C = \{c_1, c_2, \dots, c_n\}$$

ステップ1:  $k = 0$  とし、 $S_0 = S = \begin{Bmatrix} s_1^0, \dots, s_n^0 \\ p_1^0, \dots, p_n^0 \end{Bmatrix}$  とする。

ここで、 $s_i^0 = s_i, p_i^0 = p_i \quad (1 \leq i \leq n)$

ステップ2: 発生確率の大きい順に並べる。  
(添え字の順序をこの順序とする。)

44

### ステップ3: 縮退情報源

$$S_k = \begin{Bmatrix} s_1^k, \dots, s_{n-k}^k \\ p_1^k, \dots, p_{n-k}^k \end{Bmatrix}$$

に対して、確率の小さい2つの情報源記号  $s_{n-k}^k, s_{n-k-1}^k$  に対して、対応する符号の末尾に0と1を割り当てる。さらに、 $s_{n-k}^k, s_{n-k-1}^k$  を縮退して、新たな縮退情報源

$$S_{k+1} = \begin{Bmatrix} s_1^{k+1}, \dots, s_{n-k-2}^{k+1}, *s_{n-(k+1)}^{k+1} \\ p_1^{k+1}, \dots, p_{n-k-2}^{k+1}, *p_{n-(k+1)}^{k+1} \end{Bmatrix}$$

を作成する。

ここで、 $s_i^{k+1} = s_i^k, p_i^{k+1} = p_i^k \quad (1 \leq i \leq n-k-2)$

$$*s_{n-(k+1)}^{k+1} = \{s_{n-k-1}^k, s_{n-k}^k\},$$

$$*p_{n-(k+1)}^{k+1} = p_{n-k-1}^k + p_{n-k}^k \quad (i = n-k-1)$$

45

ステップ4:  $k < n-1$  である限り、 $k = k+1$  としてステップ2に戻る。

46

### ハフマン符号化例

次の符号のハフマン符号を与える。

$$S = \begin{Bmatrix} a, b, c, d \\ 0.2, 0.3, 0.1, 0.4 \end{Bmatrix}$$

符号の木を葉から構成していくと分かりやすい。

全ての点線で、縮退情報源の確率が降順になっていること。

$S_0 = \{d, b, a, c\} \quad S_1 = \{d, b, A\} \quad S_2 = \{B, d\}$

47

前のスライドの符号の木より、

$$d \rightarrow 0$$

$$b \rightarrow 11$$

$$a \rightarrow 101$$

$$c \rightarrow 100$$

平均符号長  $\bar{L}$  は、次のように計算される。

$$\bar{L} = 0.4 \times 1 + 0.3 \times 2 + 0.2 \times 3 + 0.1 \times 3$$

$$= 1.9$$

48

練習

次の符号をハフマン符号化し、  
符号の木、平均符号長、効率を求めよ。

$$S = \left\{ \begin{array}{cccc} a & , & b & , & c & , & d \\ 0.15 & , & 0.25 & , & 0.05 & , & 0.55 \end{array} \right\}$$

49

コンパクト符号

(定義)コンパクト符号

情報源に対して、符号語長が最短となる符号を  
コンパクト符号という。

コンパクト符号であっても、効率が1  
になるとは限らない。効率が1なら明  
らかにコンパクト符号である。

50

ハフマン符号のコンパクト性

ハフマン符号は、コンパクト符号である。

ハフマン符号化で得られる符号は、1  
通りとは限らない。しかし、どのような  
ハフマン符号もコンパクトとなる。

この証明のために、次の補題を示す。

51

補題

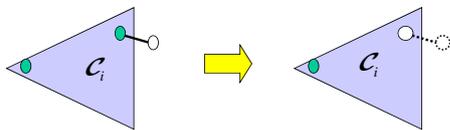
情報源  $S_i$  の最小の生成確率を持つ2記  
号を縮退して得られる情報源を  $S_{i+1}$  とする。  
情報源  $S_i$  のハフマン符号を  $C_i$  とし、情  
報源  $S_{i+1}$  のハフマン符号を  $C_{i+1}$  とする。こ  
のとき、 $C_{i+1}$  がコンパクト符号ならば、 $C_i$   
もコンパクト符号である。

証明

$C_i$  の平均符号長を  $\bar{L}_i$  とし、  
 $C_{i+1}$  の平均符号長を  $\bar{L}_{i+1}$  とする。

52

まず、符号の木の経路長が最長の葉は2つある。  
(もし、経路長が最長の葉が1つだけなら、さらに短くできる。)

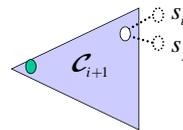


この2つの葉に発生確率最小の2記号  $s_i, s_j \in S_i$  を割り当  
てる。(もし、発生確率が最小以外の記号を割り当てるとより短い  
平均符号長が得られる。)

$$\bar{L}_i = \sum_{s_k \in S_i} p_k l_k$$

$$\leq \sum_{s_k \in S_i} p_k l_k - p_i l_i - p_j l_j + p'_i l_i + p'_j l_j$$

53



$$\bar{L}_{i+1} = \bar{L}_i - p_i l_i - p_j l_j + (p_i + p_j)(l_i - 1)$$

$$= \bar{L}_i - p_i - p_j$$

今、 $S_{i+1}$  に対する任意の符号の平均符号長を  $\bar{L}'_{i+1}$  と表す。

このとき、 $\bar{L}_{i+1}$  のコンパクト性より、以下が成り立つ。

$$\bar{L}_{i+1} \leq \bar{L}'_{i+1}$$

54

同じ値を両辺から引いているだけ。

$$\therefore \overline{L}_{i+1} - p_i - p_j \leq \overline{L}'_{i+1} - p_i - p_j$$

$$\therefore \overline{L}_i \leq \overline{L}'_i$$

情報源  $S_{i+1}$  の平均符号長と等しい。

したがって、 $C_i$  もコンパクト符号になる。

Q.E.D <sup>55</sup>

(ハフマン符号のコンパクト性の証明)

基礎  
 $i = n - 2$  のとき。  
 このとき、縮退情報源の記号の数は、  
 $|S_i| = |S_{n-2}| = 2$   
 であり、ハフマン符号は明らかにコンパクト符号である。

帰納  
 $n - 2 \geq i > 0$  のすべての  $i$  に対して、ハフマン符号が情報源のコンパクト符号だと仮定する。(帰納法の仮定)  
 先の補題より、  
 $C_i$  がコンパクト符号ならば、 $C_{i-1}$  もコンパクト符号である。  
 したがって、 $S_0$  をハフマン符号化した  $C_0$  もコンパクト符号である。

Q.E.D <sup>56</sup>