

コンパクト符号

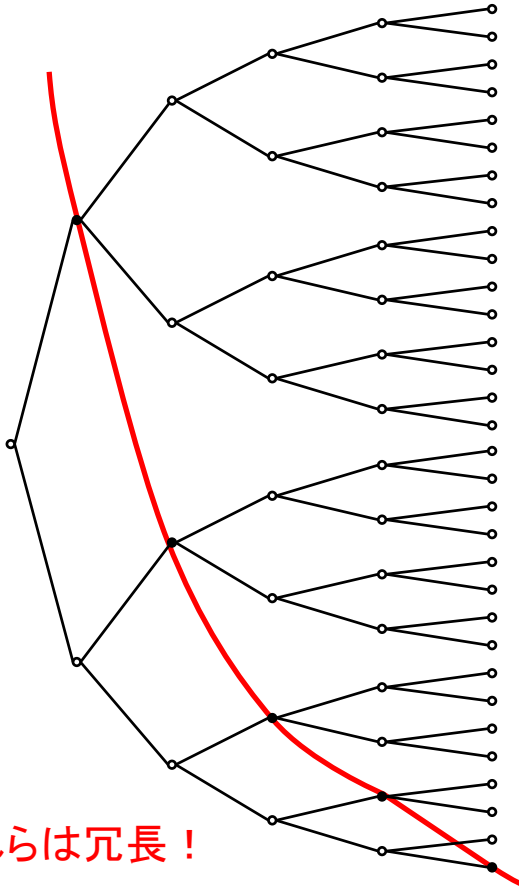
コンパクト符号

- 与えられた情報源 S から発生する情報源記号に一つずつ符号語を割り当てる符号化によってできる一意復号可能な符号のうち、平均符号長が最小となる符号をコンパクト符号という.
- マクミラン不等式とクラフト不等式が同形だから、任意の一意復号可能なコンパクト符号に対して、それと同じ符号語長バッグを持つ瞬時符号が存在する.
- 概念的には、与えられた情報源 S に対するコンパクト符号を見つけるためには、すべての瞬時符号を枚挙し、そのなかの平均符号長最小のものを選べばよい.
- 情報源 S に対するコンパクト符号が複数存在することもある.

コンパクト符号

例： 情報源記号{A, B, C, D, E}に対する2元コンパクト符号？

瞬時符号 P

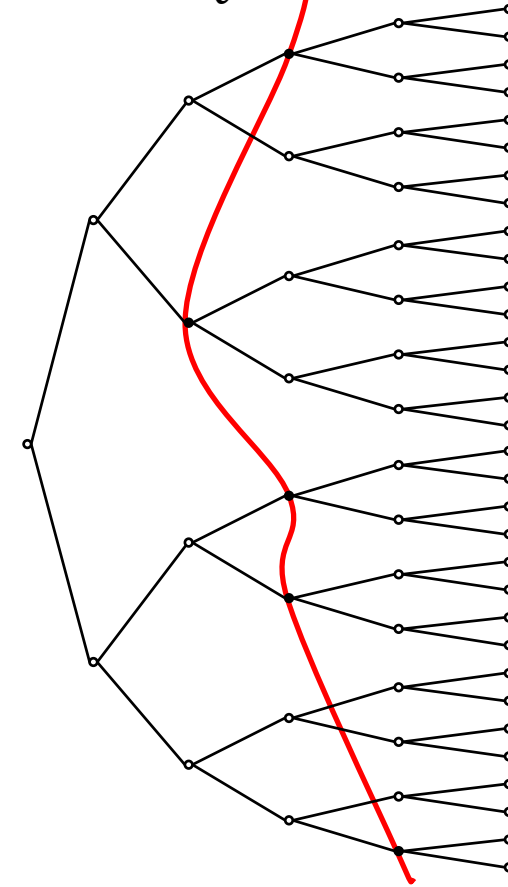


これらは冗長！

符号長バッグ: $\{1, 2, 3, 4, 5\}$

$$\rightarrow \{1, 2, 3, 4, 4\}$$

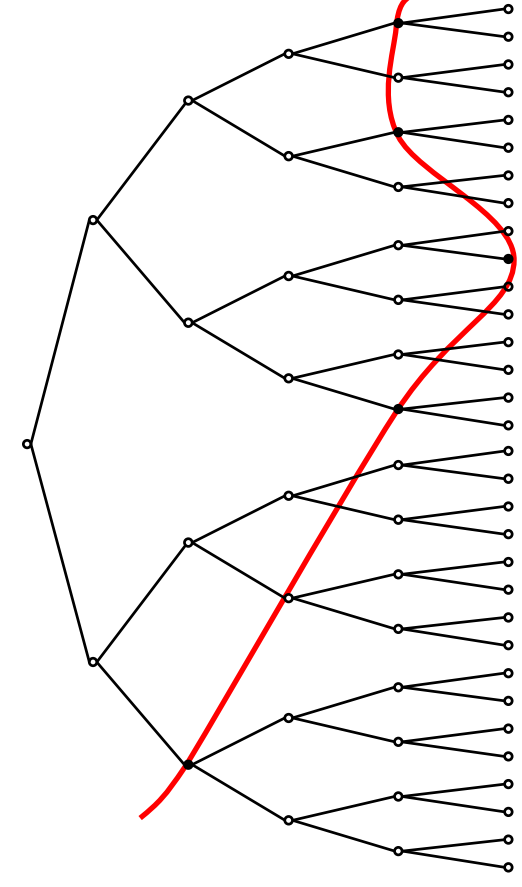
瞬時符号 Q



符号長バッグ: $\{3, 2, 3, 3, 4\}$

$$\rightarrow \{2, 2, 3, 3, 2\}$$

瞬時符号 R



符号長バッグ: $\{4, 4, 5, 4, 2\}$

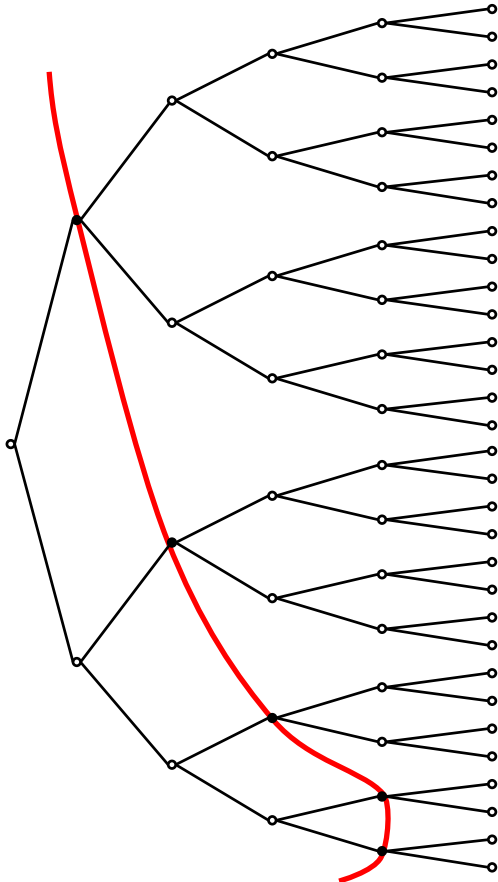
$$\rightarrow \{3, 3, 3, 3, 1\}$$

...

コンパクト符号

例： 情報源記号{A, B, C, D, E}に対する2元コンパクト符号？

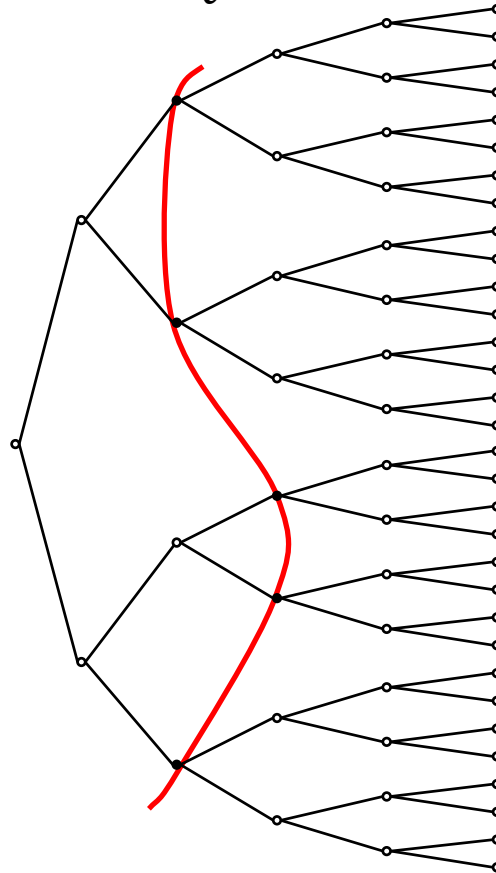
瞬時符号P



どれがいいか？

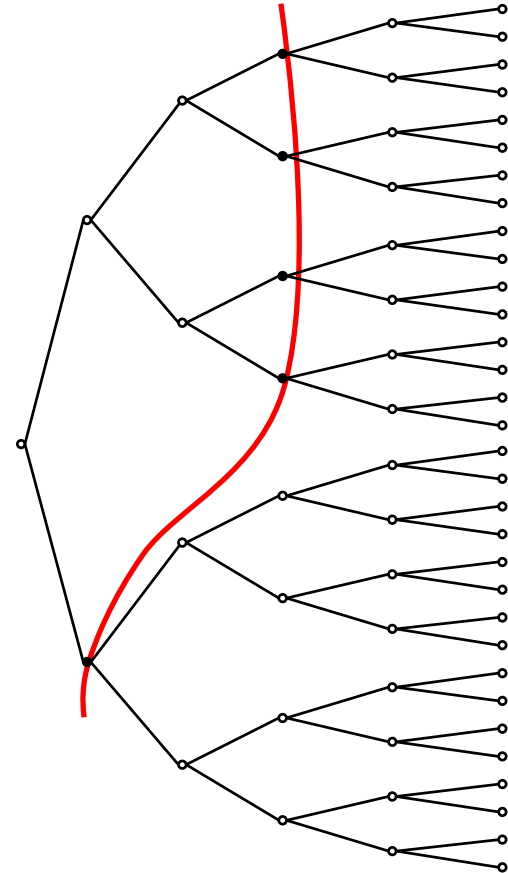
{1, 2, 3, 4, 4}

瞬時符号Q



{2, 2, 3, 3, 2}

瞬時符号R



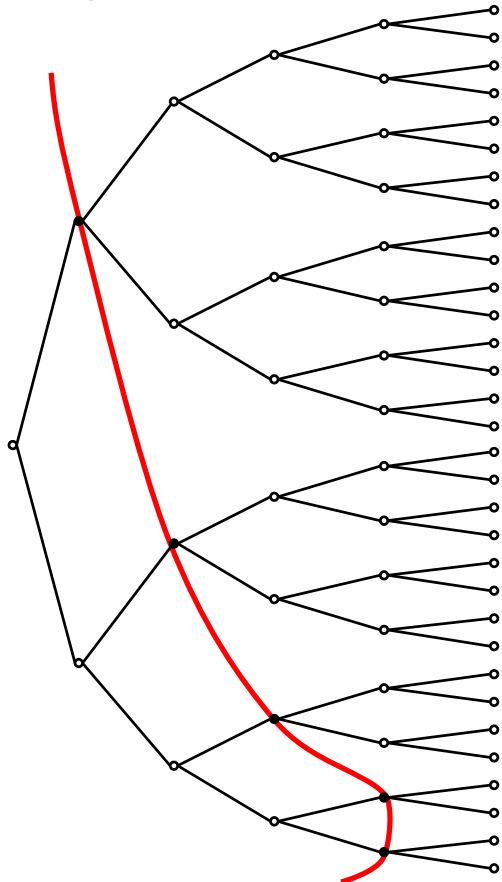
{3, 3, 3, 3, 1}

...

コンパクト符号

例： 情報源記号発生確率が $\langle A: 0.2, B: 0.2, C: 0.2, D: 0.2, E: 0.2 \rangle$ の場合

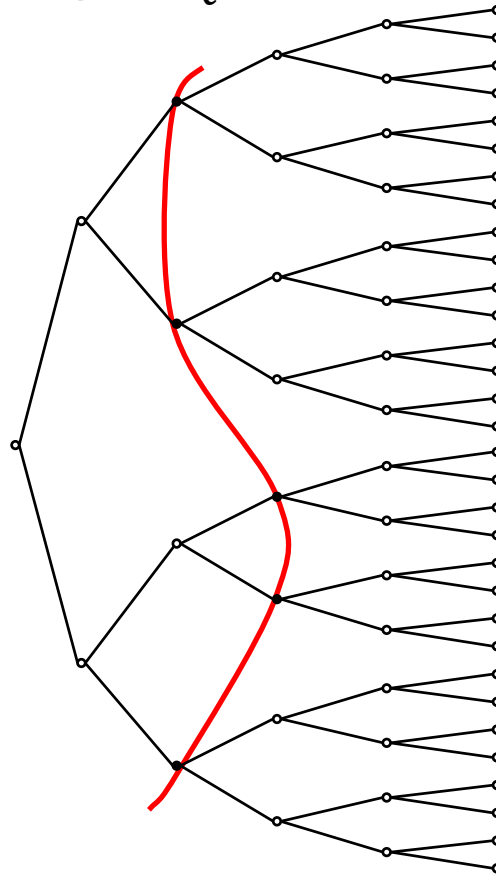
瞬時符号 P



符号長バッグ: $\{1, 2, 3, 4, 4\}$

→ 平均符号長2.8

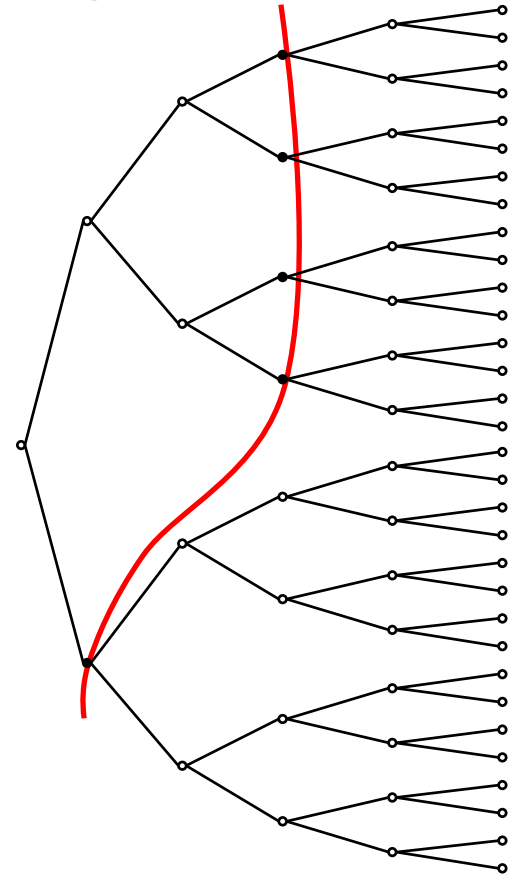
瞬時符号 Q



符号長バッグ: $\{2, 2, 3, 3, 2\}$

→ 平均符号長2.4

瞬時符号 R



符号長バッグ: $\{3, 3, 3, 3, 1\}$

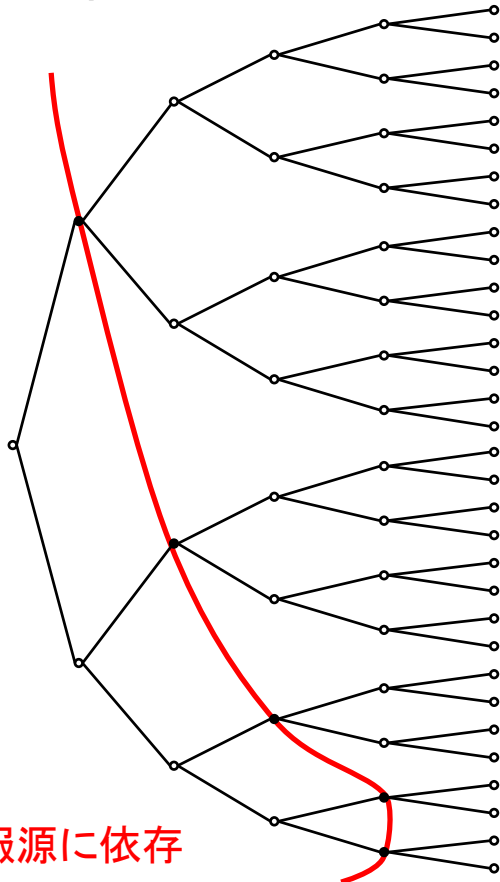
→ 平均符号長2.6

...

コンパクト符号

例： 情報源記号発生確率が〈A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03〉の場合

瞬時符号P

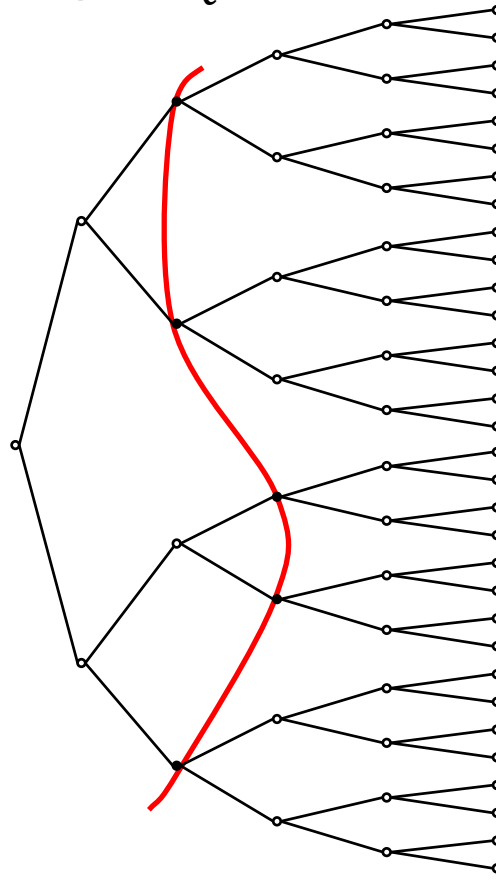


情報源に依存

符号長バッグ: {1, 2, 3, 4, 4}

→ 平均符号長1.7

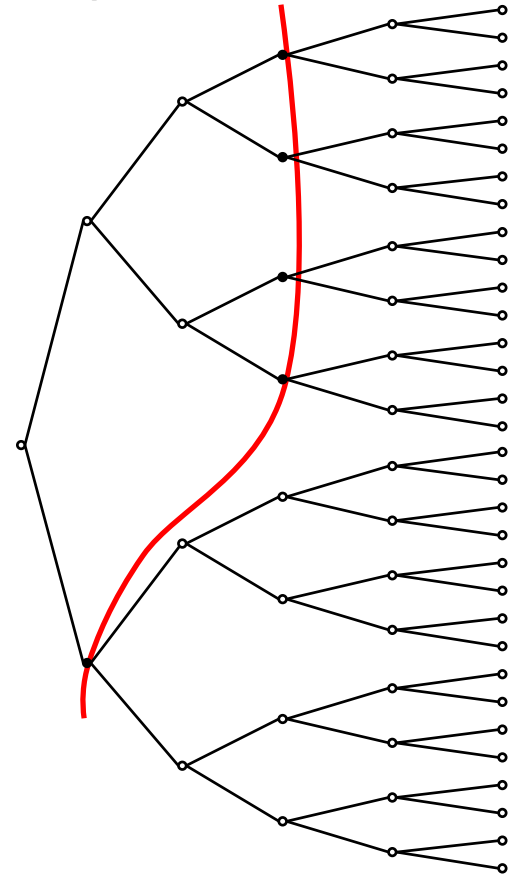
瞬時符号Q



符号長バッグ: {2, 2, 3, 3, 2}

→ 平均符号長2.17

瞬時符号R



符号長バッグ: {3, 3, 3, 3, 1}

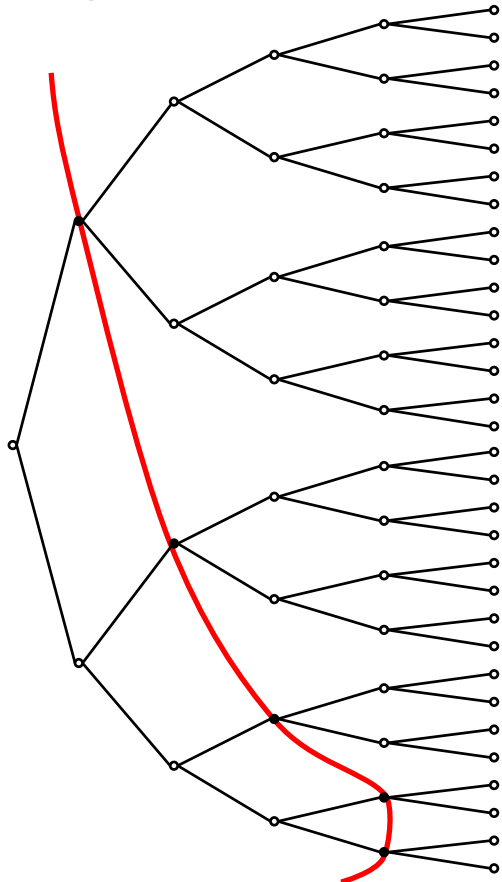
→ 平均符号長2.94

...

コンパクト符号

問題： 情報源記号発生確率 $\langle A:?, B:?, C:?, D:?, E:?\rangle$ がどのような場合？

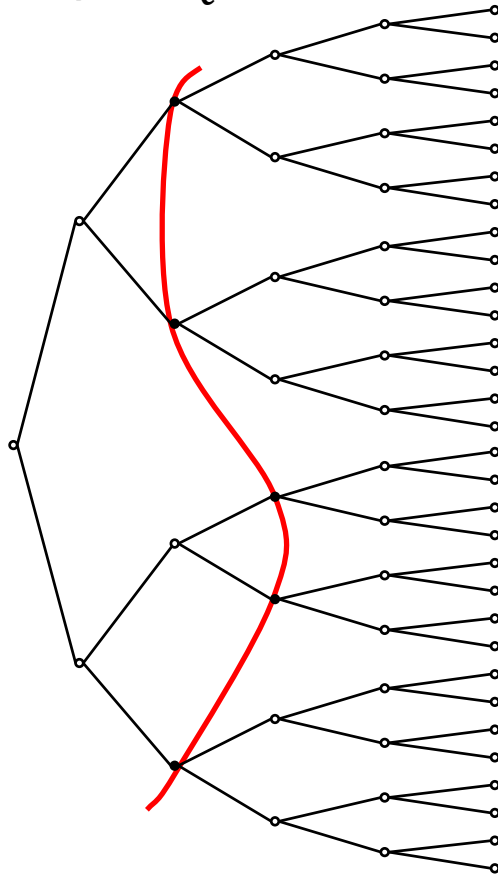
瞬時符号 P



符号長バッグ: $\{1, 2, 3, 4, 4\}$

→ 平均符号長?

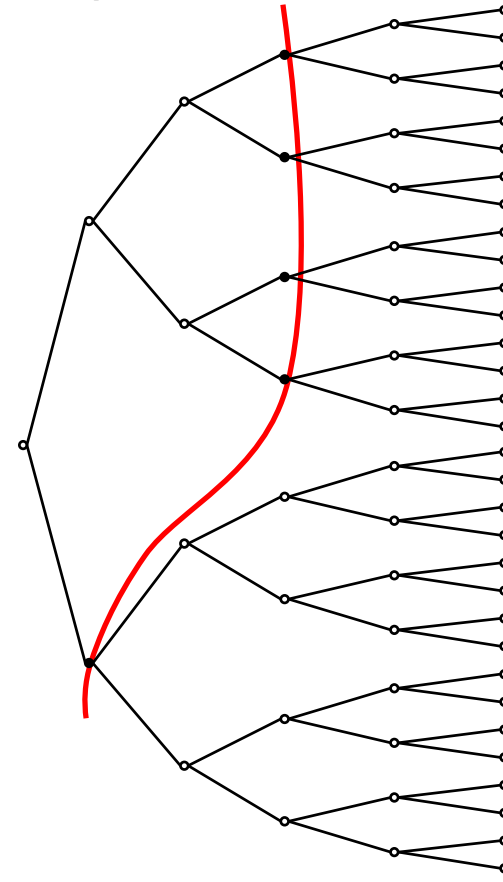
瞬時符号 Q



符号長バッグ: $\{2, 2, 3, 3, 2\}$

→ 平均符号長?

瞬時符号 R



符号長バッグ: $\{3, 3, 3, 3, 1\}$

→ 平均符号長?

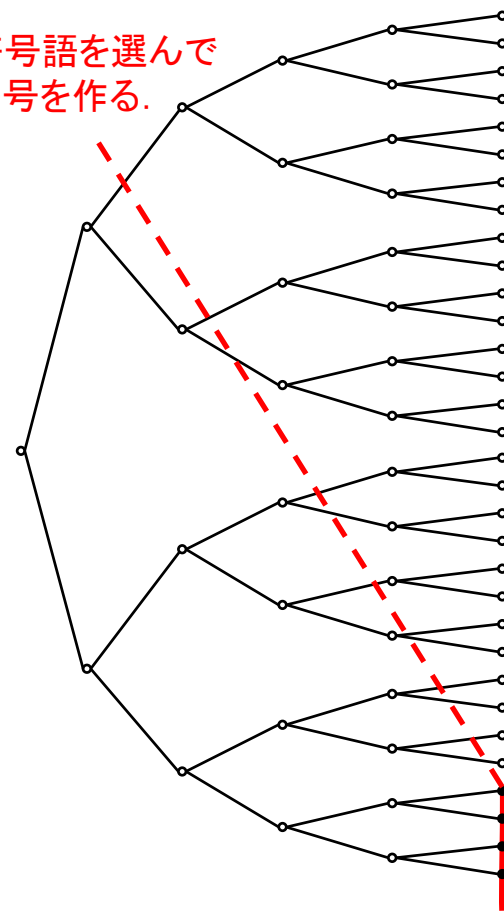
• • •

練習問題3-1

$\{c(a_1), c(a_2), c(a_3), c(a_4), c(a_5), c(a_6), c(a_7), c(a_8)\}$ が得られ,

になったという. ただし, $c(a_i)$ は情報源記号 a_i に対する符号語, $|c(a_i)|$ は符号語 $c(a_i)$ の長さを表す.

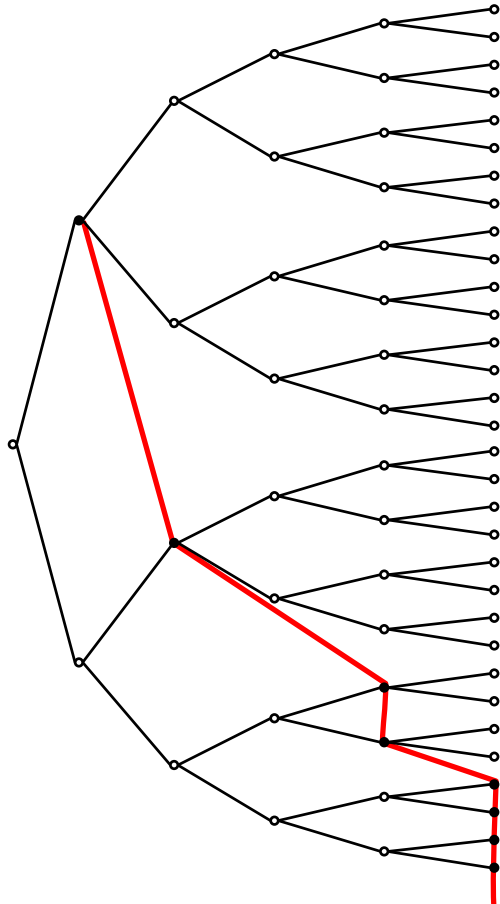
あと4つの符号語を選んで
コンパクト符号を作る.



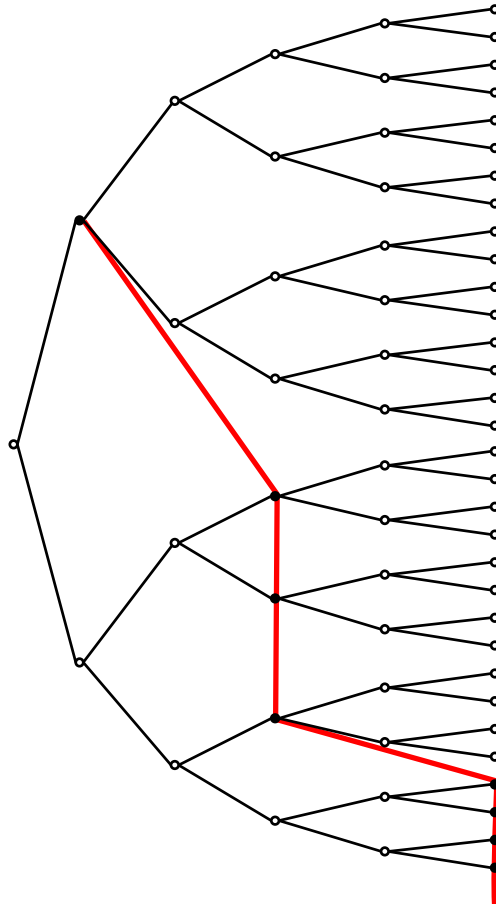
符号長バッグ: $\{?, ?, ?, ?, 5, 5, 5, 5\}$

コンパクト符号

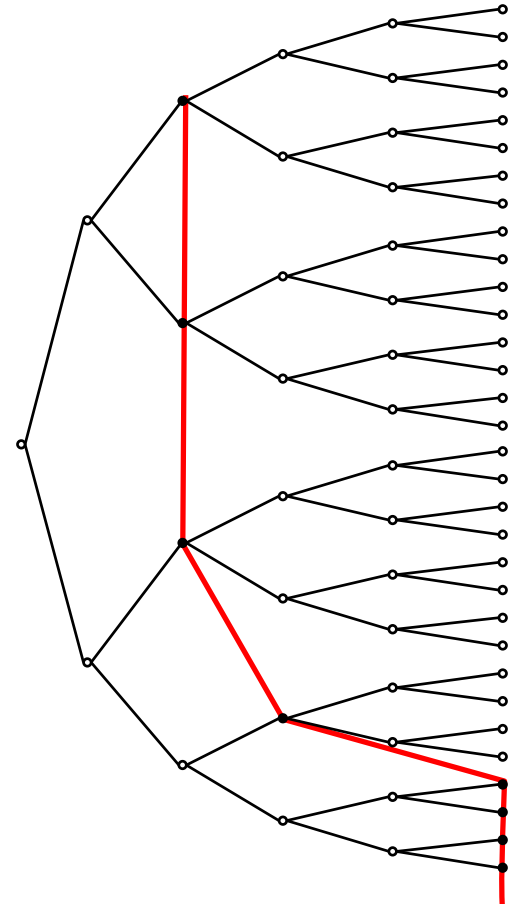
練習問題3-1



符号長バッグ: {1, 2, 4, 4, 5, 5, 5, 5}



符号長バッグ: {1, 3, 3, 3, 5, 5, 5, 5}



符号長バッグ: {2, 2, 2, 3, 5, 5, 5, 5}

ハフマン符号化:コンパクト符号を構成する

例: 情報源記号発生確率が $\langle A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03 \rangle$ の場合

情報源 $S_5: \langle A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03 \rangle$ のコンパクト符号構成法

- ① (すぐに答えが分からないので), 情報源 S_5 の発生確率の最も小さな2つの情報源記号 $D: 0.07, E: 0.03$ に着目し, それらを一つにまとめて $D' \leftarrow D + E$ とした情報源

$S_4: \langle A: 0.6, B: 0.2, C: 0.1, D': 0.1 \rangle$

に対するコンパクト符号構成を試みる.

- ② 前項と同じ: $C' \leftarrow C + D'$ として, $S_3: \langle A: 0.6, B: 0.2, C': 0.2 \rangle$ に対してトライ
③ 前項と同じ: $B' \leftarrow B + C'$ として, $S_2: \langle A: 0.6, B': 0.4 \rangle$ に対してトライ
④ S_2 に対するコンパクト符号 $\langle A \leftarrow 0, B' \leftarrow 1 \rangle$
⑤ $B' \leftarrow B + C'$ であったので, $\langle A \leftarrow 0, B \leftarrow 10, C' \leftarrow 11 \rangle$
⑥ $C' \leftarrow C + D'$ であったので, $\langle A \leftarrow 0, B \leftarrow 10, C \leftarrow 110, D' \leftarrow 111 \rangle$
⑦ $D' \leftarrow D + E$ であったので, $\langle A \leftarrow 0, B \leftarrow 10, C \leftarrow 110, D \leftarrow 1110, E \leftarrow 1111 \rangle$

なぜこれでいいのか?

符号長バッグ: $\{1, 2, 3, 4, 4\}$

→ 平均符号長1.7

情報源 S_n の発生確率の最も小さな2つの情報源記号を a_n, a_{n-1} とする.
 $a'_{n-1} \leftarrow a_n + a_{n-1}$ によってできる情報源 S_{n-1} のコンパクト符号を c_{n-1} とする.
 C_{n-1} の a'_{n-1} に対する符号語 c'_{n-1} とする. C_{n-1} の $a'_{n-1} \leftarrow c'_{n-1}$ の代わりに
 $a_{n-1} \leftarrow c'_{n-1}0, a_n \leftarrow c'_{n-1}1$ としてできる符号 C_n はコンパクト符号である.

コンパクト符号

例： 情報源記号発生確率が〈A:0.6, B:0.2, C:0.1, D:0.07, E:0.03〉の場合

情報源 S_n の発生確率の最も小さな2つの情報源記号を a_n, a_{n-1} とする。
 $a'_{n-1} \leftarrow a_n + a_{n-1}$ によってできる情報源 S_{n-1} のコンパクト符号を C_{n-1} とする。
 C_{n-1} の a'_{n-1} に対する符号語 c'_{n-1} とする。 C_{n-1} の $a'_{n-1} \leftarrow c'_{n-1}$ の代わりに
 $a_{n-1} \leftarrow c'_{n-1}0, a_n \leftarrow c'_{n-1}1$ としてできる符号 C_n はコンパクト符号である。

【性質1 (S_5 の場合)】

S_5 に対する任意のコンパクト符号 C_5 において、発生確率の最も小さな2つの情報源記号D, Eには最も長い符号語が割り当てられる。

【証明】

仮にCに対応づける符号語の長さを s , Eに対応づける符号の長さを t ,
 $s > t$ とすると,

$$\text{平均符号長 } L = \dots + s \times 0.1 + \dots + t \times 0.03$$

ここで、Cに対する符号語とEに対する符号語を交換すると

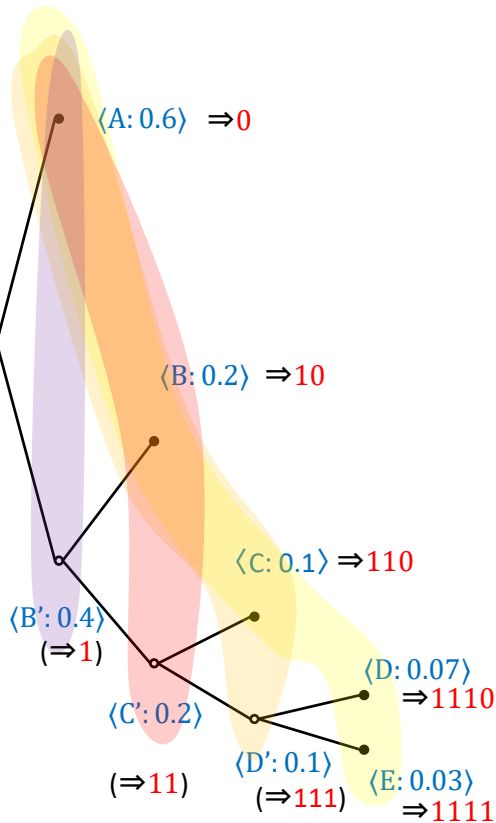
$$\text{平均符号長 } L' = \dots + t \times 0.1 + \dots + s \times 0.03$$

次のように $L' < L$ となり、 S_5 がコンパクト符号であるという前提に矛盾する。

$$L' - L = (t - s) \times (0.1 - 0.03) < 0$$

【性質1'】

S_5 に対して、D, Eに対する符号語が符号の木において兄弟ノードとなっているコンパクト符号 C'_5 が存在する。定義により C'_5 との平均符号長は等しいが、 C'_5 は C_5 と同じであるとは限らない。(証明は後述の通り)



符号長バッグ: {1, 2, 3, 4, 4}

→ 平均符号長1.7

コンパクト符号

例： 情報源記号発生確率が〈A:0.6, B:0.2, C:0.1, D:0.07, E:0.03〉の場合

情報源 S_n の発生確率の最も小さな2つの情報源記号を a_n, a_{n-1} とする。
 $a'_{n-1} \leftarrow a_n + a_{n-1}$ によってできる情報源 S_{n-1} のコンパクト符号を C_{n-1} とする。
 C_{n-1} の a'_{n-1} に対する符号語 c'_{n-1} とする。 C_{n-1} の $a'_{n-1} \leftarrow c'_{n-1}$ の代わりに
 $a_{n-1} \leftarrow c'_{n-1}0, a_n \leftarrow c'_{n-1}1$ としてできる符号 C_n はコンパクト符号である。

【性質2 (S_5 の場合)】

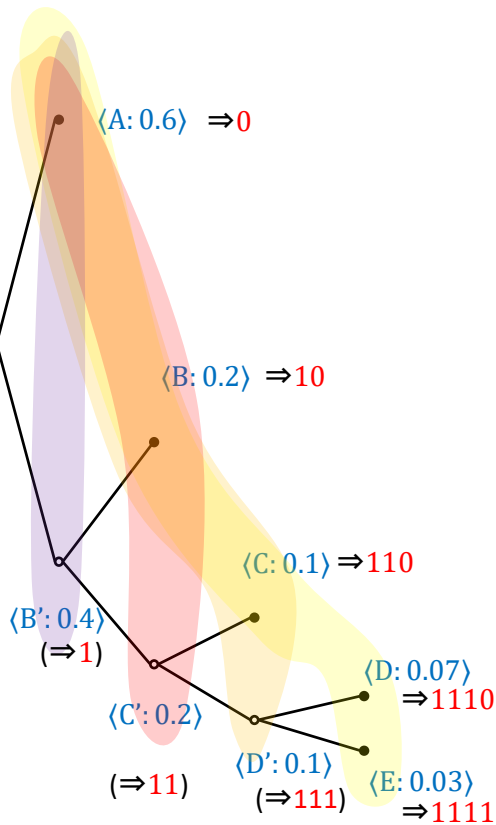
情報源 S_5 の発生確率の最も小さな2つの情報源記号D:0.07, E:0.03に着目し、それらを一つにまとめて $D' \leftarrow D + E$ とした情報源

S_4 : 〈A:0.6, B:0.2, C:0.1, D':0.1〉

に対するコンパクト符号を C_4 とする。 C_4 における D' に対する符号語 c_4 とする。 C_4 において $D' \leftarrow c_4$ の代わりに $D \leftarrow c_40, D \leftarrow c_41$ としてできる符号 C_5 はコンパクト符号である。

【証明骨子】

C_4 はコンパクト符号であるが、 C_5 はコンパクト符号でないと仮定する。 S_5 に対するコンパクト符号 C'_5 が存在し、 C'_5 の平均符号長は C_5 の平均符号長より短い。 性質1'により、 S_5 に対して、D, Eに対する符号語が符号の木において兄弟ノードとなっているコンパクト符号 C''_5 が存在し、その平均符号長は C'_5 の平均符号長に等しい。 C''_5 に対する符号木においてDとEを子とするノードを x とする。 C''_5 において、DとEに対する符号語割り当ての代わりに D' に符号語 x を割り当てる符号 C'' は、 S_4 に対する瞬時符号であり、その平均符号長は C_4 より短い。 これは、 C_4 がコンパクト符号であることに矛盾。



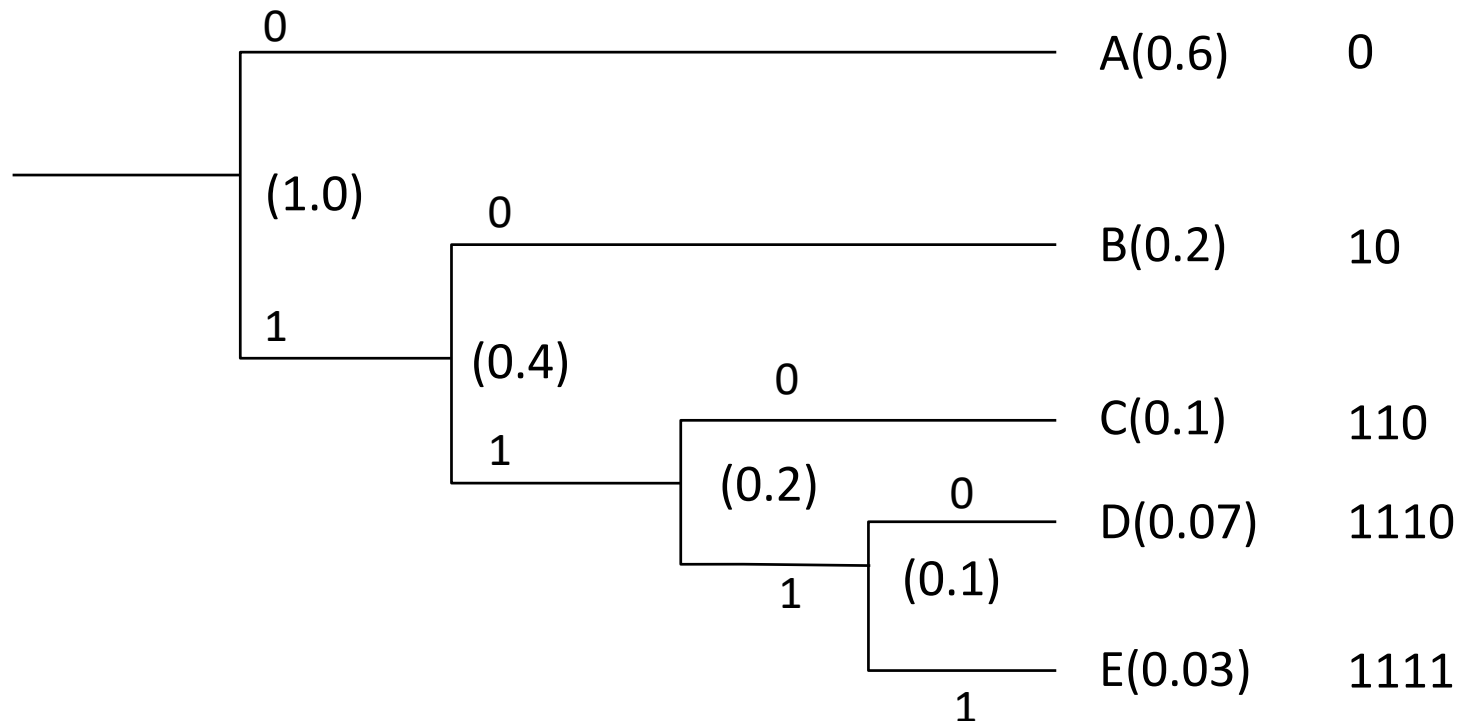
符号長バッグ: {1, 2, 3, 4, 4}

→ 平均符号長1.7

ハフマン符号化(作業のみ)

- ハフマン符号化はコンパクト符号の一つの構成法である.

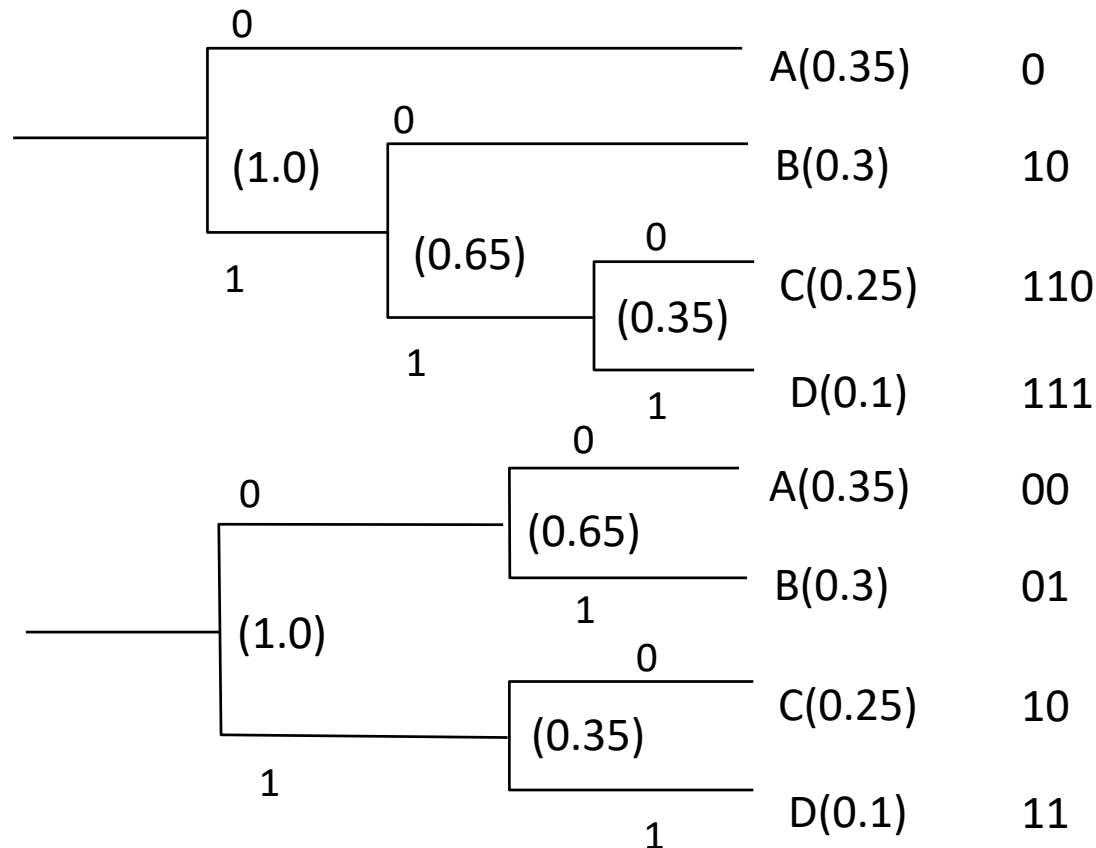
例: 情報源記号発生確率を(A:0.6, B:0.2, C:0.1, D:0.07, E:0.03)とすれば, ハフマン符号化は次のように行われる.



ハフマン符号

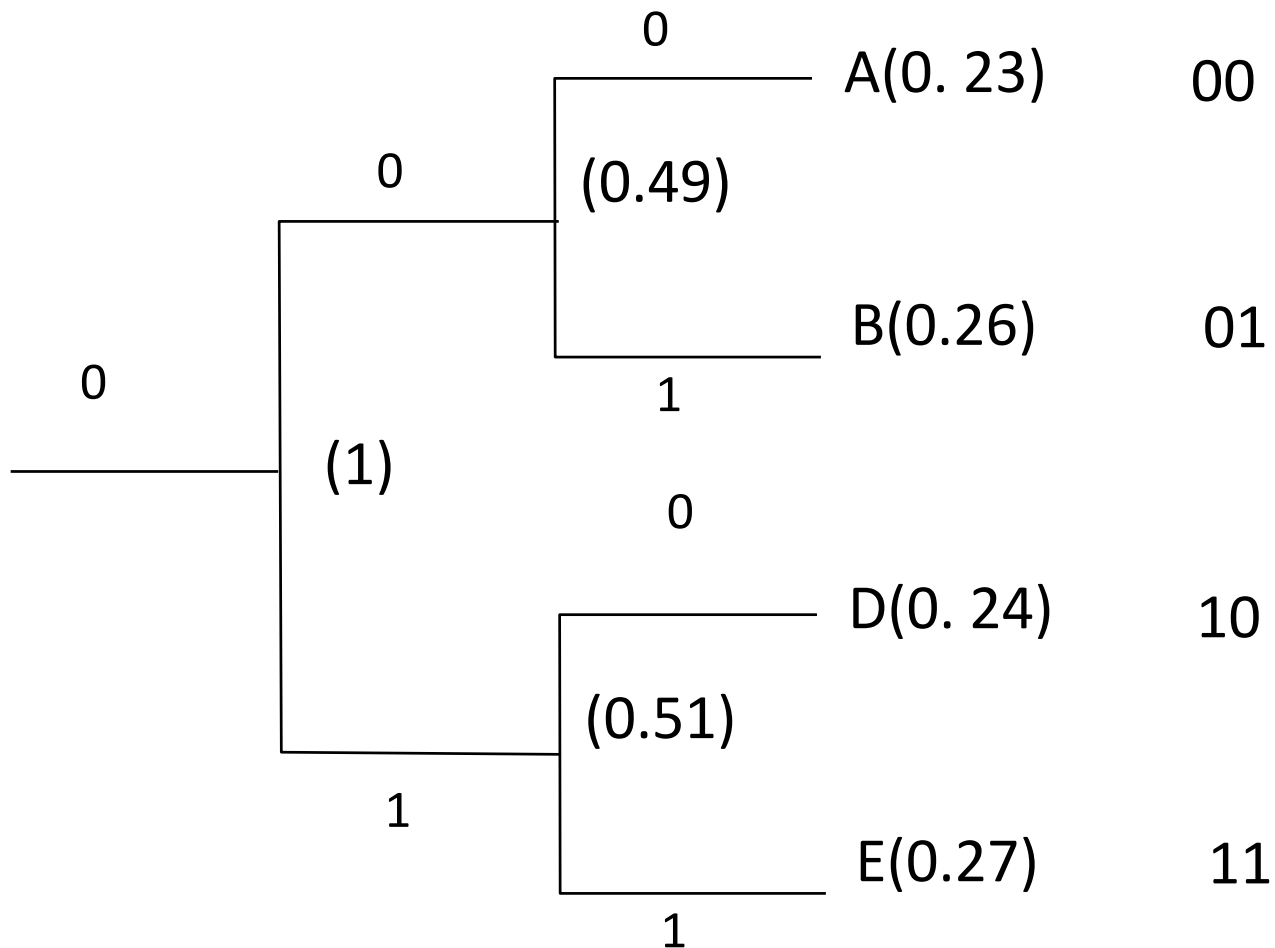
複数のハフマン符号化が可能なケース

情報源記号発生確率が〈A:0.35, B:0.3, C:0.25, D:0.1〉となっているとき



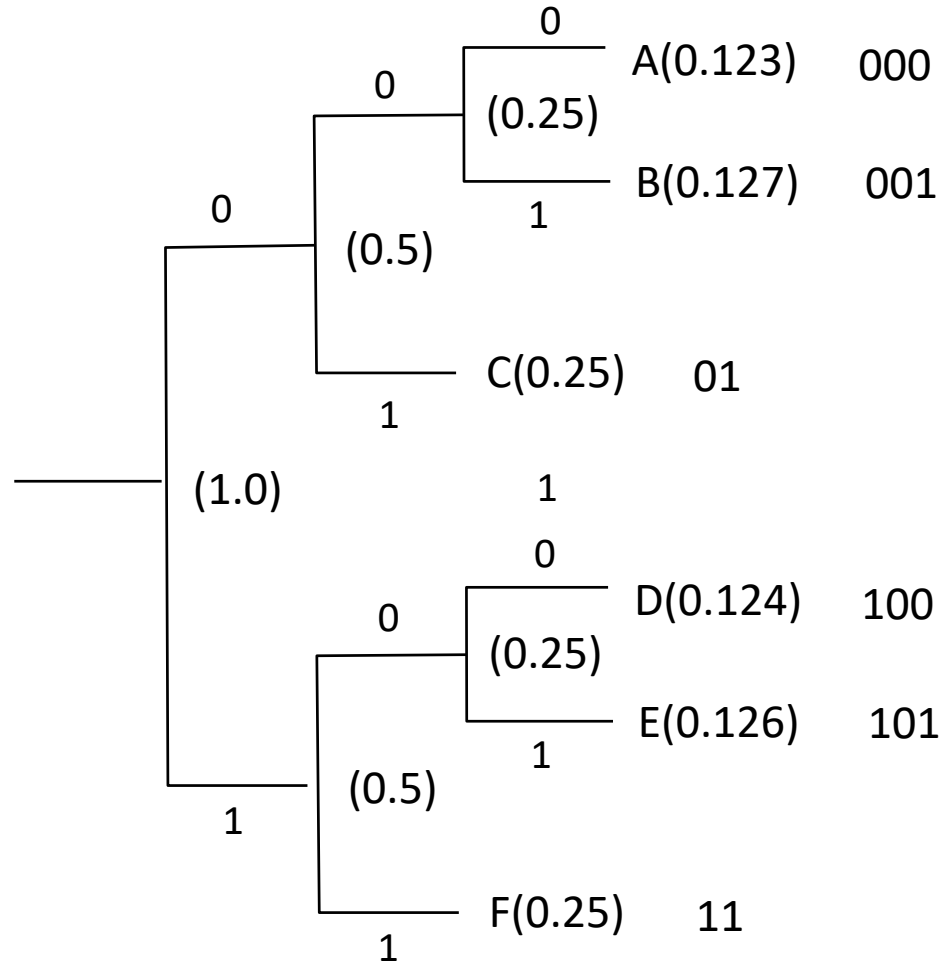
ハフマン符号

ハフマン符号化によって導き出されないコンパクト符号 その1



ハフマン符号

ハフマン符号化によって導き出されないコンパクト符号 その2



コンパクト符号

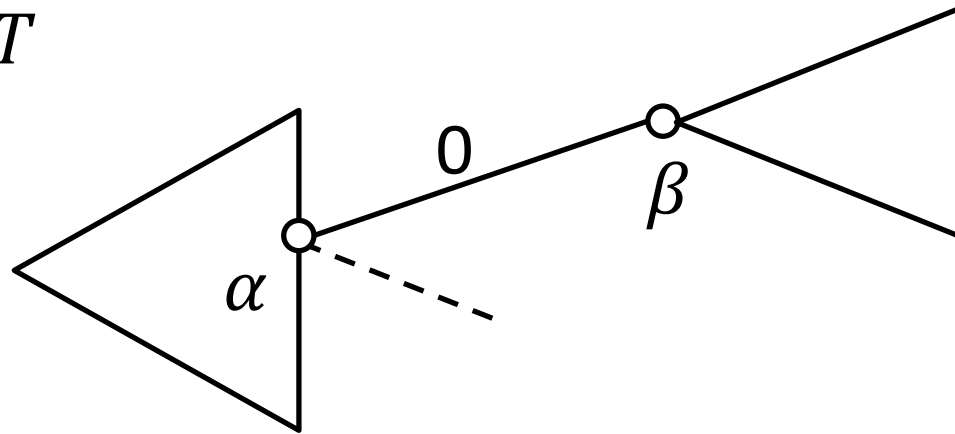
【補助定理1】（瞬時符号 C がコンパクトであるための必要条件）

2個以上の情報源記号をもつ情報源 S が与えられたとする. S に対する瞬時符号 C がコンパクトであるならば, C に対する符号の木 T において,

- (1) 葉以外の節点は必ず2個の子節点を持つ.
- (2) 情報源記号 α, β の出現確率をそれぞれ p_α, p_β とする. $p_\alpha < p_\beta$ ならば情報源記号 α, β に割り付けられる情報源記号をそれぞれ, c_α, c_β とすると, $|c_\alpha| \geq |c_\beta|$ である.

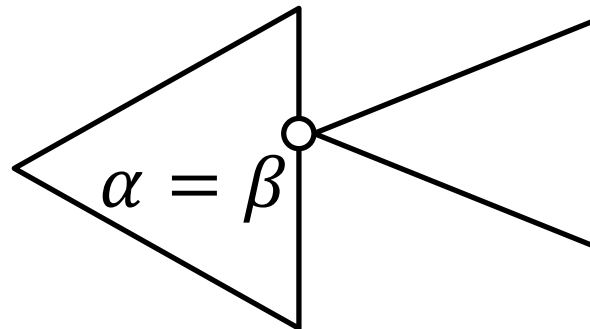
コンパクト符号

符号の木 T

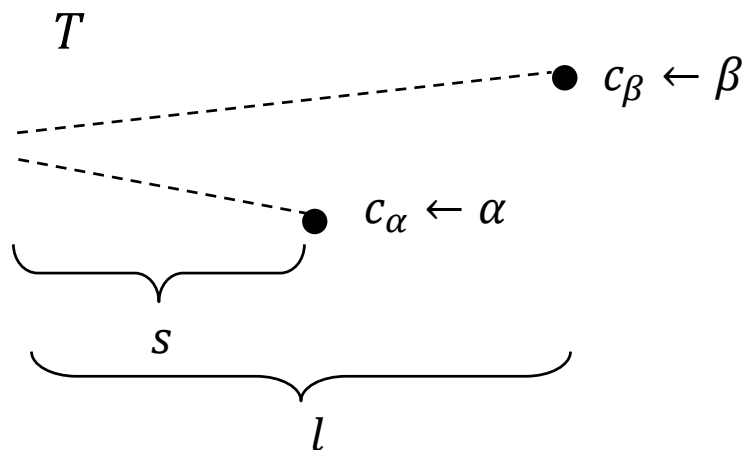


ここで、ノード α は分岐0に対応する子節点 β しか持っていないとしたら、下のようになるほうが平均符号長が短くなることは明らか.

符号の木 T'



コンパクト符号



コンパクト符号に対する符号 C の木 T において
 $p_\alpha < p_\beta$ なる出現確率を持つ情報源記号 α, β に割
 り付けられる符号語 c_α, c_β の長さ, つまり, 深さを
 $s = |c_\alpha|, l = |c_\beta|$ とすると, $s \geq l$ となる.

なぜならば, $s < l$ と仮定すると, T において情報
 源記号 α, β への符号語の割り付け方だけを入れ替
 えた符号の木 T' に対応する符号 C' の平均符号長

$$L' = \dots + sp_\beta + lp_\alpha$$

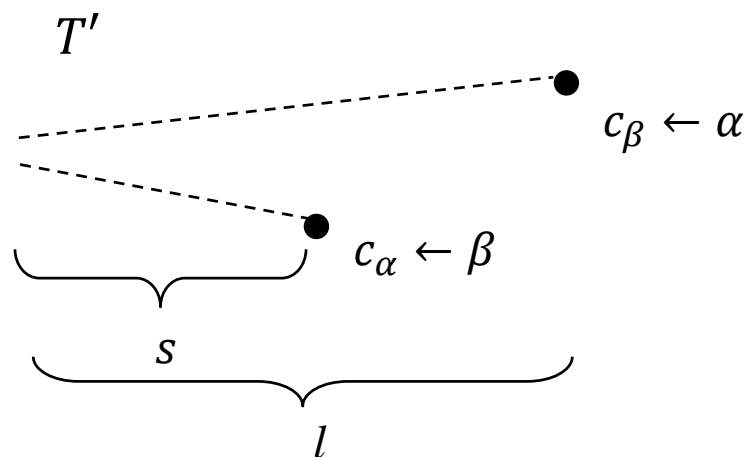
を

$$L = \dots + sp_\alpha + lp_\beta$$

と比較すると, $s < l, p_\alpha < p_\beta$ から

$$\begin{aligned} L' - L &= (sp_\beta + lp_\alpha) - (sp_\alpha + lp_\beta) \\ &= (l - s)(p_\alpha - p_\beta) < 0 \end{aligned}$$

つまり, $L' < L$ となって, C がコンパクトであるとい
 う前提に反するからである.

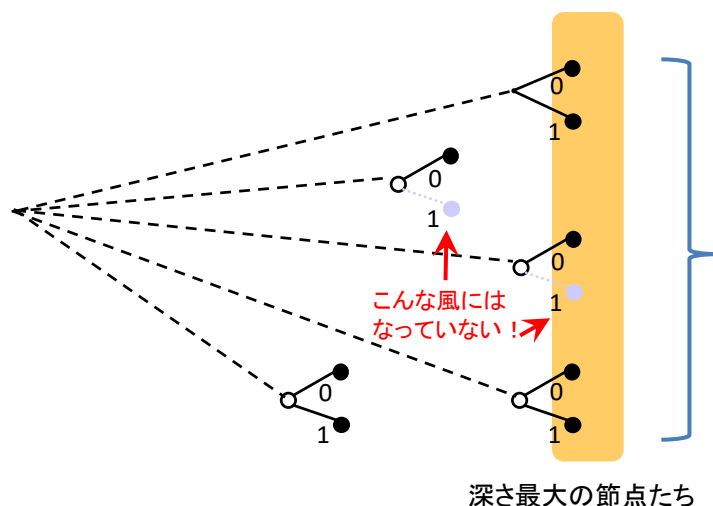


コンパクト符号

【補助定理1】（瞬時符号 C がコンパクトであるための必要条件）

2個以上の情報源記号をもつ情報源 S が与えられたとする. S に対する瞬時符号 C がコンパクトであるならば, C に対する符号の木 T において,

- (1) 深さ d が最大の葉（最高次の葉）は少なくとも2枚あり, そのうちの一つは出現確率の最も小さい情報源記号に割り付けられている.
- (2) 深さ d の葉に割り付けられている情報源記号のなかには, S の情報源記号から α を除いた情報源記号のなかで最も小さな出現確率をもつものがある.

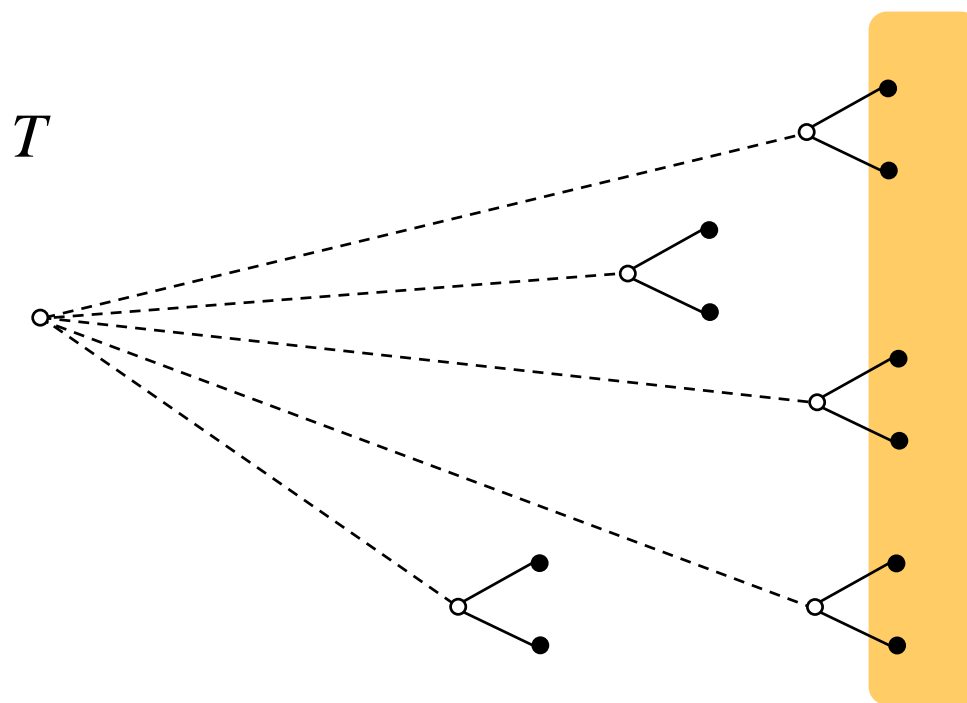


このなかのどれかは S の情報源記号のなかで出現確率の最小のもの α に割り付けられている.

それ以外の葉のどれかは, S の情報源記号から α を除いた情報源記号のなかで最も小さな出現確率をもつものに割り付けられている.

コンパクト符号

3個以上の情報源記号をもつ情報源 S の情報源記号を $A = \{a_1, \dots, a_n\}$, 各情報源記号 a_i の出現確率を p_i とする. また, すべての $1 \leq j \leq n-2$ に対して $p_j \geq p_{n-1} \geq p_n$ とする. S のコンパクトな瞬時符号 T において, a_{n-1} と a_n はどの節点に対応づけられているか?

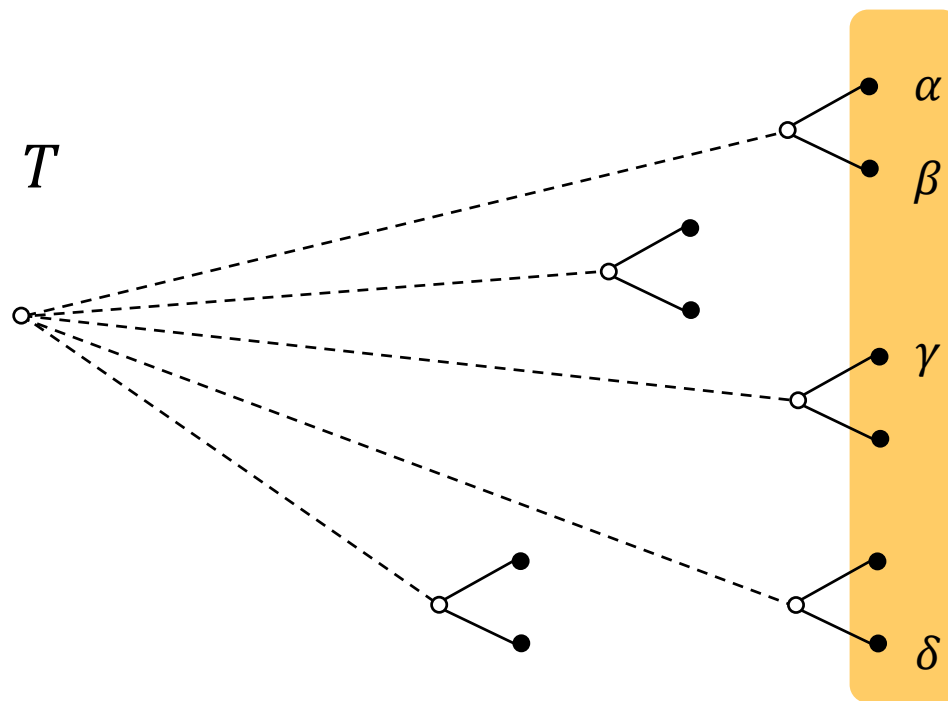


最高次の節点たち

註: 2元符号と仮定

コンパクト符号

典型的には, a_{n-1} と a_n は最高次の節点(例えば, 下図 α, β)に割り付けられる. しかし, $1 \leq \{i, j\} \leq n-2$ に対して $p_i = p_j = p_{n-1} = p_n$ となっているときは, 下図節点 α や β には a_i や a_j が割り付けられ, a_{n-1} と a_n は隣り合っていないかもしれない(例えば, γ と δ)に割り付けられることもある.

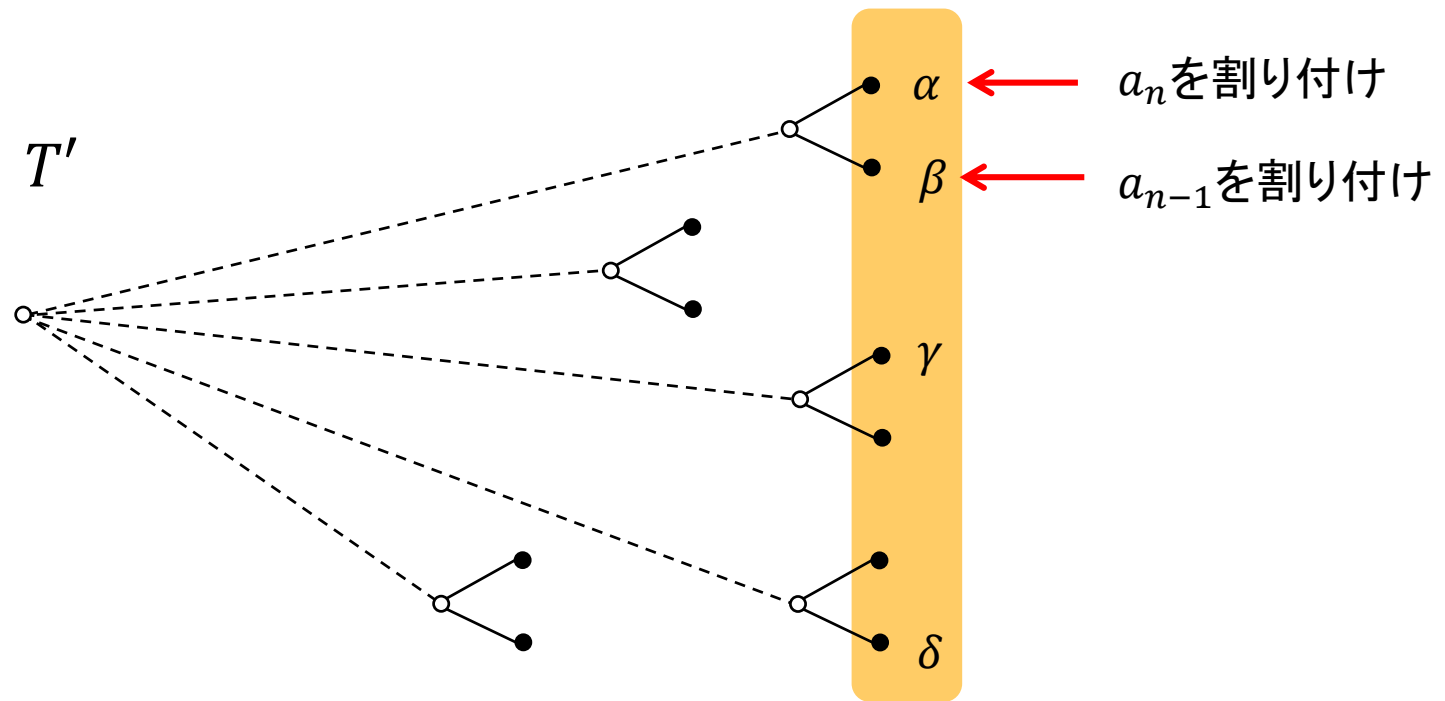


最高次の節点たち

註: 2元符号と仮定

コンパクト符号

いずれの場合でも、平均符号長を変えないように、節点への符号語の割り付けを変えることにより、下図の T' に対応するコンパクト符号を構成できる。



最高次の節点たち 註: 2元符号と仮定

コンパクト符号

【補助定理2】 （コンパクトな瞬時符号 C を構成するための十分条件）

3個以上の情報源記号をもつ情報源 S の情報源記号を
 $A = \{a_1, \dots, a_n\}$, 各情報源記号 a_i の出現確率を p_i とする. また, すべての
 $1 \leq j \leq n-2$ に対して $p_j \geq p_{n-1} \geq p_n$ とする.

このとき, a_n と a_{n-1} を一つの記号 b_{n-1} に統合して, 情報源記号の集まり
 $A' = \{a_1, \dots, a_{n-2}, b_{n-1}\}$, 各情報源記号の出現確率を $p_1, \dots, p_{n-1} + p_n$ と
する情報源 S' のコンパクトな瞬時符号 C' が得られたとする.

すると, C' において, 情報源記号 b_{n-1} に割り付けられた符号語 c_b のかわりに,
 $a_{n-1} \leftarrow c_b 0, a_n \leftarrow c_b 1$ という割り付けを加えた符号 C もまたコンパクトな
瞬時符号である.

コンパクト符号

【補助定理2】の使い方

情報源S: 情報源記号発生確率 $\langle A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03 \rangle$

情報源 S_1 : 情報源記号発生確率 $\langle A: 0.6, B: 0.2, C: 0.1, D: 0.1 \rangle$

情報源 S_2 : 情報源記号発生確率 $\langle A: 0.6, B: 0.2, C: 0.2 \rangle$

情報源 S_3 : 情報源記号発生確率 $\langle A: 0.6, B: 0.4 \rangle$

$\langle A \leftarrow 0, B \leftarrow 1 \rangle$ は情報源 S_3 のコンパクト符号

$\langle A \leftarrow 0, B \leftarrow 10, C \leftarrow 11 \rangle$ は, 情報源 S_2 のコンパクト符号

$\langle A \leftarrow 0, B \leftarrow 10, C \leftarrow 110, D \leftarrow 111 \rangle$ は, 情報源 S_1 のコンパクト符号

$\langle A \leftarrow 0, B \leftarrow 10, C \leftarrow 110, D \leftarrow 1110, E \leftarrow 1111 \rangle$ は, 情報源Sのコンパクト符号

コンパクト符号

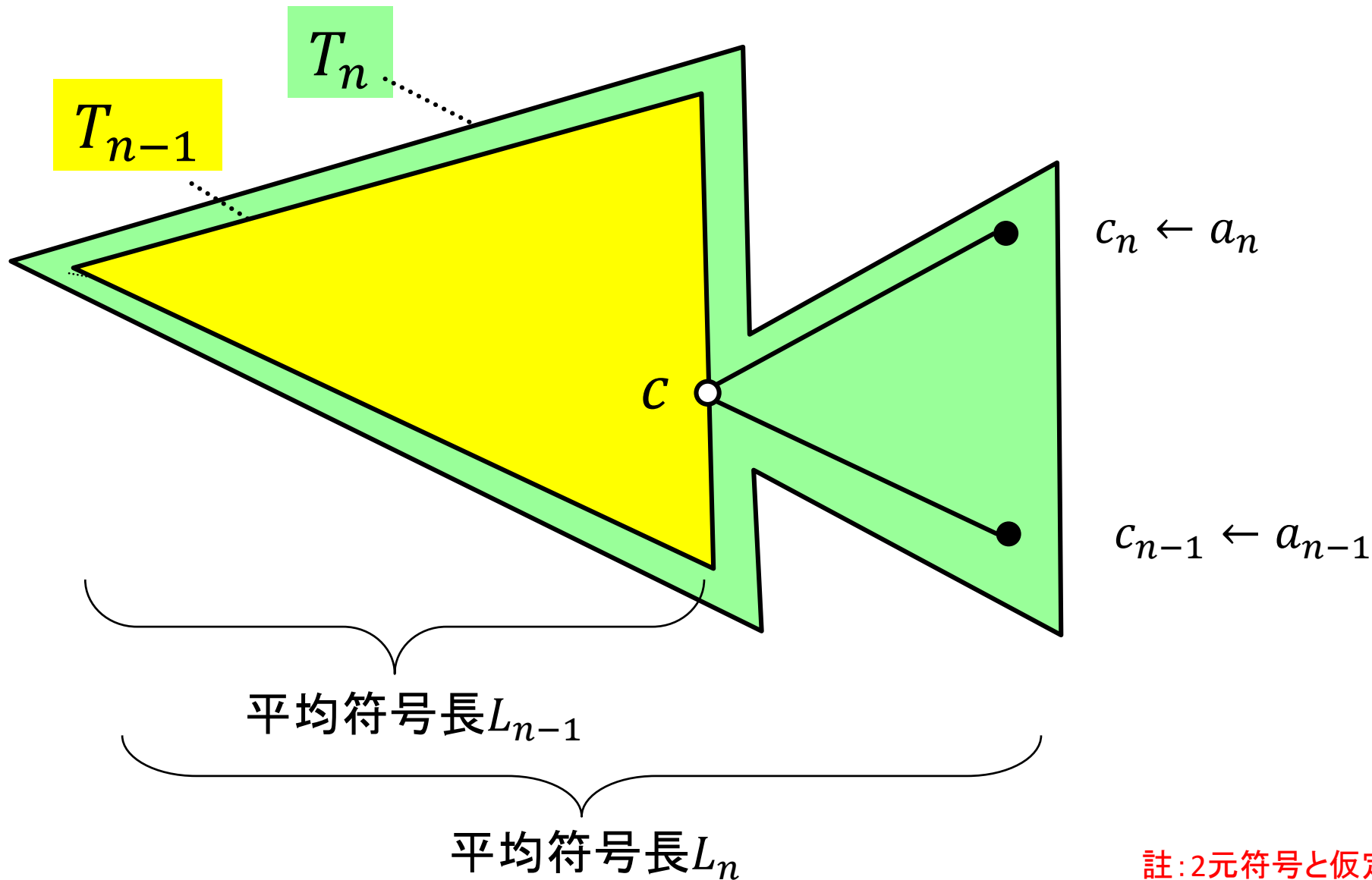
【補助定理2】の証明骨子

情報源 S_n : 情報源記号発生確率 $\langle a_1:p_1, \dots, a_{n-1}:p_{n-1}, a_n:p_n \rangle$

T_{n-1} がコンパクトであるにも関わらず, T_n がコンパクトでないと仮定する. すると, $A = \{a_1, \dots, a_n\}$ に対するコンパクト符号 C'_n が存在し, その平均符号長 L'_n は T_n の平均符号長 L_n より小さいことになる. 補助定理1に関わる議論から, C'_n と平均符号長の等しいコンパクト符号 C''_n が存在し, C''_n において, 最も深い節点には出現確率 p_n と p_{n-1} をもつ符号語が対応づけられている. C''_n において, a_n と a_{n-1} に対する符号語割り付けの代わりに, 出現確率 $p_n + p_{n-1}$ の情報源記号 b_{n-1} に対して符号語 c を割り付ける符号 C'_{n-1} を構成する. すると, C'_{n-1} は, その作り方から C_{n-1} と同じ情報源に対する瞬時符号であり, C'_{n-1} の平均符号長 $L'_{n-1} = L'_n - p_n - p_{n-1}$ が T_{n-1} の平均符号長 $L_{n-1} = L_n - p_n - p_{n-1}$ より短いことになり, T_{n-1} がコンパクトであるという前提に矛盾する.

註: 2元符号と仮定

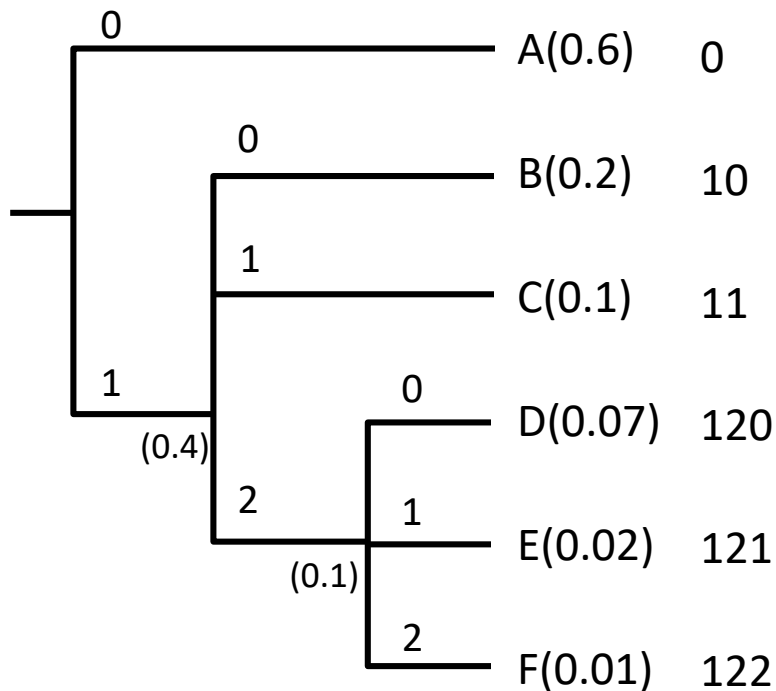
コンパクト符号



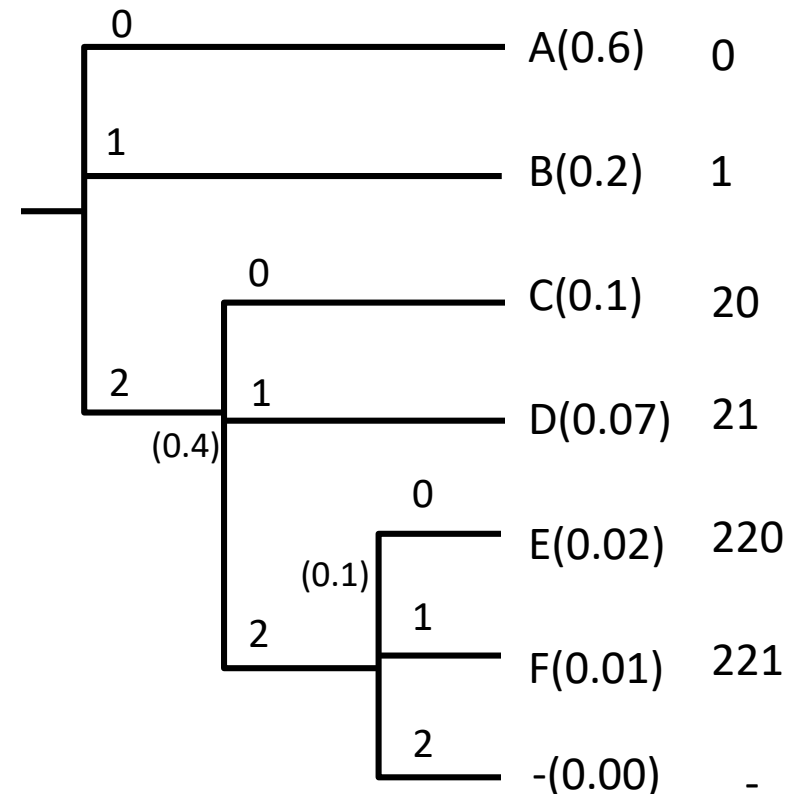
ハフマン符号

■ 3元の場合

例：情報源記号発生確率を〈A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.02, F: 0.01〉



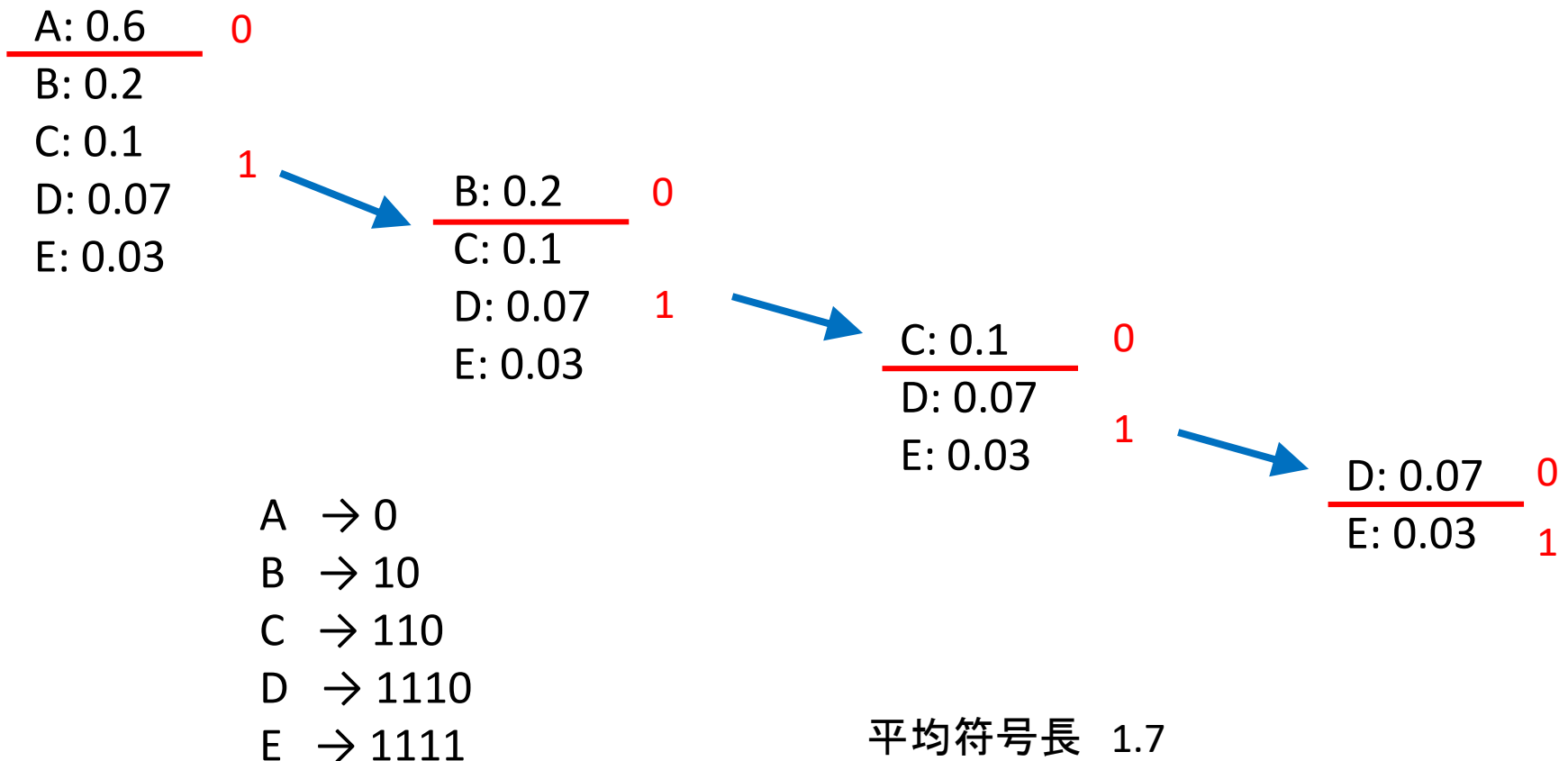
平均符号長 1.5



平均符号長 1.23

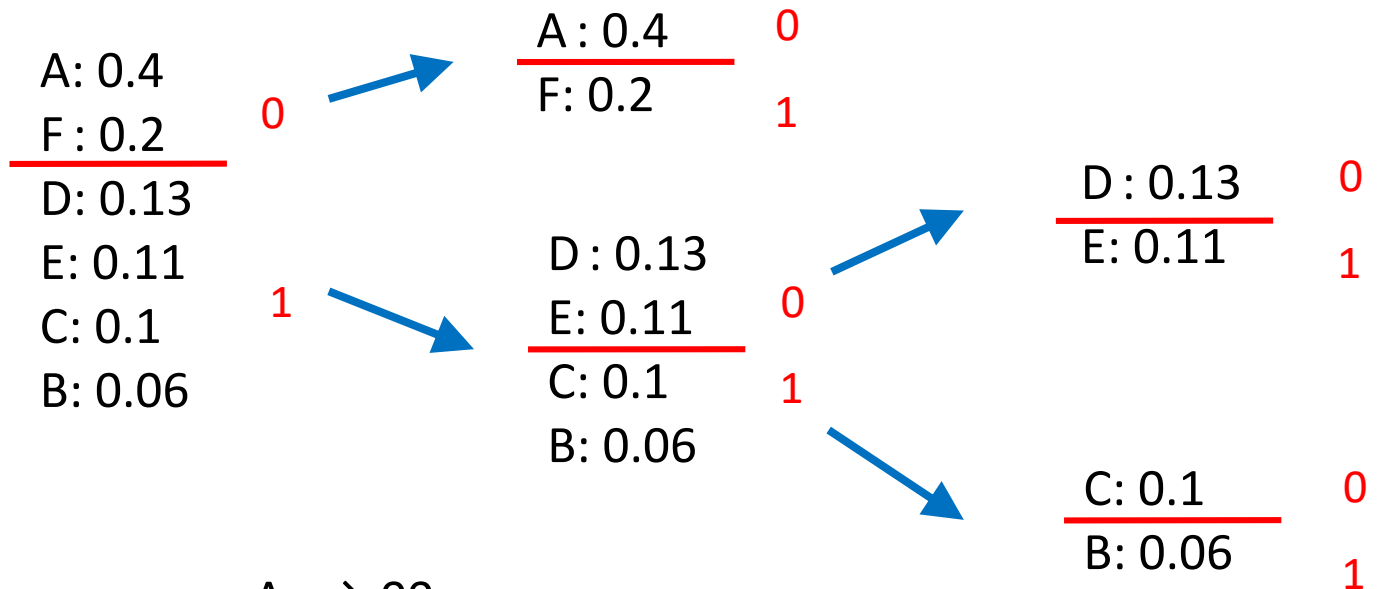
シャノン-ファノ符号

例: 情報源記号発生確率を〈A: 0.6, B: 0.2, C: 0.1, D: 0.07, E: 0.03〉



シャノン-ファノ符号

例: 情報源記号発生確率を〈A: 0.4, B: 0.06, C: 0.1, D: 0.13, E: 0.11, F: 0.2〉



A → 00
F → 01
D → 100
E → 101
C → 110
B → 111

平均符号長 2.4

課題

- シヤノン-ファノ符号化は必ずしもコンパクト符号を生成するとは限らない. どのような場合か？

まとめ

- コンパクト符号: 平均符号長の最も短い瞬時符号
- ハフマン符号: コンパクト符号構成法
- どうしてハフマン符号がコンパクト符号になるのか?
- 補助定理1: 瞬時符号がコンパクト符号であるための必要条件
- 補助定理2: コンパクトな瞬時符号の十分条件
- コンパクト符号は唯一とは限らない
- ハフマン符号で構成されないコンパクト符号
- 3元符号の場合
- シャノン-ファノ符号