



情報エレクトロニクス学科共通科目・2年次・夏ターム〔必修科目〕

講義「情報理論」

第8回

第5章 情報源符号化法

5.3 非等長情報系列の符号化

5.4 ひずみが許される場合の情報源符号化

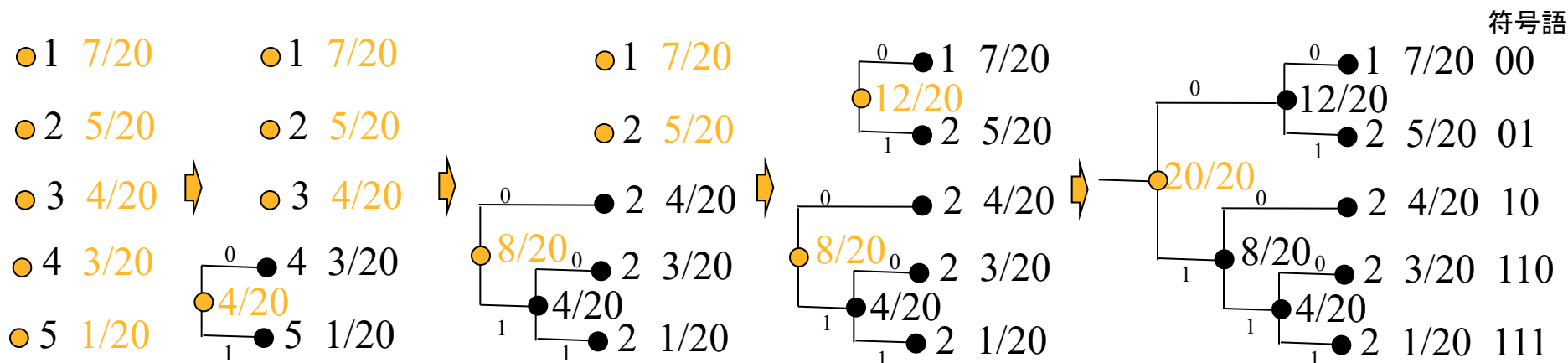


前回の演習問題

[問] 袋の中に空くじなしの1等から5等までのくじが20本入っている。各等の本数は、下表の通りとする。

1等	2等	3等	4等	5等
1	3	4	5	7

このくじを1本ずつ繰り返し引き、何等が当たったかを情報源と考え、 i 等が当たったとき情報源記号 i を発生させることにより情報源系列をつくる。ただし、引いたくじは次のくじを引く前に袋に戻すものとする。この情報源を、符号アルファベット $\{0,1\}$ を使って、情報源記号1つ1つを一意に復号可能な2元符号に符号化する場合において、コンパクト符号を1つ求めよ。また、その符号の平均符号長を求めよ。





[復習] コンパクト符号

コンパクト符号

各情報源記号を符号化した一意復号可能な符号で
平均符号長が最小の符号

情報源 S の2元コンパクト符号の平均符号長 $\geq H_1(S)$



各情報源記号を2元符号化した場合の本当の限界

代表的なコンパクト符号
ハフマン符号



[復習]ブロックハフマン符号化

情報源 S から発生する n 個の情報源記号ごとにブロック2元符号化

n 記号ごとに符号化した符号の1情報源記号当たりの平均符号長 L_n が以下の式を満たすものが存在

$$H_1(S^n)/n \leq L_n < H_1(S^n)/n + 1/n \quad \cdots \cdots \textcircled{1}$$

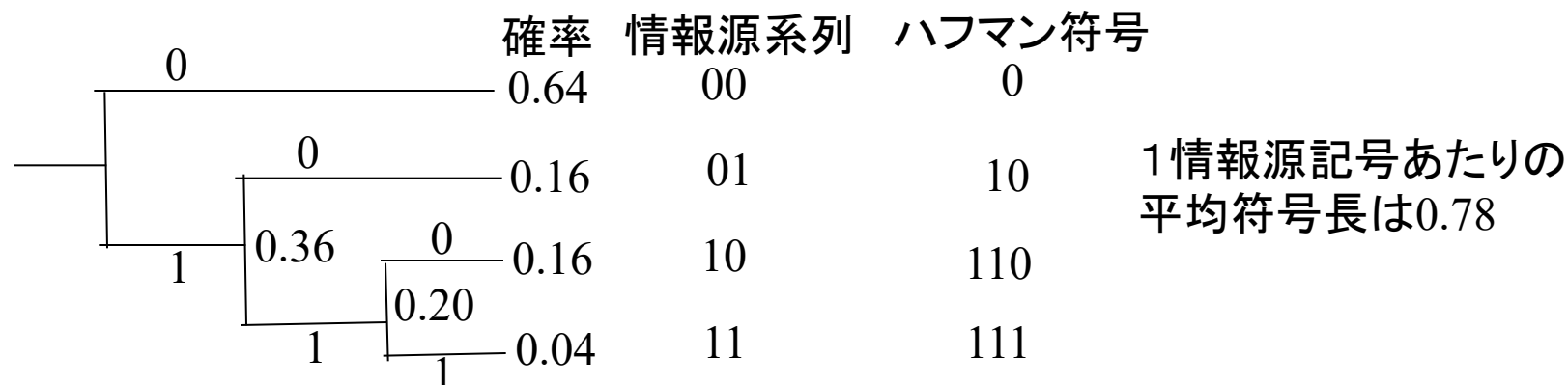
($H_1(S^n)/n$ は S の **n 次エントロピー**と呼ばれる量)

ブロックハフマン符号化

情報源 S から発生する n 個の情報源記号ごとにハフマン符号化

平均符号長 L_n は $\textcircled{1}$ 式を満たす $\rightarrow n$ を無限大に近づけると $H(S)$ に近づく。

[例: 1の出る確率が0.2の2元無記憶定常情報源を2情報源記号ごとにハフマン符号化]





ブロックハフマン符号化法の問題点

■ ブロックハフマン符号化法

- n 次の拡大情報源に対してハフマン符号化を行うやり方
- n を大きくすることにより、1情報源記号あたりの平均符号長をいくらでも下限 $H(S)$ に近づけられる。

しかし

長さ n のすべての情報源系列と符号語の巨大な対応表が必要

[例] 1, 0 の発生確率が 0.01, 0.99 の無記憶情報源 S

$$H(S) = 0.081$$

目標: 平均符号長 L を 0.089 以下

無記憶だから $L < H(S) + 1/n$ となり、 n を $1/0.008 = 125$ 以上にすれば確実。

$n = 125$ の系列は $2^{125} \doteq 4 \times 10^{37}$ 個もある！

事実上、そのようなハフマン符号を構成することはできない・・・

無記憶ならば

$$H_n(S) = H_1(S) = H(S)$$

$$1/125 = 0.008$$



非等長情報源系列の符号化

- n 記号毎符号化するハフマンブロック符号化法では、符号化すべき情報源系列の数は、 M 元情報源の場合、 M^n 個！
- 符号化すべき情報源系列を**非等長**にしてはどうだろうか？
 - 長い情報源系列と短い情報源系列を組み合わせ、
長いがよく発生する系列に、短い符号語を割り当てる
 - 符号化する情報源系列の数を減らし、
符号化のために記憶すべき表を削減する



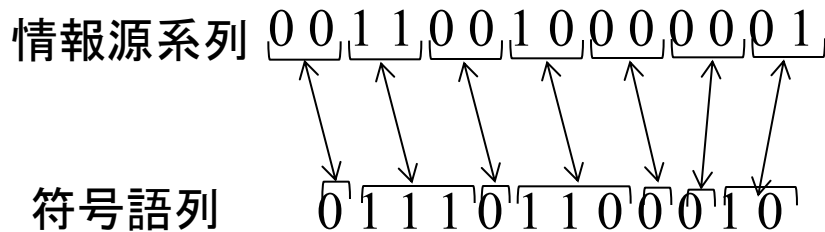
- 比較的に長いブロックで符号化したときと同じような効果を持たせられないだろうか？



[例]2元無記憶定常情報源の符号化

1の出る確率が0.2の2元無記憶定常情報源を考える

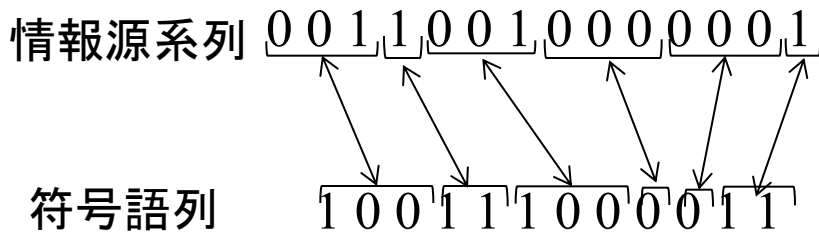
ブロック符号化



確率	情報源系列	符号語
0.64	00	0
0.16	01	10
0.16	10	110
0.04	11	111

1 情報源記号
あたりの
平均符号長は
0.78

非等長情報源系列の符号化



確率	情報源系列	符号語
0.512	000	0
0.128	001	100
0.16	01	101
0.2	1	11

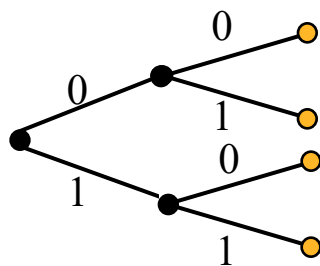
1 情報源記号
あたりの
平均符号長は
0.728

同じ数の符号語でより効率的な符号を実現

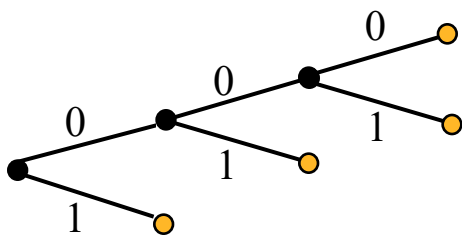


分節木

分節木



符号化する情報源系列
(橙の節点)を表す木



平均系列長が長いほど効率的

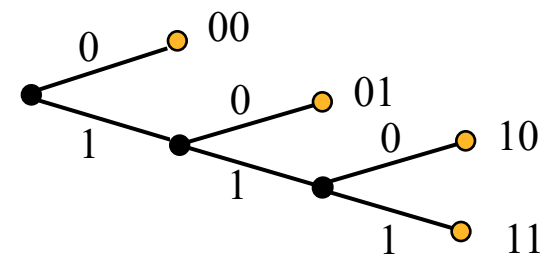
情報源系列 符号語

00	0
01	10
10	110
11	111

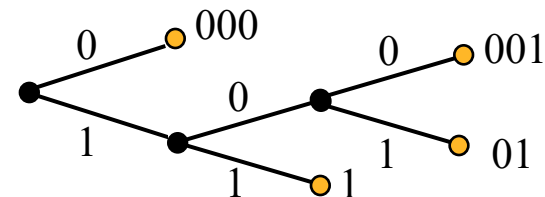
情報源系列 符号語

000	0
001	100
01	101
1	11

符号の木



符号語(橙の節点)を表す木



平均符号長が短いほど効率的



ランレングス符号化法(1)

■ ランレングス符号化法

情報源系列において同じ記号が連続する長さ
(ランレングス)を符号化して送る方法

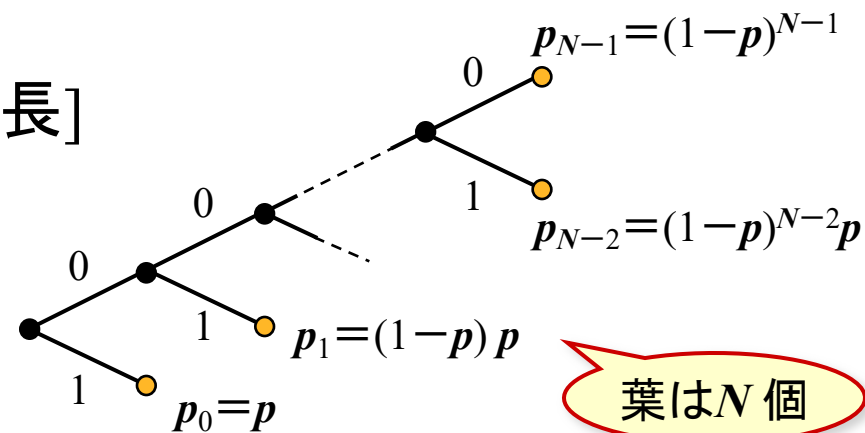
- 右上の符号は長さ3までの **0の連続(ラン)** に対するランレングス符号である。
情報源系列 1, 01, 001, 000 はそれぞれ 0のランが 0, 1, 2, 3 である。
- ランレングスをさらにハフマン符号化する方法を、
ランレングスハフマン符号化法と呼ぶ。

情報源系列	符号語
000	0
001	100
01	101
1	11

[ランレングス符号化法の平均符号長]

1, 0の発生確率が $p, 1-p$ ($p < 1-p$) の
無記憶情報源の、 $N-1$ までの 0 の
ランレングスを符号化

右の分節木で示されるような N 個の情報源系
列に対して符号化



ランレングス符号化のための分節木



ランレングス符号化法(2)

これら N 個の系列の平均系列長は

$$\begin{aligned}\bar{n} &= \sum_{i=0}^{N-2} (i+1) p_i + (N-1) p_{N-1} = \sum_{i=0}^{N-2} \underbrace{(i+1)(1-p)^i p}_{0^i 1 \text{ の系列}} + \underbrace{(N-1)(1-p)^{N-1}}_{0^{N-1} \text{ の系列}} \\ &= \frac{1 - (1-p)^{N-1}}{p}\end{aligned}$$

となる。一方、これらの系列をハフマン符号化したときの平均符号長 L_N は定理4.2から

$$L_N < -\sum_{i=0}^{N-1} p_i \log_2 p_i + 1 = H(S) \bar{n} + 1$$

を満たす。よって、1情報源記号あたりの平均符号長 L_r は

$$L_r = \frac{L_N}{n} < H(S) + \frac{1}{\bar{n}} = H(S) + \frac{p}{1 - (1-p)^{N-1}}$$

ハフマンブロック符号化の場合と比較すると

$$L_h < H(S) + \frac{1}{n} = H(S) + \frac{1}{\log_2 N}$$

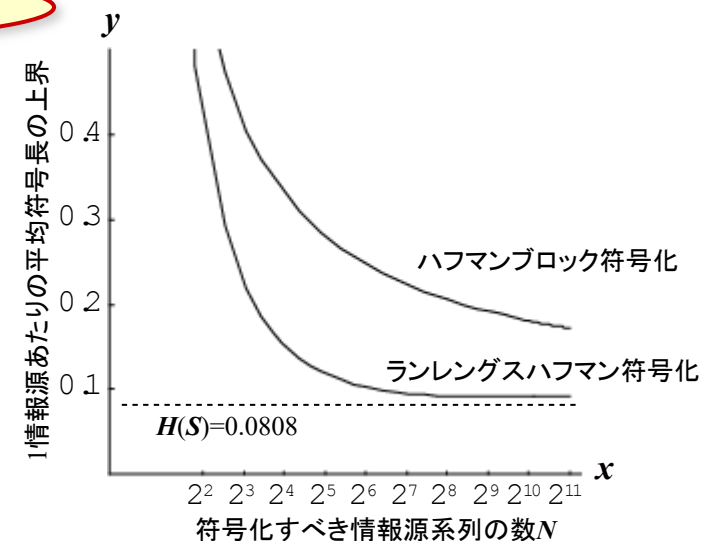


図5.11 ランレングス符号化とハフマンブロック符号化の比較($p=0.01$)

N : ハフマンブロック符号化において符号化すべき情報源系列の数



[発展] タンストール木による符号化(1)

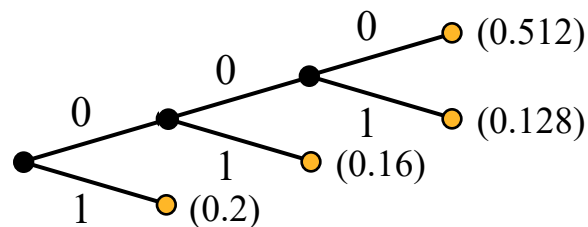
情報源系列を瞬時に(符号語を割り当てる情報源系列の境目で)分解可能
 \Leftrightarrow 分節木の葉のみに符号化すべき情報源系列が割り当たっている

N個の符号化すべき情報源系列で平均長が最大なもの

\Leftrightarrow N個の葉をもつ分節木において葉の深さの平均が最大なもの
 (平均は葉に対応する記号列の出現確率に関してとる)

タンストール木

M元情報源の場合、「最大確率が割り当てられている葉を、M個の葉を子としてもつノードに置き換える」ことを繰り返すことにより最適な分節木が得られる



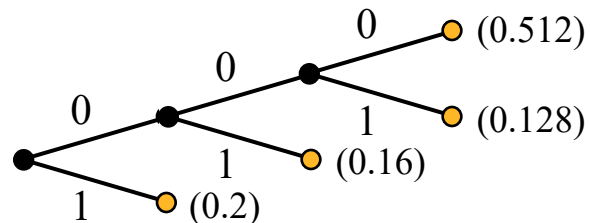
平均長 n は

$$n = 1 \times 0.2 + 2 \times 0.16 + 3 \times 0.128 + 3 \times 0.512 = 2.44$$

1の発生確率が0.2の2元無記憶定常情報源の4枚の葉をもつタンストール木



[発展] タンストール木による符号化(2)



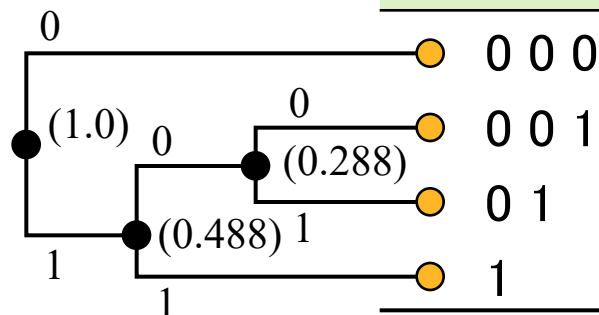
平均長 $\bar{n} = 2.44$

1の発生確率が0.2の2元無記憶定常情報源の
4枚の葉をもつタンストール木

右の符号の平均符号長 = 1.776

よって1記号あたりの
平均符号長 L は

$$L = \frac{1.776}{2.44} = 0.728$$

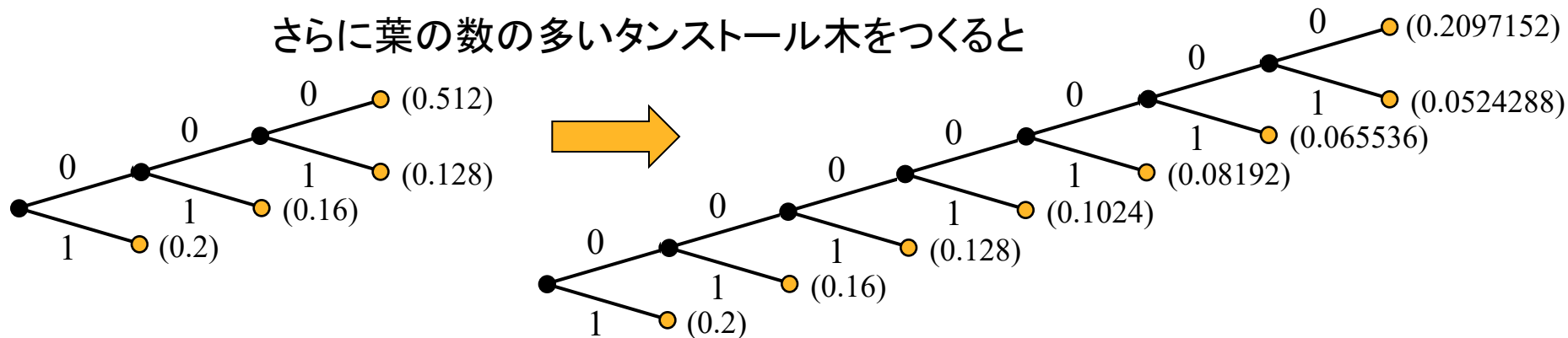


情報源系列	確率	ハフマン 符号
0 0 0	0.512	0
0 0 1	0.128	100
0 1	0.16	101
1	0.2	11



[発展] タンストール符号

さらに葉の数の多いタンストール木をつくると



最大確率
最小確率

0.512
0.128

0.384

差が縮まっていく

0.1572864

0.2097152
0.0524288

最大確率
最小確率

コンパクト符号と等長符号の平均符号長の差がなくなっていく

タンストール符号

① N枚の葉をもつタンストール木を構築 → ② 等長符号化

可変長系列-固定長符号(Variable-length-to-Fixed-length code: VF符号)

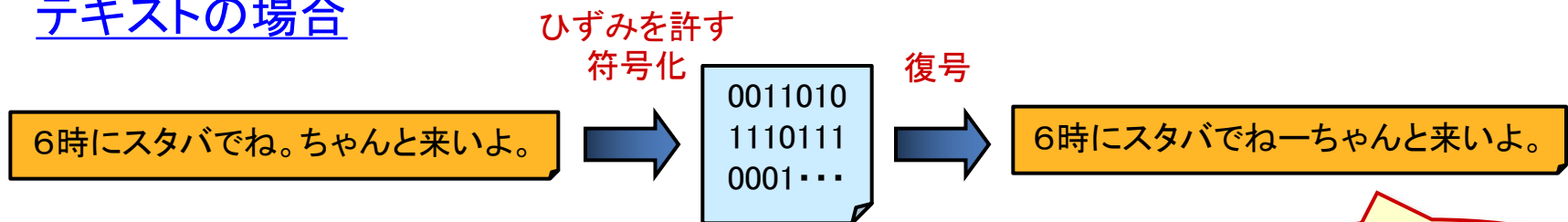
$N \rightarrow \infty$ とすれば1情報源記号あたりの平均符号長は $H(S)$ に収束



ひずみが許される場合とは？

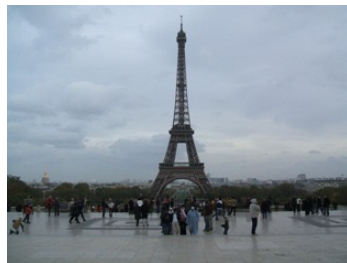
- ある程度ひずみを許しても、1情報源記号あたりの平均符号長を小さくしたい！（つまり、より小さくデータを圧縮したい）

テキストの場合



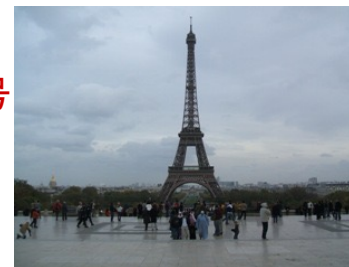
非常にまずい！

画像の場合



ビットマップファイル
(.bmp; 288 KB)

ひずみを許す符号



jpegファイル
(.jpg; 17.5 KB)

ほとんど問題ない



情報源符号化におけるひずみ

- 通信路でひずみが入るのではなく、符号化時にひずみを入れる



図：情報源出力 x と情報源復号結果 y

- **ひずみ測度**： x と y の相違を評価する関数 $d(x, y)$
 - 関数 $d(x, y)$ が大きいほど、ひずみが大きい。また次の性質を持つ。

$$d(x, y) \geq 0$$

$$x=y \text{ のとき } d(x, y)=0$$

- ひずみ測度の平均値を**平均ひずみ**と呼び、 \bar{d} で表す。

$$\bar{d} = \sum_x \sum_y d(x, y) P(x, y)$$



ひずみ測度の例

- 例1 情報源アルファベットを $A = \{0, 1\}$ とし、ひずみ測度を

$$d(x, y) = \begin{cases} 0 & ; x=y \\ 1 & ; x \neq y \end{cases}$$

とする。このとき、平均ひずみは

$$\bar{d} = \sum_x \sum_y d(x, y) P(x, y) = P(1, 0) + P(0, 1)$$

$P(1, 0)$: 入力 1 → 出力 0

$P(0, 1)$: 入力 0 → 出力 1

となる。これは、要するに、符号器の出力が元の情報源出力と異なる確率であり、通常**ビット誤り率**と呼ばれる。

- 例2 情報源アルファベットを有限個の整数または実数の集合としよう。このとき、ひずみ測度を

$$d(x, y) = |x - y|^2$$

とすれば、平均ひずみは**2乗平均誤差** (mean square error) と呼ばれる量となる。ひずみの評価量として非常によく用いられる。



ひずみが許される場合、情報源符号化定理はどうなるか？

[情報源符号化定理]

情報源 S は、任意の正数 ε に対して、1情報源記号あたりの平均符号長 L が

$$H(S) \leq L < H(S) + \varepsilon$$

となるような2元瞬時符号に符号化できる。

しかし、どのような一意復号可能な2元符号を用いても、平均符号長が $H(S)$ より小さくなるような符号化はできない。

$H(S)$ が限界！

■ ひずみを許せば、どのくらい平均符号長の限界を下げられるか？

S : 記憶の無い定常情報源、 X : 任意の時点において S から出力される情報源記号

– 1情報源記号あたりの平均符号長の下限 = 情報量 $H(S)=H(X)$



– ひずみを許した場合、出力 Y の値を知っても、元の入力 X に関して、なお平均 $H(X|Y)$ のあいまいさが残る。

– 伝えられる情報の量は $I(X;Y)=H(X)-H(X|Y)$

ひずみを許した場合の限界！



ひずみが許される場合の情報源符号化定理

- 相互情報量 $I(X;Y)$ が同じでも、平均ひずみ \bar{d} は同じとは限らない
 \Leftrightarrow 平均ひずみ \bar{d} が同じであっても、 $I(X;Y)$ は符号化の仕方で異なる
- ある与えられた値 D に対し、平均ひずみ \bar{d} が

$$\bar{d} \leq D$$

を満たす条件の下で、あらゆる情報源符号化法を考えたときの相互情報量 $I(X;Y)$ の最小値を考え、これを $R(D)$ と表す。すなわち、

$$R(D) = \min_{\bar{d} \leq D} \{I(X;Y)\}$$

これを情報源 S の **速度・ひずみ関数** (rate-distortion function) と呼ぶ。

[ひずみが許される場合の情報源符号化定理]

平均ひずみ \bar{d} を D 以下に抑えるという条件の下で、任意の正数 ε に対して、情報源 S を1情報源記号あたりの平均符号長 L が

$$R(D) \leq L < R(D) + \varepsilon$$

となるような2元符号へ符号化できる。しかし、どのような符号化を行っても、 $\bar{d} \leq D$ である限り、 L を $R(D)$ より小さくすることはできない。



速度-ひずみ関数(1)

- $R(D)$ は、 $\bar{d} \leq D$ を満たす、あらゆる情報源符号化法を考えたときの $I(X;Y)$ の最小値。では実際、これをどのようにして求めるのか？
- 情報源符号化法を変えると、条件付確率 $P(y|x)$ が変わる！
 - $I(X;Y)$ の最小化は、 $\bar{d} \leq D$ の条件の下で $P(y|x)$ を変えることで行う*
 - 下図のような通信路を考え、この通信路の特性を変えることで $I(X;Y)$ の最小化を計ると解釈できる。この仮想的通信路を試験通信路と呼ぶ。

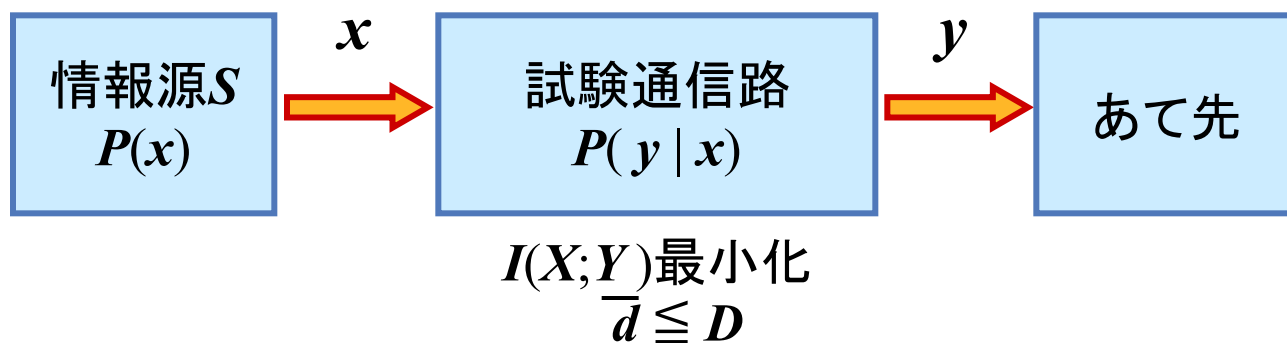


図: 試験通信路による速度・ひずみ関数 $R(D)$ の解釈

* 任意の条件付確率 $P(y|x)$ を与えるような情報源符号化・復号法が存在する



速度-ひずみ関数(2)

- 条件付確率 $P(y|x)$ を用いれば、相互情報量 $I(X;Y)$ は

$$\begin{aligned} I(X;Y) &= \sum_x \sum_y P(x,y) \log_2 \frac{P(x,y)}{P(x)P(y)} \\ &= \sum_x P(x) \sum_y P(y|x) \log_2 \frac{P(y|x)}{P(y)} \dots\dots\dots ① \end{aligned}$$

と表せる。ここに $P(y)$ は

$$P(y) = \sum_x P(x) P(y|x)$$

により計算できる。また、 $\bar{d} \leq D$ という条件は

$$\bar{d} = \sum_x P(x) \sum_y P(y|x) d(x,y) \leq D \dots\dots\dots ②$$

と表せる。また、条件付確率は

$$P(y|x) \geq 0 \dots\dots\dots ③$$

$$\sum_y P(y|x) = 1 \text{ for each } x \dots\dots\dots ④$$

を満たさなければならない。

- 式②～④の条件の下に、式①の $I(X;Y)$ を最小化すればよい！

$$P(x,y) = P(x)P(y|x)$$

$P(y|x)$ が変数
他は既知

一般には難しい



[例]記憶のない2元情報源の場合(1/3)

- 1, 0 を確率 $p, 1-p$ で発生する記憶のない2元情報源Sを考える。

また、ひずみ測度としては

$$d(x, y) = \begin{cases} 0 & ; x=y \\ 1 & ; x \neq y \end{cases}$$

を用いるものとする。このとき、平均ひずみ \bar{d} はビット誤り率となる。

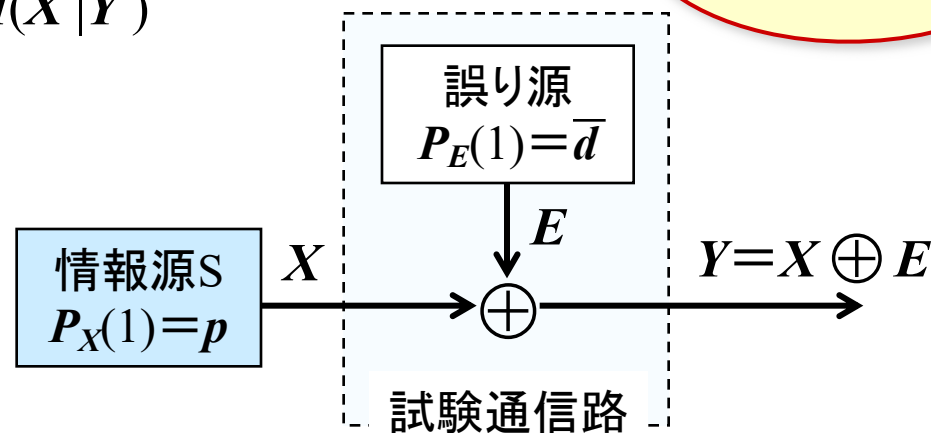
- この情報源に対して、 $0 \leq D \leq 0.5$ が与えられたとき、 $\bar{d} \leq D$ の元での速度・ひずみ関数 $R(D)$ を求めよう。

- 相互情報量 $I(X; Y) = H(X) - H(X|Y)$

$H(X) = \mathcal{H}(p)$ なので、

$H(X|Y)$ を最大化すればよい。

- ここで、 Y は左図のように、1の発生確率が \bar{d} であるような誤り源の出力 E と X の排他的論理和で表せる。



$D > 0.5$ の場合は
誤る確率の方が
大きいことを意味する

図：2元情報源に対する試験通信路



[例]記憶のない2元情報源の場合(2/3)

- $Y = X \oplus E$ であるから、 $X = Y \oplus E$ となる。したがって、

$$H(X | Y) = H(Y \oplus E | Y) = H(E | Y)$$

が得られる。

- 条件付きエントロピーの性質から $H(E | Y) \leq H(E)$ が成り立つ。さらに、誤り源に記憶がなく定常であれば、 $H(E) = \mathcal{H}(\bar{d})$ であるが、そうでなければ、 $H(E) < \mathcal{H}(\bar{d})$ であるから、

$$H(E | Y) \leq H(E) \leq \mathcal{H}(\bar{d})$$

となる。それゆえ

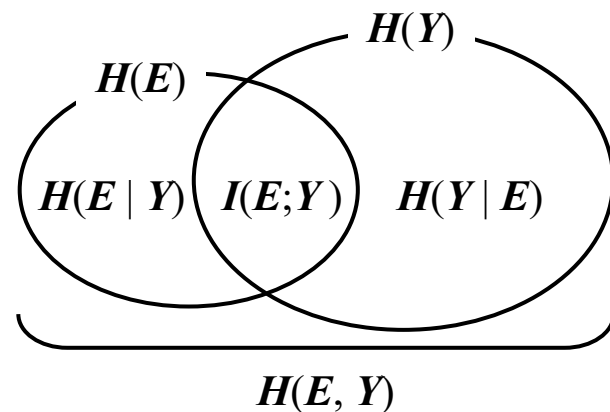
$$H(X | Y) \leq \mathcal{H}(\bar{d})$$

を得る。 $\bar{d} \leq D \leq 1/2$ なので、さらに

$$\mathcal{H}(\bar{d}) \leq \mathcal{H}(D)$$

となる。したがって、相互情報量 $I(X; Y)$ は、

$$I(X; Y) = H(X) - H(X | Y) \geq \mathcal{H}(p) - \mathcal{H}(D)$$





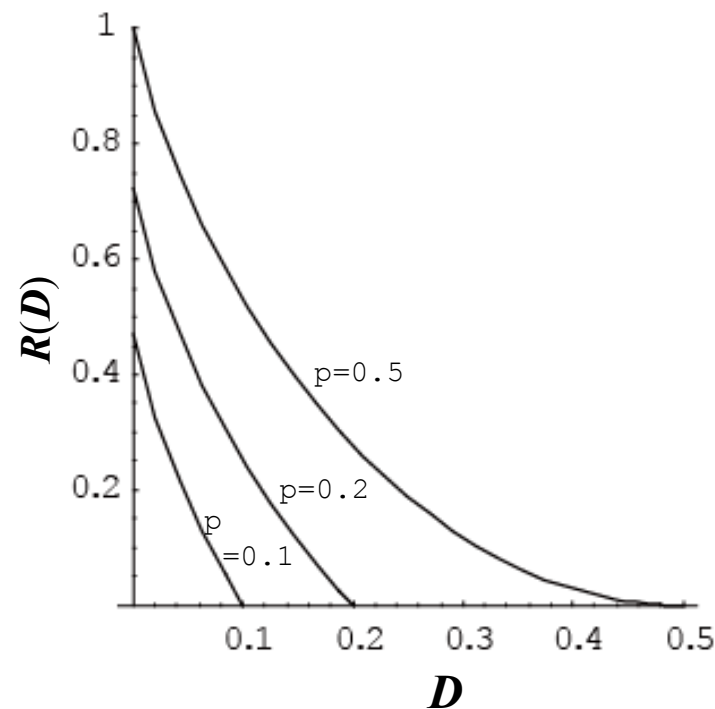
[例]記憶のない2元情報源の場合(3/3)

- $\bar{d}=D$ で誤り源が無記憶定常で情報源Sと独立であるとき、等号が成立 ($I(X;Y) = \mathcal{H}(p) - \mathcal{H}(D)$) するので、記憶のない定常2元情報源S の速度・ひずみ関数は

$$R(D) = \mathcal{H}(p) - \mathcal{H}(D)$$

で与えられることが導けた。

- 左図で分かるように、速度・ひずみ関数は、 **D に関して単調減少**であり、**下に凸な関数**である。一般の速度・ひずみ関数も同様な性質を持つことが証明されている。
- 記憶のある情報源の場合にも、
$$I(X;Y) = \lim_{n \rightarrow \infty} I(X_n;Y_n)/n$$
の最小値として、速度・ひずみ関数を定義することができる。



図：2元情報源の速度・ひずみ関数



ひずみが許される場合の情報源符号化法

〔ひずみが許される場合の情報源符号化定理〕

平均ひずみ \bar{d} を D 以下に抑えるという条件の下で、任意の正数 ε に対して、情報源 S を1情報源記号あたりの平均符号長 L が

$$R(D) \leq L < R(D) + \varepsilon$$

となるような2元符号へ符号化できる。しかし、どのような符号化を行っても、 $\bar{d} \leq D$ である限り、 L を $R(D)$ より小さくすることはできない。

- ひずみが許される場合の情報源符号化定理は、1情報源記号あたりの平均符号長を、速度・ひずみ関数 $R(D)$ にいくらでも近づく符号化法の存在を示している。では、**具体的な符号化方法はあるのか？**

– ひずみのない場合に比べて、はるかに難しい！

- ベクトル量子化
- 変換符号化
- 予測符号化
- etc.

参考図書

「情報・符号・暗号の理論」
今井秀樹 著、電子情報通信学会