

# 情報理論

## 情報量

### 英文字

アルファベット 26 文字, 空白, カンマ, ピリオド,  
アポストロフィ 計 30 文字と想定しよう.

30 文字を用いて,

100 字の文章を書く.

生成される文章の種類

$$30 \text{ の } 100 \text{ 乗通り} \quad N_{100} = 30^{100}$$

200 字の文章を書く.

$$30 \text{ の } 200 \text{ 乗通り} \quad N_{200} = 30^{200}$$

文字数が 2 倍になれば, 情報量は 2 倍と考えるのが妥当.

指数部の大きさを評価する → 対数  $\log$

## 2を底とした対数

$$\log_2(2^n) = n$$

$$\log_2(m^n) = n \times \log_2(m)$$

$$\log_2(10) = 3.3219$$

$$\Leftrightarrow 2^{3.3219} = 10$$

$$2^3 \leq 10 \leq 2^4$$

8                      16

底を2とした対数  $\log_2$  により情報量を測る.

$$N_{100} (= 30^{100}) \text{ の情報量} \\ \log_2(N_{100}) = 100 \times \log_2(30)$$

$$N_{200} (= 30^{200}) \text{ の情報量} \\ \log_2(N_{200}) = 200 \times \log_2(30)$$

情報量の比較

$$\log_2(N_{200}) = 2 \times \log_2(N_{100})$$

ちゃんと2倍になっている.

### 3ビット（2進数3桁）データの情報量

$2^3 = 8$ 通りの表現

0 0 0

0 0 1

0 1 0

0 1 1

1 0 0

1 0 1

1 1 0

1 1 1

$$\log_2(2^3) = 3$$

底を2とした対数がビット数に対応

ビット単位の情報量, ビットは情報量

$\log_2$  (表現の種類の数)

N 通りの表現が持つ情報量  $\log_2(N)$

英文字（30文字）の情報量

$$\log_2(30) = 4.907 \text{ ビット bits}$$

# 確率と情報量

サイコロ :

6通りの表現 ( $N = 6$ )

情報量  $\log_2(6) = 2.585$  bits

ある目の出る確率  $p = 1/N = 1/6$

逆に,  $N = 1/p$

確率  $p$  の事象が知れたことによる情報量は,

$$\log_2(N) = \log_2(1/p)$$

サイコロの、ある目がでたときの情報量

$$\begin{aligned}\log_2(1/p) &= \log_2(1/(1/6)) \\ &= \log_2(6) \\ &= 2.585 \text{ bits}\end{aligned}$$

対数の公式  $\log(1/x) = -\log(x)$   
より

$$\log_2(1/p) = -\log_2(p)$$

確率  $p$  の事象の情報量

$$\log_2(1/p), \quad -\log_2(p)$$

両者は等価

## おみくじ：吉の方角

東（ $1/2$ ），西（ $1/4$ ），  
南（ $1/8$ ），北（ $1/8$ ）

東と知れたことによる情報量

$$-\log_2(1/2) = 1 \text{ bit}$$

西と知れたことによる情報量

$$-\log_2(1/4) = 2 \text{ bits}$$

南と知れたことによる情報量

$$-\log_2(1/8) = 3 \text{ bits}$$

確率が低いほど，珍しいことほど，  
それが起こると情報量が多い。

# 情報エントロピー

全部で  $n$  個の事象

$i$  番目の事象の起こる確率  $p_i$

$i$  番目の事象が起きたと知れたことによる情報量

$$- \log_2 (p_i)$$

$i$  番目の事象によりもたらされると期待される情報量

$$p_i \times (- \log_2 (p_i))$$

$n$  個全部の事象により期待される情報量  $H$

(情報エントロピー)

$$p_1 \times (- \log_2 (p_1))$$

$$+ p_2 \times (- \log_2 (p_2))$$

$$+ \dots + p_n \times (- \log_2 (p_n))$$



## 方角おみくじの情報エントロピー H

4 個の事象

東 (1 / 2) , 西 (1 / 4) ,  
南 (1 / 8) , 北 (1 / 8)

$$\begin{aligned} H &= (1/2) \times (-\log_2(1/2)) \\ &+ (1/4) \times (-\log_2(1/4)) \\ &+ (1/8) \times (-\log_2(1/8)) \\ &+ (1/8) \times (-\log_2(1/8)) \\ &= (1/2) + (1/2) \\ &\quad + (3/8) + (3/8) \\ &= 1.75 \text{ bits} \end{aligned}$$

n 個の事象で、情報エントロピーが最大になるとき、

$$p_i = p = 1/n$$

$$H = -\log_2(p)$$

(例)

情報エントロピーが最大となる方角おみくじ

(n = 4)

各方角への確率が 1/4

$$H = -\log_2(1/4)$$

$$= \log_2(4)$$

$$= 2$$

n = 8 であれば、p = 1/8 であり、H = 3.

n個の事象で, ひとつの事象の確率のみが 1 で,  
他の事象の確率がゼロであるとき

$$p_1 = 1, \quad p_2 = \dots = p_n = 0$$

$$H = -1 \times \log_2(1) = -1 \times 0 = 0$$

事象間で, 確率の偏りが無いほど,  
情報エントロピは大きくなる.

# 伝達する情報の符号

## 符号化

文字	東	西	南	北
符号系A	0	1	0 1	1 1
符号系B	0	1 0	1 1 0	1 1 1

送信情報 「 東西南北 」

符号系A 0 1 0 1 1 1

東西南北？ 東西東北西？ 南東西西西？

符号系B 0 1 0 1 1 0 1 1 1

東西南北

一意解読可能符号系

瞬時解読可能符号系

情報が一意に復元できる.

順次先頭から復元できる.

# 一意解読可能, 瞬時解読可能符号系を与える ハフマンの符号化 (アルゴリズム)

- (1) 発生確率順に事象を並べる.
- (2) 小さい発生確率を持つ事象を 2 つ選び, 1 つに 0, もう 1 つに 1 を割り付ける.
- (3) ステップ (2) における 2 つの事象の確率の和を発生確率とする事象を生成し, その生成した事象を, 和を取る前の事象に代わる新たな事象として加え, あらためて全事象を発生確率順に並べる.
- (4) 事象が最後の 1 つになるまで, (2), (3) を繰り返す.

# 東西南北 (E, W, S, N) の符号系 (その1)

E (1/2)				1	← 最後の事象
			# (1/2)	0	
W (1/4)				1	
		* (1/4)		0	
S (1/8)	1				
N (1/8)	0				

## 符号系

東 (E)	1	西 (W)	0 1
南 (S)	0 0 1	北 (N)	0 0 0

任意性 : 同確率のときに, 上下どちらに配置するか

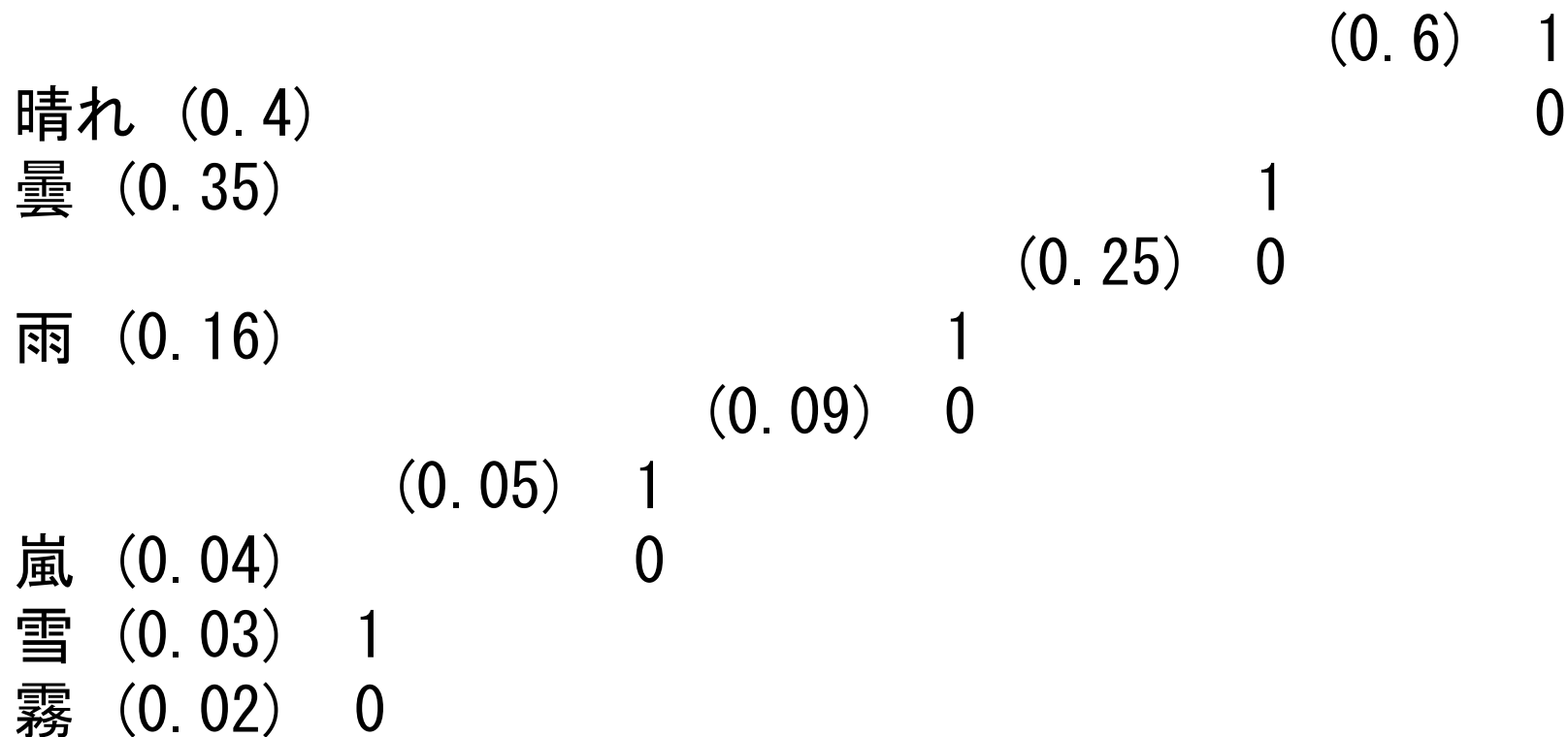
# 東西南北 (E, W, S, N) の符号系 (その2)

E (1/2)					0
				# (1/2)	1
		* (1/4)	1		
W (1/4)			0		
S (1/8)	0				
N (1/8)	1				

## 符号系

東	0			西	1 0
南	1 1 0			北	1 1 1

# 天気（晴，曇，雨，嵐，雪，霧）の符号化



## 天気符号系

晴	0	曇	1 1	雨	1 0 1	嵐	1 0 0 0
雪	1 0 0 1 1	霧	1 0 0 1 0				



## 情報量の比較

天気の符号系の平均情報量 (平均符号長)

晴 (0.4)	1 bit,	曇 (0.35)	2 bits,
雨 (0.16)	3 bits,	嵐 (0.04)	4 bits,
雪 (0.03)	5 bits,	霧 (0.02)	5 bits

$$\begin{aligned} & 0.4 \times 1 + 0.35 \times 2 + 0.16 \times 3 \\ & + 0.04 \times 4 + 0.03 \times 5 + 0.02 \times 5 \\ & = 1.99 \text{ bits} \end{aligned}$$

天気の情報エントロピー  $1.932 \text{ bits}$

6面サイコロが持つ情報エントロピー

$$-\log_2 (1/6) = 2.585 \text{ bits}$$

## ハフマンの符号化

- ・ 確率が高いほど短いコードは短い.
- ・ 平均符号長が最短の瞬時解読可能符号系