

1 符号化

1.1 アルファベットと言語

符号理論で用いるいくつかの言い回しを導入しておく。

有限集合 A に対し、その元の有限列全体の成す集合

$$A^+ := \bigcup_{n \in \mathbb{N}} A^n$$

を考える。 A の元の事を文字、 A の事をアルファベット、 A^+ の元の事を語と言う。

例 1. base, ball 等は $\{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}$ をアルファベットとする語である。

例 2. 0, 1, 01, 1011, 011011 等は $\{0, 1\}$ をアルファベットとする語である。

例 3. $\cdot, -, \dots$ 等は $\{\cdot, -, \dots\}$ をアルファベットとする語である。

語の結合を表現する二項演算子 $\cdot : A^+ \times A^+ \rightarrow A^+$ を次の様に定める。¹

$$(a_1, \dots, a_n) \cdot (b_1, \dots, b_m) := (a_1, \dots, a_n, b_1, \dots, b_m)$$

例 4. base \cdot ball = baseball, ball \cdot base = ballbase (非可換なことに注意)

例 5. 01 \cdot 1011 = 011011

例 6. $\dots \cdot - = \dots -$

利便性の為に「足しても何も変化しない語」として空語 ϕ を A^+ に追加したものを A^* で表す²:

$$A^* := \{\phi\} \cup \bigcup_{n \in \mathbb{N}} A^n$$

$u, v, w \in A^*$ について

$$w = u \cdot v$$

である時 $u(v)$ は w の接頭語 (接尾語) であると言い

空語 ϕ は任意の語 w に対し常に接頭語 (接尾語) であるが特にそれ以外のものを直接頭語 (直接尾語) と言う。

例 7. base は baseball の接頭語, ball は baseball の接尾語。

例 8. 0, 01 等は 011011 の接頭語, 11, 011 等は 011011 の接尾語。

語に含まれる文字の数

$$|(a_1, \dots, a_n)| := n, |\phi| := 0$$

を語長と呼ぶ。

例 9. baseball の語長は 8。

例 10. 011011 の語長は 6。

A^* の部分集合の事を言語と呼ぶ。

例 11. $\{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}^*$ のなかで“英単語”になるもの全を L と書けば, L は言語. baseball $\in L$ だが ballbase $\notin L$.

例 12. $\{41, 42, 43, 44, 45, 46, 47, 48, 49, 4A, 4B, 4C, 4D, 4E, 4F, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 5A\}$ は $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\}^*$ をアルファベットとする言語。

言語 L_1, L_2 が $L_1 \subset L_2$ を満たす時 L_2 は L_1 の拡大であるという。

¹この二項演算について A は (非可換な) 半群となる。

² A^* は ϕ を単位元にもつモノイドとなる。

1.2 符号化の定義

Σ, A を有限集合とする. Σ^n から A^+ への写像 φ を符号化と言う³. その像 $\varphi(L)$ は A をアルファベットにもつ言語を成す事に注意せよ. A の事を符号アルファベットと言う. $\varphi(L)$ の語を φ の符号語という.

特に $\phi: \Sigma \rightarrow A^+$ が与えられれば

$$\Sigma^n \ni (x_1, x_2, \dots, x_n) \mapsto \phi(x_1) \cdot \phi(x_2) \cdots \phi(x_n) \in A^+$$

によって任意の Σ^n からの符号化があたえられる. 以後, 特に断りの無い限り符号化は Σ からの写像を考える.

情報源文字 (Σ)	モールス信号 $A = \{ \cdot, - \}$	ASCII 符号 $A = \{0, \dots, 9, A, \dots, F\}$	冗長な符号 $A = \{A, \dots, Z\}$
a	· -	41	ALFA
b	- ...	42	BRAVO
c	- · - ·	43	CHARLIE
d	- · ·	44	DELTA
e	·	45	ECHO
f	· · - ·	46	FOXTROT
g	- - ·	47	GOLF
h	····	48	HOTE
i	··	49	INDIA
j	· - - -	4A	JULIETT
k	- - -	4B	KILO
l	· - · ·	4C	LIMA
m	- -	4D	MIKE
n	- ·	4E	NOVEMBER
o	- - -	4F	OSCAR
p	· - - ·	50	PAPA
q	- - - ·	51	QUEBEC
r	· - ·	52	ROMEO
s	···	53	SIERRA
t	-	54	TANGO
u	· · -	55	UNIFORM
v	··· -	56	VICTOR
w	· - -	57	WHISKEY
x	- · · -	58	XRAY
y	- · - -	59	YANKEE
z	- - · ·	5A	ZULU

例 13.

有限集合 Σ の元の数と同じであれば全射で自然に対応するから抽象的議論においては Σ として具体的に何を選ぶかは重要で無い場合が多い. $\#\Sigma = m$ なる符号を m 元符号と呼ぶ.

1.3 符号長

符号語の長さが全て等しい時, その符号は等長であるという. 上述の例で言えば ASCII 符号が等長符号である.

定義 1. 符号 $\varphi: \Sigma \rightarrow A^+$ と Σ を情報源アルファベットに持つ情報源 (Σ -値確率変数 X) について

$$E[|\varphi(X)|]$$

を平均符号長と言う.

明らかに等長符号の平均符号長はその長さと等しい.

例 14. $\Sigma = \{a, b, c, d\}$ を情報源アルファベットに持つ情報源

Σ	p_σ
a	0.125
b	0.125
c	0.25
d	0.5

³ Σ^n をアルファベットに持つ言語 $L \subset \Sigma^+$ からの写像に拡張することもある.

を考える。この情報源の情報量は

$$-\frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{2} \log_2 \frac{1}{2} = \frac{3}{8} + \frac{3}{8} + \frac{2}{4} + \frac{1}{2} = 1.75[\text{bit}].$$

この情報源に二種類の符号を考える。

Σ	p_σ	$\varphi_1(\sigma)$	$ \varphi_1(\sigma) $	$\varphi_2(\sigma)$	$ \varphi_2(\sigma) $
a	0.125	00	2	0001	4
b	0.125	01	2	001	3
c	0.25	10	2	01	2
d	0.5	11	2	1	1

符号 φ_1 の平均符号長は等長符号であるから 2 であり, 符号 φ_2 の平均符号長は

$$E[|\varphi_2(X)|] = 4 \times 0.125 + 3 \times 0.125 + 2 \times 0.25 + 1 \times 0.5 = 1.875$$

になって, より短い。

1.4 復号化

符号化で符号語に変換された情報は, 再び元の情報に戻せる事が期待される。

つまり

$$\psi : \varphi(\Sigma^n) \rightarrow \Sigma^n$$

で $\psi \circ \varphi$ が恒等写像であるものが存在することを期待する (ψ は復号化と呼ばれる)。その為には φ が単射であることが必要十分であり, 符号化 φ が単射でないとき特異, 単射のとき非特異という。

例 15. $\#\Sigma = m$ するとき,

$$\#\{0, 1\}^{\lceil \log_2 m \rceil} = 2^{\lceil \log_2 m \rceil} \geq 2^{\log_2 m} = m$$

であるから非特異な長さ $\lceil \log_2 m \rceil$ の 2 元等長符号が存在する。

Σ をアルファベットに持つ情報源に関して, 次の二つの事実を比較してみよう。

1. どんな分布についても, 情報量は $\log_2 m$ [bit] を超えない。
2. どんな分布についても, (平均) 符号長が $\lceil \log_2 m \rceil$ の非特異な等長符号が存在する。

この関係から, (2 元符号の) 平均符号長と (2 を底とした) 情報量が近い値であるという予想を立てることができる。この予想は後に符号化定理と言う形で正当化される。

最大の情報量を与えない情報源については, 符号長の最適化のためには, 等長でない符号が必要となる。

等長でない符号の場合は, 実用的な復号化の為に非特異なだけでは不十分である。例えば前述のモールス信号は非特異であるので

$$a \leftrightarrow \cdot -$$

の様に一文字の情報は復号できるが

$$abc \rightarrow \cdot - - \cdots - \cdot - \cdot - \rightarrow \text{pur}$$

の様に複数の文字の情報は別の文字として解釈できてしまう。