

文脈自由形言語の情報論的性質

NHK総合技術研究所 尾関 和彦

§ 1. まえがき

自然言語であれ人工言語であれ，言語の持つ最も重要な機能のひとつは情報の伝達であろう。それでは言語が伝達する情報量というものをどのように考えたらよいであろうか。

情報理論が創始されて以来，自然言語に対しても Markov 過程のモデルが考えられ，それにもとづいて一文字当りの平均情報量が計算されたこともある。Markov 過程のモデルを考えることは，その言語の文法として正規文法を考えることにほかならないが，近年の言語理論の進歩は正規文法が自然言語の文法として不十分であることを明らかにした。ここでは文法のクラスとして，自然言語の構造の基本的な部分を記述することができるといわれられている文脈自由形文法をとりあげ，それを言語能力として持つ話し手のモデルとして，プログラムシミュレーションに確率が付与された，いわゆる確率的文脈自由形

文法を考える。そしてその話し手の話す言語について情報論的議論を展開する。

§2. 言語が情報をにほう機構

ここに二人の人 A , B がいるとし、彼らは共通の言語を話すとする。彼らの言語能力のモデルとしてつぎのような文法自由形文法 G を仮定しよう：

$$G = (V_T, V_N, P, S),$$

$$V_T = \{花, 鳥, 娘, は, ととも, 美しい\},$$

$$V_N = \{S, 主部, 述部, 名詞\},$$

$$P = \{S \rightarrow 主部述部, 主部 \rightarrow 名詞は, 名詞 \rightarrow 花, \\ 名詞 \rightarrow 鳥, 名詞 \rightarrow 娘, 述部 \rightarrow 美しい, 述部 \rightarrow \\ ととも述部\}.$$

A および B が話しかつ理解するここのでできる文の全体、すなわち彼らの言語は、 $\{花は美しい, 娘は美しい, \dots, 花はととも美しい, \dots\}$ という可算無限集合である。いま A が B に

$$\text{'娘はととも美しい'} \quad (1)$$

と言ったとしよう。この文が A の伝えた情報にふさわしい、かつ B がこの文を学べるとして A の意思を知るここのでできる機構はどのようなものであろうか。(1) の生成過程を考えてみる

と、

$S \Rightarrow$ 主部述部 \Rightarrow 名詞は述部 \Rightarrow 娘は述部 \Rightarrow 娘はととも述部

\Rightarrow 娘はととも美しい (2)

というステップ^oを踏んで (1) が生成されることがわかる。この生成過程の第3ステップにおいて名詞 \rightarrow 娘というプロダクションが適用されているが、Aはここで名詞 \rightarrow 花、あるいは名詞 \rightarrow 鳥というプロダクションを選ぶこともできたはずである。つまり三つの可能なプロダクションの中からAは特に、名詞 \rightarrow 娘というプロダクションを選択したわけである。このようにして (1) は‘花’でもなく‘鳥’でもなく‘娘’は…というAの意思を伝えることができるのである。一方Bは(1)を受け取り、その生成過程の第3ステップにおいて三つの可能なプロダクションの中から特に名詞 \rightarrow 娘が選択され適用されていることを見出しAの意思を知ることができるであろう。すなわち文が情報を伝えることができるのは、その生成過程の各ステップにおいて適用するプロダクションを話し手が選択できたり、聞き手が文を受けとって情報を得ることができるのは、その文の生成過程の各ステップにおいていくつかの可能なプロダクションの中から実際にどのひとつが選択されたかを知るからにほかならない。このように文を生成過程に分析して考えれば‘情報を得る’とはいくつかの可能

性の中からどれが選択されたかを知ることである。' という情報に付する Shannon の基本的な考え方がそのまま適用できることがわかる。

ここでは話し手-聞き手の言語能力のモデルとして文脈自由形文法を考え、話し手の意志を、各プロダクションにそれが適用される確率を付与するという形でモデル化する。つまり確率的文脈自由形文法を話し手のモデルとする。それをもとにして文脈自由形言語に情報論的考察を加えていこう。

§ 3. 確率的文脈自由形言語

記法 1. 文脈自由形文法 (c.f.g.) G を通常の記法にしたがって $G = (V_T, V_N, P, S)$ と表わす。ここで V_T は末端記号の集合, V_N は非末端記号の集合, P はプロダクションの集合, $S \in V_N$ は出発記号である。

定義 1. 確率的文脈自由形文法 (s.c.f.g.) G_s は順序対 (G, φ) である。ここで $G = (V_T, V_N, P, S)$ は既約な c.f.g. であり, φ は P から実数区間 $(0, 1]$ への写像である。ただし $\varphi(P_i)$, $P_i \in P$ を同一の左辺を持つ P_i について加えたものは 1 に等しいとする。 G を G_s の台という。

記法 2. $G = (V_T, V_N, P, S)$ を cfg とするとき $\alpha \in (V_T \cup V_N)^* - V_T^*$ から出発する G による左生成過程の全体を D_α , D_α の元で文生成過程となつてゐるものの全体を E_α と書く. D_α , E_α の元で n ステップのもの全体の全体をそれぞれ $D_\alpha^{(n)}$, $E_\alpha^{(n)}$ と書く. 特に D_φ , E_φ , $D_\varphi^{(n)}$, $E_\varphi^{(n)}$ をそれぞれ D , E , $D^{(n)}$, $E^{(n)}$ と書く.

記法 3. $d \in D_\alpha^{(m)}$ の第 i ステップ ($i=1, \dots, m$) における i 段階の γ ショック P_i が適用されてゐるとき, $d = P_1 \dots P_m$ と書く.

定義 2. $G_\varphi = (G, \varphi)$ を scfg とする. そのとき写像 $\psi_\alpha: D_\alpha \rightarrow (0, 1]$ をつぎのように定義する: $d = P_1 \dots P_m \in D_\alpha$ に対して $\psi_\alpha(d) = \varphi(P_1) \dots \varphi(P_m)$.

定義 3. 写像 $\psi'_\alpha: 2^{D_\alpha} \rightarrow [0, \infty)$ をつぎのように定義する: $M \subset D_\alpha$ に対して $\psi'_\alpha(M) = \sum_{d \in M} \psi_\alpha(d)$, したがって $\psi'_\alpha(\emptyset) = 0$.

ψ'_α は ψ_α の拡張であるが, 両者とも ψ_α で表わす. ψ_φ を ψ と書く. ψ の定義域を E に制限したものを ψ で表わす.

命題 1. $(E, 2^E, \psi)$ は測度空間である.

定義4. E の元は, それによつて生成される文を対応させる写像を μ とする.

記法4. $\mu(E)$ を $L(G)$ と書く.

μ によつて E から $L(G)$ 上にひきおこされる測度を考えることができる. すなわち,

命題2. $\tilde{\psi} : 2^{L(G)} \rightarrow [0, \infty)$ を $M \subset L(G)$ に対し $\tilde{\psi}(M) = \psi(\mu^{-1}(M))$ と定義すると $(L(G), 2^{L(G)}, \tilde{\psi})$ は測度空間である.

$(L(G), 2^{L(G)}, \tilde{\psi})$ を G_{μ} に μ によつて生成される確率的文脈自由形言語 (scfl) とする.

命題3. $\tilde{\psi}(L(G)) = \psi(E)$.

命題4. $G_{\mu} = (G, \mu)$ を scfl とする. G が "あり" だけならば $(E, 2^E, \psi)$ と $(L(G), 2^{L(G)}, \tilde{\psi})$ は測度空間として同型である.

定理 1. 任意の σ -cf \mathcal{G} に対し ψ は成り立つ σ -cf $\mathcal{L}(L(\mathcal{G}), 2^{L(\mathcal{G})}, \tilde{\psi})$ により $0 < \tilde{\psi}(L(\mathcal{G})) \leq 1$ が成立する。

定理 1 を証明するためには二つの補題を用意する。

補題 1. 任意の σ -cf \mathcal{G} に対し $\psi(D^{(n)} - E^{(n)}) = \psi(D^{(n+1)})$ が成り立つ。

補題 2. 任意の σ -cf \mathcal{G} に対し $\psi(D^{(n)} \cup \bigcup_{i=1}^{n-1} E^{(i)}) = 1$ 。

(定理 1 の証明) 命題 3 により $\tilde{\psi}(L(\mathcal{G})) = \psi(E)$ であるから $0 < \psi(E) \leq 1$ を証明する。 $E = \lim_{n \rightarrow \infty} \bigcup_{i=1}^n E^{(i)}$ であり、また $\mathcal{E} = \{\bigcup_{i=1}^n E^{(i)} \mid (n=1, 2, \dots)\}$ は単調非減少な集合列であるから測度の一般論により

$$\psi(E) = \psi\left(\lim_{n \rightarrow \infty} \bigcup_{i=1}^n E^{(i)}\right) = \lim_{n \rightarrow \infty} \psi\left(\bigcup_{i=1}^n E^{(i)}\right).$$

$M_1 \subset M_2 \subset E$ ならば $\psi(M_1) \leq \psi(M_2)$ であるから $\{\psi(\bigcup_{i=1}^n E^{(i)})\}$ ($n=1, 2, \dots$) は単調非減少な数列であり、また補題 2 により $\psi(\bigcup_{i=1}^n E^{(i)}) = 1 - \psi(D^{(n+1)}) \leq 1$ であるから $\psi(\bigcup_{i=1}^n E^{(i)})$ ($n=1, 2, \dots$) は有界である。したがって $\lim_{n \rightarrow \infty} \psi(\bigcup_{i=1}^n E^{(i)})$ が存在する。この値を $\psi(E)$ とすれば、 $0 < \psi(E) \leq 1$ である。

るから $0 < \rho \leq 1$.

(証明終)

例 1. $G = (\{a, b\}, \{S\}, \{S \rightarrow a, S \rightarrow bSS\}, S)$ を台として持つ Δ cf ρ $G_\rho = (G, \rho)$, $\rho(S \rightarrow a) = p_1$, $\rho(S \rightarrow bSS) = p_2$, $p_1 + p_2 = 1$, $p_1 > 0$, $p_2 > 0$ を考える. この G_ρ により生成される Δ cf ρ により $\tilde{\Psi}(L(G))$ を定めると

$$\tilde{\Psi}(L(G)) = \begin{cases} 1; & p_1 \geq \frac{1}{2} \text{ のとき,} \\ \frac{p_1}{1-p_1}; & \frac{1}{2} > p_1 > 0 \text{ のとき.} \end{cases}$$

例 1 は Δ cf ρ が確率空間に存在する場合とそうでない場合のあることを示しているが, Δ cf ρ が確率空間に存在ということの意味についてもう少し考察を加えてみよう. Δ cf ρ を文生成オートマトンと考え, このオートマトンは文が生成されたときに停止するものとする. $D^{(n)} - E^{(n)}$ は n ステップの生成過程のうち文生成過程でないものの全体であり, $\bigcup_{i=1}^n E^{(i)}$ は n ステップ以下の文生成過程の全体であり, $\Psi(D^{(n)} - E^{(n)})$ はこのオートマトンが n ステップ以下では停止しない確率, $\Psi(\bigcup_{i=1}^n E^{(i)})$ は n ステップ以下で停止する確率であり, さで示されたように

$$\tilde{\Psi}(L(G)) = \lim_{n \rightarrow \infty} \Psi\left(\bigcup_{i=1}^n E^{(i)}\right)$$

であり, $\psi(\bigcup_{i=1}^n E^{(i)})$ は n に對して単調非減少でありから,
 ‘ $\tilde{\psi}(L(G)) = 1$ ’ と ‘任意の $\varepsilon > 0$ に対して自然数 n_0 が存
 在して $\psi(\bigcup_{i=1}^{n_0} E^{(i)}) > 1 - \varepsilon$ ’ とは同値であり. すなわち, つ
 ぎの命題が得られた.

命題 5. つぎの (i), (ii) は同値であり.

- (i) $\Delta C f g$ G_{Δ} によつて生成される $\Delta C f g$ が確率空間となり.
- (ii) 任意の $\varepsilon > 0$ に対して自然数 n_0 が存在し, オートマト
 ンとしての G_{Δ} が n_0 ステップ以下で停止する確率が $1 - \varepsilon$ より大きい.

命題 5 の (ii) を直感的にいえば ‘ G_{Δ} が殆んど確実に有限
 ステップで停止する’ と存する. 話し手のモデルとしてはその
 ような G_{Δ} のみを考えれば十分であろうから, 以下の議
 論では $\Delta C f g$ はつねに確率空間でありとする.

§ 4. 言語の情報量

4.1 文生成過程および言語のエントロピー

定義 5. $\Delta C f g$ G_{Δ} による文生成過程の全体 E のエントロ
 ピー $H(E)$ を

$$H(E) = - \sum_{d \in E} \psi(d) \log_2 \psi(d)$$

と定義する。

$H(E)$ は文生成過程 ψ とつ当りの平均情報量, すなわち P. マーカーを一つ受けと, t ときに得られる情報量の期待値である。

定義 6. $G_\delta = (G, \varphi)$ を scfg とする. G_δ により生成される scfl のエントロピー $H(L(G))$ を

$$H(L(G)) = - \sum_{w \in L(G)} \tilde{\psi}(w) \log_2 \tilde{\psi}(w)$$

と定義する。

$H(L(G))$ は G_δ により生成される文 w とつ当りの平均情報量である。

あとでわかるように一般に $H(L(G)) \leq H(E)$ であるから $H(E)$ が有限ならば $H(L(G))$ も有限である。しかし $H(E)$ は常に有限とは限らない。P. マーカーを一つ受けと, t ときに, それから無限の情報を得ることが期待できるということは一見不合理のようであるが, 以下に述べることがこの疑問を解決するであろう。

定義 7. $N: E \rightarrow \mathbb{R}$ をつきのように定義する: $d \in E$ に対し $d \in E^{(n)}$ のとき $N(d) = n$.

N は $(E, 2^E, \psi)$ 上の確率変数であり, $N = n$ とする確率 $\psi\{N = n\}$ は $\psi(E^{(n)})$ で与えられる. N の平均値を $M(N)$ と書く: $M(N) = \int_E N d\psi = \sum_{n=1}^{\infty} n \psi(E^{(n)})$.

定理 2. $H(E)$ は $M(N)$ が有限なとき, かつそのときに限り有限である.

$M(N)$ は文生成過程の平均ステップ数であるから, G_0 を話し手のモデルと考える以上 $M(N)$, したがって $H(E)$ は有限と考えてよいであろう.

4.2 言語のあいまいさ

言語にはあいまいさのある場合がある. これは同一の文を生成する複数個の生成過程が存在する場合で, 文が与えられてもそれがどの生成過程によるものか聞き手が決定できず, あいまいさが残るわけである.

ある acfg によって生成される二つの測度空間 $(E, 2^E, \psi)$, $(L(G), 2^{L(G)}, \tilde{\psi})$ によって考える. $w \in L(G)$ が生成されたと

この条件のもとでの $d \in E$ の確率 $P_w(d)$ は条件付確率の定義により $P_w(d) = \psi(\{d\} \cap h^{-1}(w)) / \psi(h^{-1}(w))$ で与えられる。

定義 8. $w \in L(G)$ のありまゝ度 $A(w)$ を

$$A(w) = - \sum_{d \in E} P_w(d) \log_2(P_w(d))$$

で、また $L(G)$ の平均ありまゝ度 $A(L(G))$ を

$$A(L(G)) = \sum_{w \in L(G)} \tilde{\psi}(w) A(w)$$

で定義する。

命題 6. $A(L(G)) = H(E) - H(L(G))$ 。

命題 7. $A(L(G)) = 0$ と存在するのは G がありまゝでないとき、かつそのときに限り。

謝辞 この研究を進めしに際しては、NHK総合技術研究所内 数理工学研究会の諸氏に討論をしていただいた。とくに例 1 の結果は上坂吉則氏および坂井徹男氏に負うところが大きいことを記して謝意を表したい。