

制御集合にもとづく線型言語の帰納的推論

5P-2

高田裕志

富士通㈱国際情報社会科学研究所

1. 導入

言語はそれを生成する文法を特定することによって同定される。ここでは、線型言語 L を同定する方法を論じる。はじめに、アルファベット Σ 上の任意の線型言語 L が、 Σ 上のある固定された線型文法 G^0 と G^0 の生成規則上の正則制御集合 C から生成されることを示す。この事実から、 G^0 に対する正則制御集合 C を同定することによって、線型言語 L を同定する。言語の構造の情報とともに、 L の所属に関する oracle と L を生成する線型文法 G の補助的な情報が与えられたとき、 L を正則集合の同定と同じ方法で同定できる。

2. 表現定理

Σ を有限のアルファベット、 Σ^* を Σ 上のすべての記号列からなる集合とする。さらに、 λ を空記号列、 $|u|$ を u の長さとする。

$M = \langle K, \Sigma, \delta, q_0, F \rangle$ を(決定性)有限状態受理機とする。ここで、 K は状態の空でない有限集合、 δ は状態遷移関数、 q_0 は初期状態、 F は最終状態の集合である。 Σ^* の部分集合 R を受理する有限状態受理機が存在するとき、 R は正則であるという。 R が正則ならば、 R を受理する状態数が最小の有限状態受理機が同型を除いて唯一存在する。この受理機を R に対する正準(canonical) 有限状態受理機という。

$G = \langle N, \Sigma, \Pi, S \rangle$ を文脈自由文法とする。ここで、 N は非終端記号の、 Π は生成規則の空でない有限集合、 S は初期記号である。 $N \cup \Sigma$ を V で表す。 G のどの生成規則も $S \rightarrow \lambda$ 、 $A \rightarrow a$ 、 $A \rightarrow aB$ または $A \rightarrow Ba$ ($S, A, B \in N, a \in \Sigma$) の形をしているとき、 G を(標準形の) 線型文法という。特に、ただ 1 つの非終端記号 S だけからなるとき、 G は最小であるという。各生成規則にはラベルがつけられているものとし、それを π で示す。また、 G は無用な非終端記号を含まないと仮定する。

【定義】 Σ を $\{a_1, a_2, \dots, a_n\}$ とする。このとき、生成規則の集合 Π^0 が

$$\begin{aligned} & (S^0 \rightarrow \lambda, S^0 \rightarrow a_1 S^0, S^0 \rightarrow a_2 S^0, \dots, S^0 \rightarrow a_n S^0, \\ & S^0 \rightarrow S^0 a_1, S^0 \rightarrow S^0 a_2, \dots, S^0 \rightarrow S^0 a_n, \\ & S^0 \rightarrow a_1, S^0 \rightarrow a_2, \dots, S^0 \rightarrow a_n,) \end{aligned}$$

である最小線型文法 $G^0 = \langle S^0, \Sigma, \Pi^0, S^0 \rangle$ を万能(universal) であるという。

アルファベット Σ に対して、万能最小線型文法 G^0 は唯一に存在する。

G を線型文法、 $x_0 \xrightarrow{\alpha} x_1 \xrightarrow{\beta} \dots \xrightarrow{\gamma} x_n$ を G における導出(derivation)とする。ここで、 x_{i-1} から x_i ($1 \leq i \leq n$) への遷移において生成規則 π_i が適用されるものとする。 $x, y \in V^*, \alpha \in \Pi^*$ に対して、 $x \xrightarrow{\alpha} y$ は G において導出 $x = x_0 \xrightarrow{\alpha_1} x_1 \xrightarrow{\alpha_2} \dots \xrightarrow{\alpha_n} x_n = y$ が存在することを意味する。ここで、 $\alpha = \pi_1 \pi_2 \dots \pi_n$ である。 α を導出 $x \xrightarrow{\alpha} y$ の付随語(associate word)という。 G によって生成された言語 $L(G)$ とは集合

$$L(G) = \{ w \in \Sigma^* \mid S \xrightarrow{\alpha} w, \alpha \in \Pi^* \}$$

である。 $L=L(G)$ となる線型文法 G が存在するならば言語 L は線型であるという。

【定義】 線型文法 $G = \langle N, \Sigma, \Pi, S \rangle$ が次の 2 つの条件を満たすとき、正準であるという：

1. $A \neq B$ である N の任意の要素 A, B に対して、 $A \xrightarrow{\alpha} w$ であるが $B \not\xrightarrow{\alpha} w$ でない記号列 $w \in L(G)$ が存在する。

2. $A \rightarrow aB, A \rightarrow aC$ または $A \rightarrow Ba, A \rightarrow Ca$ ($B \neq C$) である生成規則のペアは存在しない。

特に、条件 2 を満たすとき G は決定的であるという。

【定義】 C を Π^* の任意の正則な部分集合とする。このとき、

$$L_C(G) = \{ w \in \Sigma^* \mid S \xrightarrow{\alpha} w, \alpha \in C \}$$

を正則制御集合 C と G によって生成された言語という。

【定理 1】 任意の線型言語 L に対して、 $L=L_C(G^0)$ となる正則制御集合 C が存在する。□

証明は省略する(文献[3]を参照、以下の定理、命題についても同様である)。方針としては、 G^0 に対する正則制御集合 C を受理する有限状態受理機の各状態を、 L を生成する線型文法 G の非終端記号に対応させる。定理 1 の逆の場合も成り立つ。

【定理 2】 G^0 をアルファベット Σ 上の万能最小線型文法、 C を G^0 に対する正則制御集合とする。そのとき、 $L=L_C(G^0)$ は線型言語である。□

Inductive Inference of Linear Languages based on Control Sets

Yuji TAKADA (yuji@iias.iias.fujitsu.junet)

International Institute for Advanced Study of Social Information Science(IIAS-SIS), FUJITSU LIMITED

定理1と2から、線型文法Gが与えられたとき、 $L(G) = L_c(G^0)$ となるCを受理する有限状態受理機Mを対応させることができる。

《命題1》 任意の線型文法Gに対して、 $L(G)=L(G')$ となる決定的線型文法G'が存在する。□

《系》 任意の線型言語Lに対して、 $L=L(G)$ となる正準線型文法Gが存在する。□

任意の線型言語Lに対してしを生成する正準線型文法Gが存在する。しかし、その文法は唯一ではないので有限状態受理機における「正準」と線型文法における「正準」は必ずしも同義ではない。

《命題2》 Gを線型文法、MをC-T(M)かつ $L(G)=L_c(G^0)$ であるGに対応する有限状態受理機とする。このとき、Gが正準ならばMは正準である。□

一般に、逆の場合は成り立たない。

このように、線型言語Lを同定するには、しを生成する正準線型文法Gに対応する正準有限状態受理機Mを同定すればよい。

3. 線型言語の同定アルゴリズム

同定アルゴリズムは言語の構造に関する情報が利用可能であるとする。言語の構造はかっこ文法(parenthesis grammar)によって記述することができる。

【定義】 線型文法 $G=\langle N, \Sigma, \Pi, S \rangle$ のかっこ文法 $\langle G \rangle$ は Π のすべての生成規則 $A \rightarrow x$ ($x \in V^*$) を $A \rightarrow \langle x \rangle$ で置き換えることによって得られる。ここで、"⟨", "⟩"は Σ の要素でない特別の記号である。

W を記号列の集合とする。このとき、 $Pr(W)$ は W の要素のすべてのprefixの集合である。 $Pr(W)$ は空でないとき、 λ を含む。 X と Y が集合であるとき、 $X \oplus Y$ で X と Y の対称差を表す。

【定義】 Gを正準線型文法、 $\langle G \rangle=\langle N, \Sigma, \Pi, S \rangle$ をGのかっこ文法、 G^0 を Σ 上の万能最小線型文法、 $\langle G^0 \rangle=\langle \{S^0\}, \Sigma, \Pi^0, S^0 \rangle$ を G^0 のかっこ文法とする。 $L(\langle G \rangle)$ の表現標本(representative sample) Raとは $\langle G \rangle$ のすべての生成規則 π に対して π が導出 $S \xrightarrow{\alpha} w$ に現れるRaの要素 w が存在する $L(\langle G \rangle)$ の有限部分集合である。Raに関するCの付随表現標本Rcとは集合

$$Rc = \{\alpha \mid S^0 \xrightarrow{\alpha} w, w \in Ra\}$$

である。

言語の構造を与えることによって導出を特定することができます。したがって、 $L(\langle G \rangle)$ の要素 w に対して、 w の $\langle G^0 \rangle$ における導出を表す付隨語 α が一対一に対応する。

線型言語同定アルゴリズムLIDは表現標本Raと $L(\langle G \rangle)$ の所属に関するoracleから有限状態受理機Mを同定する。LIDはAngluin [1]によるアルゴリズムIDを線型言語版に修正したものである。

アルゴリズムLIDにおいて、Pcの各要素は可能な状態を表す。 d_0 は「落とし穴」状態を表す。関数fは任意の $\pi^0 \in \Pi^0$ に対して $f(d_0, \pi^0) = d_0$ 、各 $\mu \in Pc$ に対して $f(\mu, \pi^0) = \mu \pi^0$ と定義される。fはMの遷移を表す。関数gは付隨語 α に対して $S^0 \xrightarrow{\alpha} \langle \dots \rangle, s$ となる V^* の要素sを対応させる。LIDはPcの要素を同値類に分割する。その同値類の各要素がMの状態となる。

アルゴリズム LID

入力であるRaから次の集合を構成する

$$Rc = \{\alpha \mid S^0 \xrightarrow{\alpha} \langle \dots \rangle, w, w \in Ra\}$$

$$Pc = Pr(Rc)$$

$$Pc' = Pc \cup \{d_0\}$$

$$Tc = Tc' \cup \{f(\mu, \pi^0) \mid \mu \in Pc, \pi^0 \in \Pi^0\}$$

$$Tc = Tc' - \{d_0\}$$

ステップ 0

すべての $\mu \in Tc'$ に対して $E_0(\mu)$ を空にする
 i を1、 v_i を入にする

ステップ i

$E_i(d_0)$ を空にする

各 $\mu \in Tc$ に対して記号列 $g(\mu v_i)$ を $L(\langle G \rangle)$ のoracleに質問し、 $E_i(\mu)$ を構成する

$g(\mu v_i) \in L(\langle G \rangle)$ ならば $E_i(\mu)$ を $E_{i+1}(\mu) \cup$

$\{v_i\}$ とし、そうでなければ $E_{i+1}(\mu)$ とする

$E_i(\mu) = E_i(\mu')$ だが $E_i(f(\mu, \pi^0)) = E_i(f(\mu', \pi^0))$ となるペア $\mu, \mu' (\in Pc)$ と $\pi^0 \in \Pi^0$ がみつかったならば、 $E_i(f(\mu, \pi^0)) \oplus E_i(f(\mu', \pi^0))$ からある付隨語 ν を選び v_{i+1} を $\pi^0 \nu$ としてステップ $i+1$ に行く。そうでなければ $i=i$ として有限状態受理機Mを構成する

$M = \langle K, \Pi^0, \delta, q_0, F \rangle$ において、 $K = \{E_n(\mu) \mid \mu \in Tc\}$ 、 $q_0 = E_n(\lambda)$ 、 $F = \{E_n(\mu) \mid E_n(\mu) \ni \lambda\}$ であり、 δ は次のように定義される： $E_n(\mu)$ が空ならば、すべての $\pi^0 \in \Pi^0$ に対して $E_n(\mu)$ 自身に遷移する。そうでないならば、 $E_n(\mu) = E_n(\mu')$ となる $\mu' \in Pc$ が存在して、 π^0 に対して状態 $E_n(\mu' \pi^0)$ に遷移する。

アルゴリズムLIDはアルファベット Σ の個数とPcの要素の個数の多項式個の質問でMを同定する。LIDの正当性は文献[3]を参照。

謝辞 国際研・学習システム研究グループの諸氏に感謝する。なお、本研究は第五世代コンピュータプロジェクトの一環として行ったものである。

参考文献

- [1] Angluin, D. (1981), "A Note on the number of Queries Needed to Identify Regular Languages", Information and Control, 51, 76-87.
- [2] Ginsburg, S. and Spanier, E.H. (1968), "Control Sets on Grammars", Math. Systems Theory, 2, 159-177.
- [3] Takada, Y. (1987), "A Constructive Method of Grammatical Inference for Linear Languages based on Control Sets", in preparation.