

線形文法の推定

北九州大 商 桑次奎介

任意の線形言語 L に対し、 L に属するか否かだけが知られていない語のある有限系列を入力として $L = L(G)$ なる線形文法 G を推定するための手続きを示す。この方法は Biermann が提示した正則言語に対する学習方法⁽¹⁾ の自然な一般化である。

§1. 言語 L からの推定

アルファベット Σ 上の言語 L に対して、 $u, v \in \Sigma^*$ による L の導言語を $\{x; uxv \in L\}$ とし、 $\bar{u}L\bar{v}$ で表わす。

L に対していつも構造図 $D(L) = (N, f)$ が一意に対応する⁽³⁾。
 $N = N$ は節の集合で $\{\bar{u}L\bar{v}; u, v \in \Sigma^*\}$ を表わす。 f は節から節へのラベルつき矢印をとる写像で、 $a \in \Sigma$ に対し

$$f(X, Y) = \begin{cases} a & (Y = \bar{a}X \text{ のとき}) \\ \cdot a & (Y = X\bar{a} \text{ のとき}) \\ \text{定義されない} & (\text{その他のとき}) \end{cases}$$

と定める。 L が正則であれば $D(L)$ は L を受理する有限オートマトンをを与えることになる。そうでなければ $D(L)$ は無限構造図となるが、この場合でも L の文法的構造を充分に内包しており、それを

どうとり出すかが、ここでの中心的課題となる。

いま長さ以下の L の語全体を (L) で表わす。 $\Sigma = \{a_1, a_2, \dots, a_n\}$ として L の深さ k の構造木 $T(L, k)$ を次のように帰納的に構成する：

ステップ 0：根節 $(\varepsilon, \varepsilon)$ に $(L)_k$ をラベルせよ。

ステップ n ：現時点で得られている木において、まだチェックされていない葉のうち、最も早く生成された節 (u, v) を選べ。
同じラベルをもつ他のいかなる節からもまだ枝が張られていないなら、節 (u, v) の下に左から右へ順に節 $(ua_1, v), (u, a_1v), (ua_2, v), (u, a_2v), \dots, (ua_n, v), (u, a_nv)$ を並べ、節 (u, v) とそれぞれラベル $a_1, \cdot a_1, a_2, \cdot a_2, \dots, a_n, \cdot a_n$ をもつ枝で結べ。
新しい各節 (w, v') には $(\bar{w}L\bar{v})_k$ をラベルせよ。節 (u, v) と同じラベルをもつ他の節から既に枝が張られているなら節 (u, v) をチェックせよ。

$(\Sigma^*)_k$ が有限なので上の手続きはいつでも有限ステップで終了する。

いま、 $T(L, k)$ における節の全体を $T^*(L, k)$ とし、 $\{\bar{w}L\bar{v}; (u, v) \in T^*(L, k)\}$ を $N(L, k)$ で表わす。 $k \leq n$ ここで、 $L \subset \Sigma^*$ の深さ (k, k) の構造図 $D(L, k, k) = (N, f)$ を次のように定める：

$$N = \{(X)_k; X \in N(L, k)\}$$

$f(X, Y) = \alpha \Leftrightarrow T(L, k)$ において $(X)_k$ とラベルされた節から $(Y)_k$ とラベルされた節へラベル α をもつ枝が張られている。

この構造図 $D(L, \mathbf{f}, \mathbf{h}) = (N, f)$ から次のような線形文法

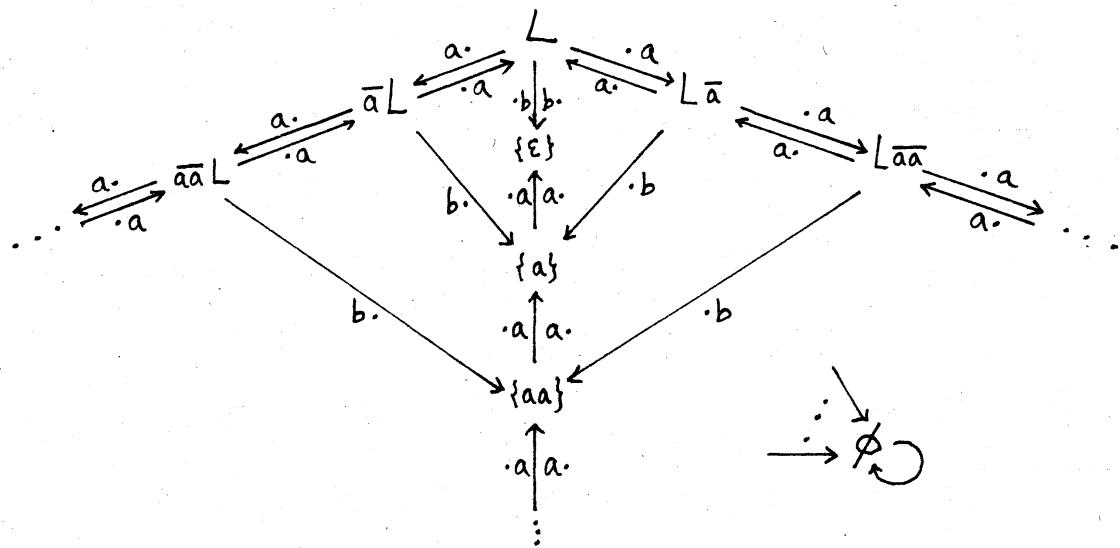
$G(L, \mathbf{f}, \mathbf{h}) = (I, \Sigma, P, S)$ を対応させることができる:

$I = N - \{\phi\}$, Σ は L のアルファベット, $S = (L)_{\mathbf{h}}$, P は
 $X, Y \in I$, $a \in \Sigma$ として $f(X, Y) = a \cdot + \varepsilon$ なら $(X \rightarrow aY) \in P$, $f(X, Y) = -a$
 なら $(X \rightarrow Ya) \in P$, さらに $\varepsilon \in X$ なら $(X \rightarrow \varepsilon) \in P$ とする。

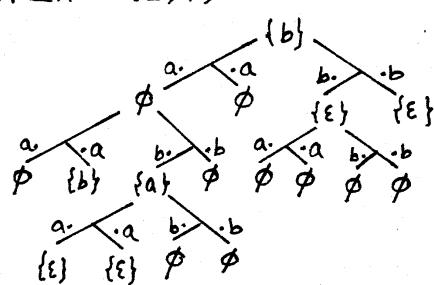
$L(G(L, \mathbf{f}, \mathbf{h}))$ を単に $L(L, \mathbf{f}, \mathbf{h})$ と書く。以上の手続きによって、
 与えられた L から一つの線形文法を得る。

例 1. $L = \{a^n b a^n ; n \geq 0\}$ とする。

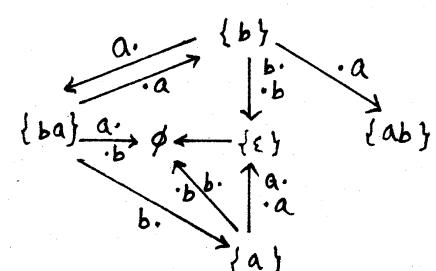
(1) 構造図 $D(L)$



(2) 構造木 $T(L, 1)$



(3) 構造図 $D(L, 1, 2)$



(4) 文法 $G(L, I, 2) = (I, \Sigma, P, S)$: $I = \{ \{b\}, \{ba\}, \{ab\}, \{\epsilon\}, \{a\} \}$, $\Sigma = \{a, b\}$, $S = \{b\}$, $P = \{ \{b\} \rightarrow a\{ba\} / \{ab\}a / b\{\epsilon\} / \{\epsilon\}b, \{ba\} \rightarrow \{b\}a / b\{a\}, \{\epsilon\} \rightarrow \epsilon, \{a\} \rightarrow a\{\epsilon\} / \{\epsilon\}a \}$.

P を簡単化すると P' と等価な規則 $P' = \{S \rightarrow aSa / b\}$ を得る。

明らかに $L(L, I, 2) = L$ である。

§2. 手続きの完全性

任意の線形言語 L に対して適当な非負整数の組 (k, ℓ) を選べば $L(L, k, \ell) = L$ とできることを示す。

いま線形文法 $G = (I, \Sigma, P, S)$ に対し, $I \cup \Sigma = V$ として

$T_G = \{(u, v); \text{導出 } S \xrightarrow{G} x \text{ には各非端末記号は } \tau \text{ か } \tau' \text{ が一度しか現われない, } x \xrightarrow{G} uyv, x, y \in V^*, u, v \in \Sigma^*\}$

とし, 言語 L の次数を

$$d(L) = \min_{L = L(G)} (\min \{k; T^*(L, k) \supseteq T_G\})$$

と定める。また

$$d(L, k) = \max_{\substack{L_1, L_2 \in N(L, k) \\ L_1 \neq L_2}} (\min \{z; (L_1)_z \neq (L_2)_z\})$$

とおく。

命題. $L \subset \Sigma^*$, $X, Y \in N(L, k)$ とするとき, $G(L, k, k) \vdash$ において $(X)_k \nRightarrow u(Y)_k v$ なる導出が存在するための必要十分条件は $Y = \bar{u} X \bar{v}$ である。 $z = i = k = d(L, k)$ とする。

証明. $D(L, k, k) = (N, f)$, $G(L, k, k) = (I, \Sigma, P, S)$

とする。 $N(L, k)$ の要素は互いにたかだかその長さで異なるので
 $X, Y \in N(L, k)$ として

$$(X)_k \rightarrow a(Y)_k \in P \Leftrightarrow f((X)_k, (Y)_k) = a \Leftrightarrow Y = \bar{a}X$$

$$(X)_k \rightarrow (Y)_k a \in P \Leftrightarrow f((X)_k, (Y)_k) = a \Leftrightarrow Y = X\bar{a}$$

が成り立つ。(かも)

$$(X)_k = (Y)_k \Leftrightarrow X = Y$$

T なので $(X)_k \xrightarrow{G} u(Y)_k v$ なる導出が存在する = $\exists Y = \bar{u}X\bar{v}$

であることが同様にわかる。□

補題 1. 任意の $L \subset \Sigma^*$ に対し, $L(L, k, d(L, k)) \subset L$.

証明. $d(L, k) = h$, $G(L, k, h) = (I, \Sigma, P, S)$ とおく。

$w \in L(L, k, h)$ とすると命題から

$$S = (L)_h \xrightarrow{*} u(\bar{u}L\bar{v})_h v \Rightarrow uv = w$$

なる導出が存在する。導出の最後には生成規則 $(\bar{u}L\bar{v})_h \rightarrow \epsilon$ が適用されており、 $\epsilon \in \bar{u}L\bar{v}$ となる。ゆえに $w = uv \in L$. □

補題 2. G をある線形文法とし, $L(G) = L$ とおく。

$T^*(L, k) \supset T_G$ ならば $L(L, k, d(L, k)) = L$ である。

証明. $G = (I, \Sigma, P, S)$ とする。 $(X \rightarrow uxv) \in P$ ($u, v \in \Sigma^*, x \in V^*$) に対し、各非端末記号がたかだか一度しか現われないような適当な導出 $S \xrightarrow{*} u'Xv'$ ($u', v' \in \Sigma^*$) が存在し、 $u'Xv' \xrightarrow{G} u'u\bar{x}v\bar{v}'$ となるので T_G の定義から (u', v') , $(u'u, v\bar{v}') \in T_G \subset T^*(L, k)$. (したがって $\bar{u}L\bar{v}, \bar{u}u\bar{L}\bar{v}'$

$\in N(L, k)$ である。そのとき命題 1 によると

$$(\bar{u}L\bar{v})_k \xrightarrow{*} u(\bar{u}\bar{u}L\bar{u}\bar{v})_k v$$

なる導出が $G(L, k, h)$ において存在する。ここで $h = d(L, k)$ とする。また、 $w \in \Sigma^*$ で $(X \rightarrow w) \in P$ に対しては、上の生成規則で $x = \epsilon$, $uv = w$ とした場合に相当し、 $\epsilon \in (\bar{u}\bar{u}L\bar{u}\bar{v})_k v$ なので $G(L, k, h)$ において

$$(\bar{u}L\bar{v})_k \xrightarrow{*} u(\bar{u}\bar{u}L\bar{u}\bar{v})_k v \xrightarrow{*} uv = w$$

なる導出が存在することになる。

以上から、 G によって生成される言語は $G(L, k, h)$ によることも生成される。一方、補題 1 から $L \supseteq L(L, k, h)$ である。 \square

補題 2 と $d(L)$ の定義からただちに次の定理を得る。

定理 1. 任意の線形言語 L に対して

$$L(L, d(L), d(L, d(L))) = L$$

§3. サンゴルからの推定

$L \subset \Sigma^*$ から $L = L(G)$ なる文法 G を推定する場合、 L 全体を入力とすることは非現実的である。実は L の適当な有限部分集合によって G を推定することが可能である。ここではそのような部分集合の条件とそれからの推定手続きを示す。

いま L の特徴集合を

$$C(L) = \{ w \in u(\bar{u}L\bar{v})_{d(L, d(L))} v ; (u, v) \in T^*(L, d(L)) \}$$

とする。 $C(L)$ は L の有限部分集合である。

定理2. L を線形言語とする。 $C(L) \subset C \subset L$ のとき、
適当な非負整数 k, h が存在して $L(C, k, h) = L$ とである。

証明。 $k = d(L)$, $h = d(L, d(L))$ とする。 $(u, v) \in T^*(L, k)$
に対して $C(L)$ の定義から $u(\bar{u}L\bar{v})_k v \subset C(L) \subset C$ を得る。
ここで $(\bar{u}L\bar{v})_k \subset \bar{u}(C\bar{v})$ となり、 $(\bar{u}L\bar{v})_k \subset (\bar{u}C\bar{v})_k$ が導か
れる。一方、 $C \subset L$ だから $(\bar{u}C\bar{v})_k \subset (\bar{u}L\bar{v})_k$ 。したがって
 $(\bar{u}L\bar{v})_k = (\bar{u}C\bar{v})_k$ を得る。 $k \leq h$ だから $(\bar{u}L\bar{v})_k = (\bar{u}C\bar{v})_k$
となって $T(L, k) \cong T(C, k)$ は全く一致し、 $D(L, k, h) = D(C, k, h)$
 $= D(C, k, h)$ すなわち

$$L(C, k, h) = L(L, k, h) = L. \quad \square$$

上の定理の応用として 線形文法 G を推定するためのアルゴリズムを図に示す。これは言語 L に属するか否かは知られていないものとして 言語の系列 w_1, w_2, \dots, w_m を次々に入力とし、 L の文法を逐次的に推定していく方法である。 I_+ は L に属する言語の集合を表すし、直接推定に用いられる。 I_- は L に属しない言語の集合を表すし推定された文法が妥当かというかの判断にのみ用いられる。

例2. $L = \{a^nba^n; n \geq 0\}$ とする。 L の特徴集合は $C(L) = \{b, aba\}$ である。 L のサンプル入力を a, b, ab, aba を二つ順に与えると、各段階で次のような文法を推定する：

(1) $I_+ = \{b\}, I_- = \{ab\}; G = G(I_+, 0, 1), L(G) = \{b\},$

$$I_+ \subset L(G), \quad I_- \cap L(G) = \emptyset$$

(2) $I_+ = \{b, aba\}, \quad I_- = \{a, ab\}$:

- ① $G = G(I_+, 0, 1)$, $L(G) = \{b\}$, $I_+ \notin L(G)$
- ② $G = G(I_+, 1, 1)$, $L(G) = \{b\}$, $I_+ \notin L(G)$
- ③ $G = G(I_+, 0, 2)$, $L(G) = \{b\}$, $I_+ \notin L(G)$
- ④ $G = G(I_+, 1, 2)$, $L(G) = L$, $I_+ \subset L(G)$, $I_- \cap L(G) = \emptyset$.

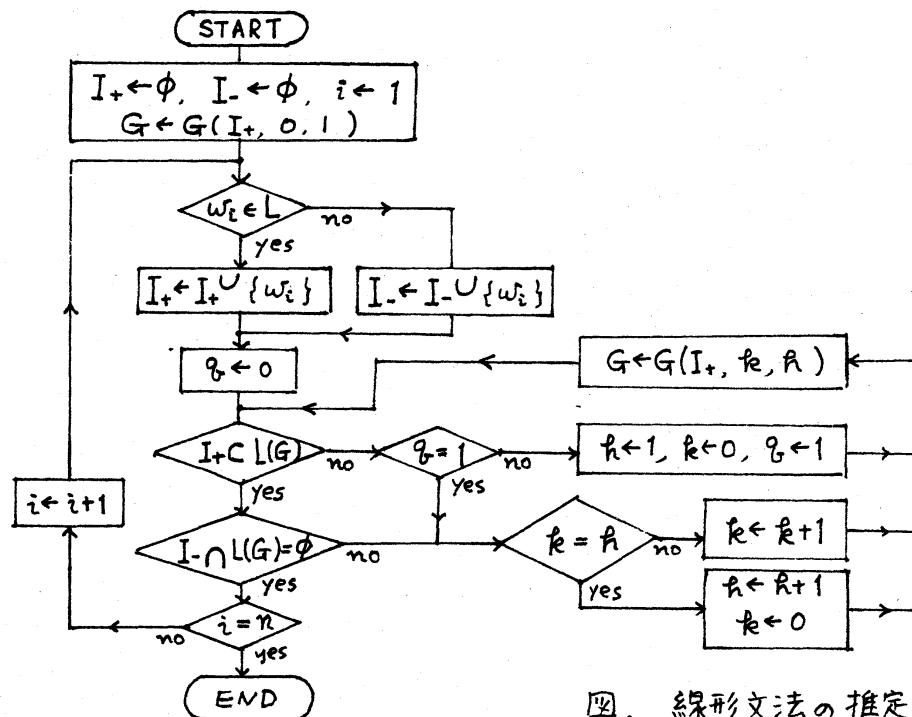


図. 線形文法の推定手続き

参考文献

- [1] A.W. Biermann, An Interactive Finite-State Language Learner, Proc. 1-st USA-JAPAN Comput. Conf. (1972).
- [2] K. Tanatsugu and S. Arikawa, On Characteristic Sets and Degrees of Finite Automata, Int. Jour. of Comput. and Inform. Sci., Vol. 6, No. 1 (1977).
- [3] 飯盛末夫・藤野精一, 自由言語の構造図を利用したフックシュタウン・オートマトンの作成について, 佐賀大学教育学部研究論文集, 第24集(II) (1976).