

模式识别中的递归结构及其句法-词义描述

路浩如

(浙江大学)

摘 要

模式的递归结构影响它的形式语言性质,成为模式结构描述的复杂问题之一。本文系统地研究了递归结构的性质,按基本递归结构划分子模式,并以句法-词义方法对它们作分层文法描述,从而建立了递归结构的有效描述体系。程序文法、属性文法以及由此派生的递归条件文法和递归属性文法,都能成功地按此体系描述各种递归结构,这些文法的描述能力也因而得到阐明。

一、引 言

句法模式识别是将一类模式作为一类形式语言而用文法加以描述。模式结构往往具有上下文有关语言的特点,分析颇为困难。即使是上下文无关语言,其文法推断问题也没有很好解决。因而人们常常采用描述能力最低的有限状态文法来分析模式的有限样本集。

近年来运用句法-词义方法的研究颇为活跃。词义信息的引入可使句法描述大大简化,从而避免了采用上下文有关文法。句法-词义方法的高度描述能力在模式识别的研究中正在受到重视^[1-3]。

模式结构复杂性的主要因素之一是它包含各种递归结构。递归结构的形式决定着模式的形式语言性质。故而研究递归结构特点及其文法描述有其十分必要的意义。

二、递归结构与文法描述

定义 1. 对于串语言(串模式) $\{x_1, x_2, \dots, x_i\}$, 设 $x_i \in \{\beta_i^{m_i} | \beta_i \in V_T^*, m_i \text{ 为变量或函数}\}$, 定义 $\{\beta_i^{m_i}\}$ 为基本递归结构; β_i 为递归子串; m_i 为递归参数。

定义 2. 根据递归子串和递归参数的内涵外连情况,基本递归结构有如下类型:

- 1) 若 $m_i = m_i(n)$, 为函数型, 否则为简单变量型;
- 2) 按 $m_i(n)$ 为线性或非线性, 而有线性型或非线性型 (简单变量型是线性型的特例);
- 3) 若 $m_i = m_i(m_j)$ 或 $m_i = m_i(n)$, $m_j = m_j(n)$, 则同一结构中的 $\{\beta_i^{m_i}\}$ 和 $\{\beta_j^{m_j}\}$

相互依赖,为相关型,否则为独立型;

4) 按 $m_i(m_j)$ 为线性或非线性,而有线性相关型或非线性相关型;

5) 若 $\{\beta_i^{m_i}\}$ 中 β_i 又含有其它基本递归结构,则成为嵌套型.

定义 3. 根据串语言是否含有非线性基本递归结构及相关型基本递归结构,定义它们的分类: 1)非递归语言; 2)线性独立递归语言; 3)线性相关递归语言; 4)非线性独立递归语言; 5)非线性相关递归语言.

对照 Chomsky 的形式语言分类,可发现以上类别 1) 和 2) 都是有限状态语言; 类别 3)包罗全部上下文无关语言(如 $\{a^n b^m c^m\}$ 及 $\{a^m b^m c^n d^n\}$)和一部分上下文有关语言(如 $\{a^n b^n c^n\}$ 及 $\{a^m b^n c^m d^n\}$);类别 4)及 5)则不含任何有限状态和上下文无关语言.

定义 4. 语言的递归复杂度定义为

$$C = [C_1, C_2, C_3]^T. \quad (1)$$

C_{1-3} 各为语言所含线性相关型、非线性独立型和非线性相关型基本递归结构的数目.

由以上分析可以从模式递归结构的特点确定其语言性质,而无须特别注意它们按 Chomsky 的分类. 再者,句法-词义文法是唯一能覆盖各类串语言的描述手段(图 1),而递归结构的性质正是确定模式结构的句法-词义描述的重要因素.

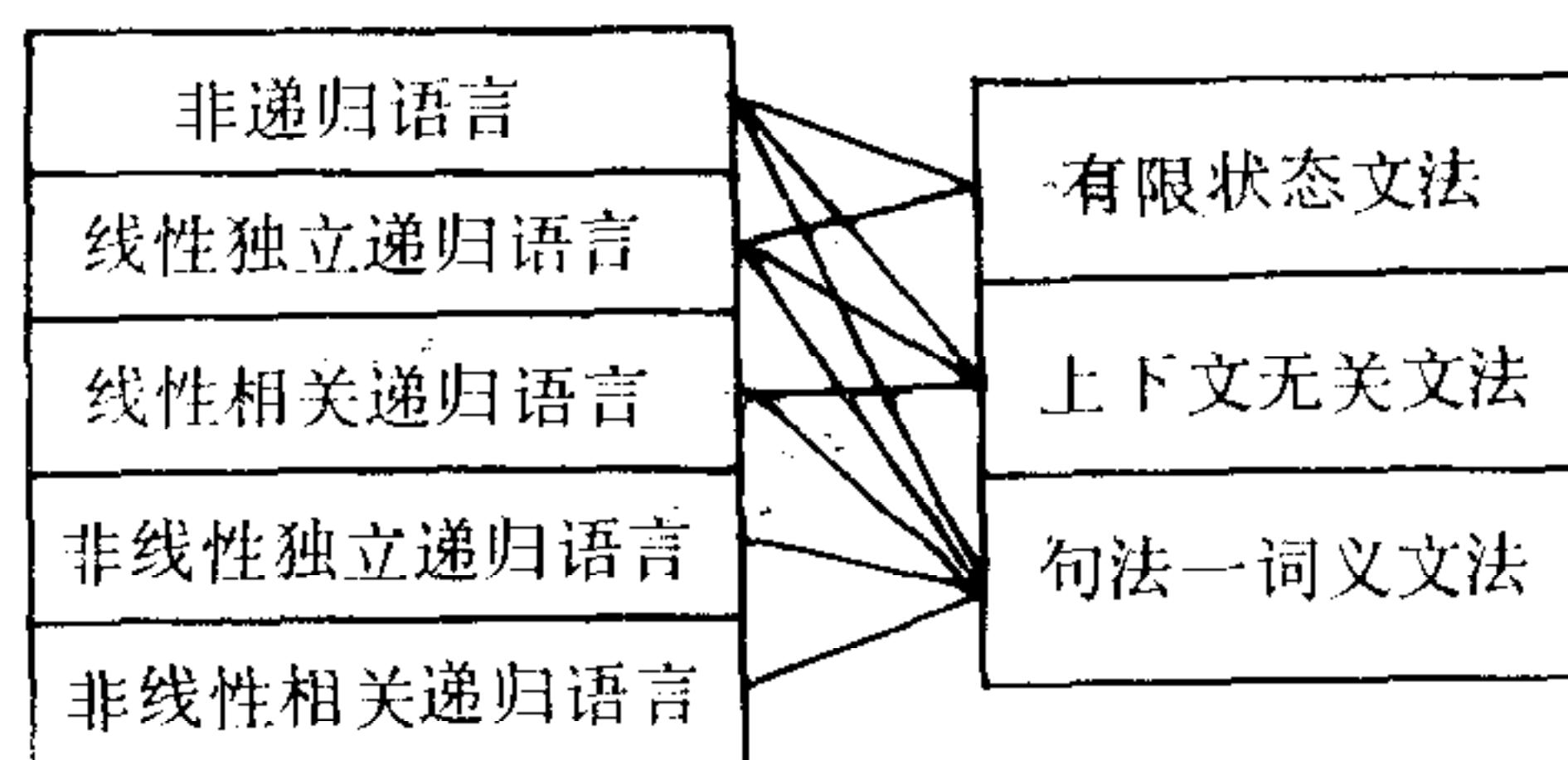


图 1 递归结构与文法描述

三、递归结构的句法-词义描述

定义 5. 模式的句法-词义描述由句法描述和词义描述两要素组成. 句法描述是各结构基元 $a \in V_T$ 在串模式中前后接续的形式语言描述; 词义描述则是基元 a 或子模式 $\gamma \in V_T^+$ 的内涵和外连特征以及相互制约关系的描述. 在句法-词义描述中生成式集 P 有如下基本形式:

$$P = \{(R_{syn}, R_{sem})_i | i = 1, 2, \dots, q\}. \quad (2)$$

R_{syn} 和 R_{sem} 各为生成式 i 的句法和词义部分.

定义 6. 模式的分层结构是由模式整体按适当方式分割成子模式,子模式又分割成更小的子模式, ..., 直至基元. 相应地,模式的分层描述是按分层结构以主干文法描述子模式所构成的模式整体,而以各个子文法描述相应的各个子模式,它们成树状结构(图 2).

本文推荐以分层文法描述作为模式句法-词义描述的规范形式,以简化分析和推断.

定义 7. 模式递归结构的句法-词义描述以基本递归结构作为子模式. 嵌套型基本递归结构的存 在将形成多级子模式分层结构. 此类子模式的词义描述之中须具备对递归性

质的制约,其实现途径有三: 1)限定运用生成式的次序; 2)限制运用生成式的次数; 3)直接限定递归参数.

途径 1) 是一种特殊方式, 2)和 3) 则有更普遍意义. 对后两种方式有如下两定理:

定理 1. 对于递归串语言 L 及 L' , 设有
 $L = \{x_1 x_2 \cdots x_t\}$, $L' = \{x'_1 x'_2 \cdots x'_t\}$,
 $x_i = \{\beta_i^{m_i}\}$, $x'_i = \{\beta_i^{m'_i}\}$, $i = 1, 2, \cdots, t$,
 $C_{L'} = [0, 0, 0]^T$, (即 L' 线性独立),
 $L \subset L'$.

若 L' 的单纯句法描述文法为 $G'(L')$, L 的句法-词义描述为 $G(L)$, 则有

$$\exists G(L) \{R_{\text{syn}}\}_{G(L)} = P_{G'(L')}. \quad (3)$$

证. 从 $L' \supset L$, 有 $x'_i \supseteq x_i$, 且 $m_i = m_i(m'_i)$. 故以 $P_{G'(L')}$ 作为 $G(L)$ 的句法部分, 再运用递归词义作为制约, 必为 L 描述之一途. 由是易证为真.

定理 2. 对于递归串语言 L 的句法-词义描述, 句法部分单独生成有限状态语言 L_{fs} 的情况与单独生成上下文无关语言 L_{cf} 的情况相比较, 前者需要较多的递归词义制约.

证. L , L_{fs} 及 L_{cf} 三者的递归复杂度各为

$$C = [C_1, C_2, C_3]^T, C_{fs} = [0, 0, 0]^T, C_{cf} = [C'_1, 0, 0]^T,$$

且有 $C_1 \geq C'_1 \geq 0$. 对定理中所指两种情况, 前者须为全部 $C_1 + C_2 + C_3$ 个基本递归结构作出递归词义制约, 后者则只对 $C_1 - C'_1 + C_2 + C_3$ 个作出递归制约. 故本定理得证.

本文推荐以单独生成线性独立递归结构的句法描述, 作为递归结构句法-词义描述中句法部分的规范形式, 使子文法描述形式一致, 便于推断和应用.

四、递归结构的程序文法描述

程序文法^[1,4-6]能限制生成式的运用次序, 本文将其列为句法-词义方法的一种特定形式.

定义 8. 程序文法 G 为五元式

$$G = (V_N, V_T, J, P, S). \quad (4)$$

V_N, V_T, J 及 $S \in V_N$ 各为终止符集、非终止符集、标号集及起始符. P 中生成式有表 1 形式. 其运用次序从 $r = 1$ 的生成式开始, 若运用成功, 转向 U 中任一标号之生成式, 若运用失败, 则转向 W 中之一. 如是, 直至去向为一空集 ϕ .

表 1

标 号	核 心	成功去向	失败去向
$r \in J$	$A \rightarrow \rho$	$U \subset J$	$W \subset J$

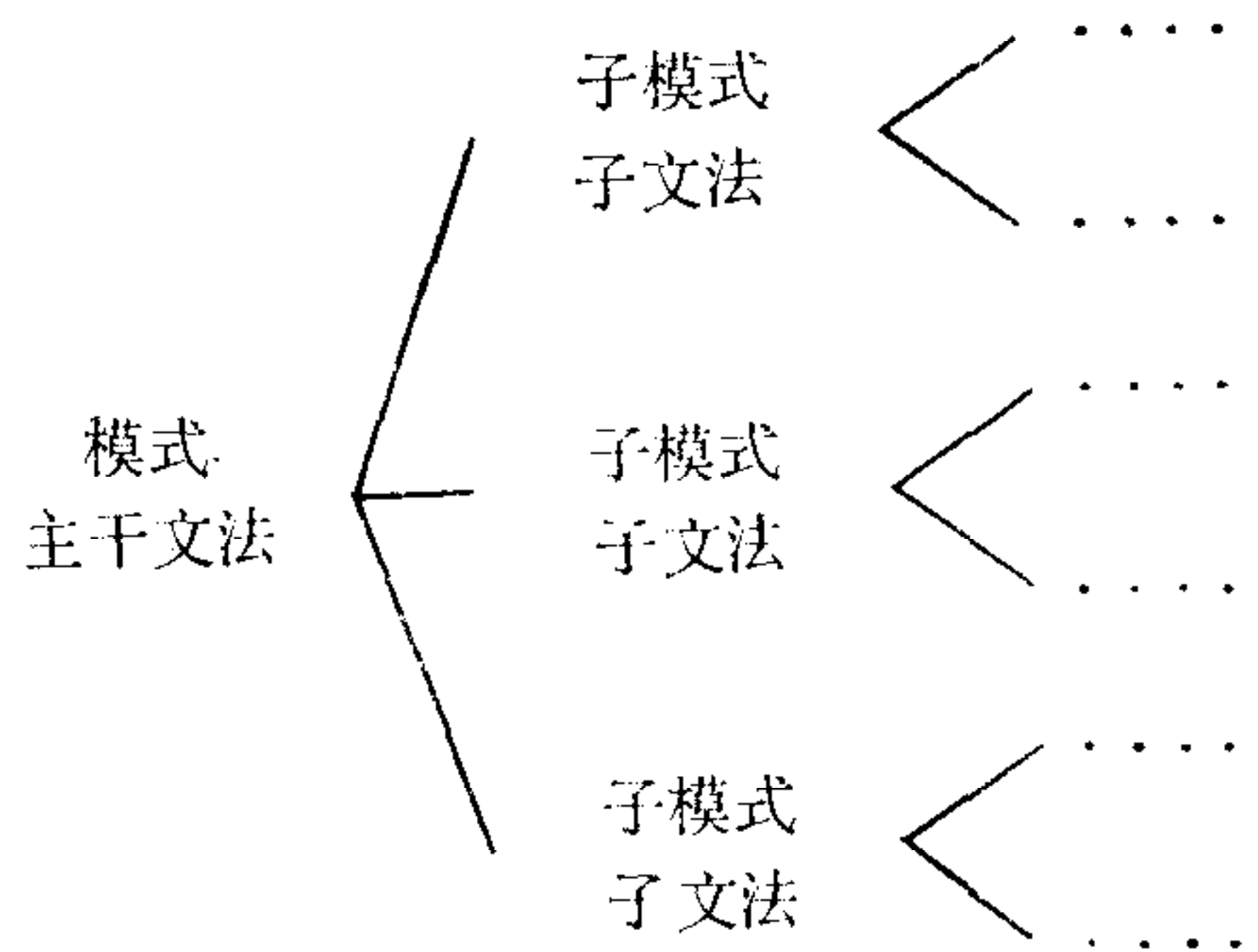


图 2 模式分层结构及分层描述

若核心有 $A \in V_N, \rho \in V^+$, 称上下文无关程序文法(本文简称为程序文法)。它能描述全部上下文无关语言和一大类上下文有关语言。

定义 9. 基本递归结构 $\beta^{f(n)}$ 的递归格式指其第 $n+1$ 次递归用第 n 次递归表示的关系式。它有两种形式: 1) 参数域递归格式, 以 $f(n)$ 表示 $f(n+1)$; 2) 生成式域递归格式, 从 A 的第 n 次递归导出第 $n+1$ 次递归。

定理 3. 递归语言 L 的每个基本递归结构 $\beta^{f(n)}$ 都有如下形式之参数域递归格式时, 必能以程序文法描述:

$$\begin{bmatrix} f(n+1) \\ f_1(n+1) \\ f_2(n+1) \\ \vdots \\ f_i(n+1) \end{bmatrix} = \begin{bmatrix} k & h & g & \cdots & p & q \\ k_1 & h_1 & g_1 & \cdots & p_1 & q_1 \\ k_2 & h_2 & g_2 & \cdots & p_2 & q_2 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ k_i & h_i & g_i & \cdots & p_i & q_i \end{bmatrix} \begin{bmatrix} f(n) \\ f_1(n) \\ f_2(n) \\ \vdots \\ f_i(n) \\ 1 \end{bmatrix} \quad (5)$$

其中 i 为有限正整数; f_1, f_2, \dots, f_i 均为导出函数; 系数矩阵各元取常值。

能够满足本定理的语言类极为广泛。应特别指出, $f(n)$ 为有限项数的多项式函数

$$f(n) = a_i n^i + a_{i-1} n^{i-1} + \cdots + a_2 n^2 + a_1 n + a_0 \quad (6)$$

时, 恒能以程序文法描述。

定理 4. 对于满足定理 3 的递归语言 L , 当且仅当其基本递归结构均有如下形式

$$f(n+1) = f(n) + q, \quad q = \text{const} \quad (7)$$

之参数域递归格式时, 其程序文法才只须成功去向 U , 而失败去向 W 恒为空集 ϕ (称 U 型程序文法)。否则, 成功去向与失败去向均为必要(称 UW 型程序文法)。

从本定理可得引论, 只有线性独立及线性相关的递归语言才能以 U 型程序文法描述。

五、递归结构的属性文法描述

属性文法^[1-3]常作为句法-词义文法的通称。本文则从比较狭义的角度给以定义。

定义 10. 属性文法 G 为四元式

$$G = (V_N, V_T, P, S), \quad (8)$$

其各个基元 $a \in V_T$ 都以给定的词义信息制约, 即基元属性 $A(a)$ 。 P 中生成式取表 2 所示两种形式之一。其中 $X \in V_N; \rho, \gamma \in V_T^+$; 形式 1) 中连接关系符 $Q \in V_T$, 而形式 2) 则将 Q 列为词义; $A(*)$ 表示 $*$ 的属性矢量; Φ 为属性转换函数。

表 2

句法规则	词义规则
1) $X \rightarrow \rho Q \gamma$	$A(X) = \Phi(A(\rho), A(Q), A(\gamma))$
2) $X \rightarrow \rho \gamma$	$Q(\rho, \gamma), A(X) = \Phi(A(\rho), A(Q), A(\gamma))$

定义 11. 基元 $a \in V_T$ 的属性矢量

$$A(a) = [\varphi_1, \varphi_2, \dots, \varphi_w]^T \quad (9)$$

是基元的特征描述。它有两种方式：

1) 绝对特征描述：取外部的某一共同基准确定基元特征。此时基元与基元的连接关系只是简单前后连接 (CAT)；

2) 相对特征描述：只从基元内部确定其特征，而不考虑其外部关系。此时基元与基元的连接一般不是简单 CAT 关系，而须引入附加的连接特征。但此法可减少基元类型，简化句法。

举直线段 a 和 b (长度 l_a, l_b , 方向 θ_a, θ_b) 为例。绝对特征描述以坐标轴为基准取它们的方向 θ_a 及 θ_b , 故有属性矢量及连接关系为

$$\begin{aligned} A(a) &= [l_a, \theta_a]^T, & A(b) &= [l_b, \theta_b]^T, \\ Q(a, b) &= CAT, & A(Q) &= [0, 0]^T, \\ aQb &= aCATb = ab. \end{aligned}$$

相对特征描述时不取直线段对坐标轴的方向角绝对值，而取线段终端方向与始端方向间的相对角度。对直线段说，该相对角为零。故有

$$\begin{aligned} A(a) &= [l_a, 0]^T, & A(b) &= [l_b, 0]^T, \\ Q(a, b) &= (CAT, \phi), & A(Q) &= [0, \phi]^T, \\ aQb &= a(CAT, \phi)b, & \phi &= \theta_b - \theta_a. \end{aligned}$$

定义 12. 模式的词义化串表述是在串表述中考虑基元之间的连接关系，其形式为

$$\{a_1 Q_1 a_2 Q_2 \dots Q_k a_{k+1}\} \text{ 或 } \{a_1 a_2 \dots a_{k+1} | Q_1, Q_2, \dots, Q_k\}. \quad (10)$$

在绝对特征描述时， $Q = CAT$ 可略而不记，成为通常的串表述 $\{a_1 a_2 \dots a_{k+1}\}$ 。

相应地，基本递归结构的词义化串表述为

$$\{\beta Q \beta Q \dots Q \beta\} \text{ 或 } \{\beta^m | Q(\beta, \beta)\}. \quad (11)$$

显然，只有在递归子串 β 与 β 之间的连接关系为相同的 $Q(\beta, \beta)$ 时，才构成基本递归结构。递归子串的最简单情况为单个基元 a 。此时结构 $\{a^m | Q(a, a)\}$ 恒可用一扩展基元 a_m 代替之，而以词义规定 a_m 的特征。

定理 5. 任何递归串结构均可用属性文法描述。

证。不失一般性，取基本递归结构 M

$$M = \{\beta^m | \beta \in V^+, Q(\beta, \beta)\}, \quad (12)$$

并以表 3 所示属性文法描述。其中 M_i 表示第 i 次导出的 M 。词义规则限定了 $i < m$ 时运用第一生成式； $i = m$ 时则运用第二生成式。在 m 为任何已知变量或函数时，恒能任意地或随机地为 m 定值，因而恒能生成基本递归结构 M 。从而任何递归串结构均可以分层形式用属性文法描述。

表 3

句法规则	词义规则
$M \rightarrow \beta M$	$Q(\beta, M) = Q(\beta, \beta),$ $A(M_i) = \Phi(A(\beta), A(Q), A(M_{i+1})),$ $i < m.$
$M \rightarrow \beta$	$A(M_i) = A(\beta)$ $i = m$

六、条件型及属性型递归词义

对递归结构的描述作递归词义制约时,程序文法和属性文法是两种方式的典型:1)条件型.按给定条件确定运用哪一个生成式,程序文法是其代表;2)属性型.指定属性关系,约束结构的生成,属性文法是其代表.引伸这一概念,递归词义信息有可能表示为更适当的形式.为此,下面定义两种派生的描述方式.

定义 13. 递归条件文法以四元式表示

$$G = (V_N, V_T, P, S), \quad (13)$$

产生基本递归结构 $\{\beta^m\}$ 的子文法生成式集,具有表 4 所示的规范形式(未列入其它词义).其中 N_M 为 M 的当前递归次数.

表 4

句法规则	词义规则
$M \rightarrow \beta M$	$N_M < m$
$M \rightarrow \beta$	$N_M = m$

定义 14. 递归属性文法以递归参数为其结构属性.表为四元式

$$G = (V_N, V_T, P, S). \quad (14)$$

基本递归结构 $\{\beta^m\}$ 的子文法生成式集成为表 5 所示的简单形式(未考虑其它词义).这里,递归参数 m_M 在词义规则中直接作为属性限定之,例如 $m_M = m$ 或 $m_M = m(n)$ 等.

表 5

句法规则	词义规则
$M \rightarrow \beta^m M$	$m_M = m$

定理 6. 任何递归串结构都可以用递归条件文法或递归属性文法描述.

证.与定理 5 的证类似.

递归条件文法和递归属性文法是递归结构描述的更为简单的形式,十分有利于模式分析和文法推断的进行.

七、结 论

模式结构描述所存在的困难问题之一是递归结构的复杂性.本文根据基本递归结构的性质对串语言进行分类,从而可以简单地分析语言结构的特点.以基本递归结构作为模式结构的子模式,据此构成模式分层结构,并以句法-词义方法作分层描述的方式被推荐为一种规范形式.这一体系的优越性在于:1)能解决任何复杂串模式的描述问题;2)十分容易进行分析;3)非常有利于文法推断过程^[7].运用这一描述体系的原则于程序文

法、属性文法、递归条件文法以及递归属性文法的研究，揭示了它们对递归结构描述的高度能力和简单形式。

本研究得到美国普渡大学傅京孙 (K. S. Fu) 教授生前热情的支持。作者在此谨表深切悼念。

参 考 文 献

- [1] Fu, K. S., *Syntactic Pattern Recognition and Applications*, Prentice-Hall, Englewood Cliffs, N. J., 1982.
- [2] Fu, K. S., *A Step Towards Unification of Syntactic and Statistical Pattern Recognition*, *IEEE Trans. Pattern Anal. Machine Intell.*, 5(1983), 200—205.
- [3] Thomason, M. G., *Syntactic/Semantic Techniques in Pattern Recognition: A Survey*, *Int. J. Comput. Inf. Sci.*, 11(1982), 75—100.
- [4] Rosenkrantz, D. J., *Programmed Grammars and Classes of Formal Languages*, *JACM* 16(1969), 117—131.
- [5] Lu, H. R. and Fu, K. S., *Inferability of Context-Free Programmed Grammars*, *Int. J. Comput. Inf. Sci.*, 13 (1984).
- [6] Lu, H. R. and Fu, K. S., *A General Approach to Inference of Context-Free Programmed Grammars*, *IEEE Trans. System, Man, Cybernetics*, 14(1984), 191—202.
- [7] 路浩如, 句法-词义模式识别中递归结构的文法推断, 自动化学报(待刊).

RECURSIVE STRUCTURES IN PATTERN RECOGNITION AND THEIR SYNTACTIC-SEMANTIC DESCRIPTIONS

LU HAORU

(Zhejiang University)

ABSTRACT

One of the complicated problems of pattern description is the effect of recursive structures on the nature of a pattern as a class of the formal language. In the present paper the properties of different recursive structures have been investigated. A strategy of hierarchical syntactic-semantic description has been studied by dividing a pattern into subpatterns according to its basic recursive structures. Thus, a systematic and effective approach for describing recursive patterns has been suggested. The approach can be successfully applied to context-free programmed grammars, attributed grammars and two special types of modifications, namely, recursion-conditioned and recursion-attributed grammars. Their descriptive power for recursive structures has been discussed.