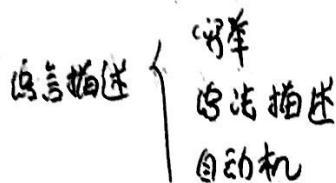


形式语言与自动机



形式语言 (代数语言): 精确地描述语言 $\alpha \rightarrow \beta$ (α, β 均为字符串)

4元组 $G = (V, \Sigma, P, S)$

V : 非终结符

P : 一组重写规则的有限集合

Σ : 终结符

$S \in V$: 句子符或初始符

推导: 设 $G = (V, \Sigma, P, S)$ 是一个文法, 在 $(V \cup \Sigma)^*$ 上定义关系 \Rightarrow 如下:

即 $\alpha \Rightarrow \beta$ 当且仅当 α 和 β 的任意因式 $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ 满足 $\alpha_1 \alpha_2 \alpha_3 \alpha_4 = \alpha$ 且 $\alpha_2 \rightarrow \alpha_3$ 是 P 的生成式, 则 $\alpha_1 \alpha_4 = \beta$

如果 $\alpha \beta \gamma$ 是 $(V \cup \Sigma)^*$ 中的字符串, 且 $\beta \rightarrow \delta$ 是 P 的生成式, 则

$$\alpha \beta \gamma \Rightarrow \alpha \delta \gamma$$

$\xRightarrow{+}_G$: 至少推导一次, 称为按非平凡方式生成 $\alpha_1 \rightarrow \alpha_2 \dots \alpha_n$ $n \geq 1$

$\xRightarrow{*}_G$: 称为生成, 表示 \Rightarrow 的自反和传递闭包

最左推导: 只改写最左边的非终结符

最右推导: 只改写最右边的非终结符 (规范推导)

例: $G = (\{E, T, F\}, \{a, +, *, (,)\}, P, E)$

$P: E \rightarrow E + T \mid T \quad T \rightarrow T * F \mid F \quad F \rightarrow (E) \mid a$

字符串 $a + a * a$ 的推导过程

最右: $E \Rightarrow E + T \Rightarrow E + F \Rightarrow E + (E) \Rightarrow E + (E + T) \Rightarrow E + (E + F) \Rightarrow E + (E + (E)) \Rightarrow E + (E + (a)) \Rightarrow E + (E + a) \Rightarrow E + T + a \Rightarrow E + F * a \Rightarrow E + (E) * a \Rightarrow E + T * a \Rightarrow E + T * F \Rightarrow E + T * (E) \Rightarrow E + T * (E + T) \Rightarrow E + T * (E + F) \Rightarrow E + T * (E + (E)) \Rightarrow E + T * (E + (a)) \Rightarrow E + T * a \Rightarrow F + a * a \Rightarrow a + a * a$

句型与句子 (推导过程中任一中间过程)

定义: 一些特殊类型的字符串为文法 $G = (V, \Sigma, P, S)$ 的句子形式(句型):

1. S 是一个句子形式

2. 如果 $\alpha \beta \gamma$ 是一个句子形式, 且 $\beta \rightarrow \delta$ 是 P 的生成式, 则 $\alpha \delta \gamma$ 也是一个句子形式

文法 G 的不含非终结符的句子形式称为 G 生成的句子。由文法 G 生成的语言, 记为 $L(G)$



$$L(G) = \{x \mid x \in \Sigma, s \xrightarrow{*} x\}$$

$$G = (N, \Sigma, P, S)$$

正则文法:

正则文法: $P: A \rightarrow Bx, \text{ 或 } A \rightarrow x$, 其中 $A, B \in N, x \in \Sigma$, 则称为左线性正则文法

正则文法或 3 型文法

$A \rightarrow xB$ 右线性正则文法

例: $G = (N, \Sigma, P, S)$

$$N = \{S, A, B\} \quad \Sigma = \{a, b\}$$

$$P: (a) S \rightarrow aA \quad (b) A \rightarrow aA \quad (c) A \rightarrow bB \quad (d) B \rightarrow bB \quad (e) B \rightarrow b$$

求 $L(G)$

可以将 (c) $A \rightarrow bB$ 改写为 $A \rightarrow bB'$ $B' \rightarrow bB$

则满足右线性正则文法

$$L(G) = \{a^n b^m\} \quad n \geq 1, m \geq 3$$

上下文无关文法: context-free grammar, (CFG) 也称为 2 型文法

$A \rightarrow \alpha$, 其中 $A \in N, \alpha \in (N \cup \Sigma)^*$

例: $G = (N, \Sigma, P, S) \quad N = \{S, A, B, C\} \quad \Sigma = \{a, b, c\}$

$$P: (a) S \rightarrow ABC \quad (b) A \rightarrow aA \mid a \quad (c) B \rightarrow bB \mid b \quad (d) C \rightarrow cC \mid c$$

$$L(G) = \{a^n b^m c^k\}, n \geq 1, m \geq 1, k \geq 0, \alpha \in \{0, 1\}$$

$$c \geq 0, \alpha = 0, \text{ 则 } \alpha = 1\}$$

上下文有关文法: (CSG) 也称为 1 型文法

$P: \alpha AB \rightarrow \alpha \gamma \beta$, 其中 $A \in N, \alpha, \beta, \gamma \in (N \cup \Sigma)^*$, 且 γ 至少包含一个非终结符

即: $|x| \rightarrow |y|, x \in (N \cup \Sigma)^+, y \in (N \cup \Sigma)^*, \text{ 并且 } |y| \geq |x|$

$| \cdot |$ 表示长度

例: $G = (N, \Sigma, P, S) \quad N = \{S, A, B, C\} \quad \Sigma = \{a, b, c\}$

$$P: (a) S \rightarrow ABC \quad (b) A \rightarrow aA \mid a \quad (c) B \rightarrow bB \mid b \quad (d) BC \rightarrow BCC$$

$$L(G) = \{a^n b^m c^k\}, n \geq 1, m \geq 1$$



无约束文法: 0型文法

$$P: \alpha \rightarrow \beta$$

$$L(G_0) \supseteq L(G_1) \supseteq L(G_2) \supseteq L(G_3)$$

任意 有限 无限 正则

语言与文法类型的约定:

如果一种语言能由几种文法所产生, 则把这种语言称为在这几种文法中当限制最少的那种文法所产生的语言。

例: $G = (S, A, B, \{a, b\}, P, S)$

$$P: S \rightarrow AB \quad S \rightarrow bA \quad A \rightarrow aS \quad A \rightarrow bAA$$

$$A \rightarrow a \quad B \rightarrow bS \quad B \rightarrow aBB \quad B \rightarrow b$$

G 为 0 型文法 $L(G) = \{\text{非空串的 } a \text{ 和 } b \text{ 构成的串}\}$

CFG 产生的语言句子的生成树表示:

1) 对于 $X \in V \cup \Sigma$ 给一个标记作为节点, S 作为树的根节点

2) 如果一个节点的标记为 A , 并且它至少有一个除它自身以外的后裔, 则 $A \in V$

3) 如果一个节点的标记为 A , 它的 $L(G)$ 个直接后裔节点按从左到右的顺序依次标记为 A_1, A_2, \dots, A_n , 则 $A \rightarrow A_1 A_2 \dots A_n$ 一定是 P 中的一个产生式

例: $G = (S, A, \{a, b\}, P, S)$

$$P: S \rightarrow bA \quad A \rightarrow bAA \quad A \rightarrow a$$

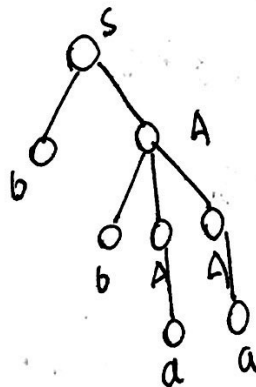
G 所生成的句子 $bbba$ 可以由下面的生成树表示:

$$S \Rightarrow bA$$

$$\Rightarrow bbAA$$

$$\Rightarrow bbba$$

$$\Rightarrow bbba$$



上下文无关文法的二义性: 一个文法 G , 如果存在某个句子有不止一棵解析树与之对应, 则称这个文法是二义的。

计算无歧义性

歧义的产生
(不同方法)



例: p : 名词 $\sim p$: 名词短语 pp : 介词短语 Avx : 动词 $\sim Avx$: 动词短语

S : 句子/短语

有限自动机与正则文法

确定的有限自动机 (DFA): 五元组 $M = (\Sigma, Q, \delta, q_0, F)$

Σ : 输入符号的有限集合

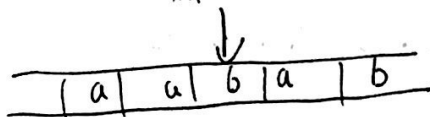
Q : Q 是状态的有限集合

q_0 : $q_0 \in Q$ 是初始状态

F : F 是终止状态集合, $F \subseteq Q$

δ : 是 Q 与 Σ 的直积 $Q \times \Sigma$ 到 Q (下一个状态) 的映射。它支配着有限状态控制的行为。有时也称为状态转移函数。

示意图:



判断, 接收, 终止

变换图: $\delta(q, a) = q'$



语言: 如果一个句子 x 使得有限自动机 M 有 $\delta(q_0, x) = p, p \in F$, 那么称句子 x 被 M 接受, $T(M)$ 即为被 M 接受的句子的集合, 即

$$T(M) = \{x \mid \delta(q_0, x) \in F\}$$

不确定的有限自动机: 五元组 $(\Sigma, Q, \delta, q_0, F)$ M

NFA

δ : Q 与 Σ 的直积 $Q \times \Sigma$ 到 Q 的幂集 2^Q 的映射。

输出不确定

区别:
 NFA $\delta(q, a)$ 是一个状态集合
 DFA $\delta(q, a)$ 是一个状态

NFA 与 DFA 的关系: (此将 FA) (因此不再区分)

定理 3.1: 设 L 是一个被 NFA 所接受的句子的集合, 则存在一个 DFA, 它恰好接受 L

正则文法与有限自动机的关系:

若 $G = (V_n, V_t, P, S)$ 是一个正则文法, 则存在一个有限自动机 $M = (\Sigma, Q, \delta, q_0, F)$

使得 $T(M) = L(G)$

