

## 第 2 章 形式语言基础

### 【前言】

计算机处理语言，首先应考虑语言的形式化、规范化，使其具有可计算性和可操作性；这就是形式语言理论研究的问题。

形式语言诞生于1956年，由chomsky创立。通常，语言研究至少涉及三个方面：**语法、语义和语用**；这里仅侧重于**语法的研究**。

形式语言的**基本观点**是：

**语言是符号串之集合！**

形式语言理论研究的基本问题是：

**研究符号串集合的表示方法、结构特性以及运算规律。**

## 第 2 章 形式语言基础

### 【内容提要】

- 2.1 形式语言是符号串集合
- 2.2 形式语言是由文法定义的
- 2.3 主要语法成分的定义
- 2.4 两类特性文法
- 2.5 文法变换方法
- 2.6 关于形式语言的分类问题

## 2.1 形式语言是符号串集合

**【形式语言】**是字母表上的符号，按一定的规则组成的所有符号串集合；其中的每个符号串称为句子。

**【名词解释】：**

- **字母表** -- 元素(符号)的非空有限集合；
- **符号串** -- 符号的有限序列；
- **符号串集合** -- 有限个或者无限个符号串组成的集合；
- **规则** -- 以某种形式表达的在一定范围内共同遵守的章程和制度；这里，指符号串的组成规则。



## 形式语言概念示例：

两个语言！

**【例2.1】**  $L_1 = \{ 00, 01, 10, 11 \}$  ;  
字母表  $\Sigma_1 = \{0, 1\}$ ,  
句子有: 00, 01, 10, 11

**【例2.2】**  $L_2 = \{ ab^m c, b^n \mid m > 0, n \geq 0 \}$   
字母表  $\Sigma_2 = \{a, b, c\}$ ,

- 句型1:  $ab^m c$  ,  
有句子:  $abc, abbc, abbbc, \dots$
- 句型2:  $b^n$  ;  
有句子:  $\varepsilon, b, bb, bbb, \dots$

**【注】** (1)  $b^0 = \varepsilon$  (空符号串),  $b^1 = b, b^2 = bb, b^3 = bbb, \dots$

(2)  $L_1$  为有限语言;  $L_2$  为无限语言。

## 2.1.1 符号串(集合)的运算

### 1. 符号串的运算

设  $\alpha$ ,  $\beta$  为两个符号串, 则:

1. 连接:  $\alpha\beta = \alpha\beta$       如  $a.b=ab$

2. 或:  $\alpha|\beta = \alpha$  (或者  $\beta$ )

3. 方幂:  $\alpha^n = \underbrace{\alpha\alpha \dots \alpha}_{n\uparrow} = \alpha\alpha^{n-1} = \alpha^{n-1}\alpha$

※  $\alpha^0 = \varepsilon$  (空符号串)

└ 什么符号也没有的符号串 !

$\alpha^1 = \alpha$  ;  $\alpha^2 = \alpha\alpha$  ; ...

4. 闭包:

$\alpha$  的正闭包:  $\alpha^+ = \alpha^1|\alpha^2|\dots|\alpha^n|\dots$

$\alpha$  的星闭包:  $\alpha^* = \alpha^0|\alpha^1|\alpha^2|\dots|\alpha^n|\dots$

## 2.1.1 符号串(集合)的运算(续1)

### ※ 符号串运算示例

【例】：

$$(1) \quad abc.de = abcde$$

$$(2) \quad ab|cd = ab \text{ (或者 } cd \text{)}$$

$$(3) \quad (a|b)^1 = (a|b) = a|b$$

$$(a|b)^2 = (a|b)(a|b) = aa|ab|ba|bb$$

...

$$(a|b)^* = (a|b)^0 | (a|b)^1 | (a|b)^2 | \dots$$

即：  $(a|b)^* = (a|b)^n, n \geq 0$

同理：  $(a|b)^+ = (a|b)^n, n > 0$

## 2.1.1 符号串(集合)的运算(续2)

### II. 符号串集合的运算

设  $A$  和  $B$  为两个符号串集合, 则:

1. 乘积:  $AB = \{xy \mid x \in A \text{ 且 } y \in B\}$

2. 和:  $A \cup B = A + B = \{x \mid x \in A \text{ 或 } x \in B\}$

3. 方幂:  $A^n = \underbrace{AA \dots A}_{n \uparrow} = AA^{n-1} = A^{n-1}A$

$\times A^0 = \{\varepsilon\};$

$A^1 = A; A^2 = AA; A^3 = AAA; \dots$

4. 闭包:

$A$  的正闭包:  $A^+ = A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$

$A$  的星闭包:  $A^* = A^0 \cup A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$

## ➤ 符号串集合运算示例：

**【例2.3】** 设  $A = \{a, b\}$ ,  $B = \{c, d\}$

则  $A+B = \{a, b, c, d\}$

则  $AB = \{xy \mid x \in A, y \in B\} = \{ac, ad, bc, bd\}$

**【例2.4】** 设  $A = \{a\}$

则  $A^* = A^0 \cup A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$

$= \{\varepsilon\} + \{a\} + \{aa\} + \{aaa\} + \dots$

$= \{\varepsilon, a, aa, aaa, \dots\}$

$= \{a^n \mid n \geq 0\}$



## ➤ 符号串集合运算示例(续):

【例2.5】 设  $A = \{a, b\}$ ,  $A^* = ?$

$$\therefore A^* = A^0 \cup A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$$

$$A^0 = \{\varepsilon\};$$

$$A^1 = A = \{a, b\};$$

$$A^2 = A \cdot A = \{a, b\} \cdot \{a, b\} = \{aa, ab, ba, bb\};$$

$$\begin{aligned} A^3 &= A \cdot A^2 = \{a, b\} \cdot \{aa, ab, ba, bb\} \\ &= \{aaa, aab, aba, abb, baa, bab, bba, bbb\}; \end{aligned}$$

...

$$\therefore A^* = \{ x \mid x = (a|b)^n, n \geq 0 \}$$

**推论:** 若  $A$  为任一字母表, 则  $A^*$  就是该字母表上的所有符号串(包括空串)的集合。

## 2.1.2 符号串集合的文法描述

长久以来，探讨符号串集合(即形式语言)的各种描述方法，一直是计算机语言处理的重要任务之一。

【例2.5】  $L = \{ ab^n c \mid n \geq 0 \}$  ,

字母表:  $\Sigma = \{a, b, c\}$  ;

展开:  $L = \{ac, abc, abbc, abbbc, \dots\}$

右图给出的表示方法——**文法规则** ;

$S \rightarrow aAc$

$A \rightarrow bA \mid \varepsilon$

- (1)  $S, A$  — 定义的对象 ( $S$  **句子**, 最大的定义对象, 又称为开始符号;  $A$  为句型  $aAc$  的**短语**) ,
- (2)  $a, b, c$  — 为字母表  $\Sigma$  中的符号;  $\varepsilon$  — 空符号串。
- (3)  $\rightarrow, |$  — 为**描述符号** ( $\rightarrow$  定义为;  $|$  或者是)

## ➤ 规则应用说明示例:

怎样利用上述**语法规则**表示语言L?

从**开始符号**出发，对符号串中的**定义对象**，采用**推导**的方法（**用其规则右部替换左部**）产生新的符号串，如此进行，直到新符号串中不再出现定义的对象为止，则最终的符号串就是一个**句子**。

【句子产生过程】( $\Rightarrow$  **推导算符**):

$$\textcircled{1} S \Rightarrow aAc \Rightarrow a \varepsilon c = ac$$

$$\textcircled{2} S \Rightarrow aAc \Rightarrow abAc \Rightarrow ab \varepsilon c = abc$$

$$\textcircled{3} S \Rightarrow aAc \Rightarrow abAc \Rightarrow abbAc \Rightarrow abbc$$

...

$$\therefore S \stackrel{+}{\Rightarrow} ab^n c, n \geq 0$$

$$S \rightarrow aAc$$

$$A \rightarrow bA \mid \varepsilon$$

## 2.2 形式语言是由文法定义的

### 2.2.1 什么是文法？

【定义】 文法(grammar)是规则的有限集, 其中的上下文无关文法可定义为四元组:

每个元素

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

$V_N$  : 非终结符集 (定义的对象集, 如: 语法成分等);  
 $V_T$  : 终结符集 (字母表);  
 $Z$  : 开始符号 (研究范畴中, 最大的定义对象);  
 $P$  : 规则集 (又称产生式集);

每个规则

$A \rightarrow \alpha$  或者  $A \rightarrow \alpha \mid \beta$

其中, 描述符号 :  $\rightarrow$  (定义为),  $\mid$  (或者是)

文法符号 :  $Z, A \in V_N, \alpha, \beta \in (V_N + V_T)^*$

## 2.2 形式语言是由文法定义的（续3）

【注意】 提供了规则集，就相当给出了一个文法：

$G(S)$ :

$S \rightarrow aAc$

$A \rightarrow bA \mid \varepsilon$

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

$V_N = \{ S, A \}; V_T = \{ a, b, c \};$

$Z = S; P :$

### 2.2.2 文法是怎样定义语言的？

设 有文法  $G(Z)$ ,  $L(G)$  为  $G$  所定义的语言；

则  $L(G) = \{ x \mid Z \stackrel{+}{\Rightarrow} x, x \in V_T^* \}$  〔2.1〕

利用  $\Rightarrow$  进行连续推导之意！

即：一个文法所定义的**语言**，就是由该文法**开始符号**推导出的所有**仅含终结符的符号串之集合**。  
其中的每个符号串皆称为**句子**。

## 【例2.6】标识符的文法

【标识符】 指字母开头的字母、数字序列。

令  $G(Z) = (V_N, V_T, Z, P)$

则  $V_N = \{I(\text{标识符}), A(\text{标识符尾})\};$

$V_T = \{\ell(\text{字母}), d(\text{数字})\};$

$Z = I;$

$P:$

$I \rightarrow \ell A \mid \ell$

$A \rightarrow \ell A \mid d A \mid \varepsilon$

同理，【无符号整数】文法 可写成：

$G(N): N \rightarrow d N \mid d$

※其四元组也可写成：  $G(N) = (\{N\}, \{d\}, N, P)$

## ※标识符文法所定义的语言求解：

上面构造的  
标识符文法属于  
正规文法 (定义在后) 类,

$$\begin{aligned} I &\rightarrow \ell A \mid \ell \\ A &\rightarrow \ell A \mid d A \mid \varepsilon \end{aligned}$$

正确性检验较容易；下面给出一个算法：

$$\text{令： } I = \ell A \mid \ell$$

$$A = \ell A \mid d A \mid \varepsilon$$

} 正规方程式

### ※ 求解 I 值步骤：

① 先求解 A:  $A = (\ell \mid d) A$ ,  $A = (\ell \mid d)^2 A$ , ...,  $A = (\ell \mid d)^n A$

代入  $A = \varepsilon$  得:  $A = (\ell \mid d)^n, n \geq 0$

②  $\because I = \ell A \mid \ell$  代入  $A = (\ell \mid d)^n, n \geq 0$

得:  $I = \ell (\ell \mid d)^n, n \geq 0$



《标识符》：字母开头的字母、数字序列；

## 【例2.7】简单算术表达式文法

令  $G(Z) = (V_N, V_T, Z, P)$

则  $V_N = \{ E(\text{算术表达式}), T(\text{项}), F(\text{因式}) \};$

$V_T = \{ i(\text{变量或常数}), +, -, *, /, (, ) \};$

$Z = E;$

$P:$

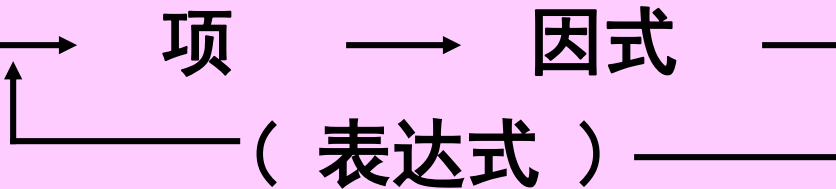
$E \rightarrow T \mid E + T \mid E - T$

$T \rightarrow F \mid T * F \mid T / F$

$F \rightarrow i \mid ( E )$

【注】此文法定义了算术表达式的层次嵌套结构：

表达式  $\longrightarrow$  项  $\longrightarrow$  因式





## ※ 算术表达式文法应用示例：

证明  $i*(i+i-i)$   
是文法 $G(E)$ 的一个句子  
(即 合法的算术表达式)：

$G(E) :$   $E \rightarrow T \mid E+T \mid E-T$   
 $T \rightarrow F \mid T*F \mid T/F$   
 $F \rightarrow i \mid (E)$

根据 语言定义式 [2.1] ,合法的算术表达式是指：

$E \xRightarrow{+} i*(i+i-i)$  成立吗？

$\therefore E \Rightarrow T \Rightarrow T*F \Rightarrow T*(E) \Rightarrow T*(E-T)$   
 $\Rightarrow T*(E+T-T) \Rightarrow F*(E+T-T) \Rightarrow i*(E+T-T)$   
 $\Rightarrow \dots \Rightarrow i*(i+i-i)$

$\therefore E \xRightarrow{+} i*(i+i-i)$

观察推导过程，可以看到：一旦产生式选择错了，会导致失败！

谢谢