

## 第二章 上下文无关文法和语言

§ 2.1 文法和语言的表示

§ 2.2 文法和语言的定义

§ 2.3 句型的分析

§ 2.4 文法的实用限制和其他表示法

§ 2.5 文法和语言的Chomsky分类

## 提要

- 所谓**形式化方法**，简单地说，就是用一整套带有严格规定的符号体系来描述问题的理论和方法，用形式化方法描述的语言（语法和语义）便是形式语言。
- 本章将从**形式语言**的角度系统地介绍什么是程序设计语言的**文法**，文法和语言的关系等问题，本章是本课程的理论基础。

## § 2.1 文法和语言的表示

语言的定义可采用下列三种方法

1. **枚举法**——把该语言的所有句子列出放在一集合内。（有限个句子时）
2. **有限条规则**——描述语言的全部句子。（有限或无限个句子），即文法表示
3. **装置**——检验和识别句子。（有限或无限个句子）即自动机

## § 2.2 文法和语言的定义

### § 2.2.1、基本概念和术语

1、字母表：元素的非空有穷集合  
元素称符号

例： $\Sigma = \{0,1\}$

2、符号串：字母表中的符号所组成的任何有穷序列

特别：空符号串  $\varepsilon$  （不包含任何符号）

### 3、字母表 $\Sigma$ 上的符号串的递归定义

(1)  $\varepsilon$  是 $\Sigma$ 上的符号串

(2) 若 $x$ 是 $\Sigma$ 上的符号串, 且 $a \in \Sigma$ ,  
则 $xa$ 或 $ax$ 是 $\Sigma$ 上的符号串

特别:  $\varepsilon x = x\varepsilon = x$

(3) 若 $y$ 是 $\Sigma$ 上的符号串,  
当且仅当 $y$ 可由 (1) 和 (2) 产生。

例:  $\Sigma = \{b, c\}$ ,

求 $\Sigma$ 上的所有符号串

根据1  $\varepsilon$ 是 $\Sigma$ 上的符号串

根据2  $\varepsilon b$ 和 $\varepsilon c$ 即  $b, c$  是 $\Sigma$ 上的符号串

$bb, bc, cb, cc$ 是 $\Sigma$ 上的符号串

$bbb, bbc, bcb, bcc, cbb, cbc, ccb, ccc$ 是 $\Sigma$ 上的符号串.....

$\Sigma$ 上的所有符号串  $\varepsilon, b, c, bb, bc, cb, cc, bbb, bbc, bcb, bcc, cbb, cbc, ccb, ccc, \dots$

- (1)  $\varepsilon$  是  $\Sigma$  上的符号串
- (2) 若  $x$  是  $\Sigma$  上的符号串, 且  $a \in \Sigma$ , 则  $xa$  或  $ax$  是  $\Sigma$  上的符号串  
特别:  $\varepsilon \cdot x = x \quad \varepsilon = x$
- (3) 若  $y$  是  $\Sigma$  上的符号串,  
当且仅当  $y$  可由 (1) 和 (2) 产生。

## 5、符号串的前缀、后缀和子串：

前缀：设 $X$ 是一符号串，从 $X$ 的尾部删去若干个（包括0个）符号之后所剩余下的部分称为 $X$ 的**前缀**；

若 $X$ 的前缀不是 $X$ 本身，则称为 $X$ 的**真前缀**。

后缀：设 $X$ 是一符号串，从 $X$ 的头部删去若干个（包括0个）符号之后所剩余下的部分称为 $X$ 的**后缀**；

若 $X$ 的后缀不是 $X$ 本身，则称为 $X$ 的**真后缀**。

子串：从一个符号串中删去它的一个前缀和一个后缀之后所剩下的部分称为此符号串的**子串**。

若 $X$ 的子串不是 $X$ 本身，则称为 $X$ 的**真子串**

例 设  $x=abc$

$x = \varepsilon a \varepsilon b \varepsilon c \varepsilon$

$x$ 的前缀:  $abc \quad ab \quad a \quad \varepsilon$

( $\varepsilon, a, ab$ 为真前缀)

$x$ 的后缀:  $abc \quad bc \quad c \quad \varepsilon$

( $\varepsilon, c, bc$ 为真后缀)

$x$ 的子串:  $abc \quad ab \quad bc \quad a \quad b \quad c \quad \varepsilon$

( $\varepsilon, a, ab, b, c, bc$ 为真子串)

用法:  $Z=X\dots\dots$  ( $x$ 为 $Z$ 的前缀)

$Z=\dots\dots X$  ( $x$ 为 $Z$ 的后缀) 感兴趣部分

$Z=\dots X \dots$  ( $x$ 为 $Z$ 的子串)



5、符号串的长度：符号串所含符号的个数

6、符号串的连接和方幂

连接：设有符号串 $x, y$ ，把 $y$ 的符号写在 $x$ 的符号之后所得的符号串，叫做 $x$ 与 $y$ 的连接，记 $xy$

方幂：设有符号串 $x$ ，则 $x$ 的 $n$ 次自身连接称为 $x$ 的 $n$ 次方幂，记为 $x^n$

特别： $x^0 = \varepsilon$

## 7、符号串集合A与B的和与积：

**和**：  $A+B=\{w|w \in A \text{ 或 } w \in B\}$

**积**：  $AB=\{xy|x \in A \text{ 且 } y \in B\}$

## 8、符号串集合的方幂：

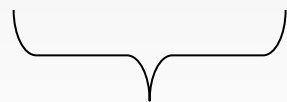
设有符号串集合A

则定义  $A^0 = \{\epsilon\}$

$$A^1=A$$

$$A^2=AA$$

$$A^n=AA\dots A$$



n个

## 9、符号串(符号)集合的正闭包

设  $A$  为符号集合，则定义  $A$  的正闭包  $A^+$  为：

例：  $\Sigma = \{b, c\}$  ,

求  $\Sigma$  上的所有符号串

$\varepsilon \cup A \cup A^2 \cup \dots$

$\varepsilon, b, c, bb, bc, cb, cc, bbb, bbc, bcb, bcc, cbb, bcb, ccb, ccc, \dots$

则  $A^+ = \{b, c\} \cup \{bb, bc, cb, cc\} \cup \dots$

$= \{b, c, bb, bc, cb, cc, bbb, bbc, bcb, bcc, cbb, bcb, ccb, ccc, \dots\}$

➤  $A^+$  是由  $A$  中的元素构成的符号串(除  $\varepsilon$ )的集合。

## 10、符号串（符号）集合的闭包 $A^*$

设 $A$ 为符号集合，则定义 $A$ 的闭包 $A^*$ 为：

$$A^* = A^0 \cup A^+ = \{\varepsilon\} \cup A^+$$

则 $A^* = \{\varepsilon\} \cup \{b, c\} \cup \{bb, bc, cb, cc\} \cup \dots$

$= \{\varepsilon, b, c, bb, bc, cb, cc, bbb, bbc, bcb, bcc, cbb, bcb, ccb, ccc, \dots\}$

➤  $A^*$ 由 $A$ 中的元素构成的所有符号串的集合。

# 问 题

- 为什么对符号、符号串、符号串集合以及它们的运算感兴趣？

若 $V$ 为某语言的字母表： $V=\{a, b, \dots 0, 1, \dots 9, +, =, (, ) \dots \text{if, else, for} \dots\}$

$X$ 为单词集： $X=\{\text{if, else, for}, \dots, \langle \text{标识符} \rangle, \langle \text{常量} \rangle, \dots\}$ , 则单词  $x \in V^*$

句子是定义在单词集 $X$ 上的符号串, 则句子  $a \in X^*$

语言是由句子组成的集合, 是字母表 $V$ 上的符号串集合

- 例如：字母表  $\Sigma = \{a, b\}$

$$\Sigma^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\}$$

那么，集合  $\{ab, aabb, aaabbb, \dots, a^n b^n, \}$  或  $\{w | w \in V^* \text{ 且 } w = a^n b^n, n \geq 1\}$

为字母表  $\Sigma$  上的一个语言

## § 2.2.2、文法和语言的形式定义

### 1、文法的形式定义

**问题：如何定义句子的合法性？**

**有穷语言：逐一列举句子**

**无穷语言：文法是以有穷的集合刻画无穷集合的工具**

**文法是对语言结构的定义和描述（或称“语法”）**

**注意：从一组规则可以推导出不同的句子，文法是在形式上对句子结构的定义和描述，而未涉及语义问题**

## § 2.2.2、文法和语言的形式定义

### 1、文法的形式定义

#### (1) 规则 (产生式)

一有序对  $(U, x)$  记为  $U ::= x$  或  $U \rightarrow x$

规则的左部:  $U$  是符号 (有的文法也可符号串)

规则的右部:  $x$  是有穷符号串

表示:  $U$  定义为  $x$

例:  $S \rightarrow abc$

$\langle \text{主函数} \rangle \rightarrow \text{main} (\langle \text{参数表} \rangle) \langle \text{参数说明} \rangle \langle \text{函数体} \rangle$



(2) 文法  $G[Z]$ : 规则的非空有穷集合.

**Z:** 开始符号 (识别符号), 至少在一条规则的左部出现。

(3) 字汇表  $V$ : 规则左右部中所有符号组成的集合

**非终结符号:** 需要定义的语法范畴  
组成非终结符号集合  $V_n$

**终结符号:** 规则中不属于  $V_n$  的符号  
组成终结符号集合  $V_t$

➤ 非终结符号以  $\langle \rangle$  括起, 但当是大写字母是时常省略

➤  $V = V_n \cup V_t$

#### (4) 文法的四元组表示

$$G=(V_n, V_t, P, Z)$$

其中  $V_n$  : 非终结符号集

$V_t$  : 终结符号集

$P$  : 规则的集合

$Z$  : 文法的开始符号

➤ 文法的BNF（巴科斯范式）表示

规则中有相同的左部时： $V \rightarrow x$

$V \rightarrow y$

...

$V \rightarrow z$

写成： $V \rightarrow x|y|\dots|z$ （' | '表达‘或’）

➤ 元符号： $\rightarrow$ ， $::=$ ， $()$ ， $|$ ， $<$ ， $>$

➤ 元语言：元符号处理的语言

➤ 由元符号组成的巴科斯范式(元语言公式)是用来描述算法语言的元语言

例:  $G=(V_n, V_t, P, S)$

其中:

$V_n=\{S, A, B\}$

$V_t=\{a, b\}$

S: 文法的开始符号

P:  $S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

文法BNF表示为

G[S]:

$S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

➤ 一般规定: 第一条规则的左部为开始符号,  
那么文法有规则的集合就完全确定了。

## 2、推导的形式定义

(1) 直接推导：如果 $U \rightarrow u$ 是 $G$ 中的一条规则，

$$x, y \in V^*,$$

则将规则 $U \rightarrow u$ 用于符号串 $r = xUy$ 上

得到符号串 $w = xuy$

记为： $xUy \Rightarrow xuy$  ( $r \Rightarrow w$ )

称符号串 $w$ 是符号串 $r$ 的**直接推导**,或符号串 $r$ 直接产生了符号串 $w$ ,也称 $w$ **直接归约**到 $r$ .

文法BNF表示为  
G[S]:

$S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

例:上述文法G[S]可进行的直接推导

①  $S \Rightarrow aB$  (规则  $S \rightarrow aB$  )

$U \Rightarrow u$

(规则  $U \rightarrow u$  ,  $x, y$  均为  $\epsilon$  )

②  $abS \Rightarrow abbA$  (规则  $S \rightarrow bA$  )

$xU \Rightarrow xu$  (规则  $U \rightarrow u$  ,  $x$  为  $ab, y$  为  $\epsilon$  )

③  $aB \Rightarrow aaBB$  (规则  $B \rightarrow aBB$  )

$xU \Rightarrow xu$  (规则  $U \rightarrow u$  ,  $x$  为  $aaB, y$  为  $\epsilon$  )

文法BNF表示为

G[S]:

$S \rightarrow aB \mid bA$

$A \rightarrow a \mid aS \mid bAA$

$B \rightarrow b \mid bS \mid aBB$

例:上述文法G[S]可进行的推导

推导过程

使用规则

$S \Rightarrow aB$

$S \rightarrow aB$

$\Rightarrow abS$

$B \rightarrow bS$

$\Rightarrow abbA$

$S \rightarrow bA$

$\Rightarrow abbbAA$

$A \rightarrow bAA$

$\Rightarrow abbbaA$

$A \rightarrow a$

$\Rightarrow abbbaa$

$A \rightarrow a$

➤只要符号串中存在非终结符号, 推导就能继续, 直至符号串全由终结符号组成, 这也是为什么称终结符和非终结符的原因

## (2) 推导(长度为n):

设 $u_0, u_1, \dots, u_n (n > 0)$ 均 $\in V^*$ , 且有

$$r = u_0 \Rightarrow u_1 \Rightarrow \dots \Rightarrow u_{n-1} \Rightarrow u_n = w$$

记为 $r \Rightarrow^+ w$  (一步或一步以上)

则称以上序列为长度 $n$ 的推导,

也称 $r$ 产生 $w$ ( $w$ 归约为 $r$ )

特例:  $r = w$ (0步推导)

➤ 长度 $n \geq 0$ 的推导记为  $r \Rightarrow^* w$ ; 长度 $n > 0$ 的推导记为  $r \Rightarrow^+ w$



### 3、语言的形式定义

(1) 句型: 设有文法  $G[Z]$ , 如果有  $Z \xRightarrow{*} x, x \in V^*$ ,  
则称  $x$  是文法  $G$  的一个句型。

□ 凡是由识别符号推导出来的字汇表  $V$  上的 **终结**  
**和非终结符号** 组成的符号串叫做 **句型**。

(2) 句子: 如有  $Z \xRightarrow{+} x, x \in V_t^*$ ,  
则称  $x$  是文法  $G$  的一个句子。

□ 由  $Z$  推导的 **终结符号** 组成的符号串为 **句子**。

(3) 语言 $L(G[Z])$ : 文法 $G[Z]$ 产生的所有句子的集合,  
称文法 $G[Z]$ 所定义的语言

$$L(G[Z]) = \{x \mid x \in V_t^* \text{ 且 } Z \Rightarrow^+ x\}$$

例: $G[\langle \text{标识符} \rangle]$

$\langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle \mid \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\mid \langle \text{标识符} \rangle \langle \text{数字} \rangle$

$\langle \text{字母} \rangle \rightarrow a \mid b \mid \dots \mid z \mid A \mid \dots \mid Z$

$\langle \text{数字} \rangle \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

问题: 符号串 “a4y” 是不是文法的句子?

结论:

a4y是文法的合法句子

例:  $G[\langle \text{标识符} \rangle]$

$\langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle \mid \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\mid \langle \text{标识符} \rangle \langle \text{数字} \rangle$

$\langle \text{字母} \rangle \rightarrow a \mid b \mid \dots \mid z \mid A \mid \dots \mid Z$

$\langle \text{数字} \rangle \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9$

## 推导过程1

$\langle \text{标识符} \rangle$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\Rightarrow \langle \text{标识符} \rangle y$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle y$

$\Rightarrow \langle \text{标识符} \rangle 4y$

$\Rightarrow \langle \text{字母} \rangle 4y$

$\Rightarrow a4y$

## 推导过程2

$\langle \text{标识符} \rangle$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle \langle \text{字母} \rangle$

$\Rightarrow \langle \text{字母} \rangle \langle \text{数字} \rangle \langle \text{字母} \rangle$

$\Rightarrow a \langle \text{数字} \rangle \langle \text{字母} \rangle$

$\Rightarrow a4 \langle \text{字母} \rangle$

$\Rightarrow a4y$

例:  $G1[A]: A \rightarrow Bb$

$B \rightarrow a$

$L(G1) = \{ab\}$

$G2[A]: A \rightarrow ab$

$L(G2) = \{ab\}$

$G1 \neq G2$  但  $L(G1) = L(G2)$

称  $G1$  和  $G2$  为等价文法 (不同文法, 相同语言).

➤ 给定文法后, 可以确定它的语言, 但由语言写出的文法是比较难的, 这里形式语言理论可以证明:

1° 给定一文法,就能从结构上唯一确定其语言.

即  $G$  唯一确定  $L(G)$

2° 给定一语言,能确定其文法,但这种文法不是唯一的。 即  $L$  确定  $G_1$ ,或 $G_2$ ...

3° 设  $G=(V_n, V_t, P, S)$  为一文法,

并设  $U \rightarrow xVy$ , 是  $P$  中一产生式,

且  $V \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$  是  $P$  中  $V$  的全部产生式

又设  $G_1=(V_n, V_t, P_1, S)$  是其中  $P_1$  从  $P$  中删去  $U \rightarrow xVy$ ,

加入  $U \rightarrow x\beta_1y, U \rightarrow x\beta_2y, \dots, U \rightarrow x\beta_ny$ ,

则  $L(G_1)=L(G)$

观察公式  $V \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$  两侧

$G3[S]: S \rightarrow A \mid S-A$

$A \rightarrow a \mid b \mid c$

$G4[S]: S \rightarrow A \mid A-S$

$A \rightarrow a \mid b \mid c$

符号串a-b-c是G3[S]、G4[S]合法句子，但语义不同。

G3解释为  $(a-b) - c$

G4解释为  $a - (b-c)$

## § 2.2.3 递归规则与递归文法

-----用有穷的规则刻划无穷的语言

### 1、递归规则

形如  $U \rightarrow xUy$      $U \in V_n, x, y \in V^*$

左右具有相同的非终结符号的规则

**特别:**  $U \rightarrow Uy$     ( $x = \varepsilon$ ) 左递归规则

$U \rightarrow xU$     ( $y = \varepsilon$ ) 右递归规则

$U \rightarrow xUy$     ( $x, y \neq \varepsilon$ ) 自嵌入递归规则

➤ 递归规则是对其左部的非终结符号进行递归定义

## 2.文法的递归性

1)直接递归性:文法中至少包含一条递归规则

2)间接递归性:文法的任一非终结符号经一步以上推导产生的递归性。

3)文法的递归性原则:文法具有直接递归性或间接递归性,否则,文法无递归性。

例1:  $G[Z]: Z \rightarrow aZb \mid ab$       具有直接递归性

例2:  $G[U]: U \rightarrow Vx$

$V \rightarrow Uy \mid z$       具有间接递归性

原因:  $U \Rightarrow Vx \Rightarrow Uyx$



- 例 :  $G[S]: S \rightarrow 0S1, S \rightarrow 01$

$$L(G) = \{ 0^n 1^n \mid n \geq 1 \}$$

- 例 :  $G[S]: (1) S \rightarrow aSBE$

$$(2) S \rightarrow aBE$$

$$(3) EB \rightarrow BE$$

$$(4) aB \rightarrow ab$$

$$(5) bB \rightarrow bb$$

$$(6) bE \rightarrow be$$

$$(7) eE \rightarrow ee$$

$$L(G) = \{ a^n b^n e^n \mid n \geq 1 \}$$

• $S \Rightarrow aSBE$	$(S \rightarrow aSBE)$
$\Rightarrow aaBEBE$	$(S \rightarrow aBE)$
$\Rightarrow aabEBE$	$(aB \rightarrow ab)$
$\Rightarrow aabBEE$	$(EB \rightarrow BE)$
$\Rightarrow aabbEE$	$(bB \rightarrow bb)$
$\Rightarrow aabb eE$	$(bE \rightarrow be)$
$\Rightarrow aabbee$	$(eE \rightarrow ee)$

- (1)  $S \rightarrow aSBE$
- (2)  $S \rightarrow aBE$
- (3)  $EB \rightarrow BE$
- (4)  $aB \rightarrow ab$
- (5)  $bB \rightarrow bb$
- (6)  $bE \rightarrow be$
- (7)  $eE \rightarrow ee$

$S \Rightarrow a^n b^n e^n \quad (n \geq 1) ?$

## 复习思考题、作业题：

39页, 2-2 (1), (2)

2-3 (1), (2)

## § 2.3 句型的分析

### § 2.3.1 规范推导和归约

- 1、最左(右)推导:在任一步推导 $V \Rightarrow w$ 中, 都是对符号串 $V$ 的最左(右) 非终结符号进行替换, 称最左 (右) 推导。
- 2、规范推导:即最右推导
- 3、规范句型: 由规范推导所得的句型。
- 4、规范归约: 规范推导的逆过程, 称规范归约或最左归约。

例:  $G[\langle \text{标识符} \rangle]$

$\langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle | \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$| \langle \text{标识符} \rangle \langle \text{数字} \rangle$

$\langle \text{字母} \rangle \rightarrow a | b | \dots | z | A | \dots | Z$

$\langle \text{数字} \rangle \rightarrow 0 | 1 | 2 | \dots | 9$

问题: 给出句子  $a4y$  的规范推导和规范归约.

给出句子  $a4y$  的最左推导.

请注意:规范推导和规范归约互为逆过程.

### 规范推导

$\langle \text{标识符} \rangle$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{字母} \rangle$

$\Rightarrow \langle \text{标识符} \rangle_y$

$\Rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle_y$

$\Rightarrow \langle \text{标识符} \rangle_{4y}$

$\Rightarrow \langle \text{字母} \rangle_{4y}$

$\Rightarrow a_{4y}$

问题:如何正确选择规则?

### 规范归约

$a_{4y}$

$\langle \neq \langle \text{字母} \rangle_{4y} \rangle$

$\langle \neq \langle \text{标识符} \rangle_{4y} \rangle$

$\langle \neq \langle \text{标识符} \rangle \langle \text{数字} \rangle_y \rangle$

$\langle \neq \langle \text{标识符} \rangle_y \rangle$

$\langle \neq \langle \text{标识符} \rangle \langle \text{字母} \rangle \rangle$

$\langle \neq \langle \text{标识符} \rangle \rangle$

问题: 如何准确选择可归约串?