

# 第二章

## 文法和语言

# 本章内容

2.1 符号、符号串及其集合的运算

2.2 文法的直观概念

2.3 文法和语言的形式定义

2.4 文法的类型

2.5 上下文无关文法及其语法树

2.6 句型的分析

2.7 文法实用中的一些说明

## 2.1 符号、符号串及其集合的运算

任何一种语言都是由该语言的基本符号组组成的符号串的集合。

- 基本符号集
- 任何语言的单词符号就是定义在它的字符集上的字符串
- 该语言的任何语句就是定义在其单词符号集上的单词串(符号串)

### 2.1.1 字母表和符号串

### 2.1.2 符号串及其集合的运算

## 2.2.1 字母表和符号串

- 字母表：是元素的非空有穷集合，把字母表中的元素称为符号，因此字母表也称符号集。例
- 符号串：字母表中的符号串组成的任何有穷序列。例

## 2.1.2 符号串及其集合的运算

- ① 符号串的长度：如果某符号串 $x$ 中有 $m$ 个符号，则其长度为 $m$ ，记为 $|x|=m$ 。例
- ② 符号串的联接：设 $x$ 和 $y$ 是符号串，它们的联接 $xy$ 是把 $y$ 的符号写在 $x$ 的符号之后得到的符号串。例
- ③ 符号串的方幂：设 $x$ 是符号串，把 $x$ 自身连接 $n$ 次得到符号串 $z$ ，即 $z=xxx\dots xx$ ，称为符号串 $x$ 的方幂。例

④ 符号串集合：若集合A中的一切元素都是某字母表上的符号串，则称A为字母表上的符号串集合。

● 语言：表示某个确定的字母表上的符号串的任何集合。

⑤ 集合的乘积： $AB = \{xy | x \in A \text{ 且 } y \in B\}$ , 例。

⑥ 集合A的闭包 $A^*$ 和正闭包 $A^+$ ：

$$\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \dots \cup \Sigma^n \dots$$

$$\Sigma^+ = \Sigma^1 \cup \dots \cup \Sigma^n \dots$$

$$\begin{cases} \Sigma^* = \Sigma^0 \cup \Sigma^+ \\ \Sigma^+ = \Sigma \Sigma^* = \Sigma^* \Sigma \end{cases}$$

## 2.2 文法的直观概念

文法的定义：对语言结构的描述和定义，即在形式上用来描述和规定语言结构的称为“文法”(或“语法”)。

例

# 规则

句子的语法结构，可以用一组规则来描述。

规则也称为“产生式”，规则中的“ $::=$ ”也常用“ $\rightarrow$ ”表示。

**注意：**文法中，描述某个特定的语法成分的规则可能不只一条。



# 由规则推导句子

如果用一组规则来描述句子的结构，就可以利用这组规则按照一定的方式去推导产生句子。

**推导方法：**使用一条规则，代替 $\Rightarrow$ 左边的某个符号，产生 $\Rightarrow$ 右端的符号串

## 2.3 文法和语言的形式定义

前面已经对规则(或产生式)的概念进行了非形式化的说明，我们已经对其有了一个直观的了解。下面将对其进行形式化说明，并在此基础上抽象地定义文法和语言。

### 2.3.1 文法的形式定义

### 2.3.2 推导的形式定义

### 2.3.3 语言的形式定义

## 2.3.1 文法的形式定义

① 规则、产生式(或重写规则)：形如  $\alpha \rightarrow \beta$  或  $\alpha ::= \beta$  的  $(\alpha, \beta)$  有序对，且  $\alpha \in V^+$ ， $\beta \in V^*$ ， $V$  为某字母表。

- $\alpha$  称为规则的左部(或产生式的左部)。
- $\beta$  称为规则的右部(或产生式的右部)。

## ②文法的定义

- 文法G定义为四元组  $(V_N, V_T, P, S)$ 
  - $V_N$ ：非终结符集
  - $V_T$ ：终结符集
  - $P$ ：产生式（规则）集合
  - $S$ ：开始符号

$$V_N \cap V_T = \phi, \quad S \in V_N$$

$V = V_N \cup V_T$ ，称为文法G的文法符号集合

# 文法的定义

- 例2.1 文法 $G = (V_N, V_T, P, S)$

$$V_N = \{ S \}, V_T = \{ 0, 1 \}$$

$$P = \{ S \rightarrow 0S1, S \rightarrow 01 \}$$

$S$ 为开始符号

# 文法的定义

- 习惯上只将产生式写出。并有如下约定：
  - 第一条产生式的左部是开始符号
  - 用尖括号括起的是非终结符，否则为终结符。或者大写字母表示非终结符，小写字母表示终结符
  - $G$ 可写成 $G[S]$ ， $S$ 是开始符号

$G: S \rightarrow aAb \quad A \rightarrow ab \quad A \rightarrow aAb \quad A \rightarrow \varepsilon$

$G[S]: A \rightarrow ab \quad A \rightarrow aAb \quad A \rightarrow \varepsilon \quad S \rightarrow aSb$

缩写形式  $G[S]: A \rightarrow ab \mid aAb \mid \varepsilon \quad S \rightarrow aSb$

注意：元符号“ $\mid$ ”读作“或”

例2.2 文法 $G = (V_N, V_T, P, S)$

$V_N = \{\text{标识符}, \text{字母}, \text{数字}\}$

$V_T = \{a, b, c, \dots, x, y, z, 0, 1, \dots, 9\}$

$P = \{ \langle \text{标识符} \rangle \rightarrow \langle \text{字母} \rangle$

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

$\quad \langle \text{标识符} \rangle \rightarrow \langle \text{标识符} \rangle \langle \text{数字} \rangle$

$\quad \langle \text{字母} \rangle \rightarrow a, \dots, \langle \text{字母} \rangle \rightarrow z$

$\quad \langle \text{数字} \rangle \rightarrow 0, \dots, \langle \text{数字} \rangle \rightarrow 9 \quad \}$

$S = \langle \text{标识符} \rangle$

## 2.3.2 推导的形式定义

- 直接推导 “ $\Rightarrow$ ”

$\alpha \rightarrow \beta$  是文法G的产生式, 若有 $v, w$ 满足:  
 $v = \gamma \alpha \delta, w = \gamma \beta \delta$ , 其中  $\gamma \in V^*, \delta \in V^*$   
则称 $v$ 直接推导出 $w$ , 记作  $v \Rightarrow w$

或 $w$ 直接归约到 $v$

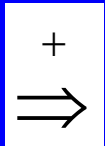
例: G:  $S \rightarrow 0S1, S \rightarrow 01$

$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000S111 \Rightarrow 00001111$

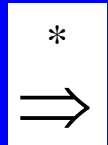
$\langle \text{程序} \rangle \Rightarrow \langle \text{分程序} \rangle. \Rightarrow \langle \text{变量说明部分} \rangle \langle \text{语句} \rangle. \Rightarrow \dots \Rightarrow \text{VAR} \langle \text{标识符} \rangle; \text{BEGIN READ}(\langle \text{标识符} \rangle) \text{END.} \Rightarrow \text{VAR A; BEGIN READ(A) END.}$



# 推导的定义



- 若存在  $v \Rightarrow w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_n = w, (n > 0)$
- 则称  $v \Rightarrow w$ ,  $v$  推导出  $w$ , 或  $w$  归约到  $v$ 。



- 若有  $v \Rightarrow w$ , 或  $v = w$ ,
- 则记为  $v \stackrel{*}{\Rightarrow} w$

# 文法的句型、句子的定义

- 句型

- 有文法G，若 $S \xRightarrow{*} x$ ，则称x是文法G的句型。

- 句子

- 有文法G，若 $S \xRightarrow{+} x$ ，且 $x \in V_T^*$ ，则称x是文法G的句子。

例：G:  $S \rightarrow 0S1, S \rightarrow 01$

$S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000S111 \Rightarrow 00001111$

– 例:  $G[E]: E \rightarrow E+T \mid T$   
 $T \rightarrow T * F \mid F$   
 $F \rightarrow (E) \mid a$

$E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow a+T \Rightarrow a+T * F$   
 $\Rightarrow a+F * F \Rightarrow a+a * F \Rightarrow a+a * a$

表示一切能用符号  $a, +, *, (和)$  构成的算术表达式

## 2.3.3 语言的形式定义

- 由文法G生成的语言记为L(G),它是文法G的一切句子的集合:

$L(G)=\{x|S \xRightarrow{+} x, \text{ 其中 } S \text{ 为文法的开始符号, 且 } x \in V_T^*\}$

例: G:  $S \rightarrow 0S1, S \rightarrow 01$

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

● 例2.3 文法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)  $eE \rightarrow ee$

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

$S \Rightarrow a S B E \quad (S \rightarrow a S B E)$

$\Rightarrow a a B E B E \quad (S \rightarrow a B E)$

$\Rightarrow a a b E B E \quad (a B \rightarrow a b)$

$\Rightarrow a a b B E E \quad (E B \rightarrow B E)$

$\Rightarrow a a b b E E \quad (b B \rightarrow b b)$

$\Rightarrow a a b b e E \quad (b E \rightarrow b e)$

$\Rightarrow a a b b e e \quad (e E \rightarrow e e)$

- G生成的每个串都在 $L(G)$ 中  
 $L(G)$ 中的每个串确实能被G生成

# 文法的等价

- 若 $L(G_1) = L(G_2)$ ，则称文法 $G_1$ 和 $G_2$ 是等价的。

如文法 $G_1[A]$ :  $A \rightarrow 0R$  与 $G_2[S]$ :  $S \rightarrow 0S1$  等价  
 $A \rightarrow 01$   $S \rightarrow 01$   
 $R \rightarrow A1$

## 练习:

- 已知语言描述, 写出文法
  - 例: 若语言由0、1符号串组成, 串中0和1的个数相同, 构造其文法。
    - $A \rightarrow 0B|1C$
    - $B \rightarrow 1|1A|0BB$
    - $C \rightarrow 0|0A|1CC$
- 已知文法, 写出语言描述
  - 例:  $G[E]: E \rightarrow E+T | T$   
 $T \rightarrow T * F | F$   
 $F \rightarrow (E) | a$



## 2.4 文法的类型

- Chomsky将文法分为四种类型：
  - 0型文法（短语文法）： $G=(V_N, V_T, P, S)$ 对任一产生式  $\alpha \rightarrow \beta$ ，都有  $\alpha \in (V_N \cup V_T)^*$ ，且至少含有一个非终结符， $\beta \in (V_N \cup V_T)^*$
  - 1型文法（上下文有关文法）：对任一产生式  $\alpha \rightarrow \beta$ ，都有  $|\beta| \geq |\alpha|$ ，仅仅  $S \rightarrow \varepsilon$  除外
  - 1型文法也可描述为  $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$
  - 2型文法（上下文无关文法）：对任一产生式  $\alpha \rightarrow \beta$ ，都有  $\alpha \in V_N$ ， $\beta \in (V_N \cup V_T)^*$
  - 3型文法（正规文法）：任一产生式  $\alpha \rightarrow \beta$  的形式都为  $A \rightarrow aB$  或  $A \rightarrow a$ ，其中  $A \in V_N$ ， $B \in V_N$ ， $a \in V_T$

# 文法的类型

- 例：1型（上下文有关）文法

文法G[S]:  $S \rightarrow aSBE$

$S \rightarrow aBE$

$EB \rightarrow BE$

$aB \rightarrow ab$

$bB \rightarrow bb$

$bE \rightarrow be$

$eE \rightarrow ee$

# 文法的类型

- 例：1型（上下文有关）文法

文法 $G[S]$ :

$S \rightarrow CD$	$Ab \rightarrow bA$
$C \rightarrow aCA$	$Ba \rightarrow aB$
$C \rightarrow bCB$	$Bb \rightarrow bB$
$AD \rightarrow aD$	$C \rightarrow \varepsilon$
$BD \rightarrow bD$	$D \rightarrow \varepsilon$
$Aa \rightarrow bD$	

$$L(G) = \{ww \mid w \in \{a, b\}^*\}$$

# 文法的类型

- 例：2型（上下文无关）文法

文法G[S]:  $S \rightarrow aB \mid bA$

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

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

文法G[S]:  $S \rightarrow 0A \mid 1B \mid 0$

$A \rightarrow 0A \mid 1B \mid 0S$

$B \rightarrow 1B \mid 1 \mid 0$

# 文法的类型

- 例：定义标识符的3型（正规）文法

文法G[I]:      $I \rightarrow iT$

$I \rightarrow i$

$T \rightarrow iT$

$T \rightarrow dT$

$T \rightarrow i$

$T \rightarrow d$

# 文法和语言

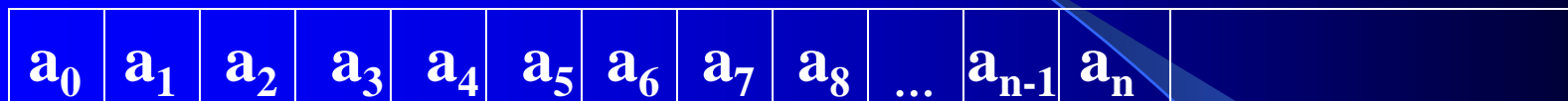
- 0型文法产生的语言称为0型语言
- 1型文法或上下文有关文法（**CSG**）产生的语言称为1型语言或上下文有关语言（**CSL**）
- 2型文法或上下文无关文法（**CFG**）产生的语言称为2型语言或上下文无关语言（**CFL**）
- 3型文法或正则（正规）文法（**RG**）产生的语言称为3型语言正则（正规）语言（**RL**）

# 文法和识别系统

- 0型文法（短语文法）的能力相当于图灵机，可以表征任何递归可枚举集，而且任何0型语言都是递归可枚举的
- 1型文法（上下文有关文法）：产生式的形式为  $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$ ，即只有A出现在  $\alpha_1$  和  $\alpha_2$  的上下文中时，才允许  $\beta$  取代A。其识别系统是线性界限自动机。

# 图灵机

带



磁头

有限控制器



# 文法的类型

- 2型文法（上下文无关文法、CFG）：产生式的形式为 $A \rightarrow \beta$ ， $\beta$ 取代 $A$ 时与 $A$ 的上下文无关。其识别系统是不确定的下推自动机。
- 3型文法（正规文法、右线性文法）：产生的语言是有穷自动机（FA）所接受的集合

## 2.5 上下文无关文法及其语法树

- 上下文无关文法有足够的描述能力描述现今程序设计语言的语法结构
  - 算术表达式
  - 语句
    - 赋值语句
    - 条件语句
    - 读语句
    - .....

## 例：算术表达式上下文无关文法表示

- 文法  $G = (\{E\}, \{+, *, I, (, )\}, P, E)$

P:  $E \rightarrow i$

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow (E)$

## 例：条件语句上下文无关文法表示

$\langle \text{条件语句} \rangle \rightarrow \text{if} \langle \text{条件} \rangle \text{then} \langle \text{语句} \rangle$   
                   $| \text{if} \langle \text{条件} \rangle \text{then} \langle \text{语句} \rangle \text{else} \langle \text{语句} \rangle$

## 2.5.1 上下文无关文法的语法树

- 用于描述上下文无关文法的 句型推导 的直观方法

句型 **aabbbaa** 的语法树（推导树）

例:  $G[S]$ :

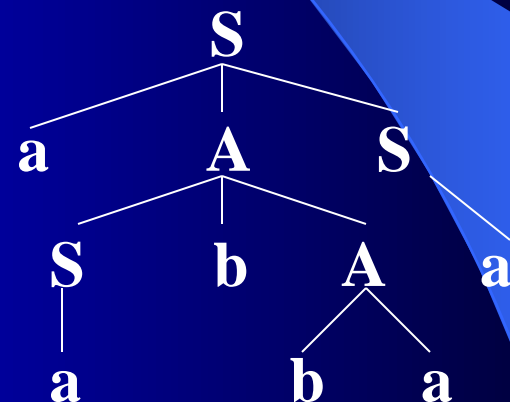
$S \rightarrow aAS$

$A \rightarrow SbA$

$A \rightarrow SS$

$S \rightarrow a$

$A \rightarrow ba$



叶子结点：树中没有子孙的结点。  
从左到右读出推导树的叶子标记，所得的句型为推导树的结果。也把该推导树称为该句型的语法树。

# 上下文无关文法的语法树

- 给定文法G，对于G的任何句型都能构造与之关联的语法树（推导树）。这棵树满足下列4个条件：
  - 1、每个结点都有一个V中的符号作标记
  - 2、根的标记是开始符号S
  - 3、若一结点n至少有一个它自己除外的子孙，并且n有标记A，则 $A \in V_N$
  - 4、如果结点n的直接子孙，从左到右的次序是结点 $n_1, n_2, \dots, n_k$ ，其标记分别为 $A_1, A_2, \dots, A_k$ ，那么 $A \rightarrow A_1 A_2 \dots A_k$ 一定是P中的一个产生式

# 上下文无关文法的语法树

- 推导过程中施用产生式的顺序

例:  $G[S]$ :

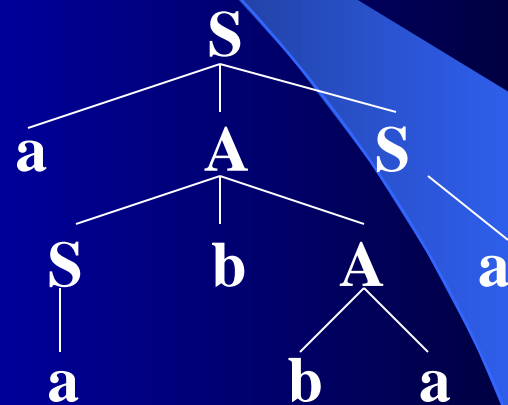
$S \rightarrow aAS$

$A \rightarrow SbA$

$A \rightarrow SS$

$S \rightarrow a$

$A \rightarrow ba$



$S \Rightarrow aAS \Rightarrow aAa \Rightarrow aSbAa \Rightarrow aSbbaa \Rightarrow aabbbaa$

$S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbbaa$

$S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aSbAa \Rightarrow aabAa \Rightarrow aabbbaa$

- 最左（最右）推导：在推导的任何一步  $\alpha \Rightarrow \beta$ ，其中  $\alpha$ 、 $\beta$  是句型，都是对  $\alpha$  中的最左（右）非终结符进行替换
- 最右推导被称为规范推导。
- 由规范推导所得的句型称为规范句型



## 2.5.2 二义性

- 问题：一个句型是否对应唯一的一棵语法树？

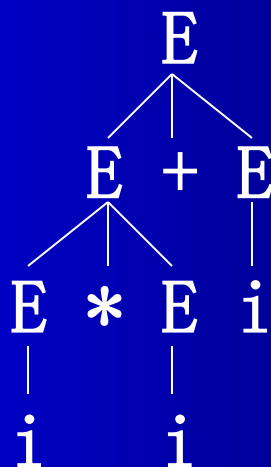
- 例：G[E]：

$E \rightarrow i$

$E \rightarrow E+E$

$E \rightarrow E * E$

$E \rightarrow (E)$



句型  $i*i+i$  的两个不同的最左推导：

推导1：  $E \Rightarrow E+E \Rightarrow E * E + E \Rightarrow i * E + E \Rightarrow i * i + E \Rightarrow i * i + i$

推导2：  $E \Rightarrow E * E \Rightarrow i * E \Rightarrow i * E + E \Rightarrow i * i + E \Rightarrow i * i + i$

# 二义文法

- 若一个文法存在某个句子对应两棵不同的语法树，则称这个文法是二义的。  
或者，若一个文法存在某个句子有两个不同的最左（右）推导，则称这个文法是二义的。
- 产生某上下文无关语言的每一个文法都是二义的，则称此语言是先天二义的。

# 二义文法

- 判定任给的一个上下文无关文法是否二义，或它是否产生一个先天二义的上下文无关语言，这两个问题是递归不可解的。但可以为无二义性寻找一组充分条件（不一定是必要的）。
- 二义文法改造为无二义文法

$G[E]: E \rightarrow i$

$E \rightarrow E+E$

$E \rightarrow E * E$

2023/5/13  $E \rightarrow (E)$

$G[E]: E \rightarrow T | E+T$

$T \rightarrow F | T * F$

$F \rightarrow (E) | i$

规定优先顺序和结合律

## 2.6句型的分析

- 句型分析 就是识别一个符号串是否为某文法的句型，是某个推导的构造过程。
- 在语言的编译实现中，把完成句型分析的程序称为分析程序或识别程序。分析算法又称识别算法。
- 从左到右的分析算法，即总是从左到右地识别输入符号串，首先识别符号串中的最左符号，进而依次识别右边的一个符号。

# 句型的分析

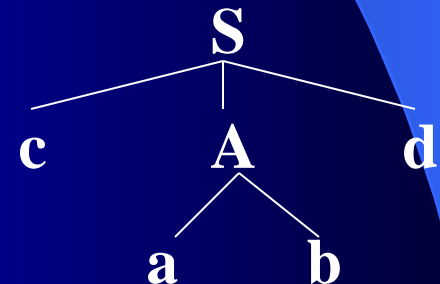
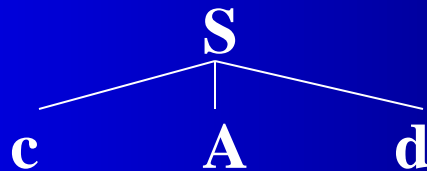
- 分析算法可分为:
- 自上而下分析法:  
从文法的开始符号出发, 反复使用各种产生式, 寻找与输入符号匹配的推导。
- 自下而上分析法:  
从输入符号串开始, 逐步进行归约, 直至归约到文法的开始符号。
- 两种方法反映了两种不同的语法树的构造过程

# 自上而下的语法分析

- 例：文法G： $S \rightarrow cAd$   
 $A \rightarrow ab$   
 $A \rightarrow a$

识别输入串  $w=cabd$  是否该文法的句子

S

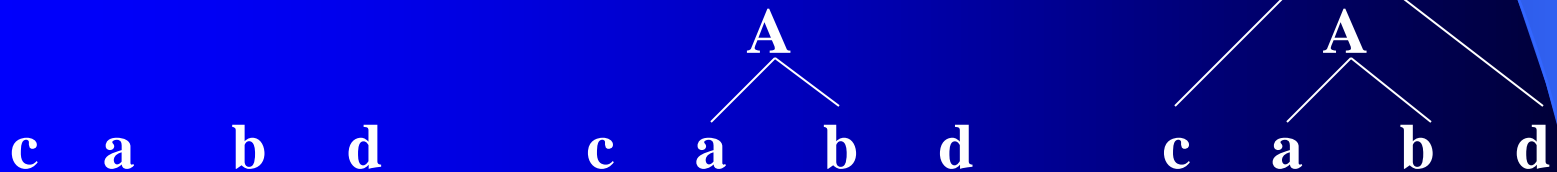


推导过程： $S \Rightarrow cAd \Rightarrow cabd$

# 自下而上的语法分析

- 例：文法G： $S \rightarrow cAd$   
 $A \rightarrow ab$   
 $A \rightarrow a$

识别输入串  $w=cabd$  是否该文法的句子



归约过程： $S \Rightarrow cAd \Rightarrow cabd$

# 句型分析的有关问题

## 1) 如何选择使用哪个产生式进行推导?

假定要被代换的最左非终结符号是 $V$ ，且有 $n$ 条规则： $V \rightarrow A1|A2|...|An$ ，那么如何确定用哪个右部去替代 $V$ ？

## 2) 如何识别可归约的串？

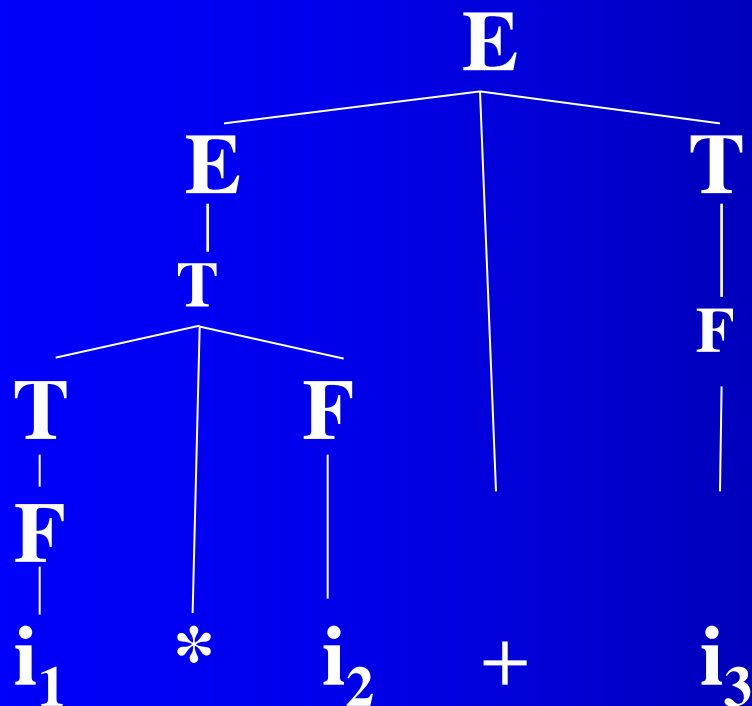
在自下而上的分析方法中，在分析程序工作的每一步，都是从当前串中选择一个子串，将它归约到某个非终结符号，该子串称为“可归约串”



# 句型分析

- 短语、直接短语、句柄的定义：
- 文法 $G[S]$ ,  $S \xRightarrow{*} \alpha A \delta$  且  $A \xRightarrow{+} b$  则称 $b$ 是句型  $\alpha b \delta$  相对于非终结符 $A$ 的短语。若有 $A \Rightarrow b$ 则称 $b$ 是句型  $\alpha b \delta$  相对于该规则 $A \rightarrow b$ 的直接短语。一个句型的最左直接短语称为该句型的句柄。

# 句型分析



$G[E]: E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow (E) \mid i$

句型:  $i * i + i$

短语:

直接短语:

句柄:

# 子树与短语、句柄

- 子树：在一棵语法树中，由某一结点及其所有分支组成的部分称为原树的一棵子树。
- 一棵子树的所有叶子自左至右排列起来形成一个相对子树根的短语。（如上例）
- 一个句型的句柄是这个句型的分析树中最左那棵只有父子两代的子树的所有叶子的自左至右排列。

- 例：设文法 $G[S]$ :
- $S \rightarrow aAS \mid a$        $A \rightarrow SbA \mid SS \mid ba$
- 则aabbbaa是该文法的一个句子，求此句子的短语和句柄

短语：

1、 a1a2b1b2a3a4

2、 a2b1b2a3

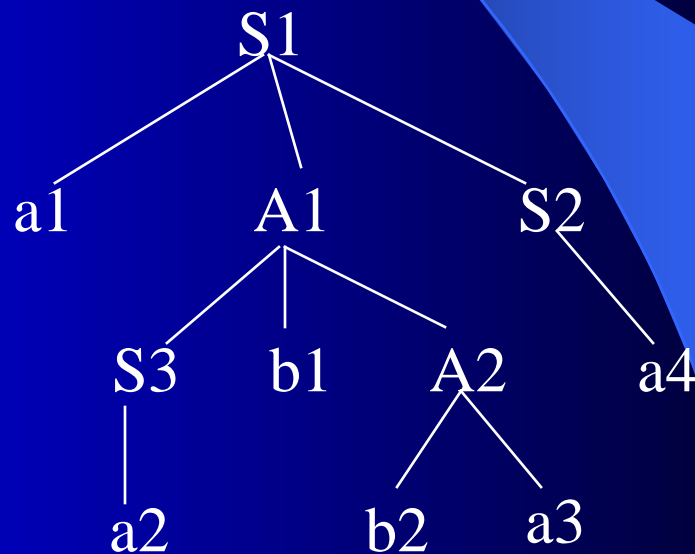
3、 a4

4、 a2

5、 b2a3

句柄：

a2



## 2.7 有关文法实用中的一些说明

- 有关文法的实用限制
- 上下文无关文法中的  $\varepsilon$  规则

# 有关文法的实用限制

- 文法中不得含有有害规则和多余规则
  - 有害规则：形如 $U \rightarrow U$ 的产生式。会引起文法的二义性
  - 多余规则：指文法中任何句子的推导都不会用到的规则
    - 1) 文法中某些非终结符不在任何规则的右部出现，该非终结符称为不可到达
    - 2) 文法中某些非终结符，由它不能推出终结符号串来，称为不可终止的

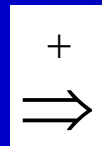
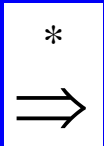
# 有关文法的实用限制

- 对于文法 $G[S]$ ，为了保证任一非终结符 $A$ 在句子推导中出现，必须满足如下两个条件：

- 1)  $A$ 必须在某句型中出现。
- 2) 必须能从 $A$ 推出终结符号串 $t$ 来。

- 即1) 文法 $G[S]$ ，对某两个符号串  $\alpha$  和  $\delta$ ：

$S \xRightarrow{*} \alpha A \delta$     2)  $A \xRightarrow{+} t, t \in V_T^*$



# 化简文法

- 例:  $G[S]$

1)  $S \rightarrow Be$

\*2)  $B \rightarrow Ce$

3)  $B \rightarrow Af$

4)  $A \rightarrow Ae$

5)  $A \rightarrow e$

\*6)  $C \rightarrow Cf$  (不可终止的)

\*7)  $D \rightarrow f$  (不可到达的)

$S \rightarrow Be$

$B \rightarrow Af$

$A \rightarrow Ae$

$A \rightarrow e$



# 上下文无关文法中的 $\varepsilon$ 规则

- 具有形式  $A \rightarrow \varepsilon$  的规则称为  $\varepsilon$  规则，其中  $A \in V_N$
- 某些著作和讲义中限制这种规则的出现。因为  $\varepsilon$  规则会使有关文法的一些讨论和证明变得复杂

# 总结

- 文法的直观概念
- 符号、符号串及其集合的运算
- 文法和语言的形式定义
- 文法的类型
- 上下文无关文法及其语法树
- 句型的分析
- 文法实用中的一些说明