

第2章

形式语言概论

本课程：文法；自动机，正规表达式

作业

- 理论: Cconline.sysu.edu.cn
- 编译器构造实验:
- 转文法
- QQ:597371232
- 实验报告提交: 1967074105@qq.com
- isscjh@mail.sysu.edu.cn

2.1 语言成分

一个语言的成分包括字母表（Alphabet），文法(Grammar)以及它的语义。

本章将主要讨论字母表、符号串、产生式文法系统以及句子分析等方面的内容。

2.1 语言成分

- 字母表与符号：字母表是元素的非空有穷集合。字母表中的元素称为符号。
- 符号串及其运算：
 - ① 符号串。
 - ② 符号串的长度。
 - ③ 符号串的连接。
 - ④ 符号串集合的乘积。
 - ⑤ 符号串的方幂。
 - ⑥ 符号串集合的方幂。
 - ⑦ 符号串集合的正闭包。
 - ⑧ 符号串集合的自反闭包。

字母表与符号串

- **字母表**：是符号的非空有穷集合，用 Σ 表示
- **符号串**：由字母表中的符号所组成的任何有穷序列被称之为该字母表上的符号串

设有字母表 $\Sigma = \{a, b, c\}$ ，那么：

序列 ab 是 Σ 上的一个符号串；

同样序列 ba ，序列 abc ；

序列 $bcca$ 等都是 Σ 上的符号串

符号串的长度 $|s|$

- 在语言的理论中，术语“句子”和“字”常常用作术语“符号串”的同义语。
- 符号串 s 的长度记作 $|s|$ ，是组成该符号串的符号的个数。
- 例如，上述 Σ 上的符号串 ab 的长度是2，记作 $|ab|=2$ 。
- 符号串 abc 的长度是3，记作 $|abc|=3$ 。
- 空符号串记作 ε ，它由零个符号组成，于是 $|\varepsilon|=0$ 。

与符号串有关的几个术语(1)

- 符号串s 的前缀：移走符号串s的尾部的零个或多个符号所得到的一个符号串。

例如ban是符号串banana的一个前缀。

- 符号串 s 的后缀：删去符号串s的头部的零个或多个符号所得到的一个符号串。

例如nana是符号串banana的一个后缀。

- 符号串s的子串：从 s 中删去一个前缀和一个后缀而得到的符号串。

与符号串有关的几个术语(2)

- 符号串 s 的真前缀、真后缀、真子串：任何非空符号串 x ，相应地，是 s 的前缀、后缀或子串，并且 $s \neq x$ 。
- 符号串 s 的子序列：从符号串 s 中删去零个或多个符号(这些符号不要求是连续的)而得到的符号串。
- 术语“语言”表示某个确定的字母表上的符号串的任何集合。

符号串的运算 (1)

- 符号串的连接：设 x, y 是符号串，则 xy 称为 x 与 y 的连接
- 符号串集合的乘积：设 A, B 是符号串集合， AB 表示 A 与 B 的乘积，则定义 $AB = \{xy \mid x \in A, y \in B\}$
- 符号串的方幂。同一符号串的连接可写成方幂形式。设 x 是一符号串，则定义

$$x^0 = \varepsilon$$

$$x^1 = x$$

$$x^2 = xx$$

$$x^3 = x^2x = xx^2 = xxx$$

例如， $x=abc$ ， $x^2=abcabc$ ， $x^3=abcabcabc$

符号串的运算 (2)

- 符号串集合的方幂。同一符号串集合的乘积也可以写成方幂形式设符号串集合A，则定义

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

$$A^1 = A$$

$$A^2 = AA$$

$$A^3 = A^2A = AA^2$$

$$A^n = A^{n-1}A = AA^{n-1}$$

例如， $A = \{a, bc\}$ ， $A^2 = AA = \{aa, abc, bca, bcbc\}$ ，
 $A^3 = A^2A = \{aaa, abca, bcaa, bcbca, aabc, abcbc, bcabc, bcbcbc\}$

符号串的运算 (3)

- 符号串集合的正闭包：设符号串集合A，A的正闭包记为 A^+ ，则有

$$A^+ = A^1 \cup A^2 \cup \dots \cup A^n \cup \dots$$

即 A^+ 为集合A上所有符号串的集合

- 符号串集合的自反闭包：符号串集合A的自反闭包记为 A^* ，则有

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

语言之上的几个重要运算

假设L和M表示两个语言

- 语言L和M的合并(union), 记作 $L \cup M$, 定义为:
$$L \cup M = \{S \mid S \text{ is in } L \text{ or } S \text{ is in } M\}$$
- 语言L和M的连接(concatenation), 记作 LM , 定义为:
$$LM = \{st \mid s \text{ is in } L \text{ and } t \text{ is in } M\}$$
- 语言L的Kleene闭包, 记作 L^* , 定义为:
$$L^* = \bigcup L^i = L^0 \cup L^1 \cup L^2 \cup L^3.$$

文法与语言的形式定义

- 定义 文法G是一个四元组:

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

V_T : 终结符集;

V_N : 非终结符集;

P : 产生式集;

S : 开始符号;

算术表达式

- 表达式是由运算符连接运算量的式子
- 表达式 \rightarrow 表达式+项|表达式-项|项
项 \rightarrow 项 \times 因子|项 \div 因子|因子
因子 \rightarrow 运算量 | (表达式)
- $G[E]=(V_N, V_T, P, E)$
 $V_N=\{E, T, F\}$
 $V_T=\{+, -, *, /, i, (,)\}$
 $P=\{E \rightarrow E+T|E-T|T,$
 $T \rightarrow T*F|T/F|F, F \rightarrow i|(E)\}$

文法及其分类

文法: $G(V_N, V_T, P, S)$

- 0型文法: $P: \alpha \rightarrow \beta$

其中: $\alpha \in (V_N \cup V_T)^+, \beta \in (V_N \cup V_T)^*$

- 1型文法 (上下文有关文法):

$P: \alpha \rightarrow \beta$ 其中: $|\alpha| \leq |\beta|$

或 $\gamma_1 A \gamma_2 \rightarrow \gamma_1 \delta \gamma_2$, $\gamma_1, \gamma_2 \in (V_N \cup V_T)^*$

$A \in V_N$ $\delta \in (V_N \cup V_T)^+$

- 2型文法 (上下文无关文法):

$P: A \rightarrow \delta$ $A \in V_N$ $\delta \in (V_N \cup V_T)^+$

- 3型文法 (右线性文法、正规文法):

$P: A \rightarrow a | aB$ $A, B \in V_N$ $a \in V_T$

0型文法

在2.2.1节中定义的文法即为0型文法(无限制的文法)。其产生式具有以下形式:

$$\alpha \rightarrow \beta$$

其中, $\alpha \in (V_N \cup V_T)^+$, $\beta \in (V_N \cup V_T)^*$

1型文法（上下文有关文法）

1型文法G的产生式具有以下形式： $\alpha \rightarrow \beta$

其中： $\alpha = \gamma_1 A \gamma_2$;

$\beta = \gamma_1 \delta \gamma_2$;

$\gamma_1, \gamma_2 \in (V_N \cup V_T)^*$;

$A \in V_N$; $\delta \in (V_N \cup V_T)^+$ 。

1型文法（上下文有关文法）举例(1)

- 条件 $P: \alpha \rightarrow \beta$ 其中: $|\alpha| \leq |\beta|$
- $G[S] = (V_N, V_T, P, S)$
 $V_N = \{S, A, B, C\}$
 $V_T = \{a, b, c\}$
 $P = \{S \rightarrow aSBC, S \rightarrow aBC, aB \rightarrow ab,$
 $bB \rightarrow bb, bC \rightarrow bc, CB \rightarrow BC, cC \rightarrow cc\}$
- 或条件 $P: \gamma_1 \mathbf{A} \gamma_2 \rightarrow \gamma_1 \mathbf{\delta} \gamma_2$
 $CB \rightarrow BC$
改为 $C\mathbf{B} \rightarrow C\mathbf{D}, \mathbf{CD} \rightarrow \mathbf{BD}, \mathbf{BD} \rightarrow \mathbf{BC}$

1型文法（上下文有关文法） 举例(2)

- $G[S] = (V_N, V_T, P, S)$

$$V_N = \{S, X, Y, Z\}$$

$$V_T = \{x, y, z\}$$

$$P = \{S \rightarrow xSYZ \mid xYZ,$$

$$xY \rightarrow xy,$$

$$yY \rightarrow yy,$$

$$yZ \rightarrow yz,$$

$$ZY \rightarrow YZ,$$

$$zZ \rightarrow zz\}$$

2型文法（上下文无关文法）

在1型文法的产生式中上下文 γ_1 和 γ_2 用空符号串 ε 代替，则有以下形式的产生式称为2型文法：

$$A \rightarrow \delta$$

其中， $A \in V_N$ ， $\delta \in (V_N \cup V_T)^+$ 。

2型文法（上下文无关文法）举例(1)

- 条件 $P: A \rightarrow \delta \quad A \in V_N, \delta \in (V_N \cup V_T)^+$

- $G[E] = (V_N, V_T, P, E)$

$$V_N = \{E, T, F\}$$

$$V_T = \{+, -, *, /, i, (,)\}$$

$$P = \{E \rightarrow E + T \mid E - T \mid T,$$

$$T \rightarrow T * F \mid T / F \mid F, F \rightarrow i \mid (E)\}$$

- 表达式是由运算符连接运算量的式子
- 表达式 \rightarrow 表达式 + 项 | 表达式 - 项 | 项

$$\text{项} \rightarrow \text{项} \times \text{因子} \mid \text{项} \div \text{因子} \mid \text{因子}$$

$$\text{因子} \rightarrow \text{运算量} \mid (\text{表达式})$$

2型文法（上下文无关文法）举例(2)

- $G_4[N]$

P: $N \rightarrow ND \mid D;$

$D \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9;$

- $G_7[S]$

P: $S \rightarrow aTd;$

$T \rightarrow bT \mid cT \mid b \mid c;$

- $G_9[S]$

P: $S \rightarrow 0S1 \mid 01 ;$

3型文法（右线性文法和正规文法）

在正规文法中， P 中的每个产生式($S \rightarrow \varepsilon$ 例外， S 为文法的开始符号)只有两种形式： $A \rightarrow a$ ， $A \rightarrow aB$ 。

其中 $A, B \in V_N$ ， $a \in V_T$ 。此外，如果 $S \rightarrow \varepsilon$ 是 P 中的一个产生式，那么 S 不能出现在任何产生式的右边。

例 正规文法 $G_5(S)$

$$S \rightarrow dB \mid +A \mid -A \mid G$$
$$A \rightarrow dB \mid G$$
$$B \rightarrow dB \mid H \mid d$$
$$G \rightarrow dH$$
$$H \rightarrow dH \mid d$$

其中 d 代表十进制数字。

3型文法（右线性文法、正规文法）举例

- 条件 P: $A \rightarrow a | aB$ $A, B \in V_N$, $a \in V_T$

- G5[S]

P: $S \rightarrow 0 | 1 | 1A | 0B$

$A \rightarrow 1A | 0B$

$B \rightarrow 0 | 1 | 0B$

- G4[N]

P: $N \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |$

$0N | 1N | 2N | 3N | 4N | 5N | 6N | 7N | 8N | 9N;$

推导和规范推导(1)

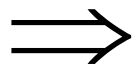
- 定义 2.3 : 文法 $G=(V_N, V_T, P, S)$, $\alpha \rightarrow \beta \in P$, $\gamma, \delta \in (V_N \cup V_T)^*$, 则称符号串 $\gamma\beta\delta$ 为 $\gamma\alpha\delta$ 用产生式 $\alpha \rightarrow \beta$ 的直接推导, 记为

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

+ 多步推导



* 0步或多步推导



推导和规范推导(2)

- **定义2.6:最左推导:**在 $xUy \Rightarrow xuy$ 的直接推导中, 若 $x \in VT^*$, $U \in VN$ —— U 是符号串 xUy 中的最左非终结符, 则称此直接推导为最左直接推导。每一步都为最左直接推导。
- **定义2.7:最右推导:**在 $xUy \Rightarrow xuy$ 的直接推导中, 若 $y \in VT^*$, $U \in VN$ —— U 是符号串 xUy 中的最右非终结符, 则称此直接推导为最右直接推导。每一步都为最右直接推导。
- 最右(直接)推导又称为规范推导

最左推导与最右推导举例



- 文法 $G[E]$

$E \rightarrow E+T | T;$

$T \rightarrow T * F | F;$

$F \rightarrow i | (E);$

- 写出 $i+(i+i)$ 的最左推导与最右推导

- 解：最左推导： $E \Rightarrow E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow i+T$

$\Rightarrow i+F \Rightarrow i+(E) \Rightarrow i+(E+T) \Rightarrow i+(T+T) \Rightarrow i+(F+T)$
 $\Rightarrow i+(i+T) \Rightarrow i+(i+F) \Rightarrow i+(i+i)$

最右推导： $E \Rightarrow E+T \Rightarrow E+F \Rightarrow E+(E) \Rightarrow E+(E+T)$
 $\Rightarrow E+(E+F) \Rightarrow E+(E+i) \Rightarrow E+(T+i) \Rightarrow E+(F+i)$
 $\Rightarrow E+(i+i) \Rightarrow T+(i+i) \Rightarrow F+(i+i) \Rightarrow i+(i+i)$

最左推导与最右推导举例

1. 确定下列文法的语言

$G[S]: S \rightarrow ABC | C;$

$A \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9;$

$B \rightarrow N | BN | \varepsilon$

$C \rightarrow 0 | 5$

$N \rightarrow A | 0$

2. 已知文法

$G[Z]: Z \rightarrow UV | VU$

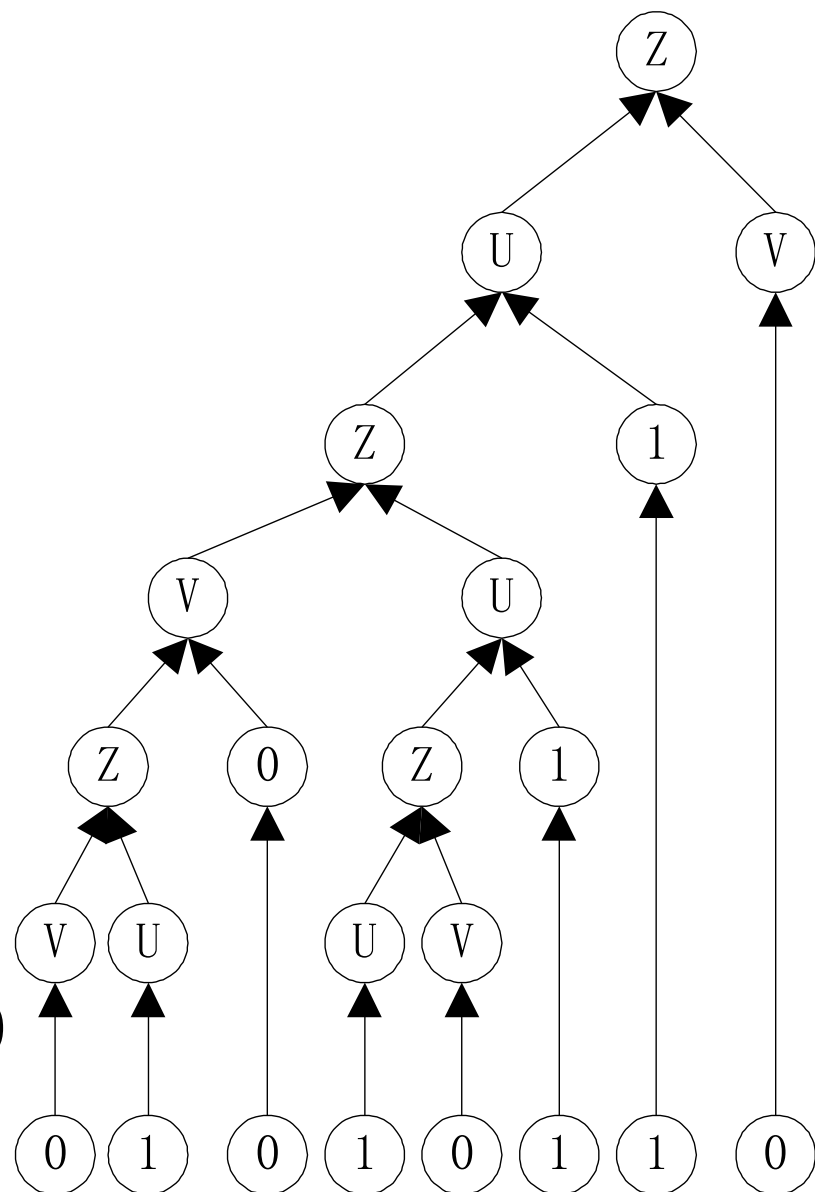
$U \rightarrow Z1 | 1$

$V \rightarrow Z0 | 0$

试写出01010110的最左推导和最右推导

- $$V \rightarrow Z0 \mid 0$$

• 解：最左推导： $Z \Rightarrow UV \Rightarrow Z1V$

$$\Rightarrow 0U0U1V \Rightarrow 010U1V \Rightarrow 010Z11V$$
$$\Rightarrow 0101011V \Rightarrow 01010110$$
$$\Rightarrow VU10 \Rightarrow VZ110 \Rightarrow VUV110 \Rightarrow VU0110$$
$$\Rightarrow V1010110 \Rightarrow 01010110$$


- $$V \rightarrow Z0 \mid 0$$

[illegible]

句型、句子和语言

• **定义2.8:** 设 S 是文法 G 的开始符号, 如果 $S \xRightarrow{*} u$,
 $u \in (V_N \cup V_T)^*$, 称 u 为文法 G 的句型。

• **定义2.9:** 设 S 是文法 G 的开始符号, 如果 $S \xRightarrow{*} u$,
 $u \in V_T^*$, 称 u 为文法 G 的句子。

• **定义2.10:** 设 S 是文法 G 的开始符号, 文法 G 的
语言

$$L(G) = \{u \mid S \xRightarrow{*} u, u \in V_T^*\}$$

短语与句柄

• 定义6.1: 设 S 是文法 G 的开始符号,

若
$$S \overset{*}{\Rightarrow} xUy \quad U \overset{+}{\Rightarrow} \alpha$$

$U \in V_N$, $x, y \in (V_N \cup V_T)^*$, 称 α 是句型 xUy 相对于 U 的短语。

又若
$$S \overset{*}{\Rightarrow} xUy \quad U \Rightarrow \alpha$$

$U \in V_N^*$, $x, y \in (V_N \cup V_T)^*$, 称 α 是句型 xUy 相对于 U 的直接短语或简单短语。

一个句型的最左直接短语称该句型的句柄。

构造下列语言的文法

- $\{ a^{3n} | n \geq 1 \};$

解: $G[S]: S \rightarrow aaaS | aaa$

- 偶数集合

解: $G[N]: N \rightarrow AB$

$$A \rightarrow DA | \varepsilon$$
$$D \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$$
$$B \rightarrow 0 | 2 | 4 | 6 | 8$$

- $\{ a^n b^m c^k | n, m, k \geq 0 \}$

解: $G[S]: S \rightarrow ABC$

$$A \rightarrow aA | \varepsilon \quad B \rightarrow bB | \varepsilon \quad C \rightarrow cC | \varepsilon$$

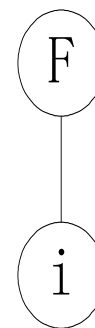
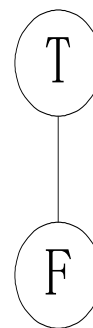
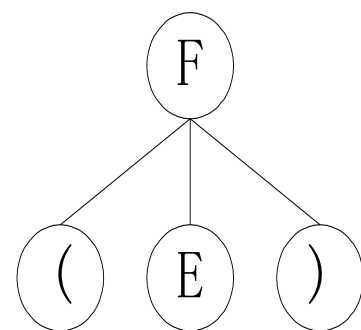
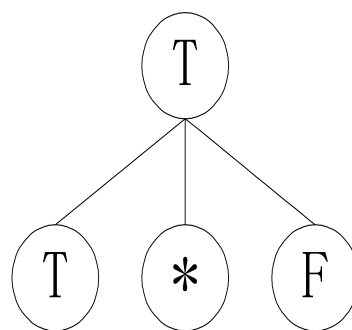
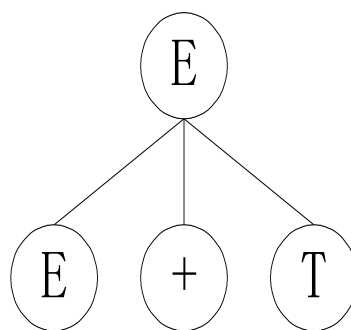
产生式树举例

- 文法 $G[E]$

$E \rightarrow E+T \mid T;$

$T \rightarrow T * F \mid F;$

$F \rightarrow i \mid (E);$



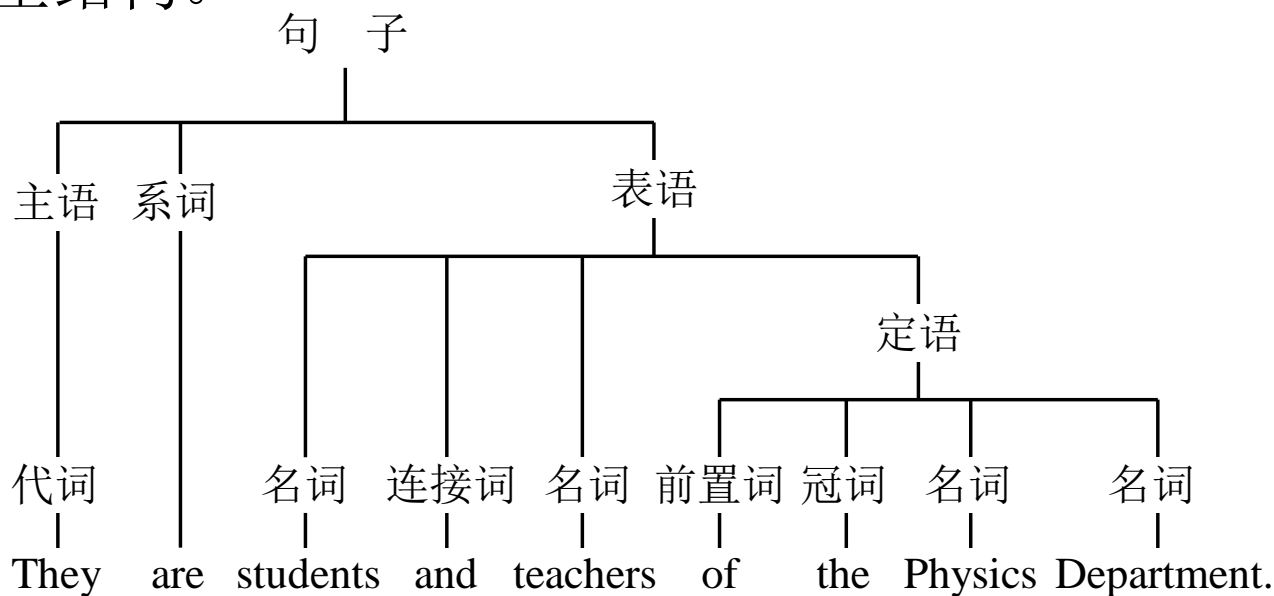
语法树

设文法 $G=(V_T, V_N, P, S)$ ，对于文法 G 的任意一个句型都存在一个相应的语法树：

- ① 树中每个结点都有一个标记，该标记是 $V_N \cup V_T$ 中的一个符号；
- ② 树的根结点标记是文法的识别符号 S ；
- ③ 若树的一个结点至少有一个叶子结点，则该结点的标记一定是一个非终结符；
- ④ 若树的一个结点有多个叶子结点，该结点的标记为 A ，这些叶子结点的标记从左到右分别是 B_1, B_2, \dots, B_N ，则 $(A, B_1, B_2 \dots B_N) \in P$ 。

语法树

在自然语言中，可通过树型表示直观地分析句子结构；在形式语言中，则是通过语法树直观地分析文法的句型结构。



句子结构

语法树举例

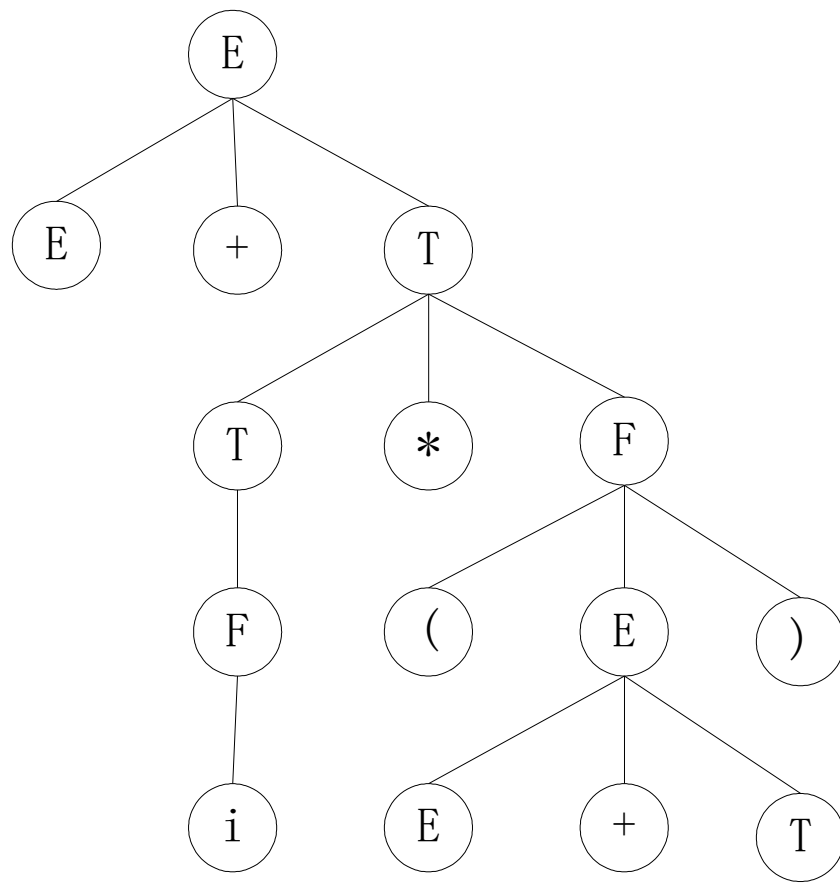
- 文法 $G[E]$

$E \rightarrow E+T \mid T;$

$T \rightarrow T*F \mid F;$

$F \rightarrow i \mid (E);$

- $E+i*(E+T)$



产生式树

- 文法的句型可依据文法的产生式树来生成相应的语法树。

产生式树

文法的句型都可依据文法的产生式来生成相应的语法树。以句型 $E+(E+T)*i$ 为例：

- ① 以文法G的识别符号E作为语法树根结点的标记。选择识别符E的一个产生式 $E \rightarrow E+T$ ，然后生成根结点E的3个分支，根结点E的3个叶子结点的标记，从左到右分别记为E、+和T。
- ② 选择产生式 $T \rightarrow T * F$ ，生成以结点T作为根结点的子树。
- ③ 选择产生式 $F \rightarrow i$ ，以②中子树最右边的叶子结点F为根结点，延伸相应的子树。
- ④ 选择产生式 $T \rightarrow F$ ，以③中子树的叶子结点T为根结点，延伸相应的子树。
- ⑤ 选择产生式 $F \rightarrow (E)$ ，以④中子树的叶子结点F为根结点，延伸相应的子树，扩充相应的子树。

无用 V_N 、不可达 V_N

- 文法中规则的个数应该是恰当的。再少，不足以完全描述一个语言，再多，又无必要。
- 程序设计语言中存在有多余规则时，往往存在有错误
- 多余规则使句子的分析复杂和增加困难度。

明显的多余规则有两种。

- 一是形如 $U \rightarrow U$ 的单规则，引起文法上的二义性。应该首先把它们删除去。
- 另一是规则 $U \rightarrow u$ ，其左都非终结符号 U 不出现在其它任何规则的右部，也应删除。（ U 例外）。

无用 V_N 、不可达 V_N 举例

例如，对于文法 $G [Z]$:

$$Z \rightarrow Be$$
$$A \rightarrow Ae \mid A \mid e$$
$$B \rightarrow Ce \mid Af$$
$$C \rightarrow Cf$$
$$D \rightarrow f$$

首先删除规则 $A \rightarrow A$;

又由于 D 不出现在任何规则右部，应删去规则 D 。

规则不多余的条件: U

*

•条件1 $S \Rightarrow xUy \quad x, y \in (V_N \cup V_T)^*$

+

•条件2 $U \Rightarrow \alpha \quad \alpha \in V_T^*$ 。

条件1要求U在句型中出现, 条件2则进一步要求能从U推导出终结符号串。

显然, 如果存在某个非终结符号U, 它不满足上述两条件, 则以U为左部的那些规则必是多余的。

二义性

- 定义2.13 一个文法，如果它的一个句子有两棵或两棵以上的语法树，则称此句子具有二义性。如果一个文法含有二义性的句子，则该文法具有二义性。

分析方法简介

➤ 一个分析器或分析自动机是这样一种系统，它能够根据给定的文法 G ，构造语言 $L(G)$ 的任意推导。分析也可看作是语法树的构造过程。我们主要讨论右线性文法和CFG的分析器。

➤ 分析的方法很多，可归纳为两类，一类是自上而下分析方法，另一类是自下而上分析方法。

2.6.1 自上而下分析方法

自上而下分析方法的基本思想是从文法的开始符号出发，利用其中的产生式，逐步推导出待分析的符号串。如果能推导出这个符号串，则表明此符号串是该文法的一个句型或句子，否则便不是。

2.6.2 确定的自上而下分析方法

- 当文法的某一个非终结符有几条产生式、而且每条产生式右部首符号都是终结符时，应保证它们是互不相同的终结符。

- 例 设文法 $G_{18}[S]$:

$$S \rightarrow aBc \mid bCd$$

$$B \rightarrow eB \mid f$$

$$C \rightarrow dC \mid c$$

试检查符号串 $aefc$ 是不是该文法的句子。

2.6.2 确定的自上而下分析方法

上例推导过程：

$$S \rightarrow aBc, S \Rightarrow aBc ;$$

$$B \rightarrow eB, S \Rightarrow aBc \Rightarrow aeBc ;$$

$$B \rightarrow f, S \Rightarrow aBc \Rightarrow aeBc \Rightarrow aefc ;$$

该例属确定的自上而下分析方法。

2.6.3 自下而上分析方法

自下而上分析方法的基本思想是从待检查的符号串出发，看最终是否能归约（推导的逆过程）到文法的识别符号。如果能归约到文法的识别符号，则表明此待检查的符号串是该文法的一个句型或句子，否则便不是。

2.6.4 文法在内存中的表示

用语法图的表格结构表示文法。一个文法的语法图由该文法所有非终结符号的定义图组成。每个非终结符号的定义图是一个结构型数据。

在自上而下分析方法中，用这种语法图表示文法有利于消除左递归，也有助于提取左因子。

分析方法简介

- 自上而下分析方法
- 确定性的自上而下分析方法
- 自下而上分析方法