

华中科技大学

编译原理

第2章 文法和语言

2019年11月4日星期一

编译原理课程组

内容摘要

本章介绍形式语言理论的基本内容，重点讨论文法和语言的概念、上下无关文法及其句型分析的基本问题。

形式语言与自动机理论是计算机科学的基础理论之一，是构造词法分析程序和语法程序的理论依据，也是学习编译原理的核心和基础内容。

关于自动机理论基本内容将在第三章讨论。

内容摘要

本章介绍形式语言理论的基本内容，重点讨论文法和语言的概念、上下无关文法及其句型分析的基本问题。

形式语言与自动机理论是计算机科学的基础理论之一，是构造词法分析程序和语法程序的理论依据，也是学习编译原理的核心和基础内容。

关于自动机理论基本内容将在第三章讨论。

重点讲解

2.1 文法的直观概念

2.2 符号和符号串

2.3 文法和语言的形式定义

2.4 文法类型

2.5 上下文无关文法及其语法树

2.6 句型分析

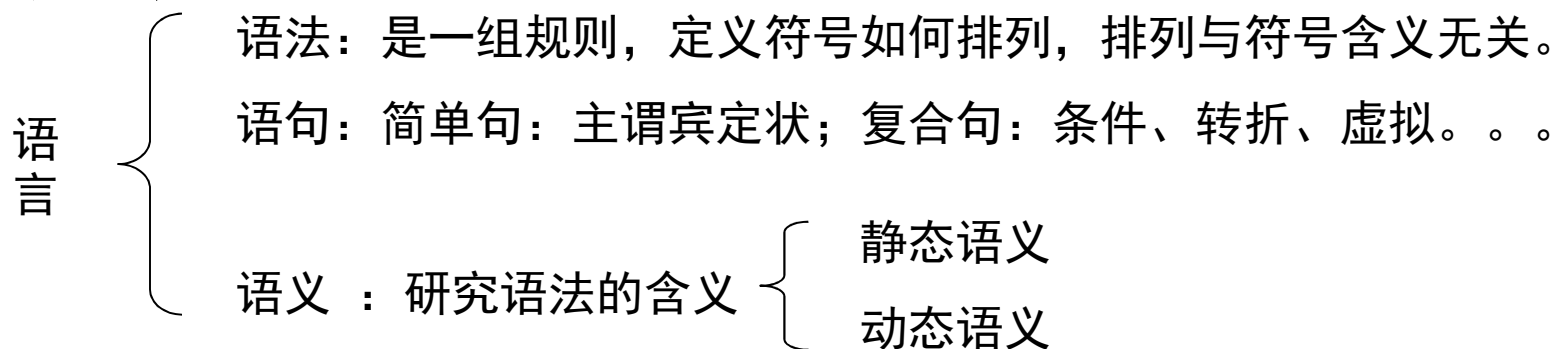
2.7 文法在实用中的一些说明

2.1 文法的直观概念

从语言结构的角度看，组成语言的基本形式是句子，句子是由单词序列构成的，单词是由语言基本符号（字母或单字）组成的。

语言既包含单词和句子这样的语言成分，又包含将这些成分组织起来的语言规则，如**词法规则**、**句法规则**等。

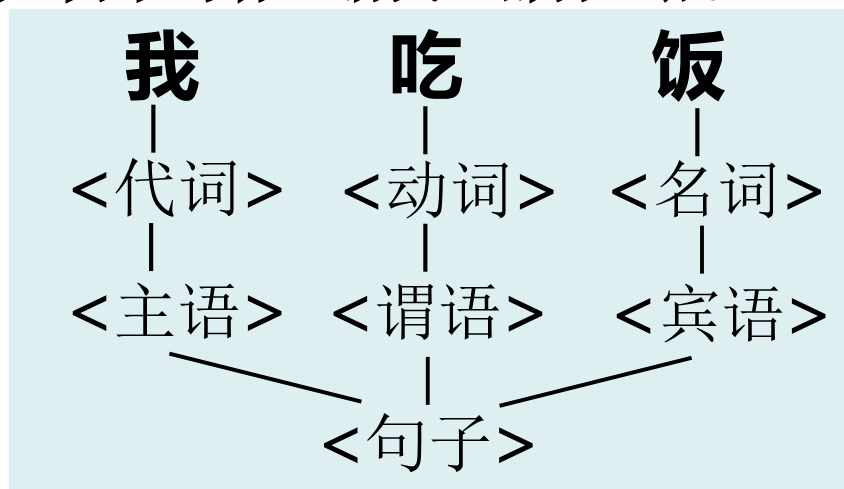
下面以自然语言为例，说明如何对语言规则进行形式化描述的基本思路。



文法是阐述语法的一个工具，语句是语法的实例

2.1 文法的直观概念

考虑句子：“我吃饭”。它的语法成分构成如下，显然是符合中文语法规则——简单句有主谓宾三部分组成。



中文语法规则可以采用如下方式陈述：

<句子>是由<主语> <谓语> <宾语>组成

<主语>是由<名词>组成

... ..

2.1 文法的直观概念

约定(1)：符号 $::=$ 的意义如下，

“... $::=$...” 表示 “...是由...组成的”。

语法规则可以形式化描述如下：

$\langle \text{句子} \rangle ::= \langle \text{主语} \rangle \langle \text{谓语} \rangle \langle \text{宾语} \rangle$
 $\langle \text{主语} \rangle ::= \langle \text{名词} \rangle$
 $\langle \text{主语} \rangle ::= \langle \text{代词} \rangle$
 $\langle \text{谓语} \rangle ::= \langle \text{动词} \rangle$
 $\langle \text{宾语} \rangle ::= \langle \text{名词} \rangle$
 $\langle \text{宾语} \rangle ::= \langle \text{代词} \rangle$
 $\langle \text{代词} \rangle ::= \text{我}$
 $\langle \text{代词} \rangle ::= \text{你}$
 $\langle \text{动词} \rangle ::= \text{吃}$
 $\langle \text{动词} \rangle ::= \text{做}$
 $\langle \text{名词} \rangle ::= \text{饭}$
 $\langle \text{名词} \rangle ::= \text{菜}$

注解：

语法概念加上 $\langle \dots \rangle$ 表示是可替换部分，目的是与构成语句的“单词”加以区别。

2.1 文法的直观概念

约定(2)：符号 | 的表示“或者”的意义。

语法规则可以简化描述如下：

<句子> ::= <主语> <谓语> <宾语>
<主语> ::= <名词> | <代词>
<谓语> ::= <动词>
<宾语> ::= <名词> | <代词>
<代词> ::= 我 | 你
<动词> ::= 吃 | 做
<名词> ::= 饭 | 菜

注解：<主语> ::= <名词> | <代词> 是两个规则的简写，
仍然表示两个规则。

2.1 文法的直观概念

约定(3): 符号 \Rightarrow 表示“推导”。

依据某一个规则, 将被替换的内容中出现的某一个与该规则左部相同的部分, 用该规则右部替换形成另一个新内容, 这个过程叫做一步“推导”。

使用推导符号 \Rightarrow , 一步推导可以描述成: 原内容 \Rightarrow 新内容
采用“多步推导过程”, 可以表示语法分析过程。

$\langle \text{句子} \rangle \Rightarrow \langle \text{主语} \rangle \langle \text{谓语} \rangle \langle \text{宾语} \rangle$
 $\Rightarrow \langle \text{代词} \rangle \langle \text{谓语} \rangle \langle \text{宾语} \rangle$
 $\Rightarrow \text{我} \langle \text{谓语} \rangle \langle \text{宾语} \rangle$
 $\Rightarrow \text{我} \langle \text{动词} \rangle \langle \text{宾语} \rangle$
 $\Rightarrow \text{我吃} \langle \text{宾语} \rangle$
 $\Rightarrow \text{我吃} \langle \text{名词} \rangle$
 $\Rightarrow \text{我吃饭}$

注解:

- ①推导过程不唯一
- ②推导起点的不同, 导致语法意义上差异的推导结果

2.1 文法的直观概念

假设将“**语法概念**”用词<句子>、<主语>、<谓语>、<宾语>、<代词>、<动词>和<名词>，分别用A、B、C、D、E、F和G表示，
“**单词**”本身：我、你、吃、做、饭和菜，分别用a、b、c、d、e和f表示，则句子、规则和推导过程，可以进一步抽象而形式化。

规则：

$A::=BCD$

$B::=G \mid E$

$C::=F$

$D::=G \mid E$

$E::=ab$

$F::=cd$

$G::=e \mid f$

推导过程：

$A \Rightarrow BCD$

$\Rightarrow ECD$

$\Rightarrow aCD$

$\Rightarrow aFD$

$\Rightarrow acD$

$\Rightarrow acG$

$\Rightarrow ace$

2.1 文法的直观概念

语法形式化方法要点：

- 语法规则的形式化
- 语法规则含有语法单位符号
- 语法规则含有构成语句的单词符号
- 特殊的语法单位符号——开始符号^[注]

综上所述，语法形式化的最终目的在于将语法分析的问题将转换成形式化的推导过程。

[注]：从不同的符号开始推导，将得到完全不同语法意义上的结果。

2.2 符号和符号串

2.2.1 基本概念

- **字母表** 字母表 Σ 是非空有穷集合，其元素称为符号。
- **符号串** 由字母表 Σ 中的符号组成的有穷序列称为（字母表 Σ 上的）符号串。特别地，不含任何符号的有穷序列称为空串，记为 ϵ 。**单词和源程序都是符号串！**

例：

设字母表 $\Sigma = \{0, 1\}$ ，则
101是 Σ 上的符号串，
201不是 Σ 上的符号串。

2.2.1 基本概念

- **符号串长度** 符号串 α 的长度是指符号串 α 中含有符号的个数，记为 $|\alpha|$ 。特别约定，空串 ε 为零，即 $|\alpha| = 0$ 。
- **符号串集合** 如果集合 A 的元素都是字母表 Σ 上的符号串，则称集合 A 为 Σ 上的符号串集合，简称串集。

例:

设字母表 $\Sigma = \{a, b, c\}$ ， $A = \{\varepsilon, a, ba, cab\}$ ， $B = \{a1, ba, cab\}$ ，则

A 是 Σ 上的符号串集合，

B 不是 Σ 上的符号串集合。

2.2.2 基本运算

- **符号串连接运算** 设 x 和 y 是字母表 Σ 上的符号串，在符号串 x 的最后一个符号之后顺序接上符号串 y 的符号得到的新符号串 z ，则称符号串 z 是由符号串 x 和符号串 y 经过连接运算的结果，记为 $z = x \cdot y$ ，其中， \cdot 是连接运算符。

例：设字母表 $\Sigma = \{a, b, c, 0, 1\}$ ，

$x = abc, y = 01cba,$

则 $z = x \cdot y = abc01cba$

2.2.2 基本运算

- **符号串方幂运算** 设 x 是字母表 Σ 上的符号串, z 是由 $n(\geq 0)$ 个 x 自身连接得到的符号串, 则称符号串 z 是由符号串 x 的 n 次方幂运算的结果, 记为 $z = x^n$ 。特别约定, $x^0 = \varepsilon, x^1 = x$ 。

讨论: $x \cdot x = x^2, x \cdot x \cdot x = x^3, \dots, \underbrace{x \cdot x \cdot \dots \cdot x}_{n \uparrow x} = x^n$

- **符号串集连接运算** 设 A, B 是字母表 Σ 上的符号串集, \cdot 是符号串集连接运算, 则 $C = A \cdot B = \{x \cdot y \mid x \in A, y \in B\}$ 。笛卡尔积

- **符号串集方幂运算** 设 A 是字母表 Σ 上的符号串集, 则 C 是由 $n(\geq 0)$ 个 A 自身连接得到的符号串集, 则称符号串集 C 是由符号串 A 的 n 次方幂运算的结果, 记为 $C = A^n$ 。特别约定, $A^0 = \{\varepsilon\}, A^1 = A$ 。

2.2.2 基本运算

- **符号串集正闭包运算** 设A是字母表 Σ 上的符号串集, A^+ 是A的正闭包, 则: $A^+ = A^1 \cup A^2 \cup A^3 \cup \dots \cup A^n \dots$.
- **符号串集闭包运算** 设A是字母表 Σ 上的符号串集, A^* 是A的闭包, 则: $A^* = A^0 \cup A^+$,

即: $A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \dots \cup A^n \dots$

例: 设字母表 $\Sigma = \{a, b\}$, $A = \{aa, bb\}$, $B = \{a\}$, $C = \{a\}$, 则

$$A \cdot B = \{aaab, aaba, bbab, bbba\},$$

$$A^2 = \{aaaa, aabb, bbaa, bbbb\},$$

$$C^+ = \{a, aa, aaa, \dots\} = \{a^n \mid n \geq 1\},$$

$$C^* = \{\epsilon, a, aa, aaa, \dots\} = \{a^n \mid n \geq 0\}.$$

语言是单
词集合的
连接

源程序就
是语言的
闭包

2.3 文法和语言的形式定义

规则是字母表 V 上形如 $\alpha ::= \beta$ 的式子，可以简写成 $\alpha \rightarrow \beta$ 。其中，符号串 $\alpha \in V^+$ 称为规则的左部，符号串 $\beta \in V^*$ 称为规则的**右部**。**规则**也称为**重写规则**、**产生式**或**生成式**。

特别地， $\alpha ::= \epsilon$ (ϵ 空串) 称为 α 的**空规则**。

对于相同左部的多个规则，可以使用符号 $|$ 简写。

如，规则 $\alpha ::= \beta$ 和 $\alpha ::= \delta$ ，简写成 $\alpha ::= \beta | \delta$ 。简写为 $\alpha \rightarrow \beta | \delta$

定义2.1 文法

文法G定义为一个四元组 (V_N, V_T, P, S) ，记为 $G = (V_N, V_T, P, S)$ 。其中，

- ① V_N 是非空有穷集合，称为**非终结符集**，其元素称为非终结符；
- ② V_T 是有穷集合，称为**终结符集**，其元素称为终结符；
- ③ P 是非空有穷集合，称为**规则集**，其元素是字母表 $V_N \cup V_T$ 上的规则， $V_N \cup V_T$ 称为文法的字母表 V ，且 $V_N \cap V_T = \Phi$ ；
- ④ $S \in V_N$ ，称为**开始符**。

2.3 文法和语言的形式定义

例2.1 定义文法G1如下:

$$\begin{aligned} G1 &= (V_N, V_T, P, S), \\ \text{其中, } V_N &= \{S\}, \\ V_T &= \{a, b\}, \\ P &= \{S \rightarrow aSb, S \rightarrow ab\} \end{aligned}$$

通常, 文法还可以采用其它形式给出定义。如文法G1可以写成其它形式如下。

$$\begin{aligned} G1 &= (\{S\}, \{a, b\}, P, S), \\ \text{其中, } P &= \{S \rightarrow aSb, S \rightarrow ab\} \end{aligned}$$

或者写成:

$$\begin{aligned} G1 &= (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S) \\ G1[S]: S &\rightarrow aSb, S \rightarrow ab \end{aligned}$$

定义2.2 直接推导、直接归约

设文法 $G = (V_N, V_T, P, S)$ ，如果 $\alpha \rightarrow \beta \in P$ ，则称 $\gamma \alpha \delta$ 推导出 $\gamma \beta \delta$ ，记为 $\gamma \alpha \delta \Rightarrow \gamma \beta \delta$ ，其中， $\gamma, \delta \in V^*$ 。

$\gamma \alpha \delta \Rightarrow \gamma \beta \delta$ 也称为直接推导或一步推导。

如果 $\gamma \alpha \delta \Rightarrow \gamma \beta \delta$ ，则也称为 $\gamma \beta \delta$ 归约到 $\gamma \alpha \delta$ ，也称为直接归约或一步归约。

例如，例3.1 定义的文法 $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$ ，推导例子有：

- (1) $S \Rightarrow aSb$ ($\alpha = S, \beta = aSb, \gamma = \varepsilon, \delta = \varepsilon$)
- (2) $aSb \Rightarrow aaSbb$ ($\alpha = S, \beta = aSb, \gamma = a, \delta = b$)
- (3) $aSb \Rightarrow aabb$ ($\alpha = S, \beta = ab, \gamma = a, \delta = b$)
- (4) $aSbSb \Rightarrow aaSbbSb$ ($\alpha = S, \beta = aSb, \gamma = a, \delta = bSb$)

定义2.3 多步推导、多步归约

设文法 $G = (V_N, V_T, P, S)$ ， $\alpha, \beta \in (V_N \cup V_T)^*$ ，如果 α, β 之间存在推导序列：

$\alpha = W_0 \Rightarrow W_1 \Rightarrow W_2 \cdots \Rightarrow W_n = \beta \quad (n \geq 1)$ ，
则称 α 经过**n步推导**出 β ，记为 $\alpha \Rightarrow^+ \beta$ 。其中， $W_i \in (V_N \cup V_T)^*$
($1 \leq i \leq n$)。 $\alpha \Rightarrow^+ \beta$ 也称**n步推导或多步推导**。

如果 $\alpha \Rightarrow^+ \beta$ ，也称为 β **归约到** α ，也称为**n步归约或多步归约**。

例如，例3.1 定义的文法 $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$ ，多步推导(\Rightarrow^+)例子有：

(1) $S \Rightarrow^+ ab \quad (\because S \Rightarrow ab)$

(2) $S \Rightarrow^+ aabb \quad (\because S \Rightarrow aSb \Rightarrow aabb)$

(3) $S \Rightarrow^+ aaaSbbb \quad (\because S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb)$

(4) $aSb \Rightarrow^+ aaabbb \quad (\because aSb \Rightarrow aaSbb \Rightarrow aaabbb)$

定义2.4 0步或0步以上推导与归约

设文法 $G = (V_N, V_T, P, S)$ ， $\alpha, \beta \in (V_N \cup V_T)^*$ ，如果有 $\alpha \rightarrow \beta$ 或 $\alpha \xRightarrow{+} \beta$ ，则称 α 经过0步或0步以上推导出 β ，记为 $\alpha \xRightarrow{*} \beta$ 。亦称 β 经过0步或0步以上归约到 α 。

例如，例3.1 定义的文法 $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$ ，0步或0步以上推导($\xRightarrow{*}$)例子有：

(1) $S \xRightarrow{*} ab$ ，因为有 $S \xRightarrow{+} ab$

(2) $S \xRightarrow{*} aabb$ ，因为有 $S \xRightarrow{+} aabb$

(3) $S \xRightarrow{*} aaabbbb$ ，因为有 $S \xRightarrow{+} aaabbbb$

(4) $aSb \xRightarrow{*} aaabbbb$ ，因为有 $aSb \xRightarrow{+} aaabbbb$

(5) $aSbSb \xRightarrow{*} aSbSb$ ，因为有 $aSbSb \xRightarrow{+} aSbSb$

定义2.5 句型、句子

设文法 $G = (V_N, V_T, P, S)$ ，如果有 $S \xRightarrow{*} \beta$ ，则称 β 是文法 G 的句型。如果有 $S \Rightarrow^* \beta$ ，且 $\beta \in V_T^*$ ，则称 β 是文法 G 的句子。

例如，例3.1 定义的文法 $G_1 = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$ ，句型和句子例子有：

(1) ab 是 G 的句子，因为有 $S \xRightarrow{*} ab$ ， $ab \in V_T^*$

(2) $aabb$ 是 G 的句子，因为有 $S \xRightarrow{*} aabb$ ， $aabb \in V_T^*$

(3) $aaaSbbb$ 是 G 的句型，因为有 $S \xRightarrow{*} aaaSbbb$ ($aaaSbbb \notin V_T^*$)

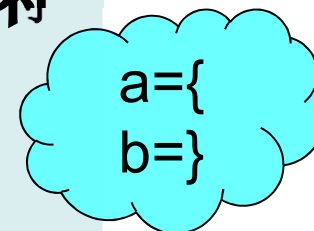
定义2.6 语言

文法 $G = (V_N, V_T, P, S)$ 的产生语言定义为文法 G 的句子集合，记为 $L(G)$ 。即：

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

例如，例3.1 定义的文法 G_1 ，其语言 $L(G)$ 是由相同个数 (≥ 1) 的 a 和 b 符号、且 b 符号均在 a 符号之后出现的符号串构成的集合。即：

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



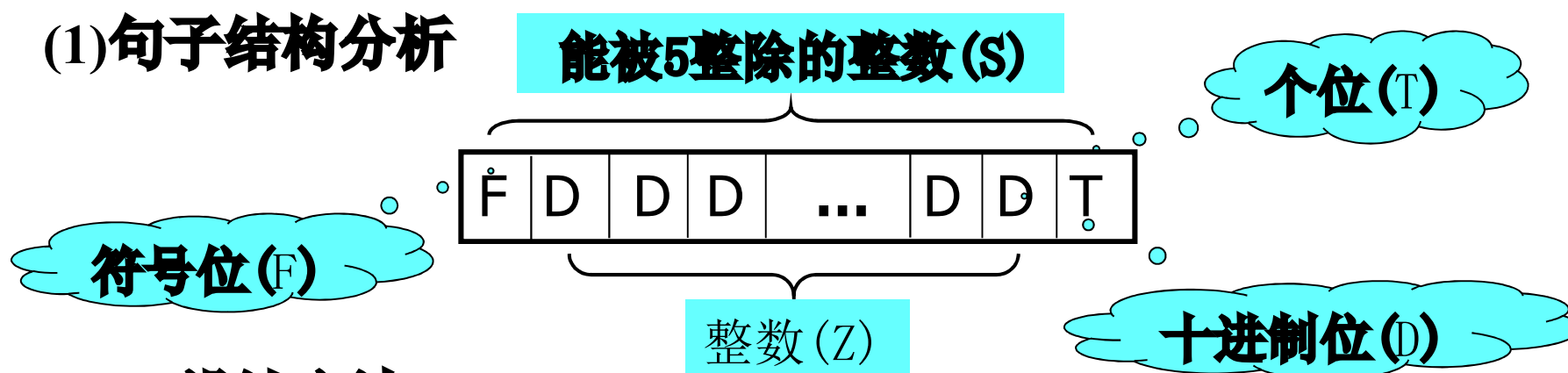
例3.1 $G_1 = (\{S\} , \{a, b\} , P, S) ,$
其中, $P = \{ S \rightarrow aSb, S \rightarrow ab \}$

2.3 文法设计举例

试设计一文法G, 使得L(G)为能被5整除的整数集。

解法 I:

(1) 句子结构分析



(2) 设计文法

G[S]: $S \rightarrow F Z T$
 $F \rightarrow + \mid - \mid \varepsilon$
 $Z \rightarrow Z D \mid \varepsilon$
 $D \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
 $T \rightarrow 0 \mid 5$

2.3 文法设计举例

解法Ⅱ:

(1) 句子结构分析

句型1

F	T
---	---

句型3

F	D	D	D	...	D	D	T
---	---	---	---	-----	---	---	---

句型2

T

句型4

D	D	D	...	D	D	T	
---	---	---	-----	---	---	---	--

(2) 设计文法

$G[S]: S \rightarrow FT \mid T \mid FZT \mid ZT$

$F \rightarrow + \mid -$

$Z \rightarrow ZD \mid D$

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

$T \rightarrow 0 \mid 5$

思考题：试设计一文法G，使得L(G)为能被3整除的整数集。

定义2.7 文法等价

设 G_1 和 G_2 是两个文法, 如果 $L(G_1)=L(G_2)$, 则称文法 G_1 和 G_2 是等价的。

例如, 下列文法 G_2 和 G_3 是等价的。因为它们产生的语言都是以字母 a 开头、字母 a 和 b 构成的符号串的集合。即 $L(G_2)=L(G_3)=\{a\}\{a, b\}^*$ 。

$G_2 = (\{ S, C \} , \{ a, b \} , P, S) ,$
其中, $P = \{ S \rightarrow aC, C \rightarrow aC, C \rightarrow bC, C \rightarrow \epsilon \} 。$

$G_3 = (\{ S \} , \{ a, b \} , P, S) ,$
其中, $P = \{ S \rightarrow Sa, S \rightarrow Sb, S \rightarrow a \} 。$

2.4 文法类型

对规则构成加以限制，可以将文法的分类为四种类型：
0型文法、1型文法、2型文法和3型文法。

(1) **0型文法** 设文法 $G = (V_N, V_T, P, S)$ ，如果任意 $\alpha \rightarrow \beta \in P$ ， α 中至少含有一个非终结符，则称文法 G 属于**0型文法**。0型文法，也称为**短语文法**。

例3.4 文法 G_4 定义如下。显然 G_4 是0型文法。 $L(G_4) = \{ \}$ 。

$$G_4 = (V_N, V_T, P, S),$$

$$\text{其中, } V_N = \{A, B, S\},$$

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

$$P = \{S \rightarrow 0AB, 1B \rightarrow 0, B \rightarrow SA \mid 01,$$

$$A1 \rightarrow SB1, A0 \rightarrow SOB\}$$

2.4 文法类型

(2) 1型文法 设文法 $G = (V_N, V_T, P, S)$ ，如果任意 $\alpha \rightarrow \beta \in P$ ， α 中至少含有一个非终结符，且除空规则之外， α 的长度不大于 β 的长度，即 $|\alpha| \leq |\beta|$ ，则称文法 G 属于**1型文法**。1型文法，也称为**上下文有关文法**。

例3.5 文法 G_5 定义如下，显然 G_5 是1型文法。

$$L(G_5) = \{a^n b^n c^n \mid n \geq 1\}。$$

$$G_5 = (V_N, V_T, P, S)，$$

$$\text{其中，} V_N = \{S, B, C\}，$$

$$V_T = \{a, b, c\}，$$

$$P = \{S \rightarrow aSBC \mid aBC, CB \rightarrow BC，$$

$$aB \rightarrow ab, bB \rightarrow bb，$$

$$bC \rightarrow bc, cC \rightarrow cc\}$$

2.4 文法类型

(3) **2型文法** 设文法 $G = (V_N, V_T, P, S)$ ，如果任意 $\alpha \rightarrow \beta \in P$ ， $\alpha \in V_N$ ，则称文法 G 属于**2型文法**。2型文法，也称为**上下文无关文法**。

例3.6 文法 G_6 定义如下，显然 G_6 是2型文法。

$$L(G_6) = \{w\$w^R \mid n \geq 0, w^R \text{ 为 } w \text{ 之逆}, w \in \{0, 1\}^*\}。$$

$$G_6 = (V_N, V_T, P, S)，$$

$$\text{其中， } V_N = \{S\}，$$

$$V_T = \{\$, 0, 1\}，$$

$$P = \{S \rightarrow 0S0 \mid 1S1 \mid \$ \}$$

2.4 文法类型

(4) 3型文法 设文法 $G = (V_N, V_T, P, S)$ ，如果任意 $\alpha \rightarrow \beta \in P$ ， $\alpha \in V_N$ ，且 β 只能是 aB 或 a （除空规则之外），则称文法 G 属于**右线性3型文法**。

设文法 $G = (V_N, V_T, P, S)$ ，如果任意 $\alpha \rightarrow \beta \in P$ ， $\alpha \in V_N$ ，且 β 只能是 Ba 或 a （除空规则之外），则称文法 G 属于**左线性3型文法**。显然 G 是3型文法。

左线性3型文法和右线性3型文法统称为3型文法，也称为正规文法。

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

$$\text{其中, } V_N = \{S, A, B\},$$

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

$$P = \{S \rightarrow A0 \mid B1, A \rightarrow 0 \mid 1, B \rightarrow 0 \mid 1\}$$

2.4 文法类型

文法分类是对规则形式逐步加以限制而得。换言之，从0型文法到1型文法、2型文法和3型文法，其规则形式逐步简单。自然，其表达力也随之逐步减弱。

如果 L_0 、 L_1 、 L_2 和 L_3 分别是0型文法、1型文法、2型文法和3型文法能产生的语言之集，则有如下关系：

$$L_0 \supseteq L_1 \supseteq L_2 \supseteq L_3.$$

在《编译原理》里，仅仅涉及到2型文法和3型文法的使用。后面，也仅仅继续讨论2型文法（或3型文法）的有关内容。

2.5 上下无关文法及其语法树

上下无关文法一个显著特征是规则左部一定有且仅有一个非终结符。利用这个特征，可以不列出 V_N 和 V_T ，给出一个上下无关文法的简洁描述方法：①文法名 G 改写成 $G[S]$ ，其中， S 表示开始符；②规则集 P ，仅书写其具体规则。

如，例2.7定义的文法 G_6 ，简化描述如下。

$G_6[S]$:

$S \rightarrow A0 \mid BS0$

$A \rightarrow 0 \mid 1$

$B \rightarrow 0 \mid 1$



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

$V_T = \{0, 1\}$

定义2.8 最左推导、最右推导

如果在推导的每一步总是选择当前句型的最左（最右）边非终结符进行推导，则称这种推导过程为**最左（最右）推导**。最右推导，也叫**规范推导**。由规范推导所得的句型，叫做**规范句型**。规范推导的逆过程，叫做**规范归约**。

例2.8 已知文法如下，试给出句子aabbbaa的推导过程。

$$\begin{aligned} G[S]: S &\rightarrow aAS \mid a \\ A &\rightarrow SbA \mid SS \mid ba \end{aligned}$$

最左推导: $S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aabAS \Rightarrow aabbaS \Rightarrow aabbbaa$

最右推导: $S \Rightarrow aAS \Rightarrow aAa \Rightarrow aSbAa \Rightarrow aSbbaa \Rightarrow aabbbaa$

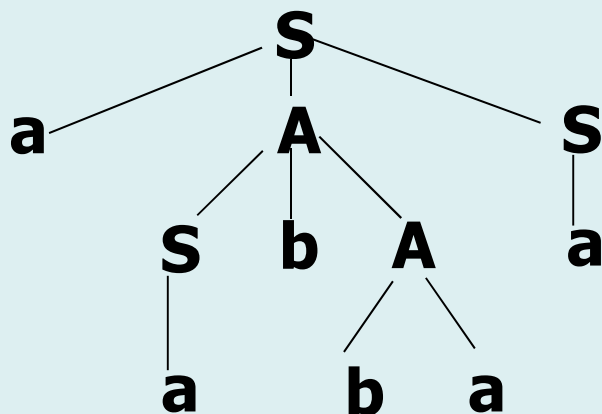
一般推导: $S \Rightarrow aAS \Rightarrow aSbAS \Rightarrow aSbAa \Rightarrow aabAa \Rightarrow aabbbaa$

定义2.9 语法树

假设文法 $G = (V_N, V_T, P, S)$ ，则文法 G 的语法树是一个满足下列条件的多叉树：

- (1) 以文法开始符 S 做为树根；
- (2) 以终结符号或非终结符号做为树的其他结点，且子树根和其孩子结点分别是某规则的左部和右部。

例如，例3.8定义的文法 $G[S]$ ，句子 $aabbbaa$ 对应的语法树如下。



例3.8 定义文法 $G[S]$:

$G[S]: S \rightarrow aAS \mid a$

$A \rightarrow SbA \mid SS \mid ba$

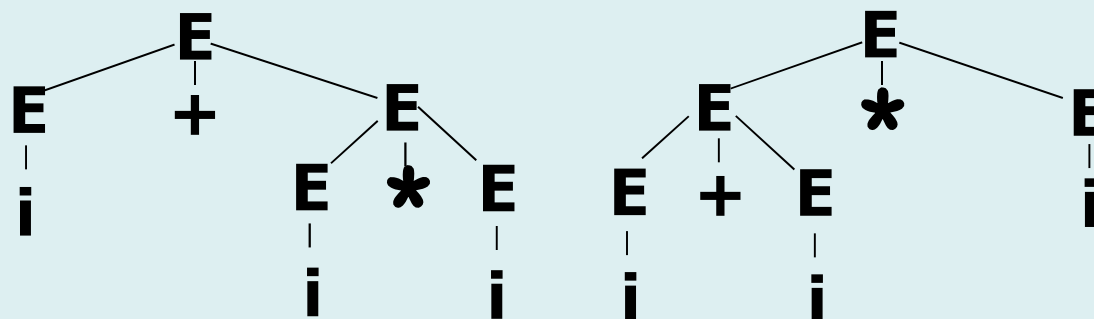
推论：①非叶子结点一定是非终结符
②全部叶子结点组成的符号串是文法的句子

定义2.10 语法二义性

如果一个文法G，某个句子存在对应的至少两棵不同的语法树，则称文法G是二义性的。

例3.9 已知文法G[E]: $E \rightarrow E + E \mid E * E \mid i$ ，证明G是二义性的。

证明: \because 句子 $i + i * i$ 存在下列两棵不同的语法树



\therefore 文法G[E]是二义性的文法

推论 ① 如果文法是无二义性的，一个句子的语法树反映了该句子的全部推导过程；② 如果文法是无二义性的，一个句子的最左（最右）推导是唯一的。

语言的先天二义性

文法的二义性，并不等同于语言的二义性，尽管两者之间可能存在非必然的联系。

因为二义性文法 G ，可能存在与之等价的无二义性的文法 G' ，即 $L(G)=L(G')$ 。

如果一个语言不存在无二义性的文法，则称该语言是**先天二义性**的。

例如，语言 $L=\{a^i b^j c^k \mid (i=j \text{ 或 } i=k), (i, j, k \geq 1)\}$ 不存在无二义性的文法，是**先天二义性**的语言。

已经证明：文法的二义性判定问题是递归不可解的。即不存在这个判定问题的算法。

2.6 句型分析

假设文法 $G[S]$ 是语言 L 之文法，即 $L(G)=L$ ，则“符号串 α 是否符合语言 L 的语法问题”被等价地转化成“推导或归约问题”，即：

$$S \xRightarrow{*} \alpha \wedge \alpha \in V_T^*$$

这样，自然地形成了**推导法和归约法**两大类分析方法。推导法和归约法，也分别称为自上而下的分析方法和自下而上的分析方法。

2.6.1 自上而下的分析方法

自上而下分析法：从文法开始符号出发，反复使用规则，寻找匹配符号串（推导）的句型，直到推导出句子或规则用遍。**进行每步推导时，存在两个选择问题：**

(1) 选择句型中哪一个非终结符进行推导

(2) 选择非终结符的哪一个规则进行推导

问题(1)可以采用最左推导解决。问题(2)通常需要穷举每一个规则的可能推导。

成功：在推到过程中一旦出现个符号串 α ，便结束穷举过程，断定符号串 α 是句子。

失败：当穷举全部可能的推导，而不存在一个符号串 α 之推导过程的时候，才可以断定符号串 α 不是句子。

2.6.1 自上而下的分析方法

(1) 输入串cabd的推导过程

(2) 输入串cabc的推导过程

序号	推导过程	输入串	说明
1	S	cabd	从S开始推导
2	$S \Rightarrow cAd$	c abd	选1号规则, c匹配成功
3	$S \Rightarrow cad$	c abd	选2号规则, a匹配成功d不成功, 回溯
4	$S \Rightarrow cAd$	c abd	选3号规则推导
5	$S \Rightarrow cabd$	cabd	abd匹配成功, 所以cabd是一个句子

G[S]:

1. $S \rightarrow cAd$
2. $A \rightarrow a$
3. $A \rightarrow ab$

这样带回溯的推导法效率很低。后面重点研究避免回溯的推导法!

2.6.2 自下而上的分析方法

自下而上分析法：从输入符号串 α 开始，逐步进行“归约”，直至归约出文法的开始符号 S ，则输入串 α 是文法 G 定义的语言的句子。否则不是。

这种分析方法在进行每步归约时，存在两个如何选择句型 α 的子串 β 进行归约的问题 ($\alpha = \beta \delta$)。

如果文法规则没有相同的右部，则在语法分析的过程中，一旦出现子串 β 与某条规则的右部相同，就可以使用这条规则进行归约，简单优先分析法就是采用此方法进行归约。

但这种限制，实际上也限制了文法的表达能力，所以通常是通过在句型中寻找所谓的“句柄”的途径解决的。

定义 2.11 短语、直接短语、句柄

设 $G[S]$ 是一文法， $\alpha \beta \delta$ 是文法 G 的句型，如果有 $S \xRightarrow{*} \alpha A \delta$ 且 $A \xRightarrow{+} \beta$ ，则称 β 是句型 $\alpha \beta \delta$ 的、相对于非终结符 A 的**短语**。

特别地，当 $A \Rightarrow \beta$ 实际是 $A \Rightarrow \beta$ 即一步推导时，则又称 β 是句型 $\alpha \beta \delta$ 的、相对于非终结符 A 的**直接短语**（或**简单短语**）。

句型的最左直接短语，称为该句型的**句柄**。

短语的理解：

“ $\alpha \beta \delta$ 是文法 G 的句型”，即 $S \xRightarrow{*} \alpha \beta \delta$

“ $S \xRightarrow{*} \alpha A \delta$ 且 $A \xRightarrow{+} \beta$ ”，即 $S \Rightarrow \underbrace{\dots}_{\geq 0} \Rightarrow \alpha A \delta \Rightarrow \underbrace{\dots}_{\geq 1} \Rightarrow \alpha \beta \delta$

这表明，如果 β 是句型 $\alpha \beta \delta$ 的、相对于 A 的短语，则至少存在一个推导，使得 $\alpha A \delta \Rightarrow \alpha \beta \delta$ ，或者 $\alpha \beta \delta \xRightarrow{+} \alpha A \delta$ 。

特别地，如果 β 是直接短语，则 $\alpha A \delta \Rightarrow \alpha \beta \delta$ ，或者 $\alpha \beta \delta \Leftarrow \alpha A \delta$ 。

短语、直接短语、句柄举例

例2.12 设文法 $G[E]$: $E \rightarrow E+T \mid T$, $T \rightarrow T * F \mid F$, $F \rightarrow (E) \mid i$, 分析句型 $i+i*i$ 的短语、直接短语和句柄。

从句型 $i+i*i$ 的如下的推导过程之一如下:

$$\begin{aligned} E &\Rightarrow E+T \Rightarrow E+T * F \Rightarrow E+T * i_3 \Rightarrow E+F * i_3 \Rightarrow E+i_2 * i_3 \\ &\Rightarrow T + i_2 * i_3 \Rightarrow F + i_2 * i_3 \Rightarrow i_1 + i_2 * i_3 \end{aligned}$$

可以看出部分短语如下:

- (1) i_3 是句型 $E+T * i_3$ 的、相对于非终结符 F 的直接短语;
- (2) i_2 是句型 $E+ i_2 * i_3$ 的、相对于非终结符 T 的短语, F 的直接短语;
- (3) i_1 是句型 $i+i*i$ 的、相对于非终结符 F 的短语、直接短语和句柄。

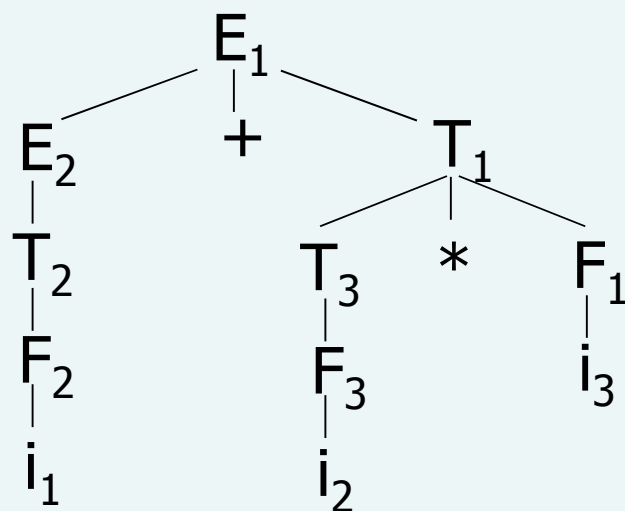
为了讨论方便起见, 句型 $i+i*i$ 改写成 $i+i_2*i_3$ 。其中, 下标1、2、3仅仅是用来标记不同 i 出现在句型中的位置。

语法树与短语、直接短语、句柄举例

从句型的语法树上，可以直观地找出句型的短语。在语法树中相对于每个非空子树的根是一个非终结符，其非空子树叶子结点组成的符号串即为短语。如果子树根与子树叶子结点之间均为父子关系，则该短语还是直接短语。

$$\begin{aligned} G[E] : E &\rightarrow E+T \mid T \\ T &\rightarrow T*F \mid F \\ F &\rightarrow (E) \mid i \end{aligned}$$

例2.12 设文法 $G[E]$ 定义如右，分析句型 $i+i*i$ 的短语、直接短语和句柄。



句型 $i+i*i$ 语法树

短语: $i_1+i_2*i_3$ (相对于 E_1)

i_1 (相对于 E_2)

i_2*i_3 (相对于 T_1)

i_1 (相对于 T_2)

i_2 (相对于 T_3)

直接短语: i_3 (相对于 F_1)

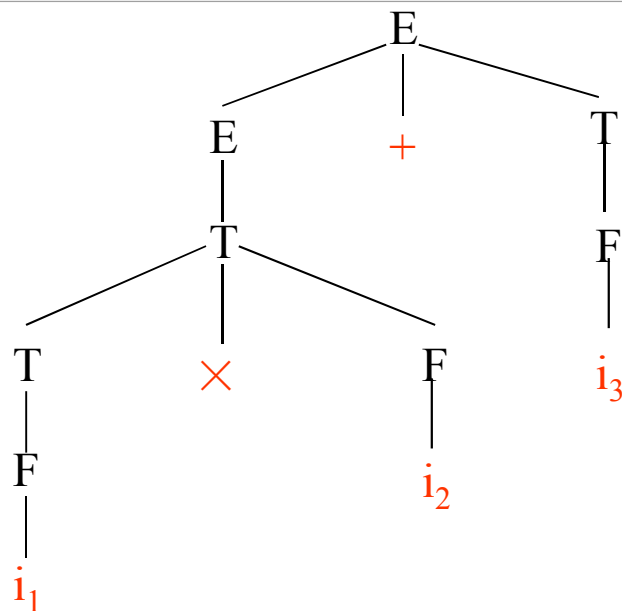
i_2 (相对于 F_3)

句柄: i_1 (相对于 F_2)

根据句柄进行归约举例

仍然以例2.12定义的文法 $G[E]$ 为例，给出通过寻找句柄，对符号串 $i+i*i$ 进行归约法分析过程，说明归约法的基本思想。

$G[E] : E \rightarrow E+T \mid T, T \rightarrow T * F \mid F, F \rightarrow (E) \mid i$



分析过程借用了语法树确定句型的句柄，显然是不符合逻辑的。因此，如何依据文法寻找句柄是归约法的关键问题。

2.7 文法在实用中的一些说明

在实际应用中，对于文法规则提出了一些限制条件，但这些并没有限制文法的语言描述能力。限制下列 3 种规则的使用：

- (1) **有害规则** 形如 $U \rightarrow U$ 的规则，称为**有害规则**。
- (2) **不可达规则** 不在任何规则右部出现的非终结符对应的规则，称为**不可达规则**。
- (3) **不可终止规则** 如果从某非终结符开始，不可能推导出任意终结串来，则该非终结符对应的规则称为**不可终止规则**。

不含有**多余规则**的文法，称为**压缩过的文法**。在后面讨论的文法时，都假设是压缩过的的文法。

例2.13 文法G定义如下,

G[S]:	$S \rightarrow Be$
	$B \rightarrow Ce \mid Af,$
	$A \rightarrow Ae \mid e$
	$C \rightarrow Cf$
	$D \rightarrow f$

显然, $C \rightarrow Cf$ 是不可终止规则, $D \rightarrow f$ 是不可达规则。除去这些规则及其相关规则之后, 得到的等价文法G' 如下。

G' [S]:	$S \rightarrow Be$
	$B \rightarrow Af$
	$A \rightarrow Ae \mid e$

ϵ 规则问题

在文法设计中，使用 ϵ 规则有时会带来方便，但会导致文法讨论和证明的复杂。

一个上下文无关文法 G 是否必须使用 ϵ 规则，完全取决于文法 G 产生的语言 $L(G[S])$ 中是否含有 ϵ 语句。

可以证明，如果 $\epsilon \notin L(G[S])$ ，则存在一个等价的文法 $G' [S']$ ，且 G' 不含 ϵ 规则。

如果 $\epsilon \in L(G[S])$ ，则存在一个等价的文法 $G' [S']$ ，且 G' 仅含 $S' \rightarrow \epsilon$ 的一个空规则。

提示：使用“代入法”，即可得到等价的文法 $G' (S')$

“代入法” 举例

已知文法 $G[S]$ 如下。容易看出 $L(G[S]) = \{0, 1, +0, +1, -0, -1\}$,
且 $\epsilon \notin L(G[S])$ 。

$G[S]$:

$$\begin{aligned} S &\rightarrow FD \\ F &\rightarrow + \mid - \mid \epsilon \\ D &\rightarrow 0 \mid 1 \end{aligned}$$

对于右部出现 F 的规则，使用 F 规则的右部替代其右部出现的 F ，之后删除 F 的规则，得到等价于 $G[S]$ 的文法 $G'[S]$ 如下。

$G[S]$:

$$\begin{aligned} S &\rightarrow +D \mid -D \mid D \\ D &\rightarrow 0 \mid 1 \end{aligned}$$

本章小结

本章重点介绍了形式化语言的基本理论、语言形式化描述方法、语言的文法分类与特性、文法与句型分析。

提出的基本概念有语言、文法、规则、推导、归约、句型、短语、句柄、语法树和语法二义性等。其中，推导是最核心的、关键的概念。重点掌握的内容是

- ① 设计一个已知语言的文法；
- ② 确定已知文法定义的语言；
- ③ 求句型的短语、直接短语和句柄。
- ④ 文法二义性判定。