

第3章 词法分析

Part I 正规文法与正规式

形式文法内容回顾

- 符号 \rightarrow 符号串 \rightarrow 句子 \rightarrow 语言
- 语言表示的两个途经：
 - 生成方式
 - 识别方式

一些基本概念

- 文法—推导—句型—短语—语法树—句柄

文法与自动机的关系

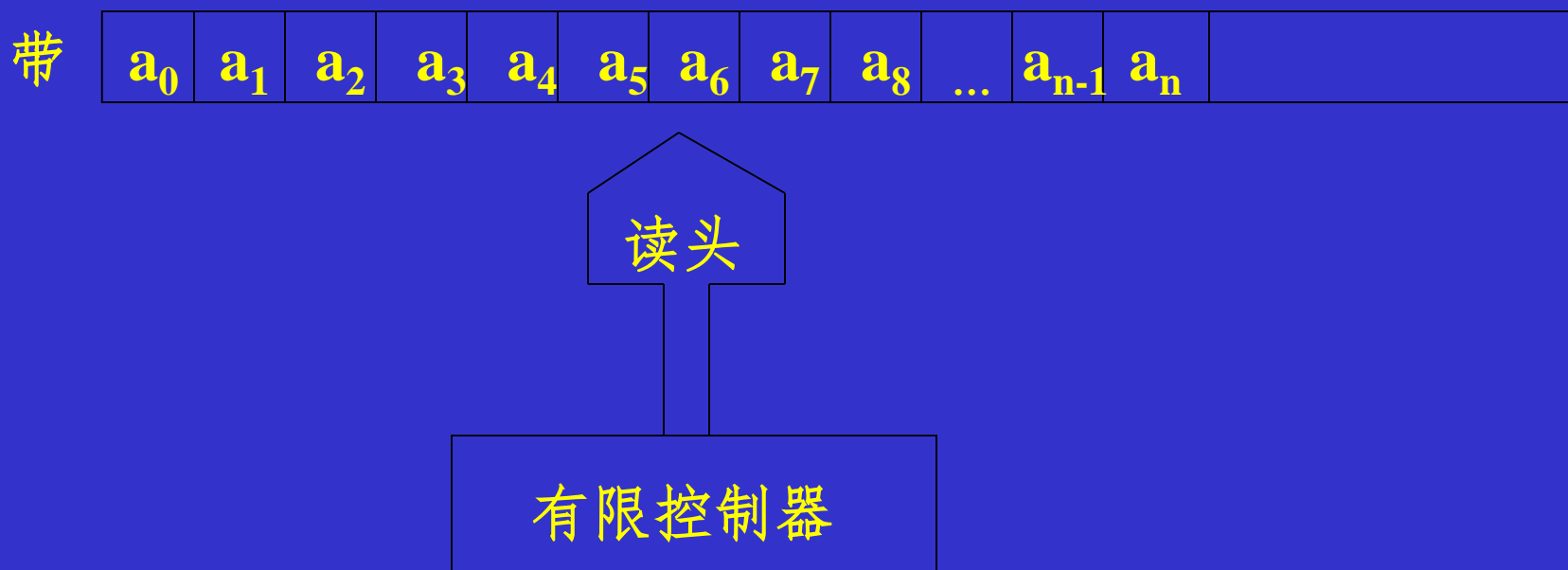
0型文法（短语结构文法）：其能力相当于图灵机，可以表征任何递归可枚举集，而且任何0型语言都是递归可枚举的

1型文法（上下文有关文法CSG）：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 α_1 和 α_2 的上下文中时，才允许 β 取代A。其识别系统是线性有界自动机

2型文法（上下文无关文法CFG）：产生式的形式为 $A \rightarrow \beta$ ， β 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RG）：产生的语言是有穷自动机（FA）所接受的集合

任何能用图灵机描述的计算都能机械地实现，
任何能在现代计算机上实现的计算都能用图灵机描述



Part I 自动机

- 正规文法与正规式(Regular Expression)
- 有限自动机(Finite Automata)
 - DFA
 - NFA
 - NFA- \rightarrow DFA
 - ϵ -FA
 - FA化简
 - FA与RG, RE的等价性
- 下推自动机

正规文法与正规式

- 单词符号结构的描述方法：
 - 正规文法（3型文法）
 - 正规式（正则表达式）

正规式

正规式也称正则表达式

正规表达式(regular expression)是说明单词的模式(pattern)的一种重要的表示法(记号)，是定义正规集的数学工具

在编译中，用以描述单词符号

定义（正规式和它所表示的正规集）：

设字母表为 Σ ，辅助字母表 $\Sigma' = \{\Phi, \varepsilon, |, \bullet, *, (,)\}$ 。

1. ε 和 Φ 都是 Σ 上的正规式，它们所表示的正规集分别为 $\{\varepsilon\}$ 和 $\{\}$ ；

2. 任何 $a \in \Sigma$, a 是 Σ 上的一个正规式, 它所表示的正规集为 $\{a\}$;
3. 假定 e_1 和 e_2 都是 Σ 上的正规式, 它们所表示的正规集分别为 $L(e_1)$ 和 $L(e_2)$, 那么, (e_1) , $e_1 \mid e_2$, $e_1 \bullet e_2$, e_1^* 也都是正规式, 它们所表示的正规集分别为 $L(e_1)$, $L(e_1) \cup L(e_2)$, $L(e_1)L(e_2)$ 和 $(L(e_1))^*$ 。
4. 仅由有限次使用上述三步骤而定义的表达式才是 Σ 上的正规式, 仅由这些正规式所表示的集合才是 Σ 上的正规集。

正规式中的符号

其中的“ $|$ ”读为“或”（也有使用“ $+$ ”代替“ $|$ ”的）；

“ \bullet ”读为“连接”；

“ $*$ ”读为“闭包”（即，任意有限次的自重复连接）。

在不致混淆时，括号可省去，但规定算符的优先顺序为“ $*$ ”、“ \bullet ”、“ $|$ ”。连接符“ \bullet ”一般可省略不写。“ $*$ ”、“ \bullet ”和“ $|$ ”都是左结合的。

例子

令 $\Sigma = \{a, b\}$, Σ 上的正规式和相应的正规集的例子有:

正规式

正规集

a

$\{a\}$

$a \mid b$

$\{a, b\}$

ab

$\{ab\}$

$(a \mid b)(a \mid b)$

$\{aa, ab, ba, bb\}$

a^*

$\{\varepsilon, a, a, \dots \dots \text{任意个 } a \text{ 的串}\}$

正规式

$(a \mid b)^*$

$(a \mid b)^*(aa \mid bb)(a \mid b)^*$

正规集

$\{\varepsilon, a, b, aa, ab, \dots\}$ 所有由a和b组成的串}

$\{\Sigma^*$ 上所有含有两个相继的a或两个相继的b组成的串}

例

令 $\Sigma=\{1, d\}$, 则 Σ 上的正规式 $r=1(1 \mid d)^*$ 定义的正规集为:
 $\{1, 11, 1d, 1dd, \dots\}$, 其中1代表字母, d代表数字, 正规式即是字母(字母|数字)*, 它表示的正规集中的每个元素的模式是“字母打头的字母数字串”, 就是Pascal和多数程序设计语言允许的标识符的词法规则

例

$\Sigma=\{d, \bullet, e, +, -\}$,

则 Σ 上的正规式 $dd^*(\bullet dd^* \mid \varepsilon)(e(+ \mid - \mid \varepsilon)dd^* \mid \varepsilon)$ 表示的是无符号数的集合。其中d为0~9的数字

程序设计语言的单词都能用正规式来定义

正规式等价

若两个正规式 e_1 和 e_2 所表示的正规集相同,则说 e_1 和 e_2 等价,写作 $e_1=e_2$ 。

例如: $e_1 = (a \mid b)$, $e_2 = b \mid a$

又如: $e_1 = b(ab)^*$, $e_2 = (ba)^*b$
 $e_1 = (a \mid b)^*$, $e_2 = (a^* \mid b^*)^*$

正规式等价变换规则

设 r, s, t 为正规式，正规式服从的代数规律有：

1. $r \mid s = s \mid r$ “或”服从交换律
2. $r \mid (s \mid t) = (r \mid s) \mid t$ “或”的可结合律
3. $(rs)t = r(st)$ “连接”的可结合律
4. $r(s \mid t) = rs \mid rt$
 $(s \mid t)r = sr \mid tr$ 分配律

$$5. \varepsilon r = r, r\varepsilon = r$$

ε 是“连接”的恒等元素
零一律

$$6. r \mid r = r$$

$$r^* = \varepsilon \mid r \mid rr \mid \dots$$

“或”的抽取律

正规文法和正规式

- 例：标识符的文法描述
 - $G = (\{ \text{DIGIT}, \text{LETTER} \}, \{ \text{idn} \}, P, \text{idn})$
 - $\text{idn} \rightarrow \text{LETTER}$
 - $\text{idn} \rightarrow \text{idn DIGIT}$
 - $\text{idn} \rightarrow \text{idn LETTER}$
- 正则式： $\text{LETTER}(\text{LETTER}|\text{DIGIT})^*$

正规式到正规文法

对 Σ 上的正规式 r , 存在一个 $RG=(V_N, V_T, P, S)$:
 $L(G)=L(r)$

初始, $V_T=\Sigma, S \in V_N$, 生成正规产生式: $S \rightarrow r$

(R.1) 对形如 $A \rightarrow r_1 r_2$ 的正规产生式: $A \rightarrow r_1 B$
 $B \rightarrow r_2 \quad B \in V_N$

(R 2) 对形如 $A \rightarrow r^* r_1$ 的正规产生式: $A \rightarrow r B$
 $A \rightarrow r_1$
 $B \rightarrow r B$
 $B \rightarrow r_1 \quad B \in V_N$

(R 3) 对形如 $A \rightarrow r_1 \mid r_2$ 的正规产生式: $A \rightarrow r_1$
 $A \rightarrow r_2$

不断应用**R**做变换, 直到每个产生式右端至多有一个 V_N

例 $r = a(a \mid d)^*$

$S \rightarrow a(a \mid d)^*$

$S \rightarrow aA \quad A \rightarrow (a \mid d)^*$

$A \rightarrow (a \mid d)B$

$A \rightarrow \varepsilon$

$B \rightarrow (a \mid d)B$

$B \rightarrow \varepsilon$

$G[S]:$

$S \rightarrow aA$

$A \rightarrow \varepsilon$

$A \rightarrow aB$

$A \rightarrow dB$

$B \rightarrow aB$

$B \rightarrow dB$

$B \rightarrow \varepsilon$

$V_T = \{a, d\}$

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

例：标识符定义的转换

- 引入 **id**

$$\mathbf{id} \rightarrow \mathbf{let} (\mathbf{let} \mid \mathbf{dig})^*$$

- 引入 **rid** 消除连接

$$\begin{aligned}\mathbf{rid} &\rightarrow (\mathbf{let} \mid \mathbf{dig})^* \\ &\rightarrow \varepsilon \mid (\mathbf{let} \mid \mathbf{dig})\mathbf{B}\end{aligned}$$

$$\mathbf{B} \rightarrow (\mathbf{let} \mid \mathbf{dig})\mathbf{B}$$

$$\mathbf{B} \rightarrow \varepsilon$$

$$\mathbf{id} \rightarrow \mathbf{let} \mathbf{rid}$$

$$\mathbf{rid} \rightarrow \varepsilon \mid \mathbf{let} \mathbf{B} \mid \mathbf{dig} \mathbf{B}$$

$$\mathbf{B} \rightarrow \mathbf{let} \mathbf{B} \mid \mathbf{dig} \mathbf{B}$$

$$\mathbf{B} \rightarrow \varepsilon$$

正规文法到正规式

对 $G=(V_N, V_T, P, S)$, 存在一个 $\Sigma = V_T$ 上的正规式 $r : L(r)=L(G)$

$$A \rightarrow xB, \quad B \rightarrow y \quad \approx \quad A = xy$$

$$A \rightarrow xA \mid y \quad \approx \quad A = x^*y$$

$$A \rightarrow x \mid y \quad \approx \quad A = x \mid y$$

例子

$$G[s]: S \rightarrow aA \mid a$$

$$A \rightarrow aA \mid a \mid dA \mid d$$

$$A \rightarrow (a \mid d)A \mid (a \mid d)$$

$$A \rightarrow (a \mid d)^*(a \mid d)$$

$$S = a(a \mid d)^*(a \mid d) \mid a$$

$$= a((a \mid d)^*(a \mid d) \mid \varepsilon)$$

$$= a((a \mid d)^+ \mid \varepsilon)$$

$$R = a(a \mid d)^*$$

请指出下列哪些字符串包含在正规式
 $ab^*c^*(a|b)c$ 所对应的正规集合中：

acac acbbc abbcac abc acc

作业

- 通读3.1

- 8