

编译程序首先是在单词级别上来分析和翻译源程序的。词法分析的任务是:从左至右逐个字符地对源程序进行扫描,产生一个个单词符号,把作为字符串的源程序改造成为单词符号串的中间程序。

因此,词法分析是编译的基础。执 行词法分析的程序称为词法分析器。

3. 1对词法分析器的要求

词法分析器功能:输入源程序,输出单词符号。

程序语言的单词符号一般分为五种:关键字,标识符,常数,运算符,界符

词类和属性

- 1. 关键字(保留字或基本字): while, if
- 2. 标识符: 用来表示各种名字
- 3. 字面常数: 256, 3.14, true, 'abc'
- 4. 运算符:如,十、一、*、/ 等等
- 5. 分界符: 如逗号, 分号, 冒号等

词法分析器输出的单词符号常常表示为二元式: (单词种别,单词符号的属性值)

词类编码原则:

界符和运算符:一符一码。

关键字可分成一类,也可以一个关键字分成一类。一字一码。

常数可统归一类,也可按类型(整型、实型、布尔型等),每个类型的常数划分成一类。一类型一码。

所有的标识符分为一类。一类一码。

对于关键字、界符、运算符来说,它们的词类编码就可以表示其完整的信息, 故对于这类单词,其单词自身的属性 值通常为空

而对于标识符, 词类编码所反映的信息不够充分, 标识符的具体特性还要通过单词自身的属性进行互相区分。 标识符的单词自身的属性常用其 在符号表中的入口指针来表示

对于常数, 其单词自身的属性常用 其在常数表中的入口指针来表示

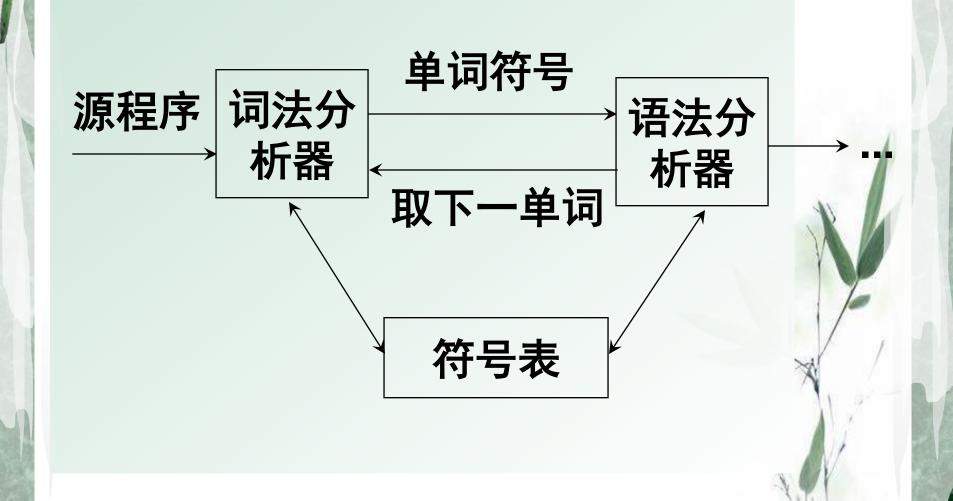
作为例子考虑下述C++代码段:

while (i>=j) i--; 经词法分析器处理后,它将转换的单词 符号序列:

< while, ->

- < (, ->
- 〈 id,指向i的符号表项的指针 〉
- <>= . ->
- 〈 id,指向j的符号表项的指针〉
- <) . ->
- 〈 id,指向i的符号表项的指针 〉
- <--->
- <:.->

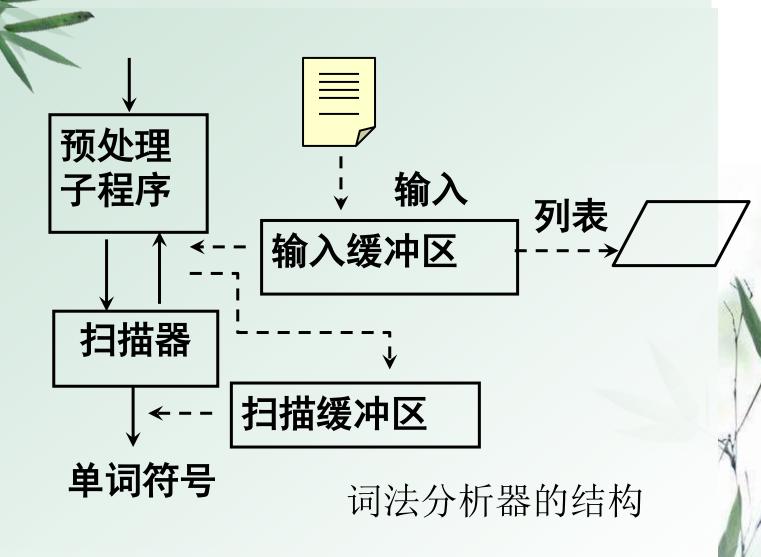
词法分析器



词法分析器作为独立子程序

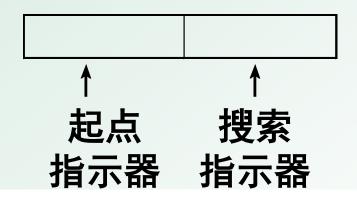
我们把词法分析器安排成一个独 立子程序,每当语法分析器需要一 个单词符号时就调用这个子程序。 每一次调用, 词法分析器就从符号 串中识别出一个单词符号,把它交 给语法分析器。这样,把词法分析 器安排成一个子程序似乎比较自然。 在后面的讨论中,我们假定词法分 析器是按这种方式进行工作的。

3.2 词法分析器的设计



~3.2.1 输入、预处理

- 输入串放在输入缓冲区中。
- 预处理子程序:剔除无用的空白、跳格、回车和换行等编辑性字符;区分标号区、捻接续行和给出句末符等
- 扫描缓冲区



3.2.2 单词符号的识别:超前搜索

1 关键字识别:

例如:

D099K=1, 10 D0 99 K = 1, 10

IF (5. EQ. M) GOTO55 IF (5. EQ. M) GOTO 55

D099K=1.10

IF(5) = 55

• 需要超前搜索才能确定哪些是关键字

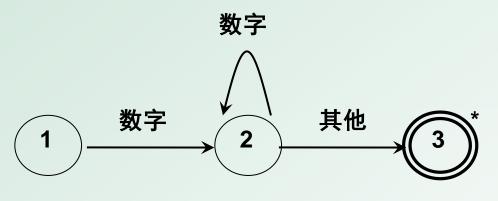
- 2 标识符识别:
- 字母开头的字母数字串,后跟界符或算符
- 3 常数识别:
- 识别出算术常数并将其转变为二进制内码表示。
 有些也要超前搜索。
 - 5. EQ. M
 - 5. E08
- 4 算符和界符的识别
- 把多个字符符合而成的算符和界符拼合成一个 单一单词符号。

3.2.3 状态转换图

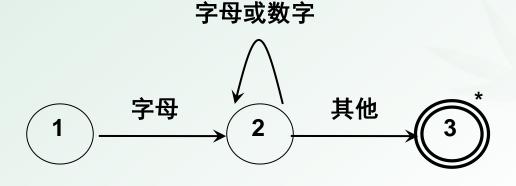
- 1 概念
- 状态转换图是一张有限方向图。
- >结点代表状态,用圆圈表示。
- 》状态之间用箭弧连结,箭弧上的标记(字符)代表射出结状态下可能出现的输入字符或字符类。
- 一张转换图只包含有限个状态,其中有一个为初态,至少要有一个终态。

一个状态转换图可用于识别(或接受)一定的字符串。

□ 终态结上打星 号表示多读进 一个不属于标 识符的字符, 应退回给输入 串。



识别整常数的状态转换图



识别标识符的状态转换图

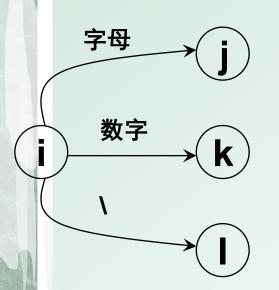
几点重要限制——不必使用超前搜索

- 所有基本字都是保留字; 用户不能用它们作自己 的标识符
- 基本字作为特殊的标识符来处理; 不用特殊的状态图来识别, 只要查保留字表。
- 如果基本字、标识符和常数(或标号)之间没有确定的运算符或界符作间隔,则必须使用一个空白符作间隔。

D099K=1, 10 要写成 D0 99 K=1, 10

3.2.4 状态转换图的实现

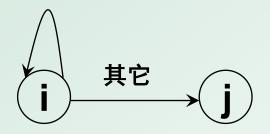
- 思想:每个状态结对应一小段程序。
- 做法:
 - 1) 对不含回路的分叉结,可用一个CASE语句或一组IF-THEN-ELSE语句实现



```
GetChar();
if (IsLetter()) {···状态j的对应程序段···;}
else if (IsDigit()) {···状态k的对应程序段···;}
else if (ch= '/') {···状态l的对应程序段···;}
else {···错误处理···;}
```

2) 对含回路的状态结,可对应一段由WHILE结构和IF语 句构成的程序.

字母或数字



3) 终态结表示识别出某种单词符号,因此,对应语句为 RETURN (C, VAL) 其中, C为单词种别, VAL为单词自身值.

3.3 正规表达式与有限自动机

• 几个概念:

- 考虑一个有穷 字母表Σ 字符集
- 其中每一个元素称为一个字符
- Σ上的字(也叫字符串) 是指由Σ中的字符所构成的 一个有穷序列
- 不包含任何字符的序列称为空字, 记为 ε
- 用 Σ*表示 Σ 上的所有字的全体, 包含空字 ε 例如: 设 Σ = {a, b}, 则 Σ*= { ε, a, b, aa, ab, ba, bb, aaa, ...}

 Σ^* 的子集U和V的连接(积)定义为 $UV = \{ \alpha\beta \mid \alpha \in U \& \beta \in V \}$ V自身的 n次积记为

Vn=**VV**---**V**

- 规定V⁰= {ε}, 令
 V*=V⁰ U V¹ U V² U V³ U ····
 称V*是V的闭包;
- 记 V⁺=VV^{*}, 称V⁺是V的正规闭包。

3. 3. 1 正规式和正规集

- 正规集可以用正规表达式(简称正规式)表示。
- 正规表达式是表示正规集一种方法。一个字集 合是正规集当且仅当它能用正规式表示。

- 正规式和正规集的递归定义: 对给定的字母表 Σ
 - 1) ϵ 和 \emptyset 都 是 Σ 上的正规式,它们所表示的正规集为 $\{\epsilon\}$ 和 \emptyset ;
 - 2) 任何 $a \in \Sigma$, $a \in \Sigma$ 上的正规式,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 集为 $L(e_1)$ 和 $L(e_2)$,则
 - i) $(e_1 | e_2)$ 为正规式,它所表示的正规集为 $L(e_1) \cup L(e_2)$,
 - ii) (e_1, e_2) 为正规式,它所表示的正规集 为 $L(e_1)L(e_2)$,
 - iii) $(e_1)^*$ 为正规式,它所表示的正规集为 $(L(e_1))^*$,
- 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,仅由这些正规式表示的字集才是Σ上的正规集。

所有词法结构一般都可以用正规式描述。

```
b(ab)*=(ba)*b
```

```
(a*b*)*=(a|b)*
```

```
L(b(ab)*)
= L(b)L((ab)*)
= L(b) (L(ab))*
= L(b) (L(a)L(b))*
= {b} {ab}*
= {b} {\epsilon}, ab, abab, ababab, \cdots}
= {b, bab, babab, bababab, \cdots}
```

```
L( (ba)*b)
= L((ba)*) L(b)
= (L(ba))*L(b)
= (L(b)L(a))* L(b)
= {ba}* {b}
= {ε, ba, baba, bababa, ···} {b}
= {b, bab, babab, bababab, ···}
```

$$\therefore L(b(ab)*) = L((ba)*b) \qquad \therefore b(ab)*=(ba)*b$$

对正规式,下列等价成立:

$$-e_1|e_2 = e_2|e_1$$
 交換律
 $-e_1|(e_2|e_3) = (e_1|e_2)|e_3$ 结合律
 $-e_1(e_2e_3) = (e_1e_2)e_3$ 结合律
 $-e_1(e_2|e_3) = e_1e_2|e_1e_3$ 分配律
 $-(e_2|e_3)e_1 = e_2e_1|e_3e_1$ 分配律

$$L(\mathbf{e}_{1} | \mathbf{e}_{2})$$
= $L(\mathbf{e}_{1}) \cup L(\mathbf{e}_{2})$
= $L(\mathbf{e}_{2}) \cup L(\mathbf{e}_{1})$
= $L(\mathbf{e}_{2} | \mathbf{e}_{1})$

 $-e\epsilon = \epsilon e = e$

3.3.2 有限自动机

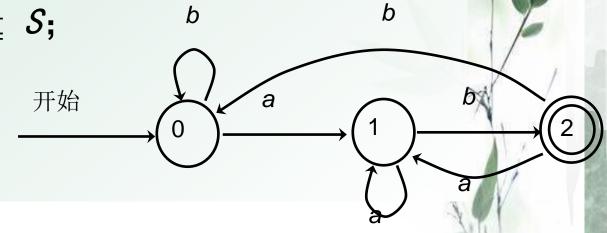
不确定的有限自动机(简称NFA)

- 一个数学模型,它包括:
- 状态集合S;
- 输入符号集合∑;
- 转换函数 $move: S \times (\Sigma \cup \{\varepsilon\}) \rightarrow P(S);$
- 状态 s_0 是开始状态;
- $F \subseteq S$ 是接受状态集合。

确定的有限自动机(简称DFA)

- 一个数学模型,包括:
- 状态集合S;
- 输入字母表Σ;
- 转换函数 $move: S \times \Sigma \rightarrow S$;
- 唯一的初态 $s \in S$;
- 终态集合 $F \subseteq S$;

识别语言 (a|b)*ab 的DFA



• 例如: DFA M=({0, 1, 2, 3}, {a, b}, f, 0, {3}), 其中: f定义如下:

$$f(0, a)=1$$
 $f(0, b)=2$
 $f(1, a)=3$ $f(1, b)=2$
 $f(2, a)=1$ $f(2, b)=3$
 $f(3, a)=3$ $f(3, b)=3$

	a	b	_ a 1
0	1	2	- a a a
1	3	2	\Rightarrow \bigcirc
2	1	3	
3	3	3	b b b
		·	2

状态转换矩阵

状态转换图

 DFA可以表示为状态转换图。假定DFA M 含有m个状态和n个输入字符,那么,这 个图含有m个状态结点,每个结点顶多 含有n条箭弧射出,且每条箭弧用Σ上 的不同的输入字符来作标记。

- 对于 Σ *中的任何字 α ,若存在一条从初态到某一终态的道路,且这条路上所有弧上的标记符连接成的字等于 α ,则称 α 为DFA M所识别(接收)
- DFA M所识别的字的全体记为L(M)。
- 可以证明: Σ 上的字集V⊆ Σ *是正规集,当且仅当存在上的DFA M,使得V=L(M)。

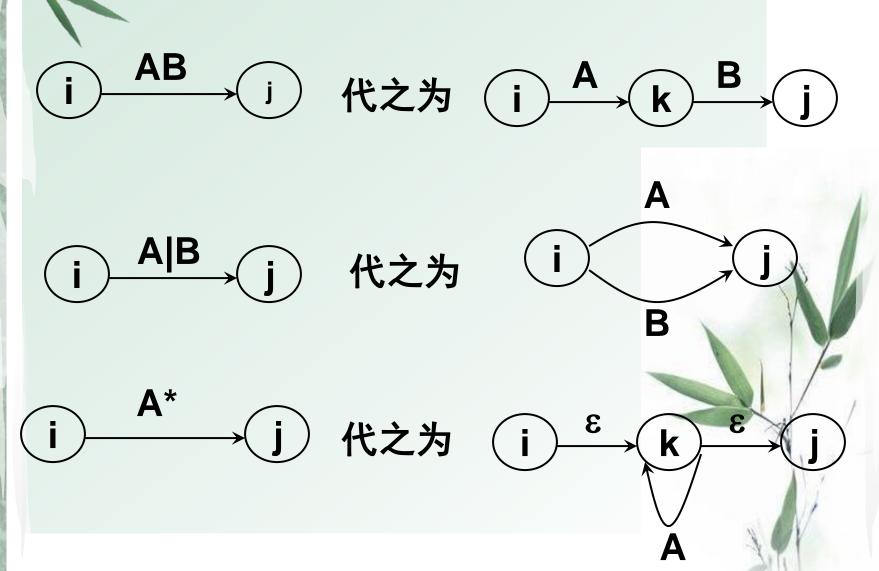
- · 从状态图中看NFA 和DFA的区别:
 - 1 弧上的标记可以是 Σ *中的一个字,而不一定是单个字符;
 - 2 同一个字可能出现在同状态射出的多条弧上。
- DFA是NFA的特例。

- 定义:对于任何两个有限自动机M和M',如果 L(M)=L(M'),则称M与M'等价。
- 自动机理论中一个重要的结论:判定两个自动机等价性的算法是存在的。
- · 对于每个NFA M存在一个DFA M',使得 L(M)=L(M')。亦即DFA与NFA描述能力相同。

证明:

- 1. 假定NFA M= $\langle S, \Sigma, \delta, S_0, F \rangle$, 我们对M的状态 转换图进行以下改造:
 - 1)引进新的初态结点X和终态结点Y, X, Y∉S, 从X到S₀中任意状态结点连一条ε箭弧, 从F中任意状态结点连一条ε箭弧到Y。
 - 2) 对M的状态转换图进一步施行替换,其中k是新引入 的状态。

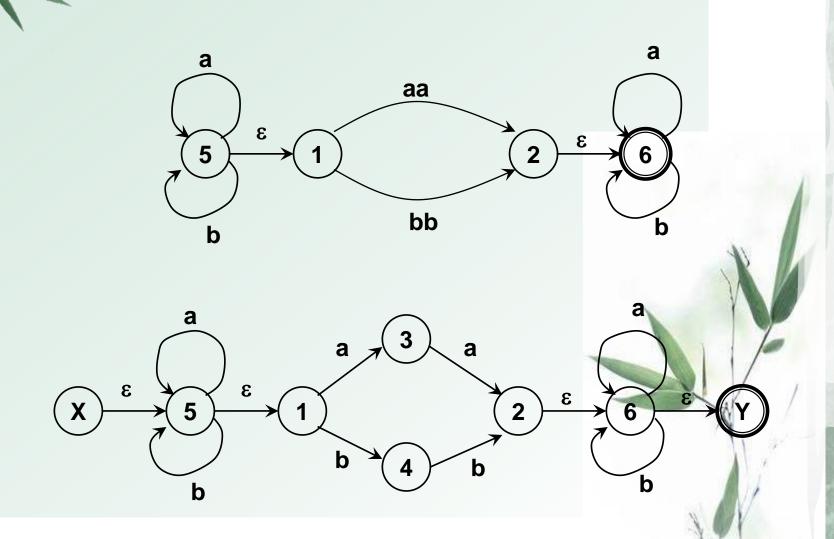
按下面的三条规则对V进行分裂:



逐步把这个图转变为每条弧只标记为Σ上的一个字符或ε,最后得到一个NFA M',显然L(M')=L(M)



识别所有含相继两个a或相继两个b的字

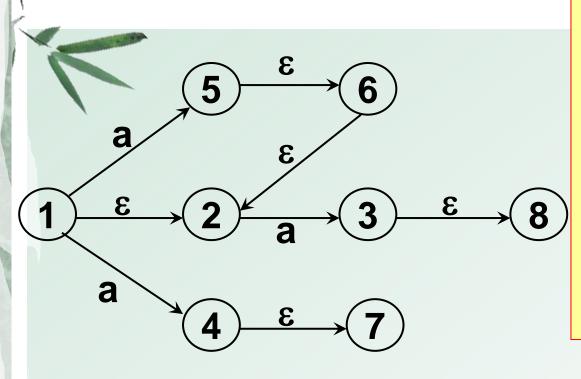


NFA变DFA: 确定化

设I是的状态集的一个子集,定义I的 ϵ -闭包 ϵ -closure(I)为:

- i) 若s∈l, 则s∈ε-closure(l);
- ii)若s∈I,则从s出发经过任意条ε弧而能到达的任何状态s'都属于 ϵ -closure(I)即

ε-closure(I)=I∪{s' | 从某个s∈I出发经 过任意条ε弧能到达s' }



- 设a是Σ中的一个 字符,定义
 - I_a= ε-closure(J) 其中,J为I中的 某个状态出发经 过一条a弧而到 达的状态集合。

•
$$\epsilon$$
-closure({1})={1, 2}=I

$$J={5, 4, 3}$$

$$I_a = \varepsilon$$
-closure(J)= ε -closure({5, 4, 3})

$$=$$
{5, 4, 3, 6, 2, 7, 8}

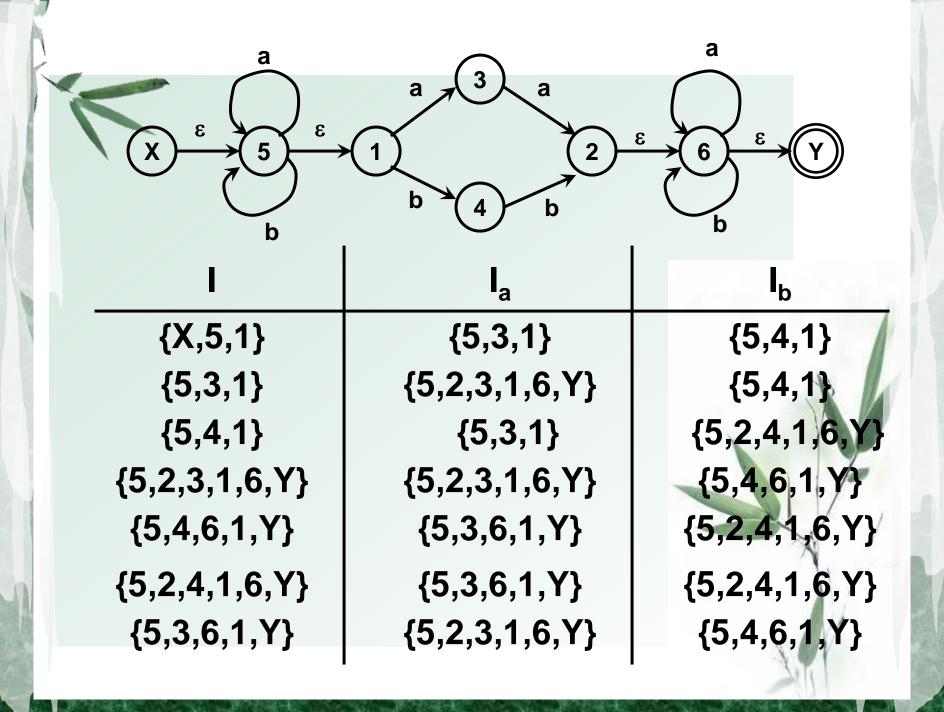
• 确定化的过程: 设字母表只包含两个a和b, 我们构造一张表:

7	7	
1	I_a	1 _b
ε-Closure({X})		
		7
		À \.×
		AV.
		1 MV

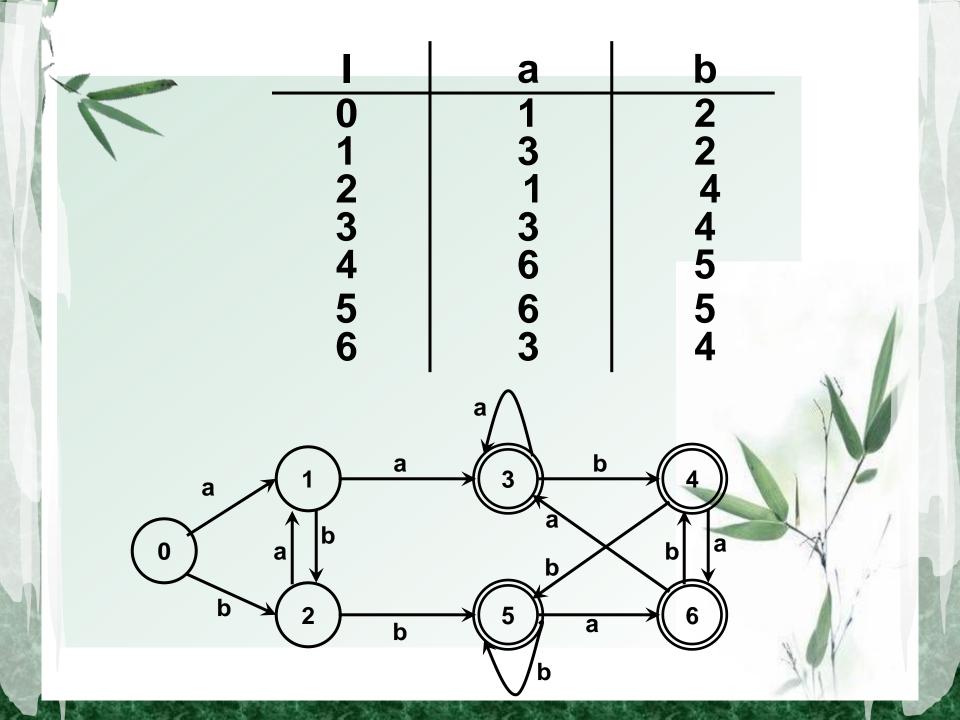
首先,置第1行第1列为 ϵ -closure({X}) 求出这一列的 I_a , I_b ;

然后,检查这两个I_a, I_b, 看它们是否已在表中的第一列中出现,把未曾出现的填入后面的空行的第1列上,求出每行第2,3列上的集合...

重复上述过程,知道所有第2,3列子集全部出现在第一列为止。



- 现在把这张表看成一个状态转换矩阵, 把其中的每个子集看成一个状态。
- 这张表唯一刻划了一个确定的有限自动机M,它的初态是ε-closure({X}),它的终态是含有原终态Y的子集。
- · 不难看出,这个DFA M与M'等价。





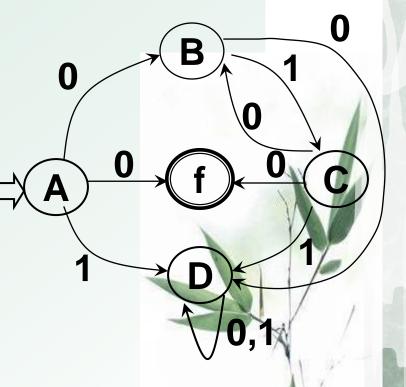
 对于正规文法G和有限自动机M,如果L(G)=L(M),则称G和M是等价的。 关于正规文法和有限自动机的等价性,有以下结论:

• 定理:

- 1. 对每一个右线性正规文法G或左线性正规文法G, 都存在一个有限自动机(FA) M, 使得L(M) =L(G)。
- 2. 对每一个FA M,都存在一个右线性正规文法 G_R 和左线性正规文法 G_L ,使得 $L(M) = L(G_R) = L(G_L)$ 。

例:

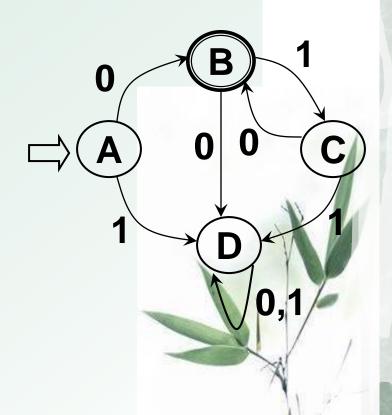
- $G_R(A)$: $A \to 0 \mid OB \mid 1D$ $B \to 0D \mid 1C$ $C \to 0 \mid OB \mid 1D$ $D \to 0D \mid 1D$
- 从G_R出发构造NFA M = 〈{A, □ \ A \ B, C, D, f}, {0, 1}, δ',
 A, {f}>, M的状态转换图 如右图所示。
- 显然 L(M) = L(G_R)。



例:

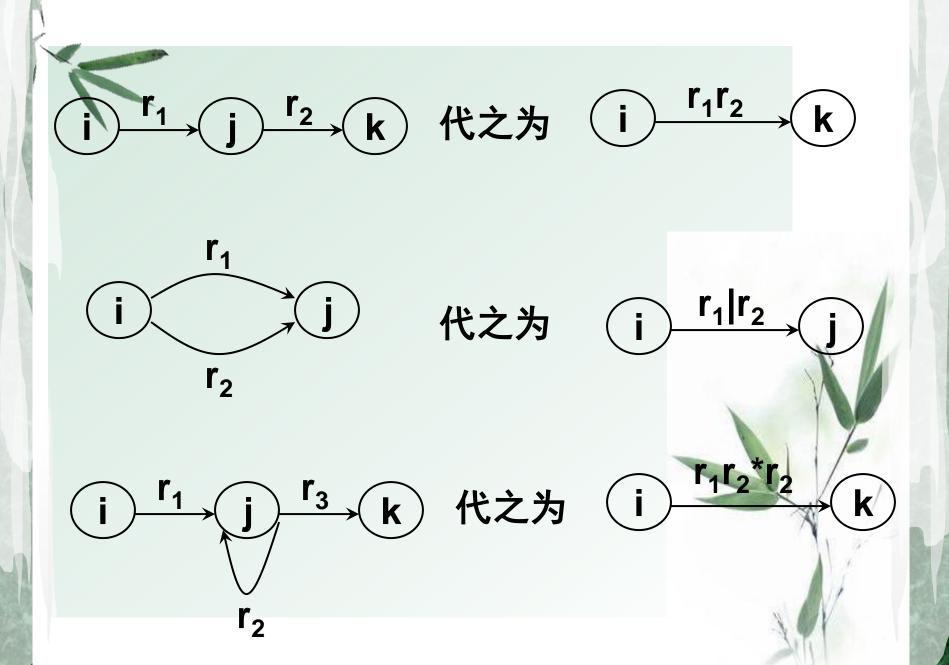
- L(M) = 0(10)*
- G_R = <{0, 1}, {A, B, C, D}, A, P>, 其中P由下列 产生式组成:

$$L(G_R) = L(M) = 0(10)*$$



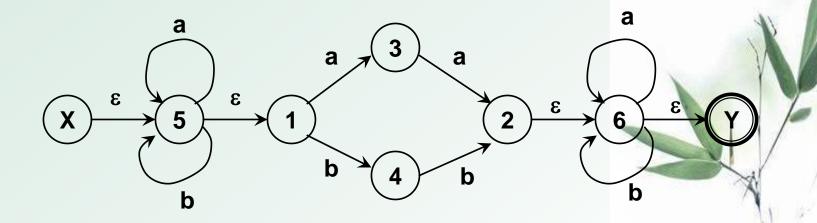
3.3.5 正规式与有限自动机的等价性

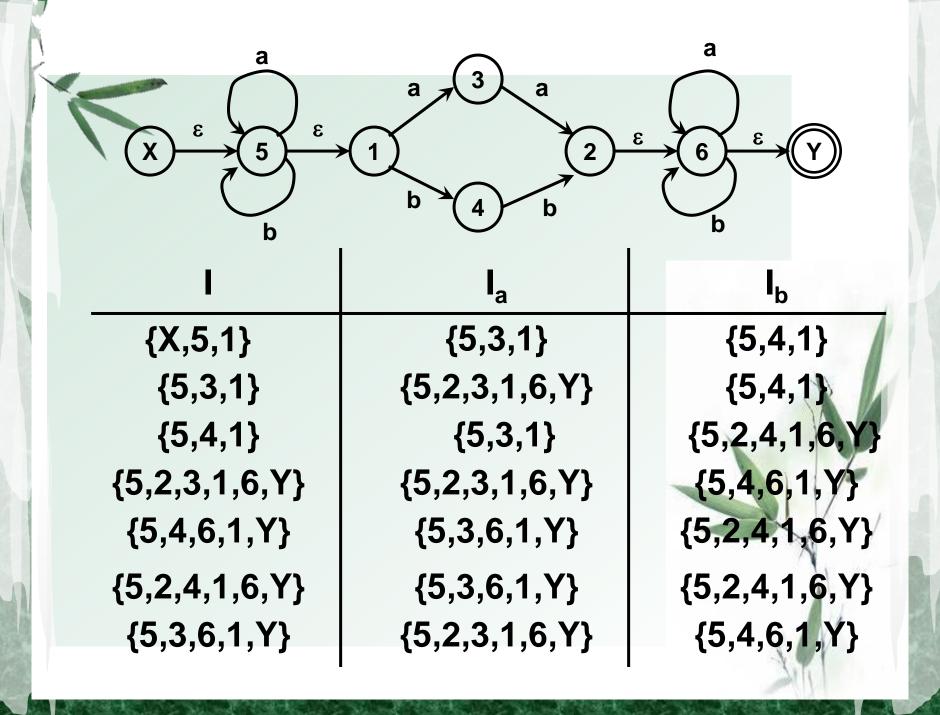
- 定理:
 - 1. 对任何FA M, 都存在一个正规式r, 使得 L(r)=L(M)。
- 2. 对任何正规式r,都存在一个FA M,使得 L(M)=L(r)。
- □ 对转换图概念拓广,令每条弧可用一个正规式作标记。(对一类输入符号)

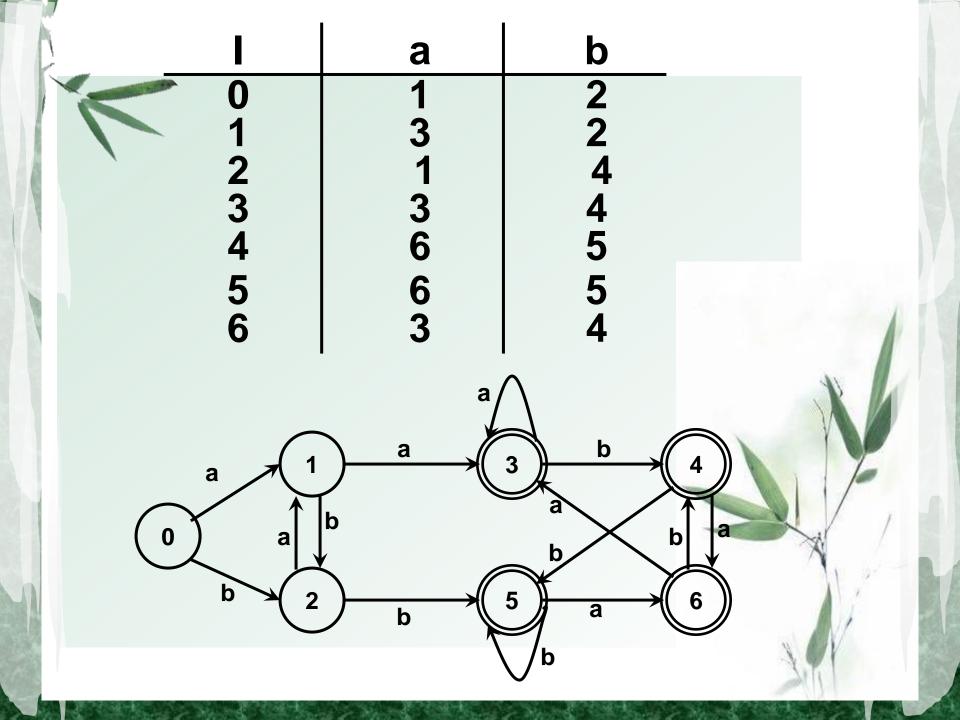


• (a|b)*(aa|bb)(a|b)*









3.3.6 确定有限自动机的化简

- 对DFA M的化简: 寻找一个状态数比M少的DFA M', 使得L(M)=L(M')
- 假设s和t为M的两个状态,称s和t等价:如果从状态s出发能读出某个字α而停止于终态,那么同样,从t出发也能读出α而停止于终态;反之亦然。
- 两个状态不等价,则称它们是可区别的。

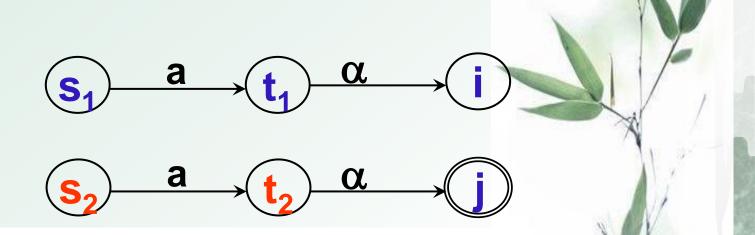
• 对一个DFA M最少化的基本思想:

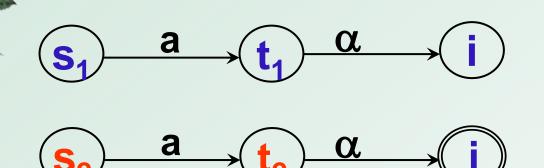
把M的状态集划分为一些不相交的子集,使得任何两个不同子集的状态是可区别的,而同一子集的任何两个状态是等价的。最后,让每个子集选出一个代表,同时消去其他状态。

具体做法:对M的状态集进行划分

- 首先,把S划分为终态和非终态两个子集,形成基本划分∏。
- 假定到某个时候, Π已含m个子集, 记为Π={I⁽¹⁾, I⁽²⁾, ..., I^(m)}, 检查Π中的每个子集看是否能进一步划分:
 - 对某个I⁽ⁱ⁾, 令I⁽ⁱ⁾={s₁, s₂, ..., s_k}, 若存在一个输入字符a使得I_a⁽ⁱ⁾ 不会包含在现行Π的某个子集I^(j)中,则至少应把I⁽ⁱ⁾分为两个部分。

例如,假定状态 \mathbf{s}_1 和 \mathbf{s}_2 经a弧分别到达 \mathbf{t}_1 和 \mathbf{t}_2 ,而 \mathbf{t}_1 和 \mathbf{t}_2 属于现行 Π 中的两个不同子集,说明有一个字 α , \mathbf{t}_1 读出 α 后到达终态,而 \mathbf{t}_2 读出 α 后不能到达终态,或者反之,那么对于字a α , \mathbf{s}_1 读出a α 后到达终态,而 \mathbf{s}_2 读出a α 不能到达终态,或者反之,所以 \mathbf{s}_1 和 \mathbf{s}_2 不等价。

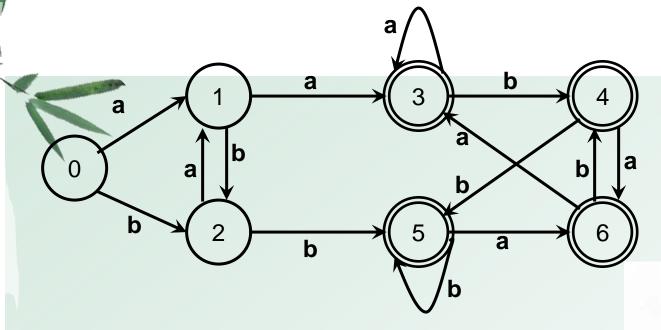




则将 I (i) 分成两半,使得一半含有 s₁:
 I (ii) = {s | s ∈ I (i) 且 s 经 a 弧 到 达 t,
 且 t 与 t₁属于现行 П 中 的 同一子集}

另一半含有s₂: | (i²)=| (i)-| (i1)

- 一般地,对某个a和I⁽ⁱ⁾,若 I_a ⁽ⁱ⁾ 落入现行 Π 中 N个不同子集,则应把I⁽ⁱ⁾划分成N个不相交的组,使得每个组J的 J_a 都落入的 Π 同一子集。这样构成新的划分。
- 重复上述过程,直到∏所含子集数不再增长。
- 对于上述最后划分∏中的每个子集,我们选取每个 子集I中的一个状态代表其他状态,则可得到化简 后的DFA M'。
- · 若 | 含有原来的初态,则其代表为新的初态,若 | 含有原来的终态,则其代表为新的终态。

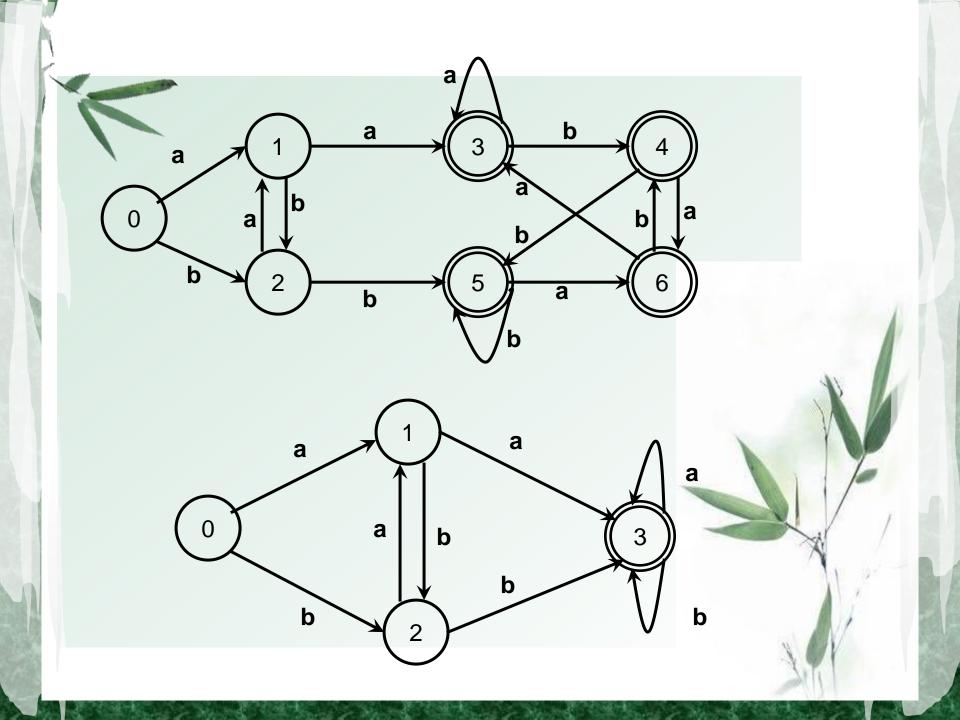


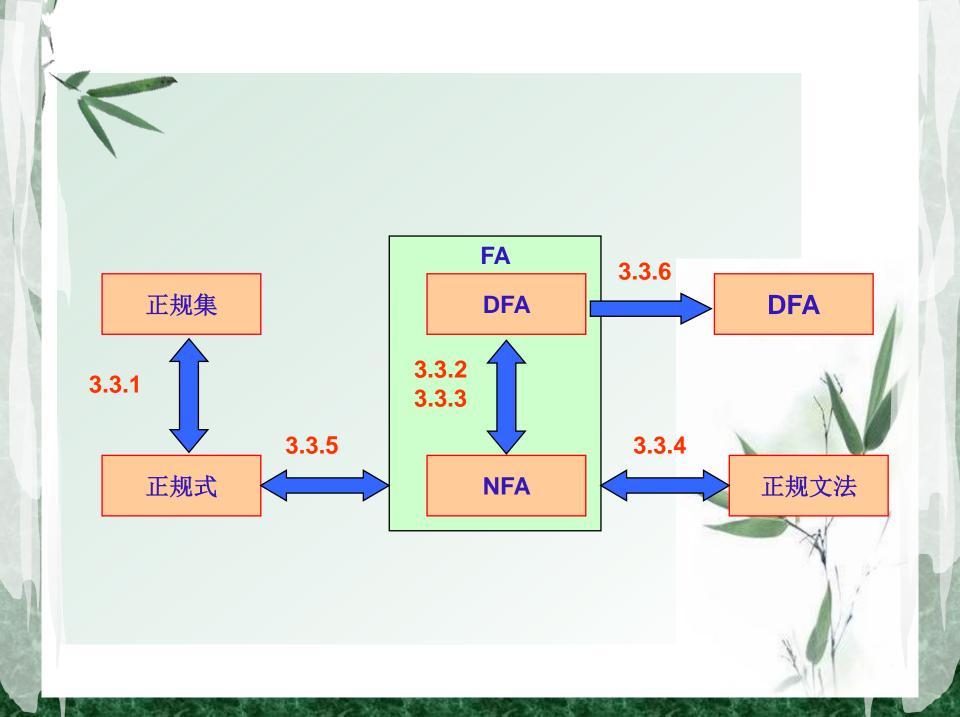
 $I_a^{(2)} = \{3, 6\} \quad I_a^{(2)} = \{4, 5\}$

$$I_{a}^{(1)} = \{0, 1, 2\} \quad I_{a}^{(2)} = \{3, 4, 5, 6\}$$

$$I_{a}^{(1)} = \{1, 3\} \quad I_{a}^{(11)} = \{0, 2\} \quad I_{a}^{(12)} = \{1\} \quad I_{a}^{(2)} = \{3, 4, 5, 6\}$$

$$I_{a}^{(11)} = \{0, 2\} \quad I_{a}^{(11)} = \{1\} \quad I_{b}^{(11)} = \{2, 5\} \quad I_{a}^{(11)} = \{0\} \quad I_{a}^{(112)} = \{2\} \quad I_{a}^{(12)} = \{1\} \quad I_{a}^{(2)} = \{3, 4, 5, 6\}$$





3.4 词法分析器的自动产生--LEX

LEX源程序

词法分析程序 自动产生器 词法分析程序L

输入串

词法分析程序L

单词符号

AUXILIARY DEFINITION

```
letter\rightarrowA|B|...|Z digit \rightarrow0|1|...|9
```



RECOGNITION RULES

```
DIM
                          { RETURN (1,-) }
                          { RETURN (2,-) }
      IF
                          { RETURN (3,-) }
      DO
                           RETURN (4,-) }
      STOP
5
                          { RETURN (5,-) }
      END
6
      letter(letter|digit) * { RETURN (6, TOKEN) }
                          { RETURN (7, DTB) }
      digit(digit)*
8
                          { RETURN (8, -) }
                           RETURN (9,-) }
                           { RETURN (10,-) }
10
                           { RETURN (11,-) }
11
      **
                          { RETURN (12,-) }
12
                          { RETURN (13,-) }
13
                          { RETURN (14.-) }
14
```

· LEX的工作过程:

- 首先,对每条识别规则P_i构造一个相应的非确定有限自动机M_i;
- 然后,引进一个新初态X,通过ε弧,将这些自动机 连接成一个新的NFA;
- 最后,把M确定化、最小化,生成该DFA的状态转换 表和控制执行程序

