

编译原理

Compiler Construction Principles



朱 青

信息学院计算机系，
中国人民大学，

zqruc2012@aliyun.com



第2章:词法分析 (Lexical Analysis)

⌘ 2.1 词法分析程序的功能

⌘ 2.2 词法分析器的设计

⌘ 2.3 正规表达式 (Regular Expression)

⌘ 2.4 有限自动机

⌘ 2.5 词法分析器的自动生成

2.4 有限自动机

⌘ 2.4.1 确定有限自动机 (DFA)

⌘ 2.4.2 非确定有限自动机 (NFA) 确定化

⌘ 2.4.3 具有 ε -转移的NFA M确定化

⌘ 2.4.4 DFA的化简

⌘ 2.4.5 正规式与有限自动机的等价性

⌘ 2.4.6 正规文法与有限自动机

2.4.1 确定有限自动机 (DFA)

有限自动机

正规表达式 \equiv 有限自动机

DFA NFA 正规式

2.4.1 确定有限自动机 (DFA)

- DFA的定义： 一个确定的有限自动机 (DFA)

M是一个五元式

$$M = (S, \Sigma, f, s_0, Z)$$

其中

- 1 S 是一个有限集，它的每个元素称为一个状态；
- 2 Σ 是一个有穷字母表，它的每个元素称为一个输入字符。

3 f 是一个从 $S \times \Sigma$ 至 S 的（单值）部分映照 $f(s, a) = s'$ 。表示：当前状态是 s ，输入字符是 a 时，下一个状态是 s' ， s' 叫 s 的后继状态。

4 $s_0 \in S$ ，是唯一的初态。

5 $Z \subset S$ ，是终态集（可空）。

•

- 一个DFA可以表示为一个矩阵。

$f(s, a)$ 值-----状态转换矩阵值。

行-----状态,

列-----输入字符。

- 一个DFA也可以表示为一张（确定的）状态转换图。

例题 2.4-1 DFA M 可用一个矩阵表示

DFA $M = (\{0, 1, 2, 3\}, \{a, b\}, f, 0, \{3\})$.

其中 f :

$$f(0,a)=1 \qquad f(0,b)=2$$

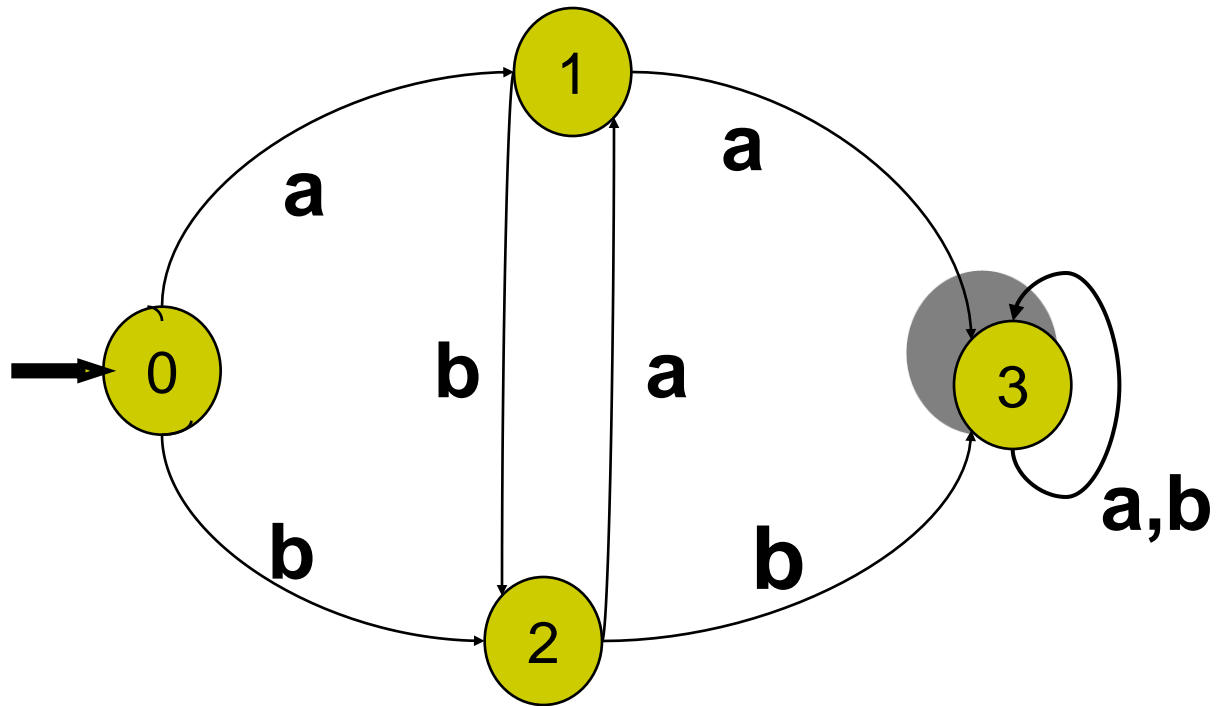
$$f(1,a)=3 \qquad f(1,b)=2$$

$$f(2,a)=1 \qquad f(2,b)=3$$

$$f(3,a)=3 \qquad f(3,b)=3$$

矩阵表示:

状态 \ 字符	字符	
	a	b
0	1	2
1	3	2
2	1	3
3	3	3



例题 2.4-2 DFA M 可用一张（确定的）状态转换图表示。

DFA $M = (\{A, B, C, D\}, \{a, b\}, \Theta, A, \{D\})$

$$\Theta(A, a) = B$$

$$\Theta(A, b) = A$$

$$\Theta(B, a) = B$$

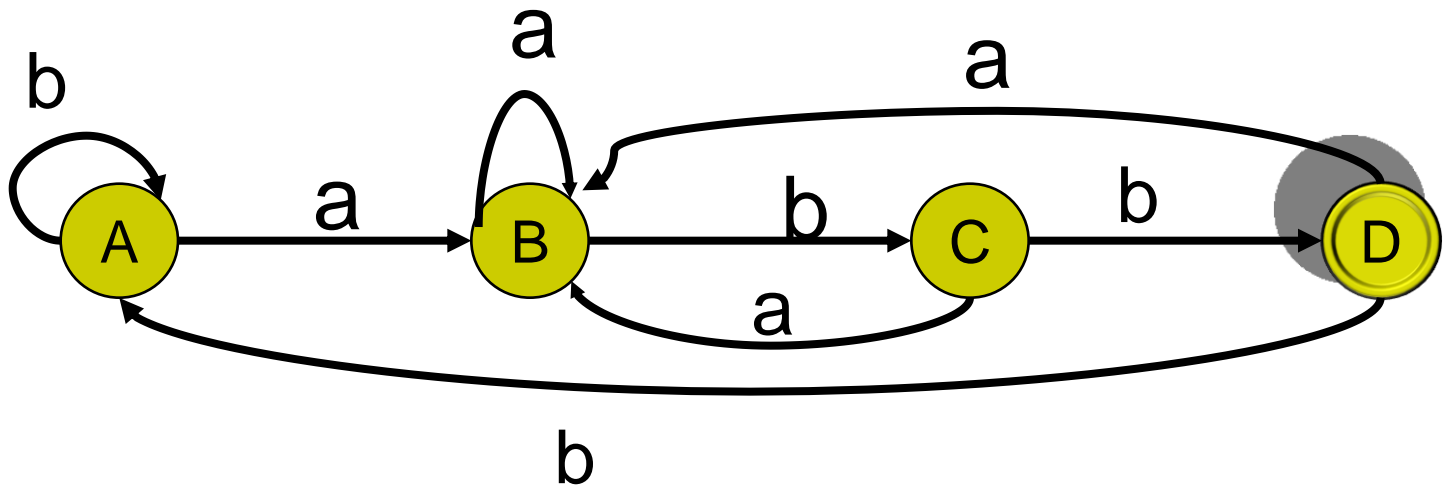
$$\Theta(B, b) = C$$

$$\Theta(C, a) = B$$

$$\Theta(C, b) = D$$

$$\Theta(D, a) = B$$

$$\Theta(D, b) = A$$



用**DFA**识别单词符号：

在**DFA**状态转换图中，存在一条从初态到终态的通路。该路上各弧的标记字符依此连结构成的字与S相同，则称S能被该**DFA**识别。

一个**DFA** M 所能识别的所有的字的集合记为： $L(M)$ 。

定理

- Σ 上的一个 字集 $V \subset \Sigma^*$ 是正规的，当且仅当存在 Σ 上的**DFA M**，使得 **$V=L(M)$** 。
- **DFA** 的确定性表现在 **f** 是一个从 **$S \times \Sigma$** 至 **S**的单值函数。即唯一确定了下一个状态。

2.4.2 非确定有限自动机 (NFA) 确定化

- **NFA的定义:**

一个非确定的有限自动机 (NFA) **M**
是一个五元式

$$\mathbf{M} = (\mathbf{S}, \Sigma, \mathbf{f}, \mathbf{S}_0, \mathbf{Z})$$

其中

1 S 是一个有限集，它的每个元素称为一个状态；

2 Σ 是一个有穷字母表，它的每个元素称为一个输入字符。

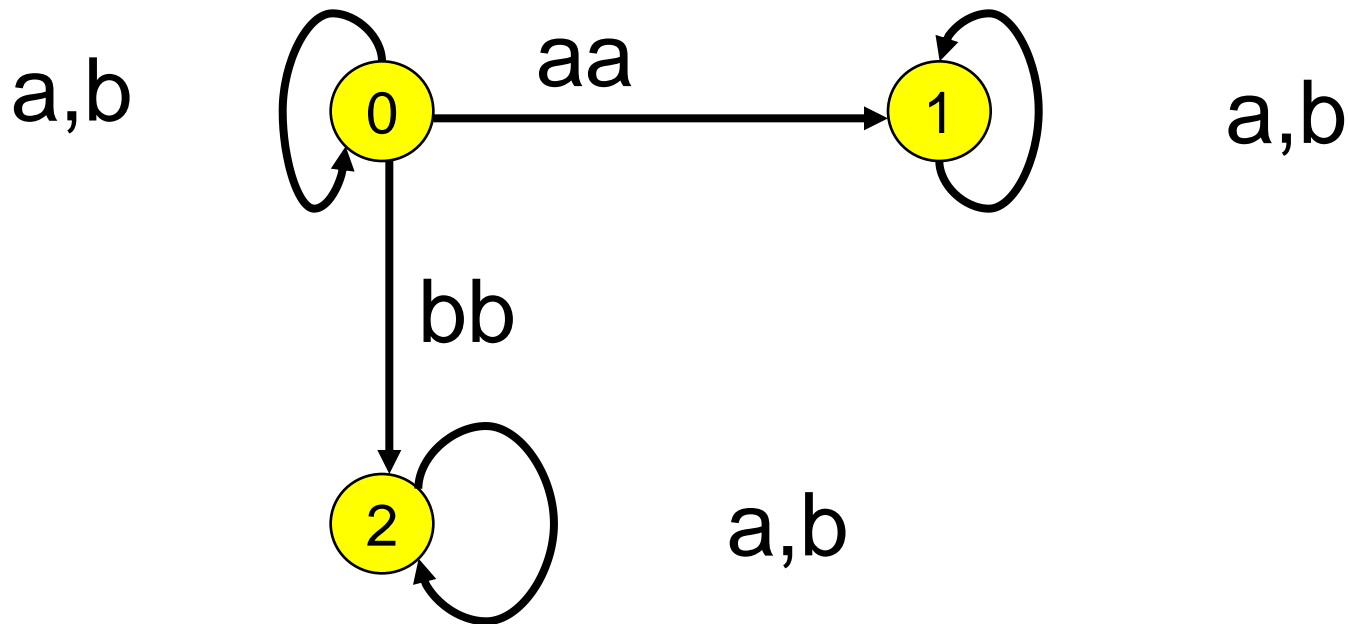
3 f 是一个从 $S \times \Sigma$ 至 S 的子集映照。

即 $f : S \times \Sigma^* \longrightarrow 2^S$ 。

4 $S_0 \subset S$ ，是一个非空初态集。

5 $Z \subset S$ ，是终态集（可空）。

NFA实例:



-
- 一个 NFA 可以表示为一个矩阵。

$f(s, a)$ 值-----状态转换矩阵。

行-----状态,

列-----输入字符

- 一个NFA也可以表示为一张（确定的）状态转换图。

例题 2.4-3 NFA M 可用一张状态转换图表示。

NFA $M = (\{A, B, C, D\}, \{a, b\}, \delta, A, \{D\})$

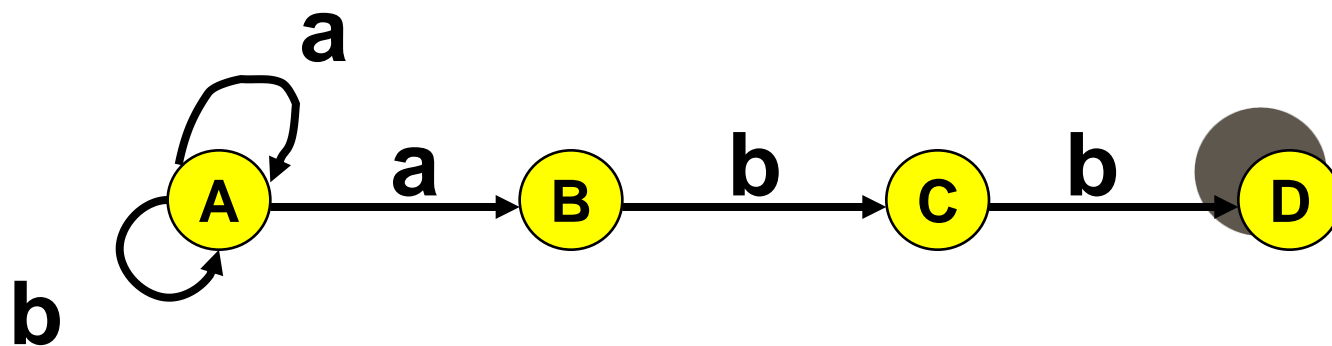
$$\delta(A, a) = \{A, B\} \quad \delta(A, b) = \{A\}$$

$$\delta(B, a) = \Phi \quad \delta(B, b) = \{C\}$$

$$\delta(C, a) = \Phi \quad \delta(C, b) = \{D\}$$

$$\delta(D, a) = \Phi \quad \delta(D, b) = \Phi$$

状态图：



例题 2.4-4 NFA M 实例。

NFA $M = (\{S, Q, U, V, Z\}, \{0, 1\}, \delta, \{S\}, \{Z\})$;

$$\delta(S, 0) = \{V, Q\} \quad \delta(S, 1) = \{U, Q\}$$

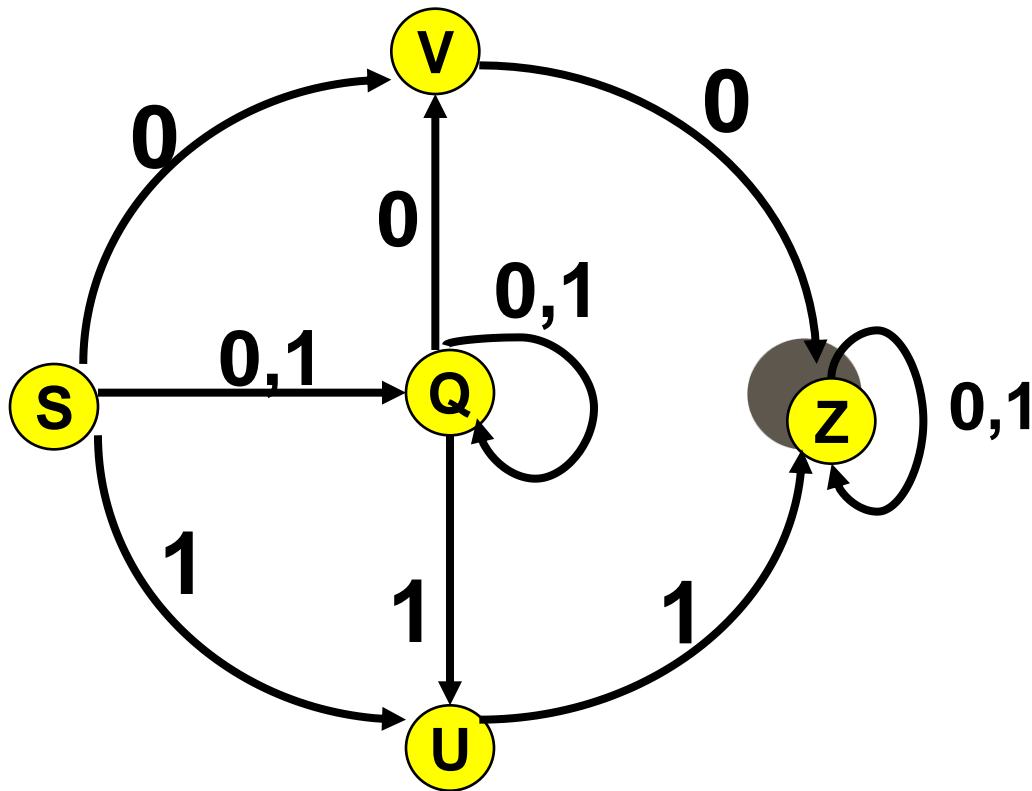
$$\delta(U, 0) = \Phi \quad \delta(U, 1) = \{Z\}$$

$$\delta(V, 0) = \{Z\} \quad \delta(V, 1) = \Phi$$

$$\delta(Q, 0) = \{V, Q\} \quad \delta(Q, 1) = \{U, Q\}$$

$$\delta(Z, 0) = \{Z\} \quad \delta(Z, 1) = \{Z\}$$

NFA的状态转换图：



图中某些状态射出两条具有相同标记的弧S,Q.

识别单词：

在NFA中识别 ϵ 有两种情况：

- (1) 终态与初态是同一点.
- (2) 从初态到终态, 所走的弧上都是 ϵ .

定理：对任何一个NFA M ，都存在一个DFA M' ，使得 $L(M') = L(M)$ 。

证明的思想是由 M 出发构造等价的 M' ，办法是让 M' 的状态对应于 M 的状态集合。

即若 $f(q,a)=\{q_1,q_2,\dots,q_k\}$ ，
我们将集合 $\{q_1,q_2,\dots,q_k\}$ 作为一个整体看作 M' 中的一个状态，即 S' 中的一个元素。

例题2.4-5： 设：

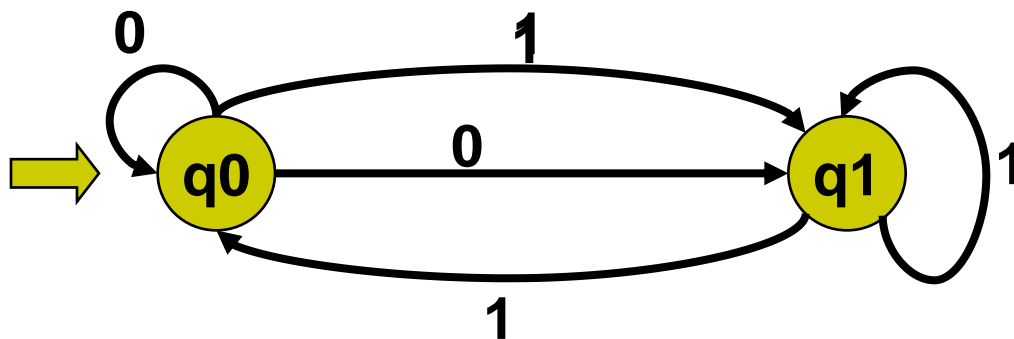
NFA $M = (\{0,1\}, \{q_0, q_1\}, f, q_0, \{q_1\})$.

$$f(q_0, 0) = \{q_0, q_1\} \quad f(q_0, 1) = \{q_1\}$$

$$f(q_1, 0) = \Phi \quad f(q_1, 1) = \{q_0, q_1\}$$

构造与M 等价的DFA M' .

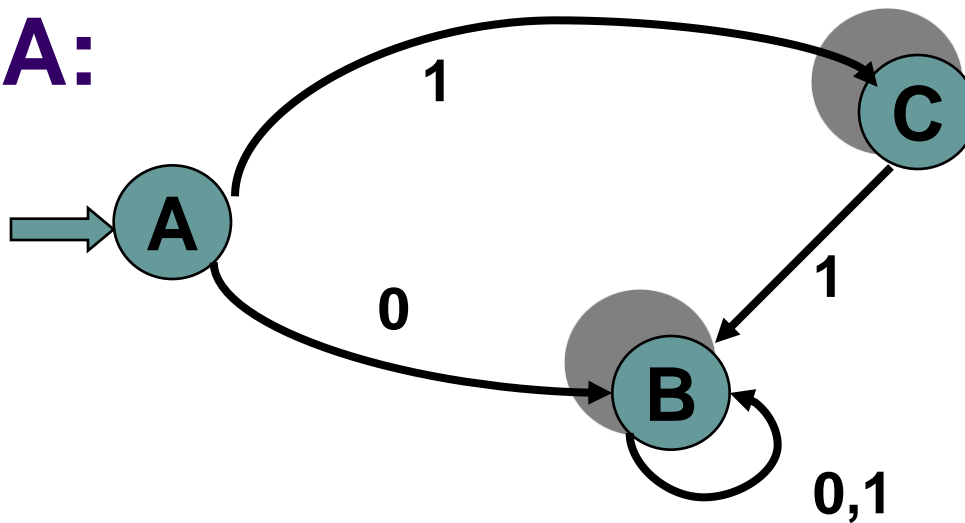
状态图:



解: 令 **DFA $M' = (\{0, 1\}, Q', f', q0', F')$**

	0	1
A.{q0}	{q0,q1}	{q1}
B.{q0,q1}	{q0,q1}	{q0,q1}
C.{q1}	Φ	{q0,q1}

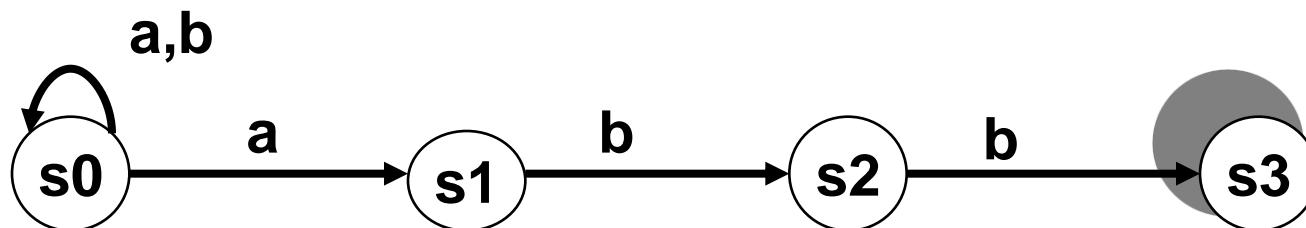
DFA:



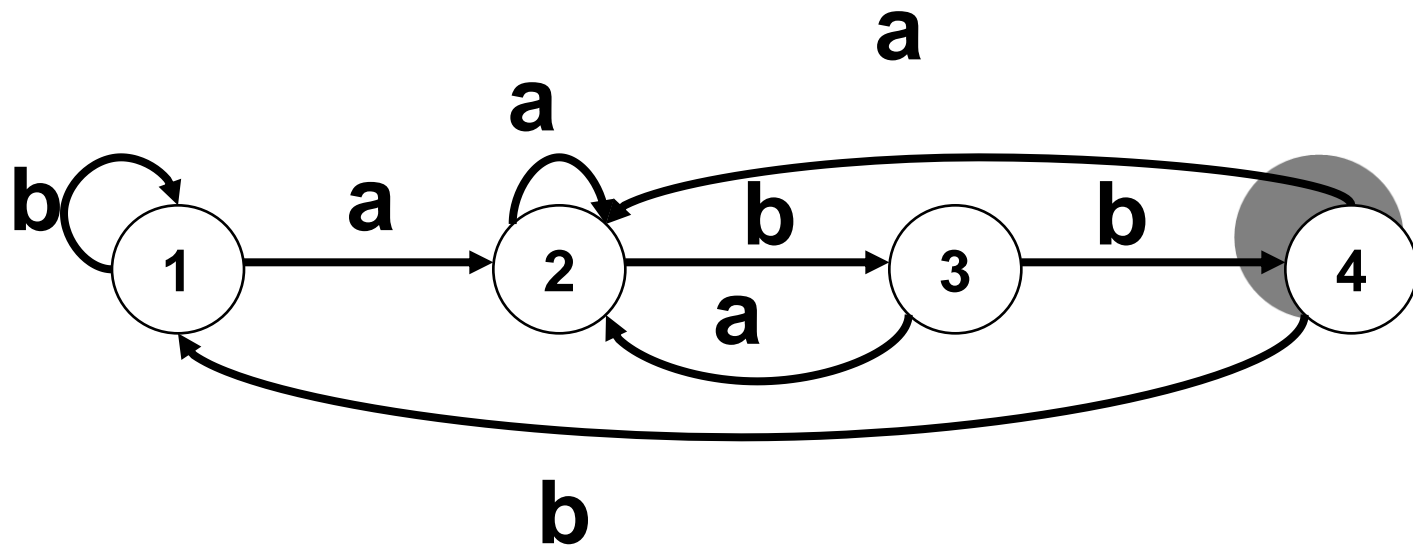
例2.4-6: 设有

NFA $M = (\{S0, S1, S2, S3\}, \{a, b\}, f, S0, \{S3\})$

状态转换图:



	a	b
1. {s0}	{s0,s1}	{s0}
2. {s0,s1}	{s0,s1}	{s0,s2}
3. {s0,s2}	{s0,s1}	{s0,s3}
4. {s0,s3}	{s0,s1}	{s0}



2.4.3 具有 ε -转移的NFA M确定化

定义1：状态集 S 的子集 I 的 ε -闭包，

即： ε -closure(I):

- (1) 若 $s \in I$ ，则 $s \in \varepsilon$ -closure(I).
- (2) 若 $s \in I$ ，则从 s 出发经过若干条 ε 弧所到达的状态 s' ，
 $s' \in \varepsilon$ -closure(I).

例题2.4-7 具有 ϵ 转移的识别正规式：
 $aa^*|bb^*$ 的非确定有限自动机.

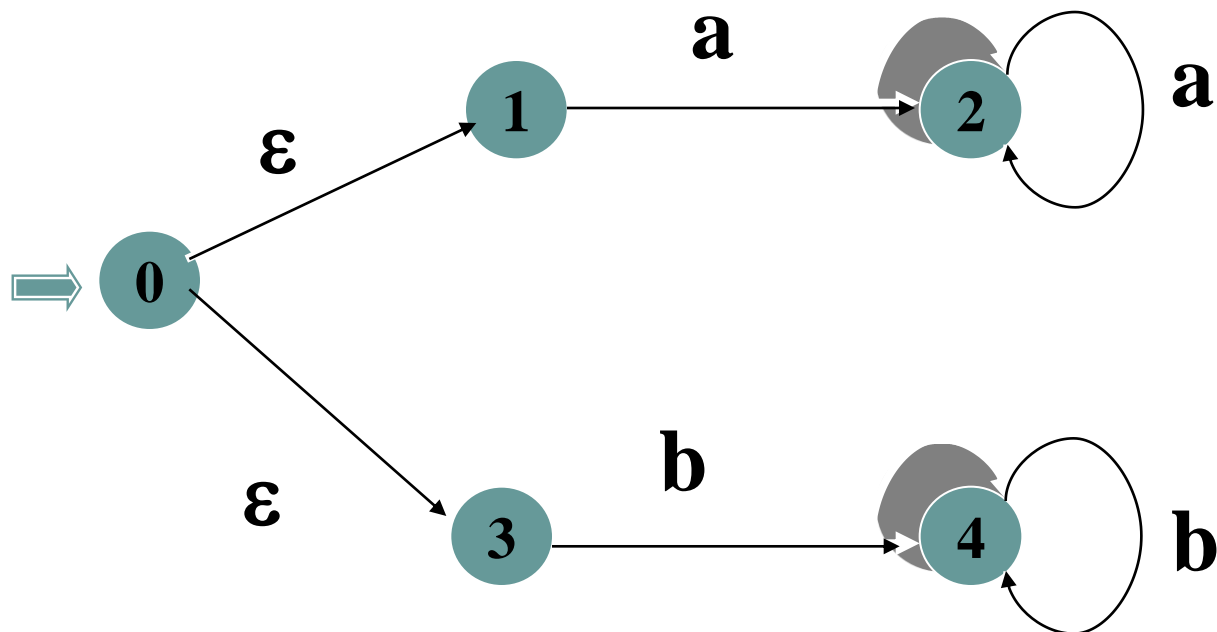
NFA $M = (\{0, 1, 2, 3, 4\}, \{a, b\}, f, \{0\}, \{2, 4\})$

$f(1, a) = \{2\}$ $f(3, b) = \{4\}$

$f(2, a) = \{2\}$ $f(4, b) = \{4\}$

$f(0, \epsilon) = \{1, 3\}$

状态转换图：



$$\varepsilon\text{-closure}(0) = \{0, 1, 3\}$$

定理1：对任何一个具有 ε -转移的NFA
M, 一定存在一个不具有 ε -转移
的 NFA M', 使
$$L(M') = L(M)。$$

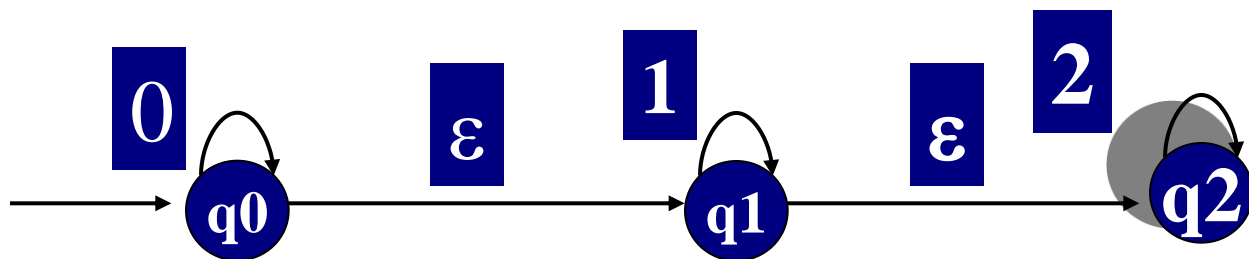
例题2.4-8:

设NFA $M = (\Sigma, Q, f, q_0, F)$

$\Sigma = \{0, 1, 2\}$, $Q = \{q_0, q_1, q_2\}$,

$F = \{q_2\}$

状态图:



从M出发构造一个不具有 ε - 转移的
NFA M' ，使得 $L(M') = L(M)$ 。

解：令 $M' = (\Sigma, Q', f', q_0, F')$ ，
其中 Σ, Q', q_0 的意义同M中完全一样。

1) F' 包含 M 的终态集 F ，其次若 M 中从 q_0 出发有一条到达某终态的 ε 道路，则将 q_0 加在 F' 中。

$$F' = \begin{cases} F \cup \{q_0\}, & \text{若 } \varepsilon\text{-closure}(q_0) \\ & \text{包含 } F \text{ 的一个状态} \\ F, & \text{否则。} \end{cases}$$

终态集 $F' = \{q_0, q_1, q_2\}$,
 ε -closure(q_0)= $\{q_0, q_1, q_2\}$
 ε -closure(q_1)= $\{q_1, q_2\}$

2)

$f'(q, a) = \{q' \mid q' \text{ 为从 } q \text{ 出发先经若干 } \varepsilon \text{ 箭弧, 接着经一个标记为 } a \text{ 的箭弧, 再经若干 } \varepsilon \text{ 箭弧组成的道路所能到达的状态。}\}$

$$f'(q_0, 0) = \{q_0, q_1, q_2\}$$

$$f'(q_0, 1) = \{q_1, q_2\}$$

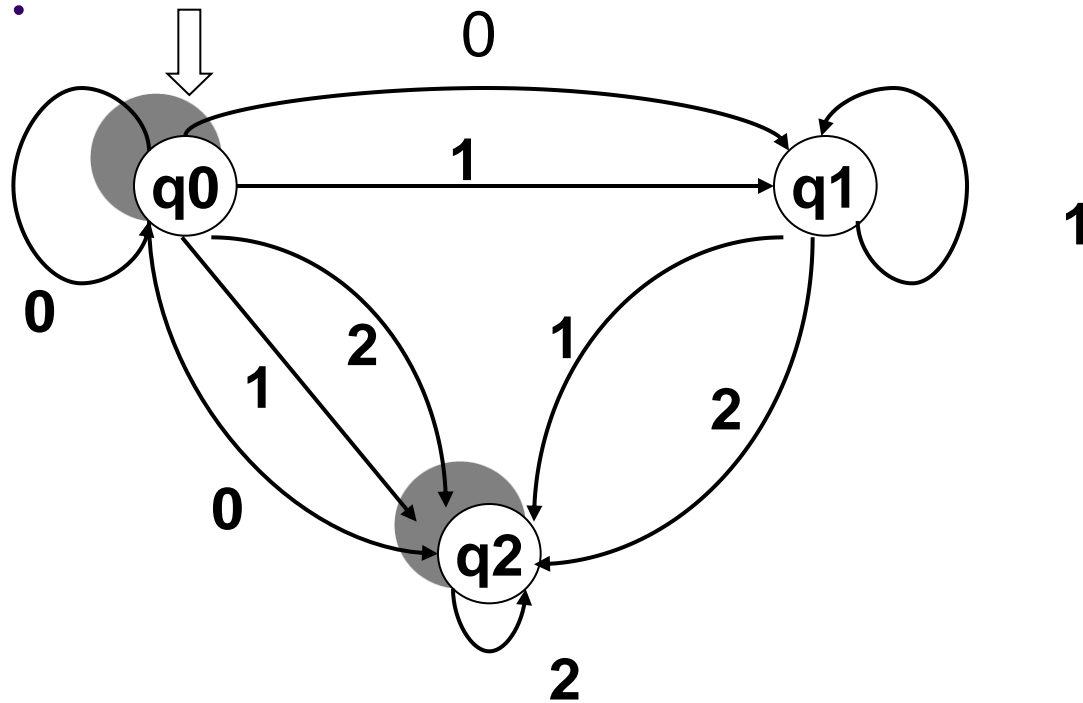
$$f'(q_0, 2) = \{q_2\}$$

$$f'(q_1, 1) = \{q_1, q_2\}$$

$$f'(q_1, 2) = \{q_2\}$$

$$f'(q_2, 2) = \{q_2\}$$

状态图:



例题2.4-7 转化为不具有 ϵ 转移的非确定有限自动机.

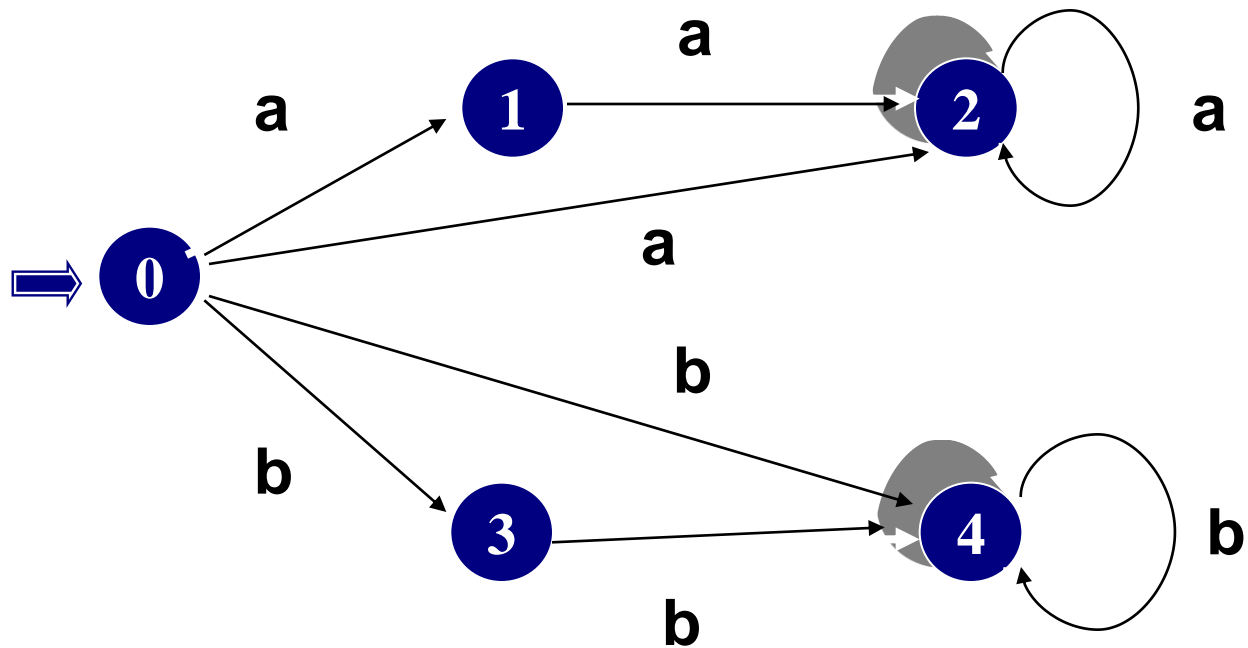
NFA $M' = (\{0, 1, 2, 3, 4\}, \{a, b\}, f', \{0\}, \{2, 4\})$

$$f'(0, a) = \{1, 2\} \quad f'(0, b) = \{3, 4\}$$

$$f'(1, a) = \{2\} \quad f'(3, b) = \{4\}$$

$$f'(2, a) = \{2\} \quad f'(4, b) = \{4\}$$

状态转换图：NFA



定理2: 对于字母表 Σ 上任何一个具有 ε -转移的NFA M , 一定存在一个的 DFA M' ,

使得: $L(M') = L(M)$ 。

将 ϵ -NFA 转化为等价的 DFA

扩充 ϵ -closure(q)。假设 T
是NFA状态的一个集合,

ϵ -closure(T):

表示所有那些可以从 T 中的元素出发经过
一条 ϵ 道路所能到达的NFA 的状态的全
体所组成的集合。

NFA $M = (S, \Sigma, f, S_0, Z')$

用构造 ϵ -closure(T) 的方法

实现DFA M' 的转换:

DFA $M = (S', \Sigma, f', q_0, Z')$

基本思想: 1)首先从 S_0 出发, 仅经过任意条 ϵ 箭弧所能到达的状态所组成的集合作为 M' 的初态 q_0 .

2) 分别把从 q_0 出发，经过对输入符号 $a \in \Sigma$ 的状态转移所能到达的状态（包括转移后再经 ε - 箭弧所能到达的状态）所组成的集合作为 M' 的状态，如此继续，直到不再有新的状态为止。

2.4.4 DFA的化简

DFA的化简是指：寻找一个状态数比M少的
DFA M' ，

使得 $L(M) = L(M')$

- 状态S与T是等价：
- 两个状态是可区分的：
- DFA M状态最少化：

最小化：是把状态集**S**分割成一些不相交的子集，不同子集中的状态是可区分的，而同一子集中的状态是等价的，用一个状态代表一个子集，并消去该子集中的其它状态，从而得到化简的**DFA**。

对M的状态集S进行分化的步骤:

1) 把S的终态与非终态分开, 生成两个子集, 形成基本分割 Π 。

2) 若某一时刻分割 Π 已包含M个子集
 $\Pi = \{I_1, I_2, \dots, I_M\}$ 。

检查每个子集 I_i , 看其是否可再分割。

$$I_i = \{S_1, S_2, \dots, S_k\}$$

设 S_1, S_2 对任意一个输入字符 a , 其后继状态 t_1, t_2 , 不完全包含在某个 I_j 中, 则有 s_1, s_2 是可区分的。

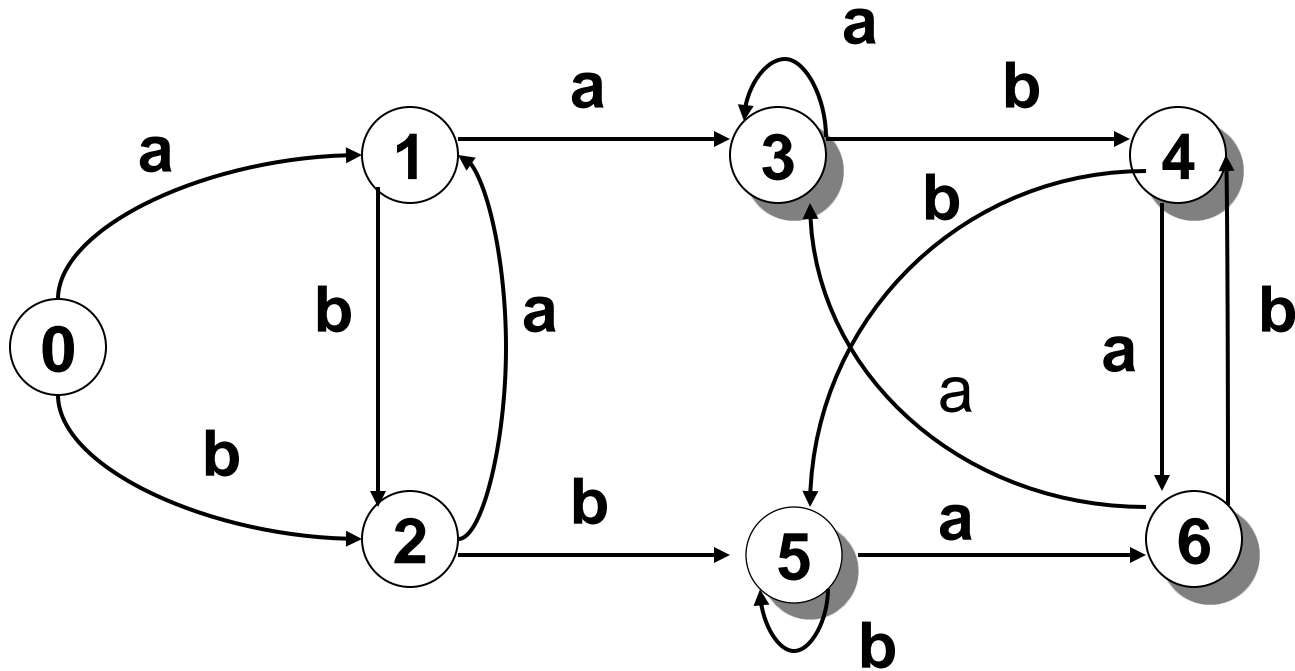
所以, 把 I_i 一分为二, 使其一半包含 S_1 , 另一半包含 S_2 。

重复上述过程，直至 Π 所含子集不再增加为止。

对于最后分割的每个子集，在该子集中选出一个状态的代表。

若该子集含有原来的初（终）态，则该状态为新的初（终）态。

例题2.4-9：未化简的DFA M（P51）



解： 1) $\{3, 4, 5, 6\}$, $\{0, 1, 2\}$
2) $\{3, 4, 5, 6\}_a = \{3, 6\}$
 $\{3, 4, 5, 6\}_b = \{4, 5\}$ 属于
 $\{3, 4, 5, 6\}$ 不能再分。
 $\{0, 1, 2\}_a = \{1, 3\}$
3) $\{1, 3\}$ 所生成的集合没有完全包含在 $\{3, 4, 5, 6\}$ 和 $\{0, 1, 2\}$ 中，故 $\{0, 1, 2\}$ 一分为二。

1 经a 弧到3, $3 \in \{3, 4, 5, 6\}$;
0, 2 经a 弧到1, $1 \in \{0, 1, 2\}$;
将 $\{0, 1, 2\}$ 分成 $\{0, 2\}$, $\{1\}$ 。

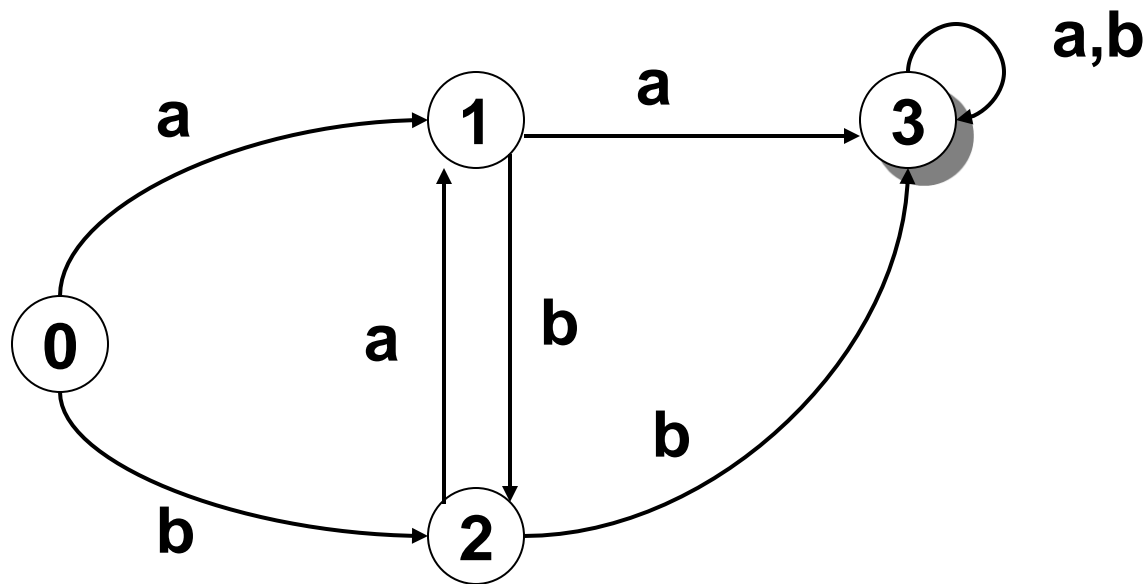
4) 再检查 $\{0, 2\}$ 。

$\{0, 2\}_b = \{2, 5\}$

$\{2, 5\}$ 不完全落在 $\{0, 2\}$, $\{3, 4, 5\}$ 中, 故 $\{0, 2\}$ 再分成 $\{0\}$, $\{2\}$ 。

最后得到: $\{0\}$, $\{1\}$, $\{2\}$, $\{3, 4, 5, 6\}$ 。

令：3代表 $\{3, 4, 5, 6\}$ ，得到：



$$M = (\{0,1,2,3\} , \{a,b\} , f , 0 , \{3\})$$

$$f(0,a) = 1$$

$$f(0,b)=2$$

$$f(1,a) = 3$$

$$f(1,b)=2$$

$$f(2,a) = 1$$

$$f(2,b)=3$$

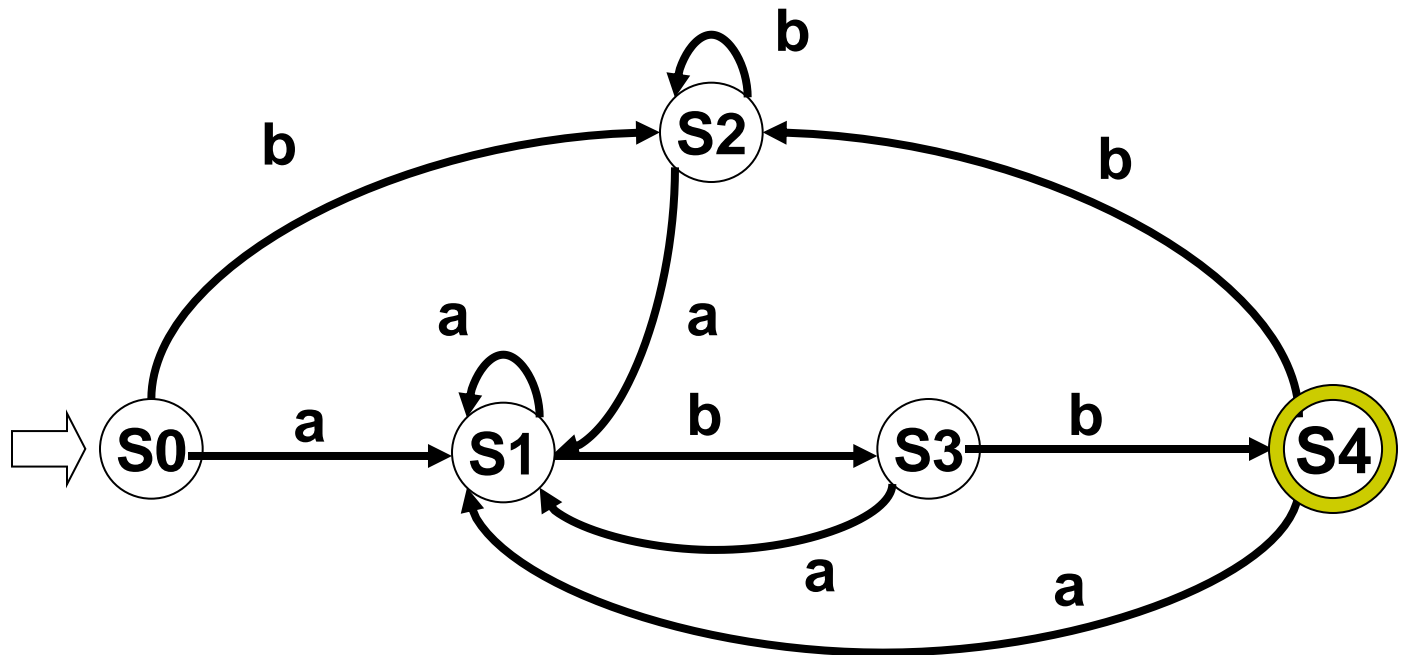
$$f(3,a) = 3$$

$$f(3,b)=3$$

例题2.4-10 将下面的DFA最小化：

$M = (\{S0, S1, S2, S3, S4\}, \{a, b\}, f, S0, \{S4\})$

	a	b
S0	S1	S2
S1	S1	S3
S2	S1	S2
S3	S1	S4
S4	S1	S2



解：

1) 初始划分： $\{S_0, S_1, S_2, S_3\}, \{S_4\}$

2) 考察 $\{S_0, S_1, S_2, S_3\}$,

$\{S_0, S_1, S_2, S_3\} \text{a} = \{S_1\} \subset \{S_0, S_1, S_2, S_3\}$

$\{S_0, S_1, S_2\} \text{b} = \{S_2, S_3\}, \{S_3\} \text{b} = \{S_4\}$

$\{S_0, S_1, S_2, S_3\}$ 不包含在同一子集中。

一分为二：

NEW: $\{S_0, S_1, S_2\}, \{S_3\}, \{S_4\}$

3)考察 $\{S_0, S_1, S_2\}$,

$\{S_0, S_1, S_2\} \text{a} = \{S_1\} \subset \{S_0, S_1, S_2\}$

$\{S_0, S_2\} \text{b} = \{S_2\}, \{S_1\} \text{b} = \{S_3\}$

$\{S_0, S_1, S_2\}$ 不包含在同一子集中。一分为二：

NEW : $\{S_0, S_2\}, \{S_1\}, \{S_3\}, \{S_4\}$

..... 直到NEW不再改变。

4) S_0 作为 $\{S_0, S_2\}$ 的代表。

