



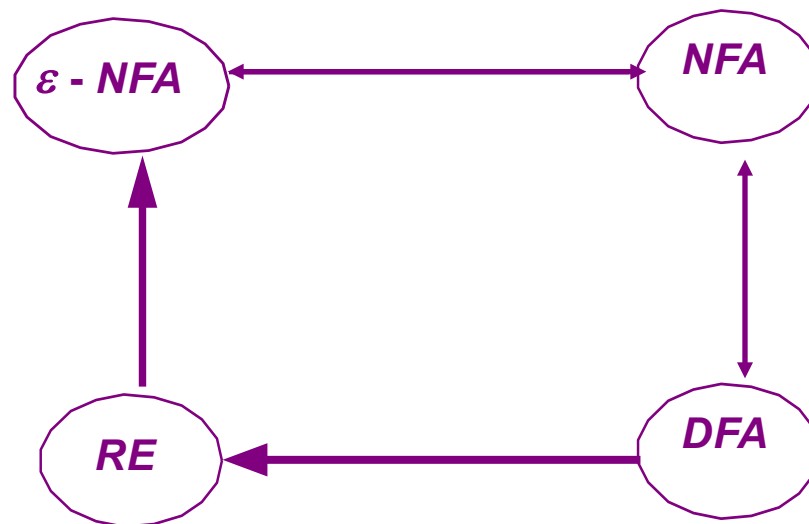
第四讲

- ✧ 正则表达式与有限自动机的关系
- ✧ 右线性语言与有限自动机的关系
- ✧ 右线性语言的性质(part1)

3.7 正则表达式与有限自动机的关系

结论：有限自动机、右（左）线性文法、正则表达式都定义了同一种语言-- 正则语言。

证明策略



RE(Regular Expression) --- 正则表达式



从有限自动机构造等价的正则表达式

(状态消去法)

思路:

(1) 扩展自动机的概念, 允许正则表达式作为转移弧的标记. 这样, 就有可能在消去某一中间状态时, 保证自动机能够接受的字符串集合保持不变.

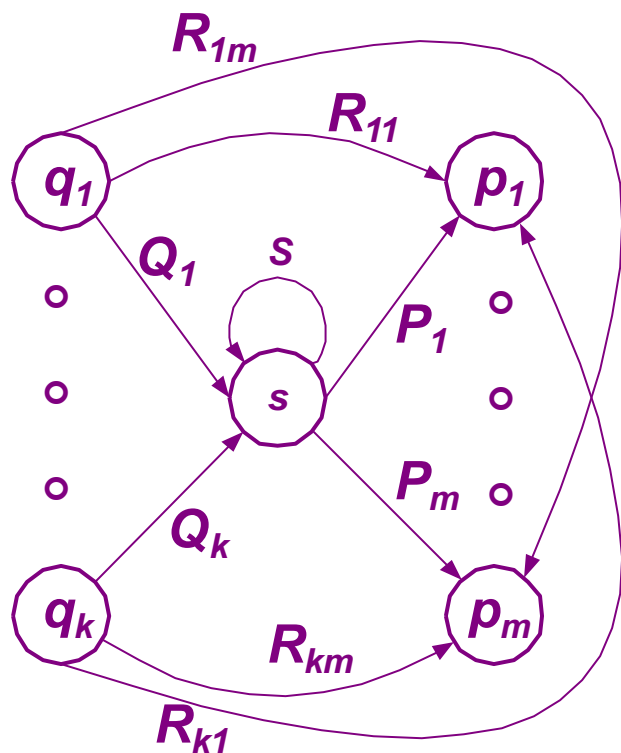
(2) 在消去某一中间状态时, 与其相关的转移弧也将同时消去, 所造成的影响将通过修改从每一个前趋状态到每一个后继状态的转移弧标记来弥补.

以下分别介绍中间状态的消去与正则表达式构造过程.

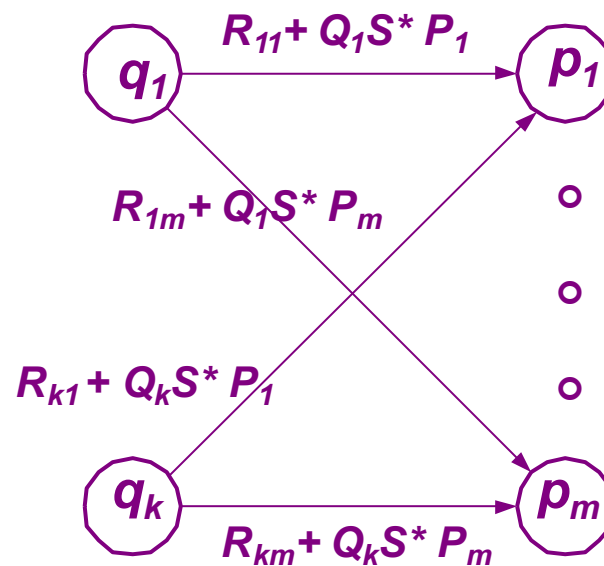
从有限自动机构造等价的正则表达式 (中间状态的消去)



以有限自动机构造等价的正则表达式 (中间状态的消去)



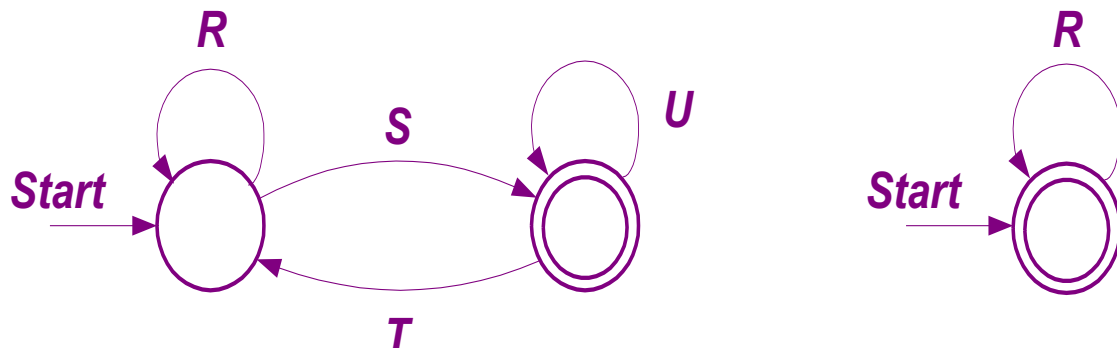
消去 s



从有限自动机构造等价的正则表达式 (状态消去法)

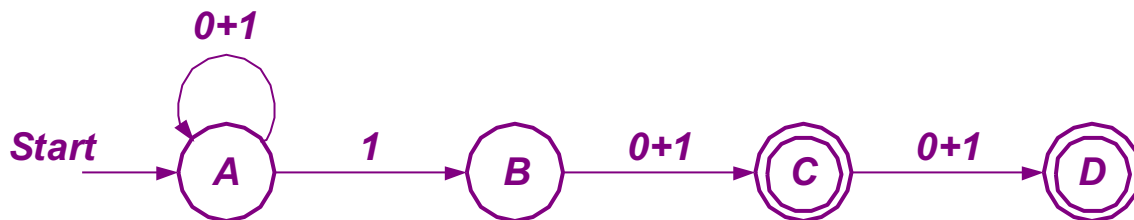
步骤:

- (1) 对每一终态 q , 依次消去除 q 和初态 q_0 之外的其它状态;
- (2) 若 $q \neq q_0$, 最终可得到一般形式如下左图两状态自动机, 该自动机对应的正则表达式可表示为 $(R+SU^*T)^*SU^*$.
- (3) 若 $q = q_0$, 最终可得到如下右图的自动机, 它对应的正则表达式可以表示为 R^* .

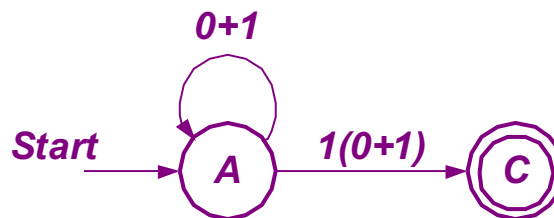
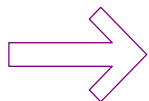


- (4) 最终的正则表达式为每一终态对应的正则表达式之和 (并) .

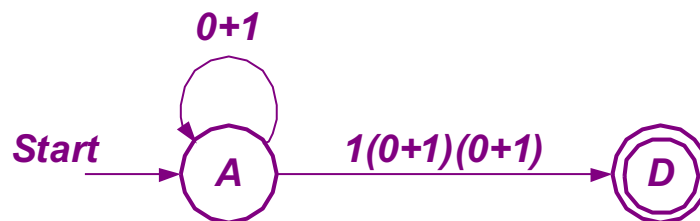
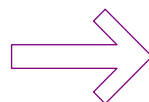
状态消去法举例



对于终态 C



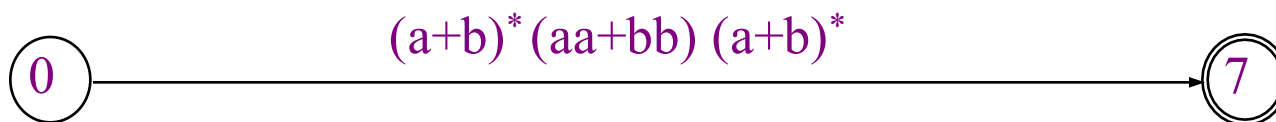
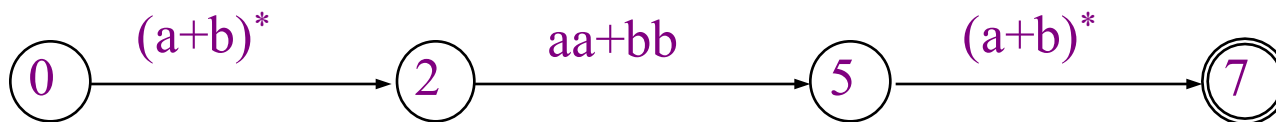
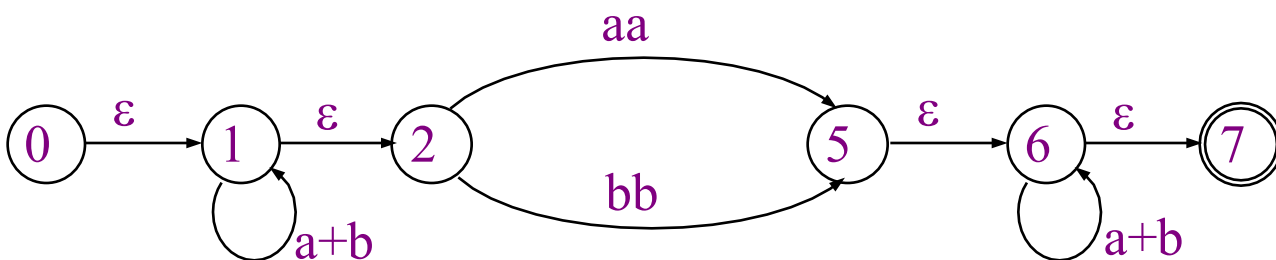
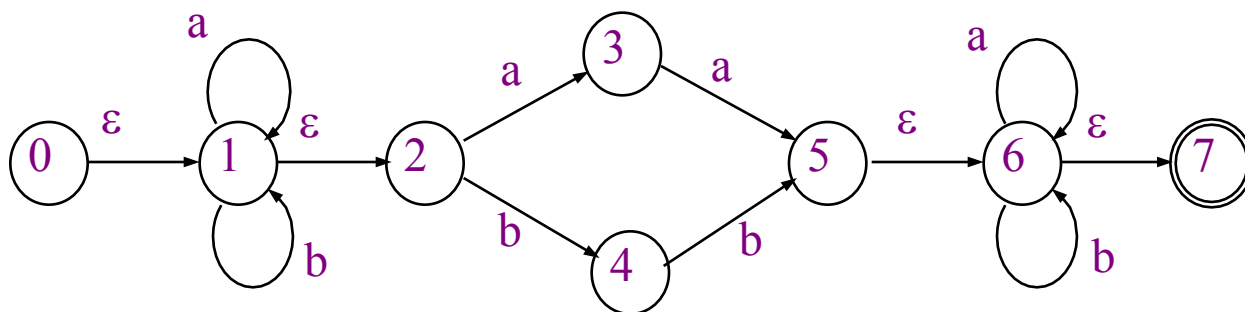
对于终态 D



等价的正则表达式

$$(0+1)^*1(0+1)+(0+1)^*1(0+1)(0+1)$$

状态消去法举例





从正则表达式构造等价的 ε - *NFA*

定理: L 是正则表达式 R 表示的语言, 则存在一个 ε - *NFA* E , 满足 $L(E) = L(R) = L$.

证明: 构造性证明. 可以通过结构归纳法证明从 R 可以构造出与其等价的, 满足如下条件的 ε - *NFA* :

- (1) 恰好一个终态;
- (2) 没有弧进入初态;
- (3) 没有弧离开终态;

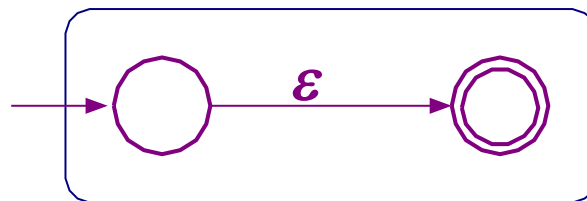


从正则表达式构造等价的 ε -NFA

(归纳构造过程)

基础:

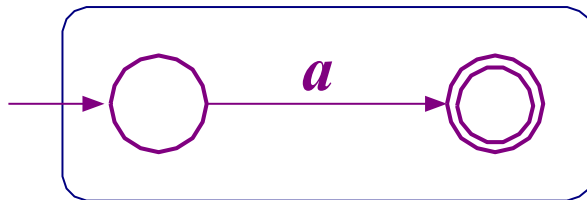
1 对于 ε , 构造为



2 对于 ϕ , 构造为



3 对于 a , 构造为

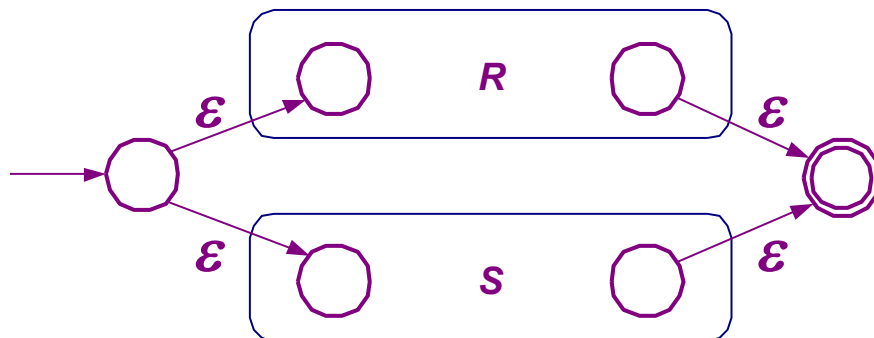


从正则表达式构造等价的 ε -NFA

(归纳构造过程)

归纳:

1 对于 $R+S$, 构造为

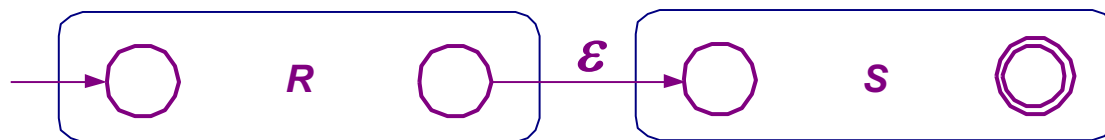


从正则表达式构造等价的 ε -NFA

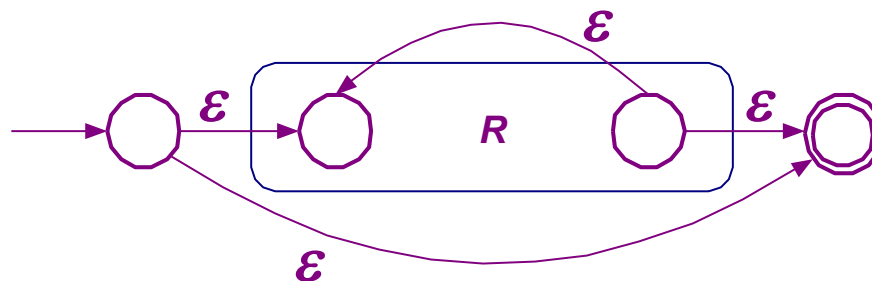
(归纳构造过程)

归纳:

2 对于 RS , 构造为

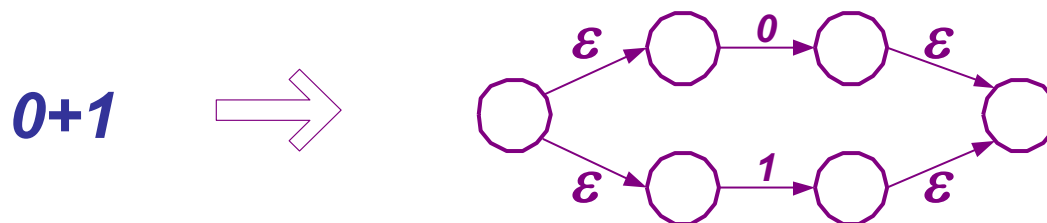
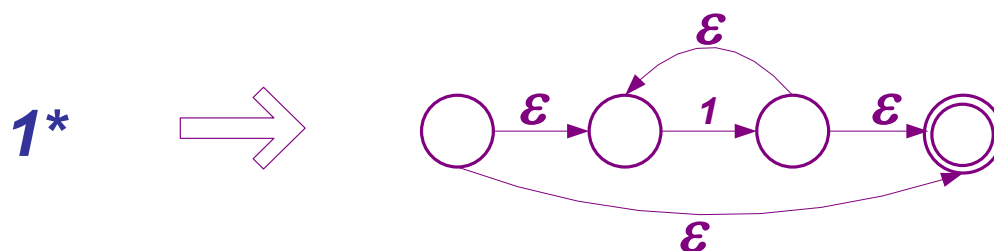


3 对于 R^* , 构造为



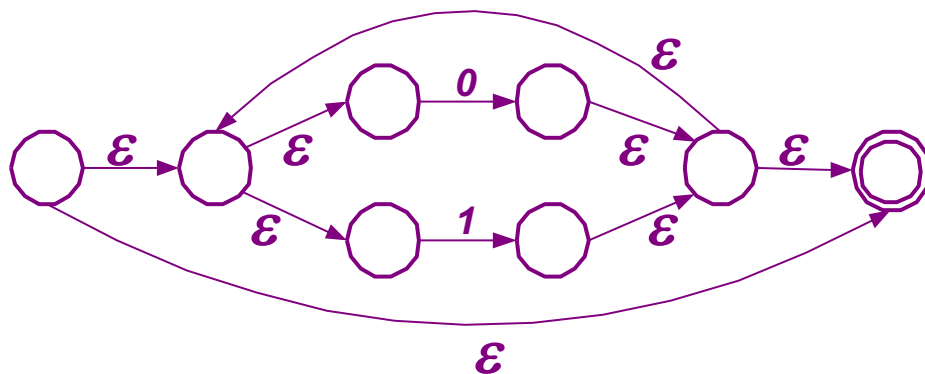
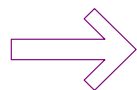
从正则表达式构造等价的 ε -NFA

举例：设正则表达式 $1^*0(0+1)^*$ ，构造等价的 ε -NFA.

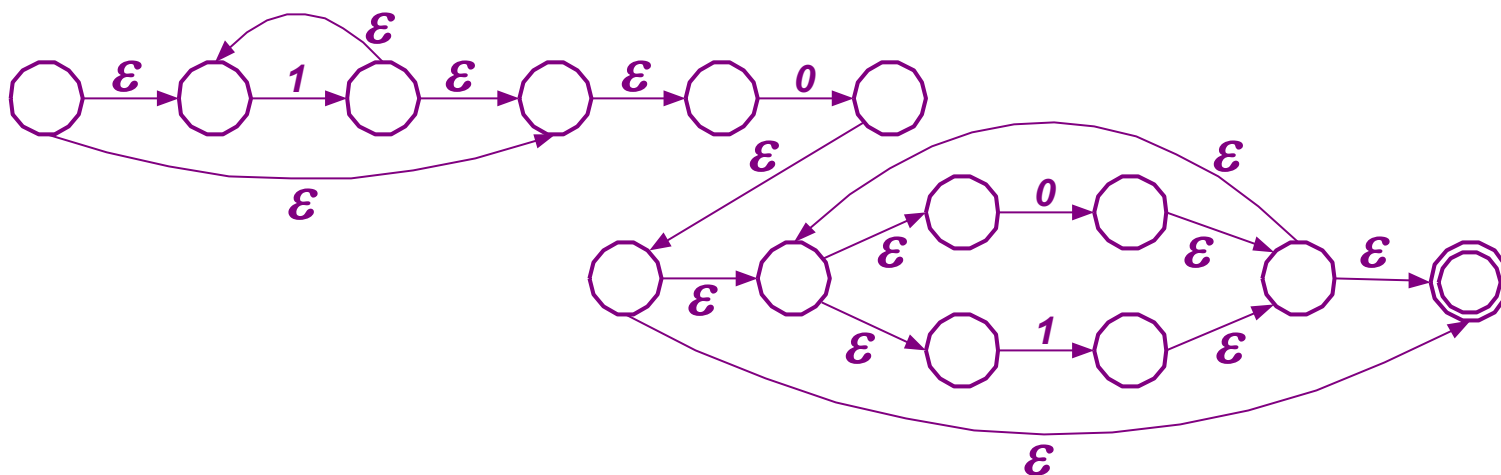


从正则表达式构造等价的 ε - NFA

$(0+1)^*$



$1^*0(0+1)^*$





3.8 右线性语言与有限自动机

至此，我们已学到正则集有三种定义方式，且这三种方式等价：

1. 正则集是含有 $\{\varepsilon\}$ ， φ ， $\{a\}$ 以及在并、连接和 $*$ 运算下封闭的语言
2. 由正规表达式定义的集合是正则集。
3. 由右线性文法生成的语言是正则集。

此外，还有第四种方式：

将正则集作为由有限自动机定义的集合。

即 正则集(右线性语言) \Leftrightarrow 有限自动机

右线性文法 \Rightarrow 有限自动机

定理3.8.1：由任意右线性文法G定义的语言必然能被一个NFA M所接受。即 $L(G) = L(M)$

证明思路（构造证明）：

设右线性文法 $G = (N, T, P, S)$ ，构造一个与G等价的有限自动机NFA $M = (Q, T, \delta, q_0, F)$ ，其中：

$Q = N \cup \{H\}$ ，H为一个新增加的状态， $H \notin N$ ， $q_0 = S$

$$F = \begin{cases} \{H, S\} & \text{当 } S \rightarrow \varepsilon \text{ 属于 } P。 \\ \{H\} & \text{否则} \end{cases}$$

δ 的定义为：

当 $A \rightarrow aB \in P$ ，则 $B \in \delta(A, a)$

当 $A \rightarrow a \in P$ ，则 $H \in \delta(A, a)$

对于任意输入， $\delta(H, a) = \varphi$ 。

右线性文法 \Rightarrow 有限自动机 (例)

例：设有右线性文法 $G=(\{S,B\}, \{a,b\}, P, S)$ ，其中

$P : S \rightarrow aB \quad B \rightarrow aB | bS | a$

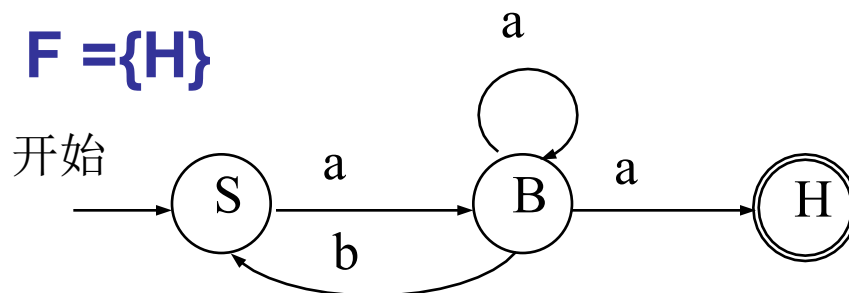
试构造与 G 等价的有限自动机 M 。

解：设 NFA $M=(Q, T, \delta, q_0, F)$

$Q=\{S,B,H\} \quad T=\{a,b\} \quad q_0 = S \quad F = \{H\}$

转换函数 δ ：

- 对于产生式 $S \rightarrow aB$ ，有 $\delta(S,a)=\{B\}$
- 对于产生式 $B \rightarrow aB$ ，有 $\delta(B,a)=\{B\}$
- 对于产生式 $B \rightarrow bS$ ，有 $\delta(B,b)=\{S\}$
- 对于产生式 $B \rightarrow a$ ，有 $\delta(B,a)=\{H\}$





右线性文法 \Rightarrow 有限自动机 (续)

求证 G 与 NFA M 两者定义了同一语言。

证明:

先证 (1) 文法 G 产生的语言 $L(G)$ 能够被 NFA M 所接收;

再证 (2) NFA M 接受的语言 $L(M)$ 可由文法 G 产生。

右线性文法 \Rightarrow 有限自动机 (续)

证明方法：通过两者定义的语言中任意一个字符串来说明。

(1) 设 $\omega = a_1a_2\dots a_n \in L(G)$ ，且 $n \geq 1$

$$\begin{aligned} \text{则有 } S &\Rightarrow a_1A_1 \Rightarrow a_1a_2A_2 \Rightarrow \dots \\ &\Rightarrow a_1a_2\dots a_{n-1}A_{n-1} \Rightarrow a_1a_2\dots a_{n-1}a_n \end{aligned}$$

则由 δ 的定义，有

$$A_1 \in \delta(S, a_1), A_2 \in \delta(A_1, a_2), \dots,$$

$$A_{n-1} \in \delta(A_{n-2}, a_{n-1}), H \in \delta(A_{n-1}, a_n), \text{ 且 } H \in \delta(S, \omega)$$

因为 $H \in F$ ，所以 ω 被NFA M 所接受。

又若 $\varepsilon \in L(G)$ ，则表明 $S \rightarrow \varepsilon \in P$ ，由 NFA M 的定义，有 $S \in F$ ，即 ε 也被NFA M 接受。

所以，由文法 G 派生的任意字符串 $\omega \in L(M)$ 。 #

右线性文法 \Rightarrow 有限自动机 (续)

(2) 再证 $L(M)$ 可由 G 产生

设 $\omega = a_1a_2\dots a_n$ 被 NFA M 接受, 即 $\omega \in L(M)$,

则必然存在状态序列 $S, A_1, A_2, \dots, A_{n-1}, H$

对 M 有转换函数为

$$A_1 \in \delta(S, a_1), A_2 \in \delta(A_1, a_2), \dots,$$

$$A_{n-1} \in \delta(A_{n-2}, a_{n-1}), H \in \delta(A_{n-1}, a_n)$$

则可规定 G 中含有产生式

$$S \Rightarrow a_1A_1, A_1 \Rightarrow a_2A_2, \dots, A_{n-1} \Rightarrow a_n$$

于是存在推导

$$S \Rightarrow a_1A_1 \Rightarrow a_1a_2A_2 \Rightarrow \dots \Rightarrow a_1a_2\dots a_{n-1}A_{n-1} \Rightarrow a_1a_2\dots a_{n-1}a_n$$

即 $a_1a_2\dots a_n$ 是文法 G 的一个句子。

也即 $\omega \in L(G)$ 。



课堂练习：

练习： 设线性文法 $G = (\{S, A, B\}, \{a, b\}, P, S)$

P: $S \rightarrow aA \mid a$

$A \rightarrow aA \mid aS \mid bB$

$B \rightarrow bB \mid b \mid a$

构造相应的 NFA M。

有限自动机 \Rightarrow 右线性文法

定理3.8.2：设有限自动机 M 接受的语言为 $L(M)$
则存在右线性文法 G ，它产生的语言 $L(G) = L(M)$ 。

证明思路：

构造一个右线性文法 G ，使它接受由 **NFA** M 定义的语言。

构造方法：

设 $M = (Q, T, \delta, q_0, F)$ ，构造一个右线性文法
 $G = (N, T, P, S)$ ，其中 $N = Q$ ， $S = q_0$

P 定义为：

若 $\delta(A, a) = B$ 且 $B \notin F$ ，则 $A \rightarrow aB$ 在 **P** 中

若 $\delta(A, a) = B$ 且 $B \in F$ ，则 $A \rightarrow a$ 和 $A \rightarrow aB$ 在 **P** 中

$L(M) \Leftrightarrow L(G)$ 的证明见书 P91（自学）。

有限自动机 \Rightarrow 右线性文法 (例)

例：设有DFA $M = (\{q_0, q_1, q_2, q_3\}, \{a, b\}, \delta, q_0, \{q_3\})$

其中转换函数如图所示，

试构造与之等价的右线性文法 G 。

解：构造右线性文法 $G = (N, T, P, S)$

$N = \{q_0, q_1, q_2, q_3\}$ $T = \{a, b\}$ $S = q_0$

产生式集合 P

$\delta(q_0, a) = q_1, \quad \therefore q_0 \rightarrow aq_1$

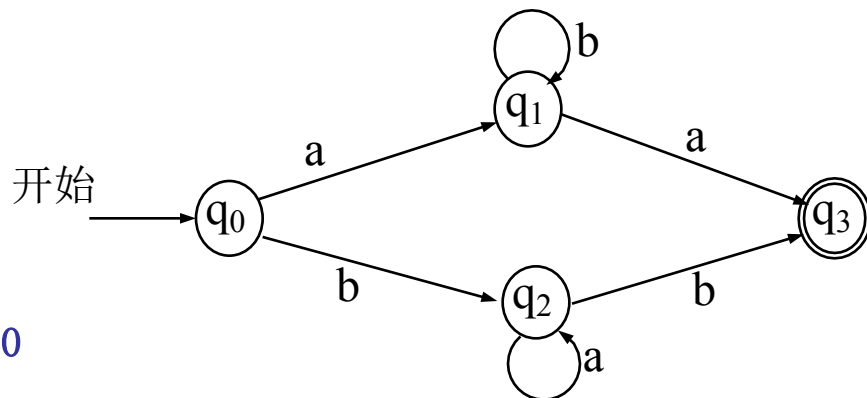
$\delta(q_0, b) = q_2, \quad \therefore q_0 \rightarrow bq_2$

$\delta(q_1, a) = q_3, \quad q_3 \in F, \quad \therefore q_1 \rightarrow a \mid aq_3$

$\delta(q_1, b) = q_1, \quad \therefore q_1 \rightarrow bq_1$

$\delta(q_2, a) = q_2, \quad \therefore q_2 \rightarrow aq_2$

$\delta(q_2, b) = q_3, \quad q_3 \in F, \quad \therefore q_2 \rightarrow b \mid bq_3$



构造的文法 G (化简 q_3) :

$G = (\{q_0, q_1, q_2\}, \{a, b\}, P, q_0)$

$P : \quad q_0 \rightarrow aq_1 \mid bq_2$

$q_1 \rightarrow a \mid bq_1$

$q_2 \rightarrow aq_2 \mid b$



3.9 右线性语言的性质

主要内容:

- DFA的极小化
- 泵浦引理
- 右线性语言的封闭性



确定有限自动机DFA的化简(极小化)

对DFA M 的极小化是找出一个状态数比 M 少的DFA $M1$, 使满足 $L(M) = L(M1)$

1. 等价和可区分的概念

设DFA $M = (Q, T, \delta, q_0, F)$

对不同的状态 $q_1, q_2 \in Q$ 和每个 $\omega \in T^*$,
如果有
 $(q_1, \omega) \vdash^* (q, \varepsilon)$ 必有 $(q_2, \omega) \vdash^* (q, \varepsilon)$ 且 $q \in F$,
则称 q_1 与 q_2 状态等价. 记为 $q_1 \equiv q_2$
否则, 称 q_1, q_2 可区分.



确定有限自动机DFA的化简

2. 不可达状态

如果不存在任何 $\omega \in T^*$, 使 $(q_0, \omega) \vdash^* (q, \varepsilon)$,
则称状态 $q \in Q$ 为不可达状态.

3. 最小化

若DFA M 不存在互为等价状态及不可达状态, 则称
DFA M 是最小化的.

一个DFA M 的最小化，是把 M 的状态集 Q 构成一个划分。

即：任何两个子集的状态都是可区分的；同一子集中的任何两个状态都是等价的。之后，每个子集用一个状态代表，并取一个状态名。

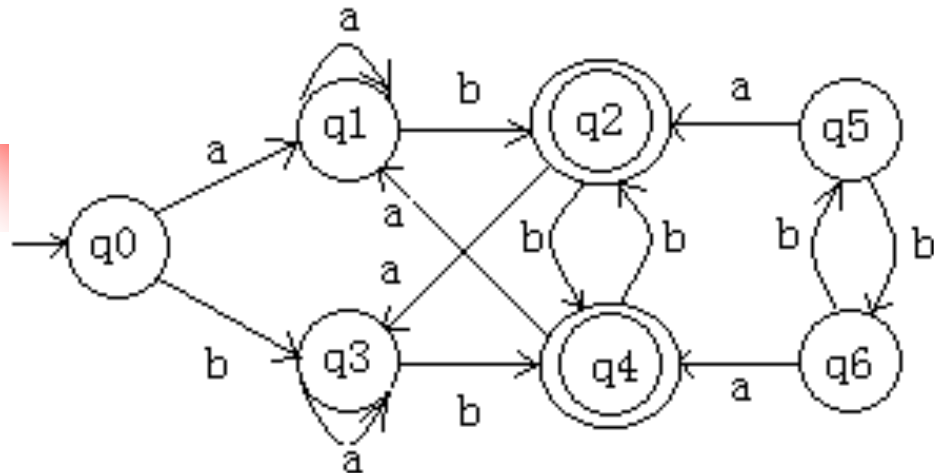
构成划分的步骤：

1. 构成基本划分 $\Pi = \{\Pi', \Pi''\}$, (Π' 为终态集, Π'' 为非终态集)
2. 细分 $\Pi = \{\Pi^1, \Pi^2, \dots, \Pi^n\}$,
 $\Pi^i = \{q_1, q_2, \dots, q_m\}$

当输入任意字符 a 时，若 Π^i 中的状态经标 a 的边可到达的状态集的元素分属于两个不同的子集中，则将 Π^i 细分为两个子集。

重复步骤(2)，直至不可再细分，得到 $M1$ 。

若 $M1$ 中有不可达状态，将其删除， $M1$ 便是最小化的。



例

(1) q_5, q_6 为不可达状态, 删除之.

(2) $Q = \{q_0, q_1, q_2, q_3, q_4\}$,

$\Pi = \{\{q_2, q_4\}, \{q_0, q_1, q_3\}\}$

构成基本划分 $\Pi = \{\Pi', \Pi''\}$

(a) 对于 $\Pi' = \{q_2, q_4\}$,

对字符 a , 有 $\delta(q_2, a) = q_3, \delta(q_4, a) = q_1$ $q_1, q_3 \in$ 同一子集.

对字符 b , 有 $\delta(q_2, b) = q_4, \delta(q_4, b) = q_2$ $q_4, q_2 \in$ 同一子集.

$\therefore \Pi' = \{q_2, q_4\}$ 不能再细分. 可用 q_2 表示 Π' 状态.

(b) 对于 $\Pi'' = \{q_0, q_1, q_3\}$

对 a , $\delta(q_0, a) = q_1, \delta(q_1, a) = q_1, \delta(q_3, a) = q_3$ $q_1, q_3 \in$ 同一子集

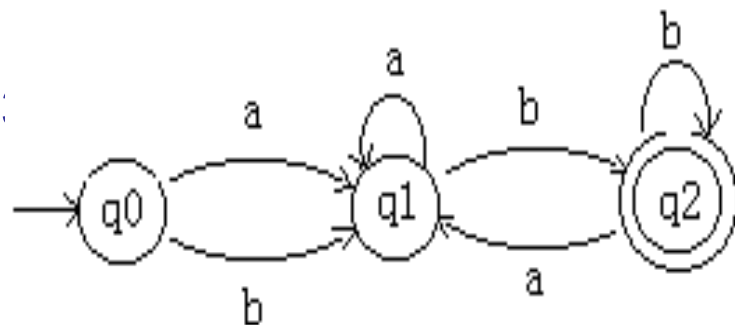
对 b , $\delta(q_0, b) = q_3, \delta(q_1, b) = q_2, \delta(q_3, b) = q_4$

$q_3, q_2, q_4 \notin$ 同一子集.

\therefore 将 Π'' 再分解. $\Pi'' = \{\{q_0\}, \{q_1, q_3\}\}$

q_1 表示

$\therefore Q = \{\{q_0\}, \{q_1\}, \{q_2\}\}$



计算状态集划分的算法——填表法

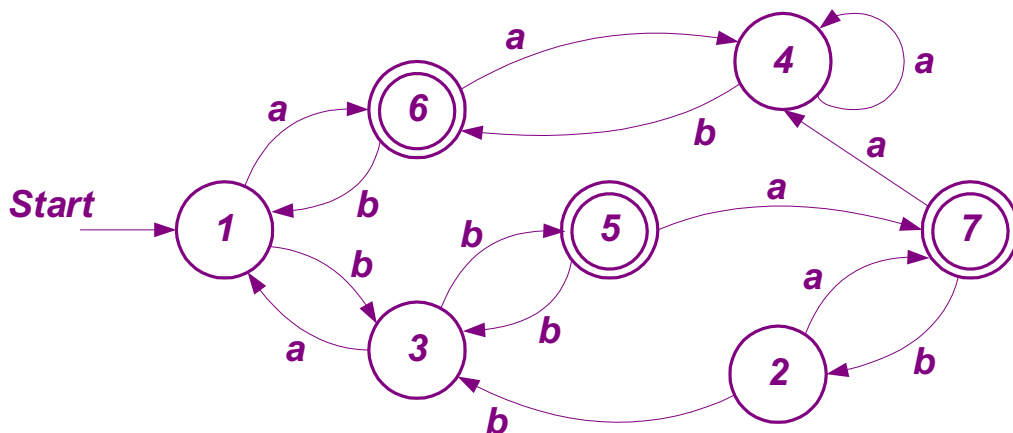
✧ 填表算法 (*table-filling algorithm*) 基于如下递归地标记可区分的状态偶对的过程:

- 基础 如果 p 为终态, 而 q 为非终态, 则 p 和 q 标记为可区分的;
- 归纳 设 p 和 q 已标记为可区分的, 如果状态 r 和 s 通过某个输入符号 a 可分别转移到 p 和 q , 即 $\delta(r, a) = p$, $\delta(s, a) = q$, 则 r 和 s 也标记为可区分的;
这是因为: 若 p 和 q 可为字符串 w 区分, 则 r 和 s 可为字符串 aw 区分.
($\because \delta'(r, aw) = \delta'(p, w)$, $\delta'(s, aw) = \delta'(q, w)$)

计算状态集划分的算法——填表法

✧ 填表算法举例

2						
3	X	X				
4	X	X	X			
5	X	X	X	X		
6	X	X	X	X	X	
7	X	X	X	X	X	
	1	2	3	4	5	6



(1) 区分所有终态和非终态

(2) 区分 (1,3), (1,4), (2,3),
(2,4), (5,6), (5,7)

(3) 区分 (3,4)

(4) 结束. 划分结果: $\{1,2\}, \{3\}, \{4\}, \{5\}, \{6,7\}$

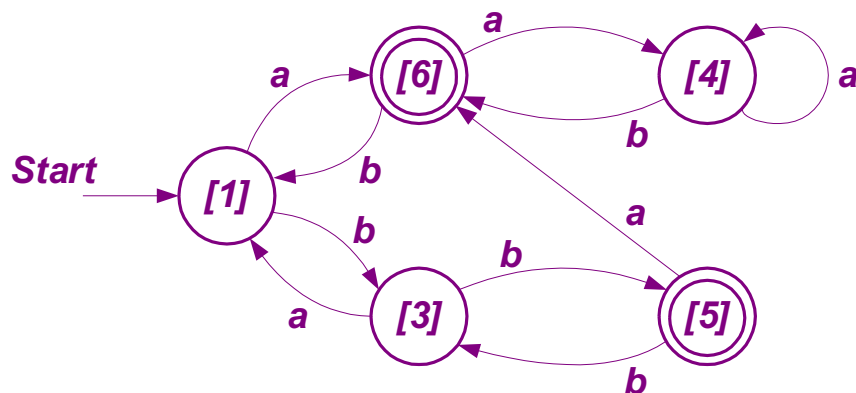
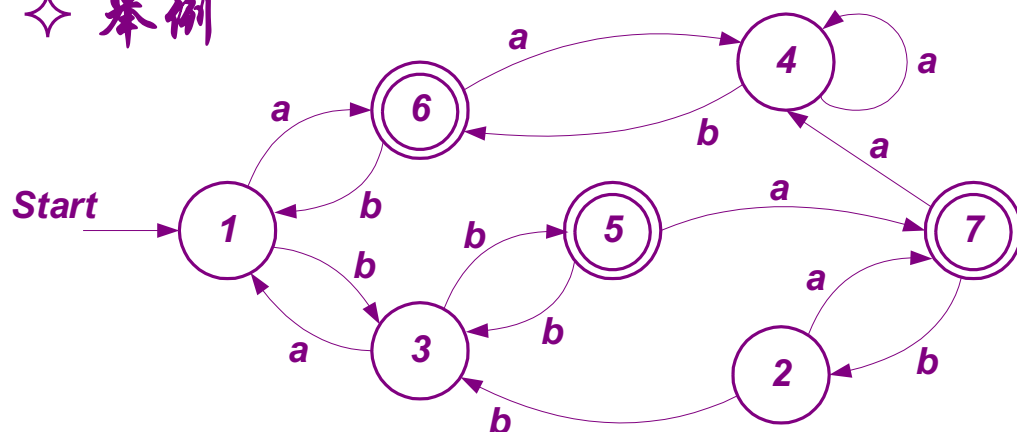
通过合并等价的状态进行 DFA 的优化

✧ 步骤

1. 删除所有从开始状态不可到达的状态及与其相关的边, 设所得到的 DFA 为 $A = (Q, T, \delta, q_0, F)$;
2. 使用填表算法找出所有等价的状态偶对;
3. 根据 2 的结果计算当前状态集合的划分块, 每一划分块中的状态相互之间等价, 而不同划分块中的状态之间都是可区分的. 包含状态 q 的划分块用 $[q]$ 表示.
4. 构造与 A 等价的 DFA $B = (Q_B, T, \delta_B, [q_0], F_B)$, 其中 $Q_B = \{ [q] \mid q \in Q \}$, $F_B = \{ [q] \mid q \in F \}$, $\delta_B([q], a) = [\delta(q, a)]$

通过合并等价的状态进行 DFA 的优化

✧ 举例



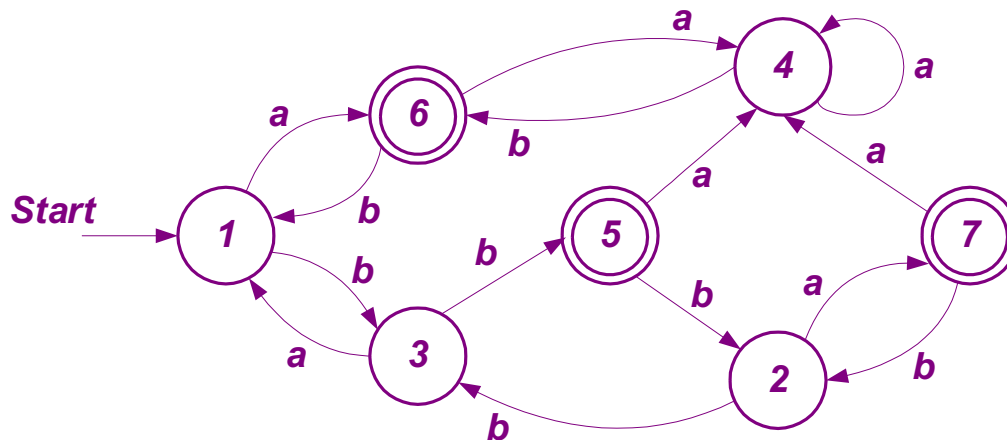
— 等价的状态偶对为：
 $(1, 2)$, $(6, 7)$

— 划分结果：
 $\{1, 2\}$, $\{3\}$, $\{4\}$,
 $\{5\}$, $\{6, 7\}$

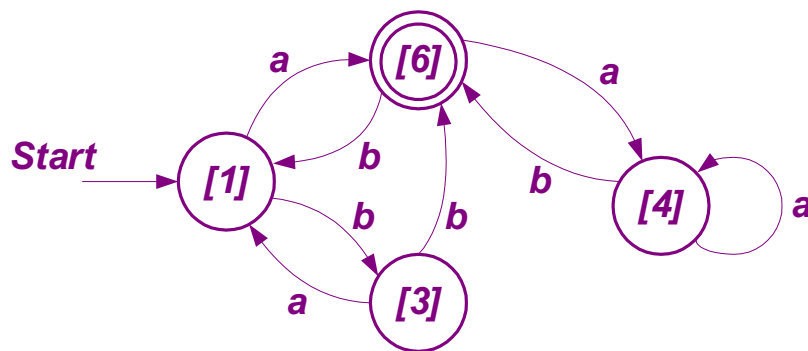
— 新的状态集合：
 $[1]$, $[3]$, $[4]$, $[5]$, $[6]$

最小化的 DFA

✧ 课堂练习 最小化下列 DFA:



✧ 参考结果



针对正则语言的 Pumping 引理

- 正则语言应满足的一个必要条件
- 用于判定给定的语言不是正则集。

物理意义：当给定一个正则集和该集合上一个足够长的字符串时，在该字符串中能找到非空的子串，并使子串重复，从而组成新的字符串。该新串必在同一个正则集内。

定理：

设 L 是正则集，存在常数 n ，对字符串 $\omega \in L$ 且 $|\omega| \geq n$ ，则 ω 可写成 $\omega_1\omega_0\omega_2$ ，其中 $|\omega_1\omega_0| \leq n$ ， $|\omega_0| > 0$ ，对所有的 $i \geq 0$ 有 $\omega_1\omega_0^i\omega_2 \in L$ 。

证明 设 L 是 DFA $D = (Q, T, \delta, q_0, F)$ 的语言，取 $n = |Q|$ 即可。

□

DFA 的 “Pumping” 特性

设 DFA $D = (Q, T, \delta, q_0, F)$, $|Q|=n$.

对于任一长度不小于 n 的字符串 $w = a_1a_2...a_m$, 其中 $m \geq n$, $a_k \in T$ ($1 \leq k \leq m$), $q \in Q$, 考察如下状态序列

$p_0 = q$
 $p_1 = \delta'(q, a_1)$
 $p_2 = \delta'(q, a_1a_2)$
...
 $p_n = \delta'(q, a_1a_2...a_n)$
 $p_{n+1} = \delta'(q, a_1a_2...a_{n+1})$
...
 $p_m = \delta'(q, a_1a_2...a_m)$

由 *pigeonhole* 原理, $p_0, p_1, p_2, \dots, p_n$ 中至少有两个状态是重复的, 即存在 i, j , $0 \leq i < j \leq n$, $p_i = p_j$.

✧ “pumping” 特性:
任一长度不小于状态数目的字符串所标记的路径上, 必然出现重复的状态.

DFA 的 “Pumping” 特性

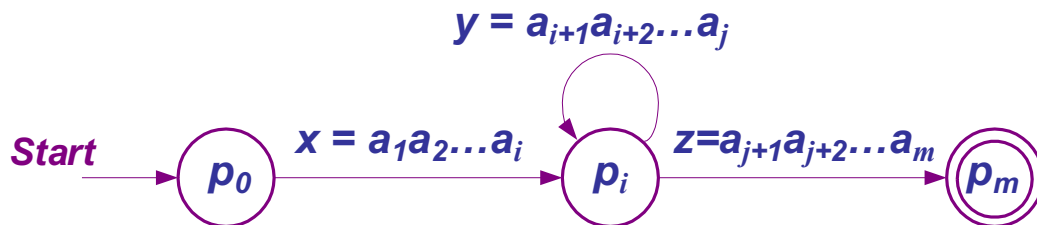
✧ “pumping” 特性：如前，设 DFA $D = (Q, T, \delta, q_0, F)$, $|Q|=n$, $w = a_1a_2\dots a_m$ ($m \geq n$), 则存在 i, j , $0 \leq i < j \leq n$, $p_i = p_j$, 其中 $p_k = \delta'(p_0, a_1a_2\dots a_k)$, $0 \leq k \leq m$.

✧ 若假定 $p_0 = q_0$, $p_m \in F$, 即 $w \in L(D)$.

令 $w = xyz$, 其中:

$$x = a_1a_2\dots a_i, y = a_{i+1}a_{i+2}\dots a_j, z = a_{j+1}a_{j+2}\dots a_m$$

则对任何 $k \geq 0$, 都有 $xy^kz \in L(D)$. (参考下图)



Pumping 引理的应用

(用于证明某个语言 L 不是正规语言)

◇ 证明步骤

1. 选任意的 n .
2. 找到一个满足以下条件的串 $w \in L$ (长度至少为 n).
3. 任选满足 $w = xyz \wedge y \neq \varepsilon \wedge |xy| \leq n$ 的 x, y, z
4. 找到一个 $k \geq 0$, 使 $xy^kz \notin L$.

◇ 举例 证明 $L = \{ a^n b^n \mid n \geq 1 \}$ 不是正则集.

◇ 证明: 由泵浦引理, 假设 L 是正则集, 则对足够大的 n , $a^n b^n$ 可写成 $\omega_1 \omega_0 \omega_2$, 其中 $0 < |\omega_0| \leq n$, $|\omega| = 2n > n$

若 $\omega_0 = a^+ \text{ 或 } b^+$, 设 $|\omega_0| = k \geq 1$, k 为常数,

取 $i=0$, 有 $\omega_1 \omega_0^0 \omega_2 = \omega_1 \omega_2 = a^{n-k} b^n \text{ 或 } a^n b^{n-k}$, 此时, a, b 字符个数不同, 即新组成的串 $\omega_1 \omega_2 \notin L$.

若 $\omega_0 = a^+ b^+$, 可取 $i=2$, 有 $\omega_1 \omega_0 \omega_0 \omega_2 = \omega_1 a^+ b^+ a^+ b^+ \omega_2 \notin L$

\therefore 与假设矛盾, 故 L 不是正则集.

Pumping 引理的应用

例 证明 $L = \{ a^{k^2} \mid k \geq 1 \text{ 的整数} \}$ 不是正则集。

证明 假设 L 是正则集，取足够大的整数 n ， $w = a^{n^2}$ 。

$$\text{有 } |w| = |\omega_1 \omega_0 \omega_2| = n^2 \geq n,$$

□

$$0 < |\omega_0| \leq n, \quad 0 < |\omega_1 \omega_0| \leq n$$

$$\text{取 } i=2, \text{ 有 } n^2 < |\omega_1 \omega_0^2 \omega_2| \leq n^2 + n < (n+1)^2$$

$$\therefore |\omega_1 \omega_0^i \omega_2| \notin L \quad (\text{串长不是整数的平方})$$

与假设矛盾。

$\therefore L$ 不是正则集。



课后练习

转换下列正则表达式为带 ϵ 转移的NFA.

a) 01^* .

b) $(0+1)01$.

c) $00(0+1)^*$.

Chap3 习题 7, 9, 20