
编译原理

- 第一章 编译程序概述
 - 第二章 PL/0编译程序的实现
 - 第三章 文法和语言（回顾）
 - 第四章 词法分析
 - 第五章 自顶向下语法分析方法
 - 第六章 自底向上优先分析方法
 - 第七章 LR分析方法
 - 第八章 语法制导翻译和中间代码生成
 - 第九章 符号表
 - 第一〇章 代码优化
 - 第十一章 代码生成
-

文法和语言

1. 概念：文法、语言、推导、句型、句子、短语、句柄的定义。
 2. 文法作为程序语言的语法的描述工具,它用规则只能陈述的是：语言的所有句子以什么样的符号串能出现.请记住文法和语言的形式定义中的“形式”的含义—只涉及语言的语法不涉及语言的语义。
 3. 符号串集合的表示法、结构及其特性，是程序设计语言语法分析研究的基础。
-

文法的概念

语言

语法：是一组规则，定义符号如何排列，**排列与符号含义无关。**

语义：研究语法的含义 { 静态语义
动态语义

文法是阐述语法的一个工具

一、文法的概念 (写出以下语言的文法)

“你是大学生” 对

“我是教师” 对

“我大学生是” 错

“我学习大学生” 对

符号和符号串

为给出语言的形式化定义，先讨论一些有关概念：

- 1、**字母表**—符号集：是字母的有穷**非空**集合。
 - 2、**符号串**—字母表的符号组成的任何有穷序列。
 - 3、**符号串的长度**：符号串 x 有 m 个符号，则长度就为 m ，表示 $|x|=m$
 - 4、**空符号串**：用 ε 表示，长度为0（不含任何符号）
-

5、符号串的运算：

(1) 符号串的头和尾

若 $z=xy$ ，则 x 是 z 的头， y 是 z 的尾。

(2) 符号串的固有头和固有尾

若 $z=xy$ 符号串， x 非空，则 y 是固有尾；

若 y 非空，则 x 是固有头。

(3) 符号串的连接：

(4) 符号串的方幂：

- (6) 闭包 (Σ^*)

- 字母表 Σ ，用 Σ^* 表示 Σ 上所有有穷长的串集合， Σ^* 称为 Σ 的闭包。

- 例：字母表 $\Sigma=\{0,1\}$

则 $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots \cup \Sigma^n$
 $= \{\epsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, \dots\}$

- (7) 正闭包(Σ^+)

$$\Sigma^+ = \Sigma^1 \cup \Sigma^2 \cup \dots \cup \Sigma^n$$

Σ^+ 称 Σ 的正闭包。

显然： $\Sigma^* = \Sigma^0 \cup \Sigma^+$

$$\Sigma^+ = \Sigma \Sigma^* = \Sigma^* \Sigma$$

三、文法和语言的形式定义

(用以上术语对文法的概念进行形式化)

1、规则 (重写规则、产生式、生成式)

形如 $\alpha \rightarrow \beta$ 或 $\alpha ::= \beta$

α 称规则的左部, β 称规则的右部。

2、文法的定义

定义为四元组 (V_N, V_T, P, S)

3、直接推导的定义

<1> $\alpha \rightarrow \beta$ 是文法 $G = (V_N, V_T, P, S)$ 的规则， γ 和 δ 是 V^* 中的任意符号，若有符号串 V, W 满足：

$$V = \gamma \alpha \delta, \quad W = \gamma \beta \delta$$

则说 V 是直接产生 W

或 W 是 V 的直接推导

或 W 直接规约到 V

记作 $V \Rightarrow W$

(2) 若存在直接推导的序列：

$$V = W_0 \Rightarrow W_1 \Rightarrow W_2 \Rightarrow \dots \Rightarrow W_n = W \quad (n > 0)$$

则称 V 推导 W (或 W 规约到 V)，记 $V \Rightarrow W$

4、句型的定义：

设 $G[S]$ 是文法，如果符号串 x 是从识别符号（开始符号）推导出来的（即 $S \xRightarrow{*} x$ ）则称 x 是文法 $G[S]$ 的**句型**。

若 x 仅由终结符号组成（ $S \xRightarrow{*} x, x \in V_T^*$ ）

则称 x 为 $G[S]$ 的**句子**。

例： $S, 0S1, 000111$ 都是文法 G 的句型， 000111 是 G 的句子。

【结论】 **句子一定是句型，句型不一定是句子。**

5、语言的定义： **$L(G)$** 表示

文法G产生的语言定义为：G产生的句子的集合

$$\{x \mid S \xRightarrow{*} x, \quad S \text{ 为文法开始符号, } x \in V_T^* \}$$

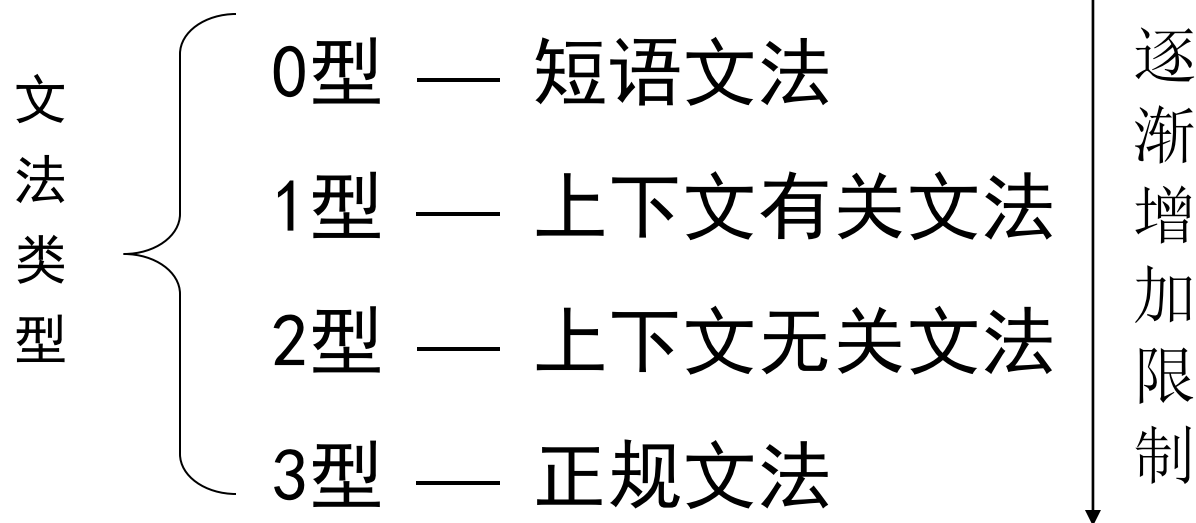
该集合为语言，用 $L(G)$ 表示。

从定义可知： x 是句型且 x 是文法G的句子。

想：例1的语言表示为什么？

四、文法的类型：

语言学家Chomsky把文法分成以下四种类型：



如果文法是正规文法 \Rightarrow 一定也是上下文无关文法

设 $G = (V_N, V_T, P, S)$ ，如果它的每一个产生式 $\alpha \rightarrow \beta$ 满足：

$\alpha \in (V_N \cup V_T)^*$ 且至少包含一个非终结符，

$\beta \in (V_N \cup V_T)^*$ ，则 G 是**0型文法**。

1型文法又称为上下文有关文法，它的每一个产生式也可描述为： $\alpha_1 A \alpha_2 \rightarrow \alpha_1 \beta \alpha_2$ ，其中 $\alpha_1 \alpha_2 \beta$ 都在 $(V_N \cup V_T)^*$ 中， $\beta \neq \varepsilon$ ， A 在 V_N 中。只有 A 出现在 $\alpha_1 \alpha_2$ 的上下文中，才允许用 β 取代 A 。

1、上下文无关文法:

设 $G = (V_N, V_T, P, S)$ ，若 P 中的每一个产生式 $\alpha \rightarrow \beta$ 满足:

α 是一非终结符, $\beta \in (V_N \cup V_T)^*$, 则此文法称为上下文无关文法 (2型文法)。(用 β 取代 α 与 α 所在的上下文无关)

例6: $G = (\{S, A, B\}, \{a, b\}, P, S)$, P 的产生式如下:

$S \rightarrow aB$

$B \rightarrow b$

$S \rightarrow bA$

$B \rightarrow bS$

$A \rightarrow bAA$

$B \rightarrow aBB$

$A \rightarrow a$

$A \rightarrow aS$

上下文无关文法

该语法可以写成:

$$S \rightarrow aB \mid bA$$

$$A \rightarrow a \mid aS \mid bAA$$

$$B \rightarrow b \mid Bs \mid aBB$$

2、正规文法:

设 $G = (V_N, V_T, P, S)$, 若 P 中的每一个产生式都是 $A \rightarrow aB$ 或 $A \rightarrow a$, A, B 都是非终结符, a 是终结符, 则 G 是正规文法。

五、上下文无关文法及其语法树

描述上下文无关文法的句型推导的直观方法。

1、语法树（推导树）的定义：

给定文法 $G = (V_N, V_T, P, S)$ ，对于 G 的任何句型都能构造与之关联的语法树，须满足条件为：

- 每个结点都有一个标记，此标记是 V 的一个符号。
- 根的标记是 S 。
- 若一结点 n 至少有一个它自己除外的子孙，并且有标记 A ，则 A 肯定在 V_N 中。
- 如果结点 n 的直接子孙，从左到右的次序是结点 n_1, n_2, \dots, n_k ，其标记分别为 A_1, A_2, \dots, A_k ，那么 $A \rightarrow A_1 A_2 \dots A_k$ 一定是 P 中的一个产生式。

► 语法树的理解：

语法树表示了**在推导过程中用到什么样的产生式和用到哪些非终结符**，并没有表明采用的顺序。

只表示推导的结果，不表示推导的过程。

2、最左推导（或最右推导）：

若在推导中的任何一步 $\alpha \Rightarrow \beta$ （ α 、 β 是句型），都是对 α 中的最左（最右）非终结符进行替换，则称为最左（最右）推导。

► 在形式语言中，最右推导常被称为规范推导。

► 由规范推导所得的句型称为规范句型。

3、文法的二义性的定义：

- ✓ 如果一个文法存在某个句子对应两棵不同的语法树，则说这个文法是二义的。
- ✓ 若一个文法中存在某个句子，它有两个不同的最左（最右）推导，则这个文法是二义的。

例10：文法 $G = (\{E\}, \{+, \times, i, (,)\}, P, E)$ ，其中 P 为：

$E \rightarrow i$

$E \rightarrow E + E$

$E \rightarrow E \times E$

$E \rightarrow (E)$

则句型 $i \times i + i$ 的推导：

六、句型的分析

❶ 句型分析：是识别一个符号串是否为某文法的句型，是整个推导的构造过程。（为一个符号串构造一个语法树/推导树）

即：识别输入符号串是否为语法上正确的程序的过程。

❷ 分析程序（识别程序）：是完成句型分析的程序。

分析算法分为

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）

自下而上分析法：从输入符号开始，逐步进行“规约”，直至文法的开始符号（归约）

4、句柄的定义：（句柄是可归约串的称呼）

令 G 是一文法， S 是文法的开始符号， $\alpha\beta\delta$ 是文法 G 的一个句型。
（为 $\alpha\beta\delta$ 确定可归约串）如果有 $S \xRightarrow{*} \alpha A \delta$ 且 $A \xRightarrow{\pm} \beta$ ，则称 β 是句型 $\alpha\beta\delta$ 相对于非终结符 A 的**短语**。

若有 $A \Rightarrow \beta$ ，则称 β 是句型 $\alpha\beta\delta$ 相对 $A \rightarrow \beta$ 的直接短语。一个句型的最左直接短语称为该句型的**句柄**。

句柄是自底向上句法分析中当前时刻需要规约的符号串。如果能够自动计算出当前的句柄，则可执行自动句法分析。

例：文法 $G[E]$:

$$E \rightarrow T \mid E+T$$
$$T \rightarrow F \mid T \times F$$
$$F \rightarrow (E) \mid i$$

求句型 $i \times i + i$ 的短语，直接短语和句柄。

文法实用中的一些说明

----化简文法

文法中**不含有**有害规则和多余规则

有害规则：形如 $U \rightarrow U$ 的产生式。会引起文法的**二义性**

多余规则：指文法中**任何句子的推导都不会用到的规则**

多余规则也可表示为：文法中**不含有不可到达和不可终止的非终结符**，分以下两种情况：

- 1) 文法中某些**非终结符不在任何规则的右部出现**，该非终结符称为**不可到达**。
- 2) 文法中某些**非终结符**，由它**不能推出终结符号串**，该非终结符称为**不可终止**。

编译原理

- 第一章 编译程序概述
 - 第二章 PL/0编译程序的实现
 - 第三章 文法和语言
 - 第四章 词法分析
 - 第五章 自顶向下语法分析方法
 - 第六章 自底向上优先分析方法
 - 第七章 LR分析方法
 - 第八章 语法制导翻译和中间代码生成
 - 第九章 符号表
 - 第一〇章 代码优化
 - 第十一章 代码生成
-

内容提要

- 词法分析程序的设计
 - 单词的描述工具
 - 有穷自动机
 - 正规式和有穷自动机的等价性
 - 正规文法和有穷自动机的等价性
 - 词法分析程序的自动构造工具
-

词法分析

词法分析是编译过程的第一个阶段。

主要任务：从左到右逐个字符地对源程序进行扫描，产生一个个单词序列，用以语法分析。

基本思路：

- 将单词符号的语法用有效的工具描述；
 - 基于该描述建立单词的识别机制；
 - 基于词法分析程序的自动构造原理，设计和实现词法分析程序
-

词法分析程序的设计

- 词法分析程序与语法分析程序的接口方式
 - 词法分析程序的输出
 - 词法分析工作分离的考虑
-

单词符号(Token)的类别

1. 关键字 (保留字) (Keyword)
 - begin, end, if, then, for, while ...
 2. 标识符 (Identifier)
 - 用来表示各种名字
 3. 常数 (Literal)
 - 256, 3.14159, 1E+3, true, "abc"
 4. 运算符 (Operator)
 - 如 +、-、*、/ 等
 5. 分界符 (Delimiter)
 - 如逗号, 分号, 冒号等
-

词法分析器的输出

- Token的基本输出格式：

<单词类别， 单词自身的属性值>

类别提供给语法分析程序使用；词法单元自身的属性值提供给语义分析程序使用。

- 单词类别可以用整数编码表示，具体的分类设计以方便语法分析程序使用为原则。
 - 标识符：1
 - 常数：2
 - 关键字：3
 - 运算符：4
 - 界符：5

词法分析示例-1

```
while i != j do
```

```
    if i > j then i := i - j else j := j - i
```

上述源程序经词法分析器的输出：

```
<3, “while” >  
<1, 指向i的符号表入口的指针 >  
<4, “!=” >  
<1, 指向j的符号表入口的指针 >  
<3, “do” >  
<3, “if” >  
<1, 指向i的符号表入口的指针 >  
.....  
<1, 指向j的符号表入口的指针 >
```

-
- 单词的描述工具
 - 正规文法
 - 正规式
 - 单词的识别机制
 - 确定有穷自动机
 - 不确定有穷自动机
 - 词法分析程序的自动构造原理
 - 正规式和有穷自动机的等价性
 - 词法分析程序的自动构造工具
-

4.2 单词的描述工具

某种程序设计语言中的所有单词构成一种语言，该语言的语法都能用正规文法表示。正规文法是描述单词的一种工具。

1、正规文法（回顾）

文法 $G=(V_N, V_T, P, S)$ ， P 中每一规则有 $A \rightarrow aB$ 或 $A \rightarrow a$ ， $A, B \in V_N$ ， $a \in V_T^*$ ，称 $G(S)$ 是正规文法。

由正规文法产生的语言称为正规集

$\langle \text{标识符} \rangle \rightarrow l | l \langle \text{字母数字} \rangle$

$\langle \text{字母数字} \rangle \rightarrow l | d | l \langle \text{字母数字} \rangle | d \langle \text{字母数字} \rangle$

$\langle \text{无符号整数} \rangle \rightarrow d | d \langle \text{无符号整数} \rangle$

l 表示 $a-z$ 中的任何英文字母， d 表示 $0-9$ 中的任何数字

2、正规式（正则表达式）

也是表示正规集的工具，是用以描述单词符号的工具。

①正规式与正规集的定义：

设字母表为 Σ ，辅助字母表 $\Sigma' = \{\phi, \varepsilon, |, \cdot, *, (,)\}$ ；

✓ ε 和 ϕ 都是 Σ 上的正规式，表示的正规集分别为 $\{\varepsilon\}$ 和 ϕ ；

✓任何 $a \in \Sigma$ 是 Σ 上的一个正规式，表示的正规集为 $\{a\}$ ；

✓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 $L(e_2)$ ，则 (e_1) ， $e_1 | e_2$ ， $e_1 \cdot e_2$ 和 e_1^* 也都是正规式，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 $(L(e_1))^*$ 。

✓ 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。

例： $\Sigma = \{a, b\}$ ， Σ 上的正规式和相应的正规集为：

- ε
- ϕ
- a
- b
- -----
- (a)
- (b)
- a|b
- ab
- a*
- b*
- -----
- (a|b)*
- a*|b*
- aba*
- (a|b)*(aa|bb)(a|b)*

一个正规式可以表示若干个符号串，
其正规集就是这些符号串的集合

- ϵ $\{\epsilon\}$
- ϕ ϕ
- a $\{a\}$
- b $\{b\}$

• -----

- (a) $\{a\}$ 一个正规式可以表示若干个符号串,
- (b) $\{b\}$ 其正规集就是这些符号串的集合
- $a|b$ $\{a, b\}$
- ab $\{ab\}$
- a^* $\{\epsilon, a, aa, aaa, aaaa, \dots\}$
- b^* $\{\epsilon, b, bb, bbb, bbbb, \dots\}$

• -----

- $(a|b)^*$ a 和 b 组成的所有串
- $a^*|b^*$ $\{\epsilon, a, aa, aaa, aaaa, \dots, b, bb, bbb, bbbb, \dots\}$
- aba^* 以 ab 开头后接若干个(包括0个) a 组成的串
- $(a|b)^*(aa|bb)$ $(a|b)^* \Sigma^*$ 上所有含有两个相继的 a 或两个相继的 b 组成的串

例：令 $\Sigma = \{d, ., e, +, -\}$ ，则 Σ 上的正规式 $d^*(.dd^*|\varepsilon)(e(+|-|\varepsilon)dd^*|\varepsilon)$ 表示的是所有无符号数。

其中 d 为 $0\sim 9$ 中的数字。

比如以下都是正规式表示集合中的元素：

2

12.59

3.6e2

471.88e-1

设 r , s , t 为正规式, 正规式服从代数规律有:

1、 $r|s=s|r$ 交换律

2、 $r|(s|t)=(r|s)|t$ 结合律

3、 $(rs)t=r(st)$ 结合律

4、 $r(s|t) = rs|rt$ 分配律

$(s|t)r = sr|tr$ 分配律

5、 $\epsilon r=r$ 零一律

$r\epsilon=r$ 零一律

程序设计语言中的单词既能用正规语法表示，又能用正规式来表示。

正规语法：

$\langle \text{标识符} \rangle \rightarrow 1 | 1 \langle \text{字母数字} \rangle$

$\langle \text{字母数字} \rangle \rightarrow 1 | d | 1 \langle \text{字母数字} \rangle | d \langle \text{字母数字} \rangle$

$\langle \text{无符号整数} \rangle \rightarrow d | d \langle \text{无符号整数} \rangle$

1表示a-z中的任何英文字母

d表示0-9中的任何数字

正规式？



程序设计语言中的单词既能用正规文法表示，又能用正规式来表示。

正规文法：

$\langle \text{标识符} \rangle \rightarrow 1 | 1 \langle \text{字母数字} \rangle$

$\langle \text{字母数字} \rangle \rightarrow 1 | d | 1 \langle \text{字母数字} \rangle | d \langle \text{字母数字} \rangle$

$\langle \text{无符号整数} \rangle \rightarrow d | d \langle \text{无符号整数} \rangle$

1表示a-z中的任何英文字母

d表示0-9中的任何数字

正规式：

标识符： $e1 = 1(l|d)^*$

无符号整数： $e2 = dd^*$

Quiz: 选择题

- 下面的正则表达式对应的语言中包含多少个不同的字符串？

$(0|1|\varepsilon)(0|1|\varepsilon)(0|1|\varepsilon)(0|1|\varepsilon)$

a) 12

b) 16

c) 31 ✓

d) 32

e) 64

f) 81

3、正规文法和正规式的等价性

一个正规语言可以由正规文法定义，也可以用正规式定义。

对于任意一个正规文法，存在一个定义同一语言的正规式。

对每一个正规式，存在一个生成同一语言的正规文法。

即正规式 \Leftrightarrow 正规文法

① 正规式 \Rightarrow 正规文法：（把正规式转换为正规文法所要求的规则形式）

将 Σ 上的一个正规式 r 转换为一个正规文法 $G = (V_N, V_T, P, S)$ 的规则：

令 $V_T = \Sigma$,

对正规式 r , 选择一个非终结符 S 生成 $S \rightarrow r$, S 为 G 的开始符号。不断拆分 r 直到符合正规文法要求的规则形式：

◆ 若 x, y 都是正规式, 对形如 $A \rightarrow xy$ 的产生式, 写成 $A \rightarrow xB, B \rightarrow y$ 。其中 $B \in V_N$

◆ 对形如 $A \rightarrow x^*y$ 的产生式, 重写为:

$A \rightarrow x A$

$A \rightarrow x B$

$A \rightarrow y$

$A \rightarrow y$

B 为新的非终结符, $B \in V_N$

$B \rightarrow x B$

对形如 $A \rightarrow x|y$ 的产生式, 重写为:

$B \rightarrow y$

$A \rightarrow x$

$A \rightarrow y$

不断利用上述规则进行变换即可。

例：将 $R = a(a|d)^*$ 变换成正规文法。令 S 是文法开始符号。

例：将 $R = a(a|d)^*$ 变换成正规文法。令 S 是文法开始符号。

解：

$$S \rightarrow a(a|d)^* \quad \left\{ \begin{array}{l} S \rightarrow aA \\ A \rightarrow (a|d)^* \end{array} \right. \rightarrow \left\{ \begin{array}{l} A \rightarrow (a|d) A \\ A \rightarrow \varepsilon \end{array} \right.$$

最后得到:

$$S \rightarrow aA$$

$$A \rightarrow aA$$

$$A \rightarrow dA$$

$$A \rightarrow \varepsilon$$

② **正规文法 \Rightarrow 正规式**: 将一个正规文法转换为正规式的规则:
转换规则:

① **$A \rightarrow xB, B \rightarrow y$** 正规式为: **$A = xy$**

② **$A \rightarrow xA|y$** 正规式为: **$A = x^*y$**

③ **$A \rightarrow x, A \rightarrow y$** 正规式为: **$A = x|y$**

不断收缩产生式规则, 直到剩下一个开始符号定义的正规式

例: 文法 $G[S]$:

$S \rightarrow aA$

$S \rightarrow a$

$A \rightarrow aA$

$A \rightarrow dA$

$A \rightarrow a$

$A \rightarrow d$

转换为正规式

$S \rightarrow aA$

$S \rightarrow a$

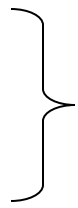


$S = aA | a$

根据上述规则3
 $A \rightarrow x, A \rightarrow y$
推出 $A = x | y$

$A \rightarrow aA$

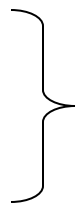
$A \rightarrow dA$



$A = aA | dA$

$A \rightarrow a$

$A \rightarrow d$



$A = a | d$

$A = (aA | dA) | (a | d) \Rightarrow$

$A = (a | d)A | (a | d) \Rightarrow$

$A = (a | d)^*(a | d)$

将它化为正规文法
变成 $A \rightarrow$
 $(a | d)A | (a | d)$

再根据上述
规则2转换
 $x = y = (a | d)$

将A代入 $S = aA|a$ 得到如下：

$$S = a((a|d)^*(a|d))|a$$

$$= a(a|d)^+|a$$

$$= a((a|d)^+|\epsilon) = a(a|d)^*$$



4.3 有穷自动机

有穷自动机：是一种自动识别装置，能正确识别正规集；是词法分析程序的工具和方法，可自动识别（且是正确识别）正规集。

分为 { 确定的有穷自动机 (**DFA**)
不确定的有穷自动机 (**NFA**)



(一) 确定的有穷自动机 **DFA** 自动识别装置

一个确定的有穷自动机 **M** 是一个五元组：

M = (K, Σ, f, S, Z)，其中：

- ① **K** 是一个有穷集，每个元素表示一个状态；
- ② **Σ** 是一个有穷字母表，每个元素是一个输入字符；
- ③ **f** 是转换函数，是在 $\mathbf{K} \times \mathbf{\Sigma} \rightarrow \mathbf{K}$ 上的映象，如：

$$\mathbf{f(K_i, a) = K_j} \quad (\mathbf{K_i, K_j \in K});$$

- ④ **S** 是初态， $\mathbf{S \in K}$ ；
- ⑤ $\mathbf{Z \subset K}$ ，是终态集。

含义：当前状态为 K_i ，输入字符 a ，转换为 K_j 状态

状态转换图

- **状态转换图(transition diagram)**
 - **状态(state)**: 表示在识别词素时可能出现的情况
 - 状态看作是已处理部分的总结
 - 某些状态为接受状态或最终状态, 表明已找到词素
 - 开始状态(初始状态): 用“**开始**”边表示
 - **边(edge)**: 从一个状态指向另一个状态; 边的标号是一个或多个符号
 - 当前状态为s, 下一个输入符号为a, 就从s离开, 沿着标号为a的边到达下一个状态
-

1、用状态图表示DFA：

方法如下：

- ◆ 初始态用 “ \Rightarrow ” 或 “ $-$ ” 表示；
- ◆ 终态点用 “ $+$ ” 或 “ \odot ” 表示；
- ◆ 若 $f(K_i, a) = K_j$ ，则从状态点 K_i 到 K_j 画弧，标记为 a 。

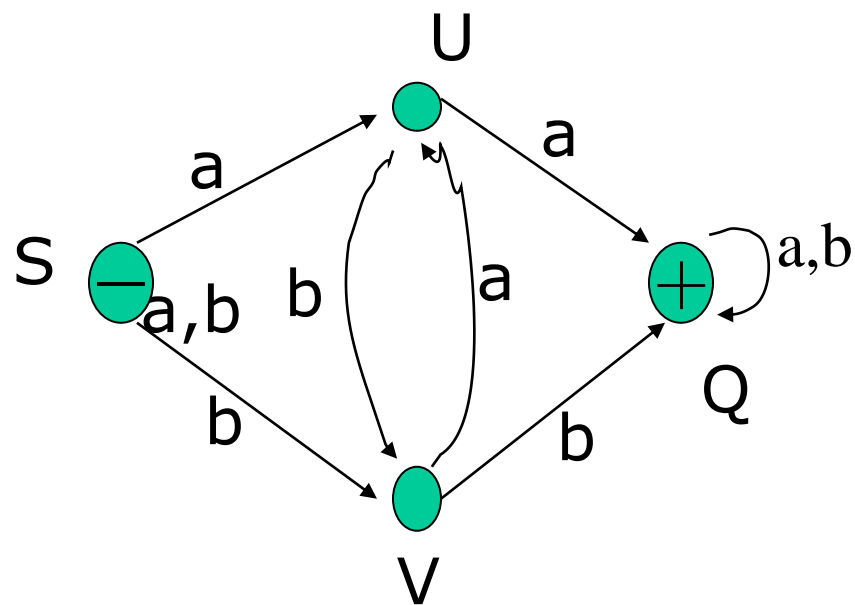
例4.6：DFA的 $M = (\{S, U, V, Q\}, \{a, b\}, f, S, \{Q\})$

其中 f 为：

$f(S, a) = U,$	$f(S, b) = V,$	$f(U, a) = Q$
$f(U, b) = V,$	$f(V, a) = U,$	$f(V, b) = Q$
$f(Q, a) = Q,$	$f(Q, b) = Q$	

[画出状态图](#)

状态图如下：



2、用矩阵表示DFA：

方法：◆ 行表示状态

◆ 列表示输入字符

◆ 元素表示相应状态行和输入字符下的新状态。



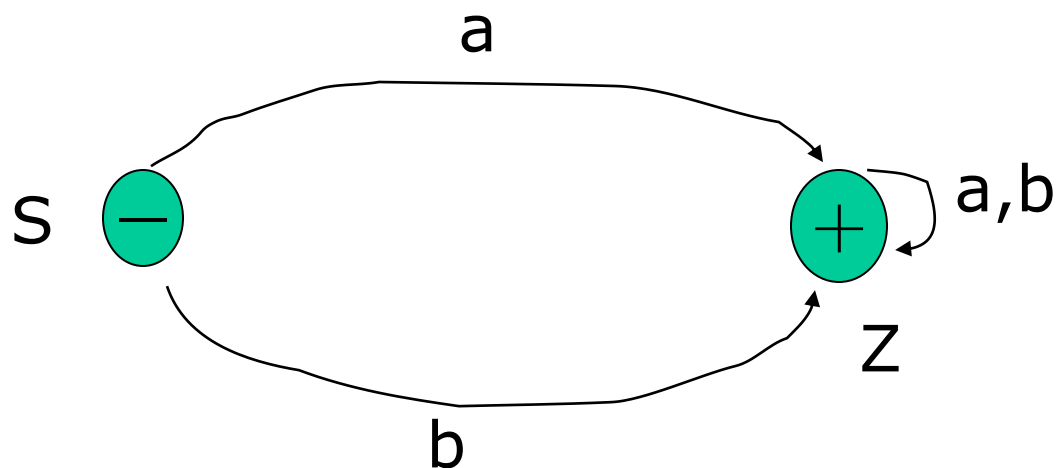
◆ “ \Rightarrow ” 标明初态，默认第一行是初态。

◆ 终态行在表右端标1，非终态标0

上例矩阵表示如下：

状态 \ 字符	a	b	
S	U	V	0
U	Q	V	0
V	U	Q	0
Q	Q	Q	1

例：



表示： $f(S,a)=Z$

$f(S,b)=Z$

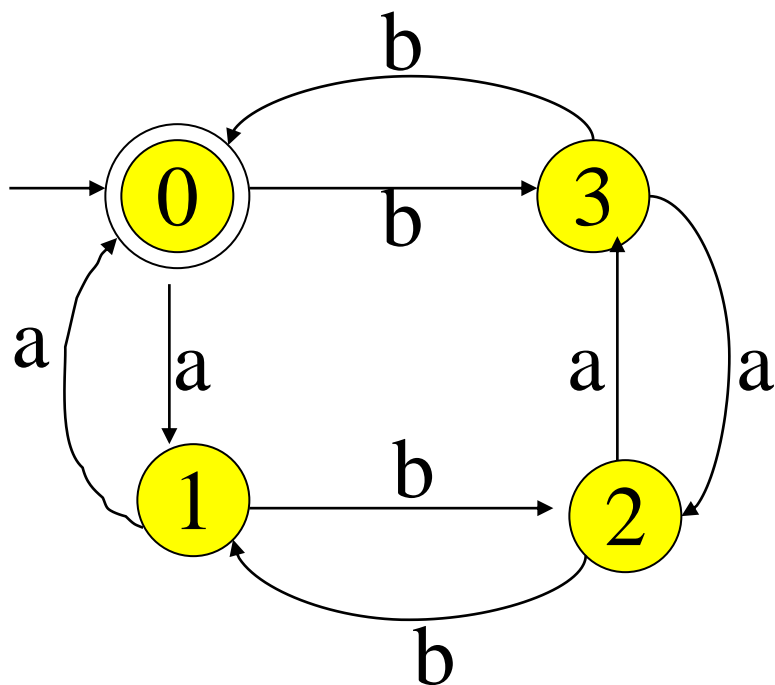
$f(Z,a)=Z$

$f(Z,b)=Z$

	a	b	
S	Z	Z	0
Z	Z	Z	1

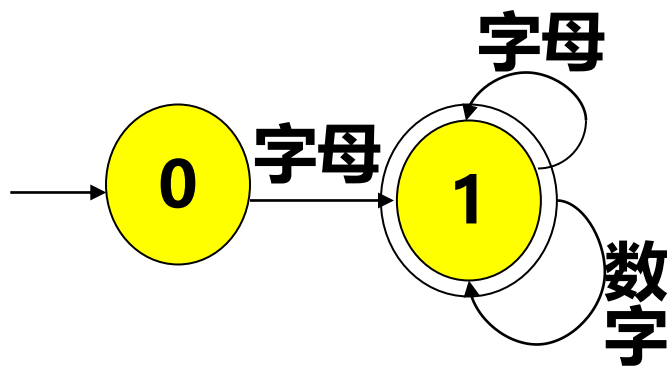
正规式是： $(a|b)(a|b)^*$

从状态转换图看，从开始状态出发，沿任一条路径到达接受状态，这条路径上的弧上的标记符号连接起来构成的符号串被DFA M接受。

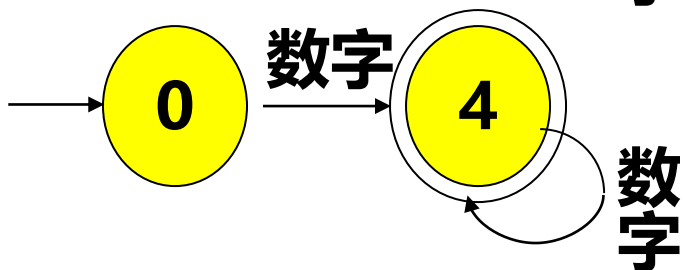


偶数个a, 偶数个b
的{a,b}串集合

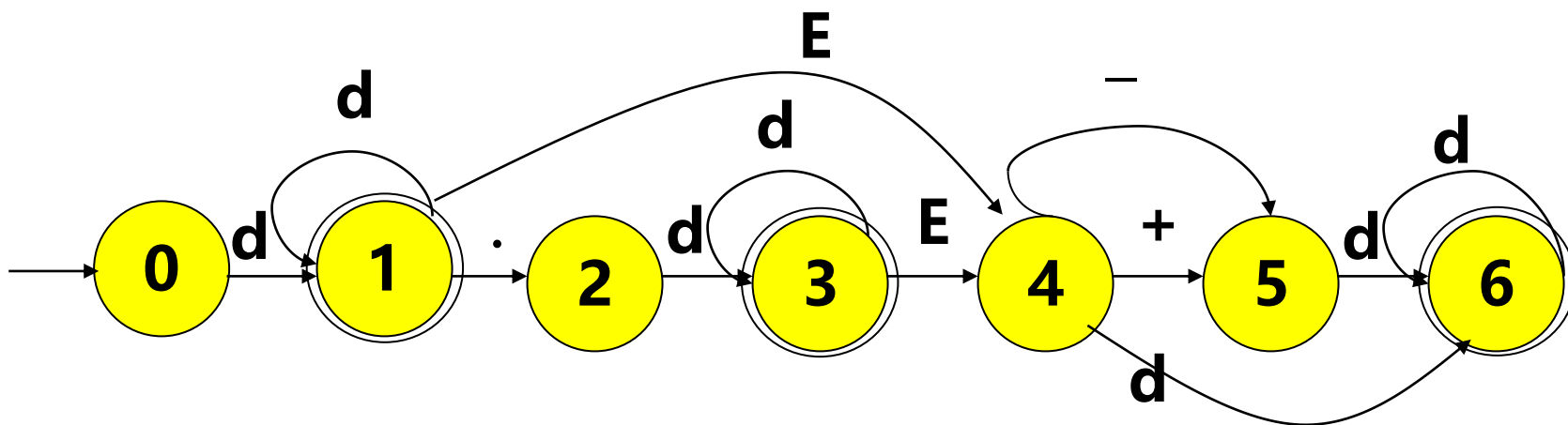
DFA示例-1



Pascal 标识符



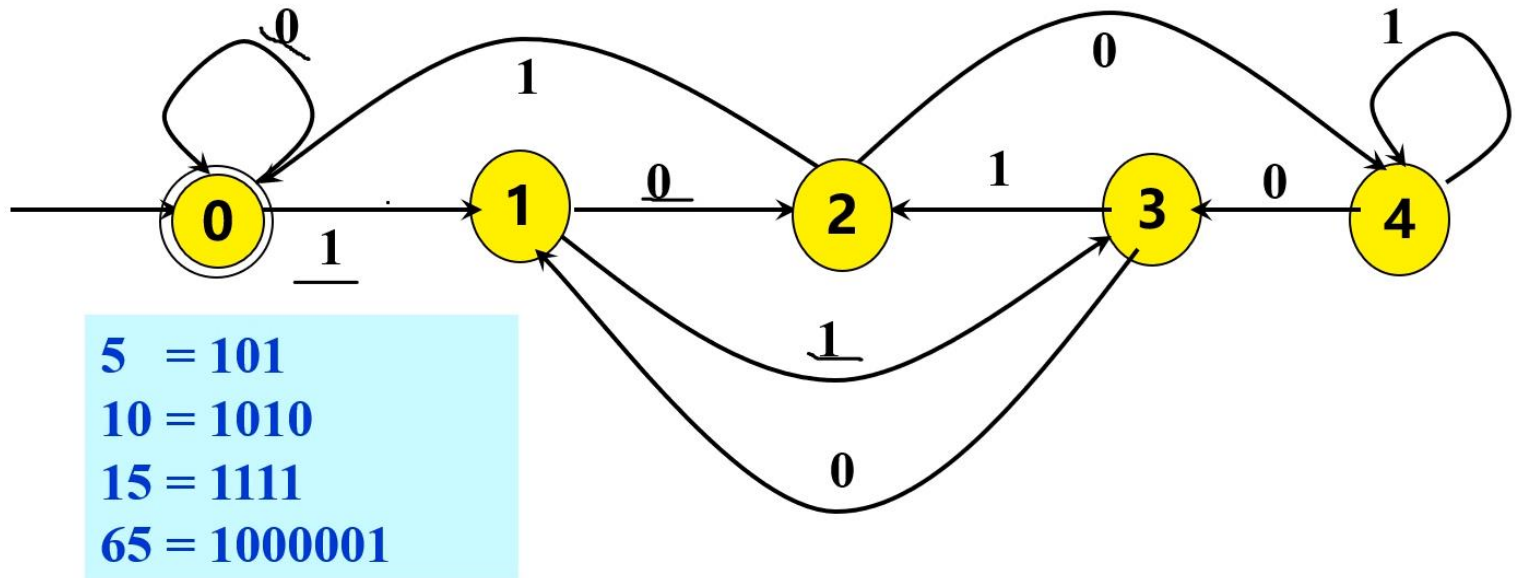
Pascal 整数和实数



DFA示例-2

识别 $\Sigma = \{0,1\}$ 上能被5整除的二进制数

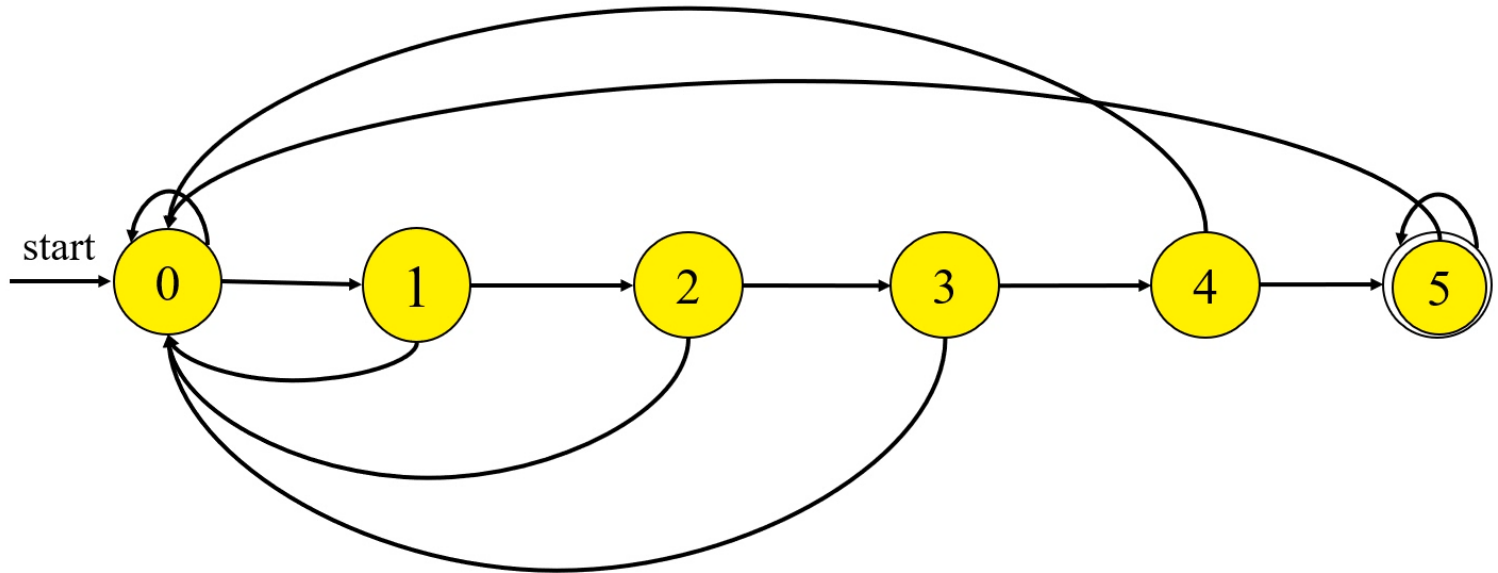
按余数跳转



Quiz

- 画一个DFA，表示所有能被32整除的二进制数。

末尾0的个数



3、接受（识别）的概念： 自动识别单词

对于 Σ^* 中的任何字符串 t ，若存在一条初态到某一终态的路，且这条路上所有弧的标记符连接成的字符串等于 t ，则称 t 可为**DFA M**所接受。

若**M**的初态同时又是终态，则空字可为**M**所接受。

4、DFA M 运行的理解：

① 设 $Q \in K$ ，函数 $f(Q, \varepsilon) = Q$ ，则输入字符串是空串，并停留在原状态上。

② 输入字符串 t （ t 表示成 Tt_1 形式， $T \in \Sigma$ ， $t_1 \in \Sigma^*$ ），在**DFA M**上运行的定义为： $f(Q, Tt_1) = f(f(Q, T), t_1)$ ，其中 $Q \in K$ 。

例如：**baab**字符串被**DFA**所接受，**DFA**见上例。

$$\begin{aligned} f(S, baab) &= f(f(S, b), aab) = f(V, aab) \\ &= f(f(V, a), ab) = f(U, ab) = f(f(U, a), b) \\ &= f(Q, b) = Q \quad (Q \text{是终态}) \end{aligned}$$

③**DFA M**所能接受的字符串的全体记为**L (M)** — 称为**语言**（也即句子的集合）

5、DFA的确定性：

当**f: $k \times \Sigma \rightarrow K$** 是一个**单值函数**，即对任何状态 **$K \in R$** ，输入符号 **$a \in K$** ， **$f(k, a)$** 唯一确定下一状态。

DFA的行为很容易用程序来模拟.

DFA $M = (K, \Sigma, f, S, Z)$ 的行为的模拟程序

- $K := S;$
- $c := \text{getchar};$
- while $c \neq \text{eof}$ do
- $\{ K := f(K, c);$
- $c := \text{getchar};$
- $\};$
- if K is in Z then return ('yes')
- else return ('no')



自动识别单词的方法：

- (1) 把单词的结构用**正规式**描述；
- (2) **把正规式转换为一个NFA**；
- (3) **把NFA转换为相应的DFA**；
- (4) **基于DFA构造词法分析程序。**

正则表达式 \Rightarrow NFA \Rightarrow DFA \Rightarrow min DFA

(二) 不确定的有穷自动机NFA

一个不确定的有穷自动机NFA M 是一个五元组：
 $M = (K, \Sigma, f, S, Z)$ ，其中：

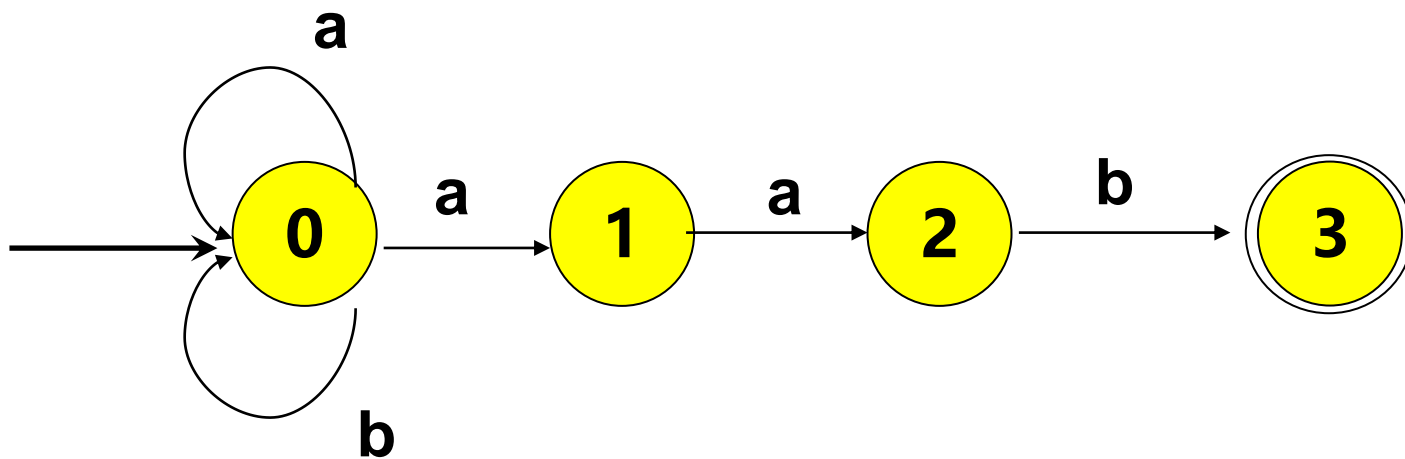
- ① K 是一个有穷集，每个元素表示一个状态；
- ② Σ 是一个有穷字母表，每个元素是一个输入字符；
- ③ f 是一个从 $K \times \Sigma^*$ 到 K 上的子集的映象；

$$f: k \times \Sigma^* \rightarrow 2^k$$

- ④ $S \subset K$ ，是一个非空初态集；
- ⑤ $Z \subset K$ ，是一个终态集。

与DFA区别：多值函数 $f(K_i, a) = K_j \quad f(K_i, a) = K_t$ ；允许输入字符为 ε

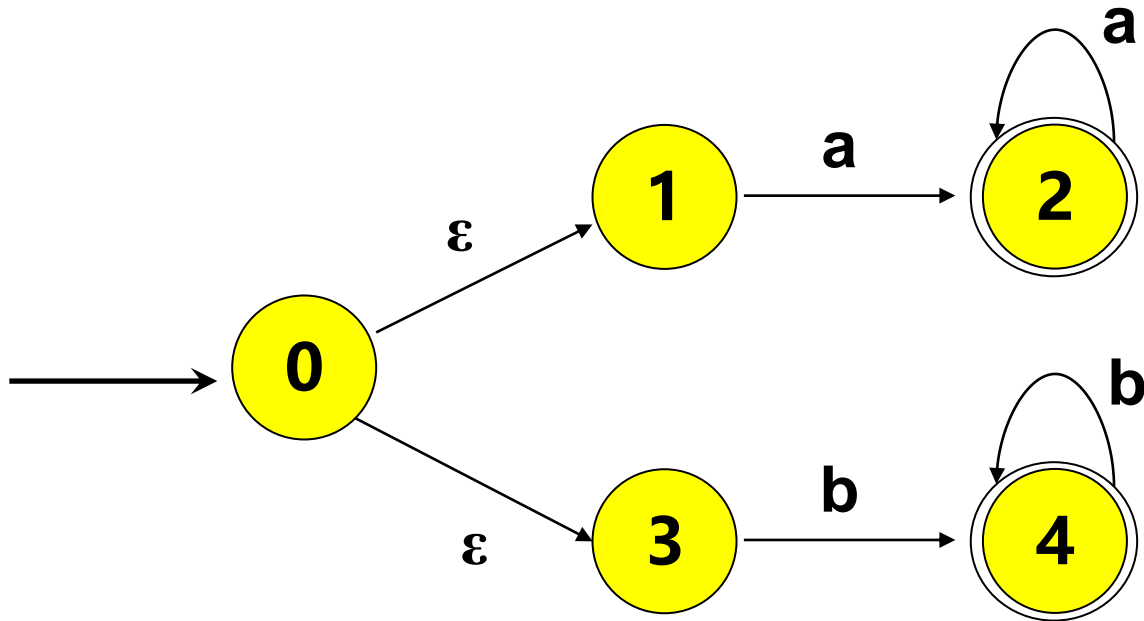
NFA示例-1



NFA 表示的语言:

$(a|b)^*aab$

NFA示例-2



NFA 表示的语言:

$aa^*|bb^*$

可不可以用DFA

来表示?

例4.7: 一个NFA,

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

其中: $f(0,a) = \{0,3\}$

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

$f(0,b) = \{0,1\}$

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

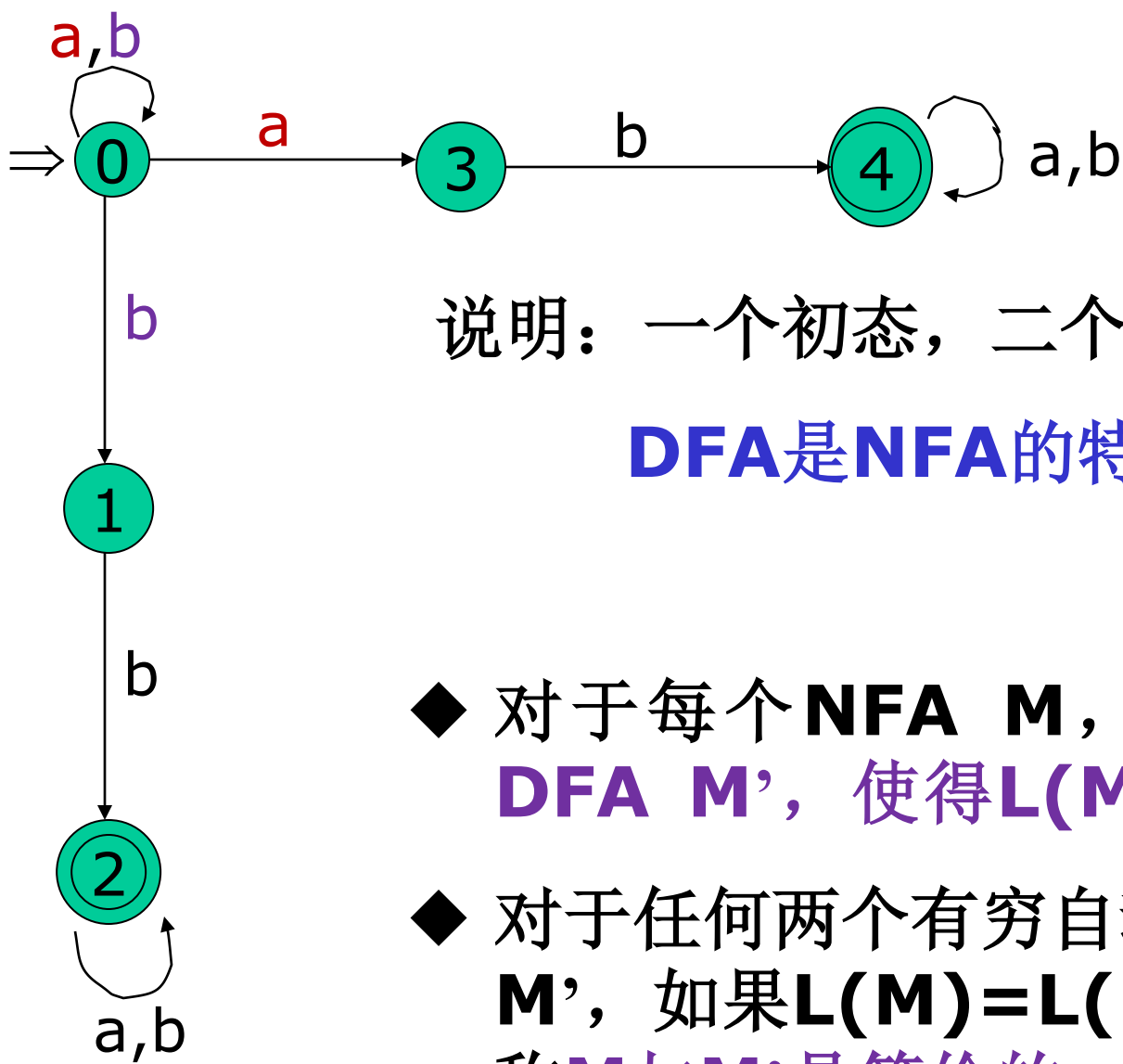
$f(1,b) = \{2\}$

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

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

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

状态图表示如下:



说明：一个初态，二个终态。

DFA是NFA的特例。

- ◆ 对于每个 **NFA M**，存在一个 **DFA M'**，使得 **$L(M) = L(M')$** 。
- ◆ 对于任何两个有穷自动机 **M** 和 **M'**，如果 **$L(M) = L(M')$** ，则称 **M** 与 **M'** 是等价的。

NFA 转换为等价的DFA

- 从NFA的矩阵表示中可以看出，表项通常是一**状态的集合**，而在DFA的矩阵表示中，表项是一个**状态**
 - NFA \Rightarrow DFA 基本思路是：**DFA的每一个状态对应NFA的一组状态**，DFA使用它的状态去记录在NFA读入一个输入符号后可能达到的所有状态
-

定义对状态集合I的几个有关运算：

1. **状态集合I的 ϵ -闭包**，表示为 ϵ -closure(I)，定义为一状态集，是状态集I中的任何状态S经任意条 ϵ 弧而能到达的状态的集合。

状态集合I的任何状态S都属于 ϵ -closure(I)

2. **状态集合I的a弧转换**，表示为move(I, a)定义为状态集合J，其中J是所有那些可从I中的某一状态经过一条a弧而到达的状态的全体。



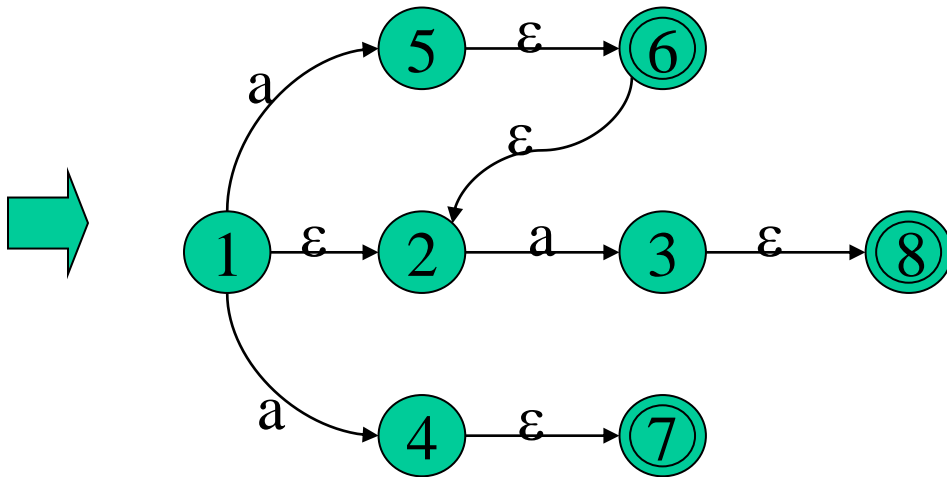
状态集合I的有关运算的例子

$I=\{1\}$, ϵ -closure(I)=?;

$I=\{5\}$, ϵ -closure(I)=?;

move($\{1,2\}$,a)=?

ϵ -closure($\{5,3,4\}$)=?;



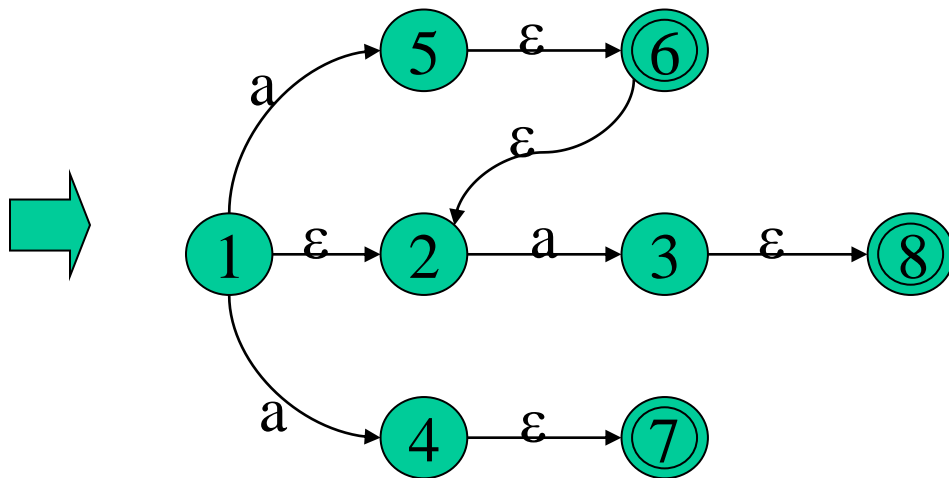
状态集合I的有关运算的例子

$I=\{1\}$, ϵ -closure(I)= $\{1,2\}$;

$I=\{5\}$, ϵ -closure(I)= $\{5,6,2\}$;

$\text{move}(\{1,2\},a)=\{5,3,4\}$

ϵ -closure($\{5,3,4\}$)= $\{2,3,4,5,6,7,8\}$;



NFA确定化算法（子集法）

NFA $N=(K, \Sigma, f, K_0, K_t)$ 按如下办法构造一个DFA $M=(S, \Sigma, d, S_0, S_t)$ ，使得

$L(M)=L(N)$ ：

1. 构造DFA M 的状态，选择NFA N 的状态的一些子集构成：

M 的状态集 S 由 K 的一些子集组成。用 $[S_1 S_2 \dots S_j]$ 表示 S 的元素，其中 S_1, S_2, \dots, S_j 是 K 的状态。并且约定，状态 S_1, S_2, \dots, S_j 是按某种规则排列的，即对于子集 $\{S_1, S_2\}=\{S_2, S_1\}$ 来说， S 的状态就是 $[S_1 S_2]$ ；

2. M和N的输入字母表是相同的，即是 Σ ；

3. 构造DFA M的转换函数，根据新构造的状态和字母表中的字符构造：

转换函数是这样定义的：

$$d([S_1 S_2 \dots S_j], a) = [R_1 R_2 \dots R_t]$$

其中 $\{R_1, R_2, \dots, R_t\} = \varepsilon\text{-closure}(\text{move}(\{S_1, S_2, \dots, S_j\}, a))$

4 $S_0 = \varepsilon\text{-closure}(K_0)$ 为M的开始状态；

5 $S_t = \{[S_i S_k \dots S_e], \text{其中}[S_i S_k \dots S_e] \in S \text{且} \{S_i, S_k, \dots, S_e\} \cap K_t \neq \Phi\}$

构造NFA N 的状态 K 的子集的算法：

把多值转换函数所转换到的多值（多状态）的集合作为一个子集，映射到DFA的一个新的状态

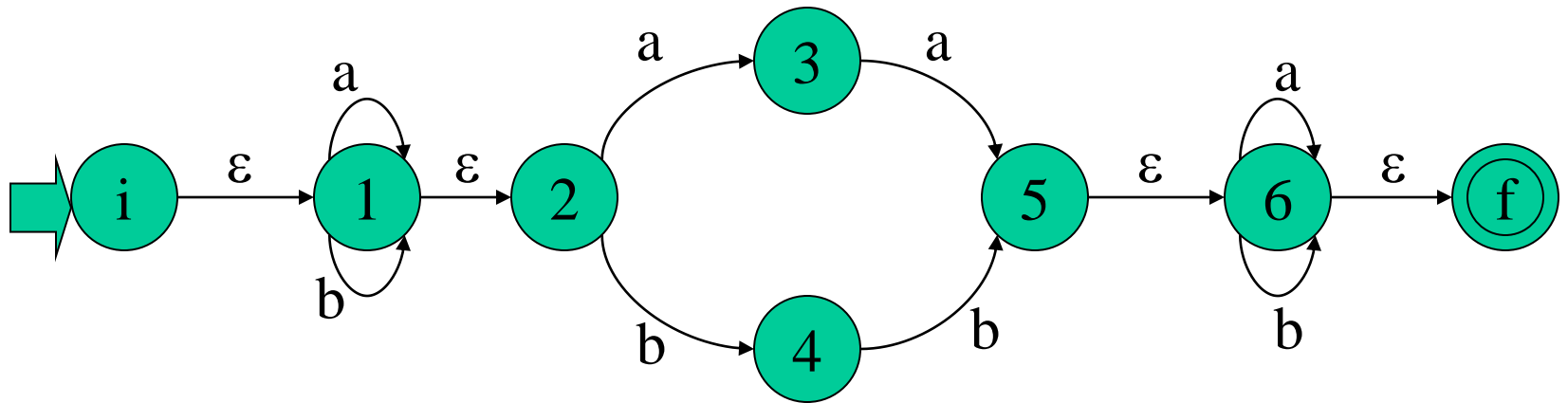
假定所构造的子集族为 C ，即 $C = (T_1, T_2, \dots, T_I)$ ，其中 T_1, T_2, \dots, T_I 为状态 K 的子集。

- 1 开始，令 ε -closure(K_0)为 C 中唯一成员，并且它是未被标记的。
-

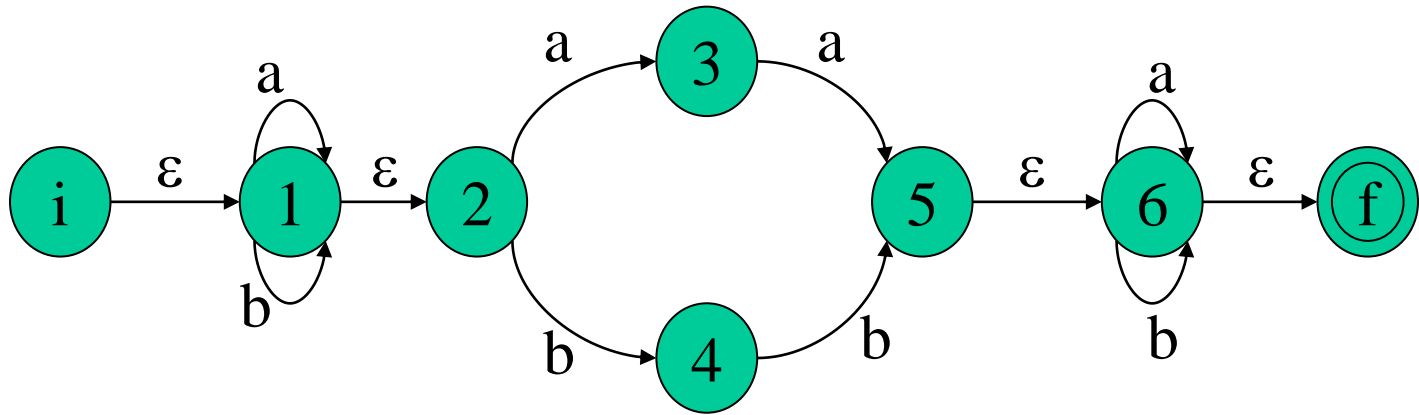
```
2 while (C中存在尚未被标记的子集T) do
{
  标记T;
  for 每个输入字母a do
  {
    U:=  $\epsilon$ -closure(move(T,a));
    if U不在C中 then
      将U作为未标记的子集加在C中
  }
}
```

NFA的确定化

例子

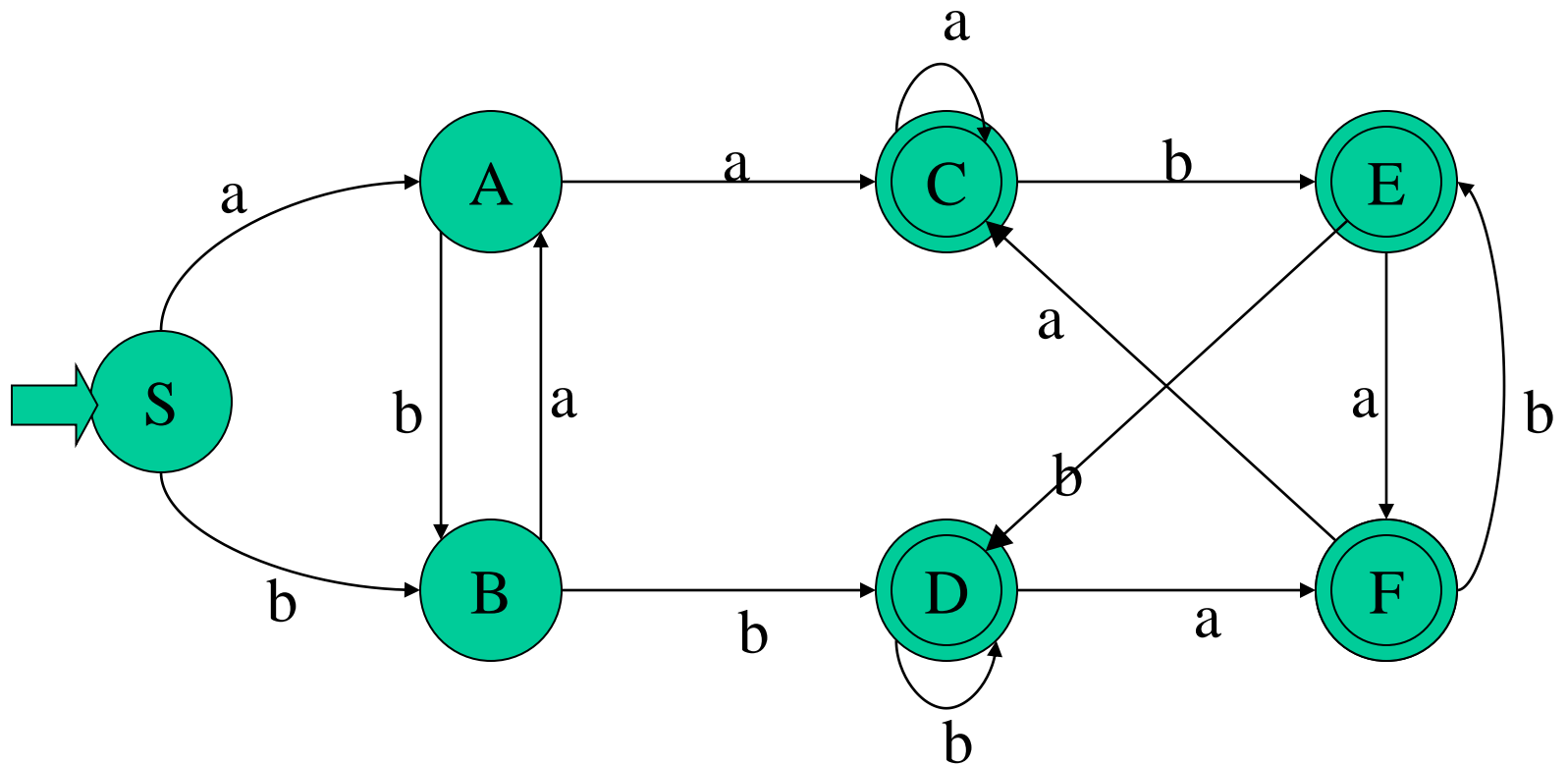


DFA还是NFA?



T	ϵ -clo (move(T,a))	ϵ -clo (move(T,b))
{i,1,2} S	{1,2,3} A	{1,2,4} B
{1,2,3} A	{1,2,3,5,6,f} C	{1,2,4} B
{1,2,4} B	{1,2,3} A	{1,2,4,5,6,f} D
{1,2,3,5,6,f} C	{1,2,3,5,6,f} C	{1,2,4,6,f} E
{1,2,4,5,6,f} D	{1,2,3,6,f} F	{1,2,4,5,6,f} D
{1,2,4,6,f} E	{1,2,3,6,f} F	{1,2,4,5,6,f} D
{1,2,3,6,f} F	{1,2,3,5,6,f} C	{1,2,4,6,f} E

等价的DFA



DFA的最小化

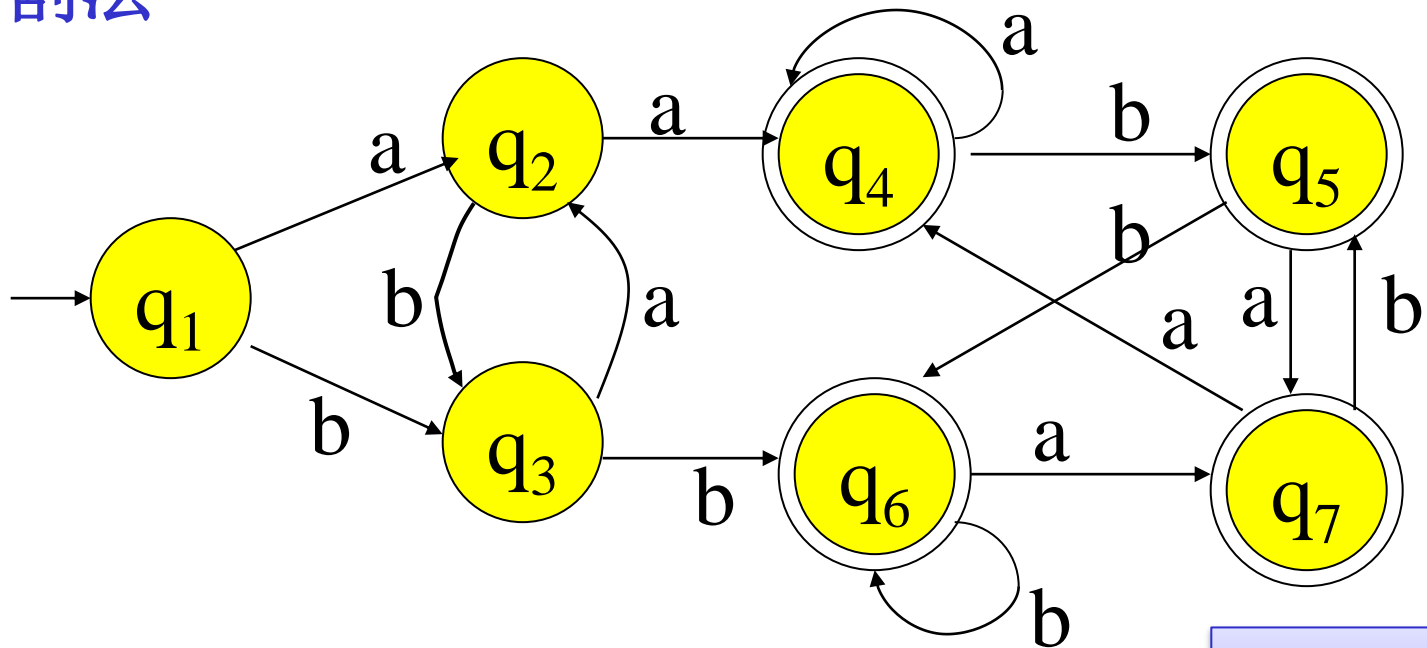
- DFA的化简

- 设有DFA M , 寻找一个状态数更少的DFA M' , 使 $L(M') = L(M)$
- 可以证明, 存在一个最少状态的DFA M' , 使 $L(M) = L(M')$ 。

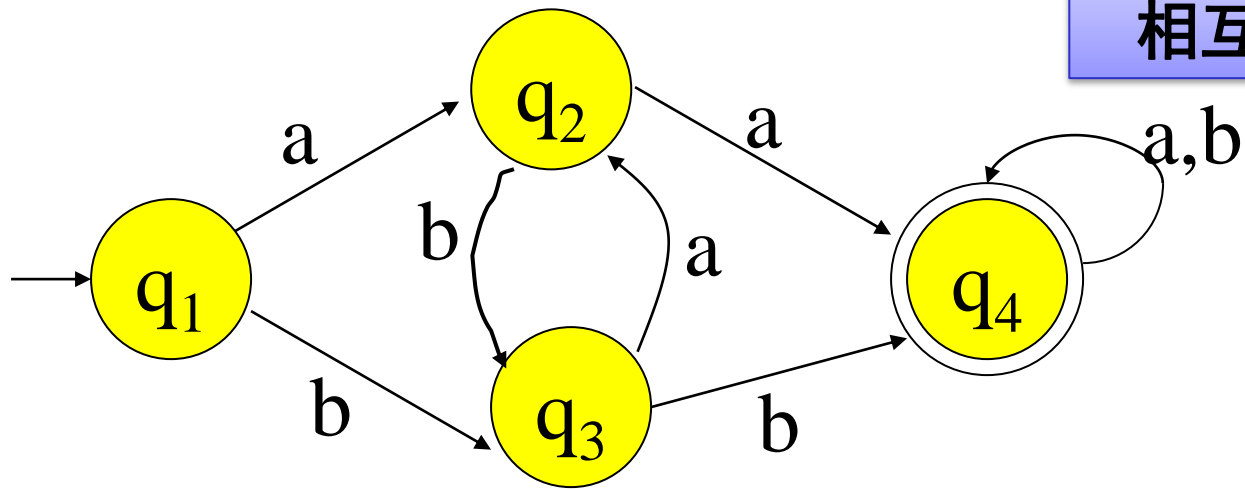
- 等价条件

- 一致性条件: 状态 s, t 必须同时为可接受状态或不可接受状态
- 蔓延性条件: 对任意输入符号, 状态 s 和 t 必须转换到等价的状态里, 否则是**可区分的**。

分割法



q4 q5 q6 q7
相互等价。



内容回顾

- 自动识别单词的方法：
 - (1) 把单词的结构用**正规式**描述；
 - (2) 把正规式转换为一个**NFA**；
 - (3) 把**NFA**转换为相应的**DFA**；
 - (4) 基于**DFA**构造词法分析程序。



课后作业

- **P72 : 1 (1) , 2, 5**
- **P73: 8**

