

有限状态自动机

乔明达

2014年2月8日

开场白

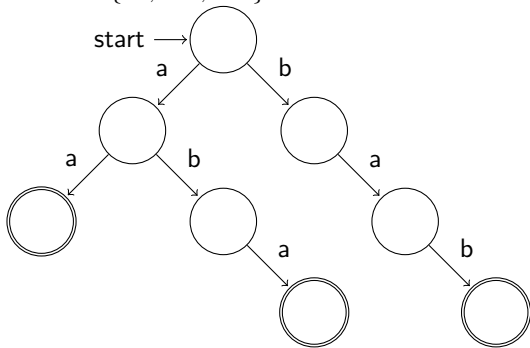
Michael Sipser 《计算理论导引》

可计算性(computability)与计算复杂性(computational complexity)

例: Trie

判断一个串 t 是否在给定的串的集合 S 中

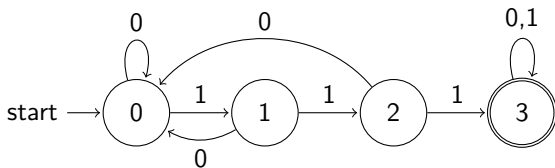
例: $S = \{aa, aba, bab\}$



例：KMP

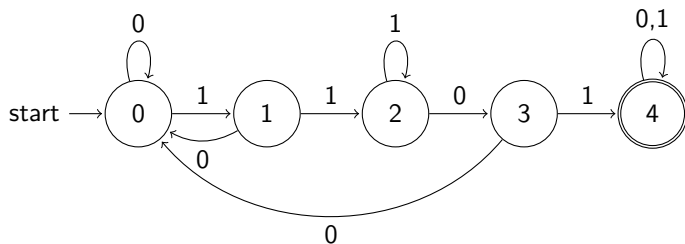
判断一个串 t 是否包含子串 s

例： $s = 111$



例：KMP

例： $s = 1101$



有限状态自动机

Finite Automaton /'faɪ,naɪt ɔ:'tɒmətən/

$(Q, \Sigma, \delta, q_0, F)$

状态集合 Q

字母表 Σ

转移函数 $\delta : Q \times \Sigma \rightarrow Q$

开始状态 $q_0 \in Q$

接受状态集合 $F \subseteq Q$

有限状态自动机

$$s = w_1 w_2 \dots w_n, w_i \in \Sigma$$

$$r_0 = q_0, r_i = \delta(r_{i-1}, w_i)$$

若 $r_n \in F$ ，则有限状态自动机接受串 s ，否则拒绝串 s

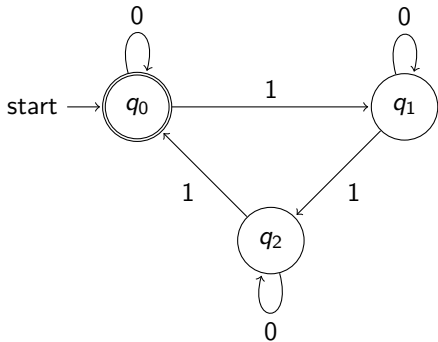
有限状态自动机

字母表 $\Sigma = \{0, 1\}$, $p \in \mathbb{N}^+$

接受包含1的个数是 p 的倍数的串

例如对于 $p = 3$, 应该接受串111, 001011, 000, 拒绝串10, 1001

有限状态自动机



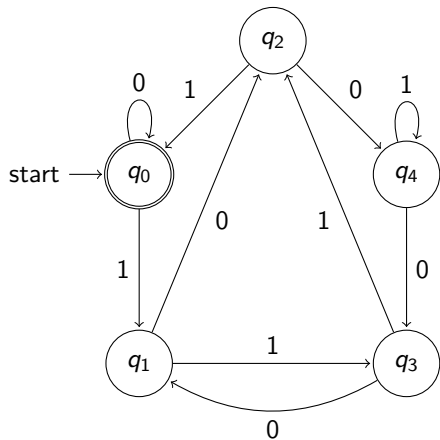
有限状态自动机

字母表 $\Sigma = \{0, 1\}$, $p \in \mathbb{N}^+$

接受 p 的倍数的二进制表示（允许前导零）

例如对于 $p = 5$, 应该接受串 101, 01010, 000, 拒绝串 11, 0001

有限状态自动机



串与语言

Σ^* 为字母表 Σ 上串的集合

例：当 $\Sigma = \{0, 1\}$ 时， $\Sigma^* = \{\epsilon, 0, 1, 00, 01, 10, 11, \dots\}$

Σ^* 的任意子集称为字母表 Σ 上的语言

有限状态自动机 M 接受的串的集合称为 M 的语言，记作 $L(M)$

可以被确定有限状态自动机识别的语言称为正则语言(regular language)

串的运算

xy

$$x^k = \underbrace{xx \cdots x}_k$$

$$x^0 = \varepsilon$$

语言的运算

交与并

补集: $\bar{A} = \{x | x \in \Sigma^*, x \notin A\}$

连接: $A \circ B = \{xy | x \in A, y \in B\}$

幂: $A^k = \{x_1 x_2 \dots x_k | x_1, x_2, \dots, x_k \in A\}$

闭包: $A^* = A^0 \cup A^1 \cup A^2 \cup \dots$

语言的运算

设有语言 $A = \{\varepsilon\}$, $B = \emptyset$

求: $A \circ A, A \circ B, B \circ B, A^*, B^*$

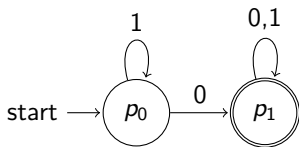
正则语言的封闭性

并集、交集：若 A 和 B 是正则语言，则 $A \cup B$ 和 $A \cap B$ 是正则语言

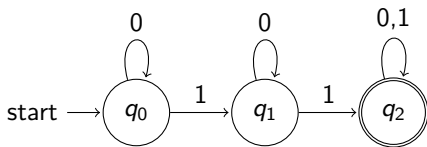
正则语言的封闭性

$$\Sigma = \{0, 1\}$$

$$A = \{x \mid x \text{ 包含至少1个0}\}$$

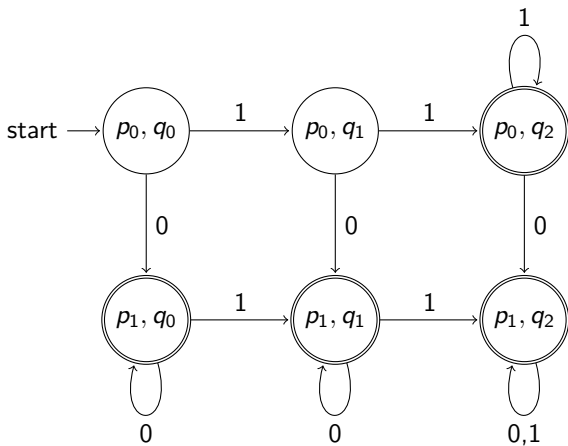


$$B = \{x \mid x \text{ 包含至少2个1}\}$$



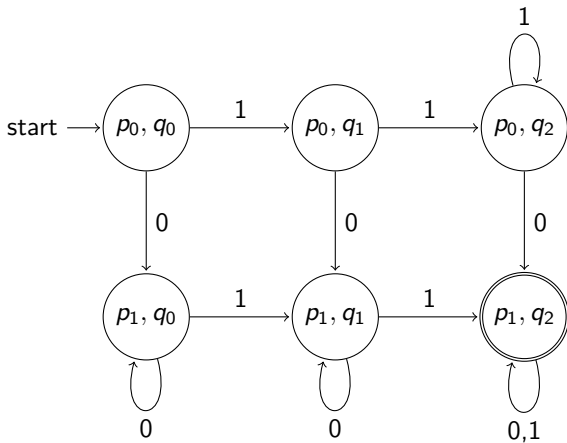
正则语言的封闭性

$A \cup B$



正则语言的封闭性

$A \cap B$



正则语言的封闭性

$M_A = (Q_A, \Sigma, \delta_A, q_A, F_A)$ 识别语言 A

$M_B = (Q_B, \Sigma, \delta_B, q_B, F_B)$ 识别语言 B

令 $M = (Q, \Sigma, \delta, q_0, F)$

其中 $Q = Q_A \times Q_B, q_0 = (q_A, q_B)$

$\delta((q_1, q_2), c) = (\delta_A(q_1, c), \delta_B(q_2, c))$

对于 $A \cup B$, 令 $F = (F_A \times Q_B) \cup (Q_A \times F_B)$

对于 $A \cap B$, 令 $F = F_A \times F_B$

正则语言的封闭性

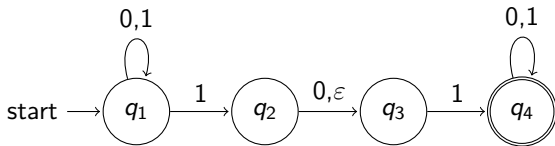
补集：若 A 是正则语言，则 \overline{A} 是正则语言

正则语言的封闭性

连接：若 A 和 B 是正则语言，则 $A \circ B$ 是正则语言

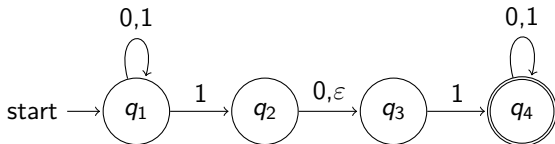
非确定有限状态自动机

DFA(Deterministic Finite Automaton)与NFA(Nondeterministic Finite Automaton)



接受串010110

非确定有限状态自动机



开始: q_1

读入0: q_1

读入1: q_1, q_2, q_3

读入0: q_1, q_3

读入1: q_1, q_2, q_3, q_4

读入1: q_1, q_2, q_3, q_4

读入0: q_1, q_3, q_4

非确定有限状态自动机

$(Q, \Sigma, \delta, q_0, F)$

令 $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$

状态集合 Q

字母表 Σ

转移函数 $\delta : Q \times \Sigma_\epsilon \rightarrow \mathcal{P}(Q)$

开始状态 $q_0 \in Q$

接受状态集合 $F \subseteq Q$

非确定有限状态自动机

若存在 $w_1, w_2, \dots, w_n, r_0, r_1, \dots, r_n$ 使得:

$$s = w_1 w_2 \dots w_n, w_i \in \Sigma \cup \{\varepsilon\}$$

$$r_0 = q_0, r_i \in \delta(r_{i-1}, w_i), r_n \in F$$

则NFA接受串 s , 否则NFA拒绝串 s

例：连接

连接：若 A 和 B 是正则语言，则 $A \circ B$ 是正则语言

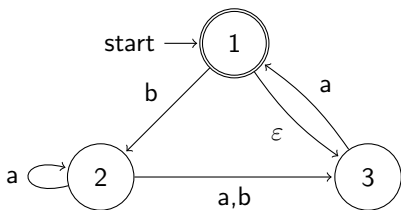
设DFAM $_A$ 和 M_B 分别识别语言 A 和语言 B

对于 M_A 中的每个接受状态，增加到 M_B 开始状态的 ϵ 转移

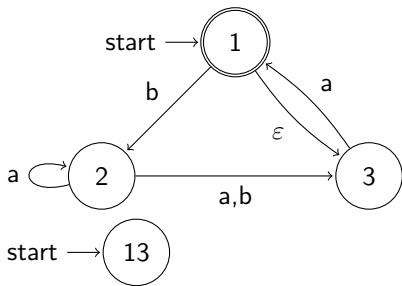
NFA与DFA的等价性

证明：对于每个NFA，存在一个DFA与它识别相同的语言

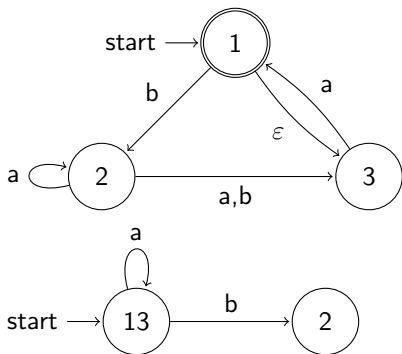
NFA与DFA的等价性



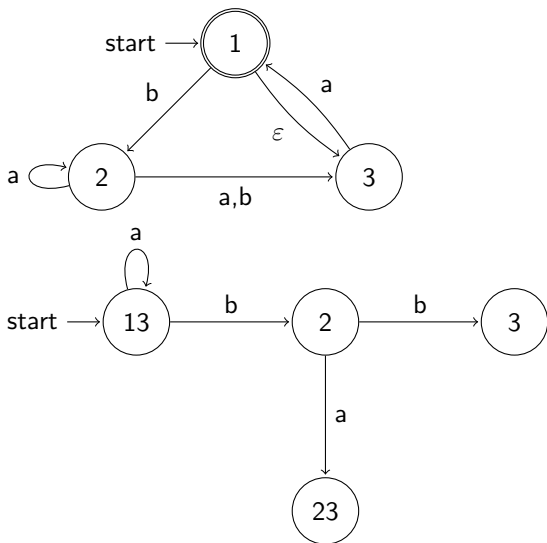
NFA与DFA的等价性



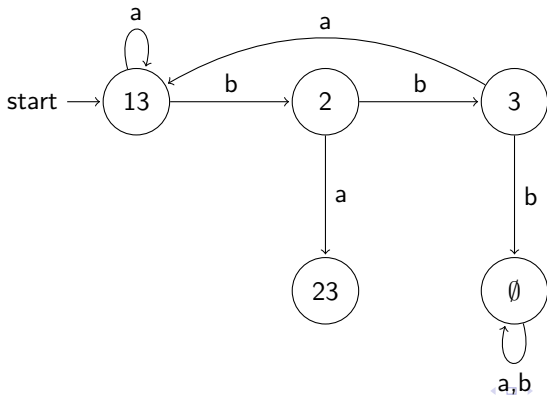
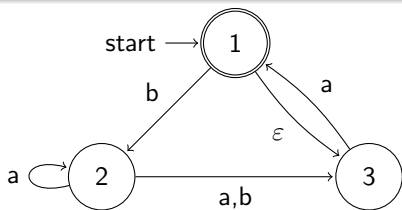
NFA与DFA的等价性



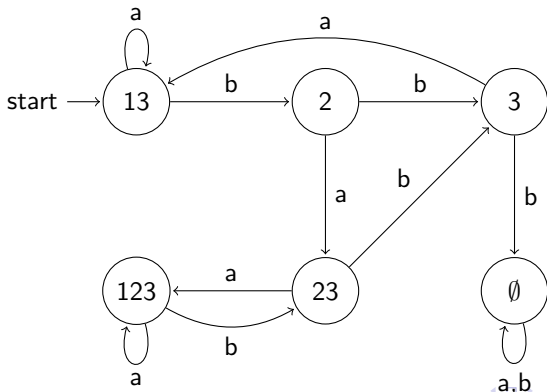
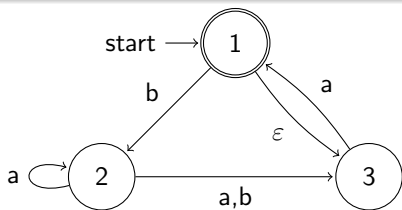
NFA与DFA的等价性



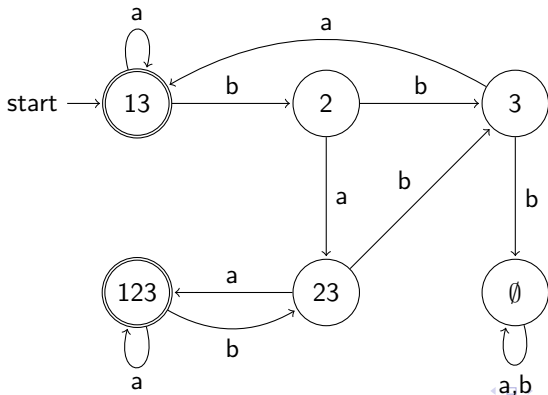
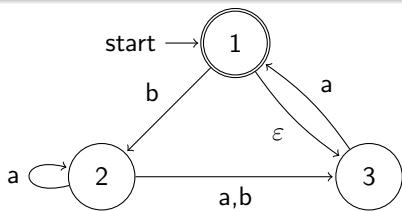
NFA与DFA的等价性



NFA与DFA的等价性



NFA与DFA的等价性



NFA与DFA的等价性

设 $N = (Q, \Sigma, \delta, q_0, F)$

定义 $E(q)$ 为：从状态 q 出发，只走 ε 转移能够到达的状态集合

构造 $M = (Q', \Sigma, \delta', E(q_0), F')$

状态集合 $Q' = \mathcal{P}(Q)$

转移函数 $\delta' : Q' \times \Sigma \rightarrow Q'$

其中 $\delta'(S, c) = \bigcup_{q \in S, q' \in \delta(q, c)} E(q')$

$F' = \{S \subseteq Q \mid S \cap F \neq \emptyset\}$

NFA与DFA的等价性

回忆定义：可以被确定有限状态自动机识别的语言称为正则语言

推论：一个语言是正则的，当且仅当存在一个NFA识别它

正则语言的封闭性

并集：若 A 和 B 是正则语言，则 $A \cup B$ 是正则语言

设DFAM $_A$ 和 M_B 分别识别语言 A 和语言 B ，开始状态分别为 q_A 和 q_B

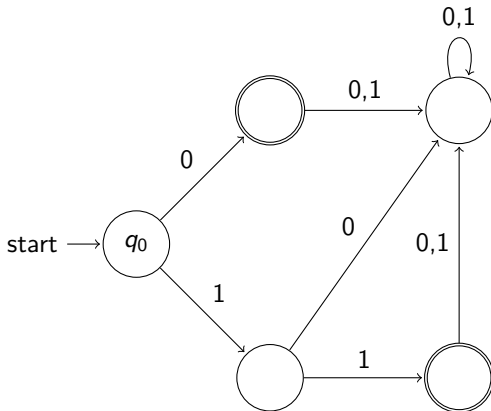
新增开始状态 q_0 ，令 $\delta(q_0, \varepsilon) = \{q_A, q_B\}$

正则语言的封闭性

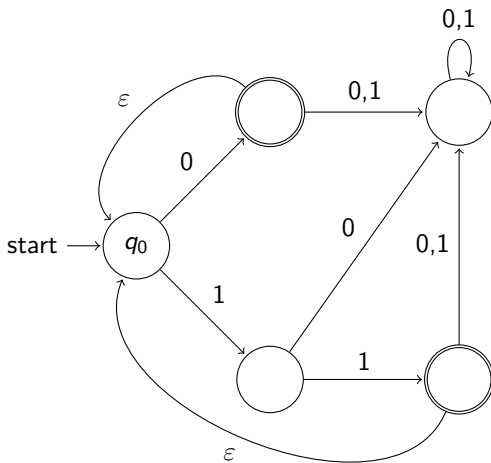
闭包：若 A 是正则语言，则 A^* 是正则语言

正则语言的封闭性

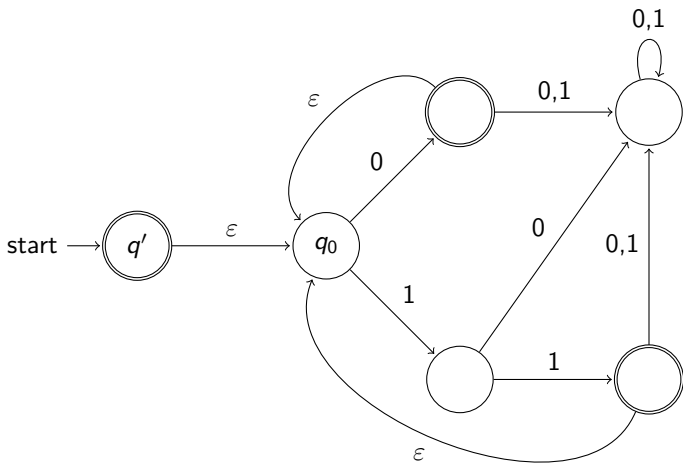
$$A = \{0, 11\}$$



正则语言的封闭性



正则语言的封闭性



正则语言的封闭性

存在DFA识别语言 A ，设其开始状态为 q_0 ，接受状态集合为 F
新增开始状态 q' ，增加 q' 到 q_0 的 ϵ 转移，并规定 q' 为接受状态
对于原来的每个接受状态 $q \in F$ ，增加 q 到 q_0 的 ϵ 转移

正则表达式

正则表达式	语言
a	$\{a\}$
ε	$\{\varepsilon\}$
\emptyset	\emptyset
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
(R_1^*)	$L(R_1)^*$

其中 R_1, R_2 为正则表达式, $L(R)$ 表示正则表达式 R 对应的语言。

正则表达式

设 $\Sigma = \{0, 1\}$

0^*10^* : 恰好包含一个1的串

$\Sigma^*1\Sigma^*$: 至少包含一个1的串

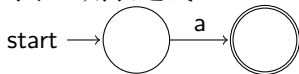
$(\Sigma\Sigma\Sigma)^*$: 长度为3的倍数的串

正则表达式与DFA的等价性

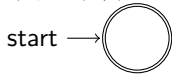
证明：一个语言是正则的当且仅当它能被某个正则表达式描述。

正则表达式与DFA的等价性

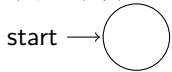
对于正则表达式 a :



对于正则表达式 ϵ :



对于正则表达式 \emptyset :



另三种情况：参考正则语言封闭性的证明

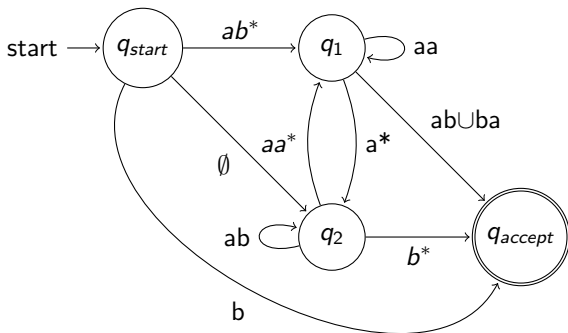
正则表达式与DFA的等价性

广义非确定有限状态自动机(Generalized Nondeterministic Finite Automaton)

$(Q, \Sigma, \delta, q_{start}, q_{accept})$

转移函数 $\delta : (Q - \{q_{accept}\}) \times (Q - \{q_{start}\}) \rightarrow \mathcal{R}$

正则表达式与DFA的等价性



接受串abbabb

$q_{start} \xrightarrow{abb} q_1 \xrightarrow{\epsilon} q_2 \xrightarrow{ab} q_2 \xrightarrow{b} q_{accept}$

正则表达式与DFA的等价性

若存在 $w_1, w_2, \dots, w_n, r_0, r_1, \dots, r_n$ 使得:

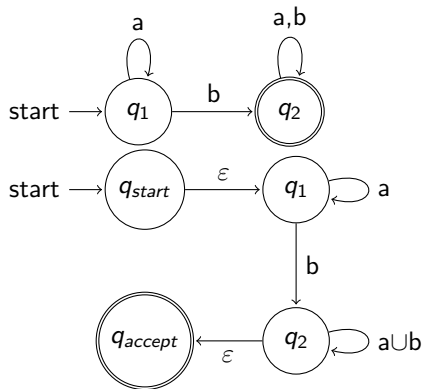
$$s = w_1 w_2 \dots w_n, w_i \in \Sigma^*$$

$$r_0 = q_{start}, r_n = q_{accept}, w_i \in L(\delta(r_{i-1}, r_i))$$

则GNFA接受串 s , 否则GNFA拒绝串 s

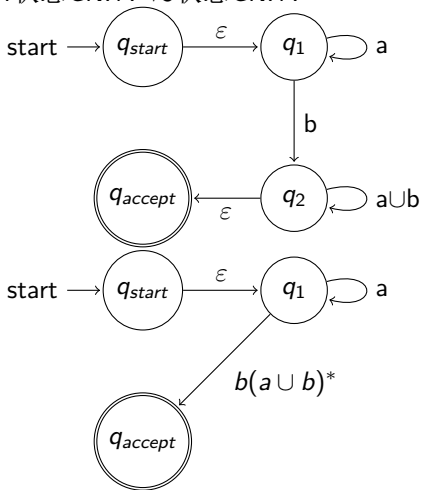
正则表达式与DFA的等价性

DFA \rightarrow GNFA



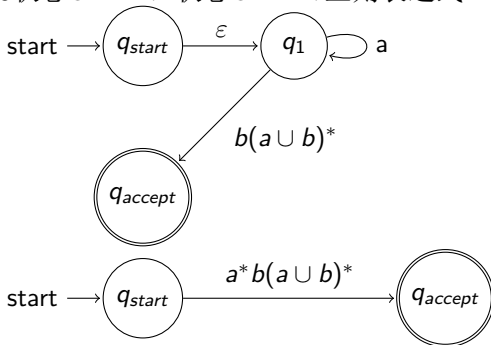
正则表达式与DFA的等价性

4状态GNFA \rightarrow 3状态GNFA



正则表达式与DFA的等价性

3状态GNFA \rightarrow 2状态GNFA \rightarrow 正则表达式



有限状态自动机的局限性

$\Sigma = \{0, 1\}$, 是否存在DFA识别: 0与1的个数相等的串?

有限状态自动机的局限性

$\Sigma = \{0, 1\}$, 是否存在DFA识别: 子串01与10出现次数相等的串?

有限状态自动机的局限性

严格证明：设存在一个包含 k 个状态的DFA识别该语言，考虑 $0^1, 0^2, \dots, 0^k, 0^{k+1}$

泵引理

正则语言的泵引理(pumping lemma for regular languages)

对于任意正则语言 A , 存在一个正整数 p 使得: 语言 A 中任意长度大于等于 p 的串 s , 都可以分拆成 $s = xyz$, 并且满足:

1. $xy^iz \in A$, 对于 $i = 0, 1, 2, \dots$
2. $y \neq \varepsilon$
3. $|xy| \leq p$

泵引理

例: $\Sigma = \{0, 1\}$, $A = \{x \mid x \text{ 含有恰好 } 2 \text{ 个 } 1\}$

取 $p = 3$, 将前三位中的0作为 $s = xyz$ 中的 y

例如对于 $100100 \in A$, 令 $x = 1, y = 0, z = 0100$

泵引理

使用泵引理证明以下语言不是正则语言：

$$A = \{0^n 1^n \mid n = 0, 1, 2, \dots\}$$

$$B = \{x \mid x \text{ 中 } 0 \text{ 与 } 1 \text{ 的个数相等}\}$$

$$C = \{ww \mid w \in \{0, 1\}^*\}$$

泵引理

设正则语言 A 可以被一个 k 状态DFA识别, 取 $p = k$

设DFA的开始状态为 q_0 , 转移函数为 δ

对于任意 $s \in A$ 且 $|s| = n \geq p$

令 $r_0 = q_0, r_i = \delta(r_{i-1}, s_i)(i = 1, 2, \dots, p)$

由抽屉原理, 必然存在 $0 \leq a < b \leq p$ 使得 $r_a = r_b$

令 $x = s_1 s_2 \cdots s_a, y = s_{a+1} s_{a+2} \cdots s_b, z = s_{b+1} s_{b+2} \cdots s_n$

从状态 $r_0 = q_0$ 出发, 输入串 x 会转移到 r_a

从状态 r_a 出发, 输入串 y 会转移回到 r_a

从状态 r_a 出发, 输入串 z 会转移到某个接受状态

思考题

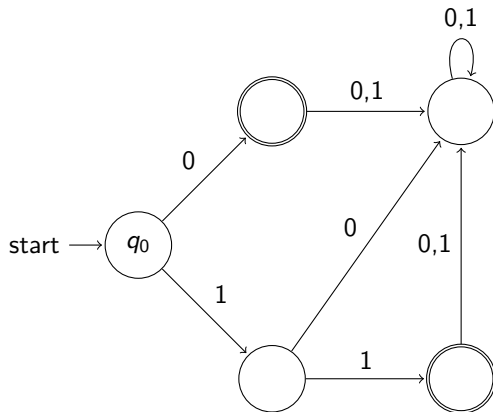
对于字母表 Σ 上的语言 A ，定义语言 $S(A) = \{xy | x \in \Sigma^*, y \in A\}$
证明：若 A 是正则语言，则 $S(A)$ 是正则语言

思考题

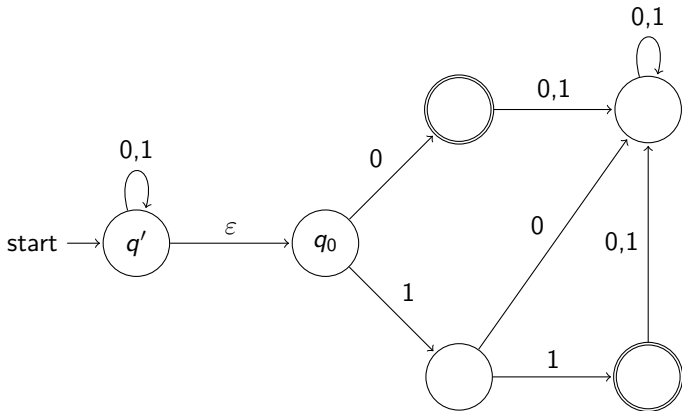
$$R \rightarrow \Sigma^* R$$

思考题

$$A = \{0, 11\}$$



思考题



例：数字搜索

来源：CTSC2004

$\Sigma = \{0, 1, \dots, 9\}$ ，给出一个正则表达式 R 和一段文本 s 。求：对于哪些 x ，存在 y 使得子串 $s[y..x]$ 被 R 匹配。

思考题

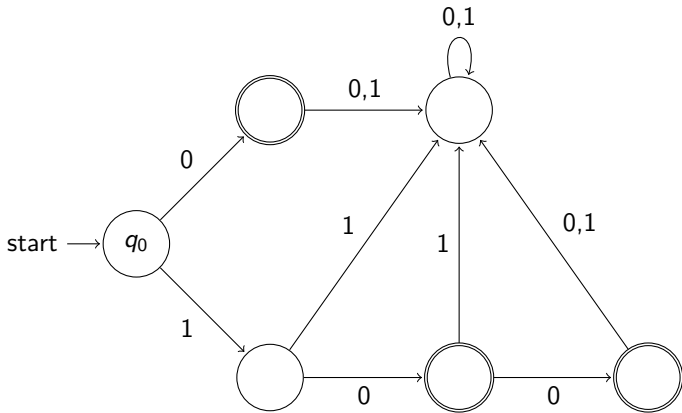
对于串 $s = s_1 s_2 \dots s_n$, 定义 $s^{\mathcal{R}} = s_n s_{n-1} \dots s_1$

对于语言 A , 定义 $A^{\mathcal{R}} = \{x^{\mathcal{R}} | x \in A\}$

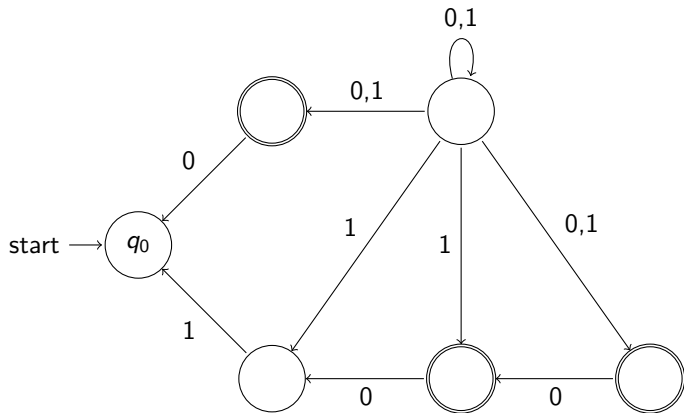
证明: 若 A 是正则语言, 则 $A^{\mathcal{R}}$ 是正则语言

思考题

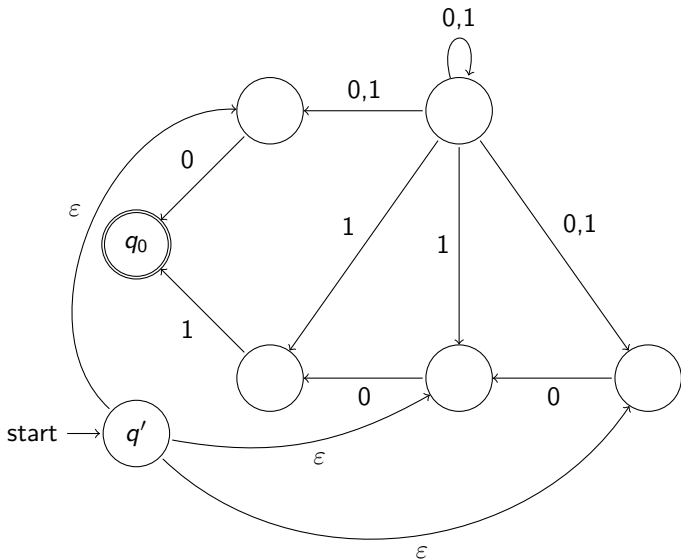
$$A = \{0, 10, 100\}$$



思考题



思考题



思考题

对于语言 $A = \{a^i b^j c^k \mid i, j, k \geq 0 \text{ 且若 } i = 1 \text{ 则 } j = k\}$

1. A在泵引理中表现得“和正则语言一样”
2. A不是正则语言

例：341在以2为底的费马小定理素性测试中表现得“和素数一样”，因为 $2^{341-1} \equiv 1 \pmod{341}$

思考题

令 $p = 3$

$$c^n = \varepsilon + c + c^{n-1}$$

$$b^n c^m = \varepsilon + b + b^{n-1} c^m (n > 1)$$

$$ab^n c^n = \varepsilon + a + b^n c^n$$

$$aab^n c^m = \varepsilon + aa + b^n c^m$$

$$aaaa^n b^m c^l = aa + a + a^n b^m c^l$$

思考题

$$A = \{a^i b^j c^k \mid i, j, k \geq 0 \text{ 且若 } i = 1 \text{ 则 } j = k\}$$

$$B = \{ab^i c^j \mid i, j \geq 0\}$$

$$C = \{ab^i c^i \mid i \geq 0\} = A \cap B$$

反设 A 为正则语言

因为 $B = L(ab^*c^*)$ 是正则语言，所以 $C = A \cap B$ 为正则语言

使用泵引理可以证明 C 不是正则语言（考虑串 $ab^p c^p$ ）

思考题

构造一个正则表达式描述 $\Sigma = \{0, 1, 2\}$ 上的语言：
 $\{x \mid x \text{ 在十进制中被 } 3 \text{ 整除}\}$

思考题

构造一个正则表达式描述 $\Sigma = \{0, 1\}$ 上的语言：
 $\{x \mid x \text{ 不包含子串 } 110\}$

总结

DFA、正则语言、正则表达式的定义

正则语言的运算封闭性

DFA、NFA、正则表达式的等价性

DFA的局限性

总结

Manufactoria
regex.alf.nu