

形式语言与自动机理论

课程简介与基础知识

王春宇

计算机科学与技术学院
哈尔滨工业大学

课程简介与基础知识

- 课程简介
- 基础知识

核心问题

计算机的基本能力和限制是什么？

- ① 究竟哪些问题, 可通过计算解决? — 可计算性理论
- ② 解决可计算的问题, 究竟需要多少资源? — 计算复杂性理论
- ③ 为了研究计算, 要使用哪些计算模型? — 形式语言与自动机理论

什么是自动机理论?

自动机理论: 研究抽象机器及其所能解决问题的理论.

- 图灵机
- 有限状态机
- 文法, 下推自动机

什么是形式语言?

形式语言: 经数学定义的语言.

语言	自然语言		形式语言		
	English	中文	化学分子式	C 语言	
	字符	A,a,B,b,...	天, 地,...	A-Z,a-z,0-9...	A-Z,a-z,0-9...
	单词	apple	苹果	H ₂ O	char
	句子	How're you?	早上好!	2H ₂ +O ₂ =2H ₂ O	char a = 10;
	语法	Grammar	语法规则	精确定义的规则	

课程内容

- 正则语言
 - 有穷自动机
 - 正则表达式
 - 正则语言的性质
- 上下文无关语言
 - 上下文无关文法
 - 下推自动机
 - 上下文无关语言的性质
- 计算导论
 - 图灵机及其扩展
 - 不可判定性

参考书

- John E. Hopcroft, Rajeev Motwani, Jeffrey D. Ullman.
Introduction to Automata Theory, Languages, and Computation.
《自动机理论、语言和计算导论》机械工业出版社
- Michael Sipser. *Introduction to the Theory of Computation.*
《计算理论导引》机械工业出版社

课程简介与基础知识

- 课程简介
- 基础知识
 - 基本概念
 - 语言和问题
 - 形式化证明

基本概念

1. **字母表**: 符号 (或字符) 的非空有穷集.

$$\Sigma_1 = \{0, 1\},$$

$$\Sigma_2 = \{a, b, \dots, z\},$$

$$\Sigma_3 = \{x \mid x \text{ 是一个汉字}\}.$$

2. **字符串**: 由某字母表中符号组成的有穷序列.

若 $\Sigma_1 = \{0, 1\}$, 那么 $0, 1, 00, 111001$ 为 Σ_1 上的字符串;

若 $\Sigma_2 = \{a, b, \dots, z\}$, 那么 $ab, xkcd$ 为 Σ_2 上的字符串.

3. **空串**: 记为 ε , 有 0 个字符的串.

字母表 Σ 可以是任意的, 但都有 $\varepsilon \notin \Sigma$.

4. 字符串的**长度**: 字符串中符号所占位置的个数, 记为 $|\square|$.
若字母表为 Σ , 可**递归定义**为:

$$|w| = \begin{cases} 0 & w = \varepsilon \\ |x| + 1 & w = xa \end{cases},$$

其中 $a \in \Sigma$, w 和 x 是 Σ 中字符组成的字符串.

★. 符号使用的一般约定:

- 字母表: $\Sigma, \Gamma, \Delta, \dots$
- 字符: a, b, c, \dots
- 字符串: \dots, w, x, y, z
- 集合: A, B, C, \dots

5. 字符串 x 和 y 的**连接**: 将首尾相接得到新字符串的运算, 记为 $x \cdot y$ 或 xy .
同样, 可递归定义为

$$x \cdot y = \begin{cases} x & y = \varepsilon \\ (x \cdot z)a & y = za \end{cases},$$

其中 $a \in \Sigma$, 且 x, y, z 都是字符串.

对任何字符串 x , 有 $\varepsilon \cdot x = x \cdot \varepsilon = x$.

连接运算的符号 “ \cdot ” 一般省略.

6. 字符串 x 的 n 次幂($n \geq 0$), 递归定义为

$$x^n = \begin{cases} \varepsilon & n = 0 \\ x^{n-1}x & n > 0 \end{cases} .$$

例如,

$$\begin{aligned} (ba)^2 &= (ba)^1ba \\ &= (ba)^0baba \\ &= \varepsilon baba \\ &= baba \end{aligned}$$

$$\begin{aligned} ba^2 &= ba^1a \\ &= ba^0aa \\ &= b\varepsilon aa \\ &= baa \end{aligned}$$

7. 集合 A 和 B 的**连接**, 记为 $A \cdot B$ 或 AB , 定义为

$$A \cdot B = \{w \mid w = x \cdot y, x \in A \text{ 且 } y \in B\}.$$

8. 集合 A 的 n 次幂($n \geq 0$), 递归定义为

$$A^n = \begin{cases} \{\varepsilon\} & n = 0 \\ A^{n-1}A & n \geq 1 \end{cases} .$$

那么, 若 Σ 为字母表, 则 Σ^n 为 Σ 上长度为 n 的字符串集合.
如果 $\Sigma = \{0, 1\}$, 有

$$\Sigma^0 = \{\varepsilon\}$$

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

$$\Sigma^2 = \{00, 01, 10, 11\}$$

$$\Sigma^3 = \{000, 001, 010, 011, 100, 101, 110, 111\}$$

\vdots

9. 克林闭包(*Kleene Closure*):

$$\Sigma^* = \bigcup_{i=0}^{\infty} \Sigma^i.$$

10. 正闭包(*Positive Closure*):

$$\Sigma^+ = \bigcup_{i=1}^{\infty} \Sigma^i.$$

显然,

$$\Sigma^* = \Sigma^+ \cup \{\varepsilon\}.$$

其他的概念如有向图, 树, 字符串的前缀, 后缀等定义这里省略.

语言

定义

若 Σ 为字母表且 $\forall L \subseteq \Sigma^*$, 则 L 称为字母表 Σ 上的语言.

- 自然语言, 程序设计语言等
- $\{0^n 1^n \mid n \geq 0\}$
- The set of strings of 0's and 1's with an equal number of each:

$$\{\epsilon, 01, 10, 0011, 0101, 1100, \dots\}$$

- \emptyset , $\{\epsilon\}$ 和 Σ^* 分别都是任意字母表 Σ 上的语言, 但注意 $\emptyset \neq \{\epsilon\}$

关于语言

唯一重要的约束就是所有字母表都是有穷的.

问题

典型问题

判断给定的字符串 w 是否属于某个具体的语言 L ,

$$w \in L?$$

- 任何所谓问题, 都可以转为语言成员性的问题
- 语言和问题其实是相同的东西

形式化证明: 演绎法, 归纳法和反证法

例 1. 若 x 和 y 是 Σ 上的字符串, 请证明 $|xy| = |x| + |y|$.

证明: 通过对 $|y|$ 的归纳来证明

① 基础: 当 $|y| = 0$, 即 $y = \varepsilon$

$$\begin{aligned}|x\varepsilon| &= |x| \\ &= |x| + |\varepsilon|\end{aligned}$$

连接的定义
长度的定义

② 递推: 假设 $|y| = n$ ($n \geq 0$) 时命题成立,
那么当 $|y| = n + 1$, 即 $y = wa$

$$\begin{aligned}|x(wa)| &= |(xw)a| \\ &= |xw| + 1 \\ &= |x| + |w| + 1 \\ &= |x| + |wa|\end{aligned}$$

连接的定义
长度的定义
归纳假设
长度的定义



形式化证明: 演绎法, 归纳法和反证法

例 1. 若 x 和 y 是 Σ 上的字符串, 请证明 $|xy| = |x| + |y|$.

证明: 通过对 y 的结构归纳来证明

① 基础: $y = \varepsilon$ 时

$$\begin{aligned}|x\varepsilon| &= |x| \\ &= |x| + |\varepsilon|\end{aligned}$$

连接的定义
长度的定义

② 递推: 假设 $y = w$ ($w \in \Sigma^*$) 时命题成立,
那么当 $y = wa$ 时

$$\begin{aligned}|x(wa)| &= |(xw)a| \\ &= |xw| + 1 \\ &= |x| + |w| + 1 \\ &= |x| + |wa|\end{aligned}$$

连接的定义
长度的定义
归纳假设
长度的定义



形式语言与自动机理论

有穷自动机

王春宇

计算机科学与技术学院
哈尔滨工业大学

有穷自动机

- 有穷状态系统
- 确定的有穷自动机
- 非确定有穷自动机
- 带有空转移的非确定有穷自动机

有穷状态系统

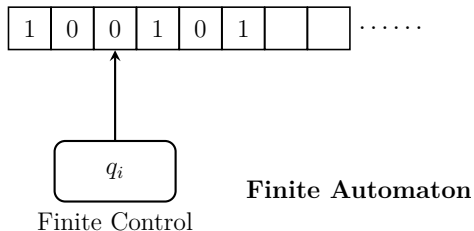
- 有限状态机: Moore Machine, Mealy Machine
- 数字电路设计
- 电脑游戏的 AI 设计
- 各种通讯协议: TCP, HTTP, Bluetooth, Wifi
- 文本搜索, 词法分析

有穷自动机

- 有穷状态系统
- 确定的有穷自动机
 - 形式定义
 - DFA 的设计举例
 - 扩展转移函数
 - DFA 的语言与正则语言
- 非确定有穷自动机
- 带有空转移的非确定有穷自动机

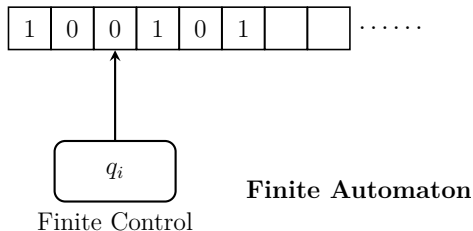
确定的有穷自动机

- 一条输入带
- 一个读头
- 一个有穷控制器



确定的有穷自动机

- 一条输入带
- 一个读头
- 一个有穷控制器



例 1. 用有穷自动机识别 $\{w \in \{0, 1\}^* \mid w \text{ 的长度 } |w| \text{ 是偶数.}\}$

确定的有穷自动机的形式定义

定义

确定的有穷自动机(*DFA*, *Deterministic Finite Automaton*) A 为五元组

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

- ① Q : 有穷状态集;
- ② Σ : 有穷输入符号集或字母表;
- ③ $\delta : Q \times \Sigma \rightarrow Q$, 状态转移函数;
- ④ $q_0 \in Q$: 初始状态;
- ⑤ $F \subseteq Q$: 终结状态集或接受状态集.

例 2. 请设计 DFA, 在任何由 0 和 1 构成的串中, 接受含有 01 子串的全部串.

例 2. 请设计 DFA, 在任何由 0 和 1 构成的串中, 接受含有 01 子串的全部串.

- ① 未发现 01, 即使 0 都还没出现过;
- ② 未发现 01, 但刚刚读入字符是 0;
- ③ 已经发现了 01.

例 2. 请设计 DFA, 在任何由 0 和 1 构成的串中, 接受含有 01 子串的全部串.

- ❶ 未发现 01, 即使 0 都还没出现过;
- ❷ 未发现 01, 但刚刚读入字符是 0;
- ❸ 已经发现了 01.

因此 DFA A 的可定义为:

$$A = (\{q_1, q_2, q_3\}, \{0, 1\}, \delta, q_1, \{q_3\})$$

其中 δ 为:

$$\delta(q_1, 1) = q_1$$

$$\delta(q_2, 1) = q_3$$

$$\delta(q_3, 1) = q_3$$

$$\delta(q_1, 0) = q_2$$

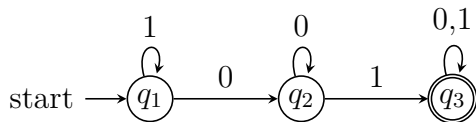
$$\delta(q_2, 0) = q_2$$

$$\delta(q_3, 0) = q_3$$

状态转移图

- ① 每个状态 q 对应一个节点, 用圆圈表示;
- ② 状态转移 $\delta(q, a) = p$ 为一条从 q 到 p 且标记为字符 a 的有向边;
- ③ 开始状态 q_0 用一个标有 start 的箭头表示;
- ④ 接受状态的节点, 用双圆圈表示.

续例 2. 含有 01 子串的全部串的状态转移图



状态转移表

- ① 每个状态 q 对应一行, 每个字符 a 对应一列;
- ② 若有 $\delta(q, a) = p$, 用第 q 行第 a 列中填入的 p 表示;
- ③ 开始状态 q_0 前, 标记箭头 \rightarrow 表示;
- ④ 接受状态 $q \in F$ 前, 标记星号 $*$ 表示.

续例 2. 含有 01 子串的全部串的状态转移表

	0	1
$\rightarrow q_1$	q_2	q_1
q_2	q_2	q_3
$*q_3$	q_3	q_3

典型问题

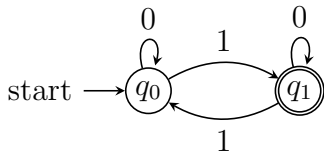
设计 DFA 使其接受且仅接受给定的语言 L .

例 3. 若 $\Sigma = \{0, 1\}$, 给出接受全部含有奇数个 1 的串 DFA.

典型问题

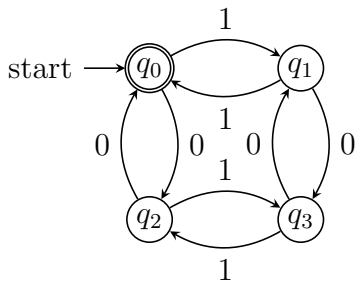
设计 DFA 使其接受且仅接受给定的语言 L .

例 3. 若 $\Sigma = \{0, 1\}$, 给出接受全部含有奇数个 1 的串 DFA.



例 4. 若 $\Sigma = \{0, 1\}$, 给出接受全部含有偶数个 0 和偶数个 1 的串 DFA.

例 4. 若 $\Sigma = \{0, 1\}$, 给出接受全部含有偶数个 0 和偶数个 1 的串 DFA.



思考题

若 $\Sigma = \{0, 1\}$

- ① 如何设计接受 \emptyset 的 DFA?
- ② 如何设计接受 Σ^* 的 DFA?
- ③ 如何设计接受 $\{\varepsilon\}$ 的 DFA?

扩展转移函数

定义

扩展 δ 到字符串, 定义扩展转移函数 $\hat{\delta} : Q \times \Sigma^* \rightarrow Q$ 为

$$\hat{\delta}(q, w) = \begin{cases} q & w = \varepsilon \\ \delta(\hat{\delta}(q, x), a) & w = xa \end{cases}$$

其中 $a \in \Sigma$, $w, x \in \Sigma^*$.

扩展转移函数

定义

扩展 δ 到字符串, 定义 **扩展转移函数** $\hat{\delta}: Q \times \Sigma^* \rightarrow Q$ 为

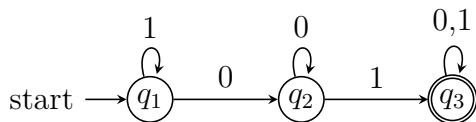
$$\hat{\delta}(q, w) = \begin{cases} q & w = \varepsilon \\ \delta(\hat{\delta}(q, x), a) & w = xa \end{cases}$$

其中 $a \in \Sigma$, $w, x \in \Sigma^*$.

那么, 当 $w = a_0a_1 \cdots a_n$, 则有

$$\begin{aligned} \hat{\delta}(q, w) &= \delta(\hat{\delta}(q, a_0a_1 \cdots a_{n-1}), a_n) \\ &= \delta(\delta(\hat{\delta}(q, a_0a_1 \cdots a_{n-2}), a_{n-1}), a_n) = \cdots \\ &= \delta(\delta(\cdots \delta(\hat{\delta}(q, \varepsilon), a_0) \cdots, a_{n-1}), a_n) \end{aligned}$$

续例 2. 接受全部含有 01 子串的 DFA, $\hat{\delta}$ 处理串 0101 的过程.



$$\begin{aligned}\hat{\delta}(q_0, 0101) &= \delta(\hat{\delta}(q_0, 010), 1) \\ &= \delta(\delta(\hat{\delta}(q_0, 01), 0), 1) \\ &= \delta(\delta(\delta(\hat{\delta}(q_0, 0), 1), 0), 1) \\ &= \delta(\delta(\delta(\delta(\hat{\delta}(q_0, \varepsilon), 0), 1), 0), 1) \\ &= \delta(\delta(\delta(\delta(q_0, 0), 1), 0), 1) \\ &= \delta(\delta(\delta(q_1, 1), 0), 1) \\ &= \delta(\delta(q_2, 0), 1) = \delta(q_2, 1) = q_2\end{aligned}$$

思考题

- ① 扩展转移函数 $\hat{\delta}$ 必须从开始状态 q_0 处理字符串吗?
- ② 对任意的串 w , $\hat{\delta}$ 能保证一定会跳转到某个状态吗?

例5. 对任何状态 q 及字符串 x 和 y , 证明 $\hat{\delta}(q, xy) = \hat{\delta}(\hat{\delta}(q, x), y)$.

证明: 对 y 使用归纳法.

❶ 当 $y = \varepsilon$ 时

$$\begin{aligned}\hat{\delta}(\hat{\delta}(q, x), \varepsilon) &= \hat{\delta}(q, x) && \hat{\delta} \text{ 的定义} \\ &= \hat{\delta}(q, x\varepsilon)\end{aligned}$$

❷ 假设 $y = w$ ($w \in \Sigma^*$) 时命题成立, 当 $y = wa$ ($a \in \Sigma$) 时

$$\begin{aligned}\hat{\delta}(q, xwa) &= \delta(\hat{\delta}(q, xw), a) && \hat{\delta} \text{ 和 连接的 定义} \\ &= \delta(\hat{\delta}(\hat{\delta}(q, x), w), a) && \text{归纳假设} \\ &= \hat{\delta}(\hat{\delta}(q, x), wa) && \hat{\delta} \text{ 的定义}\end{aligned}$$



DFA 的语言与正则语言

定义

若 $D = (Q, \Sigma, \delta, q_0, F)$ 是一个 **DFA**, 则 D **接受的语言**为

$$\mathbf{L}(D) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \in F\}.$$

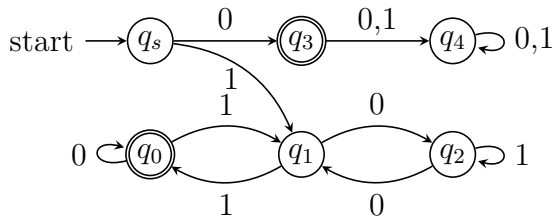
定义

如果语言 L 是某个 DFA D 的语言, 即 $L = \mathbf{L}(D)$, 则称 L 是**正则语言**.

- $\emptyset, \{\varepsilon\}$ 都是正则语言
- 若 Σ 是字母表, Σ^*, Σ^n 都是 Σ 上的正则语言

例 6. 设计 DFA 接受 $\{0,1\}$ 上的字符串 w , 且 w 是 3 的倍数的二进制表示.

例 6. 设计 DFA 接受 $\{0,1\}$ 上的字符串 w , 且 w 是 3 的倍数的二进制表示.

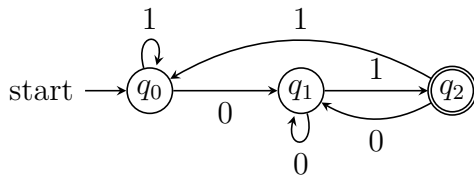


有穷自动机

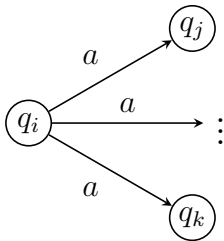
- 有穷状态系统
- 确定的有穷自动机
- 非确定有穷自动机
 - 形式定义
 - 扩展转移函数
 - NFA 的语言
 - DFA 与 NFA 的等价性
- 带有空转移的非确定有穷自动机

例 7. 由 0 和 1 构成的串中, 接受全部以 01 结尾的串, 如何设计 DFA?

例 7. 由 0 和 1 构成的串中, 接受全部以 01 结尾的串, 如何设计 DFA?

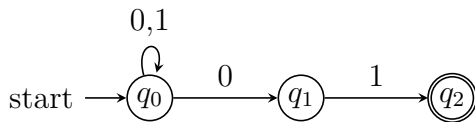
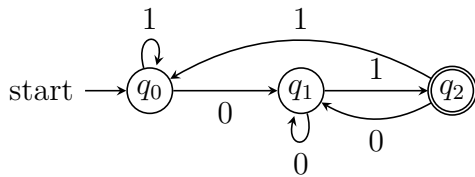


状态的非确定转移



- 同一个状态在相同的输入下,可以有多个转移状态
- 自动机可以处在多个当前状态
- 使自动机的设计更容易

续例 7. 由 0 和 1 构成的串中, 接受全部以 01 结尾的串.



思考题

有穷自动机有了非确定性, 能否增加它识别语言的能力?

非确定有穷自动机的形式定义

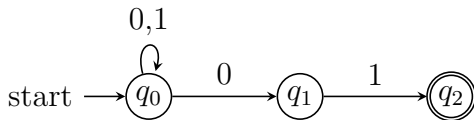
定义

非确定有穷自动机(*NFA*, *Nondeterministic Finite Automaton*) A 为五元组

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

- ① Q : 有穷状态集;
- ② Σ : 有穷输入符号集或字母表;
- ③ $\delta : Q \times \Sigma \rightarrow 2^Q$ 状态转移函数;
- ④ $q_0 \in Q$: 为初始状态;
- ⑤ $F \subseteq Q$: 为终结状态集或接受状态集.

续例7. 接受全部以 01 结尾的串的 NFA.



五元组为 $A = (\{q_0, q_1, q_2\}, \{0, 1\}, \delta, q_0, \{q_2\})$, 转移函数 δ :

$$\delta(q_0, 0) = \{q_0, q_1\}$$

$$\delta(q_1, 0) = \emptyset$$

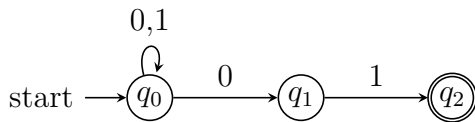
$$\delta(q_2, 0) = \emptyset$$

$$\delta(q_0, 1) = \{q_0\}$$

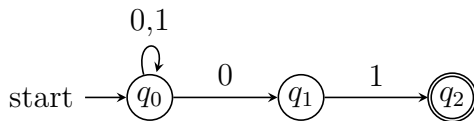
$$\delta(q_1, 1) = \{q_2\}$$

$$\delta(q_2, 1) = \emptyset$$

续例 7. 接受全部以 01 结尾的串的 NFA. 识别字符串 00101 的过程.



续例 7. 接受全部以 01 结尾的串的 NFA.



状态转移表:

	0	1
$\rightarrow q_0$	$\{q_0, q_1\}$	$\{q_0\}$
q_1	\emptyset	$\{q_2\}$
$*q_2$	\emptyset	\emptyset

扩展转移函数

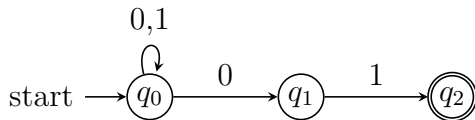
定义

扩展 δ 到字符串, 定义扩展转移函数 $\hat{\delta} : Q \times \Sigma^* \rightarrow 2^Q$ 为

$$\hat{\delta}(q, w) = \begin{cases} \{q\} & w = \varepsilon \\ \bigcup_{p \in \hat{\delta}(q, x)} \delta(p, a) & w = xa \end{cases}$$

其中 $a \in \Sigma$, $w, x \in \Sigma^*$.

续例7. 接受 01 结尾的串的 NFA, $\hat{\delta}$ 处理 00101 时每步的状态转移.



- ❶ $\hat{\delta}(q_0, \varepsilon) = \{q_0\}$
- ❷ $\hat{\delta}(q_0, 0) = \delta(q_0, 0) = \{q_0, q_1\}$
- ❸ $\hat{\delta}(q_0, 00) = \delta(q_0, 0) \cup \delta(q_1, 0) = \{q_0, q_1\} \cup \emptyset = \{q_0, q_1\}$
- ❹ $\hat{\delta}(q_0, 001) = \delta(q_0, 1) \cup \delta(q_1, 1) = \{q_0\} \cup \{q_2\} = \{q_0, q_2\}$
- ❺ $\hat{\delta}(q_0, 0010) = \delta(q_0, 0) \cup \delta(q_2, 0) = \{q_0, q_1\} \cup \emptyset = \{q_0, q_1\}$
- ❻ $\hat{\delta}(q_0, 00101) = \delta(q_0, 1) \cup \delta(q_1, 1) = \{q_0\} \cup \{q_2\} = \{q_0, q_2\}$

因为 q_2 是接受状态, 所以 NFA 接受 00101.

NFA 的语言

回顾

若 $D = (Q, \Sigma, \delta, q_0, F)$ 是一个 DFA, 则 D 接受的语言为

$$\mathbf{L}(D) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \in F\}.$$

定义

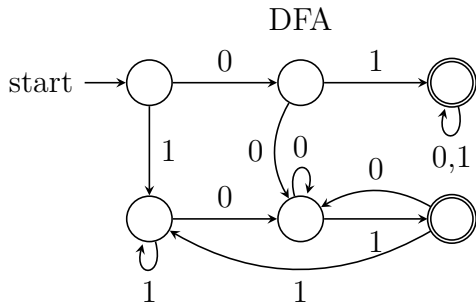
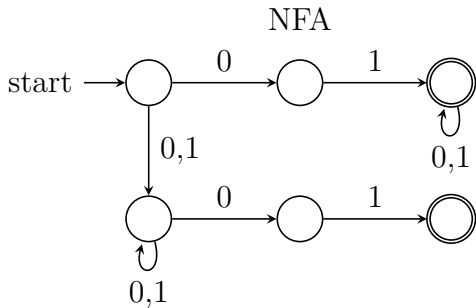
若 $N = (Q, \Sigma, \delta, q_0, F)$ 是一个 **NFA**, 则 N **接受的语言**为

$$\mathbf{L}(N) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \cap F \neq \emptyset\}.$$

例 8. 设计 $L = \{w \in \{0, 1\}^* \mid w \text{ 的首尾字符相同.}\}$ 的 NFA.

例 9. 设计 $L = \{w \in \{0, 1\}^* \mid w \text{ either begin or ends with } 01.\}$ 的 NFA.

例 9. 设计 $L = \{w \in \{0,1\}^* \mid w \text{ either begin or ends with } 01.\}$ 的 NFA.



DFA 与 NFA 的等价性

定理 1

如果语言 L 被 NFA 接受, 当且仅当 L 被 DFA 接受.

子集构造法

如果 NFA $N = (Q_N, \Sigma, \delta_N, q_0, F_N)$ 构造 DFA

$$D = (Q_D, \Sigma, \delta_D, \{q_0\}, F_D)$$

- ① $Q_D = 2^{Q_N}$;
- ② $F_D = \{S \mid S \subseteq Q_N, S \cap F_N \neq \emptyset\}$;
- ③ $\forall S \subseteq Q_N, \forall a \in \Sigma$:

$$\delta_D(S, a) = \bigcup_{p \in S} \delta_N(p, a).$$

那么有 $L(D) = L(N)$.

证明: 为证明 $L(D) = L(N)$, 对 $|w|$ 用归纳法, 往证

$$\hat{\delta}_D(\{q_0\}, w) = \hat{\delta}_N(q_0, w).$$

① 归纳基础: 当 $w = \varepsilon$ 时,

$$\hat{\delta}_D(\{q_0\}, \varepsilon) = \{q_0\} = \hat{\delta}_N(q_0, \varepsilon)$$

② 归纳递推: 假设 $w = x$ ($x \in \Sigma^*$) 时成立, 当 $w = xa$ ($a \in \Sigma$) 时

$$\begin{aligned}\hat{\delta}_N(q_0, xa) &= \cup_{p \in \hat{\delta}_N(q_0, x)} \delta_N(p, a) && \text{NFA 的 } \hat{\delta} \text{ 定义} \\ &= \cup_{p \in \hat{\delta}_D(\{q_0\}, x)} \delta_N(p, a) && \text{归纳假设} \\ &= \delta_D(\hat{\delta}_D(\{q_0\}, x), a) && D \text{ 的构造} \\ &= \hat{\delta}_D(\{q_0\}, xa) && \text{DFA 的 } \hat{\delta} \text{ 定义}\end{aligned}$$

因此上式成立.

因为

$$\hat{\delta}_D(\{q_0\}, w) = \hat{\delta}_N(q_0, w)$$

所以, 对 $\forall w \in \Sigma^*$ 有

$$w \in \mathbf{L}(N) \iff \hat{\delta}_N(q_0, w) \cap F_N \neq \emptyset$$

NFA 的语言

$$\iff \hat{\delta}_D(\{q_0\}, w) \cap F_N \neq \emptyset$$

刚证明的

$$\iff \hat{\delta}_D(\{q_0\}, w) \in F_D$$

D 的构造

$$\iff w \in \mathbf{L}(D)$$

DFA 的语言

所以

$$\mathbf{L}(D) = \mathbf{L}(N).$$

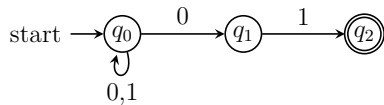


思考题

非确定性没能增加有穷自动机识别语言的能力, 原因是什么呢?

子集构造法: 构造与 NFA 等价的 DFA

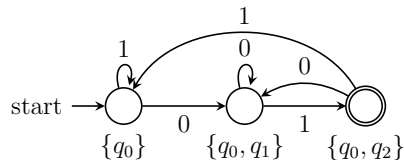
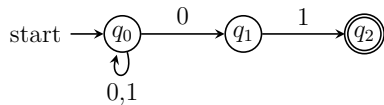
续例 7. 将接受全部以 01 结尾的串的 NFA 转换为 DFA.



	0	1
$\rightarrow q_0$	$\{q_0, q_1\}$	$\{q_0\}$
q_1	\emptyset	$\{q_2\}$
$*q_2$	\emptyset	\emptyset

子集构造法: 构造与 NFA 等价的 DFA

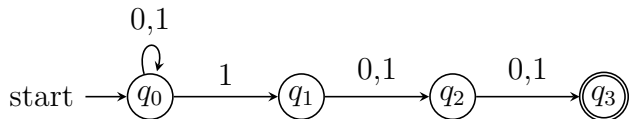
续例 7. 将接受全部以 01 结尾的串的 NFA 转换为 DFA.



	0	1
$\rightarrow q_0$	$\{q_0, q_1\}$	$\{q_0\}$
q_1	\emptyset	$\{q_2\}$
$*q_2$	\emptyset	\emptyset
$\rightarrow \{q_0\}$	$\{q_0, q_1\}$	$\{q_0\}$
$\{q_1\}$	\emptyset	$\{q_2\}$
$*\{q_2\}$	\emptyset	\emptyset
$\{q_0, q_1\}$	$\{q_0, q_1\}$	$\{q_0, q_2\}$
$*\{q_0, q_2\}$	$\{q_0, q_1\}$	$\{q_0\}$
\emptyset	\emptyset	\emptyset
$*\{q_1, q_2\}$	\emptyset	$\{q_2\}$
$*\{q_0, q_1, q_2\}$	$\{q_0, q_1\}$	$\{q_0, q_2\}$

例 10. 设计 NFA 识别 $L = \{w \in \{0, 1\}^* \mid w \text{ 倒数第 } 3 \text{ 个字符是 } 1\}$.

例 10. 设计 NFA 识别 $L = \{w \in \{0, 1\}^* \mid w \text{ 倒数第 } 3 \text{ 个字符是 } 1\}$.

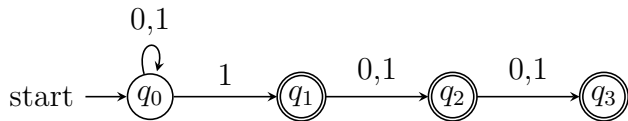


有穷自动机

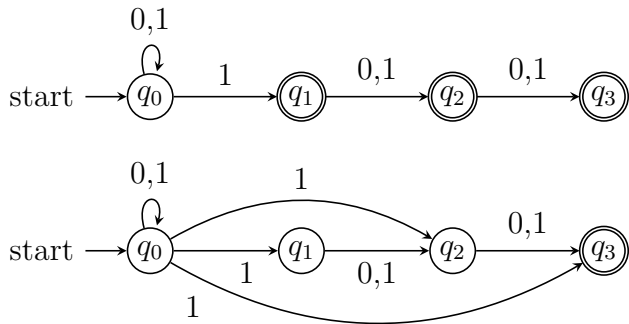
- 有穷状态系统
- 确定的有穷自动机
- 非确定有穷自动机
- 带有空转移的非确定有穷自动机
 - 形式定义
 - ε -闭包
 - 扩展转移函数
 - ε -NFA 的语言
 - ε -NFA 与 DFA 等价性

例 11. 设计 $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 NFA.

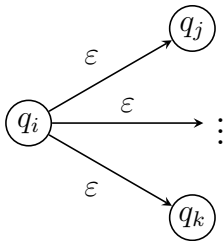
例 11. 设计 $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 NFA.



例 11. 设计 $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 3 个字符至少有一个是 } 1\}$ 的 NFA.

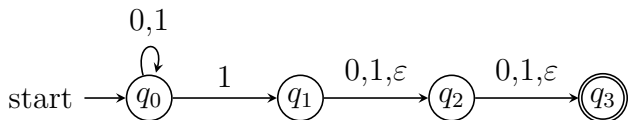
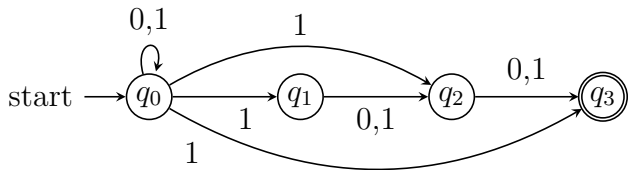
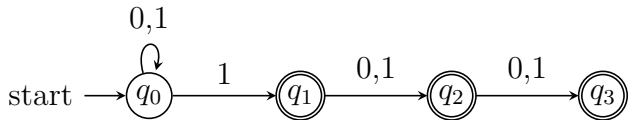


状态的 ε 转移



- 允许状态因空串 ε 而转移, 即不消耗输入字符就发生状态的改变
- 使自动机的设计更容易

续例 11. 设计 $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 NFA.



带空转移非确定有穷自动机的形式定义

定义

带空转移非确定有穷自动机(ϵ -NFA) A 为五元组

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

- ① Q : 有穷状态集;
- ② Σ : 有穷输入符号集或字母表;
- ③ $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$, 转移函数;
- ④ $q_0 \in Q$: 初始状态;
- ⑤ $F \subseteq Q$: 终结状态集或接受状态集.

ε -NFA, NFA, DFA 之间的主要区别

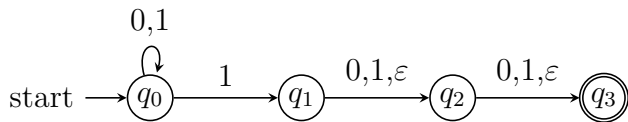
- ① 自动机在某状态, 读入某个字符时, 可能有多个转移;
- ② 自动机在某状态, 读入某个字符时, 可能没有转移;
- ③ 自动机在某状态, 可能不读入字符, 就进行转移.

注意

此后, 不再明确区分 ε -NFA 和 NFA, 而认为它们都是 NFA.

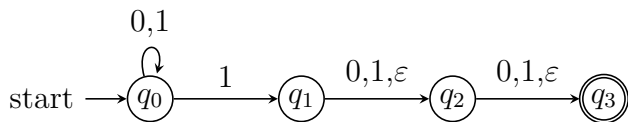
续例 11. $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 ε -NFA.

利用 ε 转移设计的有穷自动机:



状态转移表:

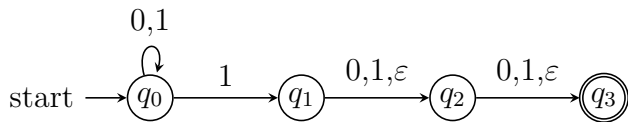
续例 11. $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 ε -NFA.
 利用 ε 转移设计的有穷自动机:



状态转移表:

	0	1	ε
$\rightarrow q_0$	$\{q_0\}$	$\{q_0, q_1\}$	\emptyset
q_1	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$
q_2	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$
$*q_3$	\emptyset	\emptyset	\emptyset

续例 11. $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 ε -NFA.
利用 ε 转移设计的有穷自动机:



当输入字符串是 011 时, ε -NFA 的状态变化.

思考题

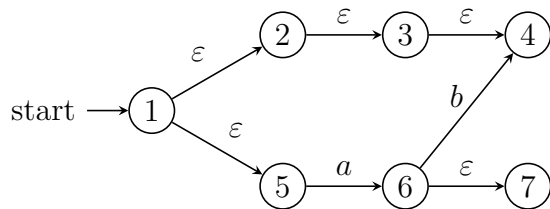
- ① 如果初始状态有 ε 转移, 第 1 个字符该如何处理?
- ② 如果最后的字符所到的状态有 ε 转移呢?

状态的 ε -闭包

定义

状态 q 的 ε -闭包(ε -Closure), 记为 $\text{ECLOSE}(q)$, 表示从 q 经过 ε 序列可达的全部状态集合, 递归定义为:

- ① $q \in \text{ECLOSE}(q)$;
- ② $\forall p \in \text{ECLOSE}(q)$, 若 $r \in \delta(p, \varepsilon)$, 则 $r \in \text{ECLOSE}(q)$.



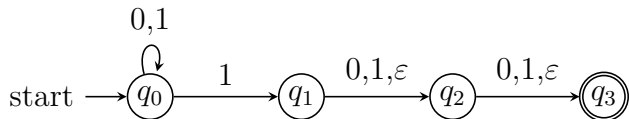
状态集合的 ε -闭包

定义

状态集 S 的 ε -闭包为

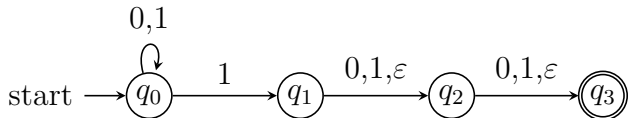
$$\text{ECLOSE}(S) = \bigcup_{q \in S} \text{ECLOSE}(q).$$

续例 11. $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 NFA.



状态转移表及每个状态的闭包:

续例 11. $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 NFA.



状态转移表及每个状态的闭包:

	0	1	ε	$\text{ECLOSE}(\sqcup)$
$\rightarrow q_0$	$\{q_0\}$	$\{q_0, q_1\}$	\emptyset	$\{q_0\}$
q_1	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$	$\{q_1, q_2, q_3\}$
q_2	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$	$\{q_2, q_3\}$
$*q_3$	\emptyset	\emptyset	\emptyset	$\{q_3\}$

扩展转移函数

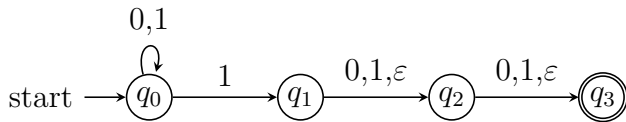
定义

扩展 δ 到字符串, 定义 **扩展转移函数** $\hat{\delta} : Q \times \Sigma^* \rightarrow 2^Q$ 为

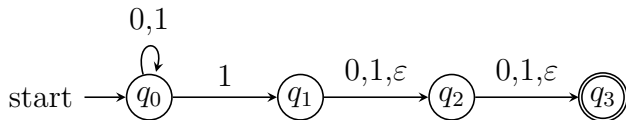
$$\hat{\delta}(q, w) = \begin{cases} \text{ECLOSE}(q) & w = \varepsilon \\ \text{ECLOSE}\left(\bigcup_{p \in \hat{\delta}(q, x)} \delta(p, a)\right) & w = xa \end{cases}$$

其中 $a \in \Sigma$, $w, x \in \Sigma^*$.

续例 11. 若 $L = \{w \in \{0, 1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 ε -NFA 如下, 求 $\hat{\delta}(q_0, 10)$.



续例 11. 若 $L = \{w \in \{0,1\}^* \mid w \text{ 倒数 } 3 \text{ 个字符至少有一个是 } 1\}$ 的 ε -NFA 如下, 求 $\hat{\delta}(q_0, 10)$.



$$\hat{\delta}(q_0, \varepsilon) = \text{ECLOSE}(q_0) = \{q_0\}$$

$$\begin{aligned} \hat{\delta}(q_0, 1) &= \text{ECLOSE}\left(\bigcup_{p \in \hat{\delta}(q_0, \varepsilon)} \delta(p, 1)\right) \\ &= \text{ECLOSE}(\hat{\delta}(q_0, 1)) = \text{ECLOSE}(\{q_0, q_1\}) = \{q_0, q_1, q_2, q_3\} \end{aligned}$$

$$\begin{aligned} \hat{\delta}(q_0, 10) &= \text{ECLOSE}\left(\bigcup_{p \in \hat{\delta}(q_0, 1)} \delta(p, 0)\right) \\ &= \text{ECLOSE}(\delta(q_0, 0) \cup \delta(q_1, 0) \cup \delta(q_2, 0) \cup \delta(q_3, 0)) \\ &= \text{ECLOSE}(\{q_0, q_2, q_3\}) = \{q_0, q_2, q_3\} \end{aligned}$$

ε -NFA 的语言

回顾

DFA $D = (Q, \Sigma, \delta, q_0, F)$ 和 NFA $N = (Q, \Sigma, \delta, q_0, F)$ 的语言分别为

$$\mathbf{L}(D) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \in F\},$$

$$\mathbf{L}(N) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \cap F \neq \emptyset\}.$$

定义

若 $E = (Q, \Sigma, \delta, q_0, F)$ 是一个 ε -NFA, 则 E 接受的语言为

$$\mathbf{L}(E) = \{w \in \Sigma^* \mid \hat{\delta}(q_0, w) \cap F \neq \emptyset\}.$$

消除空转移的子集构造法

构造方法

如果 ε -NFA $E = (Q_E, \Sigma, \delta_E, q_E, F_E)$, 构造 DFA

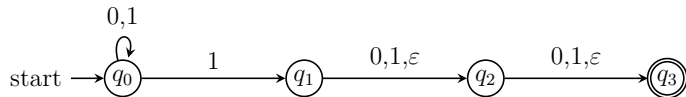
$$D = (Q_D, \Sigma, \delta_D, q_D, F_D)$$

- ❶ $Q_D = 2^{Q_E}$, 或 $Q_D = \{S \subseteq Q_E \mid S = \text{ECLOSE}(S)\}$;
- ❷ $q_D = \text{ECLOSE}(q_E)$;
- ❸ $F_D = \{S \mid S \in Q_D, S \cap F_E \neq \emptyset\}$;
- ❹ $\forall S \in Q_D, \forall a \in \Sigma,$

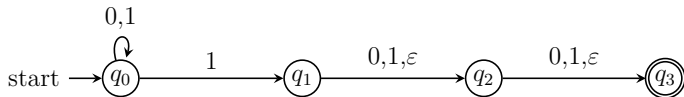
$$\delta_D(S, a) = \text{ECLOSE}\left(\bigcup_{p \in S} \delta_E(p, a)\right).$$

那么有 $\mathbf{L}(D) = \mathbf{L}(E)$.

续例 11. 将下图 L 的 ε -NFA, 转为等价的 DFA.

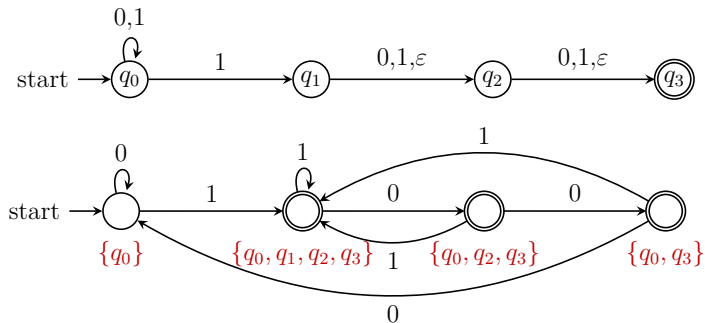


续例 11. 将下图 L 的 ε -NFA, 转为等价的 DFA.



	0	1	ε	ECLOSE()
$\rightarrow q_0$	$\{q_0\}$	$\{q_0, q_1\}$	\emptyset	$\{q_0\}$
q_1	$\{q_2\}$	$\{q_2\}$	$\{q_2\}$	$\{q_1, q_2, q_3\}$
q_2	$\{q_3\}$	$\{q_3\}$	$\{q_3\}$	$\{q_2, q_3\}$
$*q_3$	\emptyset	\emptyset	\emptyset	$\{q_3\}$

续例 11. 将下图 L 的 ε -NFA, 转为等价的 DFA.



	0	1
$\rightarrow \{q_0\}$	$\{q_0\}$	$\{q_0, q_1, q_2, q_3\}$
$*\{q_0, q_1, q_2, q_3\}$	$\{q_0, q_2, q_3\}$	$\{q_0, q_1, q_2, q_3\}$
$*\{q_0, q_2, q_3\}$	$\{q_0, q_3\}$	$\{q_0, q_1, q_2, q_3\}$
$*\{q_0, q_3\}$	$\{q_0\}$	$\{q_0, q_1, q_2, q_3\}$

ε -NFA 与 DFA 等价性

定理 2

如果语言 L 被 ε -NFA 接受, 当且仅当 L 被 DFA 接受.

证明: 必要性显然成立, 因为任何 DFA 都是 ε -NFA.

为证明充分性, 对 w 归纳, 往证 $\hat{\delta}_E(q_E, w) = \hat{\delta}_D(q_D, w)$.

① 当 $w = \varepsilon$ 时

$$\hat{\delta}_E(q_E, \varepsilon) = \text{ECLOSE}(q_E) = q_D = \hat{\delta}_D(q_D, \varepsilon).$$

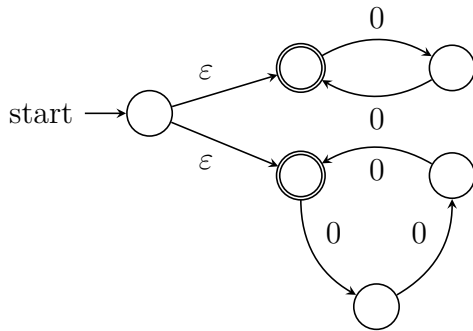
② 当 $w = xa$ 时

$$\begin{aligned}\hat{\delta}_E(q_E, xa) &= \text{ECLOSE}\left(\bigcup_{p \in \hat{\delta}_E(q_E, x)} \delta_E(p, a)\right) = \text{ECLOSE}\left(\bigcup_{p \in \hat{\delta}_D(q_D, x)} \delta_E(p, a)\right) \\ &= \delta_D(\hat{\delta}_D(q_D, x), a) = \hat{\delta}_D(q_D, xa)\end{aligned}$$



例 12. Design ε -NFA for language: $\{0^k \mid k \text{ is a multiple of 2 or 3}\}$.

例 12. Design ε -NFA for language: $\{0^k \mid k \text{ is a multiple of 2 or 3}\}$.



形式语言与自动机理论

正则表达式

王春宇

计算机科学与技术学院
哈尔滨工业大学

正则表达式

- 正则表达式
 - 语言的运算
 - 正则表达式的递归定义
 - 运算符的优先级
 - 正则表达式示例
- 有穷自动机和正则表达式
- 正则表达式的代数定律

正则表达式

- 有穷自动机
 - 通过**机器**装置描述正则语言
 - 用计算机编写相应算法, 易于实现
- 正则表达式
 - 通过**表达式**描述正则语言, 代数表示方法, 使用方便
 - 应用广泛
 - `grep` 工具 (Global Regular Expression and Print)
 - Emacs / Vim 文本编辑器
 - `lex` / `flex` 词法分析器
 - 各种程序设计语言 Python / Perl / Haskell / ...

语言的运算

设 L 和 M 是两个语言, 那么

并

$$L \cup M = \{w \mid w \in L \text{ 或 } w \in M\}$$

连接

$$L \cdot M = \{w \mid w = xy, x \in L \text{ 且 } y \in M\}$$

幂

$$L^0 = \{\varepsilon\}$$

$$L^1 = L$$

$$L^n = L^{n-1} \cdot L$$

克林闭包

$$L^* = \bigcup_{i=0}^{\infty} L^i$$

例 1. 若有语言 $L = \{0, 11\}$ 和 $M = \{\varepsilon, 001\}$, 那么

$$L \cup M = \qquad L^0 =$$

$$LM = \qquad L^1 =$$

$$ML = \qquad L^2 =$$

例 2. 对于空语言 \emptyset

$$\emptyset^0 =$$

$$\forall n \geq 1, \quad \emptyset^n =$$

$$\emptyset^* =$$

四则运算表达式的递归定义:

- ① 任何数都是四则运算表达式;
- ② 如果 a 和 b 是四则运算表达式, 那么

$$a + b, a - b, a \times b, a \div b \text{ 和 } (a)$$

都是四则运算表达式.

正则表达式的递归定义

定义

如果 Σ 为字母表, 则 Σ 上的正则表达式递归定义为:

- ① \emptyset 是一个正则表达式, 表示空语言;
 ϵ 是一个正则表达式, 表示语言 $\{\epsilon\}$;
 $\forall a \in \Sigma$, a 是一个正则表达式, 表示语言 $\{a\}$;
- ② 如果正则表达式 r 和 s 分别表示语言 R 和 S , 那么

$$r + s, rs, r^* \text{ 和 } (r)$$

都是正则表达式, 分别表示语言

$$R \cup S, R \cdot S, R^* \text{ 和 } R.$$

运算符的优先级

正则表达式中三种运算以及括号的优先级:

- ① 首先, “括号” 优先级最高;
- ② 其次, “星” 运算: r^* ;
- ③ 然后, “连接” 运算: rs , $r \cdot s$;
- ④ 最后, “加” 最低: $r + s$, $r \cup s$;

例 3.

$$\begin{aligned} 1 + 01^* &= 1 + (0(1^*)) \\ &\neq 1 + (01)^* \\ &\neq (1 + 01)^* \\ &\neq (1 + 0)1^* \end{aligned}$$

正则表达式示例

例 4.

E	$L(E)$
$\mathbf{a + b}$	$L(\mathbf{a}) \cup L(\mathbf{b}) = \{a\} \cup \{b\} = \{a, b\}$
\mathbf{bb}	$L(\mathbf{b}) \cdot L(\mathbf{b}) = \{b\} \cdot \{b\} = \{bb\}$
$\mathbf{(a + b)(a + b)}$	$\{a, b\}\{a, b\} = \{aa, ab, ba, bb\}$
$\mathbf{(a + b)^*(a + bb)}$	$\{a, b\}^*\{a, bb\} = \{a, b\}^*\{a\} \cup \{a, b\}^*\{bb\} = \{w \in \{a, b\}^* \mid w \text{ 仅以 } a \text{ 或 } bb \text{ 结尾.}\}$
$\mathbf{1 + (01)^*}$	$\{1, \varepsilon, 01, 0101, 010101, \dots\}$
$\mathbf{(0 + 1)^*01(0 + 1)^*}$	$\{x01y \mid x, y \in \{0, 1\}^*\}$

例 5. 给出正则表达式 $(aa)^*(bb)^*b$ 定义的语言.

$$\begin{aligned} \mathbf{L}((aa)^*(bb)^*b) &= \mathbf{L}((aa)^*) \cdot \mathbf{L}((bb)^*) \cdot \mathbf{L}(b) \\ &= (\{a\}\{a\})^*(\{b\}\{b\})^*\{b\} \\ &= \{a^2\}^*\{b^2\}^*\{b\} \\ &= \{a^{2n}b^{2m+1} \mid n \geq 0, m \geq 0\} \end{aligned}$$

例 6. Design regular expression for $L = \{w \mid w \text{ consists of 0's and 1's, and the third symbol from the right end is 1.}\}$

$$(\mathbf{0} + \mathbf{1})^* \mathbf{1} (\mathbf{0} + \mathbf{1}) (\mathbf{0} + \mathbf{1})$$

例 7. Design regular expression for

$L = \{w \mid w \in \{0, 1\}^* \text{ and } w \text{ has no pair of consecutive 0's.}\}$

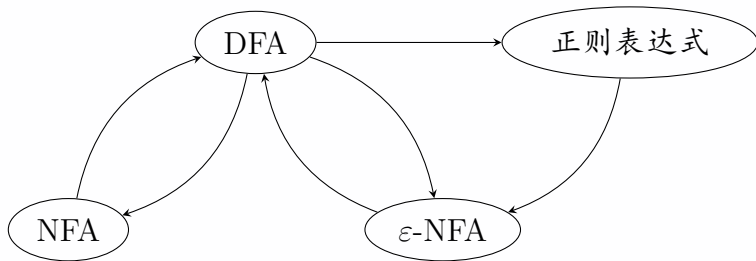
$$\mathbf{1^*(011^*)^*(0 + \varepsilon)} \text{ 或 } \mathbf{(1 + 01)^*(0 + \varepsilon)}$$

- 例. Design regular expression for $L = \{w \mid w \in \{0,1\}^* \text{ and } w \text{ contains } 01\}$.
- 例. Write a regular expression for $L = \{w \in \{0,1\}^* \mid 0 \text{ and } 1 \text{ alternate in } w\}$.
- 例. Find a regular expression for the set $\{a^n b^m \mid (n+m) \text{ is odd}\}$.
- 例. Give regular expression for the complement of $L = \{a^n b^m \mid n \geq 3, m \leq 4\}$.
- 例. Write a regular expression for the set of all C real numbers.

正则表达式

- 正则表达式
- 有穷自动机和正则表达式
 - 由 DFA 到正则表达式, 递归表达式法
 - 由 DFA 到正则表达式, 状态消除法
 - 由正则表达式到 ε -NFA
- 正则表达式的代数定律

DFA, NFA, ϵ -NFA 和正则表达式的等价性



由 DFA 到正则表达式, 递归表达式法

定理 3

若 $L = \mathbf{L}(A)$ 是某 DFA A 的语言, 那么存在正则表达式 R 满足 $L = \mathbf{L}(R)$.

由 DFA 到正则表达式, 递归表达式法

定理 3

若 $L = \mathbf{L}(A)$ 是某 DFA A 的语言, 那么存在正则表达式 R 满足 $L = \mathbf{L}(R)$.

证明: 对 DFA A 的状态编号, 令 1 为开始状态, 即

$$A = (\{1, 2, \dots, n\}, \Sigma, \delta, 1, F),$$

设正则表达式 $R_{ij}^{(k)}$ 表示从 i 到 j 但中间节点不超过 k 全部路径的字符串集:

$$R_{ij}^{(k)} = \{x \mid \hat{\delta}(i, x) = j, x \text{ 经过的状态除两端外都不超过 } k\}.$$



那么与 $A = (\{1, 2, \dots, n\}, \Sigma, \delta, 1, F)$ 等价的正则表达式为

$$\bigcup_{j \in F} R_{1j}^{(n)}$$

且递归式为

$$R_{ij}^{(k)} = R_{ij}^{(k-1)} + R_{ik}^{(k-1)} (R_{kk}^{(k-1)})^* R_{kj}^{(k-1)}$$

$$R_{ij}^{(0)} = \begin{cases} \{a \mid \delta(q_i, a) = q_j\} & i \neq j \\ \{a \mid \delta(q_i, a) = q_j\} \cup \{\varepsilon\} & i = j \end{cases}$$

下面对 k 归纳, 证明可用以上递归式求得 $R_{ij}^{(k)}$.

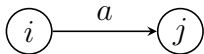
归纳基础: 当 $i \neq j$, $k = 0$ 时, 即 i 到 j 没经过任何中间节点

- 没有 i 到 j 的状态转移



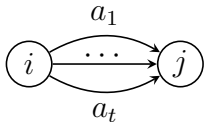
$$R_{ij}^{(0)} = \emptyset$$

- 有一个 i 到 j 的状态转移



$$R_{ij}^{(0)} = \mathbf{a}$$

- 有多个 i 到 j 的状态转移



$$R_{ij}^{(0)} = \mathbf{a}_1 + \mathbf{a}_2 + \dots + \mathbf{a}_t$$

归纳基础 (续): 当 $i = j$, $k = 0$ 时, 即从 i 到自身没经过任何中间节点

- 状态 i 没有到自己的转移



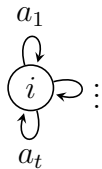
$$R_{ii}^{(0)} = \epsilon$$

- 状态 i 有一个到自身的转移



$$R_{ii}^{(0)} = \mathbf{a} + \epsilon$$

- 状态 i 有多个到自身的转移

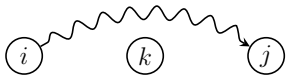


$$R_{ii}^{(0)} = \mathbf{a}_1 + \mathbf{a}_2 + \cdots + \mathbf{a}_t + \epsilon$$

归纳假设: 假设已知 $R_{ij}^{(k-1)}$, $R_{ik}^{(k-1)}$, $R_{kk}^{(k-1)}$ 和 $R_{kj}^{(k-1)}$.

归纳递推: 那么 $R_{ij}^{(k)}$ 中全部路径, 可用节点 k 分为两部分

- 从 i 到 j 不经过 k 的



$$R_{ij}^{(k)} = R_{ij}^{(k-1)}$$

- 从 i 到 j 经过 k 的

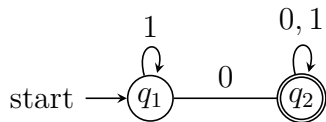


$$R_{ij}^{(k)} = R_{ik}^{(k-1)} (R_{kk}^{(k-1)})^* R_{kj}^{(k-1)}$$

因此 $R_{ij}^{(k)} = R_{ij}^{(k-1)} + R_{ik}^{(k-1)} (R_{kk}^{(k-1)})^* R_{kj}^{(k-1)}$.

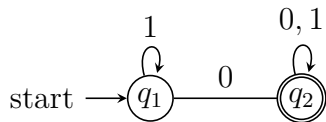


例 8. 将如图 DFA 转换为正则表达式.



- 计算 $R_{ij}^{(0)}$

例 8. 将如图 DFA 转换为正则表达式.



- 计算 $R_{ij}^{(0)}$

$R_{ij}^{(k)}$	$k = 0$
$R_{11}^{(0)}$	$\epsilon + \mathbf{1}$
$R_{12}^{(0)}$	$\mathbf{0}$
$R_{21}^{(0)}$	\emptyset
$R_{22}^{(0)}$	$\epsilon + \mathbf{0} + \mathbf{1}$

- 计算 $R_{ij}^{(1)} = R_{ij}^{(0)} + R_{i1}^{(0)}(R_{11}^{(0)})^* R_{1j}^{(0)}$

$R_{ij}^{(k)}$	$k = 0$
$R_{11}^{(0)}$	$\boldsymbol{\varepsilon} + \mathbf{1}$
$R_{12}^{(0)}$	$\mathbf{0}$
$R_{21}^{(0)}$	\emptyset
$R_{22}^{(0)}$	$\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1}$

- 计算 $R_{ij}^{(1)} = R_{ij}^{(0)} + R_{i1}^{(0)}(R_{11}^{(0)})^*R_{1j}^{(0)}$

$R_{ij}^{(k)}$	$k = 0$	$R_{ij}^{(k)}$	$k = 1$
$R_{11}^{(0)}$	$\varepsilon + \mathbf{1}$	$R_{11}^{(1)}$	$(\varepsilon + \mathbf{1}) + (\varepsilon + \mathbf{1})(\varepsilon + \mathbf{1})^*(\varepsilon + \mathbf{1})$
$R_{12}^{(0)}$	$\mathbf{0}$	$R_{12}^{(1)}$	$\mathbf{0} + (\varepsilon + \mathbf{1})(\varepsilon + \mathbf{1})^*\mathbf{0}$
$R_{21}^{(0)}$	\emptyset	$R_{21}^{(1)}$	$\emptyset + \emptyset(\varepsilon + \mathbf{1})^*(\varepsilon + \mathbf{1})$
$R_{22}^{(0)}$	$\varepsilon + \mathbf{0} + \mathbf{1}$	$R_{22}^{(1)}$	$\varepsilon + \mathbf{0} + \mathbf{1} + \emptyset(\varepsilon + \mathbf{1})^*\mathbf{0}$

- 几个基本的化简规则

如果 \mathbf{r} 和 \mathbf{s} 是两个正则表达式

$$(\epsilon + \mathbf{r})^* = \mathbf{r}^*$$

$$(\epsilon + \mathbf{r})\mathbf{r}^* = \mathbf{r}^*$$

$$\mathbf{r} + \mathbf{r}\mathbf{s}^* = \mathbf{r}\mathbf{s}^*$$

$$\emptyset \mathbf{r} = \mathbf{r} \emptyset = \emptyset$$

(零元)

$$\emptyset + \mathbf{r} = \mathbf{r} + \emptyset = \mathbf{r}$$

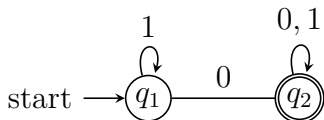
(单位元)

- 化简 $R_{ij}^{(1)}$

$R_{ij}^{(k)}$	$k = 1$	化简
$R_{11}^{(1)}$	$(\varepsilon + \mathbf{1}) + (\varepsilon + \mathbf{1})(\varepsilon + \mathbf{1})^*(\varepsilon + \mathbf{1})$	$\mathbf{1}^*$
$R_{12}^{(1)}$	$\mathbf{0} + (\varepsilon + \mathbf{1})(\varepsilon + \mathbf{1})^*\mathbf{0}$	$\mathbf{1}^*\mathbf{0}$
$R_{21}^{(1)}$	$\emptyset + \emptyset(\varepsilon + \mathbf{1})^*(\varepsilon + \mathbf{1})$	\emptyset
$R_{22}^{(1)}$	$\varepsilon + \mathbf{0} + \mathbf{1} + \emptyset(\varepsilon + \mathbf{1})^*\mathbf{0}$	$\varepsilon + \mathbf{0} + \mathbf{1}$

- 计算 $R_{ij}^{(2)} = R_{ij}^{(1)} + R_{i2}^{(1)}(R_{22}^{(1)})^*R_{2j}^{(1)}$

$R_{ij}^{(k)}$	$k = 1$	$R_{ij}^{(k)}$	$k = 2$
$R_{11}^{(1)}$	$\mathbf{1}^*$	$R_{11}^{(2)}$	$\mathbf{1}^* + \mathbf{1}^*\mathbf{0}(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})^*\emptyset$
$R_{12}^{(1)}$	$\mathbf{1}^*\mathbf{0}$	$R_{12}^{(2)}$	$\mathbf{1}^*\mathbf{0} + \mathbf{1}^*\mathbf{0}(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})^*(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})$
$R_{21}^{(1)}$	\emptyset	$R_{21}^{(2)}$	$\emptyset + (\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})^*\emptyset$
$R_{22}^{(1)}$	$\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1}$	$R_{22}^{(2)}$	$\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1} + (\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})^*(\boldsymbol{\varepsilon} + \mathbf{0} + \mathbf{1})$



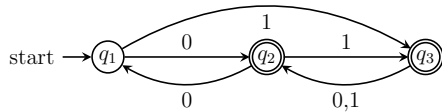
- 化简 $R_{ij}^{(2)}$

$R_{ij}^{(k)}$	$k = 2$	化简
$R_{11}^{(2)}$	$1^* + 1^*0(\epsilon + 0 + 1)^*\emptyset$	1^*
$R_{12}^{(2)}$	$1^*0 + 1^*0(\epsilon + 0 + 1)^*(\epsilon + 0 + 1)$	$1^*0(0 + 1)^*$
$R_{21}^{(2)}$	$\emptyset + (\epsilon + 0 + 1)(\epsilon + 0 + 1)^*\emptyset$	\emptyset
$R_{22}^{(2)}$	$\epsilon + 0 + 1 + (\epsilon + 0 + 1)(\epsilon + 0 + 1)^*(\epsilon + 0 + 1)$	$(0 + 1)^*$

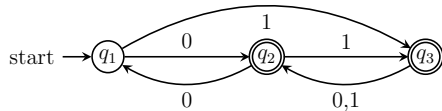
- 因只有 q_2 是接受状态, 所以该 DFA 正则表达式为

$$R_{12}^{(2)} = 1^*0(0 + 1)^*.$$

例 9. 将如图 DFA 转换为正则表达式.

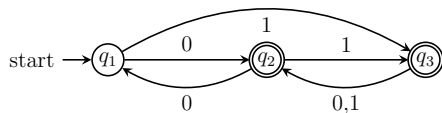


例 9. 将如图 DFA 转换为正则表达式.



	$k = 0$	$k = 1$	$k = 2$
$R_{11}^{(k)}$	ϵ	ϵ	$(00)^*$
$R_{12}^{(k)}$	0	0	$0(00)^*$
$R_{13}^{(k)}$	1	1	0^*1
$R_{21}^{(k)}$	0	0	$0(00)^*$
$R_{22}^{(k)}$	ϵ	$\epsilon + 00$	$(00)^*$
$R_{23}^{(k)}$	1	$1 + 01$	0^*1
$R_{31}^{(k)}$	\emptyset	\emptyset	$(0 + 1)(00)^*0$
$R_{32}^{(k)}$	$0 + 1$	$0 + 1$	$(0 + 1)(00)^*$
$R_{33}^{(k)}$	ϵ	ϵ	$\epsilon + (0 + 1)0^*1$

例 9. 将如图 DFA 转换为正则表达式.



仅状态 2 和 3 是接受状态:

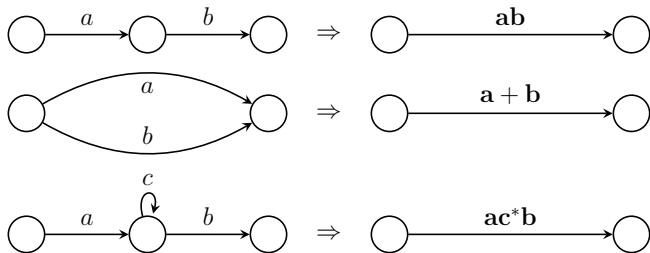
$$\begin{aligned}
 R_{12}^{(3)} &= R_{12}^{(2)} + R_{13}^{(2)}(R_{33}^{(2)})^*R_{32}^{(2)} \\
 &= 0(00)^* + 0^*1(\epsilon + (0+1)0^*1)^*(0+1)(00)^* \\
 &= 0(00)^* + 0^*1((0+1)0^*1)^*(0+1)(00)^*
 \end{aligned}$$

$$\begin{aligned}
 R_{13}^{(3)} &= R_{13}^{(2)} + R_{13}^{(2)}(R_{33}^{(2)})^*R_{33}^{(2)} \\
 &= 0^*1 + 0^*1(\epsilon + (0+1)0^*1)^*(\epsilon + (0+1)0^*1) \\
 &= 0^*1((0+1)0^*1)^*
 \end{aligned}$$

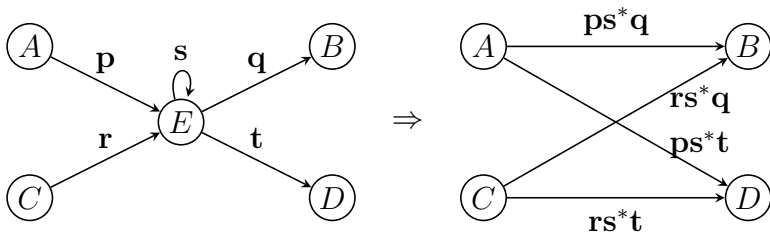
$$R_{12}^{(3)} + R_{13}^{(3)} = 0^*1((0+1)0^*1)^*(\epsilon + (0+1)(00)^*) + 0(00)^*.$$

由 DFA 到正则表达式, 状态消除法

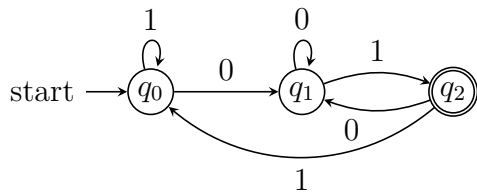
- 从 DFA 中逐个删除状态
- 用标记了正则表达式的新路径替换被删掉的路径
- 保持“自动机”等价.



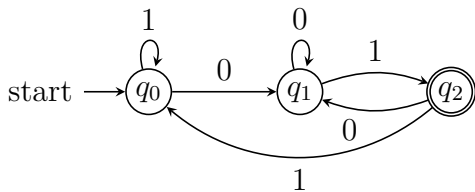
- 更一般的情况如图
- 若要删除状态 E , 需添加相应路径



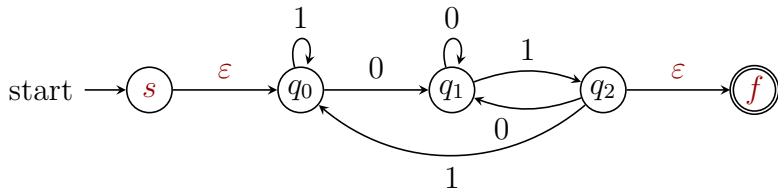
例 10. 利用状态消除法, 设计下图自动机的正则表达式.



例 10. 利用状态消除法, 设计下图自动机的正则表达式.

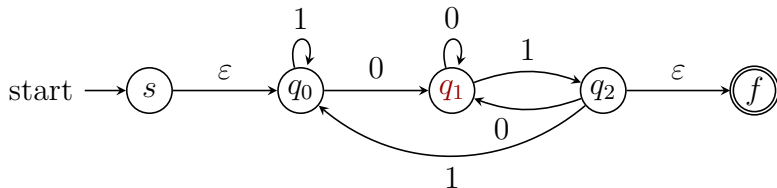


① 利用空转移, 添加新的开始 s 和结束状态 f :

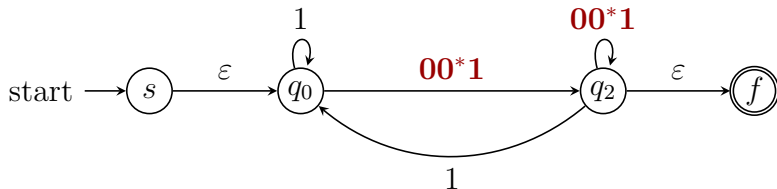


例 10. 利用状态消除法, 设计下图自动机的正则表达式.

① 利用空转移, 添加新的开始 s 和结束状态 f :

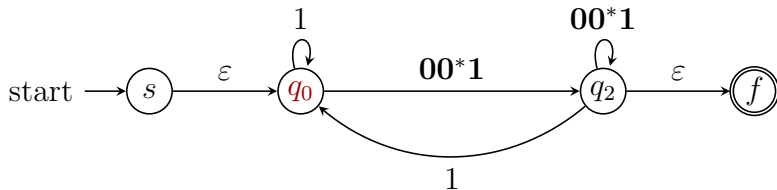


② 消除状态 q_1 , 添加路径 $q_0 \rightarrow q_2$ 和 $q_2 \rightarrow q_2$:

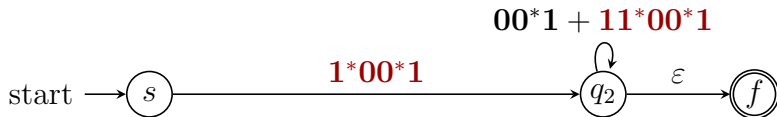


例 10. 利用状态消除法, 设计下图自动机的正则表达式.

② 消除状态 q_1 , 添加路径 $q_0 \rightarrow q_2$ 和 $q_2 \rightarrow q_2$:

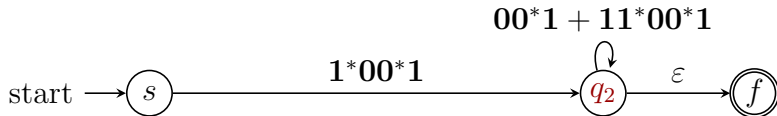


③ 消除状态 q_0 , 添加路径 $s \rightarrow q_2$ 和 $q_2 \rightarrow q_2$:



例 10. 利用状态消除法, 设计下图自动机的正则表达式.

③ 消除状态 q_0 , 添加路径 $s \rightarrow q_2$ 和 $q_2 \rightarrow q_2$:

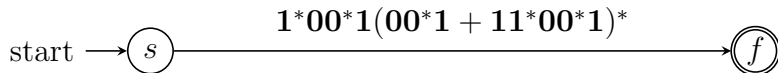


④ 消除状态 q_2 , 添加路径 $s \rightarrow f$:



例 10. 利用状态消除法, 设计下图自动机的正则表达式.

④ 消除状态 q_2 , 添加路径 $s \rightarrow f$:



⑤ 因此该自动机的正则表达式为

$$1^*00^*1(00^*1 + 11^*00^*1)^*.$$

由正则表达式到有穷自动机

定理 4

正则表达式定义的语言, 都可被有穷自动机识别.

由正则表达式构造 ε -NFA

任何正则表达式 r , 都存在等价的 ε -NFA A , 即 $L(A) = L(r)$, 并且 A 满足:

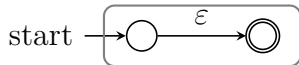
- ① 仅有一个接收状态;
- ② 没有进入开始状态的边;
- ③ 没有离开接受状态的边.

证明: 归纳基础:

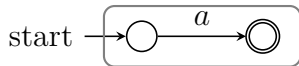
① 对于 \emptyset , 有 ε -NFA:



② 对于 ε , 有 ε -NFA:



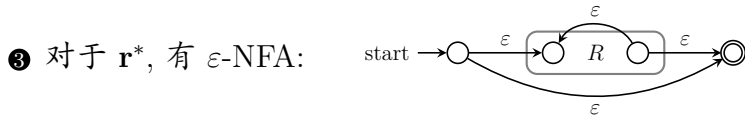
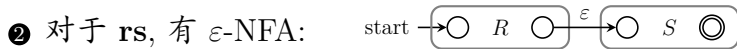
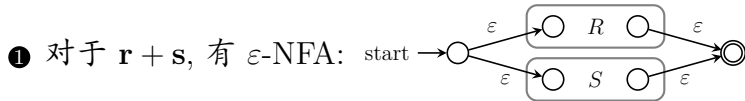
③ $\forall a \in \Sigma$, 对于 a , 有 ε -NFA:



归纳递推: 假设正则表达式 r 和 s 的 ε -NFA 分别为 R 和 S



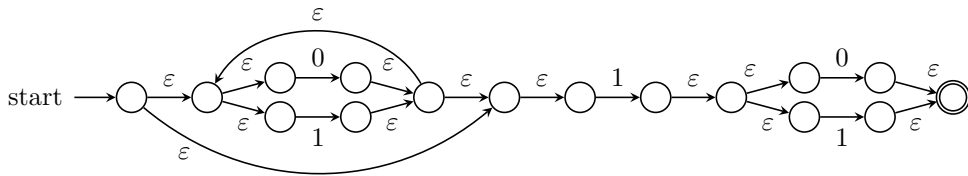
那么正则表达式 $r + s$, rs 和 r^* , 可由 R 和 S 分别构造如下:



因此任何结构的正则表达式, 都可递归构造出等价的 ε -NFA.



例 11. 正则表达式 $(0 + 1)^*1(0 + 1)$ 构造为 ε -NFA.



思考题

正则表达式到 ε -NFA 构造方法中的 3 个限制条件, 都有必要吗?

正则表达式

- 正则表达式
- 有穷自动机和正则表达式
- 正则表达式的代数定律
 - 基本的代数定律
 - 发现与验证代数定律

正则表达式的代数定律

定义

含有变量的两个正则表达式, 如果以任意语言替换其变量, 二者所表示的语言仍然相同, 则称这两个正则表达式**等价**. 在这样的意义下, 正则表达式满足一些**代数定律**.

- 并运算

$$(L + M) + N = L + (M + N) \quad (\text{结合律})$$

$$L + M = M + L \quad (\text{交换律})$$

$$L + L = L \quad (\text{幂等律})$$

$$\emptyset + L = L + \emptyset = L \quad (\text{单位元 } \emptyset)$$

- 连接运算

$$(LM)N = L(MN) \quad (\text{结合律})$$

$$\epsilon L = L\epsilon = L \quad (\text{单位元 } \epsilon)$$

$$\emptyset L = L\emptyset = \emptyset \quad (\text{零元 } \emptyset)$$

$$LM \neq ML$$

- 分配率

$$L(M + N) = LM + LN \quad (\text{左分配律})$$

$$(M + N)L = ML + NL \quad (\text{右分配律})$$

- 闭包运算

$$(L^*)^* = L^*$$

$$\emptyset^* = \epsilon$$

$$\epsilon^* = \epsilon$$

$$L^* = L^+ + \epsilon$$

$$(\epsilon + L)^* = L^*$$

发现与验证正则表达式的代数定律

检验方法

要判断表达式 E 和 F 是否等价, 其中变量为 L_1, \dots, L_n :

- ① 将变量替换为具体表达式, 得正则表达式 \mathbf{r} 和 \mathbf{s} ,
例如, 替换 L_i 为 \mathbf{a}_i ;
- ② 判断 $\mathbf{L}(\mathbf{r}) \stackrel{?}{=} \mathbf{L}(\mathbf{s})$, 如果相等则 $E = F$, 否则 $E \neq F$.

例 12. 判断 $(L + M)^* = (L^*M^*)^*$.

- ❶ 将 L 和 M 替换为 \mathbf{a} 和 \mathbf{b} ;
- ❷ $(\mathbf{a} + \mathbf{b})^* \stackrel{?}{=} (\mathbf{a}^*\mathbf{b}^*)^*$;
- ❸ 因为 $\mathbf{L}((\mathbf{a} + \mathbf{b})^*) = \mathbf{L}((\mathbf{a}^*\mathbf{b}^*)^*)$;
- ❹ 所以 $(L + M)^* = (L^*M^*)^*$.

例 13. 判断 $L + ML = (L + M)L$.

- ❶ 将 L 和 M 替换为 a 和 b ;
- ❷ 判断 $a + ba \stackrel{?}{=} (a + b)a$;
- ❸ 因为 $aa \notin a + ba$ 而 $aa \in (a + b)a$;
- ❹ 所以 $a + ba \neq (a + b)a$;
- ❺ 即 $L + ML \neq (L + M)L$.

注意

这种方法**仅限于**判断正则表达式, 否则可能会发生错误.

例 14. 若用此方法判断 $L \cap M \cap N \stackrel{?}{=} L \cap M$, 以 **a, b, c** 替换 L, M, N , 有

$$\{a\} \cap \{b\} \cap \{c\} = \emptyset = \{a\} \cap \{b\},$$

而显然

$$L \cap M \cap N \neq L \cap M.$$

形式语言与自动机理论

正则语言的性质

王春宇

计算机科学与技术学院
哈尔滨工业大学

正则语言的性质

- 证明语言的非正则性
 - 正则语言的泵引理
 - 泵引理的应用
 - 泵引理只是必要条件
- 正则语言的封闭性
- 正则语言的判定性质
- 自动机的最小化

例 1. $L = \{0^m 1^n \mid m, n \geq 0\}$ 是否是正则语言?

例 2. $L = \{0^m 1^n \mid m \geq 2, n \geq 4\}$ 是否是正则语言?

例 3. $L_{01} = \{0^n 1^n \mid n \geq 0\}$ 是否是正则语言?

正则语言的泵引理

定理 5 (正则语言的泵引理)

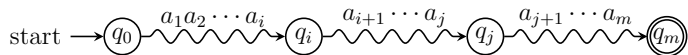
如果语言 L 是正则的, 那么存在正整数 N , 它只依赖于 L , 对 $\forall w \in L$, 只要 $|w| \geq N$, 就可以将 w 分为三部分 $w = xyz$ 满足:

- ① $y \neq \varepsilon$ ($|y| > 0$);
- ② $|xy| \leq N$;
- ③ $\forall k \geq 0, xy^kz \in L$.

证明:

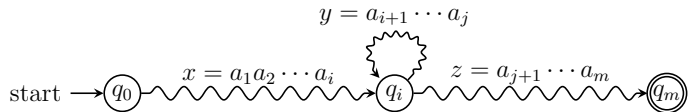
① 如果 L 正则, 那么存在有 n 个状态 DFA A 使 $\mathbf{L}(A) = L$;

② 取 $w = a_1 \dots a_m \in L$ ($m \geq n$), 定义 $q_i = \hat{\delta}(q_0, a_1 \dots a_i)$;



③ 由鸽巢原理, 必有两状态相同 $q_i = q_j$ ($0 \leq i < j \leq n$);

④ 那么 $w = xyz$ 如图, 且有 $\forall k > 0, xy^kz \in L$;



⑤ 而因为 $i < j$ 所以 $y \neq \varepsilon$ (即 $|y| > 0$), 因为 $j \leq n$ 所以 $|xy| \leq n$. □

泵引理的应用

续例3. 证明 $L_{01} = \{0^n 1^n \mid n \geq 0\}$ 不是正则语言.

证明:

- ① 假设 L_{01} 是正则的.
- ② 那么, 存在 $N \in \mathbb{Z}^+$, 对 $\forall w \in L_{01} (|w| \geq N)$ 满足泵引理.
- ③ 从 L_{01} 中取 $w = 0^N 1^N$, 显然 $w \in L_{01}$ 且 $|w| = 2N \geq N$.
- ④ 那么, w 可被分为 $w = xyz$, 且 $|xy| \leq N$ 和 $y \neq \varepsilon$.
- ⑤ 因此 y 只能是 0^m 且 $m > 0$.
- ⑥ 那么 $xy^2z = 0^{N+m} 1^N \notin L_{01}$, 而由泵引理 $xy^2z \in L_{01}$, 矛盾.
- ⑦ 所以假设不成立, L_{01} 不是正则的.



例 4. 证明 $L_{\text{eq}} = \{w \in \{0,1\}^* \mid w \text{ 由数量相等的 } 0 \text{ 和 } 1 \text{ 构成}\}$ 不是正则的.

思考题

刚刚已经证明了

$$L_{01} = \{0^n 1^n \mid n \geq 0\}$$

不是正则语言, 那么能否使用

$$L_{01} \subseteq L_{\text{eq}}$$

来说明 L_{eq} 也不是正则的呢?

续例 4. 证明 $L_{\text{eq}} = \{w \in \{0, 1\}^* \mid w \text{ 由数量相等的 } 0 \text{ 和 } 1 \text{ 构成}\}$ 不是正则的.
证明:

- ❶ 假设 L_{eq} 是正则的.
- ❷ 那么, 存在 $N \in \mathbb{Z}^+$, 对 $\forall w \in L_{\text{eq}} (|w| \geq N)$ 满足泵引理.
- ❸ 从 L_{eq} 中取 $w = 0^N 1^N$, 显然 $w \in L_{\text{eq}}$ 且 $|w| = 2N \geq N$.
- ❹ 那么, w 可被分为 $w = xyz$, 且 $|xy| \leq N$ 和 $y \neq \varepsilon$.
- ❺ 因此 y 只能是 0^m 且 $m > 0$.
- ❻ 那么 $xy^2z = 0^{N+m}1^N \notin L_{\text{eq}}$, 而由泵引理 $xy^2z \in L_{\text{eq}}$, 矛盾.
- ❼ 所以假设不成立, L_{eq} 不是正则的.



例 5. 证明 $L = \{0^i 1^j \mid i > j\}$ 不是正则的.

证明:

- ① 假设 L 是正则的.
- ② 那么, 存在 $N \in \mathbb{Z}^+$, 对 $\forall w \in L (|w| \geq N)$ 满足泵引理.
- ③ 从 L 中取 $w = 0^{N+1}1^N$, 则 $w \in L$ 且 $|w| = 2N + 1 \geq N$.
- ④ 由泵引理, w 可被分为 $w = xyz$, 且 $|xy| \leq N$ 和 $y \neq \varepsilon$.
- ⑤ 那么, y 只能是 0^m 且 $m \geq 1$.
- ⑥ 那么, $xz = xy^0z = 0^{N+1-m}1^N \notin L$, 因为 $N + 1 - m \leq N$, 而由泵引理 $xy^0z \in L$, 矛盾.
- ⑦ 所以假设不成立, L 不是正则的.



例 6. Prove $L = \{a^3b^nc^{n-3} \mid n \geq 3\}$ is not regular.

证明:

- ① 假设 L 是正则的.
- ② 那么, 存在 $N \in \mathbb{Z}^+$, 对 $\forall w \in L (|w| \geq N)$ 满足泵引理.
- ③ 从 L 中取 $w = a^3b^Nc^{N-3}$, 则 $w \in L$ 且 $|w| = 2N \geq N$.
- ④ 由泵引理, w 可被分为 $w = xyz$, 且 $|xy| \leq N$ 和 $y \neq \varepsilon$.
- ⑤ 那么, 则 y 只可能有 3 种情况 ($m > 0, r > 0, s > 0$):
 - ① $y = a^m$, 则 $xy^2z = a^{3+m}b^Nc^{N-3} \notin L$;
 - ② $y = b^m$, 则 $xy^2z = a^3b^{N+m}c^{N-3} \notin L$;
 - ③ $y = a^rb^s$, 则 $xy^2z = a^3b^sa^rb^Nc^{N-3} \notin L$.
- ⑥ 无论 y 为何种情况, xy^2z 都不可能在 L 中, 与泵引理矛盾.
- ⑦ 所以假设不成立, L 不是正则的.



思考题

- $L = \{0^n 1^n \mid 0 \leq n \leq 100\}$ 是否是正则语言?
- 有限的语言, 是否符合泵引理呢?
 - \emptyset
 - $\{\epsilon\}$
 - $\{0, 00\}$

泵引理只是必要条件

- 泵引理只是正则语言的必要条件
- 只能用来证明某个语言不是正则的
- 与正则语言等价的定理 — Myhill-Nerode Theorem

例7. 语言 L 不是正则的, 但每个串都可以应用泵引理

$$L = \{ca^n b^n \mid n \geq 1\} \cup \{c^k w \mid k \neq 1, w \in \{a, b\}^*\}$$

- 其中 $A = \{ca^n b^n \mid n \geq 1\}$ 部分不是正则的
- 而 $B = \{c^k w \mid k \neq 1, w \in \{a, b\}^*\}$ 部分是正则的
- 而 A 的任何串 $w = ca^i b^i$, 都可应用泵引理, 因为

$$w = (\varepsilon)(c)(a^i b^i)$$

重复字符 c 生成的新串都会落入 B 中

思考题

对任何正则语言 L , 在泵引理中, 与 L 相关联的正整数 N

- 与识别 L 的 DFA 的状态数 n 之间有何关系?
- 与识别 L 的 NFA 的状态数之间呢?

思考题

语言

$$L = \{0^n x 1^n \mid n \geq 1, x \in \{0, 1\}^*\}$$

是否是正则语言?

正则语言的性质

- 证明语言的非正则性
- 正则语言的封闭性
- 正则语言的判定性质
- 自动机的最小化

正则语言的封闭性

定义

正则语言经某些运算后得到的新语言仍保持正则, 称为在这些运算下封闭.

正则语言 L 和 M , 在这些运算下封闭

- 并: $L \cup M$
- 交: $L \cap M$
- 连接: LM
- 反转: $L^R = \{w^R \mid w \in L\}$
- 闭包: L^*
- 同态: $h(L) = \{h(w) \mid w \in L, \text{同态 } h: \Sigma \rightarrow \Gamma^*\}$
- 补: \bar{L}
- 逆同态:
- 差: $L - M$
- $h^{-1}(L) = \{w \in \Sigma^* \mid h(w) \in L \subseteq \Gamma^*, \text{同态 } h: \Sigma \rightarrow \Gamma^*\}$

定理 6 (并/连接/闭包的封闭性)

正则语言在并, 连接和闭包运算下保持封闭.

证明: 由正则表达式的定义得证.



定理 7 (补运算封闭性)

如果 L 是 Σ 上的正则语言, 那么 $\overline{L} = \Sigma^* - L$ 也是正则的.

证明: 设接受语言 L 的 DFA

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

即 $\mathbf{L}(A) = L$. 构造 DFA

$$B = (Q, \Sigma, \delta, q_0, Q - F)$$

则有 $\overline{L} = \mathbf{L}(B)$, 因为 $\forall w \in \Sigma^*$

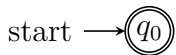
$$w \in \overline{L} \iff \hat{\delta}(q_0, w) \notin F \iff \hat{\delta}(q_0, w) \in Q - F \iff w \in \mathbf{L}(B).$$



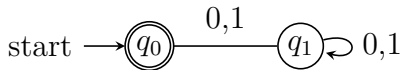
注意

使用这种方法求正则语言的补时, DFA 不能有缺失状态.

例 8. 若 $\Sigma = \{0, 1\}$, $L = \{\varepsilon\}$ 的 DFA 如图, 请给出 \bar{L} 的 DFA.



应使用完整的 DFA 去求补:



思考题

如何求正则表达式的补？

例 9. 证明 $L_{\text{neq}} = \{w \mid w \text{ 由数量不相等的 } 0 \text{ 和 } 1 \text{ 构成} \}$ 不是正则的.

证明:

- 由泵引理不易直接证明 L_{neq} 不是正则的;
- 因为无论如何取 w , 将其分为 $w = xyz$ 时, 都不易产生 L_{neq} 之外的串;
- 而证明 L_{eq} 非正则很容易;
- 由补运算的封闭性, 所以 $L_{\text{neq}} = \overline{L_{\text{eq}}}$ 也不是正则的. □

定理 8

若 DFA A_L , A_M 和 A 的定义如下

$$A_L = (Q_L, \Sigma, \delta_L, q_L, F_L)$$

$$A_M = (Q_M, \Sigma, \delta_M, q_M, F_M)$$

$$A = (Q_L \times Q_M, \Sigma, \delta, (q_L, q_M), F_L \times F_M)$$

其中

$$\delta : (Q_L \times Q_M) \times \Sigma \rightarrow Q_L \times Q_M$$

$$\delta((p, q), a) = (\delta_L(p, a), \delta_M(q, a)).$$

则对任意 $w \in \Sigma^*$,

$$\hat{\delta}((q_L, q_M), w) = (\hat{\delta}(q_L, w), \hat{\delta}(q_M, w)).$$

证明: 对 w 的结构归纳.

归纳基础: 当 $w = \varepsilon$ 时

$$\begin{aligned}\hat{\delta}((q_L, q_M), \varepsilon) &= (q_L, q_M) && \hat{\delta} \text{ 的定义} \\ &= (\hat{\delta}_L(q_L, \varepsilon), \hat{\delta}_M(q_M, \varepsilon)) && \text{同理}\end{aligned}$$

归纳递推: 当 $w = xa$ 时

$$\begin{aligned}\hat{\delta}((q_L, q_M), xa) &= \delta(\hat{\delta}((q_L, q_M), x), a) && \hat{\delta} \text{ 的定义} \\ &= \delta((\hat{\delta}_L(q_L, x), \hat{\delta}_M(q_M, x)), a) && \text{归纳假设} \\ &= (\delta_L(\hat{\delta}_L(q_L, x), a), \delta_M(\hat{\delta}_M(q_M, x), a)) && \delta \text{ 的构造} \\ &= (\hat{\delta}_L(q_L, xa), \hat{\delta}_M(q_M, xa)) && \hat{\delta} \text{ 的定义}\end{aligned}$$



定理 9 (交运算封闭性)

如果 L 和 M 是正则语言, 那么 $L \cap M$ 也是正则语言.

证明 1: 由 $L \cap M = \overline{\overline{L} \cup \overline{M}}$ 得证. □

证明 2: 由定理 8 构造识别 $L \cap M$ 的 DFA A , 则 $\forall w \in \Sigma^*$,

$$\begin{aligned}w \in L \cap M &\iff \hat{\delta}_L(q_L, w) \in F_L \wedge \hat{\delta}_M(q_M, w) \in F_M \\&\iff (\hat{\delta}_L(q_L, w), \hat{\delta}_M(q_M, w)) \in F_L \times F_M \\&\iff \hat{\delta}((q_L, q_M), w) \in F_L \times F_M \\&\iff w \in \mathbf{L}(A).\end{aligned}$$

因此 $\mathbf{L}(A) = L \cap M$, 所以 $L \cap M$ 也是正则的. □

例 10. 如果已知语言

$$L_{01} = \{0^n 1^n \mid n \geq 0\}$$

不是正则的, 请用封闭性证明语言

$$L_{\text{eq}} = \{w \in \{0, 1\}^* \mid w \text{ 由数量相等的 } 0 \text{ 和 } 1 \text{ 构成}\}$$

也不是正则的.

证明:

- ① 首先, 因为 0^*1^* 是正则语言;
- ② 而 $L_{01} = \mathbf{L}(0^*1^*) \cap L_{\text{eq}};$
- ③ 如果 L_{eq} 是正则的, L_{01} 必然也是正则的;
- ④ 因为已知 L_{01} 不是正则的, 所以 L_{eq} 一定不是正则的.



思考题

为什么又能用 L_{eq} 的子集 L_{01} 是非正则的, 来证明 L_{eq} 是非正则的呢?

例 11. 如果 L_1 和 L_2 都不是正则的, 那么 $L_1 \cap L_2$ 一定不是正则的吗?

不一定. 因为, 如果令

$$L_1 = \{0^n 1^n \mid n \geq 0\}$$

$$L_2 = \{a^n b^n \mid n \geq 0\}$$

显然两者都不是正则语言, 但

$$L_1 \cap L_2 = \{\varepsilon\}$$

是正则语言.

定理 10 (差运算封闭性)

如果 L 和 M 都是正则语言, 那么 $L - M$ 也是正则的.

证明: $L - M = L \cap \overline{M}$.



例 12. 证明正则语言在以下运算下封闭

$$\min(L) = \{w \mid w \text{ is in } L, \text{ but no proper prefix of } w \text{ is in } L\}$$

证明 1: 设 L 的 DFA 为 $A = (Q, \Sigma, \delta, q_0, F)$, 构造 $\min(L)$ 的 DFA $B = (Q, \Sigma, \delta', q_0, F)$ 其中 δ' 如下, 往证 $L(B) = \min(L)$:

$$\delta'(q, a) = \begin{cases} \delta(q, a) & \text{if } q \notin F \\ \emptyset & \text{if } q \in F \end{cases}$$

- ① $\forall w \in L(B)$, 存在转移序列 $q_0 q_1 \cdots q_n \in F$ 使 B 接受 w , 其中 $q_i \notin F$ ($0 \leq i \leq n-1$). $\therefore w \in \min(L)$.
- ② $\forall w \in \min(L)$, 有 $w \in L$, A 接受 w 的状态序列为如果 $q_0 q_1 \cdots q_n \in F$, 则显然 $q_i \notin F$ ($0 \leq i \leq n-1$), 否则 w 会有 L 可接受的前缀.
 $\therefore w \in L(B)$



例 12. 证明正则语言在以下运算下封闭

$$\min(L) = \{w \mid w \text{ is in } L, \text{ but no proper prefix of } w \text{ is in } L\}$$

证明 2:

由封闭性

$$\min(L) = L - L\Sigma^+,$$

得证.



定义

字符串 $w = a_1a_2 \dots a_n$ 的**反转**, 记为 w^R , 定义为

$$w^R = a_na_{n-1} \dots a_1.$$

定义

语言 L 的**反转**, 记为 L^R , 定义为

$$L^R = \{w^R \in \Sigma^* \mid w \in L\}.$$

定理 11 (反转的封闭性)

如果 L 是正则语言, 那么 L^R 也是正则的.

两种证明方法:

- 对正则表达式 E 的结构归纳, 往证

$$\mathbf{L}(E^R) = (\mathbf{L}(E))^R.$$

- 构造识别 L 的 NFA $A = (Q, \Sigma, \delta_A, q_0, F)$, 将其转换为识别 L^R 的 NFA

$$B = (Q, \Sigma, \delta_B, q_s, \{q_0\})$$

- ① 将 A 的边调转方向;
- ② 将 A 的初始状态 q_0 , 改为唯一的接受状态;
- ③ 新增初始状态 q_s , 且令 $\delta_B(q_s, \varepsilon) = F$.

例 13. 语言 L 及其反转 L^R 分别为

$$L = \{w \in \{0,1\}^* \mid w \text{ ends in } 01.\}$$

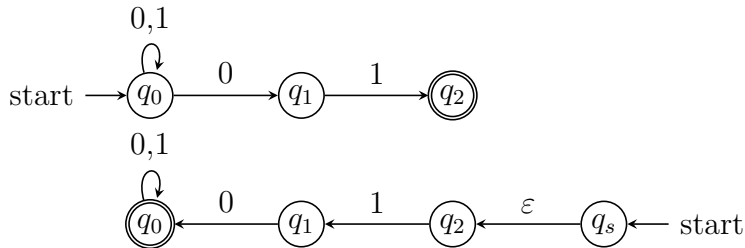
$$L^R = \{w \in \{0,1\}^* \mid w \text{ starts with } 10.\}$$

正则表达式分别为

$$L = (\mathbf{0 + 1})^* \mathbf{01}$$

$$L^R = \mathbf{10(0 + 1)^*}.$$

自动机分别为



证明: 往证如果有正则表达式 E , 则存在正则表达式 E^R 使

$$\mathbf{L}(E^R) = (\mathbf{L}(E))^R.$$

归纳基础:

- ① 当 $E = \emptyset$ 时, 有 $\emptyset^R = \emptyset$;
- ② 当 $E = \epsilon$ 时, 有 $\epsilon^R = \epsilon$;
- ③ $\forall a \in \Sigma$, 当 $E = \mathbf{a}$ 时, 有 $\mathbf{a}^R = \mathbf{a}$;

都满足 $\mathbf{L}(E^R) = (\mathbf{L}(E))^R$, 因此命题成立.

归纳递推:

① 当 $E = E_1 + E_2$ 时, 有 $(E_1 + E_2)^R = E_1^R + E_2^R$

$$(\mathbf{L}(E_1 + E_2))^R$$

$$= (\mathbf{L}(E_1) \cup \mathbf{L}(E_2))^R$$

$$= \{w^R \mid w \in \mathbf{L}(E_1) \cup w \in \mathbf{L}(E_2)\}$$

$$= (\mathbf{L}(E_1))^R \cup (\mathbf{L}(E_2))^R$$

$$= \mathbf{L}(E_1^R) \cup \mathbf{L}(E_2^R)$$

$$= \mathbf{L}(E_1^R + E_2^R)$$

正则表达式的加

语言的反转

同上

归纳假设

正则表达式的加

归纳递推:

① 当 $E = E_1 + E_2$ 时, 有 $(E_1 + E_2)^R = E_1^R + E_2^R$

② 当 $E = E_1 E_2$ 时, 有 $(E_1 E_2)^R = E_2^R E_1^R$

$$(\mathbf{L}(E_1 E_2))^R = (\mathbf{L}(E_1) \mathbf{L}(E_2))^R$$

$$= \{w_1 w_2 \mid w_1 \in \mathbf{L}(E_1), w_2 \in \mathbf{L}(E_2)\}^R$$

$$= \{(w_1 w_2)^R \mid w_1 \in \mathbf{L}(E_1), w_2 \in \mathbf{L}(E_2)\}$$

$$= \{w_2^R w_1^R \mid w_1 \in \mathbf{L}(E_1), w_2 \in \mathbf{L}(E_2)\}$$

$$= \{w_2^R \mid w_2 \in \mathbf{L}(E_2)\} \{w_1^R \mid w_1 \in \mathbf{L}(E_1)\}$$

$$= (\mathbf{L}(E_2))^R (\mathbf{L}(E_1))^R$$

$$= \mathbf{L}(E_2^R) \mathbf{L}(E_1^R) = \mathbf{L}(E_2^R E_1^R)$$

正则表达式的连接

语言的连接

语言的反转

字符串的反转

语言的连接

语言的反转

正则表达式的连接

归纳递推:

- ① 当 $E = E_1 + E_2$ 时, 有 $(E_1 + E_2)^R = E_1^R + E_2^R$
- ② 当 $E = E_1 E_2$ 时, 有 $(E_1 E_2)^R = E_2^R E_1^R$
- ③ 当 $E = E_1^*$ 时, 有 $(E_1^*)^R = (E_1^R)^*$

$$\begin{aligned} & (\mathbf{L}(E_1^*))^R \\ &= \{w_1 w_2 \dots w_n \mid n \geq 0, w_i \in \mathbf{L}(E_1)\}^R && \text{正则表达式的闭包} \\ &= \{(w_1 w_2 \dots w_n)^R \mid n \geq 0, w_i \in \mathbf{L}(E_1)\} && \text{语言的反转} \\ &= \{w_n^R w_{n-1}^R \dots w_1^R \mid n \geq 0, w_i \in \mathbf{L}(E_1)\} && \text{字符串的反转} \\ &= \{w_n^R w_{n-1}^R \dots w_1^R \mid n \geq 0, w_i^R \in \mathbf{L}(E_1^R)\} && \text{归纳假设} \\ &= \{w_1 w_2 \dots w_n \mid n \geq 0, w_i \in \mathbf{L}(E_1^R)\} && \text{变量重命名} \\ &= \mathbf{L}((E_1^R)^*) && \text{正则表达式的闭包} \end{aligned}$$

都满足 $(\mathbf{L}(E))^R = \mathbf{L}(E^R)$, 因此命题成立, 所以 L^R 也是正则语言.



同态

定义

若 Σ 和 Γ 是两个字母表, **同态** 定义为函数 $h: \Sigma \rightarrow \Gamma^*$

$$\forall a \in \Sigma, h(a) \in \Gamma^*.$$

扩展 h 的定义到字符串,

$$(1) \quad h(\varepsilon) = \varepsilon$$

$$(2) \quad h(xa) = h(x)h(a)$$

再扩展 h 到语言, 对 $\forall L \subseteq \Sigma^*$,

$$h(L) = \{h(w) \mid w \in L\}.$$

例 14. 若由 $\Sigma = \{0, 1\}$ 到 $\Gamma = \{a, b\}$ 的同态函数 h 为

$$h(0) = ab, \quad h(1) = \varepsilon.$$

则 Σ 上的字符串 0011, 在 h 的作用下

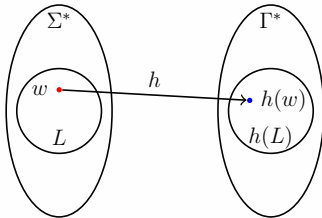
$$\begin{aligned} h(0011) &= h(\varepsilon)h(0)h(0)h(1)h(1) \\ &= \varepsilon \cdot ab \cdot ab \cdot \varepsilon \cdot \varepsilon \\ &= abab. \end{aligned}$$

语言 $L = 1^*0 + 0^*1$, 在 h 的作用下, $h(L)$ 为:

$$\begin{aligned} h(1^*0 + 0^*1) &= (h(1))^*h(0) + (h(0))^*h(1) \\ &= (\varepsilon)^*(ab) + (ab)^*(\varepsilon) \\ &= (ab)^* \end{aligned}$$

定理 12 (同态的封闭性)

若 L 是字母表 Σ 上的正则语言, h 是 Σ 上的同态, 则 $h(L)$ 也是正则的.



- 若 L 的正则表达式为 E , 即 $L = \mathbf{L}(E)$, 按如下规则构造表达式 $h(E)$

$$h(\emptyset) = \emptyset$$

$$h(\mathbf{r} + \mathbf{s}) = h(\mathbf{r}) + h(\mathbf{s})$$

$$h(\epsilon) = \epsilon$$

$$h(\mathbf{rs}) = h(\mathbf{r})h(\mathbf{s})$$

$$\forall a \in \Sigma, h(\mathbf{a}) = h(a)$$

$$h(\mathbf{r}^*) = (h(\mathbf{r}))^*$$

- 往证 $\mathbf{L}(h(E)) = h(\mathbf{L}(E))$, 而 $h(E)$ 显然也是正则表达式, 因此 $h(L)$ 正则

证明: 对 E 的结构归纳, 往证 $\mathbf{L}(h(E)) = h(\mathbf{L}(E))$.

归纳基础:

- 当 $E = \varepsilon$ 时

$$h(\mathbf{L}(\varepsilon)) = h(\{\varepsilon\}) = \{\varepsilon\} = \mathbf{L}(\varepsilon) = \mathbf{L}(h(\varepsilon))$$

- 当 $E = \emptyset$ 时

$$h(\mathbf{L}(\emptyset)) = h(\emptyset) = \emptyset = \mathbf{L}(\emptyset) = \mathbf{L}(h(\emptyset))$$

- $\forall a \in \Sigma$, 当 $E = \mathbf{a}$ 时

$$h(\mathbf{L}(\mathbf{a})) = h(\{a\}) = \{h(a)\} = \mathbf{L}(h(a)) = \mathbf{L}(h(\mathbf{a}))$$

所以命题成立.

归纳递推: 假设对正则表达式 F, G 分别有

$$\mathbf{L}(h(F)) = h(\mathbf{L}(F)), \quad \mathbf{L}(h(G)) = h(\mathbf{L}(G))$$

- 当 $E = F + G$ 时:

$$\begin{aligned} h(\mathbf{L}(F + G)) &= h(\mathbf{L}(F) \cup \mathbf{L}(G)) \\ &= h(\mathbf{L}(F)) \cup h(\mathbf{L}(G)) \\ &= \mathbf{L}(h(F)) \cup \mathbf{L}(h(G)) \\ &= \mathbf{L}(h(F) + h(G)) \\ &= \mathbf{L}(h(F + G)) \end{aligned}$$

正则表达式的加

h 作用在每个集合的串上

归纳假设

正则表达式的加

$h(F + G)$ 的定义

- 当 $E = FG$ 时: 略
- 当 $E = F^*$ 时: 略

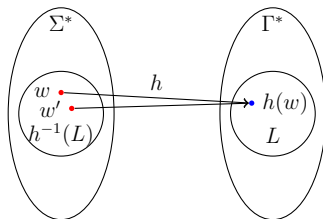


逆同态

定义

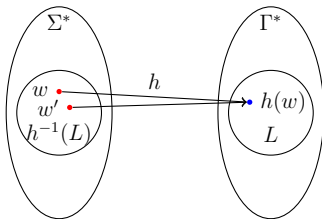
若 h 是字母表 Σ 到 Γ 的同态, 且 L 是 Γ 上的语言, 那么使 $h(w) \in L$ 的 w ($w \in \Sigma^*$) 的集合, 称为**语言 L 的 h 逆**, 记为 $h^{-1}(L)$, 即

$$h^{-1}(L) = \{w \in \Sigma^* \mid h(w) \in L\}.$$



定理 13 (逆同态的封闭性)

如果 h 是字母表 Σ 到 Γ 的同态, L 是 Γ 上的正则语言, 那么 $h^{-1}(L)$ 也是正则语言.



证明: 由 L 的 DFA $A = (Q, \Gamma, \delta, q_0, F)$, 构造识别 $h^{-1}(L)$ 的 DFA

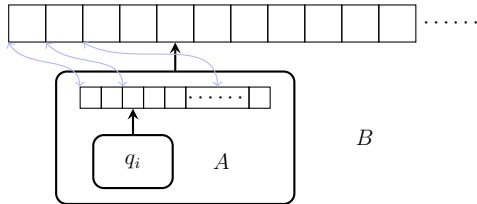
$$B = (Q, \Sigma, \delta', q_0, F),$$

证明: 由 L 的 DFA $A = (Q, \Gamma, \delta, q_0, F)$, 构造识别 $h^{-1}(L)$ 的 DFA

$$B = (Q, \Sigma, \delta', q_0, F),$$

其中

$$\delta'(q, a) = \hat{\delta}(q, h(a)).$$



为证明 $\mathbf{L}(B) = h^{-1}(L)$, 先证明 $\hat{\delta}'(q, w) = \hat{\delta}(q, h(w))$.

对 $|w|$ 归纳, 往证 $\hat{\delta}'(q, w) = \hat{\delta}(q, h(w))$.

① 归纳基础: 若 $w = \varepsilon$

$$\hat{\delta}(q, h(\varepsilon)) = \hat{\delta}(q, \varepsilon) = q = \hat{\delta}'(q, \varepsilon),$$

② 归纳递推: 若 $w = xa$

$\hat{\delta}'(q, xa) = \delta'(\hat{\delta}'(q, x), a)$	$\hat{\delta}'$ 定义
$= \delta'(\hat{\delta}(q, h(x)), a)$	归纳假设
$= \hat{\delta}(\hat{\delta}(q, h(x)), h(a))$	δ' 构造
$= \hat{\delta}(q, h(x)h(a))$	DFA 节例 5
$= \hat{\delta}(q, h(xa)).$	

所以 $\forall w \in \Sigma^*$, $\hat{\delta}'(q_0, w) = \hat{\delta}(q_0, h(w)) \in F$, 即 w 被 B 接受当且仅当 $h(w)$ 被 A 接受, B 是识别 $h^{-1}(L)$ 的 DFA, 因此 $h^{-1}(L)$ 是正则的. □

例 15. Prove that $L = \{0^n 1^{2n} \mid n \geq 0\}$ is a language not regular.

证明: 设同态 $h: \{0, 1\}^* \rightarrow \{0, 1\}^*$ 为

$$h(0) = 0,$$

$$h(1) = 11,$$

那么

$$h^{-1}(L) = \{0^n 1^n \mid n \geq 0\} = L_{01},$$

我们已知 L_{01} 非正则, 由封闭性, L 不是正则的.



例 16. 若语言 $L = (00 + 1)^*$, 同态 $h : \{a, b\} \rightarrow \{0, 1\}^*$ 为

$$h(a) = 01, \quad h(b) = 10,$$

请证明 $h^{-1}(L) = (\mathbf{ba})^*$.

证明: 往证 $h(w) \in L \iff w = (ba)^n$.

(\Leftarrow) 若 $w = (ba)^n$, 而 $h(ba) = 1001$, 因此 $h(w) = (1001)^n \in L$.

(\Rightarrow) 若 $h(w) \in L$, 假设 $w \notin (\mathbf{ba})^*$, 则只能有四种情况:

- ① w 以 a 开头, 则 $h(w)$ 以 01 开头, 显然 $h(w) \notin (00 + 1)^*$;
- ② w 以 b 结尾, 则 $h(w)$ 以 10 结尾, 显然 $h(w) \notin (00 + 1)^*$;
- ③ w 有连续的 a , 即 $w = xaay$, 则 $h(w) = z1010v$, 显然 $h(w) \notin (00 + 1)^*$;
- ④ w 有连续的 b , 即 $w = xbbby$, 则 $h(w) = z0101v$, 显然 $h(w) \notin (00 + 1)^*$;

因此 w 只能是 $(ba)^n, n \geq 0$ 的形式.



例 17. For a language L , define $\text{head}(L)$ to be the set of all prefixes of strings in L . Prove that if L is regular, so is $\text{head}(L)$.

证明. 设 L 是 Σ 上的正则语言且 $\Sigma = \{0, 1\}$, $\Gamma = \{0, 1, a, b\}$.
定义同态 $h: \Gamma \rightarrow \Sigma^*$ 和 $g: \Gamma \rightarrow \Sigma^*$ 分别为:

$$\begin{array}{llll} h(0) = 0 & h(a) = 0 & g(0) = 0 & g(a) = \varepsilon \\ h(1) = 1 & h(b) = 1 & g(1) = 1 & g(b) = \varepsilon \end{array}$$

则因为 $(\mathbf{0} + \mathbf{1})^*(\mathbf{a} + \mathbf{b})^*$ 是 Γ 上的正则语言, 所以

$$(\mathbf{0} + \mathbf{1})^*(\mathbf{a} + \mathbf{b})^* \cap h^{-1}(L)$$

是 Γ 上的正则语言, 所以

$$\text{head}(L) = g((\mathbf{0} + \mathbf{1})^*(\mathbf{a} + \mathbf{b})^* \cap h^{-1}(L))$$

是 Σ 上的正则语言, 因此 $\text{head}(L)$ 是正则的.



正则语言的性质

- 证明语言的非正则性
- 正则语言的封闭性
- 正则语言的判定性质
 - 空性, 有穷性和无穷性
 - 等价性
- 自动机的最小化

正则语言的判定性质

正则语言, 或任何语言, 典型的 3 个判定问题:

- ① 以某种形式化模型描述的语言是否为空? 是否无穷?
- ② 某个特定的串 w 是否属于所描述的语言?
- ③ 以两种方式描述的语言, 是否是相同的? — 语言的等价性

我们想知道, 要回答这类问题的具体算法, 是否存在.

空性, 有穷性和无穷性

定理 14

具有 n 个状态的有穷自动机 M 接受的集合 S :

- ① S 是非空的, 当且仅当 M 接受某个长度小于 n 的串;
- ② S 是无穷的, 当且仅当 M 接受某个长度为 m 的串, $n \leq m < 2n$.

所以, 对于正则语言:

- 存在算法, 判断其是否为空, 只需检查全部长度小于 n 的串;
- 存在算法, 判断其是否无穷, 只需检查全部长度由 n 到 $2n - 1$ 的串.

证明: 设接受正则语言 S 的 DFA 为 A .

① 必要性: 显然成立. 充分性:

❶ 如果 S 非空, 设 w 是 A 接受的串中长度最小者之一;

❷ 必然 $|w| < n$, 否则由泵引理 $w = xyz$, 接受 xz 更短.

② 必要性: 由泵引理, 显然成立. 充分性:

❶ 如果 S 无穷, 假设没有长度 n 到 $2n - 1$ 之间的串;

❷ 那么取 $w \in L(A)$ 是长度 $\geq 2n$ 中最小者之一;

❸ 由泵引理 $w = xyz$, 且 A 会接受更短的串 xz ;

❹ 于是, 或者 w 不是长度最小的, 或者长度 n 到 $2n - 1$ 之间有被接受的串, 因此假设不成立. □

正则语言的等价性

定理 15

存在算法, 判定两个有穷自动机是否等价(接受语言相同).

证明:

- ① 设 M_1 和 M_2 是分别接受 L_1 和 L_2 的有穷自动机;
- ② 则 $(L_1 \cap \overline{L_2}) \cup (\overline{L_1} \cap L_2)$ 是正则的, 所以可被某个有穷自动机 M_3 接受;
- ③ 而 M_3 接受某个串, 当且仅当 $L_1 \neq L_2$;
- ④ 由于存在算法判断 $L(M_3)$ 是否为空, 因此得证. □

正则语言的性质

- 证明语言的非正则性
- 正则语言的封闭性
- 正则语言的判定性质
- 自动机的最小化
 - DFA 状态的等价性
 - 填表算法与 DFA 最小化

DFA 状态的等价性

定义

DFA $A = (Q, \Sigma, \delta, q_0, F)$ 中两个状态 p 和 q , 对 $\forall w \in \Sigma^*$:

$$\hat{\delta}(p, w) \in F \Leftrightarrow \hat{\delta}(q, w) \in F,$$

则称这两个状态是等价的, 否则称为可区分的.

- 等价性只要求 $\hat{\delta}(p, w)$ 和 $\hat{\delta}(q, w)$ 同时在或不在 F 中, 而不必是相同状态.

填表算法

递归寻找 DFA 中全部的可区分状态对:

- ① 如果 $p \in F$ 而 $q \notin F$, 则 $[p, q]$ 是可区分的;
- ② $\exists a \in \Sigma$, 如果

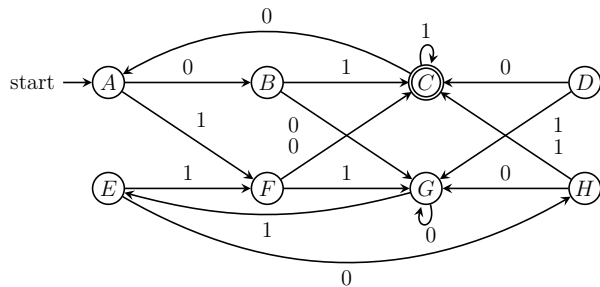
$$[r = \delta(p, a), s = \delta(q, a)]$$

是可区分的, 则 $[p, q]$ 是可区分的.

定理 16

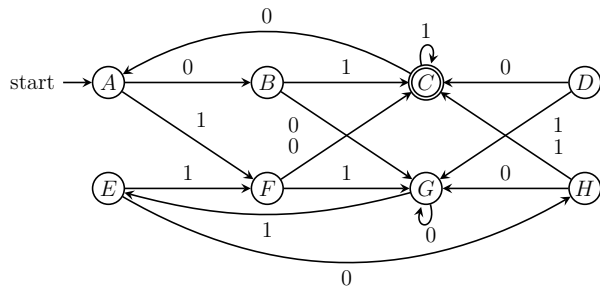
如果填表算法不能区分两个状态, 则这两个状态是等价的.

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.



<i>B</i>							
<i>C</i>							
<i>D</i>							
<i>E</i>							
<i>F</i>							
<i>G</i>							
<i>H</i>							
	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.

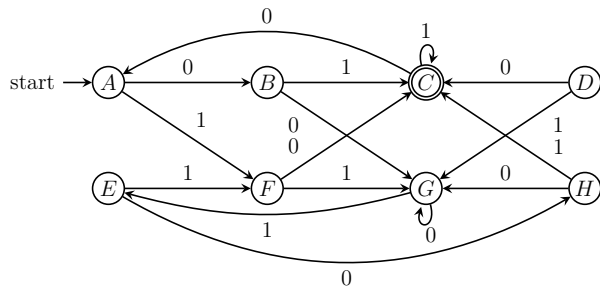


B							
C	×	×					
D			×				
E			×				
F			×				
G			×				
H			×				
	A	B	C	D	E	F	G

① 直接标记终态和非终态之间的状态对:

$$\{C\} \times \{A, B, D, E, F, G, H\}.$$

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.

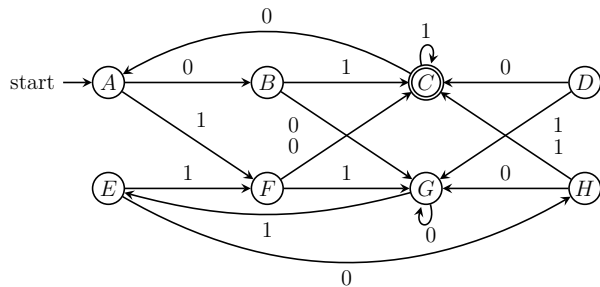


B							
C	×	×					
D	×	×	×				
E			×	×			
F	×	×	×		×		
G			×	×		×	
H			×	×		×	
	A	B	C	D	E	F	G

② 标记所有经过字符 0 到达终态和非终态的状态对:

$$\{D, F\} \times \{A, B, C, E, G, H\}.$$

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.

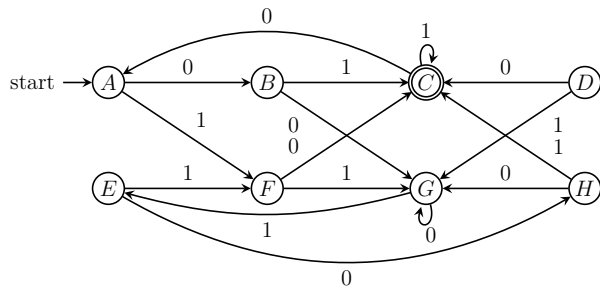


B	×						
C	×	×					
D	×	×	×				
E		×	×	×			
F	×	×	×		×		
G		×	×	×		×	
H	×		×	×	×	×	×
	A	B	C	D	E	F	G

③ 标记所有经过字符 1 到达终态和非终态的状态对:

$$\{B, H\} \times \{A, C, D, E, F, G\}.$$

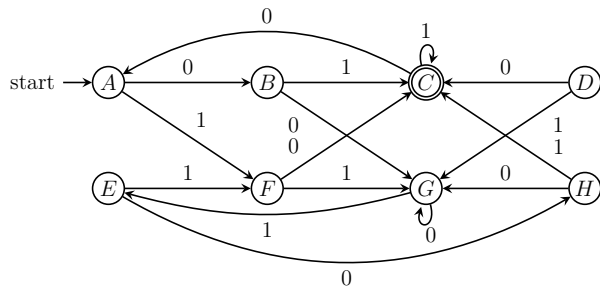
例 18. 用填表算法找到如图 DFA 中全部可区分状态对.



B	×						
C	×	×					
D	×	×	×				
E		×	×	×			
F	×	×	×		×		
G		×	×	×		×	
H	×		×	×	×	×	×
	A	B	C	D	E	F	G

④ 此时还有 $[A, E]$, $[A, G]$, $[B, H]$, $[D, F]$, $[E, G]$ 未标记, 只需逐个检查.

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.



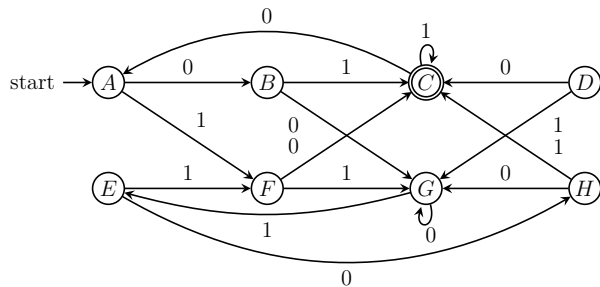
B	×						
C	×	×					
D	×	×	×				
E		×	×	×			
F	×	×	×		×		
G	×	×	×	×	×	×	
H	×		×	×	×	×	×
	A	B	C	D	E	F	G

④ 此时还有 $[A, E]$, $[A, G]$, $[B, H]$, $[D, F]$, $[E, G]$ 未标记, 只需逐个检查.

× $[A, G]$ 是可区分的, 因为经串 01 到可区分的 $[C, E]$;

× $[E, G]$ 是可区分的, 因为经串 10 到可区分的 $[C, H]$.

例 18. 用填表算法找到如图 DFA 中全部可区分状态对.



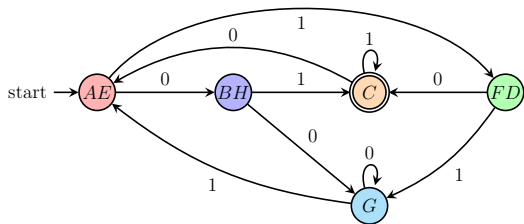
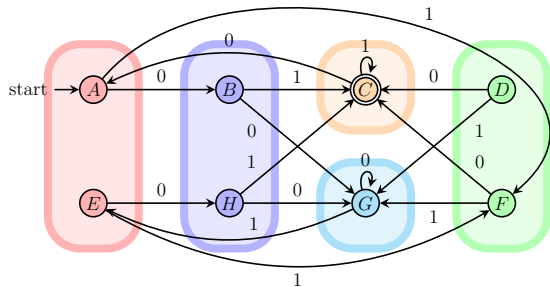
B	×						
C	×	×					
D	×	×	×				
E		×	×	×			
F	×	×	×		×		
G	×	×	×	×	×	×	
H	×		×	×	×	×	×
	A	B	C	D	E	F	G

- ⑤ 而 $[A, E]$, $[B, H]$ 和 $[D, F]$ 在经过很短的字符串后, 都会到达相同状态, 因此都是等价的.

DFA 最小化

根据等价状态, 将状态集划分成块, 构造等价的最小化 DFA.

续例 18. 构造其最小化的 DFA.



思考题

NFA 能否最小化?

形式语言与自动机理论

上下文无关文法

王春宇

计算机科学与技术学院
哈尔滨工业大学

上下文无关文法

- 上下文无关文法
 - 形式定义
 - 归约和派生
 - 最左派生和最右派生
 - 文法的语言
- 语法分析树
- 文法和语言的歧义性
- 文法的化简与范式

自然语言的文法

$\langle \text{sentence} \rangle \rightarrow \langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle$

$\langle \text{noun-phrase} \rangle \rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle \mid \langle \text{article} \rangle \langle \text{adjective} \rangle \langle \text{noun} \rangle$

$\langle \text{verb-phrase} \rangle \rightarrow \langle \text{verb} \rangle \mid \langle \text{verb} \rangle \langle \text{noun-phrase} \rangle$

$\langle \text{article} \rangle \rightarrow \text{a} \mid \text{the}$

$\langle \text{noun} \rangle \rightarrow \text{boy} \mid \text{girl} \mid \text{cat}$

$\langle \text{adjective} \rangle \rightarrow \text{big} \mid \text{small} \mid \text{blue}$

$\langle \text{verb} \rangle \rightarrow \text{sees} \mid \text{likes}$

...

自然语言的文法

使用语法规则产生句子:

$$\begin{aligned}\langle \text{sentence} \rangle &\Rightarrow \langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle \\ &\Rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle \langle \text{verb-phrase} \rangle \\ &\Rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle \langle \text{verb} \rangle \langle \text{noun-phrase} \rangle \\ &\Rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle \langle \text{verb} \rangle \langle \text{article} \rangle \langle \text{adjective} \rangle \langle \text{noun} \rangle \\ &\Rightarrow \text{the} \langle \text{noun} \rangle \langle \text{verb} \rangle \langle \text{article} \rangle \langle \text{adjective} \rangle \langle \text{noun} \rangle \\ &\Rightarrow \text{the girl} \langle \text{verb} \rangle \langle \text{article} \rangle \langle \text{adjective} \rangle \langle \text{noun} \rangle \\ &\Rightarrow \dots \\ &\Rightarrow \text{the girl sees a blue cat}\end{aligned}$$

定义

如果字符串 $w \in \Sigma^*$ 满足

$$w = w^R,$$

则称字符串 w 为回文 (*palindrome*).

定义

如果语言 L 中的字符串都是回文, 则称 L 为回文语言

$$L = \{w \in \Sigma^* \mid w = w^R\}.$$

- ε , 010, 0000, radar, racecar, drawkward
- A man, a plan, a canal — Panama
- 僧游云隐寺, 寺隐云游僧

例 1. 字母表 $\Sigma = \{0, 1\}$ 上的回文语言

$$L_{\text{pal}} = \{w \in \{0, 1\}^* \mid w = w^R\}.$$

- 很容易证明是 L_{pal} 是非正则的. 但如何表示呢?
- 可使用递归的方式来定义:
 - ① 首先 $\varepsilon, 0, 1$ 都是回文
 - ② 如果 w 是回文, $0w0$ 和 $1w1$ 也是回文
- 使用嵌套定义表示这种递归结构:

$$A \rightarrow \varepsilon \qquad A \rightarrow 0A0$$

$$A \rightarrow 0 \qquad A \rightarrow 1A1$$

$$A \rightarrow 1$$

上下文无关文法的形式定义

定义

上下文无关文法(CFG, Context-Free Grammar, 简称文法) G 是一个四元组

$$G = (V, T, P, S),$$

- ① V : 变元的有穷集, 变元也称为非终结符或语法范畴;
- ② T : 终结符的有穷集, 且 $V \cap T = \emptyset$;
- ③ P : 产生式的有穷集, 每个产生式包括:
 - i 一个变元, 称为产生式的头或左部;
 - ii 一个产生式符号 \rightarrow , 读作定义为;
 - iii 一个 $(V \cup T)^*$ 中的符号串, 称为体或右部;
- ④ $S \in V$: 初始符号, 文法开始的地方.

- 产生式 $A \rightarrow \alpha$, 读作 A 定义为 α
- 如果有多个 A 的产生式

$$A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$$

可简写为

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$$

- 文法中变元 A 的全体产生式, 称为 **A 产生式**

续例 1. 回文语言 $L_{\text{pal}} = \{w \in \{0, 1\}^* \mid w = w^R\}$ 的文法可设计为

$$G = (\{A\}, \{0, 1\}, \{A \rightarrow \varepsilon \mid 0 \mid 1 \mid 0A0 \mid 1A1\}, A).$$

字符使用的一般约定

- 终结符: $0, 1, \dots, a, b, \dots$
- 终结符串: \dots, w, x, y, z
- 非终结符: S, A, B, \dots
- 终结符或非终结符: \dots, X, Y, Z
- 终结符或非终结符组成的串: $\alpha, \beta, \gamma, \dots$

例2. 简化版的算数表达式:

- 运算只有“加”和“乘” $(+, *)$, 参数仅为标识符;
- 标识符: 以 $\{a, b\}$ 开头由 $\{a, b, 0, 1\}$ 组成的字符串.

这样的表达式集合可用文法 G_{exp} 表示

$$G_{\text{exp}} = (\{E, I\}, \{a, b, 0, 1, +, *, (,)\}, P, E),$$

其中产生式集 P 中有 10 条产生式

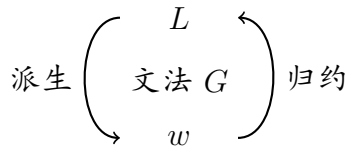
- | | | |
|--------------------------|-----------------------|------------------------|
| 1. $E \rightarrow I$ | 5. $I \rightarrow a$ | 9. $I \rightarrow I0$ |
| 2. $E \rightarrow E + E$ | 6. $I \rightarrow b$ | 10. $I \rightarrow I1$ |
| 3. $E \rightarrow E * E$ | 7. $I \rightarrow Ia$ | |
| 4. $E \rightarrow (E)$ | 8. $I \rightarrow Ib$ | |

注意, 变元 I 所定义的标识符集合, 刚好是 $(a + b)(a + b + 0 + 1)^*$.

归约和派生

非形式定义

从字符串到文法变元的分析过程, 称为递归推理或归约;
从文法变元到字符串的分析过程, 称为推导或派生.



- 归约: 自底向上, 由产生式的体向头的分析
- 派生: 自顶向下, 由产生式的头向体分析

续例2. 用算数表达式文法 G_{exp} , 将 $a * (a + b00)$ 归约的过程.

1. $E \rightarrow I$
2. $E \rightarrow E + E$
3. $E \rightarrow E * E$
4. $E \rightarrow (E)$
5. $I \rightarrow a$
6. $I \rightarrow b$
7. $I \rightarrow Ia$
8. $I \rightarrow Ib$
9. $I \rightarrow I0$
10. $I \rightarrow I1$

续例2. 用算数表达式文法 G_{exp} , 将 $a * (a + b00)$ 归约的过程.

	串归约到变元		应用产生式	重用结果	
1. $E \rightarrow I$	(1)	a	I	5. $I \rightarrow a$	—
2. $E \rightarrow E + E$	(2)	b	I	5. $I \rightarrow b$	—
3. $E \rightarrow E * E$	(3)	$b0$	I	9. $I \rightarrow I0$	(2)
4. $E \rightarrow (E)$	(4)	$b00$	I	9. $I \rightarrow I0$	(3)
5. $I \rightarrow a$	(5)	a	E	1. $E \rightarrow I$	(1)
6. $I \rightarrow b$	(6)	$b00$	E	1. $E \rightarrow I$	(4)
7. $I \rightarrow Ia$	(7)	$a + b00$	E	2. $E \rightarrow E + E$	(5), (6)
8. $I \rightarrow Ib$	(8)	$(a + b00)$	E	4. $E \rightarrow (E)$	(7)
9. $I \rightarrow I0$	(9)	$a * (a + b00)$	E	3. $E \rightarrow E * E$	(5), (8)
10. $I \rightarrow I1$					

派生和归约的形式定义

定义

若 CFG $G = (V, T, P, S)$, 设 $\alpha, \beta, \gamma \in (V \cup T)^*$, $A \in V$, $A \rightarrow \gamma \in P$, 那么称在 G 中由 $\alpha A \beta$ 可派生出 $\alpha \gamma \beta$, 记为

$$\alpha A \beta \xRightarrow{G} \alpha \gamma \beta.$$

相应的, 称 $\alpha \gamma \beta$ 可归约为 $\alpha A \beta$.

- $\alpha A \beta \xRightarrow{G} \alpha \gamma \beta$, 即用 $A \rightarrow \gamma$ 的右部 γ 替换串 $\alpha A \beta$ 中变元 A 得到串 $\alpha \gamma \beta$
- 如果语境中 G 是已知的, 可省略, 记为 $\alpha A \beta \Rightarrow \alpha \gamma \beta$

- 设 $\alpha_1, \dots, \alpha_m \in (V \cup T)^*$, $m \geq 1$, 对 $i = 1, \dots, m-1$ 如果有

$$\alpha_i \xRightarrow{G} \alpha_{i+1}$$

成立, 即 α_1 经过零步或多步派生可得到 α_m

$$\alpha_1 \xRightarrow{G} \alpha_2 \xRightarrow{G} \cdots \xRightarrow{G} \alpha_{m-1} \xRightarrow{G} \alpha_m,$$

那么, 记为

$$\alpha_1 \xRightarrow[G]{*} \alpha_m.$$

- 若 α 派生出 β 刚好经过了 i 步, 可记为

$$\alpha \xRightarrow[G]{i} \beta.$$

续例 2. 算数表达式 $a * (a + b00)$ 在文法 G_{exp} 中的派生过程.

续例 2. 算数表达式 $a * (a + b00)$ 在文法 G_{exp} 中的派生过程.

$$E \Rightarrow E * E \Rightarrow E * (E) \Rightarrow I * (E)$$

$$\Rightarrow I * (E + E) \Rightarrow I * (E + I) \Rightarrow I * (I + I)$$

$$\Rightarrow I * (a + I) \Rightarrow a * (a + I) \Rightarrow a * (a + I0)$$

$$\Rightarrow a * (a + I00) \Rightarrow a * (a + b00)$$

最左派生和最右派生

定义

为限制派生的随意性, 要求只替换符号串中最左边变元的派生过程, 称为**最左派生**, 记为

$$\Rightarrow_{\text{lm}}, \overset{*}{\Rightarrow}_{\text{lm}},$$

只替换最右的, 称为**最右派生**, 记为

$$\Rightarrow_{\text{rm}}, \overset{*}{\Rightarrow}_{\text{rm}}.$$

- 任何派生都有等价的最左派生和最右派生

$$A \Rightarrow^* w \text{ 当且仅当 } A \overset{*}{\Rightarrow}_{\text{lm}} w \text{ 当且仅当 } A \overset{*}{\Rightarrow}_{\text{rm}} w.$$

续例2. 表达式 $a * (a + a)$ 在 G_{exp} 中的最左派生和最右派生分别为:

1. $E \rightarrow I$	$E \xRightarrow{\text{lm}} E * E$	$E \xRightarrow{\text{rm}} E * E$
2. $E \rightarrow E + E$	$\xRightarrow{\text{lm}} I * E$	$\xRightarrow{\text{rm}} E * (E)$
3. $E \rightarrow E * E$	$\xRightarrow{\text{lm}} a * E$	$\xRightarrow{\text{rm}} E * (E + E)$
4. $E \rightarrow (E)$	$\xRightarrow{\text{lm}} a * (E)$	$\xRightarrow{\text{rm}} E * (E + I)$
5. $I \rightarrow a$	$\xRightarrow{\text{lm}} a * (E + E)$	$\xRightarrow{\text{rm}} E * (E + a)$
6. $I \rightarrow b$	$\xRightarrow{\text{lm}} a * (I + E)$	$\xRightarrow{\text{rm}} E * (I + a)$
7. $I \rightarrow Ia$	$\xRightarrow{\text{lm}} a * (a + E)$	$\xRightarrow{\text{rm}} E * (a + a)$
8. $I \rightarrow Ib$	$\xRightarrow{\text{lm}} a * (a + I)$	$\xRightarrow{\text{rm}} I * (a + a)$
9. $I \rightarrow I0$	$\xRightarrow{\text{lm}} a * (a + a)$	$\xRightarrow{\text{rm}} a * (a + a)$
10. $I \rightarrow I1$		

定义

CFG $G = (V, T, P, S)$ 的**语言**定义为

$$\mathbf{L}(G) = \{w \mid w \in T^*, S \xRightarrow{*}_G w\}.$$

那么符号串 w 在 $\mathbf{L}(G)$ 中, 要满足:

- ① w 仅由终结符组成;
- ② 初始符号 S 能派生出 w .

上下文无关语言

定义

语言 L 是某个 CFG G 定义的语言, 即 $L = L(G)$, 则称 L 为上下文无关语言 (CFL, Context-Free Language).

- 上下文无关是指在文法派生的每一步

$$\alpha A \beta \Rightarrow \alpha \gamma \beta,$$

符号串 γ 仅根据 A 的产生式派生, 而无需依赖 A 的上下文 α 和 β .

文法的等价性

定义

如果有两个文法 $CFG\ G_1$ 和 $CFG\ G_2$, 满足

$$L(G_1) = L(G_2),$$

则称 G_1 和 G_2 是等价的.

句型

定义

若 CFG $G = (V, T, P, S)$, 初始符号 S 派生出来的符号串, 称为 G 的句型, 即

$$\alpha \in (V \cup T)^* \text{ 且 } S \xRightarrow{*} \alpha.$$

如果 $S \xRightarrow[\text{lm}]{*} \alpha$, 称 α 为左句型. 如果 $S \xRightarrow[\text{rm}]{*} \alpha$, 称 α 为右句型.

- 只含有终结符的句型, 也称为 G 的句子
- 而 $L(G)$ 就是文法 G 全部的句子

例3. 给出语言 $L = \{w \in \{0, 1\}^* \mid w \text{ contains at least three 1s}\}$ 的文法.

解: $S \rightarrow A1A1A1A, A \rightarrow 0A \mid 1A \mid \varepsilon$

例 4. 描述 CFG $G = (\{S\}, \{a, b\}, \{S \rightarrow aSb, S \rightarrow ab\}, S)$ 定义的语言?

解: $L(G) = \{a^n b^n \mid n \geq 1\}$, 因为 $S \Rightarrow aSb \Rightarrow \cdots \Rightarrow a^{n-1} S b^{n-1} \Rightarrow a^n b^n$.

例5. 请为语言 $L = \{0^n 1^m \mid n \neq m\}$ 设计文法.

解:

$S \rightarrow AC \mid CB$	$A \rightarrow A0 \mid 0$
$C \rightarrow 0C1 \mid \varepsilon$	$B \rightarrow 1B \mid 1$

例 6. 设计 $L_{\text{eq}} = \{w \in \{0,1\}^* \mid w \text{ 中 } 0 \text{ 和 } 1 \text{ 个数相等}\}$ 的文法.

解 1: $S \rightarrow 0S1 \mid 1S0 \mid SS \mid \epsilon$, 寻找递归结构, 用变量构造递归结构;

解 2: $S \rightarrow S0S1S \mid S1S0S \mid \epsilon$, “目标串”这样构成, 由变量定义变量.

例 7. 设计 $L_{j \geq 2i} = \{a^i b^j \mid j \geq 2i\}$ 的文法.

解:
$$\begin{array}{l} S \rightarrow AB \\ A \rightarrow aAbb \mid \varepsilon \\ B \rightarrow bB \mid \varepsilon \end{array} \quad \text{或} \quad \begin{array}{l} S \rightarrow aSbb \mid B \\ B \rightarrow \varepsilon \mid bB \end{array}$$

程序设计语言的文法定义

- C — ISO C 1999 definition

...

selection-statement:

if (expression) statement

if (expression) statement else statement

switch (expression) statement

...

- Python — Full Grammar specification

...

```
try_stmt: ('try' ':' suite
          ((except_clause ':' suite)+
           ['else' ':' suite]
           ['finally' ':' suite] |
           'finally' ':' suite))
```

...

例 8. [Exe. 5.1.3] Show that every regular language is a context-free language.

例 8. [Exe. 5.1.3] Show that every regular language is a context-free language.

证明：对正则表达式 R 中运算符的个数 n 进行归纳.

归纳基础：当 $n = 0$ 时, R 只能是 ε , \emptyset 或 a ($a \in \Sigma$), 可以构造仅有一条产生式的文法 $S \rightarrow \varepsilon$, $S \rightarrow \emptyset$ 或 $S \rightarrow a$ 得到.

归纳递推：假设当 $n \leq m$ 时成立. 当 $n = m + 1$ 时, R 的形式只能由表达式 R_1 和 R_2 由连接、并或闭包形成：

- 若 $R = R_1 + R_2$, 则 R_1 和 R_2 中运算符都不超过 m , 所以都存在文法 G_1 和 G_2 , 分别开始于 S_1 和 S_2 , 只需构造新产生式和开始符号 $S \rightarrow S_1 | S_2$, 连同 G_1 和 G_2 的产生式, 构成 R 的文法;
- 若 $R = R_1 R_2$, 则同理构造 $S \rightarrow S_1 S_2$ 即可;
- 若 $R = R_1^*$, 则构造 $S \rightarrow S S_1 | \varepsilon$ 即可.

且每种构造, 文法的语言与该表达式的语言等价.



例 9. [Exe. 5.1.5] Let $T = \{0, 1, (,), +, *, \emptyset, e\}$. We may think of T as the set of symbols used by regular expressions over the alphabet $\{0, 1\}$; the only difference is that we use e for symbol ε , to avoid potential confusion in what follows. Your task is to design a CFG with set of terminals T that generates exactly the regular expressions with alphabet $\{0, 1\}$.

例 9. [Exe. 5.1.5] Let $T = \{0, 1, (,), +, *, \emptyset, e\}$. We may think of T as the set of symbols used by regular expressions over the alphabet $\{0, 1\}$; the only difference is that we use e for symbol ε , to avoid potential confusion in what follows. Your task is to design a CFG with set of terminals T that generates exactly the regular expressions with alphabet $\{0, 1\}$.

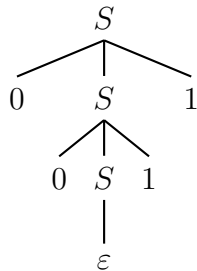
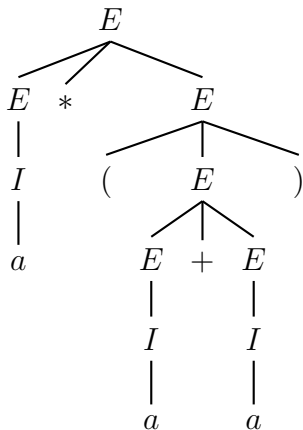
解: $S \rightarrow S + S \mid SS \mid S^* \mid (S) \mid 0 \mid 1 \mid \emptyset \mid e$.

上下文无关文法

- 上下文无关文法
- 语法分析树
 - 形式定义
 - 语法树和派生的等价性
- 文法和语言的歧义性
- 文法的化简与范式

派生或归约的过程可以表示成树形结构.

- 例2 文法 G_{exp} 中推导算数表达式 $a * (a + a)$ 的过程
- 例6 中语言 L_{eq} 的文法中推导 0011 的过程



语法分析树的形式定义

定义

CFG $G = (V, T, P, S)$ 的**语法分析树**(语法树或派生树) 为:

- ① 每个内节点标记为 V 中的变元符号;
- ② 每个叶节点标记为 $V \cup T \cup \{\varepsilon\}$ 中的符号;
- ③ 如果某内节点标记是 A , 其子节点从左至右分别为

$$X_1, X_2, \dots, X_n$$

那么

$$A \rightarrow X_1 X_2 \dots X_n \in P,$$

若有 $X_i = \varepsilon$, 则 ε 是 A 唯一子节点, 且 $A \rightarrow \varepsilon \in P$.

定义

语法树的全部叶节点从左到右连接起来, 称为该树的产物或结果.

- 如果树根节点是初始符号 S , 叶节点是终结符或 ε , 那么该树的产物属于 $L(G)$.

定义

语法树中标记为 A 的内节点及其全部子孙节点构成的子树, 称为 A 子树.

语法分析树和派生的等价性

定理 17

CFG $G = (V, T, P, S)$ 且 $A \in V$, 那么文法 G 中

$$A \xRightarrow{*} \alpha$$

当且仅当 G 中存在以 A 为根节点产物为 α 的语法树.

证明: [充分性] 对 $A \stackrel{i}{\Rightarrow} \alpha$ 的步骤数 j 归纳证明.

证明: [充分性] 对 $A \Rightarrow \alpha$ 的步骤数 j 归纳证明.

归纳基础: 当 $j = 1$ 时, 即 $A \Rightarrow \alpha$, 那么有 $A \rightarrow \alpha \in P$, 可构造 \bigwedge_{α}^A .

归纳递推: 假设 $j \leq n$ 时命题成立. 当 $j = n + 1$ 时, $A \xRightarrow{n+1} \alpha$ 的派生过程为

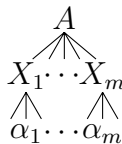
$$A \Rightarrow X_1 \cdots X_m \xRightarrow{n} \alpha_1 \cdots \alpha_m = \alpha,$$

即第 1 步一定由某产生式 $A \rightarrow X_1 X_2 \cdots X_m \in P$ 派生.

而 X_i 若非终结符, 一定有 $X_i \Rightarrow \alpha_i$ 且不超过 n 步, 由归纳假

设存在语法树 $\bigwedge_{\alpha_i}^{X_i}$. 因此可以构造以 A 为根, 以 X_i 为子树

(或叶子) 的语法树, 其产物刚好为 α .



[必要性] 对语法分析树的内节点数 j 归纳证明.

[必要性] 对语法分析树的内节点数 j 归纳证明.

归纳基础: 当 $j = 1$ 时, 即 $\begin{matrix} A \\ \swarrow \searrow \\ \alpha \end{matrix}$, A 必为根, 则 $A \rightarrow \alpha \in P$, 所以 $A \xRightarrow{*} \alpha$.

归纳递推: 假设 $j \leq n$ 时命题成立. 当 $j = n + 1$ 时, 根节点 A 的儿子依次为 X_1, X_2, \dots, X_m , 则

$$A \rightarrow X_1 \cdots X_m \in P, \text{ 且 } A \Rightarrow X_1 \cdots X_m.$$

而 X_i 子树 (或叶子) 内节点数都不超过 n , 由归纳假设有

$$X_i \xRightarrow{*} \alpha_i,$$

从左至右连接 α_i , 刚好为树的产物 α , 所以有

$$X_1 X_2 \cdots X_m \xRightarrow{*} \alpha_1 X_2 \cdots X_m \xRightarrow{*} \cdots \xRightarrow{*} \alpha_1 \alpha_2 \cdots \alpha_m = \alpha.$$

因此 $A \xRightarrow{*} \alpha$ 命题成立.



语法树唯一确定最左 (右) 派生

- 每棵语法分析树都有唯一的最左 (右) 派生
- 给定 CFG $G = (V, T, P, S)$, $A \in V$, 以下命题等价:
 - ① 通过递归推理, 确定串 w 在变元 A 的语言中.
 - ② 存在以 A 为根节点, 产物为 w 的语法分析树.
 - ③ $A \xRightarrow{*} w$.
 - ④ $A \xRightarrow[\text{lm}]{*} w$.
 - ⑤ $A \xRightarrow[\text{rm}]{*} w$.

例 10. [Exe. 5.2.2] Suppose that G is a CFG without any productions that have ε as the right side. If w is in $L(G)$, the length of w is n , and w has a derivation of m steps, show that w has a parse tree with $n + m$ nodes.

例 10. [Exe. 5.2.2] Suppose that G is a CFG without any productions that have ε as the right side. If w is in $L(G)$, the length of w is n , and w has a derivation of m steps, show that w has a parse tree with $n + m$ nodes.

证明:

- ① 派生 w 的每一步推导都对应语法树的一个内节点, 所以 w 语法树中共有 m 个内节点;
- ② 每个 w 的终结符都构成一个叶节点, 所以至少有 n 个叶节点, 而由于 G 中没有空产生式, 因此不会有标记为 ε 的叶节点, 所以只能有 n 个叶节点. 所以 w 的语法树有 $n + m$ 个节点. □

例 11. [Exe. 5.2.3] Suppose all is as in Exercise 5.2.2, but G may have some productions with ε as the right side. Show that a parse tree for a string w other than ε may have as many as $n + 2m - 1$ nodes, but no more.

例 11. [Exe. 5.2.3] Suppose all is as in Exercise 5.2.2, but G may have some productions with ε as the right side. Show that a parse tree for a string w other than ε may have as many as $n + 2m - 1$ nodes, but no more.

证明:

- ① 派生 w 的每一步推导都对应语法树的一个内节点, 所以 w 语法树中共有 m 个内节点.
- ② 每个 w 的终结符都构成一个叶节点, 所以至少有 n 个叶节点.
- ③ 推导过程中, 每次空产生式的应用, 都会增加一个标记 ε 的叶节点, 但显然不能全部的 m 步都使用空产生式, 所以最多增加 $m - 1$ 个 ε 叶节点. 因此 w 的语法树有最多 $m + n + m - 1 = n + 2m - 1$ 个节点. □

上下文无关文法

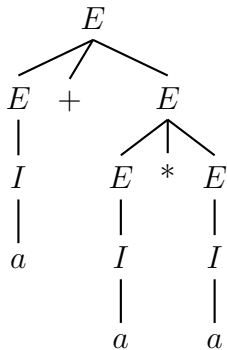
- 上下文无关文法
- 语法分析树
- 文法和语言的歧义性
 - 文法歧义性的消除
 - 语言的固有歧义性
- 文法的化简与范式

文法的歧义性

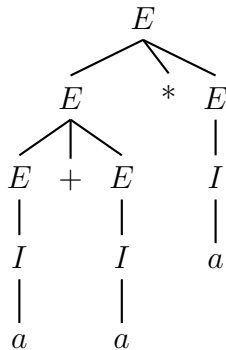
定义

如果 CFG G 使某些符号串有两棵不同的语法分析树, 称文法 G 是歧义的.

续例2. 算数表达式的文法 G_{exp} 中, 对句型 $a + a * a$ 有下面两棵语法分析树:



$$\begin{aligned}
 (1) \quad E &\Rightarrow E + E \\
 &\Rightarrow E + E * E \\
 &\stackrel{*}{\Rightarrow} a + a * a
 \end{aligned}$$



$$\begin{aligned}
 (2) \quad E &\Rightarrow E * E \\
 &\Rightarrow E + E * E \\
 &\stackrel{*}{\Rightarrow} a + a * a
 \end{aligned}$$

文法歧义性的消除

有些文法的歧义性, 可以通过重新设计文法来消除.

续例2. 文法 G_{exp} 重新设计为文法 $G_{\text{exp}'}$ 可消除歧义.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid I$$

$$I \rightarrow a$$

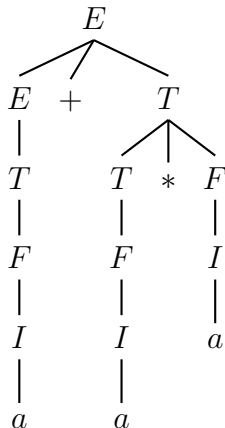
$$I \rightarrow b$$

$$I \rightarrow Ia$$

$$I \rightarrow Ib$$

$$I \rightarrow I0$$

$$I \rightarrow I1$$



语言的固有歧义性

定义

定义同样的语言可以有多个文法, 如果 CFL L 的所有文法都是歧义的, 那么称语言 L 是固有歧义的.

- 固有歧义的语言确实存在, 如语言

$$L = \{a^i b^j c^k \mid i = j \text{ or } j = k\}$$

中任何形为 $a^n b^n c^n$ 的串, 总会有两棵语法树.

- “判定任何给定 CFG G 是否歧义”是一个不可判定问题.

上下文无关文法

- 上下文无关文法
- 语法分析树
- 文法和语言的歧义性
- 文法的化简与范式
 - 消除无用符号
 - 消除 ε -产生式
 - 消除单元产生式
 - 乔姆斯基范式
 - 格雷巴赫范式

为什么要化简

- 典型问题: 给定 CFG G 和串 w , 判断 $w \in \mathbf{L}(G)$?
- 编译器设计和自然语言处理的基本问题之一
- 但文法的形式非常自由, 过于复杂不易于自动处理
- 以不改变语言为前提, 化简文法和限制文法的格式

例 12. 如下文法中, 有无意义的变元和产生式

$$S \rightarrow 0DS1D \mid B \mid \varepsilon$$

$$B \rightarrow BC1 \mid 0CBC$$

$$A \rightarrow A0 \mid A1 \mid \varepsilon$$

$$C \rightarrow D$$

$$D \rightarrow \varepsilon$$

文法的化简

- ❶ 消除无用符号: 对文法定义语言没有贡献的符号
- ❷ 消除 ε 产生式: $A \rightarrow \varepsilon$ (得到语言 $L - \{\varepsilon\}$)
- ❸ 消除单元产生式: $A \rightarrow B$

无用符号

定义

CFG $G = (V, T, P, S)$, 符号 $X \in (V \cup T)$:

- ① 如果 $S \xRightarrow{*} \alpha X \beta$, 称 X 是可达的;
- ② 如果 $\alpha X \beta \xRightarrow{*} w$ ($w \in T^*$), 称 X 是产生的;
- ③ 如果 X 同时是产生的和可达的, 即

$$S \xRightarrow{*} \alpha X \beta \xRightarrow{*} w \quad (w \in T^*),$$

则称 X 是有用的, 否则称 X 为无用符号.

消除无用符号

消除无用符号: 删除全部含有“非产生的”和“非可达的”符号的产生式

计算“产生的”符号集

- ① 每个 T 中的符号都是产生的;
- ② $A \rightarrow \alpha \in P$ 且 α 中符号都是产生的, 则 A 是产生的.

计算“可达的”符号集

- ① 符号 S 是可达的;
- ② $A \rightarrow \alpha \in P$ 且 A 是可达的, 则 α 中符号都是可达的.

定理 18

每个非空的 CFL 都能被一个不带无用符号的 CFG 定义.

注意

- 先寻找并消除全部非“产生的”符号
- 再寻找并消除全部非“可达的”符号
- 否则可能消除不完整

例 13. 消除如下文法无用符号

$$S \rightarrow AB \mid a$$

$$A \rightarrow b$$

解：先消除非产生的

$$S \rightarrow a$$

$$A \rightarrow b$$

再消除非可达的

$$S \rightarrow a$$

消除 ε -产生式

定义

文法中形如 $A \rightarrow \varepsilon$ 的产生式称为 ε -产生式.

如果变元 $A \xRightarrow{*} \varepsilon$, 称 A 是 **可空的**.

- ε -产生式在文法定义语言时, 除产生空串外没有其他帮助
- 对于 CFL L , 消除其文法中全部的 ε -产生式后, 得到语言 $L - \{\varepsilon\}$

确定“可空变元”

- ① 如果 $A \rightarrow \varepsilon$, 则 A 是可空的;
- ② 如果 $B \rightarrow \alpha$ 且 α 中的每个符号都是可空的, 则 B 是可空的.

替换产生式

将含有可空变元的一条产生式 $A \rightarrow X_1X_2\cdots X_n$,
用一组产生式 $A \rightarrow Y_1Y_2\cdots Y_n$ 代替, 其中:

- ① 若 X_i 不是可空的, Y_i 为 X_i ;
- ② 若 X_i 是可空的, Y_i 为 X_i 或 ε ;
- ③ 但 Y_i 不能全为 ε .

定理 19

任何 CFG G , 都存在一个不带无用符号和 ε -产生式的 CFG G' , 使 $\mathbf{L}(G') = \mathbf{L}(G) - \{\varepsilon\}$.

例 14. 消除 CFG $G = (\{S, A, B\}, \{a, b\}, P, S)$ 的 ε -产生式.

$$S \rightarrow AB$$

$$A \rightarrow AaA \mid \varepsilon$$

$$B \rightarrow BbB \mid \varepsilon$$

解: CFG G' 为

$$S \rightarrow AB \mid A \mid B$$

$$A \rightarrow AaA \mid Aa \mid aA \mid a$$

$$B \rightarrow BbB \mid Bb \mid bB \mid b$$

消除单元产生式

确定“单元对”

如果有 $A \Rightarrow B$, 则称 $[A, B]$ 为单元对.

- ① $A \rightarrow B \in P$, 则 $[A, B]$ 是单元对;
- ② 若 $[A, B]$ 和 $[B, C]$ 都是单元对, 则 $[A, C]$ 是单元对.

消除单元产生式

- ① 删除全部形为 $A \rightarrow B$ 的单元产生式;
- ② 对每个单元对 $[A, B]$, 将 B 的产生式复制给 A .

定理 20

每个不带 ε 的 *CFL* 都可由一个不带无用符号, ε -产生式和单元产生式的文法定义.

例 15. 消除文法的单元产生式

$$S \rightarrow A \mid B \mid 0S1$$

$$A \rightarrow 0A \mid 0$$

$$B \rightarrow 1B \mid 1$$

解: 单位对为 $[S, A]$ 和 $[S, B]$, 带入得:

$$S \rightarrow 0S1 \qquad A \rightarrow 0A \mid 0$$

$$S \rightarrow 0A \mid 0 \qquad B \rightarrow 1B \mid 1$$

$$S \rightarrow 1B \mid 1$$

文法化简的可靠顺序

- ① 消除 ε -产生式;
- ② 消除单元产生式;
- ③ 消除非产生的无用符号;
- ④ 消除非可达的无用符号.

例 16. [Exe. 7.1.2] Begin with the grammar:

$$S \rightarrow ASB \mid \varepsilon$$

$$A \rightarrow aAS \mid a$$

$$B \rightarrow SbS \mid A \mid bb$$

- ❶ Eliminate ε -productions.
- ❷ Eliminate any unit productions in the resulting grammar.
- ❸ Eliminate any useless symbols in the resulting grammar.

限制文法格式

将任意形式的文法转换为:

- ① 乔姆斯基范式 (CNF, Chomsky Normal Form)
- ② 格雷巴赫范式 (GNF, Greibach Normal Form)

乔姆斯基范式

定理 21 (乔姆斯基范式 CNF)

每个不带 ε 的 CFL 都可由这样的 CFG G 定义, G 中每个产生式都形为

$$A \rightarrow BC \text{ 或 } A \rightarrow a$$

其中 A, B 和 C 都是变元, a 是终结符.

- 利用 CNF 派生长度为 n 的串, 刚好需要 $2n - 1$ 步
- 因此存在算法判断任意字符串 w 是否在给定的 CFL 中
- 利用 CNF 的 CYK 算法 — $O(n^3)$ 时间复杂度的解析算法

CFG 转为 CNF 的方法

① 将产生式

$$A \rightarrow X_1 X_2 \cdots X_m \quad (m \geq 2)$$

中每个终结符 a 替换为新变元 C_a , 并增加新产生式

$$C_a \rightarrow a$$

② 引入新变元 D_1, D_2, \dots, D_{m-2} , 将产生式

$$A \rightarrow B_1 B_2 \cdots B_m \quad (m > 2)$$

替换为一组级联的产生式

$$A \rightarrow B_1 D_1$$

$$D_1 \rightarrow B_2 D_2$$

\dots

$$D_{m-2} \rightarrow B_{m-1} B_m$$

例 17. CFG $G = (\{S, A, B\}, \{a, b\}, P, S)$, 产生式集合 P 为:

$$S \rightarrow bA \mid aB$$

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

$$B \rightarrow aBB \mid bS \mid b$$

请设计等价的 CNF 文法.

解: CNF 为

$$S \rightarrow C_b A \mid C_a B$$

$$A \rightarrow C_a S \mid C_b D_1 \mid a \quad D_1 \rightarrow AA \quad C_a \rightarrow a$$

$$B \rightarrow C_b S \mid C_a D_2 \mid b \quad D_2 \rightarrow BB \quad C_b \rightarrow b$$

证明:

证明: 设 CFL L 不含 ε , 由定理 20, 存在不含 ε -产生式和单元产生式的等价文法 $G_1 = (V, T, P, S)$. 考虑 P 中一条产生式 $A \rightarrow X_1 X_2 \dots X_m$ ($m \geq 2$).

- ① 若某个 X_i 是终结符 a , 则引入新变元 C_a 和新产生式 $C_a \rightarrow a$, 并用 C_a 替换 X_i , 得文法 $G_2 = (V', T, P', S)$.
- ② 显然 $L(G_1) \subseteq L(G_2)$, 因为如果 $\alpha \xRightarrow{G_1} \beta$, 那么 $\alpha \xRightarrow{G_2} \beta$.
- ③ 用归纳法证明 $A \xRightarrow{G_2} w \implies A \xRightarrow{G_1} w$, 这里的 $A \in V$, $w \in T^*$.
 - i 当 $i = 1$ 时是显然的, 或者用了 P 中未修改的产生式, 或者用了被修改的产生式, 而二者都有 $A \xRightarrow{G_1} w$.
 - ii 假设当 $i \leq n$ 时命题成立. 当 $i = n + 1$ 时, $A \xRightarrow{G_2} w$ 的第 1 步, 必然使用了某个产生式 $A \rightarrow B_1 B_2 \dots B_m$, 即 $A \xRightarrow{G_2} B_1 B_2 \dots B_m \xRightarrow{G_2} w = w_1 w_2 \dots w_m$. 那么, 如果 $B_i \in V' - V$, B_i 一定是对应某个终结符 a_i 的 C_{a_i} , w_i 必然是 a_i , 令 $X_i = a_i$; 如果 $B_i \in V$, $B_i \xRightarrow{G_2} w_i$ 一定不超过 n 步, 由归纳假设, $B_i \xRightarrow{G_1} w_i$, 那么令 $X_i = B_i$. 由 P' 的结构, $A \rightarrow X_1 X_2 \dots X_m$ 是 P 的一条产生式, 所以 $A \xRightarrow{G_1} X_1 X_2 \dots X_m \xRightarrow{G_1} w_1 w_2 \dots w_m = w$.

所以 $L(G_2) \subseteq L(G_1)$.



格雷巴赫范式

定理 22 (格雷巴赫范式 GNF)

每个不带 ε 的 *CFL* 都可由这样的 *CFG* G 定义, G 中每个产生式都形为

$$A \rightarrow a\alpha$$

其中 A 是变元, a 是终结符, α 是零或多个变元的串.

- GNF 每个产生式都会引入一个终结符
- 长度为 n 的串的派生恰好是 n 步

例 18. 将以下文法转换为 GNF.

$$S \rightarrow AB$$

$$A \rightarrow aA \mid bB \mid b$$

$$B \rightarrow b$$

解: GNF 为

$$S \rightarrow aAB \mid bBB \mid bB$$

$$A \rightarrow aA \mid bB \mid b$$

$$B \rightarrow b$$

直接左递归

定义

文法中形式为 $A \rightarrow A\alpha$ 的产生式, 称为直接左递归.

消除直接左递归

① 若 A 产生式

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \cdots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m$$

其中 $\alpha_i \neq \varepsilon$, β_j 不以 A 开始;

② 引入新变元 B , 并用如下产生式替换

$$\begin{aligned} A &\rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m \mid \beta_1 B \mid \beta_2 B \mid \cdots \mid \beta_m B \\ B &\rightarrow \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_n \mid \alpha_1 B \mid \alpha_2 B \mid \cdots \mid \alpha_n B \end{aligned}$$

间接左递归

定义

文法中如果有形式为

$$A \rightarrow B\alpha \mid \dots$$

$$B \rightarrow A\beta \mid \dots$$

的产生式, 称为**间接左递归**.

- 会有 $A \Rightarrow B\alpha \Rightarrow A\beta\alpha$, 无法通过代换消除递归

消除间接左递归

- ① 将文法中变元重命名为 A_1, A_2, \dots, A_n ;
- ② 通过代入, 使产生式都形如

$$A_i \rightarrow A_j \alpha$$

$$A_i \rightarrow a \alpha$$

但要求 $i \leq j$;

- ③ 消除直接左递归 $A_i \rightarrow A_i \beta$, 再代入其他产生式.

例 19. Convert the following grammar to GNF.

$$S \rightarrow AB$$

$$A \rightarrow BS \mid b$$

$$B \rightarrow SA \mid a$$

解:

1. 重命名变元, 代换 $i > j$ 的 A_j

$$A_1 \rightarrow A_2 A_3$$

$$A_2 \rightarrow A_3 A_1 \mid b$$

$$A_3 \rightarrow a \mid \overline{A_1 A_2} \mid \overline{A_2 A_3 A_2} \mid \\ A_3 A_1 A_3 A_2 \mid b A_3 A_2$$

2. 消除直接左递归

$$A_1 \rightarrow A_2 A_3$$

$$A_2 \rightarrow A_3 A_1 \mid b$$

$$A_3 \rightarrow b A_3 A_2 \mid a \mid b A_3 A_2 B_1 \mid a B_1$$

$$B_1 \rightarrow A_1 A_3 A_2 \mid A_1 A_3 A_2 B_1$$

3. A_3 产生式代入到 A_2 , A_2 产生式代入到 A_1 , A_1 产生式代入 B_1

$$A_3 \rightarrow b A_3 A_2 \mid a \mid b A_3 A_2 B_1 \mid a B_1$$

$$A_2 \rightarrow b A_3 A_2 A_1 \mid a A_1 \mid b A_3 A_2 B_1 A_1 \mid a B_1 A_1 \mid b$$

$$A_1 \rightarrow b A_3 A_2 A_1 A_3 \mid a A_1 A_3 \mid b A_3 A_2 B_1 A_1 A_3 \mid a B_1 A_1 A_3 \mid b A_3$$

$$B_1 \rightarrow b A_3 A_2 A_1 A_3 A_3 A_2 \mid a A_1 A_3 A_3 A_2 \mid b A_3 A_2 B_1 A_1 A_3 A_3 A_2 \mid \\ a B_1 A_1 A_3 A_3 A_2 \mid b A_3 A_3 A_2 \mid b A_3 A_2 A_1 A_3 A_3 A_2 B_1 \mid a A_1 A_3 A_3 A_2 B_1 \mid \\ b A_3 A_2 B_1 A_1 A_3 A_3 A_2 B_1 \mid a B_1 A_1 A_3 A_3 A_2 B_1 \mid b A_3 A_3 A_2 B_1$$

GNF 引理 1

如果有文法 $G = (V, T, P, S)$, 设 $A \rightarrow \alpha_1 B \alpha_2$ 是 P 中的一个产生式, 且 $B \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$ 是 P 中的全部 B 产生式. 将产生式 $A \rightarrow \alpha_1 B \alpha_2$ 从 P 中删除, 并增加

$$A \rightarrow \alpha_1 \beta_1 \alpha_2 \mid \alpha_1 \beta_2 \alpha_2 \mid \cdots \mid \alpha_1 \beta_n \alpha_2$$

一组产生式, 得到文法 $G_1 = (V, T, P', S)$, 那么 $\mathbf{L}(G) = \mathbf{L}(G_1)$.

GNF 引理 1

如果有文法 $G = (V, T, P, S)$, 设 $A \rightarrow \alpha_1 B \alpha_2$ 是 P 中的一个产生式, 且 $B \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_n$ 是 P 中的全部 B 产生式. 将产生式 $A \rightarrow \alpha_1 B \alpha_2$ 从 P 中删除, 并增加

$$A \rightarrow \alpha_1 \beta_1 \alpha_2 \mid \alpha_1 \beta_2 \alpha_2 \mid \cdots \mid \alpha_1 \beta_n \alpha_2$$

一组产生式, 得到文法 $G_1 = (V, T, P', S)$, 那么 $L(G) = L(G_1)$.

证明:

- ① 显然 $L(G_1) \subseteq L(G)$, 因为 G_1 的派生中, 如果用到了 $A \rightarrow \alpha_1 \beta_i \alpha_2$, 在 G 中可以使用 $A \xRightarrow{G} \alpha_1 B \alpha_2 \xRightarrow{G} \alpha_1 \beta_i \alpha_2$.
- ② 而因为 $A \rightarrow \alpha_1 B \alpha_2$ 是唯一在 G 中而不再 G_1 中的产生式, 每当 G 的派生中用到了 $\alpha_1 B \alpha_2$ 时, 一定会在后面某一步中用到形如 $B \rightarrow \beta_i$ 的产生式来派生 B , 这两步在 G_1 中可以使用一步 $A \xRightarrow{G_1} \alpha_1 \beta_i \alpha_2$ 来代替, 所以 $L(G) \subseteq L(G_1)$. □

GNF 引理 2

如果有文法 $G = (V, T, P, S)$, 设带有直接左递归的 A 产生式为

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \cdots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m$$

其中 β_i 不以 A 开头. 在 V 中引入新的变元 B 并用以下产生式

$$A \rightarrow \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m \mid \beta_1 B \mid \beta_2 B \mid \cdots \mid \beta_m B$$

$$B \rightarrow \alpha_1 \mid \alpha_2 \mid \cdots \mid \alpha_n \mid \alpha_1 B \mid \alpha_2 B \mid \cdots \mid \alpha_n B$$

替换全部 A 产生式, 得到文法 $G_1 = (V \cup \{B\}, T, P', S)$, 那么 $\mathbf{L}(G) = \mathbf{L}(G_1)$.

证明: 在文法 G 中一系列使用 $A \rightarrow A\alpha_i$ 的最左派生, 最后必以 $A \rightarrow \beta_j$ 结束, 而这样的最左派生

$$\begin{aligned} A &\xRightarrow{\text{lm}} A\alpha_{i_1} \xRightarrow{\text{lm}} A\alpha_{i_2}\alpha_{i_1} \xRightarrow{\text{lm}} \cdots \\ &\xRightarrow{\text{lm}} A\alpha_{i_p}\alpha_{i_{p-1}} \cdots \alpha_{i_1} \\ &\xRightarrow{\text{lm}} \beta_j\alpha_{i_p}\alpha_{i_{p-1}} \cdots \alpha_{i_1}, \end{aligned}$$

在 G_1 中可以使用一系列最右派生来代替

$$\begin{aligned} A &\xRightarrow{\text{rm}} \beta_j B \xRightarrow{\text{rm}} \beta_j\alpha_{i_p} B \xRightarrow{\text{rm}} \beta_j\alpha_{i_p}\alpha_{i_{p-1}} B \xRightarrow{\text{rm}} \cdots \\ &\xRightarrow{\text{rm}} \beta_j\alpha_{i_p}\alpha_{i_{p-1}} \cdots \alpha_{i_2} B \\ &\xRightarrow{\text{rm}} \beta_j\alpha_{i_p}\alpha_{i_{p-1}} \cdots \alpha_{i_2}\alpha_{i_1}. \end{aligned}$$

而且, 相反的转换也成立, 因此 $\mathbf{L}(G) = \mathbf{L}(G_1)$.



形式语言与自动机理论

下推自动机

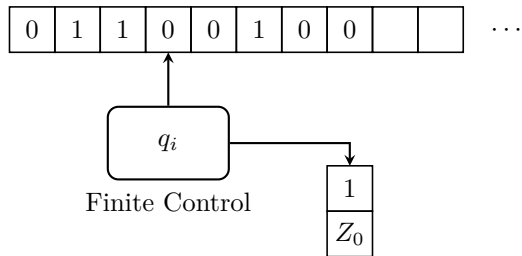
王春宇

计算机科学与技术学院
哈尔滨工业大学

下推自动机

- 下推自动机
 - 形式定义
 - 瞬时描述和转移符号
- 下推自动机接受的语言
- 下推自动机与文法的等价性
- 确定型下推自动机

下推自动机



下推自动机的形式定义

定义

下推自动机(*PDA*, *Pushdown Automata*) P 为七元组

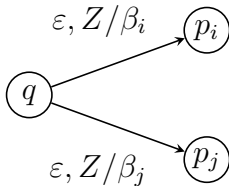
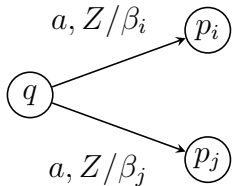
$$P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$$

- ① Q , 有穷状态集;
- ② Σ , 有穷输入符号集;
- ③ Γ , 有穷栈符号集;
- ④ $\delta : Q \times (\Sigma \cup \{\varepsilon\}) \times \Gamma \rightarrow 2^{Q \times \Gamma^*}$, 状态转移函数;
- ⑤ $q_0 \in Q$, 初始状态;
- ⑥ $Z_0 \in \Gamma - \Sigma$, 栈底符号;
- ⑦ $F \subseteq Q$, 接收状态集或终态集.

PDA 的动作和状态转移图

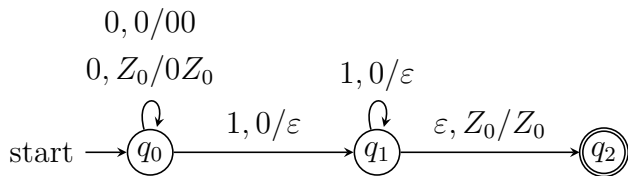
如果 $q, p_i \in Q$ ($1 \leq i \leq m$), $a \in \Sigma$, $Z \in \Gamma$, $\beta_i \in \Gamma^*$, 可以有动作:

$$\delta(q, a, Z) = \{(p_1, \beta_1), (p_2, \beta_2), \dots, (p_m, \beta_m)\}, \text{ 或}$$
$$\delta(q, \varepsilon, Z) = \{(p_1, \beta_1), (p_2, \beta_2), \dots, (p_m, \beta_m)\}.$$



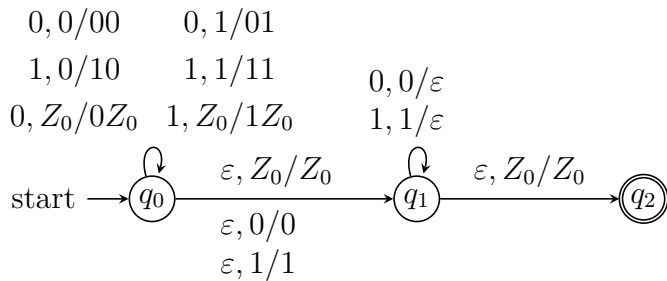
例 1. 设计识别 $L_{01} = \{0^n 1^n \mid n \geq 1\}$ 的 PDA.

例 1. 设计识别 $L_{01} = \{0^n 1^n \mid n \geq 1\}$ 的 PDA.

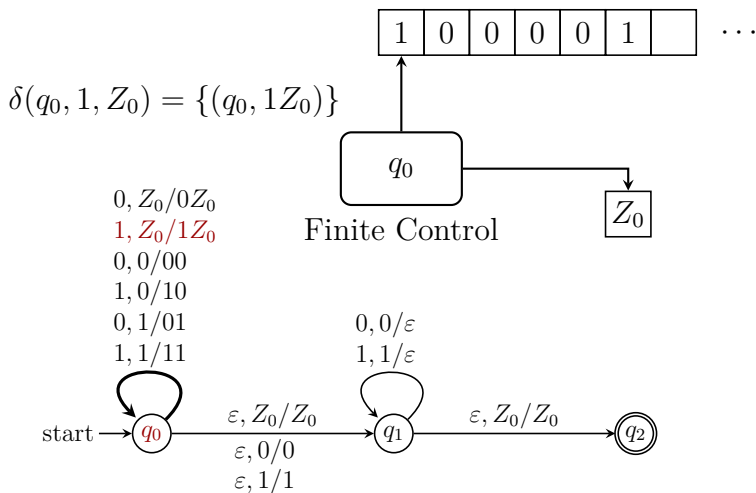


例 2. 设计识别 $L_{ww^R} = \{ww^R \mid w \in (\mathbf{0} + \mathbf{1})^*\}$ 的 PDA.

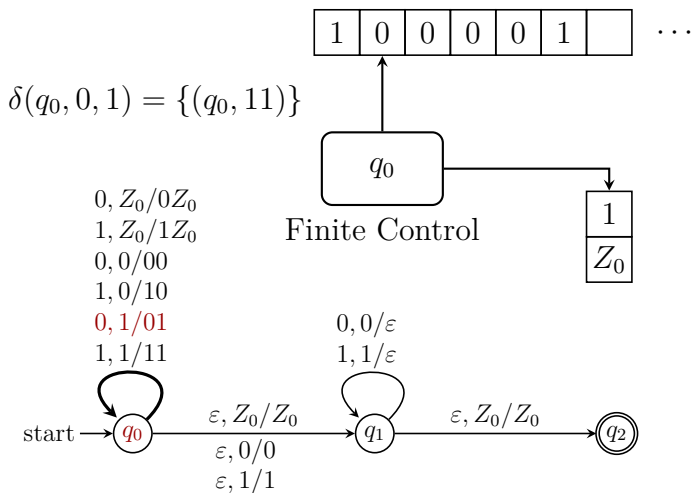
例 2. 设计识别 $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA.



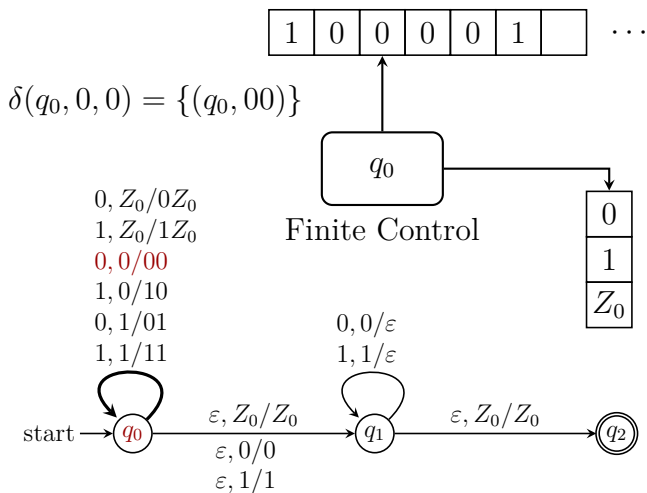
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



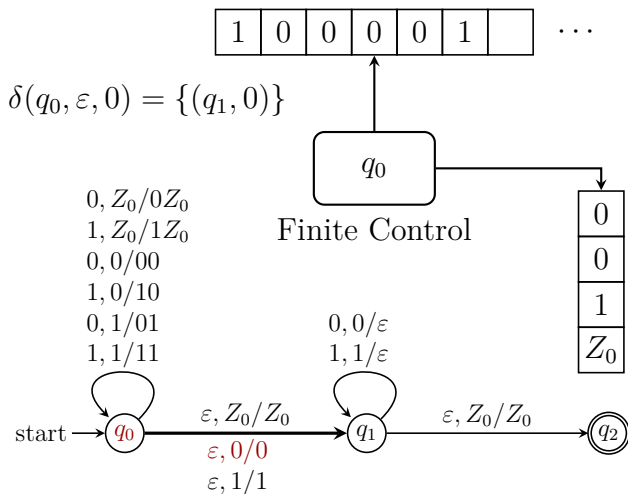
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



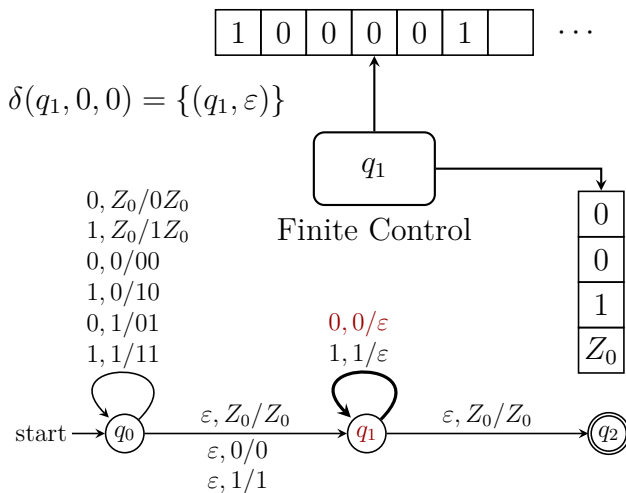
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



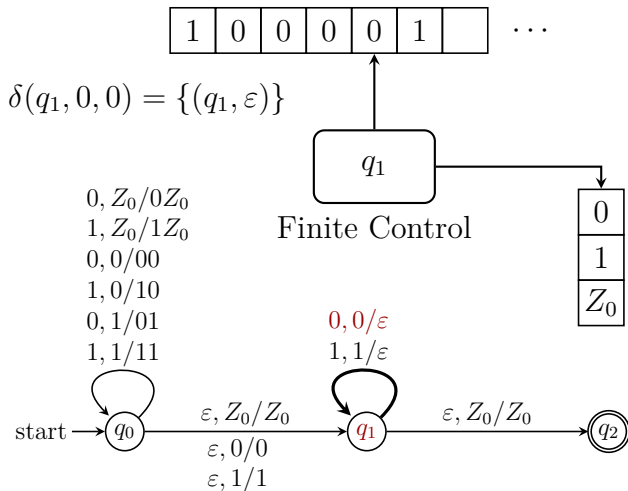
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



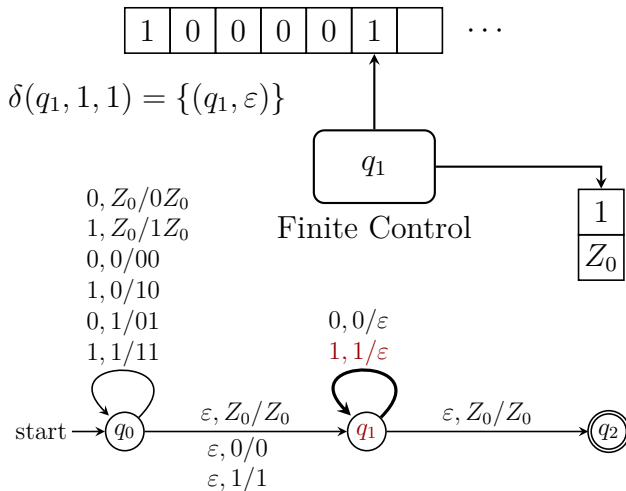
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



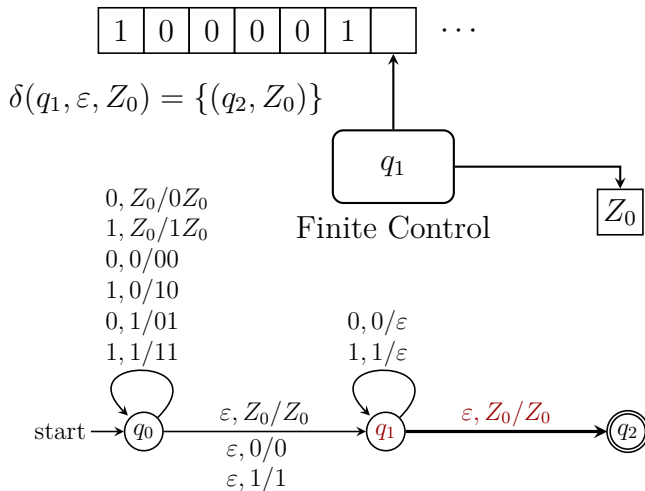
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



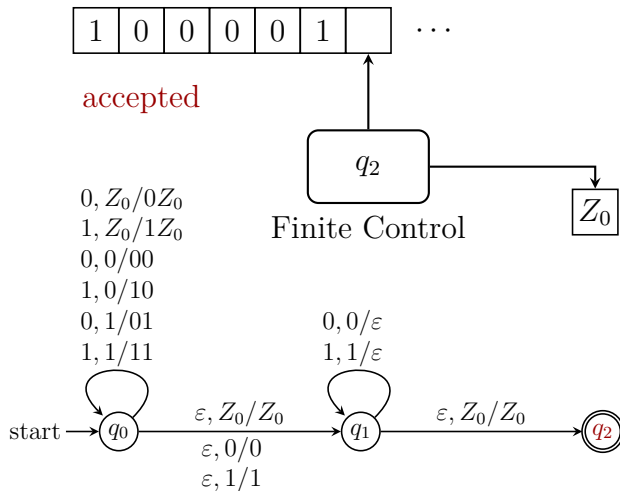
续例 2. $L_{ww^R} = \{ww^R \mid w \in (\mathbf{0} + \mathbf{1})^*\}$ 的 PDA 识别串 100001 的过程.



续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



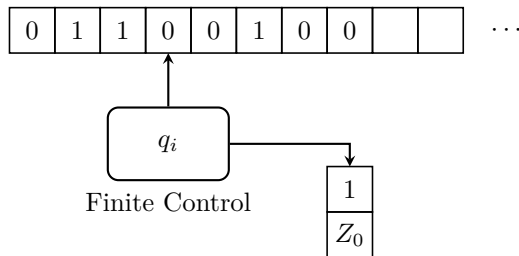
瞬时描述

定义

为描述 PDA 瞬间的格局, 定义 $Q \times \Sigma^* \times \Gamma^*$ 中三元组

$$(q, w, \gamma)$$

为**瞬时描述**(*ID*, *Instantaneous Description*), 表示此时 PDA 处于状态 q , 剩余输入串 w , 栈为 γ .



转移符号

定义

在 PDA P 中如果 $(p, \beta) \in \delta(q, a, Z)$, 由 $(q, aw, Z\alpha)$ 到 $(p, w, \beta\alpha)$ 的变化, 称为 ID 的转移 \vdash_P , 记为

$$(q, aw, Z\alpha) \vdash_P (p, w, \beta\alpha)$$

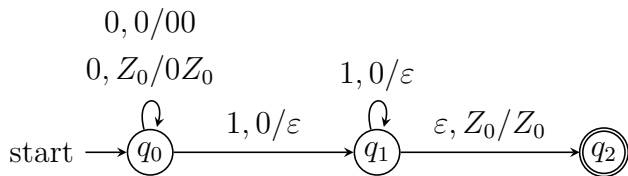
其中 $w \in \Sigma^*$, $\alpha \in \Gamma^*$.

若有 ID I, J 和 K , 递归定义 \vdash_P^* 为:

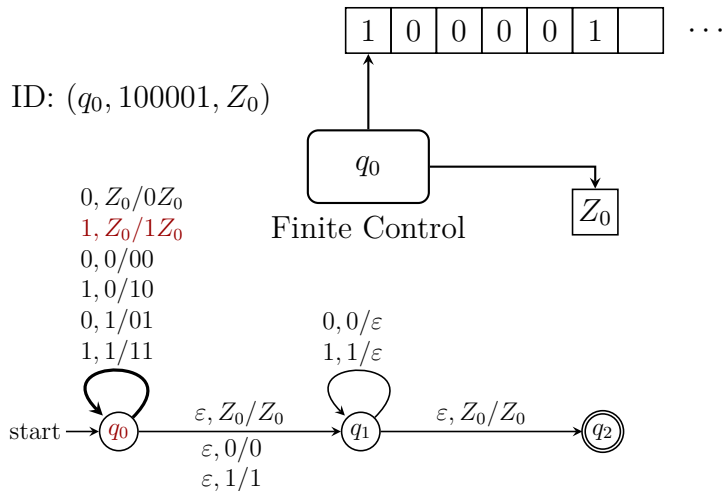
- ① $I \vdash_P^* I$;
- ② 若 $I \vdash_P J$, $J \vdash_P^* K$, 则 $I \vdash_P^* K$.

若 P 已知, 可省略, 记为 \vdash 和 \vdash^* .

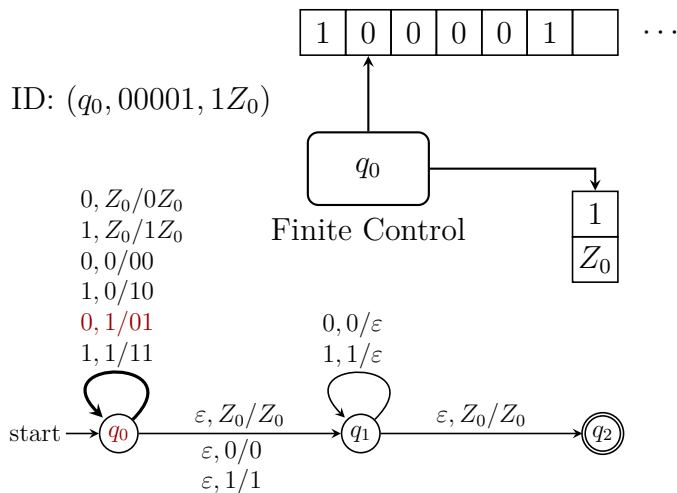
续例 1. 语言 $L_{01} = \{0^n 1^n \mid n \geq 1\}$ 的 PDA, 识别 0011 时的 ID 序列.



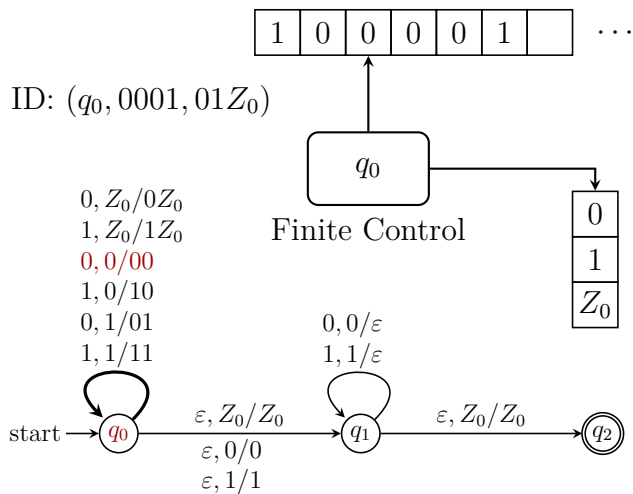
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



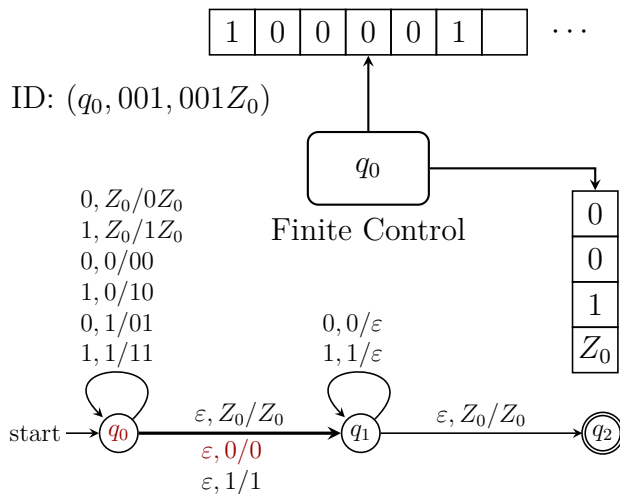
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



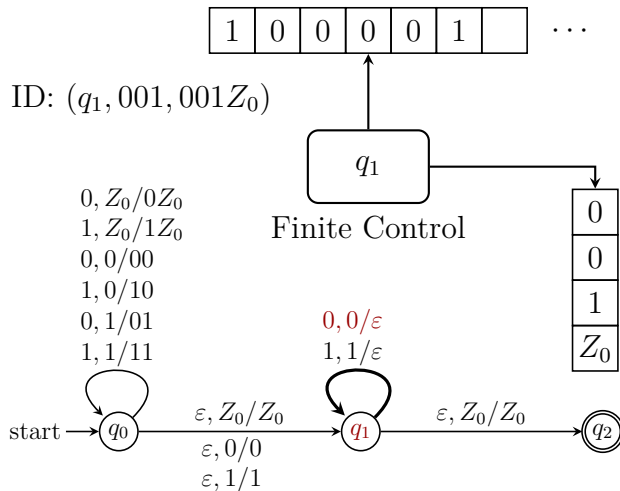
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



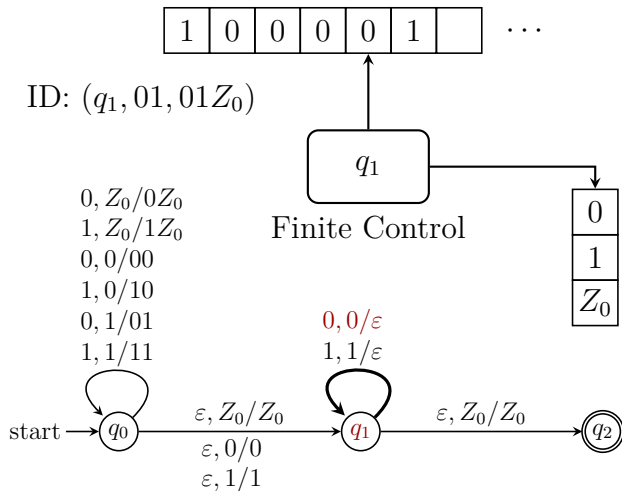
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



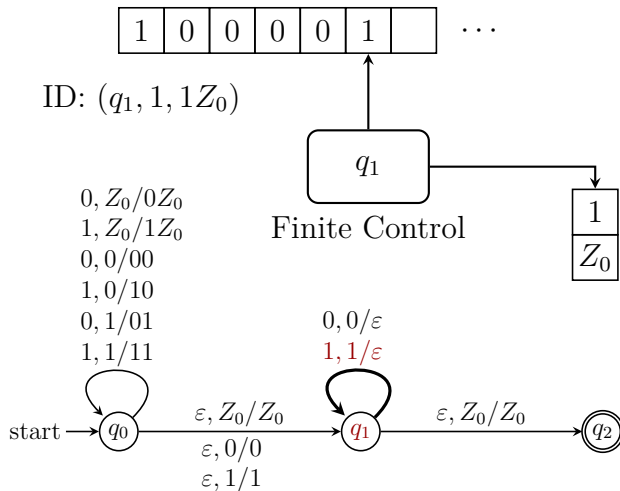
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



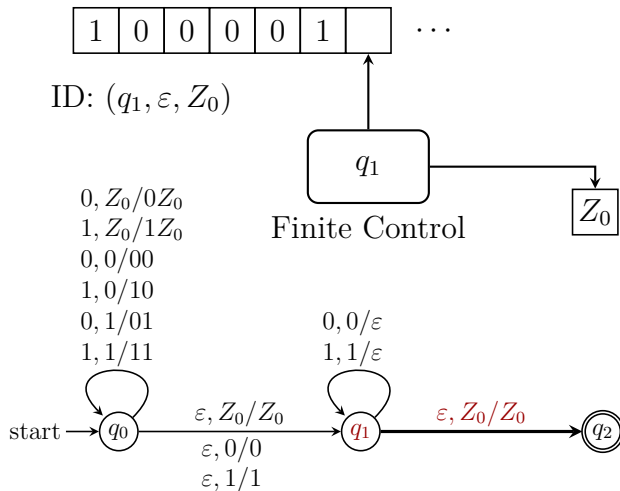
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



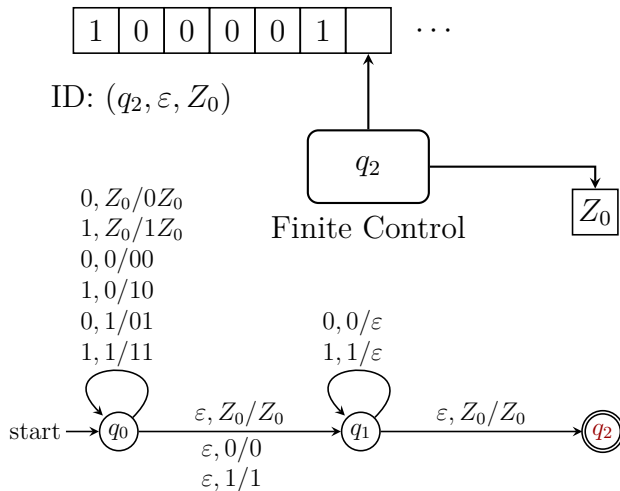
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



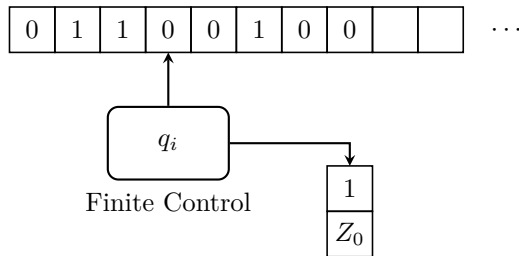
续例 2. $L_{ww^R} = \{ww^R \mid w \in (0+1)^*\}$ 的 PDA 识别串 100001 的过程.



有关 ID 的序列

对 PDA P 的一个合法 ID 序列 (计算):

- ① 把相同的字符串加到所有 ID 的输入串末尾, 所得到的计算合法;
- ② 把相同的栈符号串加到所有 ID 的栈底之下, 所得到的计算合法;
- ③ 把所有 ID 中都未消耗的部分输入串去掉, 所得到的计算合法.



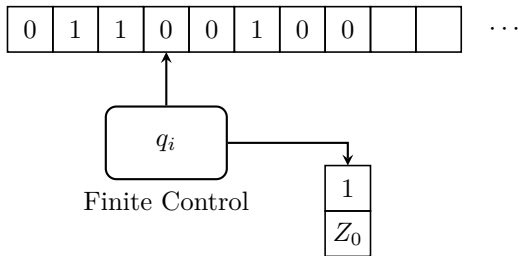
定理 23

对 $\forall w \in \Sigma^*, \forall \gamma \in \Gamma^*$, 如果

$$(q, x, \alpha) \vdash_P^* (p, y, \beta),$$

那么

$$(q, xw, \alpha\gamma) \vdash_P^* (p, yw, \beta\gamma).$$



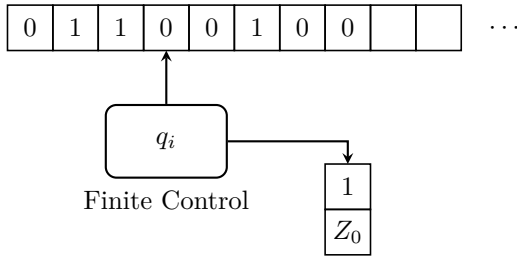
定理 24

对 $\forall w \in \Sigma^*$, 如果

$$(q, xw, \alpha) \vdash_P^* (p, yw, \beta),$$

那么

$$(q, x, \alpha) \vdash_P^* (p, y, \beta).$$



下推自动机

- 下推自动机
- 下推自动机接受的语言
 - 从终态方式到空栈方式
 - 从空栈方式到终态方式
- 下推自动机与文法的等价性
- 确定型下推自动机

下推自动机接受的语言

定义

$PDA\ P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$, 以两种方式接受语言:

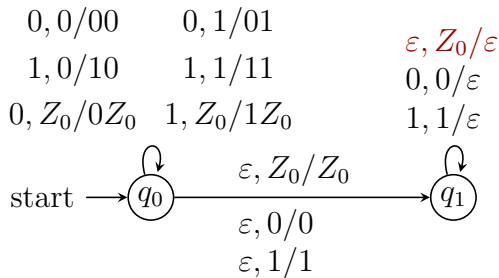
- P 以终态方式接受的语言, 记为 $\mathbf{L}(P)$, 定义为

$$\mathbf{L}(P) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \varepsilon, \gamma),\ p \in F\}.$$

- P 以空栈方式接受的语言, 记为 $\mathbf{N}(P)$, 定义为

$$\mathbf{N}(P) = \{w \mid (q_0, w, Z_0) \vdash^* (p, \varepsilon, \varepsilon)\}.$$

续例2. 识别 L_{wwr} 的 PDA P , 从终态方式接受, 改为空栈方式接受.
 用 $\delta(q_1, \varepsilon, Z_0) = \{(q_1, \varepsilon)\}$ 代替 $\delta(q_1, \varepsilon, Z_0) = \{(q_2, Z_0)\}$ 即可.



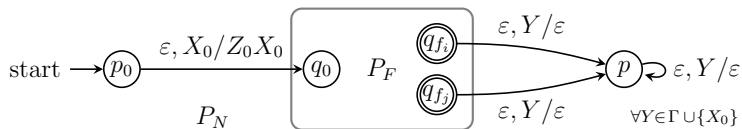
从终态方式到空栈方式

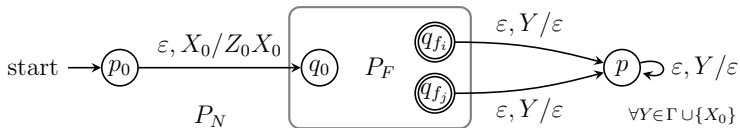
定理 25

如果 PDA P_F 以终态方式接受语言 L , 则存在 PDA P_N 以空栈方式接受 L .

证明: 设 $P_F = (Q, \Sigma, \Gamma, \delta_F, q_0, Z_0, F)$, 构造 PDA P_N ,

$$P_N = (Q \cup \{p_0, p\}, \Sigma, \Gamma \cup \{X_0\}, \delta_N, p_0, X_0, \emptyset).$$





其中 δ_N 定义如下:

- ① P_N 首先将 P_F 的栈底符号压栈, 开始模拟 P_F :

$$\delta_N(p_0, \varepsilon, X_0) = \{(q_0, Z_0 X_0)\};$$

- ② P_N 模拟 P_F 的动作: $\forall q \in Q, \forall a \in \Sigma \cup \{\varepsilon\}, \forall Y \in \Gamma$:

$\delta_N(q, a, Y)$ 包含 $\delta_F(q, a, Y)$ 的全部元素;

- ③ 从 $q_f \in F$ 开始弹出栈中符号, 即 $\forall q_f \in F, \forall Y \in \Gamma \cup \{X_0\}$:

$\delta_N(q_f, \varepsilon, Y)$ 包含 (p, ε) ;

- ④ 在状态 p 时, 弹出全部栈中符号, 即 $\forall Y \in \Gamma \cup \{X_0\}$:

$$\delta_N(p, \varepsilon, Y) = \{(p, \varepsilon)\}.$$

对 $\forall w \in \Sigma^*$ 有

$$w \in \mathbf{L}(P_F) \Rightarrow (q_0, w, Z_0) \vdash_{P_F}^* (q_f, \varepsilon, \gamma)$$

$$\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_F}^* (q_f, \varepsilon, \gamma X_0)$$

定理23

$$\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_N}^* (q_f, \varepsilon, \gamma X_0)$$

P_N 模拟 P_F

$$\Rightarrow (p_0, w, X_0) \vdash_{P_N} (q_0, w, Z_0 X_0) \vdash_{P_N}^* (q_f, \varepsilon, \gamma X_0)$$

δ_N 构造 p_0 部分

$$\Rightarrow (p_0, w, X_0) \vdash_{P_N}^* (q_f, \varepsilon, \gamma X_0) \vdash_{P_N}^* (p, \varepsilon, \varepsilon)$$

δ_N 构造 q_f 和 p 部分

$$\Rightarrow w \in \mathbf{N}(P_N)$$

即 $\mathbf{L}(P_F) \subseteq \mathbf{N}(P_N)$.

对 $\forall w \in \Sigma^*$ 有

$$\begin{aligned}
 w \in \mathbf{N}(P_N) &\Rightarrow (p_0, w, X_0) \vdash_{P_N}^* (p, \varepsilon, \varepsilon) && \text{其他状态不可能空栈} \\
 &\Rightarrow (p_0, w, X_0) \vdash_{P_N} (q_0, w, Z_0 X_0) \vdash_{P_N}^* (p, \varepsilon, \varepsilon) && \text{第一个动作必然到 } q_0 \\
 &\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_N}^* (q_f, \varepsilon, \gamma X_0) \vdash_{P_N}^* (p, \varepsilon, \varepsilon) && \text{必经 } q_f \in F \text{ 消耗完 } w \\
 &\Rightarrow (q_0, w, Z_0) \vdash_{P_N}^* (q_f, \varepsilon, \gamma) && P_N \text{ 中未用过栈底的 } X_0 \\
 &\Rightarrow (q_0, w, Z_0) \vdash_{P_F}^* (q_f, \varepsilon, \gamma) && \text{均为模拟 } P_F \\
 &\Rightarrow w \in \mathbf{L}(P_F)
 \end{aligned}$$

即 $\mathbf{N}(P_N) \subseteq \mathbf{L}(P_F)$. 所以 $\mathbf{N}(P_N) = \mathbf{L}(P_F)$.

□

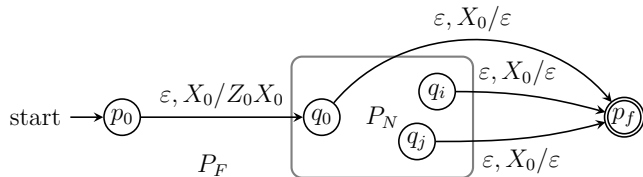
从空栈方式到终态方式

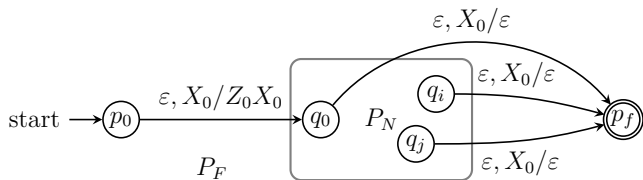
定理 26

如果 PDA P_N 以空栈方式接受语言 L , 则存在 PDA P_F 以终态方式接受 L .

证明: 设 $P_N = (Q, \Sigma, \Gamma, \delta_N, q_0, Z_0, \emptyset)$. 构造 PDA P_F ,

$$P_F = (Q \cup \{p_0, p_f\}, \Sigma, \Gamma \cup \{X_0\}, \delta_F, p_0, X_0, \{p_f\})$$





其中 δ_F 定义如下:

- ① P_F 开始时, 将 P_N 栈底符号压入栈, 并开始模拟 P_N ,

$$\delta_F(p_0, \epsilon, X_0) = \{(q_0, Z_0 X_0)\};$$

- ② P_F 模拟 P_N , $\forall q \in Q$, $\forall a \in \Sigma \cup \{\epsilon\}$, $\forall Y \in \Gamma$:

$$\delta_F(q, a, Y) = \delta_N(q, a, Y);$$

- ③ 在 $\forall q \in Q$ 时, 看到 P_F 的栈底 X_0 , 则转移到新终态 p_f :

$$\delta_F(q, \epsilon, X_0) = \{(p_f, \epsilon)\}.$$

对 $\forall w \in \Sigma^*$ 有

$$w \in \mathbf{N}(P_N) \Rightarrow (q_0, w, Z_0) \vdash_{P_N}^* (q, \varepsilon, \varepsilon)$$

$$\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_N}^* (q, \varepsilon, X_0)$$

$$\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_F}^* (q, \varepsilon, X_0)$$

$$\Rightarrow (p_0, w, X_0) \vdash_{P_F} (q_0, w, Z_0 X_0) \vdash_{P_F}^* (q, \varepsilon, X_0)$$

$$\Rightarrow (p_0, w, X_0) \vdash_{P_F}^* (q, \varepsilon, X_0) \vdash_{P_F} (p_f, \varepsilon, \varepsilon)$$

$$\Rightarrow (p_0, w, X_0) \vdash_{P_F}^* (p_f, \varepsilon, \varepsilon)$$

$$\Rightarrow w \in \mathbf{L}(P_F)$$

定理23

P_F 模拟 P_N

δ_F 构造, p_0 部分

δ_F 构造, p_f 部分

即 $\mathbf{N}(P_N) \subseteq \mathbf{L}(P_F)$.

对 $\forall w \in \Sigma^*$ 有

$$\begin{aligned}w \in \mathbf{L}(P_F) &\Rightarrow (p_0, w, X_0) \vdash_{P_F}^* (p_f, \varepsilon, \varepsilon) \\&\Rightarrow (p_0, w, X_0) \vdash_{P_F}^* (q, \varepsilon, X_0) \vdash_{P_F} (p_f, \varepsilon, \varepsilon) && \text{经 } q \text{ 才可达 } p_f \\&\Rightarrow (p_0, w, X_0) \vdash_{P_F} (q_0, w, Z_0 X_0) \vdash_{P_F}^* (q, \varepsilon, X_0) && P_F \text{ 第一个动作} \\&\Rightarrow (q_0, w, Z_0 X_0) \vdash_{P_F}^* (q, \varepsilon, X_0) && \text{即上式} \\&\Rightarrow (q_0, w, Z_0) \vdash_{P_N}^* (q, \varepsilon, \varepsilon) && P_N \text{ 与 } X_0 \text{ 无关} \\&\Rightarrow w \in \mathbf{N}(P_N)\end{aligned}$$

即 $\mathbf{N}(P_F) \subseteq \mathbf{L}(P_N)$. 所以 $\mathbf{L}(P_F) = \mathbf{N}(P_N)$.



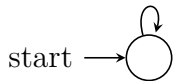
例 3. 接受 $L_{\text{eq}} = \{w \in \{0,1\}^* \mid w \text{ 中字符 } 0 \text{ 和 } 1 \text{ 的数量相同}\}$ 的 PDA.

例 3. 接受 $L_{\text{eq}} = \{w \in \{0,1\}^* \mid w \text{ 中字符 } 0 \text{ 和 } 1 \text{ 的数量相同}\}$ 的 PDA.

$0, Z_0/0Z_0$ $1, 0/10$ $0, 0/00$

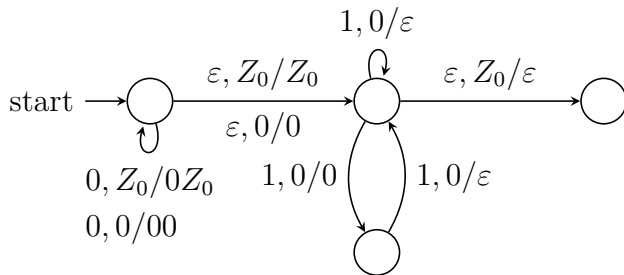
$1, Z_0/1Z_0$ $1, 1/11$ $0, 1/01$

$\varepsilon, Z_0/\varepsilon$ $1, 0/\varepsilon$ $0, 1/\varepsilon$



例 4. 接受 $L = \{0^n 1^m \mid 0 \leq n \leq m \leq 2n\}$ 的 PDA.

例 4. 接受 $L = \{0^n 1^m \mid 0 \leq n \leq m \leq 2n\}$ 的 PDA.



下推自动机

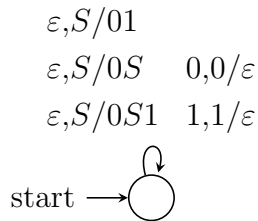
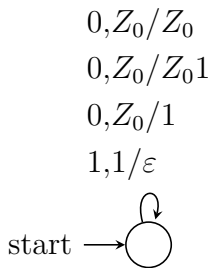
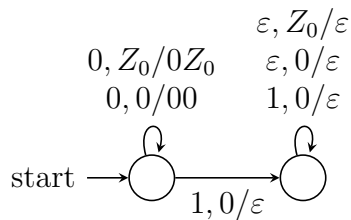
- 下推自动机
- 下推自动机接受的语言
- 下推自动机与文法的等价性
 - 由 CFG 到 PDA
 - 由 PDA 到 CFG
- 确定型下推自动机

由 CFG 到 PDA

例 5. 设计语言 $L = \{0^n 1^m \mid 1 \leq m \leq n\}$ 的 PDA.

由 CFG 到 PDA

例 5. 设计语言 $L = \{0^n 1^m \mid 1 \leq m \leq n\}$ 的 PDA.



续例 5. 设计语言 $L = \{0^n 1^m \mid 1 \leq m \leq n\}$ 的 CFG.

CFG G :

$$S \rightarrow AB$$

$$A \rightarrow 0A \mid \varepsilon$$

$$B \rightarrow 0B1 \mid 01$$

字符串 00011 的最左派生:

$$S \xRightarrow{\text{lm}} AB \xRightarrow{\text{lm}} 0AB \xRightarrow{\text{lm}} 0B \xRightarrow{\text{lm}} 00B1 \xRightarrow{\text{lm}} 00011$$

续例5. 语言 $L = \{0^n 1^m \mid 1 \leq m \leq n\}$.

用 PDA 栈顶符号的替换, 模拟文法的最左派生:

PDA		CFG	
PDA 的 ID 转移	PDA 的动作	产生式	最左派生
$(q_0, 00011, S)$			S
$\vdash (q_0, 00011, AB)$	$\varepsilon, S/AB$	$S \rightarrow AB$	$\xRightarrow{\text{lm}} AB$
$\vdash (q_0, 00011, 0AB)$	$\varepsilon, A/0A$	$A \rightarrow 0A$	$\xRightarrow{\text{lm}} 0AB$
$\vdash (q_0, 0011, AB)$	$0, 0/\varepsilon$		
$\vdash (q_0, 0011, B)$	$\varepsilon, A/\varepsilon$	$A \rightarrow \varepsilon$	$\xRightarrow{\text{lm}} 0B$
$\vdash (q_0, 0011, 0B1)$	$\varepsilon, B/0B1$	$B \rightarrow 0B1$	$\xRightarrow{\text{lm}} 00B1$
$\vdash (q_0, 011, B1)$	$0, 0/\varepsilon$		
$\vdash (q_0, 011, 011)$	$\varepsilon, B/01$	$B \rightarrow 01$	$\xRightarrow{\text{lm}} 00011$
$\vdash (q_0, 11, 11)$	$0, 0/\varepsilon$		
$\vdash (q_0, 1, 1)$	$1, 1/\varepsilon$		
$\vdash (q_0, \varepsilon, \varepsilon)$	$1, 1/\varepsilon$		

续例 5. 语言 $L = \{0^n 1^m \mid 1 \leq m \leq n\}$.

$$S \rightarrow AB$$

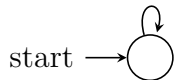
$$A \rightarrow 0A \mid \varepsilon$$

$$B \rightarrow 0B1 \mid 01$$

$$\varepsilon, S/AB$$

$$\varepsilon, A/0A \quad \varepsilon, A/\varepsilon \quad 0, 0/\varepsilon$$

$$\varepsilon, B/0B1 \quad \varepsilon, B/01 \quad 1, 1/\varepsilon$$



定理 27

任何 CFL L , 一定存在 PDA P , 使 $L = \mathbf{N}(P)$.

构造与文法等价的 PDA

如果 CFG $G = (V, T, P', S)$, 构造 PDA

$$P = (\{q\}, T, V \cup T, \delta, q, S, \emptyset),$$

其中 δ 为:

① $\forall A \in V:$

$$\delta(q, \varepsilon, A) = \{(q, \beta) \mid A \rightarrow \beta \in P'\}$$

② $\forall a \in T:$

$$\delta(q, a, a) = \{(q, \varepsilon)\}$$

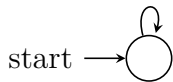
那么 $\mathbf{L}(G) = \mathbf{N}(P)$.

例 6. 为文法 $S \rightarrow aAA, A \rightarrow aS \mid bS \mid a$ 构造 PDA.

例 6. 为文法 $S \rightarrow aAA$, $A \rightarrow aS \mid bS \mid a$ 构造 PDA.

$\varepsilon, S/aAA$ $\varepsilon, A/aS$ $a, a/\varepsilon$

$\varepsilon, A/a$ $\varepsilon, A/bS$ $b, b/\varepsilon$



证明: 往证

$$S \xRightarrow{*} w \iff (q, w, S) \vdash_P^* (q, \varepsilon, \varepsilon).$$

[充分性] 往证

$$S \xRightarrow{\text{lm}}^* w \implies (q, w, S) \vdash^* (q, \varepsilon, \varepsilon).$$

设 $S \xRightarrow{\text{lm}}^* w$ 中第 i 个左句型为 $x_i A_i \alpha_i$, 其中 $x_i \in \Sigma^*$, $A_i \in V$, $\alpha_i \in (V \cup T)^*$.
并将 S 看作第 0 个左句型 $x_0 A_0 \alpha_0 = S$, 那么

$$x_0 = \varepsilon, A_0 = S, \alpha_0 = \varepsilon.$$

将 w 看作为第 n 个左句型 $x_n A_n \alpha_n = w$, 那么

$$x_n = w, A_n = \varepsilon, \alpha_n = \varepsilon.$$

再对派生步骤 i 归纳, 往证

$$S \xRightarrow{\text{lm}}^i x_i A_i \alpha_i \wedge w = x_i y_i \implies (q, w, S) \vdash^* (q, y_i, A_i \alpha_i).$$

归纳基础: 最左派生在第 0 步时, 显然成立

$$(q, w, S) \vdash^* (q, y_0, A_0\alpha_0) = (q, w, S).$$

归纳递推: 假设第 i 步时成立, 当第 $i+1$ 步时, 一定是 $A_i \rightarrow \beta$ 应用到 $x_i A_i \alpha_i$

$$S \xRightarrow{\text{lm}}^i x_i A_i \alpha_i \xRightarrow{\text{lm}} x_i \beta \alpha_i = x_{i+1} A_{i+1} \alpha_{i+1}.$$

即最左变元 A_{i+1} 一定在 $\beta \alpha_i$ 中, 设 A_{i+1} 之前的终结符为 x' , 那么由

$$\begin{aligned} x_i \beta \alpha_i &= x_i x' A_{i+1} \alpha_{i+1} = x_{i+1} A_{i+1} \alpha_{i+1} \\ x_i y_i &= x_i x' y_{i+1} = x_{i+1} y_{i+1} = w \end{aligned}$$

则有

$$\begin{aligned} \beta \alpha_i &= x' A_{i+1} \alpha_{i+1}, \\ y_i &= x' y_{i+1}. \end{aligned}$$

那么, 在 PDA 中从 ID $(q, y_i, A_i\alpha_i)$ 模拟最左派生, 用产生式 $A_i \rightarrow \beta$ 替换栈顶 A_i 后, 有

$$\begin{aligned}
 (q, w, S) &\vdash^* (q, y_i, A_i\alpha_i) && \text{归纳假设} \\
 &\vdash (q, y_i, \beta\alpha_i) && A_i \rightarrow \beta \\
 &= (q, x'y_{i+1}, x'A_{i+1}\alpha_{i+1}) \\
 &\vdash^* (q, y_{i+1}, A_{i+1}\alpha_{i+1}) && \text{弹出栈顶终结符}
 \end{aligned}$$

因此 $S \xrightarrow[n]{\text{lm}} w \implies (q, w, S) \vdash^* (q, y_n, A_n\alpha_n) = (q, \varepsilon, \varepsilon)$, 即充分性得证.

[必要性] 往证更一般的, 对任何变元 A , 都有:

$$(q, x, A) \vdash^* (q, \varepsilon, \varepsilon) \implies A \xRightarrow{*} x.$$

对 ID 转移 $(q, x, A) \vdash^i (q, \varepsilon, \varepsilon)$ 的次数 i 归纳证明.

归纳基础: 当 $i = 1$ 步时, 只能是 $x = \varepsilon$ 且 $A \rightarrow \varepsilon$ 为产生式, 所以 $A \xRightarrow{*} \varepsilon$.

归纳递推: 假设 $i \leq n$ ($n \geq 1$) 步时上式成立. 当 $i = n + 1$ 时, 因为 A 是变元, 其第 1 步转移一定是应用某产生式 $A \rightarrow Y_1 Y_2 \cdots Y_m$

$$(q, x, A) \vdash (q, x, Y_1 Y_2 \cdots Y_m)$$

其中 Y_i 是变元或终结符. 而其余的 n 步转移

$$(q, x, Y_1 Y_2 \cdots Y_m) \vdash^* (q, \varepsilon, \varepsilon)$$

中每个 Y_i 从栈中被完全弹出时, 将消耗掉的那部分 x 记为 x_i , 那么显然有

$$x = x_1 x_2 \cdots x_m.$$

而每个 Y_i 从栈中被完全弹出时, 都不超过 n 步, 所以由归纳假设,

$$(q, x_i, Y_i) \vdash^* (q, \varepsilon, \varepsilon) \implies Y_i \xRightarrow{*} x_i.$$

再由产生式 $A \rightarrow Y_1 Y_2 \cdots Y_m$, 有

$$\begin{aligned} A &\Rightarrow Y_1 Y_2 \cdots Y_m \\ &\stackrel{*}{\Rightarrow} x_1 Y_2 \cdots Y_m \\ &\stackrel{*}{\Rightarrow} x_1 x_2 \cdots Y_m \\ &\stackrel{*}{\Rightarrow} x_1 x_2 \cdots x_m = x. \end{aligned}$$

因此当 $A = S$, $x = w$ 时,

$$(q, w, S) \vdash^* (q, \varepsilon, \varepsilon) \Longrightarrow S \stackrel{*}{\Rightarrow} w$$

成立, 即必要性得证.

所以, 任何 CFL 都可由 PDA 识别.



构造与 GNF 格式文法等价的 PDA

如果 GNF 格式的 CFG $G = (V, T, P', S)$, 那么构造 PDA

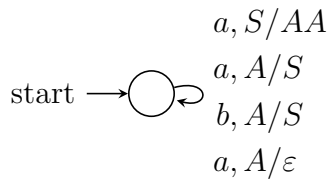
$$P = (\{q\}, T, V, \delta, q, S, \emptyset),$$

为每个产生式, 定义 δ 为:

$$\delta(q, a, A) = \{(q, \beta) \mid A \rightarrow a\beta \in P'\}.$$

续例 6. 文法 $S \rightarrow aAA, A \rightarrow aS \mid bS \mid a$ 为 GNF 格式, 构造等价的 PDA.

续例 6. 文法 $S \rightarrow aAA$, $A \rightarrow aS \mid bS \mid a$ 为 GNF 格式, 构造等价的 PDA.



由 PDA 到 CFG

定理 28

如果 PDA P , 有 $L = \mathbf{N}(P)$, 那么 L 是上下文无关语言.

构造与 PDA 等价的 CFG

如果 PDA $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, \emptyset)$, 那么构造 CFG $G = (V, \Sigma, P', S)$, 其中 V 和 P' 为

- ① $V = \{[qXp] \mid p, q \in Q, X \in \Gamma\} \cup \{S\}$;
- ② 对 $\forall p \in Q$, 构造产生式 $S \rightarrow [q_0Z_0p]$;
- ③ 对 $\forall (p, Y_1Y_2 \cdots Y_n) \in \delta(q, a, X)$, 构造 $|Q|^n$ 个产生式

$$[qXr_n] \rightarrow a[pY_1r_1][r_1Y_2r_2] \cdots [r_{n-1}Y_nr_n]$$

其中 $a \in \Sigma \cup \{\varepsilon\}$, $X, Y_i \in \Gamma$, 而 $r_i \in Q$ 是 n 次 $|Q|$ 种状态的组合;
若 $i = 0$, 为 $[qXp] \rightarrow a$.

证明: 只需证明

$$(q, w, X) \vdash^* (p, \varepsilon, \varepsilon) \iff [qXp] \xRightarrow{*} w.$$

并令 $X = Z_0$, $q = q_0$, 与开始符号 S 的产生式一起, 即可完成定理的证明.

[充分性] 对 PDA 中 $(q, w, X) \vdash^* (p, \varepsilon, \varepsilon)$ 的转移次数 i 归纳证明.

归纳基础: 当 $i = 1$ 时, P 只能消耗不超过一个的字符, 即 $w = a$

$$(q, w, X) = (q, a, X) \vdash (p, \varepsilon, \varepsilon),$$

其中 $a \in \Sigma \cup \{\varepsilon\}$ 且 $(p, \varepsilon) \in \delta(q, a, X)$, 则由文法的构造会有

$$[qXp] \rightarrow a,$$

因此 $[qXp] \xRightarrow{*} a = w$.

归纳递推: 假设当 $i \leq m$ ($m \geq 1$) 时命题成立. 当 $i = m + 1$ 时, 转移的第 1 步, 一定由某个 $(r_0, Y_1 Y_2 \cdots Y_n) \in \delta(q, a, X)$ 开始

$$(q, ax, X) \vdash (r_0, x, Y_1 Y_2 \cdots Y_n),$$

其中 $a \in \Sigma \cup \{\varepsilon\}$, $w = ax$. 而其余的 m 步为

$$(r_0, x, Y_1 Y_2 \cdots Y_n) \vdash^* (p, \varepsilon, \varepsilon).$$

而这些转移, 会从栈中依次弹出 Y_i 并消耗掉部分 x . 若分别记为 x_i , 则有

$$w = ax = ax_1 x_2 \cdots x_n.$$

若设弹出 Y_i 之前和之后的状态分别是 r_{i-1} 和 r_i , 这里 $i = 1, 2, \cdots, n$, 那么有

$$(r_{i-1}, x_i, Y_i) \vdash^* (r_i, \varepsilon, \varepsilon),$$

且转移步数都不会超过 m . 那么, 由归纳假设

$$(r_{i-1}, x_i, Y_i) \vdash^* (r_i, \varepsilon, \varepsilon) \implies [r_{i-1} Y_i r_i] \xRightarrow{*} x_i.$$

而由动作 $(r_0, Y_1 Y_2 \cdots Y_n) \in \delta(q, a, X)$ 所构造的产生式会包含

$$[qXr_n] \rightarrow a[r_0Y_1r_1][r_1Y_2r_2] \cdots [r_{n-1}Y_nr_n].$$

而显然弹出 X 后的状态 p 与弹出 Y_n 后的状态 r_n 是同一个, 所以

$$[qXp] = [qXr_n] \Rightarrow a[r_0Y_1r_1][r_1Y_2r_2] \cdots [r_{n-1}Y_nr_n] \xRightarrow{*} ax_1x_2 \cdots x_n = w$$

因此充分性得证. 那么当 $X = Z_0, q = q_0$ 时有

$$(q_0, w, Z_0) \vdash^* (p, \varepsilon, \varepsilon) \Longrightarrow [q_0Z_0p] \xRightarrow{*} w,$$

以及产生式 $S \rightarrow [q_0Z_0p]$ 有 $S \xRightarrow{*} w$, 即 PDA 接受的串可由文法派生得到.

[必要性]: 略.



例 7. 将 PDA $P = (\{p, q\}, (0, 1), \{X, Z\}, \delta, q, Z)$ 转为 CFG, 其中 δ 如下:

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

$$(2) \quad \delta(q, 1, X) = \{(q, XX)\}$$

$$(3) \quad \delta(q, 0, X) = \{(p, X)\}$$

$$(4) \quad \delta(q, \varepsilon, Z) = \{(q, \varepsilon)\}$$

$$(5) \quad \delta(p, 1, X) = \{(p, \varepsilon)\}$$

$$(6) \quad \delta(p, 0, Z) = \{(q, Z)\}$$

δ	产生式
(0)	$S \rightarrow [qZq]$ $S \rightarrow [qZp]$
(1)	$[qZq] \rightarrow 1[qXq][qZq]$ $[qZq] \rightarrow 1[qXp][pZq]$ $[qZp] \rightarrow 1[qXq][qZp]$ $[qZp] \rightarrow 1[qXp][pZp]$
(2)	$[qXq] \rightarrow 1[qXq][qXq]$ $[qXq] \rightarrow 1[qXp][pXq]$ $[qXp] \rightarrow 1[qXq][qXp]$ $[qXp] \rightarrow 1[qXp][pXp]$
(3)	$[qXq] \rightarrow 0[pXq]$ $[qXp] \rightarrow 0[pXp]$
(4)	$[qZq] \rightarrow \varepsilon$
(5)	$[pXp] \rightarrow 1$
(6)	$[pZp] \rightarrow 0[qZp]$ $[pZq] \rightarrow 0[qZq]$

消除无用符号	重命名 (可选)
$S \rightarrow [qZq]$	$S \rightarrow A$
$[qZq] \rightarrow 1[qXp][pZq]$	$A \rightarrow 1BC$
$[qXp] \rightarrow 1[qXp][pXp]$	$B \rightarrow 1BD$
$[qXp] \rightarrow 0[pXp]$	$B \rightarrow 0D$
$[qZq] \rightarrow \varepsilon$	$A \rightarrow \varepsilon$
$[pXp] \rightarrow 1$	$D \rightarrow 1$
$[pZq] \rightarrow 0[qZq]$	$C \rightarrow 0A$

下推自动机

- 下推自动机
- 下推自动机接受的语言
- 下推自动机与文法的等价性
- 确定型下推自动机
 - 正则语言与 DPDA
 - DPDA 与无歧义文法

确定型下推自动机

定义

如果 PDA $P = (Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$ 满足

- ① $\forall a \in \Sigma \cup \{\varepsilon\}, \delta(q, a, X)$ 至多有一个动作;
- ② $\exists a \in \Sigma$, 如果 $\delta(q, a, X) \neq \emptyset$, 那么 $\delta(q, \varepsilon, X) = \emptyset$.

则称 P 为确定型下推自动机(DPDA).

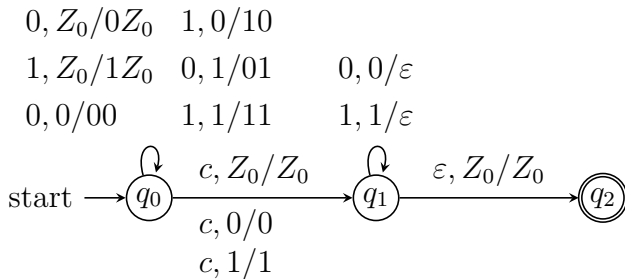
DPDA P 以终态方式接受的语言 $L(P)$ 称为 DCFL.

- DPDA 中 $\forall (q, a, Z) \in Q \times \Sigma \times \Gamma$ 满足 $|\delta(q, a, Z)| + |\delta(q, \varepsilon, Z)| \leq 1$

DPDA 与 PDA 不等价

例 8. 任何 DPDA 都无法接受 L_{ww^R} , 但是可以接受

$$L_{wcw^R} = \{wcw^R \mid w \in (\mathbf{0} + \mathbf{1})^*\}.$$



正则语言与 DPDA

定理 29

如果 L 是正则语言, 那么存在 DPDA P 以终态方式接受 L , 即 $L = \mathbf{L}(P)$.

证明: 显然, 因为 DPDA P 可以不用栈而模拟任何 DFA. □

- L_{wcur} 显然是 CFL, 所以 DCFL 语言类真包含正则语言
- DPDA 无法识别 L_{wwr} , 所以 DCFL 语言类真包含于 CFL

定义

如果语言 L 中不存在两个不同的字符串 x 和 y , 使 x 是 y 的前缀, 称语言 L 满足前缀性质.

定理 30

如果有 DPDA P 且 $L = N(P)$, 当且仅当 L 有前缀性质且存在 DPDA P' 使 $L = L(P')$.

证明: $[\Rightarrow]$ $\forall x \in N(P)$ 会弹空 P 的栈, 所以不会接受以 x 为前缀的其他串; 而转换为终态方式不改变确定性. $[\Leftarrow]$ 到达终态则弹空栈, 即可. \square

- DPDA P 的 $N(P)$ 更有限, 即使正则语言 0^* 也无法接受

DPDA 与无歧义文法

定理 31

DPDA P , 语言 $L = N(P)$, 那么 L 有无歧义的 CFG.

证明: 利用定理 28 由 P 构造的文法 G 一定无歧义, 因为:

- ① P 是确定的, 那么它接受 w 的 ID 序列也是确定的;
- ② 而由 $\delta(q, a, X) = \{(p, Y_1 \cdots Y_n)\}$ 继续弹出 Y_i 后的状态 r_i 也是确定的;
- ③ 那么由每个动作构造的一组产生式

$$[qXr_n] \rightarrow a[pY_1r_1][r_1Y_2r_2] \cdots [r_{n-1}Y_nr_n]$$

中, 仅会有一个有效的;

- ④ 那么, G 中最左派生 $S \xRightarrow{*}_{\text{lm}} w$ 就是唯一的, 所以是无歧义的.



定理 32

DPDA P , 语言 $L = \mathbf{L}(P)$, 那么 L 有无歧义的 CFG.

证明:

- ① 设符号 $\$$ 不在 L 中出现, 令 $L' = \{w\$ \mid w \in L\}$, 则 L' 具有前缀性质;
- ② 可修改 P 接受 L' , 则由定理 30, 存在 DPDA P' 使 $\mathbf{N}(P') = L'$;
- ③ 由定理 31, 存在无歧义文法 G' 使 $\mathbf{L}(G') = L'$;
- ④ 将 $\$$ 看作变元, 增加产生式 $\$ \rightarrow \varepsilon$, 修改 G' 为文法 G ;
- ⑤ 则文法 G 和 G' 一样无歧义, 且 $\mathbf{L}(G) = L$.



DCFL/DPDA 的重要应用

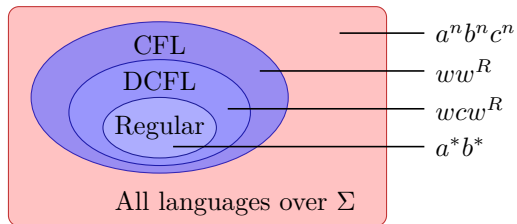
- 程序设计语言的语法分析器

如 $LR(k)$ 文法, Yacc 的基础, 解析的时间复杂度为 $O(n)$ 的算法

- 非固有歧义语言的真子集

如 L_{wwr} 有无歧义文法 $S \rightarrow 0S0 \mid 1S1 \mid \varepsilon$

语言类之间的关系



形式语言与自动机理论

上下文无关语言的性质

王春宇

计算机科学与技术学院
哈尔滨工业大学

上下文无关语言的性质

- 上下文无关语言的泵引理
 - 上下文无关语言的泵引理
 - 泵引理的应用
- 上下文无关语言的封闭性
- 上下文无关语言的判定性质
- 乔姆斯基文法体系

任何 Σ 上的所有语言是不可数的

不妨设 $\Sigma = \{a\}$, 对任何 $0 \leq x < 1$ 的实数 x , 定义语言

$$L_x = \{a^n \mid x \cdot 2^n \bmod 1 \geq 1/2\},$$

即 $a^n \in L_x$ 当且仅当 x 二进制表示的第 $n+1$ 位为 1.

- ❶ 如果 $x \neq y$, 则 x 和 y 一定有某些位不同, 所以 $L_x \neq L_y$;
- ❷ 所以 Σ 上的所有语言, 至少与 0 和 1 之间的实数一样多;
- ❸ 因此, Σ 上的所有语言是不可数的.

任何 Σ 上的所有 CFL 是可数的

任何 CFG $G = (V, \Sigma, P, S)$ 可由符号集 $V \cup \Sigma \cup \{\varepsilon, \rightarrow, |, \diamond\}$ 编码.

- 如文法 $S \rightarrow A \mid B, A \rightarrow aA \mid aC, B \rightarrow Bb \mid Cb, C \rightarrow \varepsilon \mid aCb$ 可编码为

$$S \rightarrow A \mid B \diamond A \rightarrow aA \mid aC \diamond B \rightarrow Bb \mid Cb \diamond C \rightarrow \varepsilon \mid aCb;$$

- 用 0/1 编码这些符号

$$\varepsilon \mapsto 10$$

$$a \mapsto 110$$

$$S \mapsto 1110$$

$$\rightarrow \mapsto 100$$

$$b \mapsto 1100$$

$$A \mapsto 11100$$

$$| \mapsto 1000$$

$$B \mapsto 111000$$

$$\diamond \mapsto 10000$$

$$C \mapsto 1110000.$$

- 文法编码再转换为 0/1 字符串

11101001110010001110001000011100100110111001000110
11100001000011100010011100011001000111000011001000
0111000010010100011011100001100;

- 当作二进制表示则为整数

2486025347845581444133243339142670726924.

- 而 Σ 上两个文法如果不同, 这样编码会得到不同的整数;
- 因此 Σ 上所有 CFL 至多与正整数一样多, 是可数的.
- 因此, 并非所有的语言都是 CFL.

语法分析树的大小

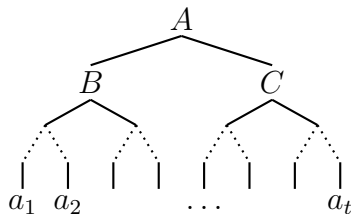
定理 33

对于乔姆斯基范式文法 $G = (V, T, P, S)$ 的语法树, 如果产物为终结字符串 w , 且树中最长路径的长度是 n , 那么 $|w| \leq 2^{n-1}$.

证明: 对最长路径的长度归纳.

基础: 为 1 时, 只能是 $\begin{smallmatrix} A \\ | \\ a \end{smallmatrix}$, 显然成立.

递推: 为 n 时根节点一定是 $A \rightarrow BC$, 而 B 和 C 子树最长路径最多为 $n-1$, 由归纳假设, 产物最长都为 2^{n-2} . 因此整棵树产物最长 $2^{n-2} + 2^{n-2} = 2^{n-1}$. \square



上下文无关语言的泵引理

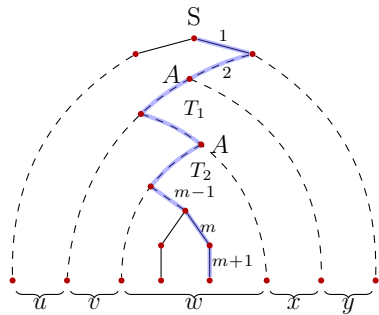
定理 34

如果语言 L 是 CFL, 那么存在正整数 N , 它只依赖于 L , 对 $\forall z \in L$, 只要 $|z| \geq N$, 就可以将 z 分为五部分 $z = uvwxy$ 满足:

- ① $vx \neq \varepsilon$ (或 $|vx| > 0$);
- ② $|vwx| \leq N$;
- ③ $\forall i \geq 0, uv^iwx^iy \in L$.

证明:

- ① 设 CNF 格式 CFG G 中变元数 $|V| = m$, 令 $N = 2^m$, 若有 $z \in L(G)$, 且 $|z| \geq N$.
- ② 则 z 的派生树内节点是二叉树, 最长路径长度至少 $m + 1$, 节点至少 $m + 2$ 个.
- ③ 该路径由下至上 $m + 1$ 个内节点中, 必有两个 T_2 和 T_1 标记了相同的变元 A .
- ④ 若记 T_2 产物为 w , 且是 T_1 的子树, T_1 的产物可记为 $vw x$, 则有 $A \Rightarrow vAx$ 和 $A \Rightarrow w$.
- ⑤ 那么 $\forall i \geq 0, A \Rightarrow v^i w x^i$. 不妨设 $z = uvwxy$, 则 $S \Rightarrow uAy \Rightarrow uv^i w x^i y$.
- ⑥ T_1 路径长不超过 $m + 1$, 那么 T_1 产物长不超过 2^m , 所以 $|vw x| \leq 2^m$.
- ⑦ T_2 必在 T_1 的左/右儿子中, 所以 v 和 x 不可能同时为空, 即 $vx \neq \varepsilon$.



□

泵引理的应用

例 1. 证明 $L = \{0^n 1^n 2^n \mid n \geq 1\}$ 不是上下文无关语言.

泵引理的应用

例 1. 证明 $L = \{0^n 1^n 2^n \mid n \geq 1\}$ 不是上下文无关语言.

证明:

- ① 假设 L 是 CFL, 那么存在整数 N , 对 $\forall z \in L (|z| \geq N)$ 满足泵引理.
- ② 从 L 中取 $z = 0^N 1^N 2^N$, 则显然 $z \in L$ 且 $|z| = 3N \geq N$.
- ③ 由泵引理, z 可被分为 $z = uvwxy$, 且有 $|vwx| \leq N$ 和 $vx \neq \varepsilon$.
- ④ 那么 vwx 只能包含一种或两种字符:
 - i 一种字符, 或为 0, 或为 1, 或为 2, 那么 $uwy \notin L$;
 - ii 两种字符, 或为 0 和 1, 或为 1 和 2, 那么也有 $uwy \notin L$;
- ⑤ 与泵引理 $uwy = uv^0wx^0y \in L$ 矛盾, 假设不成立.
- ⑥ L 不是上下文无关的.



例 2. 证明 $L = \{ww \mid w \in \{0,1\}^*\}$ 不是上下文无关的.

(错误的) 证明: 假设 L 是 CFL. 取 $z = 0^N 1 0^N 1$, 那么 $z = uvwxy$ 为

$$z = \overbrace{\underbrace{00 \cdots 00}_u \underbrace{0}_v \underbrace{1}_w \underbrace{0}_x \underbrace{00 \cdots 01}_y}^{0^N 1} \overbrace{\quad}^{0^N 1}$$

则对任意 $i \geq 0$, 有 $uv^iwx^iy \in L$, 满足泵引理. □

(正确的) 证明: 假设 L 是 CFL. 取 $z = 0^N 1^N 0^N 1^N$, 将 z 分为 $z = uvwxy$ 时

- ① 若 vw 在 z 中点的一侧, uv^0wx^0y 显然不可能属于 L ;
- ② 若 vw 包括 z 中点, 那么 uv^0wx^0y 为 $0^N 1^i 0^j 1^N$, 也不可能属于 L .

所以假设不成立, L 不是 CFL. □

CFL 的泵引理同样只是必要条件

有些非 CFL, 泵引理对它们没有什么作用. 例如

$$L = \{a^i b^j c^k d^l \mid i = 0 \text{ 或 } j = k = l\}$$

不是上下文无关的.

- 如果选 $z = b^j c^k d^l$, 则可以让 $z = uvwxy$ 的 vwx 只含有 b , 那么对任何 m , 都有 $uv^mwx^my \in L$;
- 如果选 $z = a^i b^j c^j d^j$, 则可以让 v 和 x 只包含 a , 那么对任何 m , 都有 $uv^mwx^my \in L$.

所以无法使用泵引理证明 L 非 CFL.

Ogden 引理 (的较弱形式)

如果语言 L 是 CFL, 那么存在正整数 N , 它只依赖于 L , 对 $\forall z \in L$, 在 z 中至少 N 个任意位置作**标记**后, 就可以将 z 分为五部分 $z = uvwxy$ 满足:

- ① v 和 x 一起至少含有一个标记位置;
- ② vwx 中至多有 N 个标记位置;
- ③ $\forall i \geq 0, uv^iwx^iy \in L$.

上下文无关语言的性质

- 上下文无关语言的泵引理
- 上下文无关语言的封闭性
 - 代换的封闭性
 - 并/连接/闭包/同态/逆同态/反转的封闭性
 - 交和补运算不封闭
 - 封闭性的应用
- 上下文无关语言的判定性质
- 乔姆斯基文法体系

代换

定义

两个字母表 Σ 到 Γ 的函数 $s: \Sigma \rightarrow 2^{\Gamma^*}$ 称为代换. Σ 中的一个字符 a 在 s 的作用下为 Γ 上的一个语言 L_a , 即

$$s(a) = L_a.$$

扩展 s 的定义到字符串,

$$s(\varepsilon) = \{\varepsilon\}$$

$$s(xa) = s(x)s(a)$$

再扩展 s 到语言, 对 $\forall L \subseteq \Sigma^*$,

$$s(L) = \bigcup_{x \in L} s(x).$$

定理 35

如果有 Σ 上的 CFL L 和代换 s , 且每个 $a \in \Sigma$ 的 $s(a)$ 都是 CFL, 那么 $s(L)$ 也是 CFL.

构造方法

设 CFL L 的文法 $G = (V, T, P, S)$, 每个 $s(a)$ 的文法 $G_a = (V_a, T_a, P_a, S_a)$. 那么 $s(L)$ 的文法可以构造为

$$G' = (V', T', P', S) :$$

- ① $V' = V \cup \left(\bigcup_{a \in T} V_a \right)$
- ② $T' = \bigcup_{a \in T} T_a$
- ③ P' 包括每个 P_a 和 P 中产生式, 但是要将 P 的产生式中每个终结符 a 均替换为文法 G_a 的开始符号 S_a .

证明: 对 $\forall w \in s(L)$, 那么一定存在某个 $x = a_1 a_2 \cdots a_n \in L$ 使

$$w \in s(x) = s(a_1)s(a_2) \cdots s(a_n).$$

那么 w 可以分为 $w = w_1 w_2 \cdots w_n$ 且 $w_i \in s(a_i)$, 即

$$S_{a_i} \xrightarrow[G_{a_i}]{*} w_i.$$

由于 $S \xrightarrow[G]{*} x = a_1 a_2 \cdots a_n$, 所以

$$S \xrightarrow[G]{*} S_{a_1} S_{a_2} \cdots S_{a_n} \xrightarrow[G]{*} w_1 w_2 \cdots w_n = w,$$

所以 $w \in \mathbf{L}(G')$, 即 $s(L) \subseteq \mathbf{L}(G')$.

因为 G' 的终结符仅能由每个 S_a 派生, 因此对 $\forall w \in \mathbf{L}(G')$ 有

$$S \xrightarrow[G]{*} \alpha = S_{a_1} S_{a_2} \cdots S_{a_n} \xrightarrow[G]{*} w.$$

因为 G' 中的每个 S_a 在 G 中是终结符 a , 所以

$$S \xrightarrow[G]{*} a_1 a_2 \cdots a_n = x \in L$$

又因为 $\alpha = S_{a_1} \cdots S_{a_n} \xrightarrow[G]{*} w = w_1 \cdots w_n$, 所以 $S_{a_i} \xrightarrow[G]{*} w_i$, 即 $w_i \in s(a_i)$. 那么

$$w = w_1 w_2 \cdots w_n \in s(a_1) s(a_2) \cdots s(a_n) = s(a_1 a_2 \cdots a_n) = s(x) \subseteq s(L),$$

所以 $w \in s(L)$, 即 $\mathbf{L}(G') \subseteq s(L)$. 因此 $\mathbf{L}(G') = s(L)$. □

例3. 设 $L = \{w \in \{a, b\}^* \mid w \text{ 有相等个数的 } a \text{ 和 } b\}$, 代换

$$s(a) = L_a = \{0^n 1^n \mid n \geq 1\}$$

$$s(b) = L_b = \{ww^R \mid w \in (0 + 1)^*\}$$

求 $s(L)$ 的文法.

解: 设计 L 的文法为: $S \rightarrow aSbS \mid bSaS \mid \varepsilon$

L_a 的文法为: $S_a \rightarrow 0S_a1 \mid 01$

L_b 的文法为: $S_b \rightarrow 0S_b0 \mid 1S_b1 \mid \varepsilon$

那么 $s(L)$ 的文法为: $S \rightarrow S_aSS_bS \mid S_bSS_aS \mid \varepsilon$

$$S_a \rightarrow 0S_a1 \mid 01$$

$$S_b \rightarrow 0S_b0 \mid 1S_b1 \mid \varepsilon$$

CFL 对并/连接/闭包/同态封闭

定理 36

上下文无关语言在并, 连接, 闭包, 正闭包, 同态下封闭.

CFL 对并/连接/闭包/同态封闭

定理 36

上下文无关语言在并, 连接, 闭包, 正闭包, 同态下封闭.

证明 1: 设 $\Sigma = \{1, 2\}$, L_1, L_2 是任意 CFL. 定义代换

$$s(1) = L_1, \quad s(2) = L_2.$$

语言 $\{1, 2\}$, $\{12\}$, $\{1\}^*$ 和 $\{1\}^+$ 显然都是 CFL, 那么

- ① 由 $s(\{1, 2\}) = s(1) \cup s(2) = L_1 \cup L_2$, 所以并运算封闭;
- ② 由 $s(\{12\}) = s(12) = s(\varepsilon)s(1)s(2) = L_1L_2$, 所以连接运算封闭;

③ 闭包和正比包运算封闭, 因为

$$\begin{aligned} s(\{1\}^*) &= s(\{\varepsilon, 1, 11, \dots\}) \\ &= s(\varepsilon) \cup s(1) \cup s(11) \cup \dots \\ &= \{\varepsilon\} \cup s(1) \cup s(1)s(1) \cup \dots \\ &= L_1^*. \end{aligned}$$

若 h 是 Σ 上的同态, L 是 Σ 上的 CFL, 对 $\forall a \in \Sigma$ 令代换 $s'(a) = \{h(a)\}$, 则

$$h(L) = \{h(w) \mid w \in L\} = \bigcup_{w \in L} \{h(w)\} = \bigcup_{w \in L} s'(w) = s'(L),$$

所以同态运算封闭.



证明 2: 用文法证明并/连接/闭包的封闭性. 设 CFL L_1 和 L_2 的文法分别为

$$G_1 = (V_1, T_1, P_1, S_1), G_2 = (V_2, T_2, P_2, S_2)$$

那么, 分别构造

① $L_1 \cup L_2$ 的文法为

$$G_{\text{union}} = (V_1 \cup V_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1 \mid S_2\}, S);$$

② $L_1 L_2$ 的文法为

$$G_{\text{concat}} = (V_1 \cup V_2 \cup \{S\}, T_1 \cup T_2, P_1 \cup P_2 \cup \{S \rightarrow S_1 S_2\}, S);$$

③ L_1^* 的文法为

$$G_{\text{closure}} = (V_1 \cup \{S\}, T_1, P_1 \cup \{S \rightarrow S_1 S \mid \varepsilon\}, S).$$

再证明所构造文法的正确性, 略.



CFL 对反转封闭

定理 37

如果 L 是 CFL, 那么 L^R 也是 CFL.

证明:

设 L 的文法 $G = (V, T, P, S)$, 构造文法

$$G' = (V, T, \{A \rightarrow \alpha^R \mid A \rightarrow \alpha \in P\}, S),$$

则 $L(G') = L^R$. 证明略.



CFL 对逆同态封闭

定理 38

如果 L 是字母表 Δ 上的 CFL , h 是字母表 Σ 到 Δ^* 的同态, 那么 $h^{-1}(L)$ 也是 CFL .

CFL 对逆同态封闭

定理 38

如果 L 是字母表 Δ 上的 CFL, h 是字母表 Σ 到 Δ^* 的同态, 那么 $h^{-1}(L)$ 也是 CFL.

证明:

设 PDA $P = (Q, \Delta, \Gamma, \delta, q_0, Z_0, F)$, $\mathbf{L}(P) = L$. 构造 $\mathbf{L}(P') = h^{-1}(L)$ 的 PDA

$$P' = (Q', \Sigma, \Gamma, \delta', [q_0, \bar{\epsilon}], Z_0, F \times \{\bar{\epsilon}\}).$$

在 P' 的状态中, 使用缓冲, 暂存字符 $a \in \Sigma$ 的同态串 $h(a)$ 的后缀.

① $Q' \subset Q \times \Delta^*$: 状态 $[q, \bar{x}]$ 中的 \bar{x} 为缓冲;

② 设 $q \in Q$, 那么 δ' 定义如下:

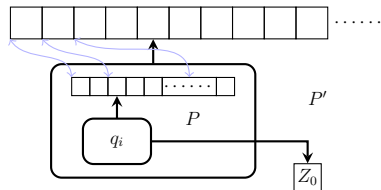
❶ $\forall [q, \bar{\varepsilon}] \in Q \times \{\bar{\varepsilon}\}, \forall a \in \Sigma, \forall X \in \Gamma$

$$\delta'([q, \bar{\varepsilon}], a, X) = \{([q, h(a)], X)\}$$

❷ 若 $\delta(q, \bar{a}, X) = \{(p_1, \beta_1), (p_2, \beta_2), \dots, (p_k, \beta_k)\}$, 则

$$\delta'([q, \overline{a\bar{x}}], \varepsilon, X) = \{([p_1, \bar{x}], \beta_1), ([p_2, \bar{x}], \beta_2), \dots, ([p_k, \bar{x}], \beta_k)\}$$

这里 $\bar{a} \in \Delta \cup \{\bar{\varepsilon}\}$, \bar{x} 是某个 $h(a)$ 的后缀.



□

CFL 对交/补运算不封闭

CFL 对交运算不封闭

因为语言

$$L_1 = \{0^n 1^n 2^i \mid n \geq 1, i \geq 1\}$$

$$L_2 = \{0^i 1^n 2^n \mid n \geq 1, i \geq 1\}$$

都是 CFL, 而

$$L_1 \cap L_2 = \{0^n 1^n 2^n \mid n \geq 1\}$$

不是 CFL.

CFL 对补运算不封闭

因为

$$L_1 \cap L_2 = \overline{\overline{L_1} \cup \overline{L_2}}.$$

定理 39

若 L 是 CFL 且 R 是正则语言, 则 $L \cap R$ 是 CFL.

证明: 设 DFA $D = (Q_1, \Sigma, \delta_1, q_1, F_1)$ 且 $\mathbf{L}(D) = R$,
PDA $P = (Q_2, \Sigma, \Gamma, \delta_2, q_2, Z_0, F_2)$ 且 $\mathbf{L}(P) = L$, 构造 PDA

$$P' = (Q_1 \times Q_2, \Sigma, \Gamma, \delta, [q_1, q_2], F_1 \times F_2)$$

其中 δ 为:

定理 39

若 L 是 CFL 且 R 是正则语言, 则 $L \cap R$ 是 CFL.

证明: 设 DFA $D = (Q_1, \Sigma, \delta_1, q_1, F_1)$ 且 $\mathbf{L}(D) = R$,
PDA $P = (Q_2, \Sigma, \Gamma, \delta_2, q_2, Z_0, F_2)$ 且 $\mathbf{L}(P) = L$, 构造 PDA

$$P' = (Q_1 \times Q_2, \Sigma, \Gamma, \delta, [q_1, q_2], F_1 \times F_2)$$

其中 δ 为:

$$\delta([p, q], a, Z) = \begin{cases} \{([p, s], \beta) \mid (s, \beta) \in \delta_2(q, a, Z)\} & a = \varepsilon \\ \{([r, s], \beta) \mid r = \delta_1(p, a) \wedge (s, \beta) \in \delta_2(q, a, Z)\} & a \neq \varepsilon \end{cases}$$

再往证 $\mathbf{L}(P') = L \cap R$, 略.



封闭性的应用

例 4. 请证明语言 L 不是 CFL

$$L = \{w \in \{a, b, c\}^* \mid n_a(w) = n_b(w) = n_c(w)\},$$

其中 $n_a(w)$ 表示 w 中 a 的个数.

证明:

- ① 因为 $\mathbf{a^*b^*c^*}$ 是正则语言,
- ② 而 $L \cap \mathbf{a^*b^*c^*} = \{a^n b^n c^n \mid n \geq 0\}$ 不是 CFL,
- ③ 由 CFL 与正则语言的交还是 CFL, 所以 L 不可能是 CFL.



上下文无关语言的性质

- 上下文无关语言的泵引理
- 上下文无关语言的封闭性
- 上下文无关语言的判定性质
- 乔姆斯基文法体系

上下文无关语言的判定性质

可判定的 CFL 问题

- 空性: 只需判断文法的开始符号 S 是否为非产生的.
- 有穷性和无穷性:
 - ① 用不带无用符号的 CNF 的产生式画有向图;
 - ② 变元为顶点, 若有 $A \rightarrow BC$, 则 A 到 B 和 C 各画一条有向边;
 - ③ 检查图中是否有循环.
- 成员性: 利用 CNF 范式, 有 CYK 算法检查串 w 是否属于 L .

CYK¹算法

- CNF $G = (V, T, P, S)$, 以 $O(n^3)$ 时间检查 “ $w = a_1a_2 \cdots a_n \in \mathbf{L}(G)$?”
- 以动态规划方式, 在表中由下至上逐行计算 X_{ij} , 再检查 “ $S \in X_{1n}$?”

$$X_{ij} = \{A \in V \mid A \Rightarrow a_i a_{i+1} \cdots a_j, 1 \leq i \leq j \leq n\},$$

- 计算首行

$$X_{ii} = \{A \mid A \rightarrow a_i \in P\}$$

- 计算其他

$$X_{ij} = \left\{ A \mid \begin{array}{l} i \leq k < j, \\ BC \in X_{ik} X_{k+1,j}, \\ A \rightarrow BC \in P \end{array} \right\}$$

$$\begin{array}{ccccc} & & & & X_{15} \\ & & & & X_{14} \quad X_{25} \\ & & & X_{13} & X_{24} \quad X_{35} \\ & & X_{12} & X_{23} & X_{34} \quad X_{45} \\ X_{11} & X_{22} & X_{33} & X_{44} & X_{55} \\ \hline a_1 & a_2 & a_3 & a_4 & a_5 \end{array}$$

¹J. Cocke, D. Younger, T. Kasami 分别独立发现了算法的基本思想

例 5. CNF G 如下, 用 CYK 算法判断 $bbabaa \in \mathbf{L}(G)$?

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$

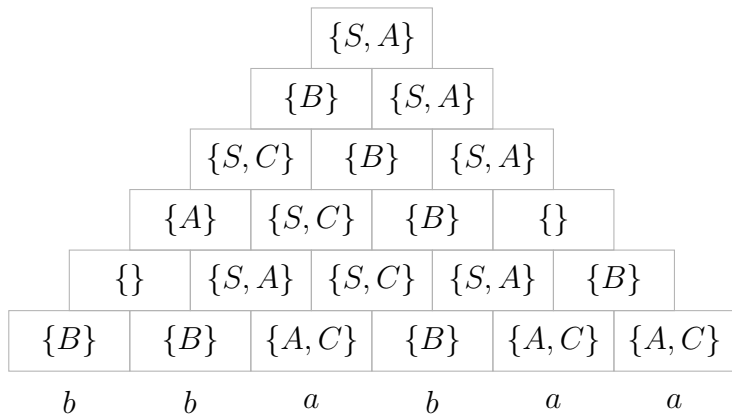
例 5. CNF G 如下, 用 CYK 算法判断 $bbabaa \in \mathbf{L}(G)$?

$$S \rightarrow AB \mid BC$$

$$A \rightarrow BA \mid a$$

$$B \rightarrow CC \mid b$$

$$C \rightarrow AB \mid a$$



因为 $S \in X_{16} = \{S, A\}$, 所以 $bbabaa \in \mathbf{L}(G)$.

上下文无关语言的判定性质

不可判定的 CFL 问题

- ① 判断 CFG G 是否歧义的?
- ② 判断 CFL 是否固有歧义的?
- ③ 两个 CFL 的交是否为空?
- ④ 两个 CFL 是否相同?
- ⑤ 判断 CFL 的补是否为空? 尽管有算法判断 CFL 是否为空
- ⑥ 判断 CFL 是否等于 Σ^* ?

上下文无关语言的性质

- 上下文无关语言的泵引理
- 上下文无关语言的封闭性
- 上下文无关语言的判定性质
- 乔姆斯基文法体系

乔姆斯基文法体系

如果文法 $G = (V, T, P, S)$, 符号串 $\alpha \in (V \cup T)^* V (V \cup T)^*$, $\beta \in (V \cup T)^*$, 产生式都形如

$$\alpha \rightarrow \beta$$

即每个产生式的左部 α 中至少要有一个变元, 那么:

- ① 称 G 为 0 型文法或短语结构文法;
- ② 若 $|\beta| \geq |\alpha|$, 称 G 为 1 型文法或上下文有关文法;
- ③ 若 $\alpha \in V$, 称 G 为 2 型文法或上下文无关文法;
- ④ 若 $\alpha \rightarrow \beta$ 都形如 $A \rightarrow aB$ 或 $A \rightarrow a$, 称 G 为 3 型文法或正则文法.

形式语言与自动机理论

图灵机

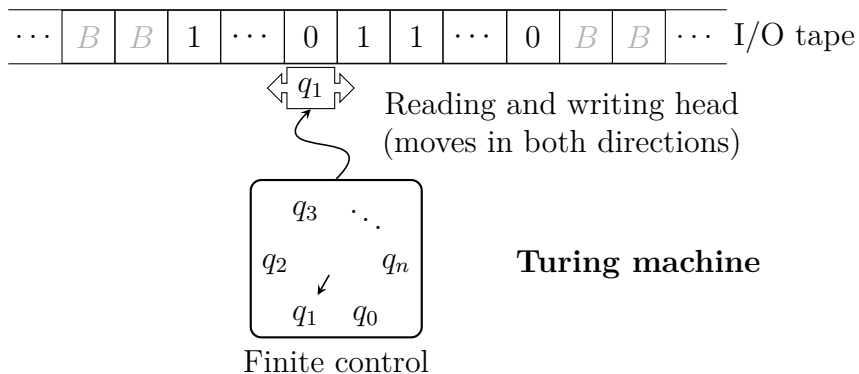
王春宇

计算机科学与技术学院
哈尔滨工业大学

图灵机

- 图灵机
 - 形式定义
 - 瞬时描述及其转移
 - 语言与停机
 - 整数函数计算器
- 图灵机的变形

图灵机



图灵机的形式定义

定义

图灵机(*TM*, *Turing Machine*) M 为七元组

$$M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$$

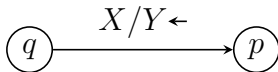
- ① Q : 有穷状态集;
- ② Σ : 有穷输入符号集;
- ③ Γ : 有穷带符号集, 且总有 $\Sigma \subset \Gamma$;
- ④ $\delta: Q \times \Gamma \rightarrow Q \times \Gamma \times \{L, R\}$ 转移函数;
- ⑤ $q_0 \in Q$: 初始状态;
- ⑥ $B \in \Gamma - \Sigma$: 空格符号;
- ⑦ $F \subseteq Q$: 终态集或接受状态集.

图灵机的动作及状态转移图

有穷控制器处于状态 q , 带头所在单元格为符号 X , 如果动作的定义为

$$\delta(q, X) = (p, Y, L),$$

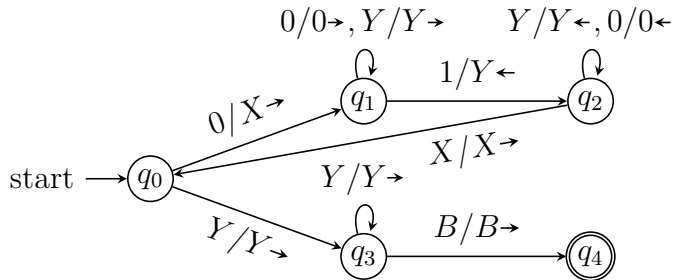
表示状态转移到 p , 单元格改为 Y , 然后带头向左移动一个单元格.



因为每个动作都是确定的, 因此是“确定的图灵机”.

例 1. 设计识别 $\{0^n 1^n \mid n \geq 1\}$ 的图灵机.

例 1. 设计识别 $\{0^n 1^n \mid n \geq 1\}$ 的图灵机.



续例 1. 设计识别 $\{0^n 1^n \mid n \geq 1\}$ 的图灵机.

$$M = (\{q_0, q_1, q_2, q_3, q_4\}, \{0, 1\}, \{0, 1, X, Y, B\}, \delta, q_0, B, \{q_4\})$$

δ	0	1	X	Y	B
q_0	(q_1, X, R)	—	—	(q_3, Y, R)	—
q_1	$(q_1, 0, R)$	(q_2, Y, L)	—	(q_1, Y, R)	—
q_2	$(q_2, 0, L)$	—	(q_0, X, R)	(q_2, Y, L)	—
q_3	—	—	—	(q_3, Y, R)	(q_4, B, R)
q_4	—	—	—	—	—

瞬时描述

定义

图灵机虽有无穷长的带, 但非空内容总是有限的. 因此用带上全部的非空符号、当前状态和带头位置, 同时定义**瞬时描述**(ID) 为

$$X_1X_2\cdots X_{i-1}qX_iX_{i+1}\cdots X_n$$

- ① 图灵机的当前状态 q ;
- ② 带头在左起第 i 个非空格符上;
- ③ $X_1X_2\cdots X_n$ 是最左到最右非空格内容.
- ④ 为避免混淆, 一般假定 Q 和 Γ 不相交.

转移符号

定义

图灵机 M 中, 如果 $\delta(q, X_i) = (p, Y, L)$, 定义 ID 转移为

$$X_1 \dots X_{i-1} q X_i \dots X_n \vdash_M X_1 \dots X_{i-2} p X_{i-1} Y X_{i+1} \dots X_n$$

如果 $\delta(q, X_i) = (p, Y, R)$ 那么

$$X_1 \dots X_{i-1} q X_i \dots X_n \vdash_M X_1 \dots X_{i-1} Y p X_{i+1} \dots X_n$$

若某 ID 是从另一个经有限步 (包括零步) 转移而得到的, 记为 \vdash_M^* .

若 M 已知, 简记为 \vdash 和 \vdash^* .

续例 1. 设计识别 $\{0^n 1^n \mid n \geq 1\}$ 的图灵机.

接受 0011 的 ID 序列 (M 的一个计算) 为

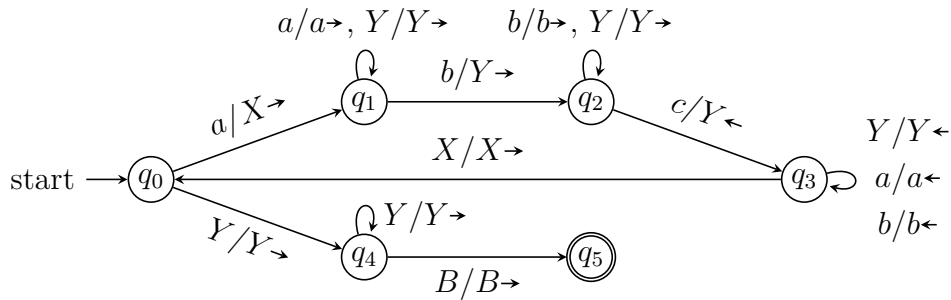
续例 1. 设计识别 $\{0^n 1^n \mid n \geq 1\}$ 的图灵机.

接受 0011 的 ID 序列 (M 的一个计算) 为

$$\begin{array}{lll}
 q_0 0011 \vdash X q_1 011 & \vdash X 0 q_1 11 & \vdash X q_2 0 Y 1 \\
 \vdash q_2 X 0 Y 1 & \vdash X q_0 0 Y 1 & \vdash X X q_1 Y 1 \\
 \vdash X X Y q_1 1 & \vdash X X q_2 Y Y & \vdash X q_2 X Y Y \\
 \vdash X X q_0 Y Y & \vdash X X Y q_3 Y & \vdash X X Y Y q_3 B \\
 \vdash X X Y Y B q_4 B
 \end{array}$$

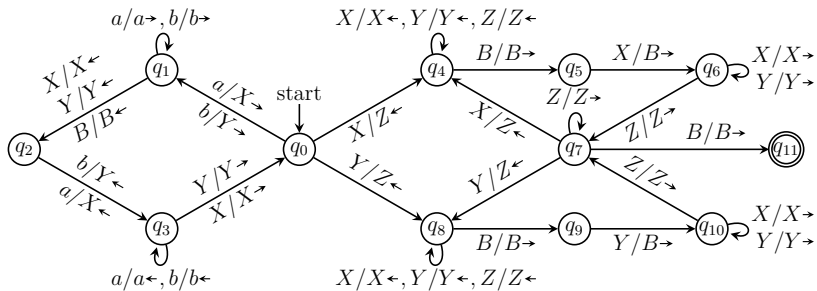
例 2. 设计接受 $L = \{a^n b^n c^n \mid n \geq 1\}$ 的图灵机.

例 2. 设计接受 $L = \{a^n b^n c^n \mid n \geq 1\}$ 的图灵机.



例 3. 设计接受 $L = \{ww \mid w \in \{a, b\}^+\}$ 的图灵机.

例3. 设计接受 $L = \{ww \mid w \in \{a, b\}^+\}$ 的图灵机.



思考

- ① DFA 和 TM 的主要区别?
- ② 计算机, 究竟是 TM 还是 DFA?

图灵机的语言

定义

如果 $M = (Q, \Sigma, \Gamma, \delta, q_0, B, F)$ 是一个图灵机, 则 M 接受的语言为

$$\mathbf{L}(M) = \{w \mid w \in \Sigma^*, q_0 w \vdash^* \alpha p \beta, p \in F, \alpha, \beta \in \Gamma^*\}.$$

定义

如果 L 是图灵机 M 的语言, 即 $L = \mathbf{L}(M)$, 则称 L 是递归可枚举语言.

- 一般假定, 当输入串被接受时, 图灵机总会停机;
- 然而, 对于不接受的输入, 图灵机可能永远不停止.

定义

对接受和不接受的输入, 都保证停机的图灵机, 所接受的语言称为**递归语言**.

算法的形式化

保证停机的图灵机, 正是算法的好模型, 即算法概念的形式化.

- λ -calculus — Alonzo Church, Stephen Kleene
- Partial recursive functions — Kurt Gödel
- Post machines — Emil Post
- Turing machines — Alan Turing

整数函数计算器

- 传统的方法, 把整数 $i \geq 0$ 写为 1 进制, 用字符串 0^i 表示;
- 若计算 k 个自变量 i_1, i_2, \dots, i_k 的函数 f , 用

$$0^{i_1} 1 0^{i_2} 1 \dots 1 0^{i_k}$$

作为 TM M 的输入;

- M 停机, 且输入带上为 0^m , 表示 $f(i_1, i_2, \dots, i_k) = m$.
- M 计算的 f , 不必对所有不同的 i_1, i_2, \dots, i_k 都有值.

定义

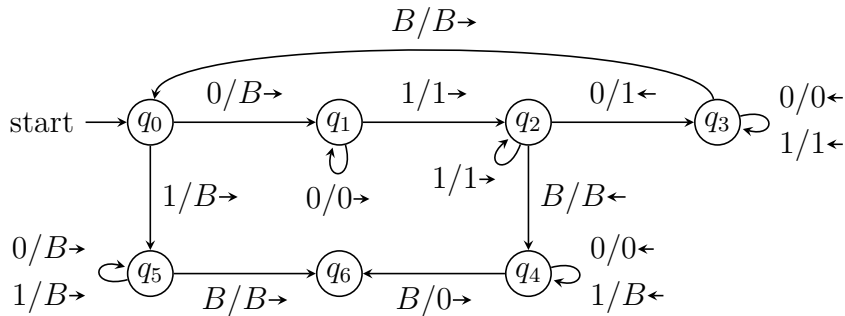
如果 $f(i_1, i_2, \dots, i_k)$ 对所有 i_1, i_2, \dots, i_k 都有定义, 称 f 为全递归函数.
被图灵机计算的函数 $f(i_1, i_2, \dots, i_k)$ 称作部分递归函数.

例 4. 给出计算整数真减法 ($\dot{-}$) 的图灵机, 其定义为

$$m \dot{-} n = \begin{cases} m - n & m \geq n \\ 0 & m < n \end{cases}.$$

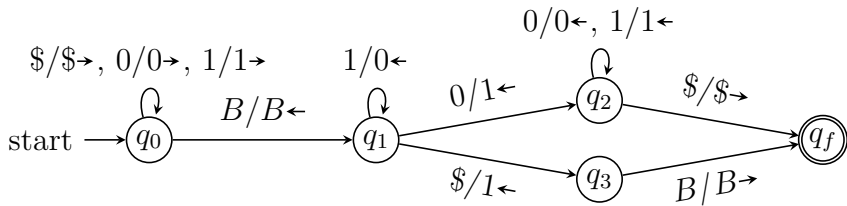
例 4. 给出计算整数真减法 ($\dot{-}$) 的图灵机, 其定义为

$$m \dot{-} n = \begin{cases} m - n & m \geq n \\ 0 & m < n \end{cases}.$$



例 5. 二进制数的加 1 函数, 使用符号 \$ 作为数字前的占位标记.
例如 $q_0\$10011 \vdash^* \q_f10100 , $q_0\$111 \vdash^* q_f1000$.

例 5. 二进制数的加 1 函数, 使用符号 \$ 作为数字前的占位标记.
 例如 $q_0 \$10011 \vdash^* \$q_f 10100$, $q_0 \$111 \vdash^* q_f 1000$.



图灵机

- 图灵机
- 图灵机的变形
 - 扩展的图灵机
 - 受限的图灵机

状态中存储

有限控制器中可以存储有限个符号的图灵机:

$$M' = (Q', \Sigma, \Gamma, \delta, q'_0, B, F')$$

其中 $Q' = Q \times \Gamma \times \cdots \times \Gamma$, $q'_0 = [q_0, B, \cdots, B]$.

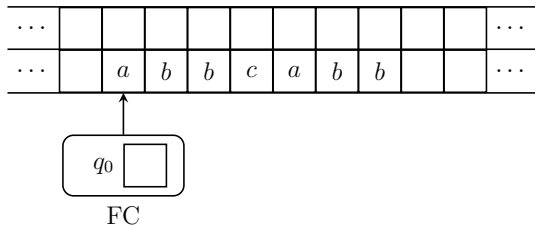
多道

多道图灵机:

$$M' = (Q, \Sigma, \Gamma', \delta, q_0, B', F)$$

其中 $\Gamma' = \Gamma \times \Gamma \times \cdots \times \Gamma$.

例 6. 利用状态中存储与多道设计 TM 识别 $L = \{w cw \mid w \in \{a, b\}^*\}$.



子程序

设计 TM 的一部分作为一个子程序:

- 具有一个指定的初始状态;
- 具有一个指定的返回状态, 但暂时没有定义动作;
- 可以具有参数和返回值.

通过进入子程序的初始状态, 实现调用; 通过返回状态的动作, 实现返回.

例 7. 设计 TM 实现全递归函数“乘法”.

扩展的图灵机

多带图灵机

有穷控制器、 k 个带头和 k 条带组成. 每个动作, 根据状态和每个带头符号:

- ① 改变控制器中的状态;
- ② 修改带头单元格中的符号;
- ③ 每个带头独立的向左或右移动一个单元格, 或保持不动.

开始时, 输入在第 1 条带上, 其他都是空的, 其形式定义非常繁琐.

定理 40

如果语言 L 被一个多带图灵机接受, 那么 L 能够被某个单带图灵机接受.

证明方法:

- ① 用 $2k$ 道的单带图灵机 N 模拟 k 带图灵机 M ;
- ② N 用两道模拟 M 一带, 一道放置内容, 另一道标记带头;
- ③ 模拟 M 的一个动作, N 需要从左至右, 再从右至左扫描一次;
- ④ 第一次扫描搜集当前格局, 第二次扫描更新带头和位置.

非确定图灵机 (NTM)

在每个状态 q 和每个带符号 X 的转移, 可以有有限个选择的图灵机, 即

$$\delta(q, X) = \{(q_1, Y_1, D_1), (q_2, Y_2, D_2), \dots, (q_k, Y_k, D_k)\}.$$

图灵机增加了非确定性, 并未改变图灵机接受语言的能力.

定理 41

如果 L 被非确定图灵机接受, 那么 L 被图灵机接受.

证明方法:

- ① 同样用多带技术, 用确定的 TM M 模拟 NTM N ;
- ② M 用第 1 条带存储 N 未处理的 ID, 用第 2 条带模拟 N ;
- ③ M 从第 1 条带取 N 的当前 ID 放到第 2 带;
- ④ 若不接受, 把当前 ID 可能的 k 个转移 ID 复制到第 1 条带的最末端;
- ⑤ 然后循环, 继续从第 1 带取下一个 ID 去模拟.

思考题

为什么非确定性没有改变图灵机识别语言的能力？

多维图灵机

- ① 这种装置具有通常的有穷控制器;
- ② k 维阵列组成的带, 在 $2k$ 个方向上都是无限的;
- ③ 根据状态和读入符号改变状态, 并沿着 k 个轴的正和负向移动;
- ④ 开始时, 输入沿着某一个轴排列, 带头在输入的左端.

同样, 这样的扩展也没有增加额外的能力, 仍然等价于基本的图灵机.

受限的图灵机

半无穷带图灵机

图灵机的输入输出带只有一侧是无穷的.

定理 42

半无穷带图灵机, 与图灵机等价.

证明方法:

一侧无穷的图灵机带, 可使用多道技术, 模拟双侧无穷的图灵机带.

多栈机

基于下推自动机的扩展, k 栈机器是具有 k 个栈的确定型下推自动机.

定理 43

如果图灵机接受 L , 那么双栈机接受 L .

证明方法:

- ① 一个堆栈保存带头左边内容, 一个堆栈保存带头右边内容;
- ② 带头的移动用两个栈分别弹栈和压栈模拟;
- ③ 带头修改字符 A 为 B , 用一个栈弹出 A 而另一个压入 B 来模拟;
- ④ 开始时输入在双栈机的输入带, 但先将输入扫描并压入一个栈, 再依次弹出并压入另一个栈, 然后开始模拟图灵机.

例 8. 利用双栈机器接受 $L = \{a^n b^n c^n \mid n \geq 0\}$ 和 $L = \{a^n b^n c^n d^n e^n \mid n \geq 0\}$.

形式语言与自动机理论

不可判定性

王春宇

计算机科学与技术学院
哈尔滨工业大学

不可判定性

- 不可判定性
- 非递归可枚举的语言
- 递归可枚举但非递归的语言
- 语言类的关系

不可判定问题

典型问题

给定语言 $L \subseteq \Sigma^*$ 和字符串 $w \in \Sigma^*$, 判断是否 $w \in L$ 的问题, 称为语言 L 上的一个判定性问题.

(非形式) 定义

如果一个问题, 不存在能解决它的程序, 则称为不可判定的.

不可判定问题

典型问题

给定语言 $L \subseteq \Sigma^*$ 和字符串 $w \in \Sigma^*$, 判断是否 $w \in L$ 的问题, 称为语言 L 上的一个判定性问题.

(非形式) 定义

如果一个问题, 不存在能解决它的程序, 则称为不可判定的.

是否存在不可判定的问题?

- ① $\{L \mid L \in \Sigma^*\}$ 是不可数的;
- ② $\{P \mid P \text{ 是一个程序}\}$ 是可数的;
- ③ 问题显然比程序多, 必然存在不可判定的问题.

hello-world 问题

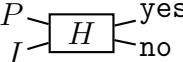
判断带有给定输入的任意给定的程序, 是否以 hello, world 为其输出的前 12 个字符.

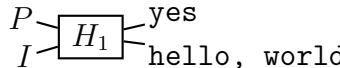
定理

hello-world 问题是不可判定的.

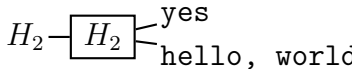
- “具有这个输入的这个程序是否显示 hello, world?”
- 解决这样问题的通用程序是不存在的.

(非形式) 证明: 反证法. 假设这样的程序 H 存在.

① H 检查程序 P 在输入 I 时的输出, 并回答 yes 或 no: 

② 修改 H , 在回答 no 时, 输出 hello, world: 

③ 再修改 H_1 , 将程序 P 作为 P 的输入: 

④ 那么, 当程序 H_2 以 H_2 为输入时: 

⑤ H_2 的输出会出现矛盾 (悖论), 所以 H 不可能存在. □

问题的归约

如何证明问题是不可判定的？

- ① 归谬法 (反证法)
- ② 问题的归约

不可判定问题

- 递归可枚举语言 — 图灵机所识别
- 递归语言 — 保证停机的图灵机所识别

定义

一个问题, 如果它的语言是递归的, 就称为**可判定问题**, 否则称为**不可判定问题**.

不可判定的问题

- 不存在保证停机的图灵机识别该问题的语言
- 不存在解决该问题的算法

不可判定性

- 不可判定性
- 非递归可枚举的语言
 - 第 i 个串
 - 图灵机编码与第 i 个图灵机
 - 对角化语言 L_d
- 递归可枚举但非递归的语言
- 语言类的关系

可判定吗?

“图灵机 M 接受输入 w 吗?”

第 i 个串 w_i

定义

将全部 $(0+1)^*$ 中的字符串按长度和字典序排序,
那么第 i 个串就是 w_i . 且刚好有

$$\text{binary}(i) = 1w_i.$$

比如:

i	1	2	3	4	5	6	7	8	9	...
$\text{binary}(i)$	1ε	10	11	100	101	110	111	1000	1001	...
w_i	ε	0	1	00	01	10	11	000	001	...

图灵机编码

将 $\Sigma = \{0, 1\}$ 上的全部图灵机, 用二进制字符串编码

$$M = (Q, \Sigma, \Gamma, \delta, q_1, B, F)$$

- ① $Q = \{q_1, q_2, \dots, q_{|Q|}\}$, 开始状态为 q_1 , 终态为 q_2 且停机;
- ② $\Gamma = \{X_1, X_2, \dots, X_{|\Gamma|}\}$, 总有 $X_1 = 0, X_2 = 1, X_3 = B$;
- ③ 设带头移动方向 $D_1 = L, D_2 = R$;
- ④ 任意的转移 $\delta(q_i, X_j) = (q_k, X_l, D_m)$ 编码为

$$C = 0^i 10^j 10^k 10^l 10^m;$$

- ⑤ 则全部 n 个转移的编码合并在一起, 作为图灵机 M 的编码:

$$C_1 11 C_2 11 \dots C_{n-1} 11 C_n.$$

第 i 个图灵机 M_i

定义

如果图灵机 M 的编码为第 i 个串 w_i , 则称 M 是第 i 个图灵机 M_i .

- 任意图灵机 M 都对应一个字符串 w
- 任意的字符串 w 都可以看作图灵机的编码
- 如果编码不合法, 将其看作接受 \emptyset 且立即停机的图灵机

非递归可枚举的语言

定义

使第 i 个串 w_i 不属于第 i 个图灵机 M_i 的语言 $L(M_i)$ 的所有 w_i 的集合, 称为**对角化语言** L_d , 即

$$L_d = \{w_i \mid w_i \notin \mathbf{L}(M_i), i \geq 1\}.$$

Diagram illustrating a matrix M_i with a shaded diagonal band. The matrix is indexed by i (rows) and j (columns). The diagonal band is shaded gray. The matrix is shown with rows 1 to 6 and columns 1 to 6, with ellipses indicating continuation. The diagonal band is defined by the condition $|i - j| \leq 1$.

定理 44

L_d 不是递归可枚举语言, 即不存在图灵机接受 L_d .

证明: 反证法.

假设存在识别 L_d 的图灵机 M , 那么 M 也可被编码, 不妨设第 i 个图灵机 $M_i = M$, 即 $\mathbf{L}(M_i) = L_d$.

那么, 考虑第 i 个串 w_i 是否会被 M_i 识别:

- ① 如果 $w_i \in \mathbf{L}(M_i) = L_d$, 那么由 L_d 的定义, 又有 $w_i \notin \mathbf{L}(M_i)$;
- ② 如果 $w_i \notin \mathbf{L}(M_i)$, 那么由 L_d 的定义, 又有 $w_i \in L_d = \mathbf{L}(M_i)$.

无论如何都会矛盾, 因此假设不成立, 不存在接受 L_d 的图灵机.



不可判定性

- 不可判定性
- 非递归可枚举的语言
- 递归可枚举但非递归的语言
 - 递归语言的封闭性
 - 通用语言与通用图灵机
- 语言类的关系

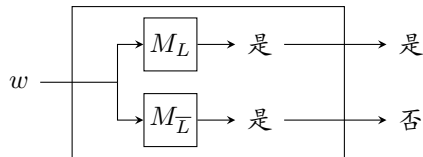
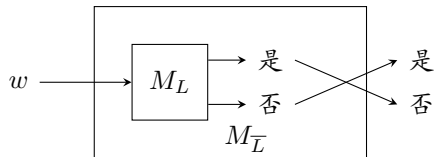
递归语言的封闭性

定理 45

如果 L 是递归的, 那么 \bar{L} 也是递归的.

定理 46

如果语言 L 和 \bar{L} 都是递归可枚举的, 那么 L 是递归的.



通用语言与通用图灵机

定义

如果图灵机 M 接受串 w , 那么由 $M111w$ 表示的有序对 (M, w) 构成的语言 L_u , 称为通用语言

$$L_u = \{M111w \mid w \in \mathbf{L}(M)\}.$$

定义

构造图灵机 U , 当输入 $M111w$ 时, 利用多带技术模拟 M 处理串 w 的过程. 因为 M 接受 w 时会停机, 因此 U 可以识别 L_u , 图灵机 U 称为通用图灵机.

递归可枚举但非递归的语言

定理 47

通用语言 L_u 是递归可枚举的, 但不是递归的.

证明: L_u 是递归可枚举的. 用反证法证明 L_u 不是递归的.

通用图灵机 U 使用 3 条带分别:

(1) 装载 M 的编码; (2) 放置 w , 模拟 M 的带; (3) 存储 M 的状态.

假设存在算法 A 识别 L_u , 那么可如下得到识别对角化语言 L_d 的算法 B .

将 B 的输入 $w = w_i$ 转换为 M_i111w_i 交给 A 判断:

- 当 A 接受, 表示 $w_i \in L(M_i)$, 则 B 拒绝;
- 当 A 拒绝, 表示 $w_i \notin L(M_i)$, 则 B 接受.

而由于 L_d 不是递归的, 所以 B 不可能存在, 所以 L_u 不可能是递归的. □

通用图灵机的重要意义

- 识别 L_u 的通用图灵机 U , 可以模拟任意图灵机
- 冯·诺伊曼通用数字电子计算机体系结构设计思想的灵感来源
- 抽象理论的先期发展可以对实际问题有很大帮助

不可判定性

- 不可判定性
- 非递归可枚举的语言
- 递归可枚举但非递归的语言
- 语言类的关系

语言类的关系

