

编译原理 (南京大学软件学院)

Spring 2024

第 3 节: 词法分析

Lecturer: 魏恒峰

Scribes: 张桃玮

§1 词法分析

1.1. 语言. 语言是字符串的集合.

定义 1.1 (字母表). 字母表 Σ 是一个有限的符号集合.

现在, 符号没有任何的含义. 符号是什么意思是语义干的事情.

定义 1.2 (字母表). 字母表 Σ 上的串 (s) 是由 Σ 中符号构成的一个有穷序列.其中特殊的串是空串, $|\epsilon| = 0$ 定义 1.3 (串上的连接运算). 例如 $x = \text{dog}, y = \text{house}$ $xy = \text{doghouse}$, $s\epsilon = \epsilon s = s$

定义 1.4 (串上的指数运算).

$$s^0 \triangleq \epsilon$$

$$s^i \triangleq ss^{i-1}, i > 0$$

定义 1.5. 语言是给定字母表 Σ 上一个任意的可数的串集合。

注:

1. 注意 \emptyset 和 $\{\epsilon\}$. 后者是只有空串的语言.
2. 例如 $\text{ws} : \{\text{blank}, \text{tab}, \text{newline}\}$
3. 这就可以通过集合操作构造新的语言.

定义 1.6. 我们可以通过如下的方式构造新语言.

运算	定义和表示
L 和 M 的并	$L \cup M := \{s \mid s \text{ 属于 } L \text{ 或者 } s \text{ 属于 } M\}$
L 和 M 的连接	$LM := \{st \mid s \text{ 属于 } L \text{ 且 } t \text{ 属于 } M\}$
L 的 Kleene 闭包	$L^* := \cup_{i=0}^{\infty} L^i$
L 的正闭包	$L^+ := \cup_{i=1}^{\infty} L^i$

注: $L^* (L^+)$ 允许我们构造无穷集合例子 1.1. 假设 $L = \{A, B, \dots, Z, a, b, \dots, z\}, D = \{0, 1, \dots, 9\}$. 那么

- $L \cup D$ = 所有的大小写字母和所有的一位数码构成的集合.
- $|LD| = 52 \times 10 = 520$.
- L^4 = 所有长度为 4 的由字母构成的语言的集合.

- L^* = 所有字母构成的字符串 (包含空串) 构成的集合.
- D^+ = 所有数字构成的字符串 (不包含空串) 构成的集合.
- $L(L \cup D)^*$ = 以字母开头, 跟上 0 个或者若干个字母或者数字的字符串构成的集合. (也就是类似于 id).

1.2. 正则表达式. 注意区分语法和语义. 例如正则表达式中, 每个正则表达式 r 对应一个正则语言 $L(r)$. 正则表达式是语法, 正则语言是语义.

定义 1.7 (正则表达式). 给定字母表 Σ , Σ 上的正则表达式由且仅由以下规则定义:

1. ϵ 是正则表达式;
2. $\forall a \in \Sigma, a$ 是正则表达式;
3. 如果 r 是正则表达式, 则 (r) 是正则表达式;
4. 如果 r 与 s 是正则表达式, 则 $r | s, rs, r^*$ 也是正则表达式.

运算优先级: $() > * > \text{连接} > |$

例如

$$(a) | ((b)^*(c)) \equiv a | b^*c.$$

每个正则表达式 r 对应一个正则语言 $L(r)$, 这代表了它的语义.

定义 1.8 (正则表达式对应的正则语言).

$$L(\epsilon) = \{\epsilon\}$$

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

$$L((r)) = L(r)$$

$$L(r | s) = L(r) \cup L(s) \quad L(rs) = L(r)L(s) \quad L(r^*) = (L(r))^*$$

例子 1.2. 如果 $\Sigma = \{a, b\}$, $L(a | b) = \{a, b\}$. 那么

- $L(a^*) = \{\epsilon, a, aa, aaa, \dots\}$.
- $L((a | b)^*)$ = 由 a 和 b 构成的任意长度的字符串 (包括空串).
- $L(a | a^*b) = \{a, b, ab, aab, aaab, \dots\}$.

实际生活中, 我们的正则表达式的列表如下:

表达式	匹配	例子
c	单个非运算符字符 c	a
$\backslash c$	字符 c 的字面值	$\backslash *$
$"s"$	串 s 的字面值	$"**"$
$.$	除换行符以外的任何字符 (看环境, ANTLR4 中可以匹配换行符)	$a.*b$
$^$	一行的开始	abc
$\$$	行的结尾	$abc\$$
$[s]$	字符串 s 中的任何一个字符	$[abc]$
$[\^s]$	不在串 s 中的任何一个字符	$^\wedge[abc]$
r^*	和 r 匹配的零个或多个串连接成的串	a^*
r^+	和 r 匹配的一个或多个串连接成的串	a^+
$r^?$	零个或一个 r	$a^?$
$r\{m,n\}$	最少 m 个, 最多 n 个 r 的重复出现	$a\{1,5\}$
r_1r_2	r_1 后加上 r_2	ab
$r_1 \mid r_2$	r_1 或 r_2	$a \mid b$
(r)	与 r 相同	$(a \mid b)$
r_1/r_2	后面跟有 r_2 时的 r_1	$abc/123$

有一些简单记录方法 (在 Vim 中)

- $\backslash w$ 表示所有大小写字母, 数字, 以及下划线
- $\backslash W$ 表示除去所有大小写字母, 数字, 以及下划线
- $\backslash d$ 表示所有数码
- $\backslash D$ 表示除去所有数码

例子 1.3.

$$(0 \mid (1(01^*0)^*1))^*$$

表示的二进制的 3 的倍数. 可以从 regex.com 中观察一些例子.

接下来介绍自动机.

定义 1.9 (NFA). 非确定性有穷自动机 \mathcal{A} 是一个五元组 $\mathcal{A} = (\Sigma, S, s_0, \delta, F)$:

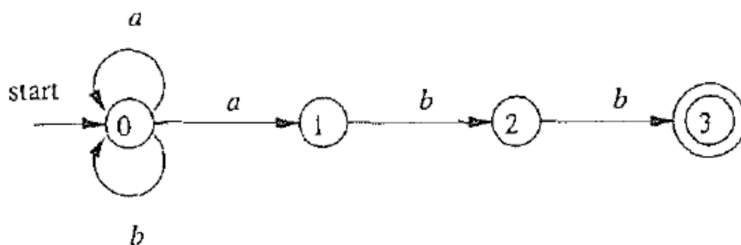
1. 字母表 $\Sigma (\epsilon \notin \Sigma)$;
2. 有穷的状态集合 S ;
3. (唯一) 的初始状态 $s_0 \in S$;
4. 状态转移函数 δ .

$$\delta : S \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^S$$

5. 接受状态集合 $F \subseteq S$

注:

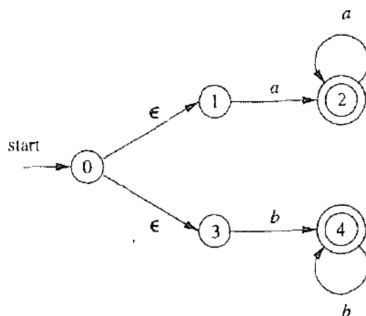
1. 非确定性指的是其“出路”可能不唯一. 此外, 可以接受空串然后进行转移 (称为 ϵ 转移).
2. 所有没有对应出边的字符默认指向“空状态” \emptyset . 此状态无论接受什么字符都到自身. 一旦进入将无法出去.



下面来看 NFA 的语义. (非确定性) 有穷自动机是一类极其简单的计算装置, 它可以识别 (接受/拒绝) Σ 上的字符串.

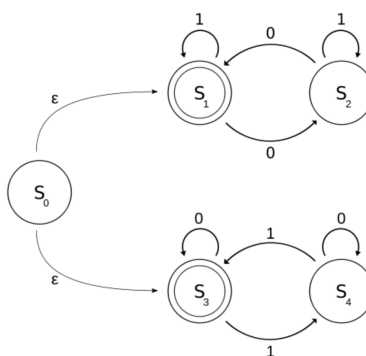
定义 1.10 (接受 (Accept)). (非确定性) 有穷自动机 \mathcal{A} 接受字符串 x , 当且仅当存在一条从开始状态 s_0 到某个接受状态 $f \in F$ 、标号为 x 的路径.

例子 1.4. 考虑下面的自动机:



对于上面的一个状态图, $L(\mathcal{A}) = L((aa^* | bb^*))$.

例子 1.5. 考虑下面的自动机:



实际上 $L(\mathcal{A}) = \{ \text{包含偶数个 1 或偶数个 0 的 01 串} \}$.

因此, \mathcal{A} 定义了一种语言 $L(\mathcal{A})$: 它能接受的所有字符串构成的集合.

关于自动机 \mathcal{A} 的两个基本问题:

1. Membership 问题: 给定字符串 $x, x \in L(\mathcal{A})$?
2. $L(\mathcal{A})$ 究竟是什么?

定义 1.11 (DFA). 确定性有穷自动机 \mathcal{A} 是一个五元组 $\mathcal{A} = (\Sigma, S, s_0, \delta, F)$:

- 字母表 $\Sigma (\epsilon \notin \Sigma)$
- 有穷的状态集合 S
- 唯一的初始状态 $s_0 \in S$
- 状态转移函数 δ

$$\delta : S \times \Sigma \rightarrow S$$

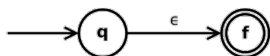
- 接受状态集合 $F \subseteq S$

上述约定 (所有没有对应出边的字符默认指向一个“死状态”) 同样适用于这里.

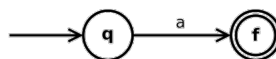
NFA 简洁易于理解, 便于描述语言 $L(\mathcal{A})$ DFA 易于判断 $x \in L(\mathcal{A})$, 适合产生词法分析器. 下面我们来走 $RE \Rightarrow NFA \Rightarrow DFA \Rightarrow$ 词法分析器的流程.

a) 正则表达式到 NFA 要求 $L(N(r)) = L(r)$. 可以使用 Thompson 构造法. 使用四种基础的归纳.

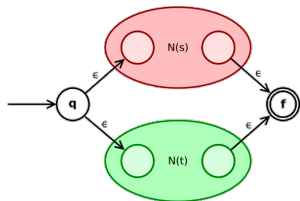
ϵ 是正则表达式。



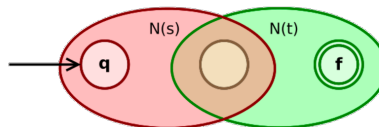
$a \in \Sigma$ 是正则表达式。



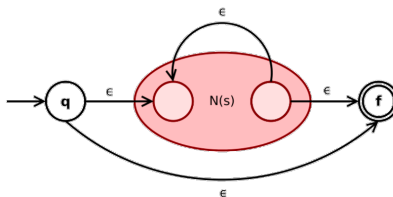
如果 s, t 是正则表达式, 则 $s|t$ 是正则表达式。



如果 s, t 是正则表达式, 则 st 是正则表达式。



如果 s 是正则表达式, 则 s^* 是正则表达式。

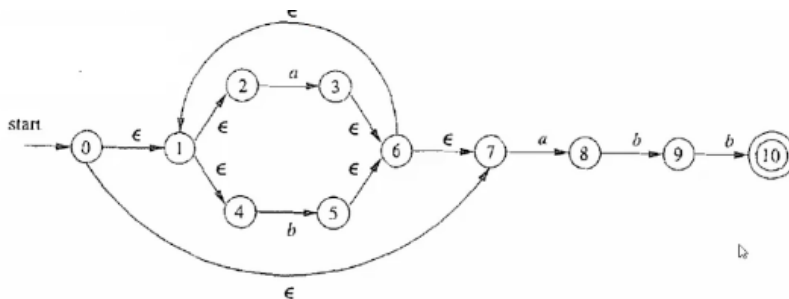


$N(r)$ 的性质以及 Thompson 构造法复杂度分析

1. $N(r)$ 的开始状态与接受状态均唯一。
2. 开始状态没有人边, 接受状态没有出边。
3. $N(r)$ 的状态数 $|S| \leq 2 \times |r|$ 。
4. 每个状态最多有两个 ϵ -入边与两个 ϵ -出边。
5. $\forall a \in \Sigma$, 每个状态最多有一个 a -入边与一个 a -出边。 ($|r|$: r 中运算符与运算分量的总和)

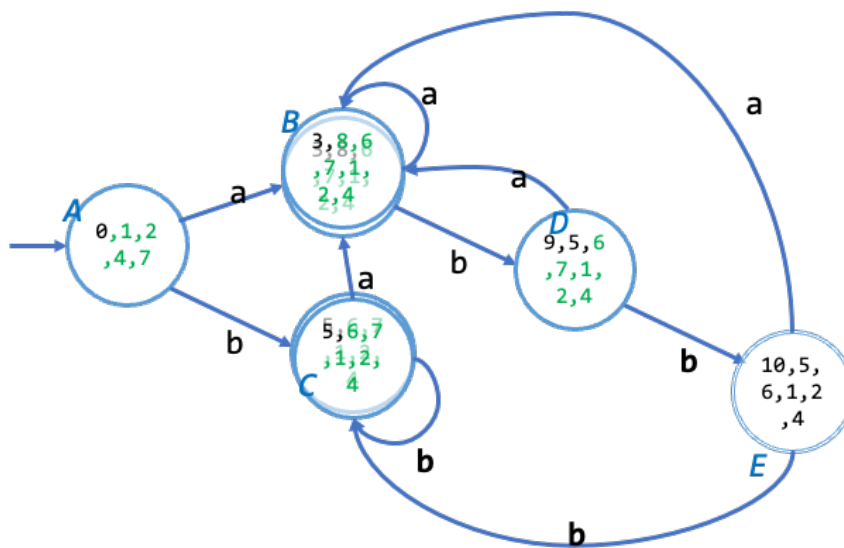
b) 从 NFA 到 DFA 的转换 主要思想: 使用 DFA 模拟 NFA – 子集构造法

假设我们想要把下图的 NFA(用正则语言表示就是 $L(\mathcal{A}) = L((a|b)^*abb)$) 转换为 DFA, 我们可以



- 找到 NFA 的初始状态. 也就是从起始点开始, 仅通过 ϵ 转移可以到达的所有的状态的集合. 如 0, 1, 2, 4, 7 号.
- 从这个当前的“等价状态集合”每个元素出发, 每次消耗字母表中的一个相同字符, 每一个状态往前走一步, 看一看到达哪些状态. 然后进行 ϵ 转移. 然后把等价的状态集合合并.
- 直到不会再出现新的状态了.

上图经过转换的等价的 DFA 是



用表格表示就是

NFA 状态	DFA 状态	a	b
{0, 1, 2, 4, 7}	A	B	C
{1, 2, 3, 4, 6, 7, 8}	B	B	D
{1, 2, 4, 5, 6, 7}	C	B	C
{1, 2, 4, 5, 6, 7, 9}	D	B	E
{1, 2, 4, 5, 6, 7, 10}	E	B	C

回顾 NFA 的接受是存在一条路径可以接受, 那么根据状态的等价性, 等价的状态都是可以接受的. 因此状态 E 是接受状态.

将刚刚的内容形式化:

定义 1.12 (ϵ 闭包). 从状态 s 开始, 只通过 ϵ -转移可达的状态集合

$$\epsilon\text{-closure}(s) = \{t \in S_N \mid s \xrightarrow{\epsilon^*} t\}$$

并且我们可以构建一个状态集合的 ϵ 闭包, 就是看每一个元素的 ϵ 闭包之并. 也就是

$$\epsilon\text{-closure}(T) = \bigcup_{s \in T} \epsilon\text{-closure}(s)$$

定义 1.13 (move 描述). 那么当前在的状态往下读入一个字符, 比如 a , 定义做一次 move, 定义为

$$\text{move}(T, a) = \bigcup_{s \in T} \delta(s, a)$$

那么子集构造法对于一个 NFA 转移为 DFA 的原理就是假设有 NFA 和 DFA, 分别叫做 N 和 D ,

$$N : (\Sigma_N, S_N, n_0, \delta_N, F_N)$$

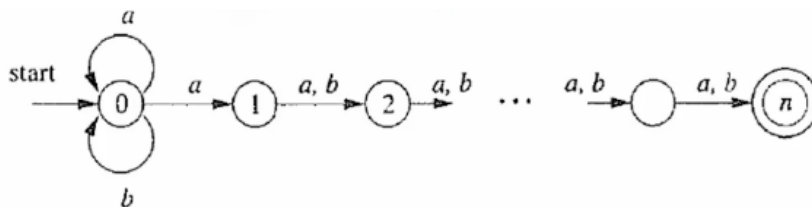
$$D : (\Sigma_D, S_D, d_0, \delta_D, F_D)$$

我们构造的 D 满足

1. 字母表: $\Sigma_D = \Sigma_N$
2. 状态集合: $S_D \subseteq 2^{S_N} \quad (\forall s_D \in S_D : s_D \subseteq S_N)$
3. 初始状态: $d_0 = \epsilon\text{-closure}(n_0)$
4. 转移函数: $\forall a \in \Sigma_D : \delta_D(s_D, a) = \epsilon\text{-closure}(\text{move}(s_D, a))$
5. 接受状态集: $F_D = \{s_D \in S_D \mid \exists f \in F_N. f \in s_D\}$

下面分析构造复杂度. 假设 $|S_N| = n$, 构造出来的 DFA 的状态数不会超过 2^n . 这从状态集合的定义可以立即得出. 也就是 $|S_D| = \Theta(2^n)$. 实际上, 遗憾的事情是, 构造出的下界也是 2^n : 存在一个用例, 它恰好生成一个 2^n 个状态的 DFA. 也就是 $|S_D| = O(2^n) \cap \Omega(2^n)$. 实际上, 可以证明, 对于任何算法, 最坏情况下, $|S_D| = \Omega(2^n)$. 这说明这个复杂度是问题本身导致的, 而不是取决于发明算法人的聪明程度.

下面给出这个用例: “长度为 $m \geq n$ 个字符的 a, b 串, 且倒数第 n 个字符是 a ”. 其中 n 是给定的整数. 这用正则语言写作 $L_n = (a \mid b)^* a (a \mid b)^{n-1}$. 用 NFA 表示就是如下图.

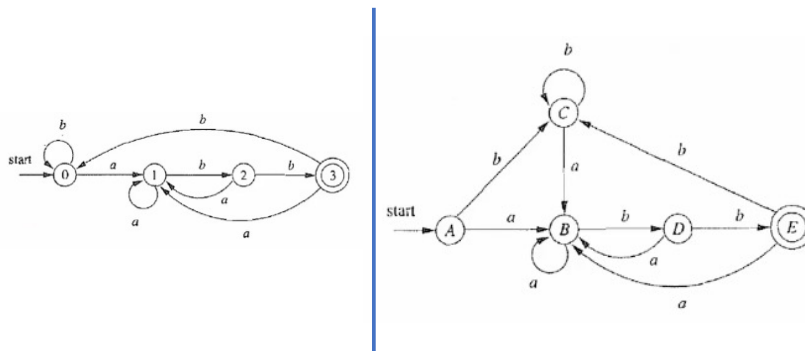


其对应的 DFA 是



下面回顾闭包的概念. 我们通常会说一个 “ f -closure (T)”, 这表示我们在初始的集合 T 上对于每一个元素重复作用函数 f , 并且不断地把新得到的元素放到 T 中. 直到继续应用 f 之后没有新的元素加入集合为止. 换句话说, 也就是构造一系列的序列, $T \Rightarrow f(T) \Rightarrow f(f(T)) \Rightarrow f(f(f(T))) \Rightarrow \dots$, 直到找到 x 使得 $f(x) = x$ (x 称为 f 的不动点) 的时候结束.

c) DFA 最小化 比如语言 $L(\mathcal{A}) = L((a|b)^*abb)$, 下图的左边和右边同样识别了我们的语言, 但是左边的更简洁. 我们希望让形如右图的 DFA 简化为左图的. 这个过程叫做 DFA 的最小化.

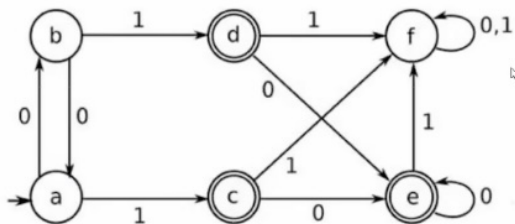


其核心思想是 “等价的状态可以合并”. 但是定义等价状态绝非易事. 首先的尝试是如果状态 s 和 t 等价, 那么

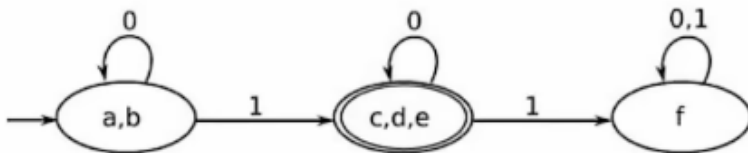
$$s \sim t \iff \forall a \in \Sigma. \left(\left(s \xrightarrow{a} s' \right) \wedge \left(t \xrightarrow{a} t' \right) \right) \implies (s' = t')$$

但是这个定义行不通. 我们考察上图右侧的例子. 可以发现 $A \sim C \sim E$. (通过 a 都可以到达状态 B, 通过 b 都可以到达 C) 但是他们显然不等价! 因为 A, C 是非结束状态, E 是结束状态. 我们把不等价的状态认做了等价的状态.

下面的例子展示了我们这样的定义同样可以把等价的状态认做不等价的状态. 如下图



例如 $a \approx b$, 通过相同的字符到达的状态都不是同一个状态. 但是这两个状态是等价的. 如下图所示.



这个定义需要错的地方在于: 经过同一个字符, 我们应当允许到达不同的状态, 但是到达的这个不同的状态要等价. 也就是原来的定义太“短视”了. 它只考虑了一步的情况, 忽略了很多步最后等价的情况. 因而, 只要取一个反就好了.

定义 1.14. DFA 中的两个状态不等价, 定义如下:

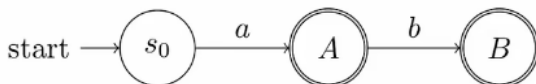
$$s \approx t \iff \exists a \in \Sigma. \left(s \xrightarrow{a} s' \right) \wedge \left(t \xrightarrow{a} t' \right) \wedge (s' \not\approx t')$$

也就是可以像这样定义两个状态的等价:

$$s \sim t \iff \forall a \in \Sigma. \left(\left(s \xrightarrow{a} s' \right) \wedge \left(t \xrightarrow{a} t' \right) \right) \implies (s' \sim t')$$

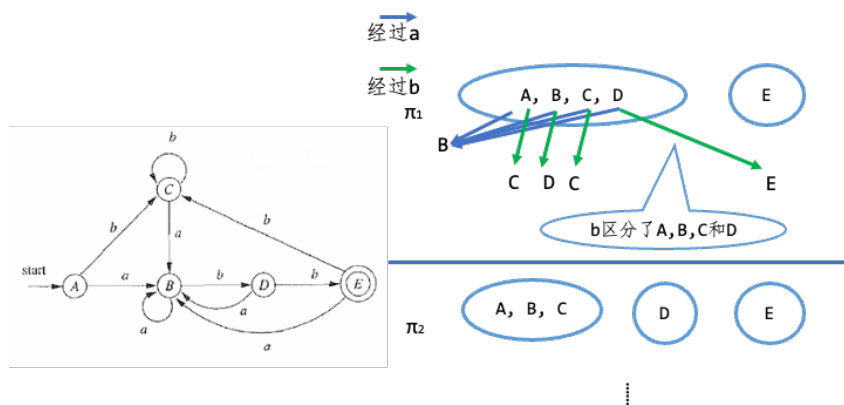
也就是存在一个字符 a. 可以“区分”两个状态, 使他们在相同的字符的驱使下到达了不同的状态, 也就是 s' 和 t' 不等价. 这样的定义就要求我们, 不断合并等价的状态, 直到无法合并为止. 但是, 这是一个递归定义, 从哪里开始呢?

首先, 所有接受状态都是等价的吗? 如下图的自动机, A 可以接受的是 a, B 可以接受的是 ab, 他们显然不等价. 因为如果我们可以把他俩合并的话, 整个自动机就剩下 2 个状态了, 自然不可能接受长度为 2 的字符串了.



我们的基础情况就出现问题了. 既然正面很难入手, 考虑问题的反面. 我们来划分状态, 而非合并. 首先假设所有状态都是等价的, 然后根据状态不等价的定义进行“分裂”. 这样, 我们首先知道所有的接受状态和非接受状态肯定不等价, 于是我们现在的划分是 $\Pi = \{F, S \setminus F\}$. 这也就是我们递归的基础. 然后, 根据不等价的定义 $s \approx t \iff \exists a \in \Sigma. (s \xrightarrow{a} s') \wedge (t \xrightarrow{a} t') \wedge (s' \approx t')$ 继续分裂.

用下图的例子, 我们首先将 A, B, C, D 归为一类, E 分为一类. 接下来对于每一类, 寻找一类里面是不是存在“这个类里面的元素受相同的字符的驱使, 随后却到达不同的状态”的情况. 如果是这样, 我们就做分裂. 然后继续这样做. 下图展示了第一次操作.



最后的迭代过程就是

$$\Pi_0 = \{\{A, B, C, D\}, \{E\}\}$$

$$\Pi_1 = \{\{A, B, C\}, \{D\}, \{E\}\}$$

$$\Pi_2 = \{\{A, C\}, \{B\}, \{D\}, \{E\}\}$$

$$\Pi_3 = \Pi_2 = \Pi_{\text{final}}$$

实际上这就是寻找闭包的过程. 因此原来最少可以有 4 个状态, 将等价的状态合并即可. 我们在合并的时候, 会不会让它成了一个 NFA? 简而言之, 实际上是不会的. 如果有的话, 说明分裂不彻底, 还可以再分裂.

初始状态、接受状态是哪些? 原先的初始状态就是初始状态, 原来的结束状态就是结束状态.

注:

1. 注意考虑把“死状态”单独分为一类.
2. 这是 DFA 最小化算法. 不适用于 NFA 最小化; NFA 最小化问题是 PSPACE-complete 的 (复杂度很高).

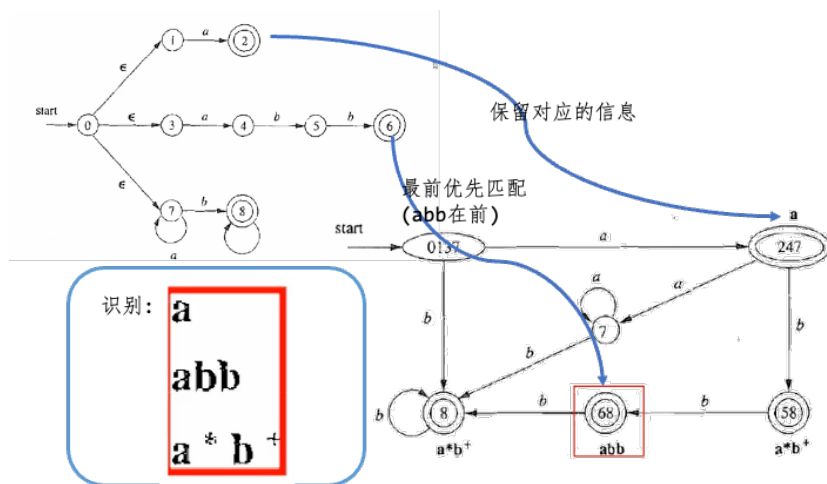
补充资料: 如何分析 DFA 最小化算法的复杂度? 为什么 DFA 最小化算法是正确的? 最小化 DFA 是唯一的吗?

定义 1.15 (可区分的 (Distinguishable); 等价的 (Equivalent)). 如果存在某个能区分状态 s 与 t 的字符串, 则称 s 与 t 是可区分的; 否则, 称 s 与 t 是等价的.

定义 1.16 (字符串 x 区分了 s 和 t). 如果分别从 s 与 t 出发, 沿着标号为 x 的路径到达的两个状态中只有一个是接受状态, 则称 x 区分了状态 s 与 t .

然后我们得到了最小的 DFA 就可以写词法分析器了.

d) **构造词法分析器** 记得我们匹配的原则: 最前优先匹配: *abb* (比如关键字), 以及最长优先匹配: *aabbbb*. 比如我们要构造下图的词法分析器, 就要保留各个 NFA 的接受状态信息 (识别出来了个什么), 并采用最前优先匹配 (两个接受状态合并的时候, 接受状态信息要是靠前的那一个).

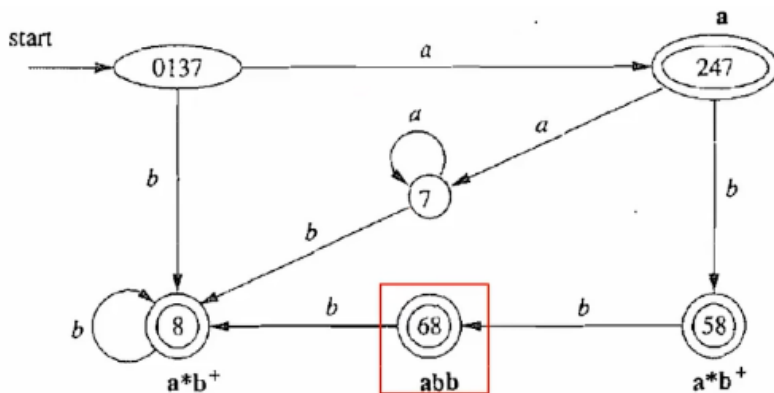


另外要注意, 需要消除“死状态”, 避免词法分析器徒劳消耗输入流. 这就要保证 $\forall a \in \Sigma \cdot \delta(\emptyset, a) = \emptyset$.

做了这样的准备, 我们就可以模拟运行该 DFA, 直到无法继续为止 (输入结束或状态无转移). 假设此时状态为 s .

1. 若 s 为接受状态, 则识别成功;
2. 否则, 回溯 (包括状态与输入流) 至最近一次经过的接受状态, 识别成功;
3. 若没有经过任何接受状态, 则报错 (忽略第一个字符), 这个字符程序不被允许.

例子 1.6. 对于下图的 DFA:



- $x = a$, 输入结束; 接受; 识别出 a
- $x = abba$ 状态无转移; 回溯成功; 识别出 abb ;
- $x = aaaa$ 输入结束; 回溯成功; 识别出 a
- $x = cabb$ 状态无转移; 回溯失败; 报错 c

如果希望对于 DFA 最小化, 需要做一些小改动. 初始划分需要考虑不同的词法单元. 一开始我们只是分为了接受状态, 拒绝状态, (以及死状态); 初始划分需要考虑不同的词法单元. 也就是上述例子中一开始我们应该这样分组: $\Pi_0 = \{\{0137, 7\}, \{247\}, \{8, 58\}, \{68\}, \{\emptyset\}\}$.