

# 扫描/词法分析

词法分析 (lexical analysis) 又称扫描 (scanning)

## 1、scanner的功能✓

词法分析程序 (scanner) 从源代码中读取字符形成若干个记号 (token)。

## scanner的任务

从源程序中读取字符并将其分成一个个的逻辑单元token, ~~每个~~ token是源程序中表示信息单元的字符序列。

## scanner的输出形式

词法分析程序的输出形式是一个二元式: (token类型, token的属性值)。

## 2、token的分类✓

### 关键字 (keywords)

例如 if 和 while, 它们是固定的字母串

### 标识符 (identifiers)

由用户定义的串, 通常由字母和数字组成且由字母开头

### 特殊符号 (special symbols)

如算术符号 + 和 \*; 一些多字符符号, >=

## token

token有若干种, 在程序语言中通常定义为枚举类型。比如在C语言中:

```
1 typedef enum{
2     IF, THEN, ELSE, PLUS, MINUS, NUM, ID, ...
3 } TokenType;
```

### 保留字 (reserved words)

比如 IF、THEN, 代表了 if 和 then

### 特殊符号 (special words)

比如算术符号 PLUS 和 MINUS, 代表 + 和 -

### 可表示多个字符串的token

token { keywords  
identifiers  
special symbols

can represent multiple strings, 比如 `NUM` 和 `ID`, 分别代表数字和标识符 (identifiers)

## token的attribute

### token的string value

token必须和它所代表的字符串区分开来, 比如 `IF` 必须和它代表的 `if` 区分开。

为了使这个区别更明显, 由token代表的字符串有时被称作串值 (string value) 或它的词义 (lexeme)。

某些token只有一个lexeme, 比如reserved words;

有些token有多个词义: 比如identifiers, 多个字符串 (根据用户定义得出) 都是 `ID`, 编译器会把这些标识符都记录在symbol table中。

### token的attribute

任何与token相关联的值都被叫做token的attribute, string value就是attribute的一种。

token还可以有其他的attribute:

比如 `NUM`, 它除了串值 `32767`, 还有一个由其串值得到的数字值 `32767`

比如 `PLUS`, 它除了串值 `+`, 还有一个真正的算术操作 `+`

实际上, 可以将token看作其attributes的结构体。

### scanner计算attribute

scanner必须计算每一个token的若干必要的attribute。

比如 `NUM` 的串值比如被计算, 而其数字值不必要被计算, 因为其可由串值计算得到; 另一方面, 如果其数字值被计算出, 那其串值可能被丢弃。

## token record

既然需要计算token的attribute, 那将其attribute收集起来则是十分有用的, 这种数据结构叫做token record, 其在C语言中定义可以如下:

```
1 typedef struct{
2     TokenType tokenval;
3     char* stringval;
4     int numval;
5 } TokenRecord;
```

## scanner的工作模式

scanner很少一次性将整个源程序转换为多一个个token。

实际上, scanner是在分析程序 (parser, 语法分析) 的控制下进行操作, 它通过一个函数从输入中返回有关命令的下一个token, 该函数有如下声明:

```
1 TokenType getToken(void);
```

该函数没有参数, 参数保存在缓冲区或由系统输入设备提供。



# 3. 正则表达式 (regular expressions)

## 定义

正则表达式表示字符串的格式，用来匹配相应的字符串，它完全由其所匹配的字符串来定义，一般用 $r$ 来表示

$L(r)$

$L(r)$ :  $r$ 匹配的字符串的集合，称为正则表达式 $r$ 生成的语言 (language generated by the regular expression)

该语言与程序设计语言无关 (至少在此是这样)，其依赖于适用的字符集 (character set)，通常是ASCII字符的集合或其子集。

$\Sigma$

$L(r)$ 的合法字符集称为字母表 (alphabet)，记作 $\Sigma$  (Sigma)，其中的元素称为符号 (symbol)

## patterns

在正则表达式中，所有的symbol都象征 (indicate) 着patterns.

## 元字符

元字符 (metacharacters or metasymbols) 有着特殊的含义，它可能不是 $\Sigma$ 中的字符，或者我们区分不出它是作为元字符适用还是作为 $\Sigma$ 中的普通字符适用 (为了区分其作用，可以用转义字符 (escape character) 实现)

## 基本正则表达式

基本正则表达式是字母表中的单个字符且与自身匹配

- 若 $a \in \Sigma$ ，则 $L(a) = \{a\}$
- 空串 $\epsilon$ ， $L(\epsilon) = \{\epsilon\}$
- 空集 $\emptyset$ ， $L(\emptyset) = \{\}$

## 正则表达式的三个基本运算及优先级

三个基本运算的优先级顺序为：\* 优先级最高，连接 其次，| 优先级最低，可以用 $()$ 改变这个优先级顺序

在该定义中， $\emptyset$ 、 $\epsilon$ 、 $|$ 、 $*$ 、 $()$ 这6个符号都有元字符的含义。

- 选择 (choice) :  $r|s$ 
  - $L(r|s) = L(r) \cup L(s)$ ，例如 $L(a|b) = L(a) \cup L(b) = \{a, b\}$
- 连接 (concatenation) :  $rs$ 
  - $L(rs) = L(r)L(s)$ ，例如 $L((a|b)c) = L(a|b)L(r) = \{a, b\}\{c\} = \{ab, bc\}$
- 重复 (repetition or closure) :  $r^*$ 
  - 有时称为Kleene闭包、克林闭包 (Kleene closure)，写作 $r^*$ ， $r$ 的0次或多次连接
  - 例如 $S^* = \{\epsilon\} \cup S \cup SS \cup SSS \cup \dots$

基本正则表达式

$a$
$\epsilon$
$\emptyset$

重复 repetition :  $r^*$   
连接 concatenation :  $rs$   
选择 choice :  $r|s$

$$\circ L(r^*) = L(r)^*$$

## 正则表达式的名字

我们可以给一个较长的正则表达式一个名字，以便使用

例如:  $digit = 0|1|2|\dots|9$

## 正则表达式的两个事实

- ① 不同的正则表达式可以生成相同的语言 (在实际中从未尝试着证实已找到了最简单的, 这有两个原因)
- ② 并非用简单术语描述的所有串都可由正则表达式产生, 我们将正则表达式可匹配的串的集合称为正则集合 (regular set)

## 正则表达式的扩展 (sf)

- 一个或多个重复:  $r^+$
- 任意字符:  $.$
- 字符范围:  $[a-z], [0-9], [a-zA-Z]$
- 非:  $\sim r$
- 可选的子表达式 (即0个或1个重复):  $r?$

## 最长子串定理

当串可以是单个token也可以是多个token的序列时, 则通常解释为单个token.

principle of longest substring: 可组成单个token的字符的最长串在任何时候都是假设为下一个记号, 待考究

# 有限自动机 (finite automata)

它是对由正则表达式给出的串格式的识别算法

## 几个概念

- 状态 (state)
  - 在图中用圆圈表示
- 转换 (transition)
  - 在图中有带有箭头的线表示由一个状态向另一个状态的转换
- 开始状态 (start state)
  - 识别过程开始的状态, 开始状态表示为一个不来自任何地方且指向它的箭头线无标识
- 接受状态 (accepting state)
  - 接受状态指匹配成功
- 出错转换 (error transition)
  - 来自出错状态的所有转换都要回到其本身
  - 出错状态是非接受的, 因此一旦发生一个出错, 则无法从这个出错状态逃出, 而且再也不能接受串了

## 确定有限自动机 (deterministic finite automata)



# 4. DFA定义和组成(部分)

DFA: 下一个状态由当前状态和当前输入字符唯一确定的自动机

一个确定有限自动机M由5部分组成:

- 字母表 $\Sigma$
- 状态集合 $S$
- 转换函数 $T: S \times \Sigma \rightarrow S$ 
  - $S \times \Sigma$ 指的是 $S$ 和 $\Sigma$ 的笛卡尔积或叉积: 集合对 $(s, c)$ , 其中 $s \in S, c \in \Sigma$ .
  - 例如:  $T(s, c) = s'$
- 初始状态 $s_0$ 
  - $s_0 \in S$
- 接受状态集合 $A$ 
  - $A \subset S$

## ε-转换

ε-transition 是无需考虑输入串就有可能发生的转换, 它可以看作一个空串的“匹配”。

有两个好处:  
①方便设计DFA  
②描述空串的匹配

- 使最初的确定性有限自动机保持完整并且只添加一个新的初始状态就可以将几个确定性有限自动机合并
- 有了它就可以清晰地描述出空串的匹配

## 非确定性有限自动机 (nondeterministic finite automata)

NFA和DFA有两点不同:  
①字母表 $\Sigma$   
②转移函数 $T$

- 字母表 $\Sigma$ 扩充为 $\Sigma \cup \epsilon$
- 转移函数 $T: S \times (\Sigma \cup \{\epsilon\}) \rightarrow P(S)$ , 其中 $P(S)$ 为 $S$ 的幂集 (power set)

任意个 $\epsilon$ 都可能在任一状态上引入到串中, 并与NFA中 $\epsilon$ 转换的数量相对应。因此, NFA并不表示算法, 但是却可以通过在每个非确定性选择中回溯的算法来模拟它

## 从正则表达式到DFA

将正则表达式翻译成DFA可以这样做:

1. 将正则表达式翻译成NFA
2. 将NFA翻译成DFA

选择重复者只会增加两个状态, 并不会增加状态, 只是首尾连接原NFA的状态不能少

## 从正则表达式到NFA

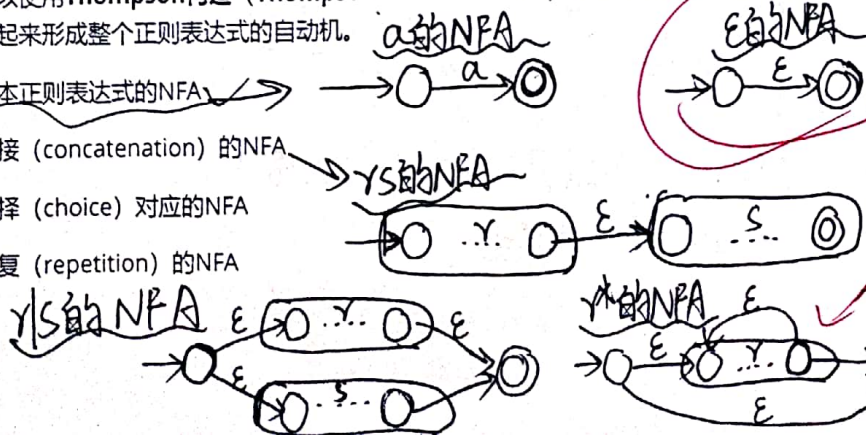
可以使用Thompson构造 (Thompson's construction), 其利用 $\epsilon$ 转换将正则表达式对应的自动机连接起来形成整个正则表达式的自动机。

基本正则表达式的NFA

连接 (concatenation) 的NFA

选择 (choice) 对应的NFA

重复 (repetition) 的NFA



记得这个形状



## 从NFA到DFA

### 状态 $s$ 的 $\epsilon$ -闭包 ( $\epsilon$ -closure)

我们将单个状态 $s$ 的 $\epsilon$ -closure定义为: 状态 $s$ 可经过0个或多个 $\epsilon$ -closure达到的状态的集合, 并将这个集合记作 $\bar{s}$ 。

一个状态的 $\epsilon$ 总是包含着该状态本身

### 状态集合的 $\epsilon$ -closure

我们将一个状态集合的 $\epsilon$ -closure定义为集合中每个状态的 $\epsilon$ -closure的并。

记状态集合为 $S$ , 则 $\bar{S} = \bigcup_{s \in S} \bar{s}$

### 子集构造 (subset construction)

设有NFA  $M$ , 则其DFA为 $\bar{M}$ 。

①  $\bar{M}$ 开始状态的闭包即 $\bar{M}$ 的初始状态

② 求每个状态在每字符 $a$ 上的转换, 即求 $Sa \rightarrow \bar{S}a$

解法 A. 求 $Sa$ 。若 $S$ 中有状态 $t$ 的转换( $t \xrightarrow{a} u$ ), 则 $u \in Sa$ 。

解法 B. 求 $\bar{S}a$ 。求 $\bar{S}$ 的 $\epsilon$ -闭包。

状态转换

$\bar{M}$ 中包含 $M$ 接受状态的状态为 $\bar{M}$ 的接受状态。

注意检查集合间的重复  
一直计算新集合

## 9. DFA状态数最小化

定理 对于任意的DFA, 都有一个含有最少量的状态的等价DFA, 且这个最少量状态的DFA是唯一的。

状态在同一集合内代表着等价

方法:

① 将DFA状态分为两个集合: 非接受状态集合和非接受状态集合

② 对于字符集上的每个字符 $a$

检查每个集合中是否有 $a$ 区分状态 $s$ 和 $t$ , 若有, 则拆分, 并重复步骤2

若无, 则定义了该集合到其自身的 $a$ 转换

区分的定义: 若 $a$ 区分 $s$ 和 $t$ ,

即 $s$ 和 $t$ 在 $a$ 上有不同的转换, 有2种情况。

① 转换至不同集合

② 存在错误转换和空转换



# 1. 词法分析的功能任务

## 2. token的分类(3类)

## 正则表达式(regular expression)

作用

正则表达式生成的语言

字母表(alphabet)  $\Sigma$  (sigma)

符号(symbol)

元字符(metacharacter/metasybol) "" \

基础正则表达式(3个)  $a \in \Sigma$

正则表达式的名字

正则表达式的两个事实

正则表达式的扩展: 5个:  $\gamma^+$ ,  $[a-z]$ ,  $\sim \gamma$ ,  $\gamma?$

最长子串处理  $\rightarrow$  歧义性

odd 奇数个

even 偶数个

3个基本运算  $\gamma s$   $\gamma | s$   $\gamma^*$   
及优先级:  $\gamma^* > \gamma s > \gamma | s$

## 4. DFA定义和组成 $\rightarrow$ (部分) 初始状态, 接受状态, 状态集, 转移函数, 字母表

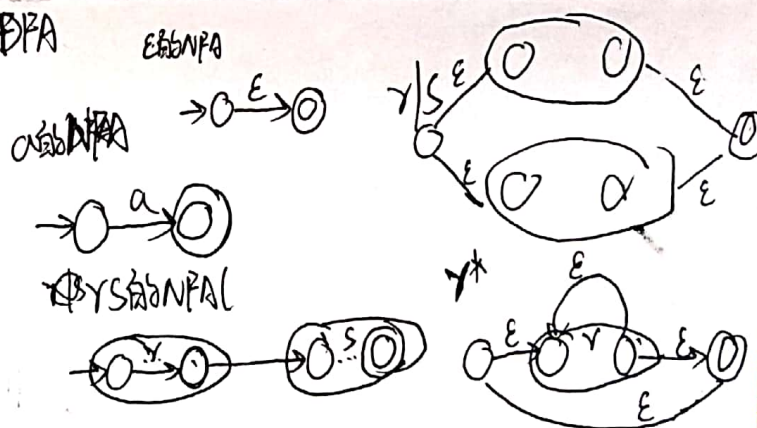
## 5. $\epsilon$ 转换及其好处: 简化匹配, 合并DFA

## 6. NFA与DFA的两点不同

## 从正则表达式到NFA

基础正则表达式的NFA

连接, 选择, 重复的NFA



## NFA $\rightarrow$ DFA

状态  $Q$  的  $\epsilon$  闭包  $\bar{Q}$ : 包括  $Q$

状态集  $S$  的  $\epsilon$  闭包  $\bar{S}$ : 所有  $\epsilon$  闭包的并

子集构造法

①  $M$

② ① 求  $S_a$

② ② 求  $\bar{S}_a$

2步: 1. 遍历状态集和字符  
2. 求开始状态的  $\epsilon$  闭包作为初始状态:  $\bar{A}$   
可结合在一起  
得  $S \xrightarrow{a} \bar{S}_a$  bar表示闭包

① 求  $M$  的初始状态

② 求每个状态集在每个字符  $a$  上的转换

即  $S \xrightarrow{a} \bar{S}_a$

A. 求  $S_a$

B. 求  $\bar{S}_a$

## DFA 状态数最小化 (即状态等价, step 1, step 2)

定理: 唯一性

方法: ① 先分成可接受状态集和不可接受状态集

② 遍历状态集和字符寻找区分: 若无区分, 则定义在  $a$  上的转换