

CH2 词法分析

2.1 实现方式

- 1. 先由自然语言规则通过人工描述将其变成 REX
  - 2. 把 REX 使用 Thompson 算法以自动化的形式转为 NFA
  - 3. 通过子集构造法（也是自动化的形式）转成 DFA
  - 4. 通过 Hopcroft 最小化算法简化 DFA
  - 5. 组合不同 token 类型的 DFA, 再通过一个二维数组判断状态转换 (state [0] 是判定为不接收的状态), 一个一维 finality 数组实现终止状态类型判断实现 DFA 的自动执行来(不是很重要)
- 生成词法分析器的代码

2.2 正则表达式

优先级

闭包>连接>选择

简写方式

- [abcd] 代表 a|b|c|d
- [b-g] 代表 [bcdefg]
- [b-gm-Qkr] 代表 [bcdefgmnOPQkr]
- e 可以被省略
- M? 代表一个或没有
- M+ 代表至少一个
- . 代表任意一个 character (除newline)
- "" 代表引用, 只能一个一个字符
- 精准匹配引用里面的内容

REX命名

通过  $r_1 \rightarrow r_2$  的方式给  $r_1$  命名成  $d_1$ ,

并且之后可以将  $d_1$  视为一个 character 一样

在 REX 中使用, 比如:

```
1 digit -> [0-9]
2 number -> digit+
```

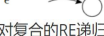
歧义处理

\*如果所有语句都只对应一个语法树, 那么这个文法没有歧义\*

- 1. longest match 最长匹配
  - 实现方式: 通过 last final 序列 (细节见词法分析 2.3.1)
- 2. rule priority 规则优先 - 根据规则和规则之间的优先级关系, 越前面越高

2.3 REX 转 NFA

- 1. 单个 NFA 的构造
- 对基本的 RE 逆递归构造:  $\epsilon, a$



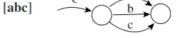
对复合的 RE 递归构造:  $st, s|t, s^*$



对简化表示方式的转换:

$M^+$  constructed as  $M \cdot M^*$

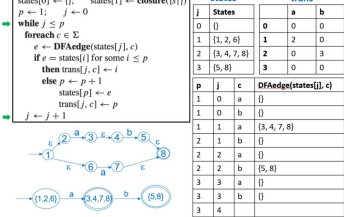
$M^?$  constructed as  $M | \epsilon$



"abc" constructed as  $a \cdot b \cdot c$

2.4 NFA 转 DFA

- 1. 首先列出所有状态的闭包;
- 2. 将初始状态的闭包作为新的初始状态;
- 3. 计算在每个新状态下在各个字符上的转移的闭包作为新的状态, 转移自然成为新的转移;
- 4. 包含原接受状态的所有新状态都是接受状态。



2.5 DFA 最小化

- 1. 划分等价类
- 1. 初始划分: 接受状态组和非接受状态组  $P = \{S - F, F\}$
- 2. 利用“可区分”的概念, 反复分割划分中的组 G

for P 中的每个元素/集合 G {  
  细分 G, 使得 G 中的 s, t 仍然在同一组中 iff  
  对任意输入 a, s, t 都到达 P 中的同一组;  
   $P_{new}$  = 将 P 中的 G 替换为细分得到的小组;  
}

集合 G 中的每个状态读入同一字符后, 都落入相同的某个集合, 那么就不用细分 S

- 3. 直到不可再分裂 (如果  $P_{new} = P$ , 令  $P_{final} = P$ , 算法完成; 否则  $P = P_{new}$ , 转步骤 2)

- 2. 重建 DFA: 从划分得到的等价类中选取代表, 并重建 DFA

第二章例题: 自然语言到正则语言的转换

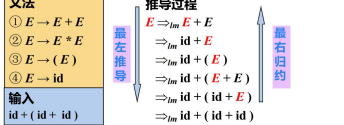
- a. Strings over the alphabet  $\{a, b, c\}$  where the first a precedes the first b.
- b. Strings over the alphabet  $\{a, b, c\}$  with an even number of a's.
- c. Binary numbers that are multiples of four.
- d. Binary numbers that are greater than 101001.
- e. Strings over the alphabet  $\{a, b, c\}$  that don't contain the contiguous substring  $baa$ .
- f. The language of nonnegative integer constants in C, where numbers beginning with 0 are octal constants and other numbers are decimal constants.
- g. Binary numbers  $n$  such that there exists an integer solution of  $a^2 + b^2 = c^2$ .
- a.  $c^*a(a|c)^*b(a|b|c)^*$
- b.  $(b|c)^*(b|c)a(b|c)^*a(b|c)^*$
- c.  $1(0|1)^*00$
- d.  $(1(0|1)^*(0|1)(0|1)(0|1)(0|1)(0|1)) | (1011(0|1)(0|1)) | (11(0|1)(0|1)(0|1)(0|1)(0|1))$
- e.  $(a|c)^*(b|ba)^*(c(a|c)^*)^*$
- f.  $(0|1-7|0-7)^* | (1-9|0-9)^* | 00 | 0$
- g.  $1|10$

CH3 语法分析

3.1 定义

EOF 的解决: 添加一个新的开始符号 S' 和一个新的规则  $S' \rightarrow SS$

推导和规约:



句型、句子和语言:

- 句型: 起始符号推导出来的任意串, 既可以包含终结符, 又可以包含非终结符, 也可能是空串
- 句子: 仅含终结符的句型
- 语言: 由文法推导出的所有句子构成的集合

编程语言的文法

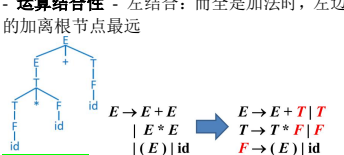
编程语言的文法由于为了语法分析的高效而添加了一定的处理和限制, 是任意 CFG 的子集。

消除二义性

“给定 CFG 是否无二义性”是不可判定问题, 不存在通用的算法。

通常用于解决这个问题的方式是分层:

- 运算优先级 - 根据算符不同的优先级, 引入新的非终结符; 越接近开始符号 S 的文法符号优先级越低
- 运算结合性 - 左结合: 而全是加法时, 左边的加高根节点最近



消除左递归  
 $A \rightarrow A\alpha | \beta$ , 其中  $\alpha \neq \epsilon, \alpha, \beta$  不以 A 开头

$A \rightarrow \beta A'$   
 $A' \rightarrow \alpha A' | \epsilon$

把左递归转成了右递归

提左公因子

对形如  $P \rightarrow \alpha a | \alpha y$  的一对产生式, 用如下产生式替换

$P \rightarrow \alpha Q$   
 $Q \rightarrow \beta | y$

Q 为新增加的未出现过的非终结符

3.2 Top-Down - LL(1)

采用最左推导 (每次优先替换最左边的非终结符), 代表文法是 LL(1) 文法。

3.2.1 First、Follow 和 Nullable

Nullable 集 - 一个文法对应单个集合, 定义是集合内的每一个元素可以生成空串

定义:

- Base case:  $X \rightarrow \epsilon$
- Inductive case:  $X \rightarrow Y_1 \dots Y_n$

if  $Y_1 \dots Y_n$  是 n 个非终结符且都可以生成空串

算法: 根据定义递归生成寻找, 结束条件是集合不再增长。

First 集 - 一个符号 (终结符和非终结符) 对应单个集合, 定义是可以由这个符号推导得到的串的首终结符号的集合

Base case:

- If X is a terminal: First(X) = {X}

Inductive case:

- If  $X \rightarrow Y_1 Y_2 \dots Y_n$ 
  - First(X) U= First( $Y_1$ ) 等价于 First(X)  $\supseteq$  First( $Y_1$ )
  - If  $Y_1 \in \text{Nullable}$ , First(X) U= First( $Y_2$ )
  - If  $Y_1, Y_2 \in \text{Nullable}$ , First(X) U= First( $Y_3$ )
  - ...

上述规则似乎是关于非终结符的, 但是 First 是关于文法符号串 (如产生式右部) 的, 规则如 inductive case。

对于非终结符对应单个集合, 定义是从起始符号出发, 可能在推导过程中跟在这个符号右边的终结符号的集合

Base case:

- Follow(A) = {}

Inductive case:

- If  $B \rightarrow s_1 A s_2$  for any  $s_1$  and  $s_2$ 
  - Follow(A) U= First( $s_2$ )
  - If  $s_2$  is Nullable, Follow(A) U= Follow(B)

对于一些特殊产生式的 follow 集合推论:

Production Constraints

$T' \rightarrow TS$  {S}  $\subseteq \text{FOLLOW}(T)$

$T \rightarrow R$  FOLLOW(T)  $\subseteq \text{FOLLOW}(R)$

$T \rightarrow aTc$  {c}  $\subseteq \text{FOLLOW}(T)$

$R \rightarrow \epsilon$

$R \rightarrow RbR$  {b}  $\subseteq \text{FOLLOW}(R)$

LL(1) 算法的预测分析表

行 A: 对应一个非终结符

列 a: 对应某个终结符或输入结束符 \$

项 M(A,a): 针对非终结符 A, 当下一个输入 Token 为 a 时, 可选的产生式

对文法 G 的每个产生式  $X \rightarrow y$

• if  $t \in \text{First}(y)$ : enter  $(X \rightarrow y)$  in row X, col t

• if y is Nullable and t  $\in \text{Follow}(X)$ : enter  $(X \rightarrow y)$  in row X, col t

如果表格内有冲突, 则代表其不是 LL(1) 文法, 我们需要通过消除左递归和提左公因子将其转换为 LL(1) 文法。

而保证产生式唯一性的条件如下 (对任何两个产生式  $A \rightarrow \alpha | \beta$ ):

- 1. First( $\alpha$ )  $\cap$  First( $\beta$ ) =  $\emptyset$

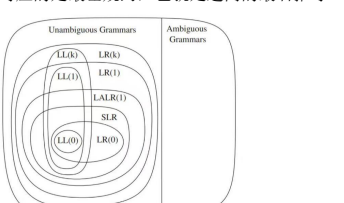
$\alpha$  和  $\beta$  推导不出以同一个单词为首的串

- 2. 若  $\beta \neq \epsilon$ , 那么  $\alpha \neq \epsilon$ , 且 First( $\alpha$ )  $\cap$  Follow( $A$ ) =  $\emptyset$

$\alpha$  和  $\beta$  不能同时推出  $\epsilon$ : First( $\alpha$ ) 不应在 Follow(A) 中

3.3 Bottom-Up - LR(0), SLR(1), LR(1) 和 LALR

对应的是最左规约, 也就是逆向的最右推导



关注三条线:

- LR(0)  $\subset$  SLR(1)  $\subset$  LALR(1)  $\subset$  LR(1)  $\subset$  LR(k)

- LR(0)  $\subset$  LL(1)  $\subset$  LL(k)

- $\forall k, LL(k) \subset LR(k)$ , 但是  $LL(1) \not\subset LALR(1)$

3.3.1 LR(0)

状态项

产生式加一个  $\cdot$  记录当前识别进度, 当  $\cdot$  在产生式最后的时候是可以使用这个产生式规约的起始状态

- 初始状态为  $S' \rightarrow \cdot S$

- 终止状态为  $S' \rightarrow S \cdot$

项集闭包

两个循环, 先遍历闭包中每一个 item, 再遍历每一个以 item 的  $\cdot$  后面的 symbol 为起始项的转换关系; 直至不再增长 (不动点思想)。

Closure(I) = repeat

for any item  $A \rightarrow \alpha \cdot \beta$  in I

for any production  $X \rightarrow \gamma$

if  $\alpha \cdot \beta$  is a prefix of  $X$

add  $A \rightarrow \alpha X \cdot \beta$  to I

until I does not change

return I

转换关系

goto 函数最开始是空集, 对于闭包里面每一个项, 如果后面是 X, 则将其加入 goto 的结果, 并且计算其 closure。

Goto(I, X) =

set J to the empty set

for any item  $A \rightarrow \alpha \cdot X \beta$  in I

add  $A \rightarrow \alpha X \cdot \beta$  to J

return Closure(J)

整体 DFA 构造算法

对每一个状态里面的每一项计算 goto 的结果并加入 DFA 直至不再增长。

Initialize T to {Closure(S'  $\rightarrow$   $\cdot$  SS)}

Initialize E to empty.

repeat

for each state I in T

for each item A  $\rightarrow \alpha \cdot X \beta$  in I

let J be Goto(I, X)

T  $\leftarrow T \cup \{J\}$

E  $\leftarrow E \cup \{J\}$

until E and T did not change

不断增加新状态直到不动点

语法分析表

意义:

以状态为行, 以 Action 的终结符和 GOTO 的非终结符为列, 表项内容为操作:

- Action 表项:

- shift+goto 下一状态: 例如 s2 就是 shift 并且去到第二个状态

- 当状态存在以  $\cdot$  为结尾的项时 (有一个项满足就整体都选择规约), 选择 index 为 k 的产生式进行 reduce

- GOTO 表项:

- 列是每次规约所用产生式左边的 nonterminal A, 表示每次规约之后去到的下一个状态是什么

通过 DFA 构造:

- 1. s 和 g 都由状态转移关系决定, 如果是非终结符就是 g, 终结符就是 s, 然后当前状态对应行, 转移读的符号是列, 去到的状态是 g/s 后面跟的数字;

- 2. r 就是如果这个状态包括在最后的产生式, 那么使用这个产生式规约, 这一整行的 Action 都加上 r+产生式对应的序号

这一项。

例子:

0: S'  $\rightarrow$  S\$

1: S  $\rightarrow$  x S

2: S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

S  $\rightarrow$  x S

S  $\rightarrow$  y

S'  $\rightarrow$  S\$

Stack (states)	Stack (symbols)	Input	Action
1		x y S	shift 1
1 2	x	x y S	shift 2
1 2 2	x x	y S	shift 3
1 2 2 3	x x y	S	reduce 2 (S $\rightarrow$ y)
1 2 2 5	x x S	S	reduce 1 (S $\rightarrow$ x S)
1 2 5	x S	S	reduce 1 (S $\rightarrow$ x S)
1 4	S	S	accept

3.3.2 SLR(1)

在 LR(0) 的基础上在规约条件中增加了 next token 是否在产生式左边的 follow 集合内的判断:

What are the valid next tokens for r2?

- Idea: we can choose the reduce action only if the next input token  $t \in \text{Follow}(E) = \{S\}$

- Only T|3, S| can be r2!

通过 DFA 构造:

- 1. s 和 g 都由状态转移关系决定, 如果是非终结符就是 g, 终结符就是 s, 然后当前状态对应行, 转移读的符号是列, 去到的状态是 g/s 后面跟的数字;

- 2. r 就是如果这个状态包括在最后的产生式, 那么使用这个产生式规约, 这一整行的 Action 都加上 r+产生式对应的序号

这一项。

例子:

S'  $\rightarrow$  E\$

E  $\rightarrow$  T+E

T  $\rightarrow$  x

S'  $\rightarrow$  E\$

E  $\rightarrow$  T+E

T  $\rightarrow$  x

S'  $\rightarrow$  E\$

E  $\rightarrow$  T+E

## AC