

CH2 词法分析

2.1 扫描处理

1.某些记号只有一个词义:保留字;某些记号有无限多个语义:标识符 ID 表示.

2.2 正则表达式

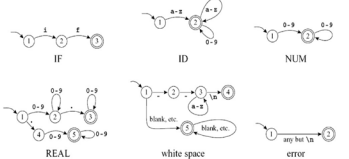
1.一些 notation

- a An ordinary character stands for itself.
- ε The empty string.
- Another way to write the empty string.
- M | N Alternation, choosing from M or N.
- M · N Concatenation, an M followed by an N.
- MN Another way to write concatenation.
- M* Repetition (zero or more times).
- M+ Repetition, one or more times.
- M? Optional, zero or one occurrence of M.
- [a - zA - Z] Character set alternation.
- . A period stands for any single character except newline.
- "a +*" Quotation, a string in quotes stands for itself literally.

2.RE 匹配优先匹配保留字;最长字符串优先

2.3 有穷自动机

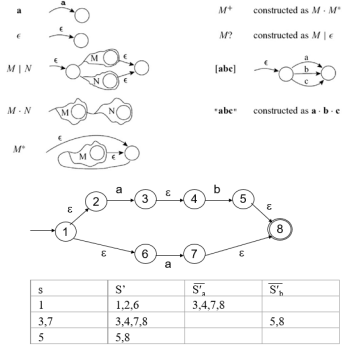
- 1.DFA:M 由字母表Σ、状态集 S、转换函数 T:S×Σ→S、初始状态 S₀∈S 以及接受状态 A⊂S.
- 2.错误状态默认不画,但是存在;错误状态下的任何转移均回到自身,永远无法进入接受.
- 3.NFA:M 由字母表Σ、状态集 S、转换函数 T:S×(Σ∪{ε})→P(S)、初始状态 S₀及接受状态 A 的集合.
- 一些给出的基本 FA 图:



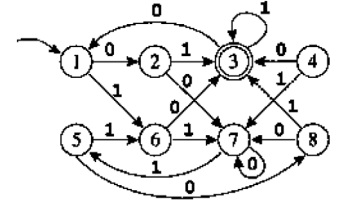
2.4 正则表达式到 DFA

1.子集构造的过程:首先列出所有状态的ε闭包;然后将初始状态的ε闭包作为新的初始状态;然后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;包含原接受状态的所有新状态都是接受状.

PS:ε闭包首先包含自身.



4.DFA 状态数最小化:一种正向思路:从 DFA 的终态出发逆着找,看出边、终点一致的就是等价态.以 2.6 为例:{4, 6} {2, 8} {1, 5}等价.



5. NFA-DFA 的步骤:给出一个表头如下:新状态编号 新加字符 新加旧状态 新状态代号;

CH3 上下文无关文法分析

3.1 CFG

1.左递归:定义 A 的推导式的右边第一个出现的是 A;右递归:定义 A 的推导式右边最后一个出现的是 A;

3.2 分析数和抽象语法树

- 1.二义性 (ambiguous):同一个串存在多个推导即多个分析树
- 2.分析树 (concrete syntax tree)是一个作了标记 labeled 的树,内部节点是非终结符,树叶是终结符;

对一个内部节点运用推导时,推导结果从左到右依次成为该内部节点的子节点
3.最左推导和前序编号对应,最右推导后序
4.AST(syntax tree)去除了终结符和非终结符信息,仅保留了语义信息;一般用左孩子右兄弟

3.3 Ambiguity 二义性

- 1.定义:带有两个不同的分析树的串的文法
- 2.解决方法①设置消歧规则 disambiguating rule,在每个二义性情况下指出哪个是对的.无需对文法进行修改,但是语法结构就不是单纯依赖文法了,还需要规则②修改文法.
- 4.修改文法时需要同时保证优先级和结合律 precedence and associativity
- 5.在语法树中,越接近根,越高,优先级越低;左递归导致左结合,右递归会右结合
- 6.相同优先级的运算符组叫 precedence cascade 优先级联
- 7.期中错题: A grammar is ambiguous if it has only one parse trees for all sentences.

Answer: T

自顶向下分析: LL(k)

第一个 L 是从左到右处理,第二个 L 是最左推导,1 代表仅使用 1 个符号预测分析方向
递归下降:1.将一个非终结符 A 的文法规则看作将识别 A 的一个过程的定义.递归下降需要使用 EBNF:将可选[]翻译成 if,将重复{}翻译成

while 循环

- LL(1)1.动作:①生成 (generate),利用文法将栈顶的 N 替换成串,串反向进栈
- ②匹配 (match):将栈顶的记号和下一个输入记号匹配
- ③接受 (accept):接受字符串
- ④错误 (error)
- 2. LL(1)文法是无二义性的,对任意规则 A → α1|α2, First(α1)∩First(α2) 为空,否则不是 LL(1).
- 3. LL(1)面对重复和选择的解决方法:消除左递归 left recursion removal 和提取左因子 left factoring.
- 4.简单直接左递归: A → Aα|β, αβ 是 N,且 β 不以 A 开头. A → βA', A' → αA'|ε
- 5.普遍直接左递归: A → Aα1|Aα2|...|Aαn|β1|β2|...|βm A → β1A'|β2A'|...|βmA' A → α1A'|α2A'|...|αnA'|ε
- 6.提取左因子: A → αβ|αγ. A → αA', A' → β|γ

first follow sets:
1.Fisrt 定义:令 X 为一个 T 或 N 或 ε, Fisrt(X) 由 T 或 ε 组成.①若 X 为 T 或 ε, Fisrt(X)={X}②若 X 为 N,对于每个产生式 X→X1 X2 ...Xn, First(X) 都包含了 First(X1) - {ε}.

若对于某个 i<n,所有的 Fisrt(X1),...,First(Xi) 都含有 ε,则 First(X) 也包括了 First(Xi+1) - {ε}.若所有 Fisrt(X1),...,First(Xn) 都含有 ε,则 First(X) 也包含 ε.

2.定理:A non-terminal A is nullable if and only if First(A)

contains ε
3.Follow 定义:若 A 是一个 N,那么 Follow(A) 由 T 和 \$ 组成.①若 A 是 \$,直接进入 Follow(A) ②若存在产生式 B→αAy,则 First(y) - {ε} 在 Follow(A) 中 ③若存在产生式 B→αAy,且 ε 在 First(y) 中,则 Follow(A) 包括 Follow(B)
PS:③更常见的情况是 B→αA,那么 Follow(A) 包括 Follow(B)
4.First 关注点在产生式左边,Follow 在右边
期中错题:Given production A→BαC, we have:Follow(A)⊂Follow(C), First(B)⊂First(A)
5.LL(1) 分析表 M[N,T] 的构造算法:为每个非终结符 A 和产生式 A→α 重复以下:
①对于 First(α) 中的每个记号 a,都将 A→α 添加到项目 M[A,a] ②若 α 可空 (nullable),则 Follow(A) 中的每个元素 a (包括 \$),都将 A→α 添加到 M[A,a]
例:下列文法构造 LL(1)

Table with 3 columns: nullable, FIRST, FOLLOW. Rows for X, Y, Z.

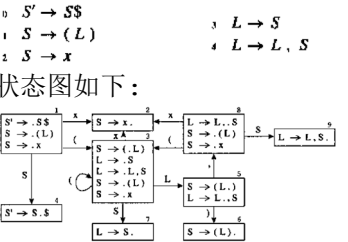
Table with 3 columns: a, c, d. Rows for X, Y, Z.

6.LL(1) 判别:A grammar in BNF is LL(1) if the

following conditions are satisfied. ① For every production Ai→ α1|α2|...|αn, First(αi)∩First(αj) is empty for all i and j, 1 ≤ i, j ≤ n, i ≠ j ② For every non-terminal A such that First(A) contains ε, First(A)∩Follow(A) empty.

自底向上分析: LR(k)

- Yacc 基于 LALR(1).
- 1.动作为①shift,将 T 从输入开头移到栈顶②reduce 使用产生式 A→α 将栈顶的 α 规约成 A③accept 分析栈为开始符号,输入栈为空时的动作④error ⑤转移 (goto):吃 non-terminal
- 2.加一个 S':新的开始符号
- 3.LR 分析器需要构造和使用一张 LR 分析表;而 LR 分析表的构造需要一个生成器.对于 start symbol 而言,需要新增 S'→S\$ 规则来引入.
- 4.LR(0) 的状态中包含带 ' . ' 的规则, ' . ' 左边表示已读,右边表示未读;如果圆点 ' . ' 右边的第一个非终结符有规则,那么也写到这一状态中.
- 例如:如下文法构造 LR(0)



由上面的状态图构造表如下:

	()	x	.	\$	S	L
1	s3		s2			g4	
2	r2	r2	r2	r2			
3	s3		s2			g7	g5
4							
5			s6		s8		
6	r1	r1	r1	r1	r1		
7	r3	r3	r3	r3	r3		
8	r4	r4	r4	r4	r4		
9						g9	

5.s-r conflict:表中同时出现了 sn 和 rn,即构成 s-r 冲突

6.r-r conflict:表中同时出现两个及以上的 rn,则构成 r-r

7.LR(0) 文法不可能二义的
8.A grammar is LR(0) if and only if ①Each state is a shift state (a state containing onlyshift items) or a reduce state (containing a single complete item).
这里的 r1 都是用 s'→s·规约,应该写成 accept

SLR:

1.SLR 算法定义:移进规则不变；规约时要求输入必须在属于 follow(A) 的终结符的项中,而不是全放.就是比 LR(0)改进
 $R \leftarrow \{ \}$
for each state I in T
 for each item $A \rightarrow \alpha .$ in I
 for each token X in FOLLOW(A)
 $R \leftarrow R \cup \{ (I, X, A \rightarrow \alpha) \}$

2.SLR 不可能是二义性
3.自底向上右递归可能引起栈溢出,需要避免

4.SLR 中的两种冲突,sr 冲突使用消歧规则:优先移进；rr 冲突基本是设计出问题
LR(1) and LALR(1)

1.LR(1) items:[$A \rightarrow \alpha \cdot \beta, a$] 前面是 LR(0) 项,后面是 lookahead token
2.LR(1) 的起始状态 $[S' \rightarrow \cdot S, \$]$ 的闭包
3.LR(1) 中如何获得规则状

态的闭包？难点在找 look ahead symbol.使用如下算法：
Closure(I) = repeat

for I 中任意项 ($A \rightarrow \alpha.X\beta,z$)
 for 任意产生式 $X \rightarrow y$
 for 任意 $w \in \text{First}(\beta z)$
 $I := I \cup \{ (X \rightarrow y, w) \}$
until I 没有变化

$S' \rightarrow \cdot S \& \quad ?$	$S' \rightarrow \cdot S \& \quad ?$
$S \rightarrow \cdot V = E \quad \$$	$S \rightarrow \cdot V = E \quad \$$
$S \rightarrow \cdot E \quad \$$	$S \rightarrow \cdot E \quad \$$
$E \rightarrow \cdot V \quad \$$	$E \rightarrow \cdot V \quad \$$
$V \rightarrow \cdot x \quad \$$	$V \rightarrow \cdot x \quad \$$
$V \rightarrow \cdot * E \quad \$$	$V \rightarrow \cdot * E \quad \$$
$V \rightarrow \cdot x \quad =$	$V \rightarrow \cdot x \quad =$
$V \rightarrow \cdot * E \quad =$	$V \rightarrow \cdot * E \quad =$

4.LR(1) 表的构造如下:在圆点到末尾的位置发生规约,对应项中填入 rn(n 是规则编号)；读入终结符则在对应项中写入 sn;读入非终结符则在对应项中写入 gn;读入 s 后在标记符位置填 a 表示 accept

5.LR(1) 文法不可能二义性

6. LR(1) 分析表示例

	x	w	=	S	T	V'
1	s8	s6		g2	g5	g3
2			a			
3			s4	r3		
4	s11	s13			g9	g7
5						
6	s8	s6			g10	g12
7			r4	r3		
8			r4	r3		
9			r1	r1		
10			r5	r5		
11			r5	r4		
12			r3	r3		
13	s11	s13			g14	g7
14			r5			

8.LALR(1)将 look ahead symbol 进行合并:除了 look ahead symbol 不同以外全都相同的规则合并成一个,目的在于减少表项。
9.A grammar is an LALR(1) grammar if no parsing conflicts arise in the LALR(1)

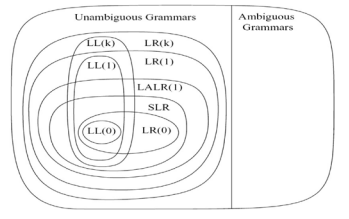
parsing algorithm.

10.如果文法是 LR(1),那么 LALR(1)中必然没有 sr 冲突,但是可能有 rr 冲突.

$[A \rightarrow \alpha \cdot a]$	$[A \rightarrow \alpha \cdot b]$	$[A \rightarrow \alpha \cdot a/b]$
$[B \rightarrow \alpha \cdot b]$	$[B \rightarrow \alpha \cdot a]$	$[B \rightarrow \alpha \cdot a/b]$

11.如果文法是 SLR(1),那么必然是 LALR(1)。

12. 各类文法的层次如下：



13.悬挂 else(dangling else)的处理:悬挂 else 会导致 s-r 冲突,冲突文法如下:

$S \rightarrow \text{if } E \text{ then } S \text{ else } S$
 $S \rightarrow \text{if } E \text{ then } S$
 $S \rightarrow \text{other}$

解决办法:引入新非终结符用于 if-else 的匹配问题.M 是匹配的 else;U 是未匹配的 else

$S \rightarrow M$
 $S \rightarrow U$
 $M \rightarrow \text{if } E \text{ then } M \text{ else } M$
 $M \rightarrow \text{other}$
 $U \rightarrow \text{if } E \text{ then } S$
 $U \rightarrow \text{if } E \text{ then } M \text{ else } U$

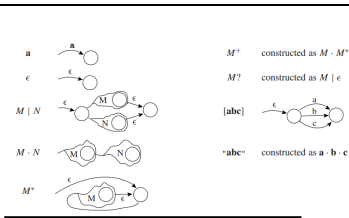


FIGURE 2.6. Translation of regular expressions to NFA.

0	$S' \rightarrow \cdot S$	5	$S \rightarrow \cdot S \& S$
1	$S \rightarrow \cdot$	6	$S \rightarrow \cdot (S)$
2	$S \rightarrow \cdot x$	7	$S \rightarrow \cdot \text{WORD}$
3	$S \rightarrow \cdot \backslash \text{begin} (\text{WORD})$	8	$S \rightarrow \cdot \backslash \text{begin}$
4	$S \rightarrow \cdot \backslash \text{end} (\text{WORD})$	9	$S \rightarrow \cdot \text{end}$
10	$S \rightarrow \cdot \backslash \text{WORD}$		

	nullable	FIRST	FOLLOW
S'	False	$\backslash \{ S, \text{WORD}, \text{begin}, \text{end} \}$	
S	True	$\backslash \{ \text{WORD}, \text{begin}, \text{end} \}$	S, \backslash
B	False	\backslash	$\backslash \{ \text{WORD}, \text{begin}, \text{end} \}$
E	False	\backslash	$\backslash \{ \text{WORD}, \text{begin}, \text{end}, \backslash, S \}$
X	False	$\backslash \{ \text{WORD}, \text{begin}, \text{end} \}$	$\backslash \{ \text{WORD}, \text{begin}, \text{end}, \backslash, S \}$

	\backslash	begin	end	WORD	{	}	S
S'	$S' \rightarrow \cdot S S$	$S' \rightarrow \cdot S$	$S' \rightarrow \cdot S$	$S' \rightarrow \cdot S S$	$S' \rightarrow \cdot S S$	$S' \rightarrow \cdot S$	$S' \rightarrow \cdot S$
S	$S \rightarrow \cdot \epsilon, S \rightarrow \cdot X S$	$S \rightarrow \cdot X S$	$S \rightarrow \cdot X$	$S \rightarrow \cdot X S$	$S \rightarrow \cdot X$	$S \rightarrow \cdot S$	$S \rightarrow \cdot \epsilon$
B	$B \rightarrow \cdot \backslash \text{begin} (\text{WORD})$						
E	$E \rightarrow \cdot \text{end} (\text{WORD})$						
X	$X \rightarrow \cdot B S E, X \rightarrow \cdot \backslash \text{WORD}$	$X \rightarrow \cdot \text{begin}$	$X \rightarrow \cdot \text{end}$	$X \rightarrow \cdot \text{WORD}$	$X \rightarrow \cdot \{$	$X \rightarrow \cdot \}$	

	nullable	FIRST	FOLLOW
S	False	ϵ	
B	False	w	v, y, x, x
D	True	y, x	x
E	True	y	x, x
F	True	x	x

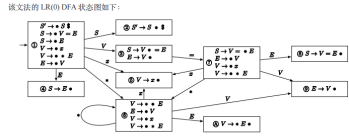
构造 LL(1) 分析表, 给出证据表明该文法不是 LL(1) 文法, 尽管少量地修改文法使他成为一个接收相同语言的 LL(1) 文法.

	nullable	FIRST	FOLLOW
S	False	ϵ	
B	False	w	v, y, x, x
D	True	y, x	x
E	True	y	x, x
F	True	x	x

	u	v	w	x	y	z
S	$S \rightarrow \cdot u B D x$					
B			$B \rightarrow \cdot B v, B \rightarrow w$			
D				$D \rightarrow \cdot \epsilon F$	$D \rightarrow \cdot \epsilon F$	$D \rightarrow \cdot \epsilon F$
E				$E \rightarrow \cdot \epsilon$	$E \rightarrow \cdot y$	$E \rightarrow \cdot \epsilon$
F				$F \rightarrow \cdot x$		$F \rightarrow \cdot \epsilon$

因为在上面的 LL(1) 分析表中有两个规则包含到同一个表项 (B, w) 中, 所以不是 LL(1) 文法, 修改目标是将 $B \rightarrow B v$ 与 $B \rightarrow w$ 分拆, 将这两个规则替换为 $B \rightarrow w B, B \rightarrow v B' \{ \epsilon \}$ 即可.

	x	w	=	S	T	V'
0	s1	s2		g2	g5	g3
1	s2	s2		g4	g5	g3
2	s2	s2		g4	g5	g3



构造 SLR 分析表还需要计算非终结符的 Follow 集: Follow(S) = { $\$$ } / Follow(V) = { $\epsilon, \$$ } / Follow(E) = { $\epsilon, \$$ }. 构造 SLR 分析表如下 (前四行 Action, 后三行 GOTO) :

	x	*	=	\$	S	V	E
1	s5	s6			g2	g3	g4
2					accept		
3			s7, r3	r3			
4				r2			
5		r4	r4				
6	s5	s6			g2	g4	
7	s5	s6			g2	g4	
8			r1				
9		r3	r3				
10		r5	r5				

因为单元格 (3, *) 中包含了 s7, r3 两个 action, 出现冲突.

构造下面文法的 LR(0) DFA, 并分析是否是 LR(0) / SLR(1) 文法, 给出证据.

0	$S \rightarrow \cdot S$	2	$S \rightarrow \cdot (S)$
1	$S \rightarrow \cdot$	3	$S \rightarrow \cdot \epsilon + (S)$



id	()	+	\$	S	E
1	s3				g2	
2			g4	accept		
3	r1	s5, r1	r1	r1	r1	
4	s6					
5	s3					g8
6	r3	r3	r3	r3	r3	
7			s8	g4		
8	r2	r2	r2	r2	r2	

因为状态 2 中同时存在 shift 和 accept, 并且 3 (单元格中 s5, r1 冲突, 所以不是 LR(0) 文法.
因为 Follow(E) = { $\$, +$ }, 所以在 SLR 分析表中 3 (单元格中不会填入 r1, 没有冲突, 而且 SLR(1) 可以 lookahead 一个 token, 所以状态 2 中可以通过 token 区分 g4 和 accept, 所以是 SLR 文法, 同时因为 LR(1) 优于 LR(0), 所以也是 LR(1) 文法.

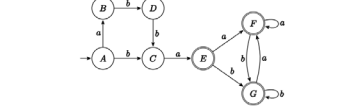
写出 (a+b)(a+b)* 的 NFA, 并写它的 regular DFA.



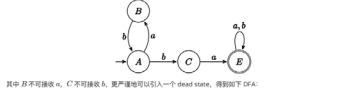
NFA 转化为 DFA
利用子集构造法, 从初始状态开始构造:
• $A = \epsilon\text{-closure}(0) = \{0, 1, 6\}$
• A 接收 a 转移到 $B = \epsilon\text{-closure}(2) = \{2\}$
• A 接收 b 转移到 $C = \epsilon\text{-closure}(5) = \{5\}$
• D 接收 b 转移到 $D = \epsilon\text{-closure}(3) = \{1, 3, 4\}$
• C 接收 a 转移到 $E = \epsilon\text{-closure}(6) = \{6, 7, 8, 10, 13\}$
• D 接收 a 转移到 $F = \epsilon\text{-closure}(9) = \{7, 8, 9, 10, 12, 13\}$
• F 接收 b 转移到 $G = \epsilon\text{-closure}(11) = \{7, 8, 10, 11, 12, 13\}$
• D 接收 c 转移到 G , 接收 b 转移到 C , F 和 G 接收 a 和 b 转移到 F 和 G

状态	接收 a	接收 b
A	B	C
B	/	D
C	E	/
D	B	C

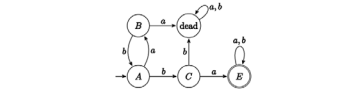
状态	接收 a	接收 b
E	F	G
F	F	G
G	F	G



最小化 DFA
初始划分分为 {A, B, C, D}, {E, F, G}, 根据上述划分转移可知 A, D 同属 E, F, G 同属不可区分, 最终的划分为 {A, D}, {B}, {C}, {E, F, G}, 每组内选一代表状态, 得到最小化 DFA 如下:



其中 B 不可接收, C 不可接收, 更严谨地可以引入一个 dead state, 得到如下 DFA:



10. _____ is commonly used to make operations left associative.
单选题 (5.0 分) (难易度:中)
A. Right recursion
B. Indirect recursion
C. Left recursion
D. Left factoring
正确答案: C
答案解析: 暂无

1. Given a legal string of tokens for a CFG, there must be a unique parsing tree to derivate the string.
判断题 (4.0 分) (难易度:中)
A. True
B. False
正确答案: B
答案解析: 暂无

2. Left-recursive CFG cannot be LL(1).
判断题 (4.0 分) (难易度:中)
A. True
B. False
正确答案: A
答案解析: 暂无

3. The language $L = \{ a^n b^m | n > 1 \}$ can't be generated by any CFG.
判断题 (4.0 分) (难易度:中)
A. True
B. False
正确答案: A
答案解析: 暂无

5. Both DFA and NFA can recognize regular set.
判断题 (4.0 分) (难易度:中)
A. True
B. False
正确答案: A
答案解析: 暂无

6. Which of the following string can be defined by the regular expression $(ab)^* ab^* (ab)^* c^*$
单选题 (5.0 分) (难易度:中)
A. Match
B. Generate
C. Accept
D. Reduce
正确答案: C
答案解析: 暂无

8. The output of the parser is _____.
单选题 (5.0 分) (难易度:中)
A. token
B. syntax tree
C. target code
D. intermediate code
正确答案: B
答案解析: 暂无

9. _____ reflects the derivation steps of a string.
单选题 (5.0 分) (难易度:中)
A. Syntax tree
B. Parse tree
C. Binary tree
D. Quadtree
正确答案: B
答案解析: 暂无

10. _____ is commonly used to make operations left associative.
单选题 (5.0 分) (难易度:中)
A. Right recursion
B. Indirect recursion
C. Left recursion
D. Left factoring
正确答案: C
答案解析: 暂无