

编译原理

2020年4月





编译原理 第四章

上海交通大学 张冬茉

Email:zhang-dm@cs.sjtu.edu.cn

2020年3月



括:

第4章 语法分析



本章将介绍编译器中使用的典型语法分析方法。主要包

- 介绍基本概念
- 讨论适合于手工实现的技术,
- 介绍自动生成工具中所使用的算法
- 扩展语法分析方法,使之能从常见的错误中得以恢复。



4.1 引论



上下文无关文法可以给出一个程序设计语言的精确易懂的语法描述,本节将通过一个算术表达式的典型文法,说明语法分析的本质,因为处理表达式的语法分析技术可以用于处理程序设计语言的大部分结构。



4.1.1 语法分析器的作用



在本书的编译器模型中, 语法分析器接受词法分析器提供的记号串, 检查它们是否能由源程序语言的文法产生, 如图4-1所示。我们希望词法分析器能用易于理解的方式提示语法错误信息, 并能从常见的错误中恢复过来, 以便后面的输入能继续处理下去。

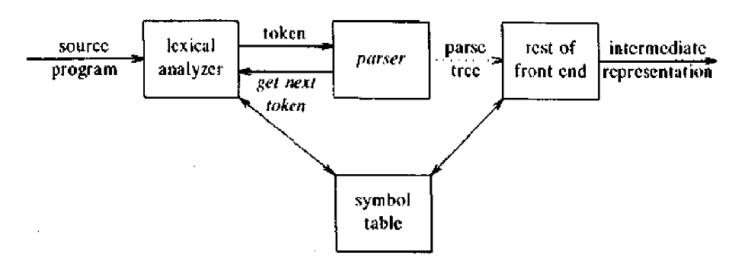
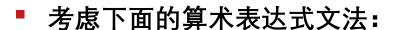


Fig. 4.1. Position of parser in compiler model,

- 典型的文法的语法分析器有三类
 - 通用的语法分析方法
 - ■自顶向下的语法分析方法
 - ■自底向上的语法分析方法
- 本节的剩余部分将考虑语法错误的性质和错误恢复的一般策略。我们在讨论各种语法分析方法时将详细介绍两种错误恢复策略:
 - 应急模式恢复策略
 - ■短语级恢复策略



4.1.2 代表性文法



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

$$T \rightarrow T * F | F$$

<mark>消除E和T的直接左递归</mark>,可以得到:

E -> TE'

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

• 考虑下面的算术表达式文法:

$$E \rightarrow E + E \mid E * E \mid (E) \mid id$$

(4.3)

(4.1)

(4.2)



4.1.3 语法错误的处理



- 我们知道程序可能包含不同级别的错误。下边是一些程序错误的示例:
 - 词法错误, 如标识符、关键字或操作符的拼写错误。
 - 语法错误,如算术表达式的括号不配对。
 - 语义错误,如操作符作用于不相容的操作数。
 - 逻辑错误,如无限的递归调用。
- 语法分析器中出错处理程序的基本目标是:
 - 清楚而准确地报告错误的出现。
 - 迅速地从每个错误中恢复过来,以便能继续检查后面的错误。
 - 不能过分降低正确程序的处理速度。





- 有些分析方法,如LL方法和LR方法,可以尽快地检查出语法错误。更准确地说,它们具有可行前缀性(viable-prefix property),即当它们一旦发现一个输入字符串的前缀不是该语言任何字符串的前缀时就能检测出错误。
- 60%的被编译程序语法和语义都是正确的。即使有错误出现,错误的 种类也相当少。80%的出错语句只含一个错误,只有13%的出错语句含有 两个错误。多数错误是微不足道的,90%的错误是单个记号错。
- 多数错误可简单地分为以下几类: 60%是标点符号错, 20%是操作符或运算对象错, 15%是关键字错, 剩下的5%是其他类型的错误。多数标点号错都属于分号的不正确使用。





出错处理程序应该怎样报告错误呢?很多编译器普遍采用的方法是:显示错处的程序行,用指针指出检测到错误的位置。如果能够知道实际错误可能是什么,编译器还会显示附带的诊断信息,如"此处遗漏了分号"等。

一但检测出错误,语法分析器将如何恢复这个错误呢?有很多一般性的策略,但没有哪种策略占绝对优势。



4.1.4 错误恢复策略



我们将主要介绍下面几种策略:

- 应急模式恢复策略。最容易实现。当发现错误时,语法分析器开始 抛弃输入记号,每次抛弃一个记号,直到发现某个指定的同步记号 为止。比较简单,不会陷入死循环。对一个语句中出现错误较少时 比较合适。
- 短语级恢复策略。发现错误时,语法分析器对剩余的输入字符串作局部纠正,即用一个能使语法分析器继续工作的字符串来代替剩余输入的前缀。可能引起死循环。被用于自顶向下语法分析中。主要缺点是难以应付实际错误出现在诊断点之前的情况。





- 出错产生式策略。如果对经常遇到的错误有很清楚地了解,我们可以扩充语言的文法,增加产生错误结构的产生式。然后用由这些错误产生式扩充的文法构造语法分析器。
- 全局纠正策略。我们总是希望一个理想的编译器在处理不正确的 输入字符串时尽可能少的改动。如果给定错误输入串×和文法G, 这 些算法会发现y的一棵分析树, 以便使用最少的符号插入、删除和修 改操作把×变换成正确的输入字符串y。不幸的是, 实现这些算法的时间和空间开销太大, 目前只作理论上的探讨。



4.2 上下文无关文法

在第二章我们已经介绍了上下文无关文法(简称文法)的例子。

例如: if-else语句的构造规则可以表达为:

stmt -> if (expr) stmt else stmt

这里,箭头"->"可以读作"可以具有形式"。这样的规则称为产生式 (production)。在一个产生式中,像关键字if, else和括号这样的词法元素称为记号(token)或终结符(terminal),像*expr*和*stmt*这样的变量表示一个记号序列,并称之为非终结符(nonterminal)。



4.2.1 上下文无关文法的正式定义

- 上下文无关文法由终结符、非终结符、开始符号和产生式组成。其中:
 - 终结符是组成字符串的基本符号。在讨论程序设计语言的文法时,"记号"和"终结符"是同义词。
 - 非终结符是表示字符串集合的语法变量。非终结符所 定义的字符串集合有助于定义该文法所产生的语言。非终 结符强加给语言一种层次结构,这种层次结构对 语法分析和翻译都非常有用。
 - 有一个非终结符被定义为开始符号。
 - 文法的产生式说明了终结符和非终结符组合成串的方式。

例4.1



■ 具有下述产生式的文法定义了简单的算术表达式。

$$expr \rightarrow expr \ op \ expr$$

 $expr \rightarrow (expr)$
 $expr \rightarrow -expr$
 $expr \rightarrow id$
 $op \rightarrow +$
 $op \rightarrow -$
 $op \rightarrow +$
 $op \rightarrow /$
 $op \rightarrow \uparrow$

■ 在该文法中,终结符包括id、+、-、*、/、↑、(和),非终结符包括 expr和op, expr是开始符号。



4.2.2 符号的使用约定

- 1、下列符号是终结符:
 - 字母表中比较靠前的小写字母,如a、b、c等
 - 操作符号,如+、-等
 - 标点符号 如括号、逗号等
 - 数字0, 1, …, 9
 - <mark>黑体串</mark>,如id、if等
- 2、下列符号是非终结符:
 - 字母表中比较靠前的大写字母,如A、B、C等
 - 字母S,它常常代表开始符号
 - <mark>小写斜体名字</mark>,如*expr、stmt*等



- 3、字母表中比较靠后的大写字母,如X、Y、Z等,表示文法符号,也就是说,可以是非终结符也可以是终结符。
- 4、字母表中<mark>比较靠后的小写字母</mark>,如u、v、·····、z等,表示 终结符号的串。
- 5、小写希腊字母,如α、β、γ等,表示文法符号的串。
- 6、如果A $\rightarrow \alpha_1$ 、A $\rightarrow \alpha_2$ 、……、A $\rightarrow \alpha_k$ 是所有以A为左部的产生式(称为A产生式),则可以把它们写成: $A\rightarrow \alpha_1 |\alpha_2| \dots |\alpha_k,$

我们将 α_1 、 α_2 、.....、 α_k 称为A的候选式。

• 7、除非特别说明,否则第一个产生式左部的符号是开始符号。



例4.2



使用上述简写约定,例4.1的文法可以写成:

$$E -> E + T | E - T | T$$
 $T -> T * F | T / F | F$
 $F -> (E) | id$

根据上述约定,E,T和F是非终结符,E是开始符,其他符号都是终结符。



4.2.3 推导



考虑下面的算术表达式文法:

$$E -> E + E | E * E | (E) | -E | id$$
 (4-3)

则称: E => -E => -(E) => -(id) 是 E 到-(id)的推导。

• 更抽象地,我们说αAβ=>αγβ,如果A->γ是产生式,

<mark>而且α和β是任意的文法符号的串</mark>。如果:

$$\alpha_1 = > \alpha_2 = > \cdots = > \alpha_n$$
,则说 α_1 推导出 α_n 。

符号=>表示"一步推导"。通常我们用 👣 表示"零步或多步推导",因此,

- 对任何串α, α *> α。
- 如果α *> β, 而且β=>γ, 则α *> γ。





- 类似地,我们用 ⇒ 表示"一步到多步的推导"。当且 仅当S ⇒ w时,我们说终结符w在L(G)中。终结符串 w成为G的句子。由上下文无关文法产生的语言称为 上下文无关语言。如果两个产生相同的语言,则称 这两个文法等价。
- 对于开始符号S的文法G,如果S ≛ α,则称α为G的 句型,其中α可能含有非终结符。句子是不含非终结符的句型。



例4.3



■ 字符串-(id+id)是文法(4-3)的句子, 因为存在如下推导:

$$E = -E = -(E) = -(E+E) = -(id+E) = -(id+id)$$
 (4-4)

- 出现在这个推导中的字符串E、-E、-(E)、···、-(id+id)都是该文法的句型。 我们用E *> -(id+id)表示-(id+id)可以由E推导出来。
- 按推导长度进行归纳,我们可以证明文法(4-3)产生的语言中的每个句子都是由二元操作符+和*、一元操作符-、括号以及运算对象id组成的算术表达式。同样按算术表达式长度进行归纳,我们也可以证明这样的算术表达式都可以由文法(4-3)产生。因此,文法(4-3)正好产生所有包括二元操作符+和*、一元操作符-、括号以及操作数id的算术表达式的集合。



术语:



- 每一步都替代最左非终结符的推导, 叫做最左推导。
- ▶ 为了强调α通过最左推导推导出β这一事实,我们写α 🔭 β。
- 如果 S = α,则称 α 是该文法的左句型。
- 我们可以类似地定义最右推导,即每步推导都替代最右非终结符的推导。
- ■最右推导有时也称为规范推导。



4.2.4 分析树和推导

■ 分析树可以看成是推导的图形表示, <mark>但它</mark>不能显示出替代顺序的选择。回 顾2.2节, 分析树的每个内节点都标以 <mark>某个非终结符A</mark>。A的子节点从左到 右分别被用来替换A所使用的产生式 右部的各符号标记。分析树的叶节点 用非终结符或终结符来标记,它们从 左到右构成一个句型, 称为树的边界 <mark>或果实</mark>。例如,-(id+id)的推导过程如 (4-4)所示,其分析树如图4-2所示。



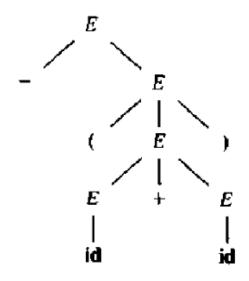


Fig. 4.2. Parse tree for -(id + id).





从该推导所构造出的分析树序列如图4-3所示。推导的第一步是E=>-E。 为了模拟这一步,我们为最初的分析树的根节点E增加两个子节点,分 别标记为-和E。推导第二步是-E=>-(E),所以为第二个分析树的标记为 E的叶结点增加三个子节点,分别标记为(、E和),从而获得带有结果-(E) 的第三棵树。如此继续,我们将得到第六棵树所示的这个分析树。

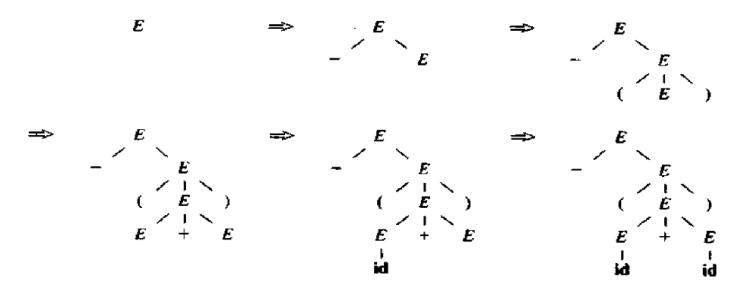


Fig. 4.3. Building the parse tree from derivation (4.4).



例4.4



- 每棵分析树都有一个与之对应的唯一的最左推导和唯一的最右推导。然而,每一个句子不一定只有一个分析树,或者说不一定只有一个最左推导或最右推导。
- 再次考虑算术表达式文法(4-3)。句子id+id*id有两个不同的推导:

$$E \Rightarrow E + E$$

$$\Rightarrow id + E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

$$E \Longrightarrow E*E$$

$$\Longrightarrow E+E*E$$

$$\Longrightarrow id+E*E$$

$$\Longrightarrow id+id*E$$

$$\Longrightarrow id+id*id$$



4.2.5 二义性



- 给定一个文法G,如果L(G)中存在一个具有两棵或两棵以上分析树的句子,则称G是二义性的。我们也可以如下定义二义性:如果L(G)中存在一个具有两个或两个以上最左(或最右)推导的句子,则G是二义性文法。
- 例4.4表明算术表达式文法:
 E -> E + E | E * E | (E) | -E | id
 就是一个二义性文法。



//4.2.6 验证文法所产生的语言



对"文法G产生语言L"的证明包括两部分:

- (1) 我们必须证明由G产生的每个字符串都在L中;
- (2) 反之, L中的每个字符串都能由G产生。



例4.5



文法G: $S \rightarrow (S)S \mid \mathbf{E}$ 能而且仅能产生所有的配对的括号串证明 $L(G) = \{s \mid s$ **是一个配对的括号串** $\}$ 的方法:

- (1) 首先证明从S推导出的所有句子都是配对的括号串,
- (2) 然后证明每个配对的括号串都可以从S推导出来。
- 对推导的步数使用数学归纳法,证明从S推导出的所有句子都是配对的括号串。从S经过一步推导得出的终结符串只有空串。空串可以视为配对的括号串。显然,当推导步数为1时命题正确。
- 现在假设所有少于n步的推导所产生的句子都是配对的括号串考察一个n 步最左推导。这个推导一定具有如下形式:

$$S => (S)S \stackrel{*}{=} (x)S \stackrel{*}{=} (x)y$$

由于从S推导出x和y的步数少于n步,根据归纳假设,x和y都是配对括号串。因此,(x)y也一定是配对括号串。





- 因此我们已经证明了任何从S中推导出来的串都是配对的括号串。下边, 我们对括号串的长度使用数学归纳法,证明每个配对括号串都可以从S 推导出来。空串是配对的括号串,并且使用产生式S->€由S推导出来。 于是,长度为0的括号串可以从S推导出来。
- 假设每个长度小于2n的配对括号串都可以从S推导出来,我们来考虑长度为2n(n >= 1)的配对括号串w。可以肯定,w是由左括号开始的,令(x)是w的具有相同个数的左右括号的最短前缀。那么w可以写作(x)y,其中x和y都是配对的括号串。既然x和y的长度小于2n,由归纳假设,它们可以从S推导出来。于是,我们可以找出如下形式的推导:



4.2.7 正规表达式和上下文无关文法的比较

我们可以用下列规则机械地<mark>把一个NFA转换成一个等价的上下文无关文法</mark>:

- 对NFA的每一个状态i,创建一个非终结符Ai
- 如果状态i遇见输入符号a转换到状态j,则引入产生式A_i->aA_i
- 如果状态i遇见输入符号€转换到状态j,则引入产生式A_i-> A_i
- 如果状态i是接受状态,则引入产生式A_i-> €
- 如果状态i是开始状态,则A_i是文法的开始符号





既然<mark>正规集都是上下文无关语言</mark>,我们可能要问:为什么要用正规表达式而不用上下文无关文法来定义语言的词法? 理由如下:

- 语言的词法规则通常都非常简单,不必动用强大的文法来解决
- 对于记号,正规表达式比上下文无关文法提供了更简洁且易于理解的定义
- 从正规表达式可以自动地构造出有效的词法分析器,从任何文法都很难构造词法分析器
- 把语言的语法结构分成词法和非词法两部分为编译器前端的模块 划分提供了方便的途径



4.3 设计文法



文法能够描述程序设计语言的大部分语法成分,但不能描述程序设计语言的全部语法成分。语法分析以后的各编译阶段必须分析语法分析器的输出,以保证输入字符串符合语法分析器无法检查的那些规则。



4.3.1词法分析和语法分析

- 词法分析器和语法分析器的分工明确。
- 词法规则的识别可以满足快速进行词法分析的要求。
- 每种语法分析方法只能处理一种形式的文法。为了适应所选择的分析方法,我们常常不得不改写初始文法。



4.3.2 消除二义性



• 有些二义性文法可以通过改写来消除二义性。作为一个例子, 我们来消除下面的"不匹配else"文法的二义性:

这里, other代表任何其他语句。



如图4-6。



• 文法(4-9)是具有二义性的,因为串
if E₁ then if E₂ then S₁
else S₂
有两棵分析树,

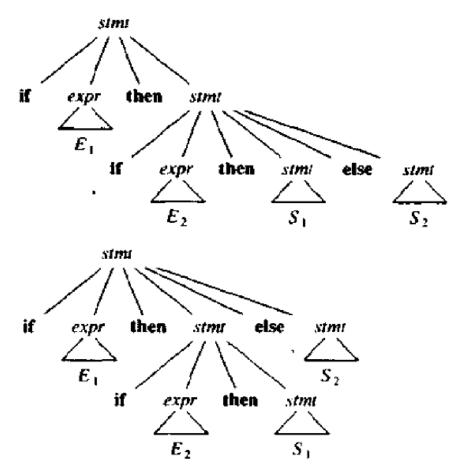


Fig. 4.6. Two parse trees for an ambiguous sentence.



例4.6



所有包含这种条件语句的程序设计语言都使用第一种分析树。

一般规则是,"每个else和前面最近的没有配对的then配对"。

这条避免二义性的规则可以直接并入文法中。

例如:可以把文法(4-9)改写成下面的无二义性文法,

其基本思想是:出现在then和else之间的语句必须是"配对"的,即它不

能以一个未配对的then后面跟随任意的非else语句结束,

于是else被迫与这个未配对的then配对。配对的语句是一

个不包含不配对语句的if-then-else语句或者任何非条件语句。

改写如下:

stmt -- matched_stmt

unmatched_stmt

matched_stmt - if expr then matched_stmt else matched_stmt

other

unmatched_stmt → if expr then stmt

if expr then matched_stmt else unmatched_stmt



4.3.3 消除左递归



- 左递归的定义: 如果文法G具有一个非终结符A使得对某个字符串α存在
 - 推导A =→ Aα,则称G是左递归的文法。
- 消除左递归的原因: 在最左推导中会造成无限推导的情况。
- 消除左递归的方法:对于文法中形如: $A \rightarrow A\alpha \mid \beta$ 直接左递归形式的产生式

其中 β 不以A为前缀。它产生的符号串为 $\beta\alpha$ *结构。引进非终结符A',令:

$$A \to \beta A'$$

$$A' \to \alpha A' \mid \varepsilon$$

它同样产生βα*形符号串,所以与等价,但是它不是左递归的。





• 考虑下面的算术表达式文法:

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

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

(4-10)

$$F -> (E) \mid id$$

消除E和T的直接左递归,可以得到:

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

(4-11)

$$T' \rightarrow *FT' \mid E$$

$$F -> (E) | id$$



消除左递归的一般方法



无论有多少A产生式,我们都可以用下面的技术来<mark>消除直接左递归</mark>。 首先,把A产生式放在一起:

 $A \rightarrow A\alpha_1 | A\alpha_2 | \cdots | A\alpha_m | \beta_1 | \beta_2 | \cdots | \beta_n$

其中: 每个βi都不以A开头。

然后用下面的产生式代替A产生式:

 $A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \cdots \mid \beta_n A'$

 $A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \cdots \mid \alpha_m A' \mid \epsilon$

但不能消除包括两步或多步推导的左递归,如:

S -> Aa | b

 $A \rightarrow Ac \mid Sd \mid E$

(4-12)



算法4.1 消除左递归

输入:无循环推导和€产生式的文法G

输出: 与G等价的无左递归文法

方法:对文法G应用图4-7的算法。

注意,得到的非左递归文法可能包含有C产生式。

1. Arrange the nonterminals in some order A_1, A_2, \ldots, A_n .

```
2. for i := 1 to n do begin
for j := 1 to i - 1 do begin
replace each production of the form A<sub>i</sub> → A<sub>j</sub>γ
by the productions A<sub>i</sub> → δ<sub>1</sub>γ | δ<sub>2</sub>γ | ···· | δ<sub>k</sub>γ,
where A<sub>j</sub> → δ<sub>1</sub> | δ<sub>2</sub> | ···· | δ<sub>k</sub> are all the current A<sub>j</sub>-productions;
end
eliminate the immediate left recursion among the A<sub>j</sub>-productions
end
```

Fig. 4.7. Algorithm to eliminate left recursion from a grammar.





文法(4-12)应用消除左递归算法后的过程为:

• 令非终结符的次序是S、A。在S产生式中没有直接左递归,所以在算法第2步,对于i=1,什么也没做。i=2时,用S产生式替换A->Sd中的S,得到下面的A产生式:

A -> Ac | Aad | bd | €

消除A产生式中的直接左递归,产生下面的文法:

 $S \rightarrow Aa \mid b$

 $A \rightarrow bdA' \mid A'$

A' -> cA' | adA' | €



4.3.5 提取左因子



算法4.2 提取左因子

输入: 文法G

输出: 一个等价的提取了左因子的文法

方法:对每个非终结符A,找出它的两个或更多候选式的最长公共前缀 α 。

如果 $\alpha \neq \mathbf{C}$,既有一个非平凡的公共前缀,则用下面的产生式代替

所有A产生式A-> $\alpha\beta_1 | \alpha\beta_2 | \cdots | \alpha\beta_n | \gamma$,其中 γ 表示所有不以 α 开头的

候选式: A -> αA' | γ

 $A' \rightarrow \beta_1 | \beta_2 | \cdots | \beta_n$

其中: A'是一个新的非终结符。反复应用这种变换,直到任一

非终结符都没有两个候选式具有公共前缀为止。





(4-13)

▶ 下面是条件语句抽象出来的文法:

这里的i、t和e分别表示if、then和else, E和S表示表达式和语句。提取左因子后,该文法变为

S -> iEtSS' | a

$$S' \rightarrow eS \mid \mathbf{\epsilon}$$
 (4-14)

 $E \rightarrow b$

于是,如果输入是i,我们可以将S扩展到iEtSS',等到iEtS出现后,再决定是将S'扩展到eS还是E。当然由于文法(4-13)和(4-14)都是具有二义性,所以对于输入e,我们不清楚应该选择S'的哪个候选式。



4.3.6 非上下文无关语言的结构

- L1是所有由c隔开的两个相同的a、b串组成的字母串集合。
- 例如: aabcaab。这个语言是检查程序中标识符的声明应先于其引用的抽象,即wcw中的第一个w表示标识符w的声明,第二个w表示它的引用。可以证明该语言不是上下文无关语言。(上下文无关文法不能描述两个相同的a、b串组成的字母串集合)

这个例子意味着C++和Java等程序设计语言都不是上下文无关语言,因为它们要求标识符的声明先于引用,并且允许标识符任意长。

 由于上述原因,描述C++和Java语法的文法并不定义标识符中的字符, 而只用文法中id这样的记号代表所有的标识符。在这类语言的编译器中, 语义分析阶段检查标识符的声明是否先于引用。





- 同理语言L₂={anbmcndm | n≥1且m≥1}不是上下文无关语言。
- L₂是由正规表达式a*b*c*d * 所表示的语言的子集合,在它的每个句子中,a和c的个数相同,b和d的个数相同。它是"过程声明中的形参个数和过程引用的实参个数应该一致"的抽象。
- 注意,过程定义和引用的语法并不涉及到参数的个数,通常 在语义分析阶段检查实参个数是否正确。



作业

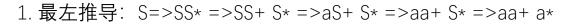


对下列每一个文法和对应的串完成:

- 给出串的一个最左推导;
- 给出串的一个最右推导;
- 给出串的一棵语法分析树;
 - (1) S->SS+|SS*|a 和串aa+a*
 - (2) S->0S1|01 和串000111
 - (3) S->aSbS|bSaS| € 和串aabbab

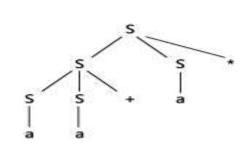


答案:



最右推导: S=>SS* =>Sa* =>SS+ a* =>Sa+ a* =>aa+ a*

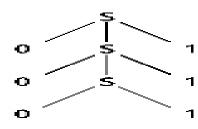
语法分析树:



2. 最左推导: S=>0S1=>00S11=>000111

最右推导: S=>0S1=>00S11=>000111

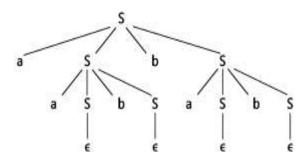
语法分析树:



3. 最左推导: S=>aSbS=>aaSbSbS=>aabSbS=>aabbS=>aabbaSbS=>aabbabS=>aabbab

最右推导: S=>aSbS=>aSbaSbS=>aSbaSb=>aSbab=>aaSbSbab=>aaSbbab=>aabbab

语法分析树:





4.4 自顶向下语法分析



- 语法分析的方法通常有两类,即自上而下分析方法和自下而上分析方法,或称为自顶向下(top-down)和自底向上(bottom-up)分析方法。它们的目标都是判断一输入串是否为一合法句子。
- 自上而下分析就是能否找到从文法开始符开始的推导序列,使得推导出的句子恰为输入串。或者说,能否从根结点出发,向下生长出一棵语法分析树,其叶结点组成的句子恰为输入串。显然语法分析树的每一步生长(每一步推导)都以能否与输入串匹配为准。



4.4.1 递归下降语法分析法

递归下降分析法是一种自顶向下的分析方法, 文法的每个非终结符,对应于一个递归过程。分析 就是从方法开始符出发执行一组递归过程,向下推 导,直到推导出句子。或者说从根结点出发,自上 而下为输入串寻找一个最左匹配序列,建立一棵语 法分析树。





考虑下述文法: S -> cAd

A -> ab | a

若输入字符串为w=cad

为了自顶向下地为w建立分析树,我们首先建立只有标记为S的单个节点的树。输入指针指向w的第一个符号c。然后,我们用S的第一个产生式来扩展该树,图4-9a所示的树。

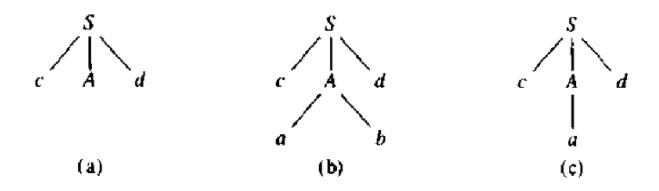


Fig. 4.9. Steps in top-down parse.





- 最左边的叶子标记为c, 匹配w的第一个符号。现在, 我们将输入指针移到w的第二个符号a。考虑下一个标记为A的叶子。用A的第一个候选式扩展A, 得到图4-9b所示的树。现在我们已经匹配了第二个输入符号a, 再将输入指针移到第三个输入符号d, 把它和下一个标记为b的叶结点进行比较。因为b和d不匹配, 报告失败, 回到A, 看是否还有别的候选式可试。
- 回到A时,我们必须将输入指针重置到第二个符号,即第一次进入A时的位置。这意味着A的程序必须将输入指针保存在一个局部变量中。现在尝试A的第二种候选式,得到图4-9c所示的分析树。叶子a匹配w的第二个符号,叶子d匹配w的第三个符号。因为已经产生了w的分析树,我们停止分析。
- 注意: (1) 左递归文法可能会使其进入无限循环。
- (2) 递归下降分析法可能需要回溯,即需要重复地扫描输入。



4.4.2 预测语法分析器



- 在许多情况下,通过仔细地编写文法,消除左递归,提取左因子, 我们可以获得一个有效的文法,这个文法可以用于不带回溯的递归下降 语法分析器。
- 为了构造预测语法分析器,对给定的当前输入符号a和将要扩展的非终结符A,我们必须知道,在A的所有可选产生式A->α₁ | α₂ | ··· | α_n中,哪个候选式是唯一能推导出以a开头的串。
- 预测语法分析器能够通过观察候选式所推导出的第一个符号,确定 正确的候选式。这种方法可以检测出多数程序设计语言中具有不同关键 字的控制流结构。



4.4.3 预测语法分析器的状态转换图

- 词法分析器的状态转换图和预测语法分析器的状态转换图 具有明显的区别。
- 对于预测语法分析器,每个非终结符都对应一个状态转换图,边上的标记是记号和非终结符。
- 记号(终结符)上的转换意味着如果该记号是下一个输入符号,就应 该进行转换。
- 非终结符A上的转换是对与A对应的过程的调用。
- 为了由文法构造预测语法分析器的状态转换图,首先消除文法中的左递归,然后提取左因子,并对每个非终结符A执行如下操作:
 - 创建一个开始状态和一个终态
 - 对每一个产生式 $A->X_1X_2\cdots X_n$,创建一条从开始状态到终止状态的路径,边上的标记分别为 X_1 , X_2 ,…, X_n 。



4.4.2 预测语法分析器



- 基于状态转换图的预测分析程序试图进行终结符和输入的匹配,并且,当它经过标记为非终结符的边时,进行潜在的递归过程调用。一种非递归的实现方法是,当在状态S上有一个标记为非终结符的指向其他状态的转换时,则将状态S压入栈中,当到达该非终结符的终止状态时,将状态S弹出栈。
- 如果给定的状态转换图是确定的,即一个状态对于一个输入仅有一个转换,则上述方法是有效的。如果出现二义性,可以用下边例4.10中的方法解决。如果不能消除不确定性,我们就不能构造预测语法分析器。但我们可以构造递归下降语法分析器,用回溯的方法尝试所有可能的情况。



图4-10给出了文法(4-11)所对用的状态转换图。

E -> TE'

E' -> +TE' | ϵ T -> FT'

T' -> *FT' | ϵ

F -> (E) | id

唯一的二义性在于确定是否经过Є边。如果我们把E'的开始状态的出边解释为:对E'的开始状态,如果下一个输入是+,则选择+上的转换,否则选择Є上的转换。对T'也作同样的假设,我们就消除了二义性。

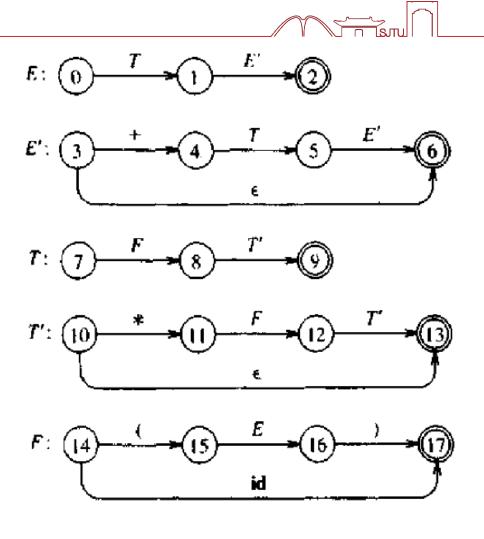


Fig. 4.10. Transition diagrams for grammar (4.11),



状态转换图的变换化简



通过图4-11的变换化简状态转换图。用C实现的与图4-12对 应的预测语法分析器要比图4-10所对应的快20%~25%。

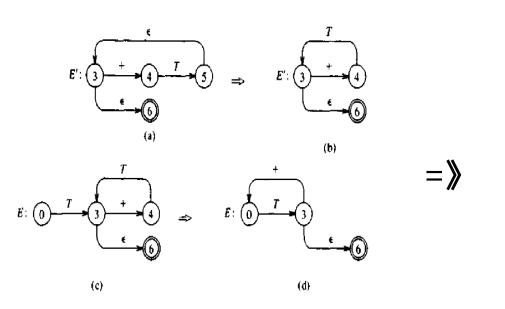


Fig. 4.11. Simplified transition diagrams.

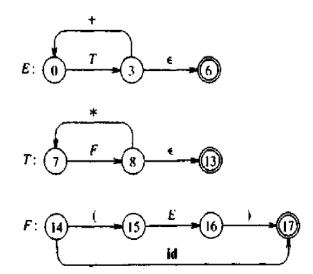


Fig. 4.12. Simplified transition diagrams for arithmetic expressions.



状态转换图的实现

```
void E()
{ T(); while(lookahead=='+'){match('+');T();}
void T()
{ F(); while(lookahead=='*'){ match('*');F();T'();}
void F()
  if(lookahead=='i'){match('i');}
   else if (lookahead=='(')
          {match('(');E();if(lookahead== ')')
                          {match(')');}
                          else {error();}}
       else{error();}
```

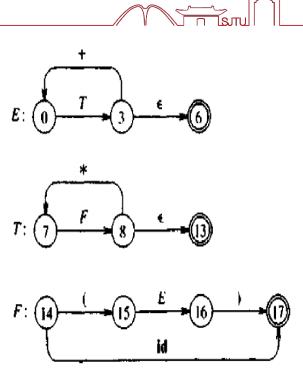


Fig. 4.12. Simplified transition diagrams for arithmetic expressions.



4.4.4 非递归的预测分析

通过显示地维护一个状态栈,而不是通过隐式的递归调用,我们可以构造非递归的预测语法分析器。<mark>预测分析的关键问题是确定用于扩展非终结符的产生式</mark>,图4-13中的非递归语法分析器通过查分析表来选取产生式。

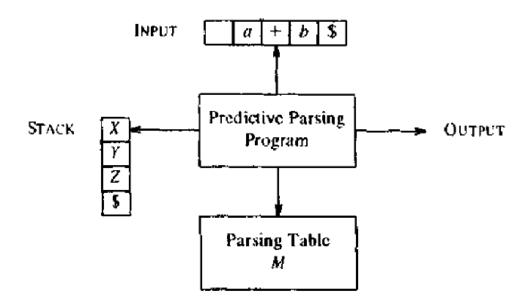


Fig. 4.13. Model of a nonrecursive predictive parser.



语法分析器工作方式

- 语法分析器由一个按如下方式工作的程序控制:程序根据<mark>栈顶当前的符号×和当前的输入符号a</mark>决定语法分析器的动作:
 - 如果X=a=\$, 则语法分析器宣告分析成功并停止。
 - 如果X=a≠\$,则语法分析器弹出栈顶符号X,并将输入 指针移到下一个输入符号上。
 - 如果X是非终结符,则程序访问分析表M的M[X,a]项。 M[X,a]项是文法的一个X产生式或者是出错 信息。



算法4.3 非递归的预测分析

输入: 串w和文法G的分析表M。

输出:如果w属于L(G),则 输出w的最左推导, 否则报告错误。

方法: 开始时,语法分析器的格局是\$S在栈里(其中S是G的开始符号且在栈顶),

图4-14是程序。

```
set ip to point to the first symbol of w$:
repeat
      let X be the top stack symbol and a the symbol pointed to by ip;
      if X is a terminal or $ then.
           if X = a then
                pop X from the stack and advance ip
           else error()
      else
                /* X is a nonterminal */
           if M(X, a) = X \rightarrow Y_1 Y_2 \cdots Y_n then begin
                pop X from the stack;
                push Y_k, Y_{k-1}, ..., Y_1 onto the stack, with Y_1 on top;
                output the production X \rightarrow Y_1 Y_2 \cdots Y_k
            end
            else error()
until X = S
                 /* stack is empty */
```

Fig. 4.14. Predictive parsing program.





考虑例4.6中的文法(4-11):

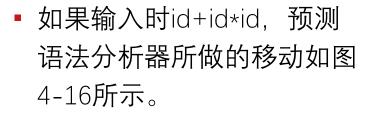
(4-11)

该文法预测分析表如图4-15所示。表中空白表项表示出错,非空白表项表示一个产生式。

NONTER-	INPUT SYMBOL						
	id	+	*	()	\$	
E	E→TE'			E→TE'			
\boldsymbol{E}^{r}		$E' \rightarrow +TE'$			E'→€	E' →€	
T	T→FT′		}	T →FT ′			
T		T'→€	T' →*FT'		Τ"→ε	Τ' →ε	
F	F→id			F→(E)	ļ		

Fig. 4.15. Parsing table M for grammar (4.11).





- 输入指针指向输入栏中符号 串最左边的符号。语法分析 器跟踪的是输入的最左推导, 即产生式输出的正好是最左 推导中使用的那些产生式。
- 已经扫描过的输入符号加上 栈中的文法符号(从顶到底) 构成该推导的左句型。

STACK	INPUT	Оптрит
\$ <i>E</i>	id + id * id\$	
\$E'T	id + id * id\$	$E \rightarrow TE'$
\$E'T'F	id + id * id\$	$T \rightarrow FT'$
\$ <i>E'T'</i> id	id + id * id\$	$F \rightarrow id$
\$ <i>E'T'</i>	+ id * id\$	
\$ <i>E</i> '	+ id * id\$	T' → €
\$E'T +	+ id * id\$	$E' \rightarrow +TE'$
\$E'T	id * id\$	
\$E'T'F	id * id\$	$T \rightarrow FT'$
\$E'T'id	id * id\$	F → id
\$ <i>E'T'</i>	* id\$	
\$E'T'F*	* id\$	$T' \rightarrow *FT'$
\$E'T'F	id\$	
\$E'T'id	id\$	$F \rightarrow id$
\$E'T'	\$	
\$ <i>E</i> '	s	$T' \rightarrow \epsilon$
\$	s	$E' \rightarrow \epsilon$

Fig. 4.16. Moves made by predictive parser on input id+id*ic



4.4.5 FIRST和FOLLOW



- 构造文法G的分析表需要两个与G有关的函数FIRST和 FOLLOW。
- 我们可以用这<mark>两个函数来填写G的分析表的表项</mark>。
- 由FOLLOW函数产生的记号集合还可用作紧急方式错误恢 复期间的同步记号。



FIRST集合的构造方法



为了计算文法符号X的FIRST(X),我们可以应用下列规则,直到没有终结符或E可加到某个FIRST集合为止:

- 如果X是终结符,则FIRST(X)是{X}。
- 如果X->€是一个产生式,则将€加到FIRST(X)中。
- 如果X是非终结符,且X->Y₁Y₂...Y_k是一个产生式,

则: 若对于某个i, a属于FIRST(Y_i)且€属于FIRST(Y₁), ···,FIRST(Y_{i-1}),

即 $Y_1 \cdots Y_{i-1} \stackrel{*}{=} \epsilon$,

则将a加入FIRST(X)中;

若对于所有的j=1, 2, ···, k, €在FIRST(Y_i)中,

则将加到FIRST(X)中。例如: FIRST(Y_1)中的每个元素确实都在FIRST(X)中。如果 Y_1 不能推导出 \mathfrak{E} ,则不再往FIRST(X)中增加新的符号,如果 $Y_1 \stackrel{*}{=} > \mathfrak{E}$,则将FIRST(Y_2)加到FIRST(X)中,以此类推。



FOLLOW集合的构造方法



- 为计算所有非终结符A的后继符号集合FOLLOW(A), 我们可以应用如下规则, 直到每个FOLLOW集合都不能再加入任何符号或\$为止:
 - 将\$放入FOLLOW(S)中, 其中S是开始符号, \$是输入串的结束符。
 - 如果存在产生式A->αBβ,则将FIRST(β)中除€以外的符号都放入 FOLLOW(B)中。
 - 如果产生式A->αB, 或A->αBβ, 其中: FIRST(β)中包含€
 (即β *> €) , 则将FOLLOW(A)中的所有符号都放入
 FOLLOW(B)中。





我们再来看文法(4-11):

```
E -> TE'
```

T' -> *FT' |
$$\epsilon$$

那么:

$$FIRST(E) = FIRST(T) = FIRST(F) = \{(, id)\}$$

$$FIRST(E') = \{+, \mathbf{\epsilon}\}$$

FIRST(T') =
$$\{*, \ \mathbf{\epsilon}\}$$

$$FOLLOW(E) = FOLLOW(E') = {}, {}$$

$$FOLLOW(T) = FOLLOW(T') = \{+, \}$$

$$FOLLOW(F) = \{+, *, \}$$

(4-11)



4.4.6 预测分析表的构造



算法4.4 构造预测分析表。

• *输入*: 文法G。

■ *输出*:分析表M。

• 方法:

- 1、对于文法中的每个产生式A->α, 执行第2和第3步。
- 2、对FIRST(α)中的每个终结符a,将A->α加入到M[A,a]中。
- 3、若ε在FIRST(α)中,则对FOLLOW(A)的每个终结符b, 将A->α加入到M[A,b]中;若ε在FIRST(α)中,且\$在 FOLLOW(A)中,则将A->α加入到M[A,\$]中。
- 4、将M中每个没定义的表项均设置为error。





NONTER- MINAL	INPUT SYMBOL							
	id	+	*	()	\$		
E	$E \rightarrow TE'$			E →TE'				
\boldsymbol{E}^{i}		$E' \rightarrow +TE'$			E' →€	E'→€		
T	$T \rightarrow FT'$		1	T→FT'				
T'		T'→€	T'→*FT'		T' → €	Τ' →ε		
F	F→id			F→(E)	ļ			

Fig. 4.15. Parsing table M for grammar (4.11).

- FIRST(TE')=FIRST(T)={(, id},
- ∴ 产生式E->TE'被填入M[E, (]和M[E, id] 表项中.
- FIRST(+TE')={+},
- ∴ 产生式E'->+TE'被填入M[E', +]表项中
- :: FIRST(€)={ € }, 而需要求FOLLOW(E')={), \$}
- ∴ 产生式E'-> € 被填入M[E',)]和M[E', \$]表项含有E'->€。

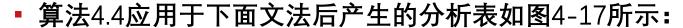




消除下列文法中左递归, 并构造一个预测分析表



4.4.7 LL(1)文法



S -> iEtSS' | a S' -> eS | € E -> b

NONTER- MINAL	INPUT SYMBOL					
	a	b	e	i	t	\$
S	S→a			S→iEtSS'		
S'			S' → ∈ S' → eS			S, →
E		E→b				

Fig. 4.17. Parsing table M for grammar (4.13).

- ∵ FOLLOW(S')={e, \$},
- ∴ M[S', e]同时包含S'->eS和S'->€。该文法具有二义性的。
- LL(1)文法的定义:分析表中没有多重定义表项的文法叫做LL(1)文法。算法4.4可以为任何LL(1)文法G产生分析表。这个分析表能分析G的所有句子,而且只能分析G的句子。



LL(1)文法的性质

- LL(1)文法有一种特殊的性质。它不是二义性的,也不含左递归。 可以证明,G是LL(1)文法,当且仅当G的任何两个不同的产生式 A->α|β满足下面的条件:
 - 不存在这样的终结符a,使得α和β导出的串都以a开始。
 - α和β中至多有一个能导出空串。
 - 如果 $\beta \stackrel{*}{=} \epsilon$,那么 α 不能导出以FOLLOW(A)中的终结符开始的任何字符串。
 - 注意: 如果一个文法经过消除左递归和提取所有可能的左因子后也不能产生LL(1)文法。一般来说,没有一个普遍适用的规则可以用来删除多重定义的表项,使其成为单值而不影响语法分析器所识别的语言。



作业



对下列每一个文法完成:

- 提取左公因子;
- 消除左递归;
- 构造预测分析表;
- 判断文法是否是LL(1)文法;
 - (1) S->SS+|SS*|a
 - (2) S->0S1|01
 - (3) S->aSbS|bSaS| €



作业答案



(1)
$$S->SS+|SS*|a$$

• 提取左公因子: S ->SSS'| a

• 消除左递归: S->aS"

• 构造预测分析表: FIRST(aS")={a}, FIRST(SS'S")={a},



$$FIRST(+)=\{+\}$$
, $FIRST(*)=\{*\}$, $FIRST(\mathbf{C})=\{\mathbf{C}\}$

$$FOLLOW(S")=FOLLOW(S)=\{+,*,\$\}$$

	+	*	а	\$
S			S->aS"	
S"	S"-> €	5"-> €	S"->SS'S"	5'-> €
S'	S'->+	S'->*		

表中没有多重定义的表项,所以是 LL(1)文法。



作业答案



- (2) S->0S1|01
- 提取左公因子: S->0S'

$$S' - > S1|1$$

■ 构造预测分析表: FIRST(OS)={0}

$$FIRST(S1)=\{0\}$$

	0	1	\$
S	S->0S'		
S'	S'->S1	S->1	

■ 表中没有多重定义的表项, 所以是 LL(1) 文法。



作业答案



- (3) S->aSbS|bSaS| €
- 构造预测分析表: FIRST(aSbS)={a}, FIRST(bSaS)={b}

$$FIRST(\mathbf{E}) = {\mathbf{E}}, FOLLOW(S) = {a,b,\$}$$

	а	b	\$
S	S->aSbS	S->bSaS	
	S-> €	S-> E	S-> €

■ 表中有多重定义的表项, 文法有二义性, 不是 LL (1) 文法

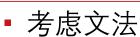


4.5 自底向上语法分析

- 本节介绍一种比较常用的自底向上分析方法, 称为移动归约分析法。
- 移动归约分析法为输入串构造分析树时从叶结点开始, 像根节点前进。
- 我们可以把该过程看成是把输入串w"归约"成文法开始符号的过程。在每一步归约中,如果一个子串和某个产生式的右部匹配,则用该产生式的左部符号代替该子串。如果每一步都能恰当地选择子串了,我们就可以得到最右推导的逆过程。



例4.15



```
S -> aABe
A \rightarrow Abc \mid b
B \rightarrow d
句子abbcde可按下述步骤归约到S:
abbcde
aAbcde
aAde
aABe
S
```



4.5.1 归约



- 非形式定义: 一个符号串的"句柄"是和一个产生式 右部匹配的子串,而且把它'归约'到该产生式左部的非 终结符代表了最右推导逆过程的一步。
- 形式定义: 右句型γ的句柄是一个产生式A->β以及 γ的一个位置, 在该位置可以找到串β, 而且用A代替β 可以得到γ的最右推导的前一个右句型, 即如果S ** αAw ** αβw, 那么紧跟在α后面的

A->β是αβw的句柄。

注意: 如果文法是具有二义性的,则句柄不一定唯一,因为可能有不止一个αβw的最右推导。只有文法没有二义性时,它的每个右句型才有一个句柄。在上面的例子中,abbcde是右句型,句柄是在位置2的A->b。



例4.16



(3)
$$E \to (E)$$
 (4-16)

(4)
$$E \rightarrow id$$

■ 最右推导:
$$E \underset{\text{rin}}{\Rightarrow} E + E$$
 $\underset{\text{rin}}{\Rightarrow} E + E * E$
 $\underset{\text{rin}}{\Rightarrow} E + E * id_3$
 $\underset{\text{rin}}{\Rightarrow} E + id_2 * id_3$
 $\underset{\text{rin}}{\Rightarrow} id_1 + id_2 * id_3$

方便起见,我们给id加以下标,并给每个右句型的句柄加上下划线。例如,id₁是右句型id₁+id₂*id₃的句柄。注意句柄右边的串中仅含终结符。



■ 因为文法(4-16)是具有二义性,存在 $id_1+id_2*id_3$ 的另一个右推导:

$$E \underset{rm}{\Longrightarrow} E * \underline{E}$$

$$E * \underline{id}_3$$

$$E + \underline{E} * \underline{id}_3$$

$$E + \underline{id}_2 * \underline{id}_3$$

$$E + \underline{id}_2 * \underline{id}_3$$

$$E + \underline{id}_2 * \underline{id}_3$$

■ 考虑右句型E+E*id₃,在该推导中E+E是E+E*id₃的句柄,而在上一个 推导中id₃是该右句型的句柄



4.5.2 句柄裁剪



- 通过"裁剪句柄"可以得到最右推导的逆过程。
- 考虑例4.16的文法(4-16)的文法和输入串id₁+id₂*id₃。它的规约序列如图 4-21所示。容易看出,该右句型序列正好是例4.16中第一个最右推导序 列的逆序。

RIGHT-SENTENTIAL FORM	HANDLE	REDUCING PRODUCTION
$\mathbf{id}_1 + \mathbf{id}_2 * \mathbf{id}_3$	id,	E → id
$E + id_2 * id_3$	id ₂	E → id
$E + E * id_3$	id ₃	E → id
E + E * E	E * E	$E \rightarrow E * E$
E + E	E + E	$E \rightarrow E + E$
E		

Fig. 4.21. Reductions made by shift-reduce parser.



4.5.3 移动-归约语法分析技术

- 语法分析器将零个或多个输入符号压入栈,直到句柄β 在栈顶出现为止,语法分析器再把β归约成某个恰当的产生式右部。语法分析器重复此过程,直到它发现错误或者栈中只含有开始符号并且输入串为空: STACK SS
- 进入这个格局后, 语法分析器停止并宣告分析成功。



例4.17

让我们逐步看一 下移动归约语法分析 器在分析输入串 id₁+id₂*id₃ 时的动作, 文法是(4-16), 使用 例4.16的第一种推导。 动作序列如图4-22所 示。注意,由于文法 (4-16)对该输入有两 种最右推导,所以语 法分析器还可能采取 另一个动作序列。

	Stack	INPUT	Action
(1)	\$	$id_1 + id_2 * id_3$ \$	shift
(2)	\$id,	+ id ₂ * id ₃ \$	reduce by $E \rightarrow i\mathbf{d}$
(3)	\$ <i>E</i>	$+ id_2 * id_3$	shift
(4)	\$ <i>E</i> +	id ₂ * id ₃ \$	shift
(5)	$E + id_2$	* id,\$	reduce by $E \rightarrow id$
(6)	E + E	* id ₃ \$	shift
(7)	E + E *	id,\$	shift
(8)	$SE + E * id_3$	\$	reduce by $E \rightarrow id$
(9)	E + E * E	\$	reduce by $E \to E * E$
(10)	E + E	\$	reduce by $E \rightarrow E + E$
(11)	\$ <i>E</i>	s	accept

Fig. 4.22. Configurations of shift-reduce parser on input $id_1 + id_2 * id_3$.



移动归约语法分析器的基本动作

- 移动归约语法分析器的基本动作是移动和归约的,但实际上有四种可能的动作:移动、归约、接受、出错。
 - ▶ 移动: 把下一个输入符号移动到栈顶。
 - 归约: 语法分析器知道句柄的右端已在栈顶。它必须在栈中确定句柄的左端,并选择正确的非终结符替代句柄。
 - 接受: 语法分析器宣告分析成功。
 - 出错: 语法分析器发现了一个语法错误, 并调用错误恢复处理程序进行错误处理。



4.5.4 可行前缀 (活前缀)

- 出现在移动归约语法分析器<mark>栈中的右句型的前缀集合称为可行前缀。</mark>
- 等价的定义为: 可行前缀是右句型的前缀,而且其右端不会超过该 句型的最右边句柄的末端。

即:若 $A \rightarrow \alpha\beta$ 是文法的一个产生式,S是文法开始符,

并有: $S \underset{R}{\Rightarrow} \delta A \omega \underset{R}{\Rightarrow} \delta \alpha \beta \omega$

则 $\delta \alpha \beta$ 的任何前缀都是规范句型 $\delta \alpha \beta \omega$ 的可行前缀。

在上述定义中, $\alpha\beta$ 是句型 $\delta\alpha\beta\omega$ 关于A的直接短语, 并且是一最右推导(规范推导), 所以 $\alpha\beta$ 是最左直接短语, 是一句柄,

因此 $\delta\alpha\beta$ 的任何前缀不含句柄后的任何符号,特别,句柄是 $\alpha\beta$ 的后缀 $\delta\alpha\beta$ 是分析栈栈顶的符号串。



4.5.5 移动归约分析过程的冲突

有些上下文无关文法不能使用移动归约分析方法进行分析。这种文法的每一个移动归约语法分析器会形成这样的格局: 根据栈中的内容和下个输入符号不能决定是移动还是归约(移动-归约冲突),或不能决定按哪一个产生式进行归约(归约-归约冲突)。

二义性文法一定不能使用移动归约分析方法进行分析。



4.6 LR语法分析技术介绍: 简单LR技术

- 本节介绍一种有效的自底向上的语法分析技术。我们也将给出三种构造LR分析表的方法:
- 第一种方法称为<mark>简单LR方法(简称SLR</mark>),它最容易实现,但功能最弱,对某些文法可能失败。
- 第二种方法称为<mark>规范的LR方法</mark>,它功能最强,代价也 最高。
- 第三种方法叫做<mark>向前看的LR方法</mark>(简称LALR),其功 能和代价介于前两者之间。



4.6.1 为什么使用 LR语法分析技术



- LR分析富有吸引力的原因如下:
 - LR语法分析器能识别几乎所有能用上下文无关文法描述的程序设计语言的结构。
 - LR分析法是已知的最一般的无回溯移动归约语法分析法, 而且可以和其他移动归约分析法一样被有效地实现。
 - LR分析法分析的文法类是预测分析法能分析的文法类的真超集。
 - 在自左向右扫描输入符号串时,LR语法分析器能及时发现语法错误。
- 主要缺点是: 对典型的程序设计语言文法, 手工构造 LR语法
 - 分析器的工作量太大,需要专门的工具。如Yacc。



4.6.2 LR语法分析算法

- LR语法分析器的模型如图4-29所示。
- 驱动程序对所有的LR语法分析器都是一样的,不同的语法分析器只是语法分析表不同。

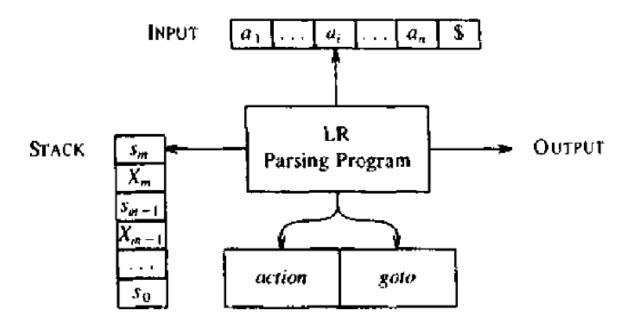


Fig. 4.29. Model of an LR parser.



- 语法分析表由动作函数action和转移函数goto两部分组成。驱动LR语法分析器的程序工作如下: 它根据<mark>栈顶状态s_m和当前输入符号a_i访问action[s_m, a_i],即s_m和a_i所对应的语法分析表项的动作部分,可能的动作如下:</mark>
 - ■移动状态s进栈。
 - 按文法产生式A->β归约。
 - 接受
 - 出错
- 函数goto以状态和文法符号为参数,产生一个状态。我们将看到,用SLR、规范的LR和LALR方法从文法G构造的语法分析表的goto函数是识别G的活前缀的确定的有穷自动机的转换函数。



■ 算法4.7 LR分析算法 *输入*: 文法G的LR语法 分析表和输入串w。 输出:如果w属于L(G), 则输出w的自底向上 分析. 否则报错。 方法: 首先, 把初始状 态Sn放在语法分析器 的栈顶,把w\$放在输 入缓存区;然后。语 法分析器执行图4-30 的程序, 直到遇见接

受或出错动作为止。

set ip to point to the first symbol of w\$; repeat forever begin let s be the state on top of the stack and a the symbol pointed to by ip; if action(s, a) = shift s' then begin push a then s' on top of the stack; advance ip to the next input symbol else if $action[s, a] = reduce A \rightarrow \beta$ then begin pop $2 * |\beta|$ symbols off the stack; let s' be the state now on top of the stack; push A then goto[s', A] on top of the stack; output the production $A \rightarrow \beta$ end else if action[s, a] = accept thenreizem else error()

Fig. 4.30. LR parsing program.



例4.18



- 下面是含有二元操作符
 - +和*的算术表达式文法:

$$(1) \quad E \to E + T$$

- $(2) \quad E \to T$
- $(3) \quad T \to T * F$
- $(4) \qquad T \to F$
- $(5) \quad F \to (E)$
- (6) $F \rightarrow id$
- 图4-31给出了文法G的LR语法分析表,包括动作函数和转移函数

Con a rec	action					goto			
\$tate	id	+	*	()	\$	E	T	F
0	s5		· <u> </u>	s 4			1	2	3
1		s6				acc			
2		τ2	s7		r2	£2	İ		
3	{	r4	г4		r4	г4	1		
4	s5			s4			8	2	3
5		r6	r6		r6	r6	l		
6	s5			s4			}	9	3
7	s5			s4			ļ		ŧ0
8		só			s11				
9	}	ri	s7		ιj	rl			
10	}	r3	r3		г3	r3			
11		£ 5	r5		r 5	r5			

Fig. 4.31. Parsing table for expression grammar.





- 各个动作的编码的含义是:
 - 1、s_i表示把当前输入符号和状态i压进栈。
 - 2、r_i表示按第j个产生式归约。
 - 3、acc表示接受。
 - 4、空白表示出错。
- 注意,对于终结符a,状态转移动作goto[s, a]可以 在表的动作表项action[s, a]中找到,所以转移域中 仅给出非终结符A的goto[s, A]。另外,我们还没有 解释图4-31的表项都是怎么确定的,这个问题稍后 讨论。





- 输入字符串: id*id+id,
- 栈和输入内容的变化序 列如图4-32所示。
- 例如第一行,状态0,id 是第一个输入符号,对 应的动作域是s5,则id进 栈,状态5压入栈顶。
- 然后,*为当前输入符号, 状态5,按F->id归约:弹 出栈顶两个符号,栈顶 状态0。因为goto[0,F]是 3,则把F和3压入栈。
- 其余类推。

	Stack	INPUT	ACTION
(1)	0	id * id + id \$	shift
(2)	0 id 5	* id + id \$	reduce by $F \rightarrow id$
(3)	0 F 3	* id + id \$	reduce by $T \to F$
(4)	0 T 2	* id + id \$	shift
(5)	0 T 2 * 7	id + id \$	shift
(6)	0 T 2 * 7 id 5	+ id \$	reduce by $F \rightarrow id$
(7)	0.72 * 7F10	+ id \$	reduce by $T \to T *F$
(8)	0 T 2	+ id \$	reduce by $E \to T$
(9)	0 <i>E</i> 1	+ id \$	shift
(10)	0EI + 6	id \$	shift
(11)	0E1 + 6 id 5	\$	reduce by $F \rightarrow id$
(12)	0E1 + 6F3	s	reduce by $T \to F$
(13)	0E1 + 6T9	s	$E \rightarrow E + T$
(14)	0 E i	s	accept

Fig. 4.32. Moves of LR parser on id * id + id.



4.6.3 LR文法

- 给定文法G,如果我们能为G构造出LR语法分析表,则称G是LR文法。很多上下文无关文法不是LR文法。但 是典型的程序设计语言的结构一般都可以避免非LR文法。
- 直观地,为了使一个文法是LR文法,只要保证在句柄出现在栈顶时,自左向右扫描的移动归约分析器能够及时地识别它。
- 能够用来帮助LR语法分析器做出移动归约决定的另一个信息源是下k个输入符号。例如图4-31是LR(1)。每步需要向前看k个符号的LR语法分析器所分析的文法叫做LR(k)文法。



LL文法和LR文法之间区别

- 对于LR(k)文法,我们必须通过向前看k个输入符号就能够知道一个产生式的右部所能推导出的所有字符串,进而识别出这个产生式右部的出现。这个要求远远不如LL(k)那么严格。
- LL(k)文法,在决定是否使用某个产生式时,要求只能向前看该产生式右部推出的前k个符号。所以LR文法比LL文法描述的语言更多。



4.6.4 构造SLR语法分析表

■ 文法G的LR(0)项目是在G的产生式右部的某处加点的 产生式。例如,产生式A->XYZ可以生成如下四个项目:

 $A \rightarrow XYZ$

 $A \rightarrow X \cdot YZ$

 $A \rightarrow XY \cdot Z$

 $A \rightarrow XYZ$

- · SLR方法的主要思想是首先从文法构造出识别活前缀的确定有穷自动机。
- 我们把项目划分成一组集合,这些集合对应SLR语法分析器的状态。项目可以看成识别活前缀的NFA的状态,而"项目分组"就是3.6节中讨论的子集构造法。





- 被称为规范LR(0)项目集族是构造SLR语法分析表的基础。为了构造文法的规范LR(0)项目集族,我们需要定义拓广文法的概念,并引入闭包(closure)运算和转移函数(goto)。
- 如果文法G的开始符号是S,那么G的拓广文法G'是在G的基础上增加一个新的开始符号S'和产生式S'->S。新产生式的目的是用来指示语法分析器什么时候应该停止分析并宣布接受输入,即当且仅当语法分析器执行归约S'->S时,分析成功。



4.6.4.1 闭包运算closure

■ 如果 I 是文法G的项目集,那么closure(I)是

从I出发由下面两条规则构造的项目集:

- 1、初始时,把I的每个项目都加入到closure(I)中。
- 如果A->α·Bβ在closure(I)中,且存在产生式
 B->γ,若B->·γ不在closure(I)中,则将其加入到
 closure(I)中。反复运用这条规则,直到没有更多的项目可以加到closure(I)为止。



例4.19



• 考虑拓广的表达式文法:

$$E' \rightarrow E$$

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

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

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

$$(4-19)$$

• 如果I是一个项目集合{[E' -> ·E]}, 那么closure(I)包含下

列项目:
$$E' \rightarrow \cdot E$$
 $E \rightarrow \cdot E + T$
 $E \rightarrow \cdot T$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$





■ 例计算函数closure的算法如图4-33所示。

```
function closure ( I );

begin

J := I;

repeat

for each item A \to \alpha \cdot B\beta in J and each production

B \to \gamma of G such that B \to \gamma is not in J de

add B \to \gamma to J

until no more items can be added to J;

return J

end
```

Fig. 4.33. Computation of closure.

- 核心项目:初始项S'->S和所有点不在左端的项目
- 非核心项目: 点在左端的非初始项目。



4.6.4.2 goto函数

- 第二个非常有用的函数是goto(I, X), 其中I是项目集, X是文法符号。goto(I, X)定义为所有项目集 $[A->\alpha X\cdot\beta]$ ($[A->\alpha \cdot X\beta]$ 在I中)的闭包。
- 例4.35 若I是两个项目的集合{[E'->E], [E-> E·+T]}, 则goto(I, +)包含下列项目:

$$E \rightarrow E + \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot id$$



4.6.4.3 项目集的构造



```
procedure items(G');
begin

C := {closure({[S' → ·S]})};
repeat
    for each set of items I in C and each grammar symbol X
        such that goto(I, X) is not empty and not in C do
        add goto(I, X) to C
. until no more sets of items can be added to C
end
```

Fig. 4.34. The sets-of-items construction.



例4.20

考虑拓广的表达 式文法:

$$E' \rightarrow E$$

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

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

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

规范LR(0)项目集 族如图4-35所示。

$$I_{0}: E' \rightarrow \cdot E$$

$$E \rightarrow \cdot E + T$$

$$E \rightarrow \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot id$$

$$I_{1}: E' \rightarrow E \cdot E \rightarrow E \cdot + T$$

$$T \rightarrow T * F$$

$$T \rightarrow T * F$$

$$E \rightarrow E \cdot + T$$

$$T \rightarrow T * F$$

$$T \rightarrow \cdot T *$$

Fig. 4.35. Canonical LR(0) collection for grammar (4.19).

 $F \rightarrow \cdot id$



• 考虑拓广的表达式文法:

$$E' \rightarrow E$$

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

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

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

· 这个项目集的 goto函数被示为 图4-36中有穷自动 机D的转换状态图。

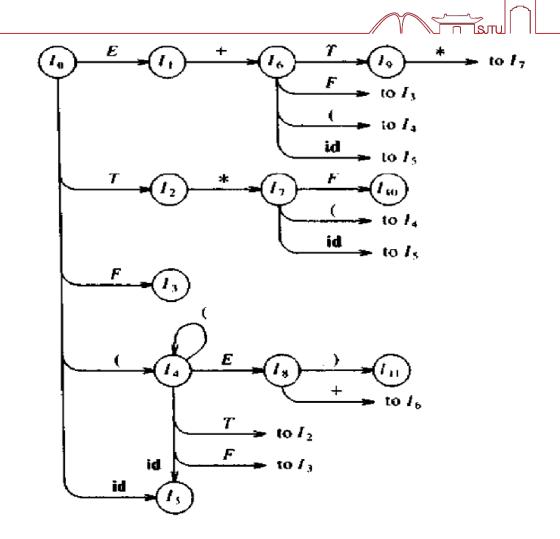


Fig. 4.36. Transition diagram of DFA D for viable prefixes.





• 有效项目:如果存在一个推导S' 💺 αAw 💺 α β₁ β₂w,

则:项目A-> β_1 · β_2 对可行前缀 $\alpha\beta_1$ 是有效的。

- 一般而言,同一个项目可能对多个可行前缀有效。
- $A -> \beta_1 \cdot \beta_2$ 对前缀 $\alpha \beta_1$ 有效这个事实告诉我们,在发现 $\alpha \beta_1$ 在分析栈时是 移进还是归约,特别是:

如果: $\beta_2 \neq \mathbf{C}$, 它暗示句柄还没有完全进栈,

动作应该是移进。

如果: $β_2$ = $\mathbf{\epsilon}$, 那么A-> $β_1$ 是句柄,

应该用这个产生式归约。





- 注意 (1) 当然,同一个可行前缀的两个有效项目可能告诉我们做不同的事情,有些这样的冲突可以通过向前看下一个输入符号来解决,其他一些冲突可以用下一节的方法解决。当LR方法用于构造任意文法的语法分析表时,不能保证所有冲突都能解决。
- 注意 (2) : 如果A是一个由规范项目集族构造、以goto函数为转换函数的DFA,则一个活前缀γ的有效项目集正好从A的初态出发,沿着标记为γ的路径所能到达的那些项目的集合。这是LR分析理论的一条基本定理。



4.6.4.4 SLR语法分析表的构造算法

- *输入*: 拓广文法G'
- *输出*: G'的SLR语法分析表函数action和goto
- 方法: 1、构造C= $\{I_0, I_1, \cdots, I_n\}$,即G'的规范LR(0)项目集族。
 - 2、从I_i构造状态i,它的分析动作确定如下:
 - (a) 如果[A->α·aβ]在 I_i 中,并且goto(I_i ,a)= I_j ,则置action[i,a]为"移动j进栈",这里a必须是终结符。
 - (b) 如果[A->α·]在I_i中, 则对FOLLOW(A)中的所有a, 置action[i,a]为"归约A->α",这里A不能是S'。
 - (c) 如果[S'->S]在I;中,则置action[I,\$]为"接受"。

如果:由上面的规则产生的动作有冲突,则G不是SLR(1)文法, 就构造不出语法分析器。

- 3、对所有的非终结符A,使用下面的规则构造状态i的goto函数: 如果goto(I_i ,A)= I_i ,则goto[i,A]=i。
- 4、不能由规则(2)和(3)定义的表项都置为"出错"。
- 5、语法分析器的初始状态是从包含[S'->·S]的项目集构造出的状态。



• 让我们为文法(4-19)构造SLR表。它的规范LR(0)项目集族已在图4-35中给出。 考虑 I_0 :

$$E' \rightarrow \cdot E$$

$$E \rightarrow \cdot E + T$$

$$E \rightarrow \cdot T$$

$$T \rightarrow \cdot T * F$$

$$T \rightarrow \cdot F$$

$$F \rightarrow \cdot (E)$$

$$F \rightarrow \cdot id$$

• 项目F->·(E)使得action[0,(]="移动4进 栈",项目F->·id使得action[0,id]="移动5 进栈"。 I_0 其他项目不产生动作。其他项 目集一次类推,可得到图4-31的分析表。

\$tate	L		æ	tion				goto	
	id	+	*	()	\$	E	T	F
0	s 5			s4			1	2	3
1)	s6				acc			
2		г2	s7		r2	£2	ļ		
3	{	r4	г4		r4	г4	1		
4	s5			s4			8	2	3
5		r6	r6		r6	r6			
6	s5			s4			1	9	3
7	s5			s4					ŧ0
8		só			s11				
9	}	ri	s7		ri	rl			
10	}	r3	r3		г3	r3			
11		£ 5	r5		г\$	r5			

Fig. 4.31. Parsing table for expression grammar.



注意:

- 每个SLR(1)文法都不是二 义的,
- 但有很多非二义的文法不 是SLR(1)文法。
- ▶ 考虑下面的文法:

$$S \rightarrow L=R$$

$$S \rightarrow R$$

(4-20)

$$L \rightarrow id$$

文法(4-20)的规范LR(0) 项目集族如图4-37所示。

$$I_0: S' \rightarrow \cdot S$$

 $S \rightarrow \cdot L = R$

$$S \rightarrow R$$

$$L \rightarrow \cdot *R$$

$$L \rightarrow -id$$

$$R \rightarrow \cdot L$$

$$I_1: S' \rightarrow S$$

$$l_2: S \rightarrow L \cdot = R$$

 $R \rightarrow L \cdot$

$$I_3: S \rightarrow R$$

$$I_4: L \rightarrow *R$$

$$R \rightarrow L$$

$$L \rightarrow *R$$

$$L \rightarrow id$$

$$I_5$$
: $L \rightarrow id$

$$l_6: S \rightarrow L = R$$

$$R \rightarrow L$$

$$L \rightarrow R$$

$$L \rightarrow R$$

$$L \rightarrow R$$

$$l_2: L \rightarrow *R$$

$$l_8: R \rightarrow L$$

$$I_9: S \rightarrow L = R$$





- 考虑项目集I₂。该集合的第一项使得action[2,=]为 "移动6进栈",因为FOLLOW(R)包含=。而第二项使 得action[2,=]为"归约R-> L"。于是,action[2,=]有 多重定义,存在移动归约冲突。但文法(4-20)不是 二义性的。
- 部分这类有冲突的无二义性文法可以用后面要讨论的规范LR文法和LALR文法来消除。



4.7 更强大的LR语法分析器

构造LR语法分析表的最一般技术: 在SLR表中

- 如果项目集I_i包含项目[A->α·]并且a在FOLLOW(A)中, 则状态i调用A->α归约。
- 然而,在某些情况下,当状态i出现在栈顶时,栈内的可行前缀βα使得在任何右句型中,βA不能跟随a。因此在这种情况下,用A->α进行归约无效。



让我们重新考虑下面的文法:

$$S -> L = R$$

 $S -> R$
 $L -> *R$ (4-20)
 $L -> id$
 $R -> I$

文法(4-20)的规范LR(0)项目集 族如图4-37所示。

- 状态2中项目R->L·对应于 上面的A->α, a对应于符号=, =在FOLLOW(R)中。于是, 当SLR语法分析器处于状态2而且 下一个输入符号为=时, SLR语法 分析器调用R->L归约。然而, 例 4.39的文法没有以R=···开始的右句 型。因此, 仅于活前缀L相对应的状态2不该进行归约。

$$I_0: S' \rightarrow \cdot S$$

$$S \rightarrow \cdot L = R$$

$$S \rightarrow \cdot R$$

$$L \rightarrow \cdot *R$$

$$L \rightarrow \cdot *id$$

$$R \rightarrow \cdot L$$

$$I_1: S' \to S$$

$$I_2: S \rightarrow L \cdot = R$$
 $R \rightarrow L$

$$I_3: S \rightarrow R$$

$$I_4: \quad L \to *\cdot R$$

$$R \to \cdot L$$

$$L \to \cdot *R$$

$$L \to \cdot id$$

$$I_5$$
: $L \rightarrow id$

$$l_6: S \rightarrow L = R$$

$$R \rightarrow L$$

$$L \rightarrow R$$

$$L \rightarrow \mathbf{id}$$

$$l_2: L \rightarrow *R$$

$$l_8: R \rightarrow L$$

$$I_9: S \rightarrow L = R$$

Fig. 4.37. Canonical LR(0) collection for grammar (4.20).



4.7.1 规范LR(1)项目



- 通过重新定义项目,使之包含一个终结符作为第二个分量,可以把更多的信息并入状态中。项目的一般形式也就变成了[A->α·β,a],其中A->αβ是产生式,a是终结符或\$。这样的对象叫做LR(1)项目。1是第二个分量的长度,这个分量叫做项目的搜索符。
- 形式地,我们说LR(1)项目[A->α·β,a]对活前缀γ有效, 如果存在推导S \rightarrow δAW $\stackrel{\sim}{\Longrightarrow}$ δαβW, 其中:
 - 1, γ=δα
 - 2、a是w的第一个符号,或者w是E且a是\$。





▶ 让我们考虑如下的文法:

S -> BB

 $B \rightarrow aB \mid b$

它有一个最右推导S aaBab aaaBab。在上面的定义中,令 δ =aa,A=B,w=ab, α =a, β =B,我们可以看到,项目[B->a·B,a]对活前缀 γ =aaa是有效的。这个文法的另一个最右推导是S BaB BaaB。从这个推导可以看出,项目[B->a·B,\$]对于活前缀Baa是有效的。



4.7.2 LR(1)项目集的构造算法

- **输入**: 拓广文法G'。
- 输出: LR(1)项目集, 它们是对G'的一个 或多个活前缀有效 的项目集。
- 方法:构造项目集的过程closure和goto及主例程items如图4-38所示。

```
function closure (i);
begin
       repeat
             for each item (A \rightarrow \alpha \cdot B \beta, a) in I,
                   each production B \rightarrow \gamma in G',
                   and each terminal b in FIRST(\beta a)
                   such that \{B \rightarrow y, b\} is not in I do
                         add (B \rightarrow \gamma, b) to I;
       until no more items can be added to I:
        return i
 end:
 function goto(1, X);
 begin
        let J be the set of items [A \rightarrow \alpha X \cdot \beta, a] such that
              [A \rightarrow \alpha \cdot X\beta, a] is in I;
        return closure (J)
 end:
 procedure items(G');
 begin
        C := \{closure(\{[S' \rightarrow \cdot S, \$]\})\};
         repeat
              for each set of items I in C and each grammar symbol X
                    such that goto(I, X) is not empty and not in C do
                          add goto(I, X) to C
         until no more sets of items can be added to C
  end
```

Fig. 4.38. Sets of LR(1) items construction for grammar G'.





▶ 考虑下面的拓广文法:

(4-21)

 $C \rightarrow cC \mid d$

首先计算{[S'->S·,\$]}的闭包,我们用项目[S'->S·,\$]来匹配过程closure中的项目[A->αB·β,a],即令A=S',α=€,B=S,β=€,a=\$。在函数closure中,要求对每个产生式B->γ和FIRST(βa)的每个终结符b,把[B->·γ,b]加到闭包中。根据当前文法,B->·γ只能是S->CC,而且因为β=€,a=\$,b也只能是\$,因此将[S->·CC,\$]加入到闭包中。





接下来,对于FIRST(C\$)中的b,将项目[C->·γ,b]加入闭包。
 也就是说,用项目[S->·CC,\$]匹配项目[A->αB·β,a],
 其中A=S,α=€,B=C,β=C,a=\$。

因为C不会推导出空串,所以FIRST(C\$)=FIRST(C)。

又因为FIRST(C)包含终结符c和d, 所以将以下项目加入闭包:

 $[C->\cdot cC,c],\quad [C->\cdot cC,d],\quad [C->\cdot d,c],\quad [C->\cdot d,d]_\circ$

因为再没有紧跟在点右面的非终结符,所以我们完成了对一个LR(1)项目集的计算。这个初始项目集为:

$$I_0: S' \rightarrow \cdot S, \$$$

$$S \rightarrow \cdot CC, \$$$

$$C \rightarrow \cdot cC, c/d$$

$$C \rightarrow \cdot d, c/d$$



- 为方便起见,我们 把LR(1)项目的方括号省略了, 而且用[C->·cC,c|d]表示 [C->·cC,c]和[C->·cC,d]的缩写。
- 现在对不同的X值计算 goto(I₀,X)。对X=S,因为点在 右端,项目集中只有项目 [S'->S,\$],因此第二个项目集 为: I₁: S'-> S·,\$ 以此类推,得到图4-39所示的。

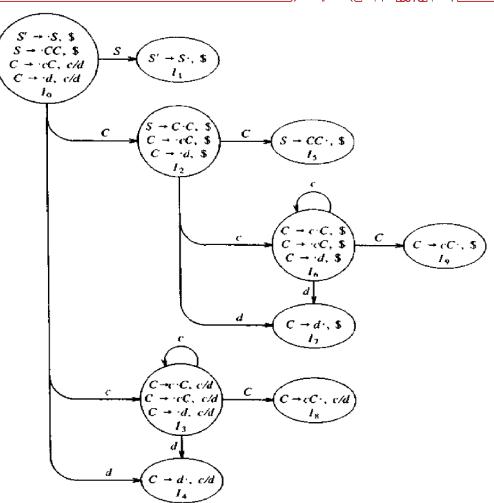


Fig. 4.39. The goto graph for grammar (4.21).



4.7.3 规范LR(1)语法分析表的构造算法

- *输入*: 拓广文法G'。
- 输出: 文法G'的规范LR语法分析表函数action和goto。
- 方法:
- 1、构造G'的LR(1)项目集规范族C= $\{I_0, I_1, \dots, I_n\}$ 。
- 2、从 构造语法分析器的状态i, 状态i的分析动作确定如下:
 - (a) 如果[A->α·aβ,b]在 I_i 中,并且goto(I_i ,a)= I_j , 则置action[i,a]为 s_i ,即"移动j进栈",这里a必须是终结符。
 - (b) 如果[A->α·,a]在I_i中且A≠S', 则置action[i,a]为r_i,即按r_i归约,其中j是产生式A->α的序号。
 - (c) 如果[S'->S·,\$]在 I_i 中,则置action[I_i ,\$]=acc,表示"接受"。 如果由上面的规则产生的动作有冲突,则G不是LR(1)文法,就构造不出语法分析器。
- 3、对所有的非终结符A,使用下面的规则构造状态i的goto函数: 如果goto(I_i ,A)= I_i ,则goto[i,A]=i。
- 4、不能由规则(2)和(3)定义的表项都置为"出错"。
- 5、语法分析器的初始状态是从包含[S'->·S,\$]的项目集构造出的状态。



文法: S'->S

 $S \rightarrow CC$

 $C \rightarrow cC \mid d$

的规范LR语法分析表如 图4-40所示,产生式1、2 和3分别是S->CC、C->cC 和C->d。

■ 注意:每个SLR(1)文法都是LR(1)文法,但对于SLR(1)文法,规范LR(1)文法可能具有更多的状态。上面的例子中的文法是SLR文法,它的SLR语法分析器只有7个状态,而图4-40却有10个状态。

CT + TE		action				
STATE	C	d	\$	S	C	
0	s3	s4		1	2	
1			acc			
2	s6	s7			5	
3	s3	84			8	
4	т3	τ3		}		
5			rl			
6	s6	s7			9	
7			r3	1		
8	r2	r2				
9			r2	ł		

Fig. 4.40. Canonical parsing table for grammar (4.21).



4.7.4 构造LALR语法分析表

- 就语法分析器大小而言, SLR表和LALR表对同一个文法 具有同样多的状态。构造SLR表和LALR表比构造规范LR表 要经济得多。
- 我们可以寻找同心的(即第一分量相同)LR(1)项目集, 并把这些同心的项目集合并成一个项目集。例如,图4-39 中,I₄和I₇、I₃和I₆以及I₈和I₉。注意,一般而言,心是相 应文法的一个LR(0)项目集;另外,LR(1)文法可能产生多个 同心的项目集合。
- 同心集的合并不会引起新的移动-归约冲突;但可能产生新的归约-归约冲突。



考虑如下文法:

 $S' \rightarrow S$

S -> aAd | bBd | aBe | bAe

 \forall -> C

B -> c

它只产生四个串: acd, ace, bcd和bce。通过构造 该文法的LR(1)项目集,可以看出没有冲突,它是LR(1)文法。 在它的项目集中,对活前缀ac有效的项目集为:

{[A->c·,d], [B->c·,e]}, 对bc有效的项目集为:

 $\{[A->c\cdot,e], [B->c\cdot,d]\},$

这两个集合都没有产生冲突且是同心的,然而他们合并后

 $A \rightarrow c$, de

 $B \rightarrow c$, de

产生归约-归约冲突。



一个简易但耗空间的LALR表构造算法

- *输入*: 拓广文法G'。
- 输出: G'的LALR语法分析表的action函数和goto函数。
- 方法:
- 1、构造文法的LR(1)项目集规范族 $C=\{I_0,I_1,\cdots,I_n\}$ 。
- 2、对出现在LR(1)项目集中的每个心,找出所有与之同心的项目集,用它们的并集代替它们。
- 3、令C'={J₁,J₂,···,J_n}是合并后的LR(1)项目集族。按照规范LR语法分析表的构造算法的方式从J_i来构造状态i的动作。如果分析动作出现冲突,算法无法产生分析表,说明该文法不是LR(1)文法。
- 4、goto表的构造如下:

如果J是一个或多个LR(1)项目集的并,即J= I_1 U I_2 U···U I_k ,那么goto(I_1 ,X), goto(I_2 ,X),···,goto(I_k ,X)也同心。 记K为所有与goto(I_1 ,X)同心的项目集的并,则goto(I_k ,X)=K。



再次考虑文法文法: S'->S 它的转移图如图 4-39所示。

$$C \rightarrow cC \mid d$$

正如我们之前提到的,有3对项目集可以合并 I_4 和 I_7 、 I_3 和 I_6 、 I_8 和 I_9 、结果如图4-41所示。

Sm. mr		golo			
STATE	c	đ	\$	S	C
0	s36	s47		L	2
1			acc	ì	
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3	ነ	
5			гÌ		
89	г2	г2	τ2		

G		action				
STATE	Ċ	d	5	S	c	
0	s3	s4	,	1	2	
- 1			acc			
2	s6	s7			5	
3	s3	84			8	
4	т3	t3		}		
5			rl			
6	s6	s7			9	
7			r3			
8	r2	r2				
9			r2	ł		

Fig. 4.41. LALR parsing table for grammar (4.21).

Fig. 4.40. Canonical parsing table for grammar (4.21).



- 当输入串存在错误时,LALR语法分析器可能比LR语法分析器多做一些不必要的归约,而LR语法分析器则能立即报错。但LALR不会比LR多移进一个符号。
- 例如,若输入串是ccd(针对文法(4-21))并跟随以\$时,图4-34的LR语法分析器将把0c3c3d4压入栈,并在状态4发现错误,因为状态4面临\$的动作是"出错"。

■ 对于同一输入串,LALR语法分析器将产 生相应的动作,即把

0 c 36 c 36 d 47

压入栈。但状态47面临\$的动作是 归约C->d, 栈的内容成为 0 c 36 c 36 C 89

现在状态89面临\$的动作是归约C->cC,

这时栈的内容改为

0 c 36 C 89

在经一次类似的归约,获得的栈为 0C2

最后, 状态2面临\$的动作是"出错"。

作业



• (1) 构造下面拓广文法的LR(0)和SLR(1)分析表:

$$C \rightarrow cC \mid d$$

■ (2) 构造下面拓广文法的规范LR(1)分析表:



作业答案



(1) 构造下面拓广文法的LR(0) 和SLR (1) 分析表:

S -> CC

C -> cC | d

• 对各条产生式编号:

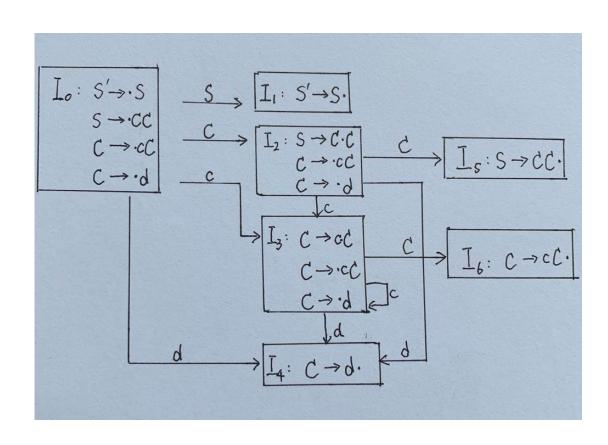
0: S' -> S

1: S -> CC

2: C -> cC

3: C -> d

 然后,构造给定文法的 LR(0) 自 动机,结果如图所示:







■ 则 LR(0) 分析表结果如图所示:

状态	ACTION			G	ОТО
	С	d	\$	S	С
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		





- 在此基础上,构造 SLR 分析表,则还需要计算文法的 Follow 集:
- Follow(S')={\$}, Follow(S)={\$}, Follow(C)={\$,c,d}
- 则 SLR (1) 分析表为:

状态	ACTION			GO	ОТО
	С	d	\$	S	С
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5			r1		
6	r2	r2	r2		





(2) 构造下面拓广文法的

规范LR(1)分析表: S->AS|b

A->SA|a

对各条产生式编号:

0: S' -> S

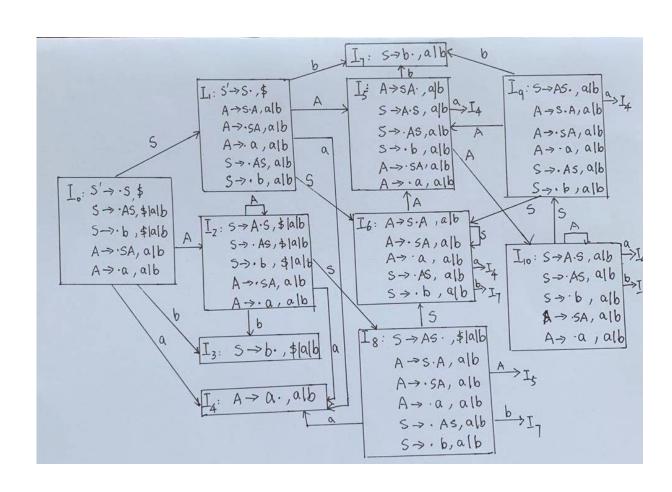
1: S -> AS

2: S -> b

3: A -> SA

4: A -> a

- 先计算文法的 First 集合,
- 结果为: First (A)={a,b), First(S)={a,b}
- 然后,构造给定文法的 LR(1) 自动机,计算过程用到了 First 集合,最终结果如图所示:







- 则 LR(1) 分析 表结果如图所 示:
- 从表中可以看到,存在冲突项,所以该文法不是LR(1)的。

状态	A	CTION	GO	ТО	
	а	b	\$	S	A
0	s4	s3		1	2
1	s4	s7	acc	6	5
2	s4	s3		8	2
3	r2	r2	r2		
4	r4	r4			
5	S4/r3	s7/r3		9	10
6	s4	s7		6	5
7	r2	r2			
8	S4/r1	s7/r1	r1	6	5
9	S4/r1	s7/r1		6	5
10	s4	s7		9	10



附加: 算符优先分析法



一、算符优先关系表

■ 算符文法:上下文无关文法G,如果它不含 ε - 产生式,并且产生式右部均不含相连的非终结符,即没有 $P \to \varepsilon$,或 $P \to \cdots QR \cdots$ 形式的产生式,则G是一个算符文法。

例:
$$E \rightarrow E + T | T$$
 $T \rightarrow T * F | F$
 $F \rightarrow (E) | i$

单击此处编辑母版标题样式

优先关系定义: 算符文法G,对任何 $a,b \in V_T$,定义a与b之间的优先关系如下:

- (1) a = b, G中有 $P \rightarrow \cdots ab \cdots$ 或 $P \rightarrow \cdots aQb \cdots$ 形式的产生式。
- (2) a \lessdot b, G中有 $P \rightarrow \cdots aQ \cdots$ 的产生式,且 $Q \Rightarrow b \cdots$ 或 $Q \Rightarrow Rb \cdots$ 。
- (3) a > b, G中有 $P \rightarrow \cdots Qb \cdots$ 的产生式,且 $Q \stackrel{\text{\tiny T}}{\Rightarrow} \cdots a$ 或 $Q \stackrel{\text{\tiny T}}{\Rightarrow} \cdots aR$ 。

对
$$E \rightarrow E + T$$
,

由
$$T \stackrel{\tau}{\Rightarrow} i$$
 得 $+ \stackrel{\epsilon}{\circ} i$

由
$$T \stackrel{+}{\Rightarrow} T * F$$
 得 $+ \lessdot *$

由
$$T \stackrel{+}{\Rightarrow} (E)$$
 得 $+ \lessdot ($

由
$$E \Rightarrow i$$
 得 $i > +$

由
$$E \stackrel{+}{\Rightarrow} (E)$$
 得 $) > +$

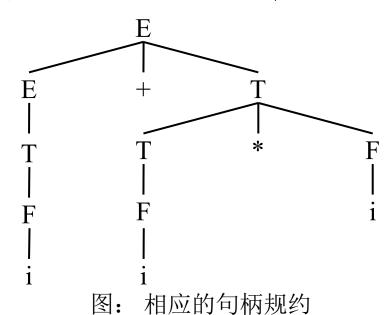
由
$$E \Rightarrow E + T$$
 得 $+ > +$

	+	*	i	()	\$
+	>	<	<	<	>	>
+ *	> > >	>	<			>
i	>	>			>	>
(<	<	<	<	<u> </u>	
	> <	> <			>	>
) \$	<	<	<	<		<u>o</u>



二、算符优先分析方法

这也是一种表驱动的分析方法,除了优先 关系表作为分析表外,还需要一个分析栈 和一个分析控制程序,协同地进行语法分 析。文法句子i+i*i的LR分析过程如下:



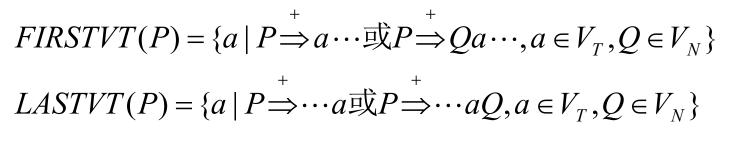


例子文法句子i+i*i的算法优先分析过程

分析栈	输入串	动作
\$	i+i*i\$	\$≪i
\$i	+i*i\$	i > + 归约
\$F	+i*i\$	\$ < +
\$F+	i*i\$	+ < i
\$F+i	*i\$	i > * 归约
\$F+F	*i\$	+ < *
\$F+F*	i\$	* < i
\$F+F*i	\$	i > \$ 归约
\$F+F*F	\$	* > \$ 归约
\$F+T	\$	+ > \$ 归约
\$E	\$	结束



三、优先关系表的构造



构造FIRSTVT(P)可遵照下列两条规则:

- 若有产生式 $P \rightarrow a \cdots$ 或 $P \rightarrow Qa \cdots$, 则 $a \in FIRSTVT(P)$;
- 若有产生式 $P \to Q \cdots$ 则 $FIRSTVT(Q) \subseteq FIRSTVT(P)$ 。 构造LASTVT(P)有类似的两条规则:
- 若有产生式 $P \rightarrow \cdots a$ 或 $P \rightarrow \cdots aQ$, 则 $a \in LASTVT(P)$;
- 若有产生式 $P \rightarrow \cdots Q$ 则 $LASTVT(Q) \subseteq LASTVT(P)$ 。



例



右面文法的FIRTVST集和LASTVT集, 是这样得到的:

 $\pm F \rightarrow (E)|i$,

 $FIRSTVT(F) = \{(, i), LASTVT(F) = \{(, i), i\}\}$

由T→T*F|F,

FIRSTVT(T)={*, (, i }, LASTVT(T)={*,) , i }

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

FIRSTVT(E)={+, *, (, i}, LASTVT(E)={+, *,), i}

参照前节FIRST集的成员表表示法,本例的 FIRST集、LASTVT集也可用成员表表示,其

中用 √表示FIRST集的成员,○表示LASTVT

集的成员,则有下图:

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

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

	+	*	()	i
$E \rightarrow E + T \mid T$	√0	O	4	0	√ O
T→T*F F		√0	4	0	√ O
F→(E) i			1	0	√ O





于是,由
$$E \rightarrow E + T$$

$$\pm T \to T * F$$

$$\pm F \rightarrow (E)$$

	+	*	i	()	\$
+	>	<	<	<	>	*
*	>		<<		>	>
i	>	>			>	>
(<	<	<	<	<u> </u>	
	>	>			>	>
) \$	<	<	<	<		0



大作业



- 测试文法: 1. $E \rightarrow E + T \mid T$ $T \rightarrow T * F \mid F$ $F \rightarrow (E) \mid i$
 - 2. E -> E + E | E * E | (E) | id
- 提交方式:canvas平台
- 提交时间:2020年6月21日



本章小结



- 语法分析器
- 上下文无关文法和文法设计,
- 推导、分析树和二义文法
- LL(1)语法分析器
- SLR语法分析器
- 规范LR语法分析器和LL(1)语法分析器
- 算符优先法

谢谢!

