

第六讲 自顶向下语法分析

2021-10

1. 基本思想

对于给定语言的文法和一个单词符号串（终结符串），一般的自顶向下分析过程是：从文法开始符号进行推导，每一步推导都获得文法的一个句型，直到产生一个句子，恰好是所期望的单词符号串。每一步推导是对当前句型中剩余的某个非终结符进行展开，即使用以该非终结符为左部的某个产生式的右部替换该非终结符。如果不存在这样的推导，则表明该单词符号串存在语法错误。

例如，给定如下文法 $G[S]$ ：

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aA \mid \varepsilon \\ B &\rightarrow b \mid bB \end{aligned}$$

针对单词符号串 $aaab$ 的一个自顶向下分析过程为：

$$\begin{aligned} S & \quad // \text{使用产生式 } S \rightarrow AB \\ \Rightarrow AB & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aAB & \quad // \text{使用产生式 } B \rightarrow b \\ \Rightarrow aAb & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aaAb & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aaaAb & \quad // \text{使用产生式 } A \rightarrow \varepsilon \\ \Rightarrow aaab & \end{aligned}$$

2. 带回溯的自顶向下分析

一般的自顶向下分析过程存在两类非确定性。其一是在每一步推导中，选择对哪一个非终结符进行展开。其二是如果选定的非终结符是多个产生式的左部，那么应该选择使用哪一个产生式。这种非确定性导致推导过程需要不断地进行试探和回溯。例如，对于文法 $G[S]$ 单词符号串 $aaab$ ，以下是一个试探的推导过程：

$$\begin{aligned} S & \quad // \text{使用产生式 } S \rightarrow AB \\ \Rightarrow AB & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aAB & \quad // \text{使用产生式 } B \rightarrow bB \\ \Rightarrow aAbB & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aaAbB & \quad // \text{使用产生式 } A \rightarrow aA \\ \Rightarrow aaaAbB & \quad // \text{使用产生式 } A \rightarrow \varepsilon \\ \Rightarrow aaabB & \quad // \text{必须回溯} \end{aligned}$$

回溯会带来很高的复杂度，这在编译程序的设计中是不现实的。试想一个实际中需要编

译的程序单位可能会包含多少个单词符号，就会得出这样的结论。那么在自顶向下分析过程中，如何避免回溯呢？解决办法就是消除上述两类非确定性。

很容易想到，如果我们只允许最左推导或最右推导，那么就可以避开上述的第一类非确定性。由于通常都是从左到右读入单词符号串，所以我们可以规定在自顶向下分析过程只使用最左推导。例如，对于文法 $G[S]$ 单词符号串 $aaab$ ，以下试探推导过程的每一步总是对最左边的非终结符进行展开：

S	// 使用产生式 $S \rightarrow AB$
$\Rightarrow AB$	// 使用产生式 $A \rightarrow aA$
$\Rightarrow aAB$	// 使用产生式 $A \rightarrow aA$
$\Rightarrow aaAB$	// 使用产生式 $A \rightarrow \varepsilon$
$\Rightarrow aaB$	// 必须回溯

这个试探推导过程虽然出现了必须回溯的情形，但我们能够确信这是由于第二类非确定性造成的。也就是说，在前面的推导步骤中选错了产生式。

3. 自顶向下预测分析

如果在自顶向下分析过程中，我们能够同时消除两类非确定性，即能够保证非终结符选择和产生式选择都是确定的，那么这样的分析过程就是**确定的自顶向下分析**。

为了解决消除第二类非确定性的问题，通常采取的策略是向前查看确定数目的单词符号，然后确定应该选择哪一个产生式进行最左推导。这种方法称为**自顶向下预测分析**。成功分析的结果是一个唯一的最左推导。

我们先来看一个可以进行自顶向下预测分析的例子。

对于本讲前面所讨论的文法 $G[S]$ ，以及任意的单词符号序列 $a^n b^m$ ($n \geq 0, m > 0$)，总存在如下的最左推导：

S	// 使用产生式 $S \rightarrow AB$
$\Rightarrow AB$	// 使用产生式 $A \rightarrow aA$
$\Rightarrow aAB$	// 使用产生式 $A \rightarrow aA$
.....	
$\Rightarrow a^n AB$	// 使用产生式 $A \rightarrow \varepsilon$
$\Rightarrow a^n B$	// 使用产生式 $B \rightarrow bB$
$\Rightarrow a^n bB$	// 使用产生式 $B \rightarrow bB$
.....	
$\Rightarrow a^n b^{m-1} B$	// 使用产生式 $B \rightarrow b$
$\Rightarrow a^n b^m$	

不难看出，对于给定的 $a^n b^m$ 来说，这是唯一的最左推导。然而，问题是当最左的非终结符为 A 时，我们如何在产生式 $A \rightarrow aA$ 和 $A \rightarrow \varepsilon$ 之间做出选择；同样，当最左的非终结符为 B 时，我们如何在产生式 $B \rightarrow bB$ 和 $B \rightarrow b$ 之间做出选择。

我们来分析一下向前查看确定数目的单词符号进行预测分析的可能性。首先，在需要展开的最左非终结符为 A 时，我们通过向前查看 1 个单词符号，就可以在产生式 $A \rightarrow aA$ 和 $A \rightarrow \varepsilon$ 之间作出选择。当查看到下一个单词符号为 a 时，选择 $A \rightarrow aA$ ；而当查看到下一

个单词符号为 b 时, 则选择 $A \rightarrow \varepsilon$ 。然而, 在需要展开的最左非终结符为 B 时, 当查看到下一个单词符号为 b 时, 并不足以在 $B \rightarrow bB$ 和 $B \rightarrow b$ 之间做出选择。我们将向前查看单词符号的数目变为 2。通过分析, 当后两个单词符号为 bb 时, 我们应该选择 $B \rightarrow bB$; 而当后两个单词符号中只有一个 b 时(后跟一个单词符号序列的结束符号), 应该选择 $B \rightarrow b$ 。

由于 $L(G[S]) = \{ a^n b^m \mid n \geq 0, m > 0 \}$, 所以我们可以得出结论: 只要向前查看 2 个单词符号, 就可预测分析 $L(G[S])$ 中的所有句子。

我们先来看一个不能进行自顶向下预测分析的例子。设有如下文法 $G'[S]$:

$$\begin{aligned} S &\rightarrow Sa \\ S &\rightarrow b \end{aligned}$$

对于 $L(G'[S])$ 中的任意的单词符号序列 ba^n ($n \geq 0$), 总存在如下的最左推导:

$$\begin{aligned} S & & // \text{使用产生式 } S \rightarrow Sa \\ \Rightarrow Sa & & // \text{使用产生式 } S \rightarrow Sa \\ \Rightarrow Saa & & // \text{使用产生式 } S \rightarrow Sa \\ \dots\dots & & \\ \Rightarrow San & & // \text{使用产生式 } S \rightarrow b \\ \Rightarrow ba^n & & \end{aligned}$$

然而, 在第 1 步推导时, 在向前查看到第 2 个单词符号为 a 时, 我们才会选择产生式 $S \rightarrow Sa$; 在第 2 步推导时, 在向前查看到第 3 个单词符号为 a 时, 我们才会选择产生式 $S \rightarrow Sa$; 这样, 在第 n 步推导时, 我们就需要向前查看 $n+1$ 个单词。这样, 无论向前查看单词符号的数确定为多少, 都无法满足对 $L(G)$ 中所有句子进行预测分析的需求。

可见, 不是所有文法都是可以成功实施自顶向下预测分析的。在实践中, 通常我们可以对文法的设计进行适当限制, 以满足预测分析的要求。比如, 我们将在下一节讨论的 LL(1) 文法是一种满足这种要求的文法, 并且足以满足多数程序设计语言的文法描述需求。

4. LL(1) 分析

LL(1) 分析是应用较普遍的一种自顶向下预测分析方法。

LL(1) 中的第一个“L”代表从左 (Left) 向右扫描单词符号, 第二个“L”代表产生的是最左 (Leftmost) 推导, “1”代表向前查看 (lookahead) 一个单词符号。

LL(1) 分析方法仅适用于 LL(1) 文法。为给出 LL(1) 文法的定义, 我们首先介绍两个重要概念: First 集合和 Follow 集合。这些概念在介绍自底向上分析方法时也要用到。

4.1 First 集合和 Follow 集合

我们先来看 First 集合的定义:

设上下文无关文法 $G = (V_N, V_T, P, S)$ 。对 $\alpha \in (V_N \cup V_T)^*$,

$$\text{First}(\alpha) = \{ a \mid \alpha \Rightarrow^* a\beta, a \in V_T, \beta \in (V_N \cup V_T)^*, \text{或者 } \alpha \Rightarrow^* \varepsilon \text{ 时 } a = \varepsilon \}$$

直观来理解，一个句型 α 若可以推导出另一个以终结符 a 开头的句型，那么 a 属于 $\text{First}(\alpha)$ ；若 α 可以推导出 ε ，那么 ε 属于 $\text{First}(\alpha)$ 。

由于任何推导都可以对应一个最左推导（可以参考2.2.5的结果直接得出该结论），所以 First 集合的定义也可以基于最左推导出：

$$\text{First}(\alpha) = \{a \mid \alpha \Rightarrow_{lm}^* a\beta, a \in V_T, \beta \in (V_N \cup V_T)^*, \text{ 或者 } \alpha \Rightarrow_{lm}^* \varepsilon \text{ 时 } a = \varepsilon\}$$

在实际应用中，对于给定文法 $G = (V_N, V_T, P, S)$ ，我们通常会关心下列集合 X_G 中句型的 First 集合：

$$X_G = V_N \cup V_T \cup \{\varepsilon\} \cup \{v \mid A \rightarrow u \in P, \text{ 且 } v \text{ 是 } u \text{ 的后缀}\}$$

对所有 $x \in X_G$ ，可以通过如下过程计算 $\text{First}(x)$ ：

- 初始时，对 $x \in V_T \cup \{\varepsilon\}$ ，置 $\text{First}(x) = \{x\}$ ；对其它 x ，置 $\text{First}(x) = \Phi$ ；
- 重复如下步骤，直到所有 First 集合没有变化为止：
 - (1) 对于 $y_1y_2\dots y_k \in \{v \mid A \rightarrow u \in P, \text{ 且 } v \text{ 是 } u \text{ 的后缀}\}$ ，其中 $k \geq 1$ ， $y_j \in V_N \cup V_T$ ($1 \leq j \leq k$)，若 $\forall j: 1 \leq j \leq i-1 (\varepsilon \in \text{First}(y_j)) \wedge \varepsilon \notin \text{First}(y_i)$ ，其中 $1 \leq i \leq k$ ，则令

$$\text{First}(y_1y_2\dots y_k) = \text{First}(y_1) \cup \text{First}(y_2) \cup \dots \cup \text{First}(y_i) - \{\varepsilon\}$$
 否则，若 $\forall j: 1 \leq j \leq k (\varepsilon \in \text{First}(y_j))$ ，则令

$$\text{First}(y_1y_2\dots y_k) = \text{First}(y_1) \cup \text{First}(y_2) \cup \dots \cup \text{First}(y_k)$$
 - (2) 若有 $A \rightarrow y_1y_2\dots y_k \in P$ ，则置 $\text{First}(A) = \text{First}(A) \cup \text{First}(y_1y_2\dots y_k)$ 。

例 1 设上下文无关文法 $G = (V_N, V_T, P, S)$ ，其中 P 为

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow Da \mid \varepsilon \\ B &\rightarrow cC \\ C &\rightarrow aADC \mid \varepsilon \\ D &\rightarrow b \mid \varepsilon \end{aligned}$$

试计算下列集合 X_G 中句型的 First 集合：

$$X_G = V_N \cup V_T \cup \{\varepsilon\} \cup \{v \mid A \rightarrow u \in P, \text{ 且 } v \text{ 是 } u \text{ 的后缀}\}$$

解 集合 X_G 为：

$$X_G = \{\varepsilon, S, A, B, C, D, a, b, c, AB, Da, cC, aADC, ADC, DC\}$$

初始时，置 $\text{First}(a) = \{a\}$ ， $\text{First}(b) = \{b\}$ ， $\text{First}(c) = \{c\}$ ， $\text{First}(\varepsilon) = \{\varepsilon\}$ ，而对其它的 $x \in X_G$ ，置 $\text{First}(x) = \Phi$ 。

第一轮应用以上三个步骤之后， X_G 中句型的 First 集合更新情况为：

$$\begin{aligned} \text{First}(cC) &= \{c\} \\ \text{First}(aADC) &= \{a\} \\ \text{First}(A) &= \{\varepsilon\} \\ \text{First}(B) &= \{c\} \end{aligned}$$

$$\text{First}(C) = \{\varepsilon, a\}$$

$$\text{First}(D) = \{\varepsilon, b\}$$

第二轮应用以上三个步骤之后， X_G 中句型的 **First** 集合更新情况为：

$$\text{First}(AB) = \{c\}$$

$$\text{First}(Da) = \{b, a\}$$

$$\text{First}(ADC) = \{\varepsilon, b, a\}$$

$$\text{First}(DC) = \{\varepsilon, b, a\}$$

$$\text{First}(S) = \{c\}$$

$$\text{First}(A) = \{\varepsilon, b, a\}$$

第三轮应用以上三个步骤之后， X_G 中句型的 **First** 集合更新情况为：

$$\text{First}(AB) = \{b, a, c\}$$

$$\text{First}(S) = \{b, a, c\}$$

$$\text{First}(C) = \{\varepsilon, a\}$$

第四轮应用以上三个步骤之后， X_G 中句型的 **First** 集合无变化。我们将最后结果整理如下：

$$\text{First}(\varepsilon) = \{\varepsilon\}$$

$$\text{First}(S) = \{b, a, c\}$$

$$\text{First}(A) = \{\varepsilon, b, a\}$$

$$\text{First}(B) = \{c\}$$

$$\text{First}(C) = \{\varepsilon, a\}$$

$$\text{First}(D) = \{\varepsilon, b\}$$

$$\text{First}(a) = \{a\}$$

$$\text{First}(b) = \{b\}$$

$$\text{First}(c) = \{c\}$$

$$\text{First}(AB) = \{b, a, c\}$$

$$\text{First}(Da) = \{b, a\}$$

$$\text{First}(cC) = \{c\}$$

$$\text{First}(aADC) = \{a\}$$

$$\text{First}(ADC) = \{\varepsilon, b, a\}$$

$$\text{First}(DC) = \{\varepsilon, b, a\}$$

我们再来看 **Follow 集合** 的定义：

设上下文无关文法 $G = (V_N, V_T, P, S)$ 。对每个 $A \in V_N$,

$$\text{Follow}(A) = \{ a \mid S \Rightarrow^* \alpha A \beta \# \text{ 且 } a \in \text{First}(\beta \#), \alpha, \beta \in (V_N \cup V_T)^* \}$$

其中， $\#$ 代表输入单词符号序列的结束符。

显然，一定有 $\# \in \text{Follow}(S)$ 。

直观来理解，若 G 中存在一个包含子串 Xa 的句型，那么 a 属于 $\text{Follow}(X)$ ； G 中存在一个以 X 结尾的句型，那么 $\#$ 属于 $\text{Follow}(X)$ 。

对所有 $A \in V_N$ ，我们可以通过如下过程计算 $\text{Follow}(A)$ ：

- 初始时，置 $\text{Follow}(S) = \{\#\}$ ；对其它 $A \in V_N$ ，置 $\text{Follow}(A) = \Phi$ 。
- 重复如下步骤，直到所有 Follow 集合不再变化为止：

若有 $A \rightarrow \alpha B \beta \in P$ ，其中 $\alpha, \beta \in (V_N \cup V_T)^*$ ，则

(1) 令 $\text{Follow}(B) = \text{Follow}(B) \cup (\text{First}(\beta) - \{\epsilon\})$ ；

(2) 若 $\epsilon \in \text{First}(\beta)$ ，则令 $\text{Follow}(B) = \text{Follow}(B) \cup \text{Follow}(A)$ 。

例 2 设上下文无关文法 $G = (V_N, V_T, P, S)$ ，其中 P 为

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow Da \mid \epsilon \\ B &\rightarrow cC \\ C &\rightarrow aADC \mid \epsilon \\ D &\rightarrow b \mid \epsilon \end{aligned}$$

试计算 V_N 中所有非终结符的 Follow 集合。

解 初始时，置 $\text{Follow}(S) = \{\#\}$ ；对其它 $X \in V_N$ ，置 $\text{Follow}(X) = \Phi$ 。

在例 4.1 的结果基础上，我们进行每一轮的迭代计算。

依次考虑每一个右端含有非终结符的产生式，第一轮计算为：

由于 $S \rightarrow AB \in P$ ，所以置 $\text{Follow}(A) = \text{Follow}(A) \cup (\text{First}(B) - \{\epsilon\}) = \{c\}$ ；

由于 $S \rightarrow AB \in P$ ， $\epsilon \in \text{First}(\epsilon)$ ，所以置 $\text{Follow}(B) = \text{Follow}(B) \cup \text{Follow}(S) = \{\#\}$ ；

由于 $A \rightarrow Da \in P$ ，所以置 $\text{Follow}(D) = \text{Follow}(D) \cup (\text{First}(a) - \{\epsilon\}) = \{a\}$ ；

由于 $B \rightarrow cC \in P$ ， $\epsilon \in \text{First}(\epsilon)$ ，所以置 $\text{Follow}(C) = \text{Follow}(C) \cup \text{Follow}(B) = \{\#\}$ ；

由于 $C \rightarrow aADC$ ，所以置 $\text{Follow}(A) = \text{Follow}(A) \cup (\text{First}(DC) - \{\epsilon\}) = \{c, b, a\}$ ；

由于 $C \rightarrow aADC$ ， $\epsilon \in \text{First}(DC)$ ，所以置 $\text{Follow}(A) = \text{Follow}(A) \cup \text{Follow}(C) = \{c, b, a, \#\}$ ；

由于 $C \rightarrow aADC$ ，所以置 $\text{Follow}(D) = \text{Follow}(D) \cup (\text{First}(C) - \{\epsilon\}) = \{a\}$ ；

由于 $C \rightarrow aADC$ ， $\epsilon \in \text{First}(C)$ ，所以置 $\text{Follow}(D) = \text{Follow}(D) \cup \text{Follow}(C) = \{a, \#\}$ ；

再依次考虑这些产生式，进行第二轮计算，结果无变化。我们将最后结果整理如下：

$$\begin{aligned} \text{Follow}(S) &= \{\#\} \\ \text{Follow}(A) &= \{c, b, a, \#\} \\ \text{Follow}(B) &= \{\#\} \\ \text{Follow}(C) &= \{\#\} \\ \text{Follow}(D) &= \{a, \#\} \end{aligned}$$

4.2 LL (1) 文法

下面我们来定义 **LL(1) 文法**。为方便，我们先引入预测集合的概念。

设上下文无关文法 $G = (V_N, V_T, P, S)$ 。对任何产生式 $A \rightarrow \alpha \in P$ ，其预测集合 $PS(A \rightarrow \alpha)$ 定义如下：

- 如果 $\varepsilon \notin \text{First}(\alpha)$ ，那么 $PS(A \rightarrow \alpha) = \text{First}(\alpha)$ ；
- 如果 $\varepsilon \in \text{First}(\alpha)$ ，那么 $PS(A \rightarrow \alpha) = (\text{First}(\alpha) - \{\varepsilon\}) \cup \text{Follow}(A)$

文法 G 是 LL(1) 的，当且仅当对于 G 中任何两个有相同左部的不同产生式 $A \rightarrow \alpha$ 和 $A \rightarrow \beta$ ，都满足： $PS(A \rightarrow \alpha) \cap PS(A \rightarrow \beta) = \Phi$ 。

例 3 设上下文无关文法 $G = (V_N, V_T, P, S)$ ，其中 P 为

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow Da \mid \varepsilon \\ B &\rightarrow cC \\ C &\rightarrow aADC \mid \varepsilon \\ D &\rightarrow b \mid \varepsilon \end{aligned}$$

试指出该文法是否 LL(1) 文法。

解 根据 LL(1) 文法的定义，我们只需要考虑产生式 $A \rightarrow Da$ 和 $A \rightarrow \varepsilon$ ， $C \rightarrow aADC$ 和 $C \rightarrow \varepsilon$ ，以及 $D \rightarrow b$ 和 $D \rightarrow \varepsilon$ 的预测集合之间是否相交。

先分别计算出这些产生式的预测集合：

$$\begin{aligned} PS(A \rightarrow Da) &= \{b, a\} \\ PS(A \rightarrow \varepsilon) &= \{c, b, a, \#\} \\ PS(C \rightarrow aADC) &= \{a\} \\ PS(C \rightarrow \varepsilon) &= \{\#\} \\ PS(D \rightarrow b) &= \{b\} \\ PS(D \rightarrow \varepsilon) &= \{a, \#\} \end{aligned}$$

可见： $PS(A \rightarrow Da) \cap PS(A \rightarrow \varepsilon) \neq \Phi$ 。因此，该文法不是 LL(1) 文法。

4.3 LL（1）分析的实现

在这一小节里，我们介绍常用的两种 LL（1）分析的实现方法。

4.3.1 递归下降 LL（1）分析程序

在递归下降 LL（1）分析程序的设计中，每个非终结符都对应一个分析子程序，分析程序从调用文法开始符号所对应的分析子程序开始执行。非终结符对应的分析子程序根据下一个单词符号可确定自顶向下分析过程中应该使用的产生式，根据所选定的产生式，分析子程序的行为依据产生式右端依次出现的符号来设计：

- 每遇到一个终结符，则判断当前读入的单词符号是否与该终结符相匹配，若匹配，则继续读取下一个单词；若不匹配，则进行错误处理。
- 每遇到一个非终结符，则调用相应的分析子程序。

例如，设有如下产生式：

$$\langle \text{function} \rangle \rightarrow \text{FUNC ID } (\langle \text{parameter_list} \rangle) \langle \text{statement} \rangle$$

其中， $\langle \text{function} \rangle$ ， $\langle \text{parameter_list} \rangle$ ，和 $\langle \text{statement} \rangle$ 是非终结符，而 FUNC 和 ID 是终结符。若它是左部为 $\langle \text{function} \rangle$ 的唯一产生式，那么非终结符 $\langle \text{function} \rangle$ 对应的分析子程序 ParseFunction() 的设计可描述为：

```
void ParseFunction()
{
    MatchToken(T_FUNC);      // T_FUNC 为匹配终结符 FUNC 的单词符号
    MatchToken(T_ID);        // T_ID 为匹配终结符 ID 的单词符号
    MatchToken(T_LPAREN);    // T_LPAREN 为匹配终结符 '(' 的单词符号
    ParseParameterList();
    MatchToken(T_RPAREN);    // T_RPAREN 为匹配终结符 ')' 的单词符号
    ParseStatement();
}
```

其中，ParseParameterList() 和 ParseStatement() 分别为非终结符 $\langle \text{parameter_list} \rangle$ 和 $\langle \text{statement} \rangle$ 对应的分析子程序，而函数 MatchToken() 则是用于匹配当前终结符和正在扫描的单词符号（若匹配，则调用词法分析程序取下一个单词符号；否则报告词法错误信息）。函数 MatchToken() 的一种简单的设计为：

```
void MatchToken(int expected)
{
    if (lookahead != expected)    // 判别当前扫描的单词符号是否与期望的终结符
匹配
    {
        printf("syntax error \n"); // 若不匹配，则报告出错信息，跳出
        exit(0);
    }
    else                          // 若匹配，消费掉当前单词符号
        lookahead = getToken();  // 并向词法分析程序申请并读入下一个单词符号
}
```

其中，lookahead 为全局量，存放当前扫描的单词符号。

为方便，我们在随后的讨论以及例子中，将继续使用这个 MatchToken() 函数。

一般情况下，设 LL(1) 文法中某一非终结符 A 对应的所有产生式的集合为：

$$A \rightarrow u_1 \mid u_2 \mid \dots \mid u_n$$

那么相对于非终结符 A 的分析子程序 ParseA() 可以具有如下形式的一般结构：

```
void ParseA()
{
    switch (lookahead) {
        case PS(A→u1):
            ..... /*根据 u1 设计的分析过程*/
    }
```



```

        break;
    case PS( $A \rightarrow u_2$ ):
        ..... /*根据  $u_2$  设计的分析过程*/
        break;
    ...
    case PS( $A \rightarrow u_n$ ):
        ..... /*根据  $u_n$  设计的分析过程*/
        break;
    default:
        printf("syntax error\n");
        exit(0);
}
}

```

值得注意的是：由于是 LL(1) 文法，所以产生式 $A \rightarrow u_1, A \rightarrow u_2, \dots$ ，以及 $A \rightarrow u_n$ 的预测集合是两两互不相交的，故上述选择语句中的各个选择之间是互斥的。

例 4 设文法 $G[S]$ 为

$$\begin{aligned}
 S &\rightarrow AaS \mid BbS \mid d \\
 A &\rightarrow a \\
 B &\rightarrow \varepsilon \mid c
 \end{aligned}$$

试验证该文法是 LL(1) 文法，并设计一个基于该文法的递归下降分析程序。

解 先验证该文法是 LL(1) 的。为此，我们先计算出所需要的 First 集合和 Follow 集合：

$$\begin{aligned}
 \text{First}(AaS) &= \{a\} \\
 \text{First}(BbS) &= \{c, b\} \\
 \text{First}(d) &= \{d\} \\
 \text{First}(a) &= \{a\} \\
 \text{First}(\varepsilon) &= \{\varepsilon\} \\
 \text{First}(c) &= \{c\} \\
 \text{Follow}(S) &= \{\#\} \\
 \text{Follow}(A) &= \{a\} \\
 \text{Follow}(B) &= \{b\}
 \end{aligned}$$

再计算出各产生式的预测集合：

$$\begin{aligned}
 PS(S \rightarrow AaS) &= \{a\} \\
 PS(S \rightarrow BbS) &= \{c, b\} \\
 PS(S \rightarrow d) &= \{d\} \\
 PS(A \rightarrow a) &= \{a\} \\
 PS(B \rightarrow \varepsilon) &= \{b\} \\
 PS(B \rightarrow c) &= \{c\}
 \end{aligned}$$

因为 $PS(S \rightarrow AaS)$ ， $PS(S \rightarrow BbS)$ 以及 $PS(S \rightarrow d)$ 互不相交， $PS(B \rightarrow \varepsilon)$ 和 $PS(S \rightarrow d)$ 不相交，所以， $G[S]$ 是 LL(1) 文法。

这样，开始符号 S 对应的分析子程序可以设计为：

```
void ParseS()  
{  
    switch (lookahead) {  
        case a:  
            ParseA();  
            MatchToken(a);  
            ParseS();  
            break;  
        case b,c:  
            ParseB();  
            MatchToken(b);  
            ParseS();  
            break;  
        case d:  
            MatchToken(d);  
            break;  
        default:  
            printf("syntax error \n")  
            exit(0);  
    }  
}
```

这里，我们假设了下一个单词符号与该文法的终结符号具有相同的表示形式，下同。

终结符 A 对应的分析子程序可以设计为：

```
void ParseA()  
{  
    if (lookahead==a) {  
        MatchToken(a);  
    }  
    else {  
        printf("syntax error \n");  
        exit(0);  
    }  
}
```

终结符 B 对应的分析子程序可以设计为：

```
void ParseB()  
{  
    if (lookahead==c) {  
        MatchToken(c);  
    }  
    else if (lookahead==b) {
```

```

    }
    else {
        printf("syntax error \n");
        exit(0);
    }
}

```

在实践中，这种递归下降分析程序的设计思想也常用于其它形式的文法描述。比如，在 PL/0 语法分析程序时采用了 EBNF 形式的语法描述。对于这种语法描述，我们将会得到更加简洁和高效的分析子程序，有利于语法分析程序的手工构造。

EBNF 形式的语法描述中，每一条规则右部的语法成分之间除了可以有连接算符之外，还包含其它一些算符，主要有选择、重复、任选以及优先括号等。这种更加丰富的表达方式有利于精简分析子程序的设计，从而提高递归下降分析程序的效率。在递归下降分析子程序的设计中，针对不同算符，可选择不同的处理语句：

- $X_1 | X_2 | \dots | X_m$ 表示多个成分之间的选择，可对应到选择语句；
- $\{ X \}$ 表示成分 X 的重复（0 到多次），可对应到循环语句；
- $[X]$ 表示成分 X 的任选（0 或 1 次），可对应到 if-then 语句；
- (X) 表示成分 X 的优先处理，可对应到复合语句。

在处理多个成分之间的选择时，也需要保证各成分的预测集合之间互不相交。计算预测集合可以采取类似于 4.1 节叙述的方法，先计算出必要的 First 集合和 Follow 集合。由于本课不涵盖有关 EBNF 形式定义的内容，所以这里不讨论这些集合的严格定义的算法，而是在随后介绍的实例中进行必要的讨论。为帮助理解这些实例，我们列出关于 First 集合的一些性质：

- $\text{First}(X_1 | X_2 | \dots | X_m) = \text{First}(X_1) \cup \dots \cup \text{First}(X_m)$
- $\text{First}(\{ X \}) = \text{First}(X) \cup \{\epsilon\}$
- $\text{First}([X]) = \text{First}(X) \cup \{\epsilon\}$
- $\text{First}((X)) = \text{First}(X)$

4.3.2 表驱动 LL (1) 分析程序

递归下降分析程序比较直观，容易设计。但不足之处就是递归调用可能带来的效率问题。这一小节里，我们介绍另外一种 LL (1) 分析程序的实现方法，称为表驱动的方法。这一方法用到了两个重要的数据结构（对象）：一个二维表和一个栈，称前者为预测分析表，后者为下推栈。

我们先介绍预测分析表的结构和构造算法。

设上下文无关文法 $G = (V_N, V_T, P, S)$ 。G 的**预测分析表**是一个二维表 M ，表的每一行 $A \in V_N$ 对应一个非终结符，表的每一列 $a \in V_T \cup \{\#\}$ 对应某个终结符或输入结束符 $\#$ 。表中的项 $M(A, a) \subseteq P$ 是一个产生式集合，它由如下过程得到：

- 对文法 G 的每个产生式 $A \rightarrow \alpha \in P$ ，若它的预测集合 $PS(A \rightarrow \alpha)$ 中包含 $a \in V_T \cup \{\#\}$ ，那么将 $A \rightarrow \alpha$ 加入 $M[A, a]$ 。

在预测分析表中，当 $M(A, a)$ 不含产生式时，就对应一个出错位置。

可以证明：一个文法 G 的预测分析表中每个表项最多只包含一个产生式，当且仅当 G 是 LL(1) 文法。对于 LL(1) 文法，预测分析表也称为 LL(1) 分析表。

例 5 设文法 $G[S]$ 为

$$\begin{aligned} S &\rightarrow AaS \mid BbS \mid d \\ A &\rightarrow a \\ B &\rightarrow \varepsilon \mid c \end{aligned}$$

试给出该文法的预测分析表。

解 根据例 4 的计算结果，各产生式的预测集合为：

$$\begin{aligned} PS(S \rightarrow AaS) &= \{a\} \\ PS(S \rightarrow BbS) &= \{c, b\} \\ PS(S \rightarrow d) &= \{d\} \\ PS(A \rightarrow a) &= \{a\} \\ PS(B \rightarrow \varepsilon) &= \{b\} \\ PS(B \rightarrow c) &= \{c\} \end{aligned}$$

这样，我们得到如图 1 的预测分析表。可以看出，表中每个表项最多只包含一个产生式。因此，文法 $G[S]$ 为 LL(1) 文法，即图 1 是一个 LL(1) 分析表。

	a	b	c	d	$\#$
S	$S \rightarrow AaS$	$S \rightarrow BbS$	$S \rightarrow BbS$	$S \rightarrow d$	
A	$A \rightarrow a$				
B		$B \rightarrow \varepsilon$	$B \rightarrow c$		

图 1 预测分析表

表驱动 LL(1) 分析程序的工作原理为：初始时，下推栈只包含 $\#$ ；首先将文法开始符号入栈；然后执行如下步骤：

- (1) 若栈顶为终结符，则判断当前读入的单词符号是否与该终结符相匹配，若匹配，再读取下一单词符号继续分析；若不匹配，则进行错误处理。
- (2) 若栈顶为非终结符，则根据该终结符和当前输入单词符号查 LL(1) 分析表，若相应表项中是产生式（唯一的），则将此非终结符出栈，并把产生式右部符号从右至左入栈；若表项为空，则进行错误处理。
- (3) 重复 (1) 和 (2)，直到栈顶为 $\#$ 同时输入也遇到结束符 $\#$ 时，分析结束。

图 2 给出了该分析过程的一个更具体的描述。

```
初始时 # 入栈，然后文法开始符号入栈；首个输入符号读进 a；
flag =TRUE；
while (flag) do {
    栈顶符号出栈并放在X中；
    if (X≠V1) {
        if (X==a)
            把下一个输入符号读进a；
        else ERROR；
    }
    else if (X=='#') {
        if (a=='#') flag = FALSE；
        else ERROR； /* 退出，转出错处理 */
    }
    else if (M[X,a] == {X0X1X2...Xk}) Xk,Xk-1,...,X1依次进栈；
    else ERROR； /* 退出，转出错处理 */
}
/*分析成功，过程完毕*/
```

图 2 表驱动 LL（1）分析过程

例如，基于图 1 中的 LL（1）分析表，对于输入符号串 *aabd* 的分析过程如图 3 所示。初始时，下推栈内容为 *#S*，*S* 为栈顶；余留符号串内容为 *aabd#*，其中，*#* 为输入结束符，最左边的 *a* 为下一个输入符号。由于当前栈顶为非终结符，所以通过查图 1 中的分析表得到下一步骤需要用到的产生式为 *S* → *AaS*（*S* 所在行，*a* 所在列）。下一步动作是将栈顶的 *S* 出栈，将 *AaS* 中的符号从右至左依次压栈。对栈顶为非终结符的情形依此类推。当栈顶为终结符时，如步骤 3），检查到该终结符 *a* 与下一个输入符号是匹配的，然后将栈顶的 *a* 出栈，下推栈内容变为 *#Sa*，同时余留符号串内容变为 *abd#*。假若在某一步发现栈顶的终结符与下一个输入符号不匹配，那末就要转到错误处理过程。重复类似的过程，直到步骤 10），下推栈栈顶和余留符号串内容都是 *#*，分析成功并返回。假若下推栈栈顶和余留符号串的内容只有一个为 *#*，则需要转错误处理过程。

步骤	下推栈	余留符号串	下一步动作
1)	#S	a a b d #	应用产生式 S ® AaS
2)	#S a A	a a b d #	应用产生式 A ® a
3)	#S a a	a a b d #	匹配栈顶和当前输入符号
4)	#S a	a b d #	匹配栈顶和当前输入符号
5)	#S	b d #	应用产生式 S ® BbS
6)	#S b B	b d #	应用产生式 B ® e
7)	#S b	b d #	匹配栈顶和当前输入符号
8)	#S	d #	应用产生式 S ® d
9)	#d	d #	匹配栈顶和当前输入符号
10)	#	#	返回：分析成功

图 3 例：表驱动 LL（1）分析过程

5. 文法变换：消除左递归、提取左公因子

如果所设计的文法不是 LL(1) 的，那么就不方便采用 LL(1) 分析方法。有时对文法进行适当的变换，有利于构造出所期望的文法类型。比如，消除左递归和提取左公因子。

当的变换，有利于构造出所期望的文法类型。

5.1 消除左递归

若文法中含有形如 $P \rightarrow P_1\alpha_0, P_1 \rightarrow P_2\alpha_1, \dots, P_n \rightarrow P\alpha_n$ 的一组产生式，其中 $n \geq 0$ ，那么称该文法是左递归的。若 $n=0$ ，则文法中含有形如 $P \rightarrow P\alpha$ 的产生式，称为直接左递归。相应地，若 $n > 0$ ，我们称之为间接左递归。

在实际设计的文法中，包含左递归和左公因子的文法一般不是 LL(1) 的，大家可以讨论一下其中缘由。所以如果想要获得 LL(1) 文法，多数情况下需要消除文法中的左递归和左公因子。

下面我们介绍消除左递归的一般方法。先考虑直接左递归的情形。

对于下列构成直接左递归的产生式组合：

$$P \rightarrow P\alpha \mid \beta$$

其中， $\alpha \neq \varepsilon$ ，以及 β 不以 P 打头，我们可将它们替换为：

$$\begin{aligned} P &\rightarrow \beta Q \\ Q &\rightarrow \alpha Q \mid \varepsilon \end{aligned}$$

其中 Q 为新增加的非终结符。对于如下具有更一般形式直接左递归的一组产生式：

$$P \rightarrow P\alpha_1 \mid P\alpha_2 \mid \dots \mid P\alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

其中 α_i ($1 \leq i \leq m$) 不为 ε ， β_j ($1 \leq j \leq n$) 不以 P 打头，我们可将这一组产生式替换为：

$$\begin{aligned} P &\rightarrow \beta_1 Q \mid \beta_2 Q \mid \dots \mid \beta_n Q \\ Q &\rightarrow \alpha_1 Q \mid \alpha_2 Q \mid \dots \mid \alpha_m Q \mid \varepsilon \end{aligned}$$

其中 Q 为新增加的非终结符。

例 6 设文法 $G[E]$ 为

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \end{aligned}$$

试变换该文法，得到一个等价的 不含左递归的文法。

解 消除直接左递归后得到等价的文法 $G'[E]$ ：

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \\ E' &\rightarrow \varepsilon \end{aligned}$$

$$\begin{aligned}
T &\rightarrow FT' \\
T' &\rightarrow *F T' \\
T' &\rightarrow \varepsilon \\
F &\rightarrow (E) \\
F &\rightarrow a
\end{aligned}$$

对于无回路（不存在满足 $A \Rightarrow^+ A$ 的非终结符 A ）、无 ε -产生式（不存在形如 $A \rightarrow \varepsilon$ 的产生式）的文法，通过下列步骤可消除一般左递归（包括直接和间接左递归）：

(1) 以某种顺序将文法非终结符排列为 A_1, A_2, \dots, A_n 。

(2) for $i = 1, n$ do {

for $j = 1, i-1$ do {

用 $A_i \rightarrow \alpha_1 r \mid \alpha_2 r \dots \mid \alpha_k r$ 反复替代形如 $A_i \rightarrow A_j r$ 的产生式，

其中 $A_j \rightarrow \alpha_1 \mid \alpha_2 \dots \mid \alpha_k$ 是关于 A_j 的全部产生式；

}

消除关于 A_i 的直接左递归；

}

(3) 化简由 (2) 得到的文法。

例 7 设文法 $G[S]$ 为

$$\begin{aligned}
S &\rightarrow PQ \mid a \\
P &\rightarrow QS \mid b \\
Q &\rightarrow SP \mid c
\end{aligned}$$

试变换该文法，得到一个等价的不含左递归的文法。

解 这是一个含有间接左递归的文法，并且无回路和 ε -产生式。我们按照上述消除一般左递归的方法，首先将非终结符排序为 S, P 和 Q ，然后进行如下变换：

先考虑关于 S 的产生式。不含直接左递归，所以不发生变化。

再考虑关于 P 的产生式。因为两个产生式都没有 S 打头，所以执行内层 for 循环后不发生变化。同样，关于 P 的产生式不含直接左递归，所以仍不发生变化。

最后考虑关于 Q 的产生式。执行内层 for 循环的过程可描述为

$$Q \rightarrow SP \mid c$$

替换为 $Q \rightarrow PQP \mid aP \mid c$

再替换为 $Q \rightarrow QSQP \mid bQP \mid aP \mid c$

然后再消去关于 Q 的直接左递归，得到

$$Q \rightarrow bQPR \mid aPR \mid cR$$

$$R \rightarrow SQPR \mid \varepsilon$$

我们列出最终结果：

$$\begin{aligned} S &\rightarrow PQ \mid a \\ P &\rightarrow QS \mid b \\ Q &\rightarrow bQPR \mid aPR \mid cR \\ R &\rightarrow SQPR \mid \varepsilon \end{aligned}$$

这里应该指出，如果对非终结符的不同排序方式，依上述步骤消除一般左递归后所得到的文法从形式上看可能会有区别，但它们之间是相互等价的。如，我们将非终结符排序为 Q , P 和 S ，则变换后所得到的文法为：

$$\begin{aligned} S &\rightarrow cSQR \mid bQR \mid aR \\ P &\rightarrow SPS \mid cS \mid b \\ Q &\rightarrow SP \mid c \\ R &\rightarrow PSQR \mid \varepsilon \end{aligned}$$

5.2 提取左公因子

$$\begin{aligned} P &\rightarrow \alpha Q \\ Q &\rightarrow \beta \mid \gamma \end{aligned}$$

这种文法变换称为提取左公因子。

在实际设计中，一个含有左公因子的文法通常不是 LL(1) 的（大家可以想一下其中的原因）。所以为了设计 LL(1) 文法时，一般需要提取文法中的所有左公因子。下面我们给出提取左公因子的一般变换规则。

设有如下形式的一组产生式

$$P \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_m \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_n$$

其中， α 称为 $\alpha\beta_1$, $\alpha\beta_2$, ... 和 $\alpha\beta_m$ 之间的左公因子，而每个 γ 不以 α 开头。提取左公因子 α 后，这一组产生式改写为：

$$\begin{aligned} P &\rightarrow \alpha Q \mid \gamma_1 \mid \gamma_2 \mid \dots \mid \gamma_n \\ Q &\rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m \end{aligned}$$

其中 Q 为新增加的非终结符。

例 8 设文法 $G[S]$ 为

$$\begin{aligned} S &\rightarrow \underline{if} \ C \ \underline{then} \ S \mid \underline{if} \ C \ \underline{then} \ S \ \underline{else} \ S \mid a \\ C &\rightarrow \underline{true} \mid \underline{false} \end{aligned}$$

其中， \underline{if} , \underline{then} , \underline{else} , \underline{true} 和 \underline{false} 为终结符。试变换该文法，得到一个等价的不含左公因子的文法。

解 提取 $S \rightarrow \underline{if} \ C \ \underline{then} \ S \mid \underline{if} \ C \ \underline{then} \ S \ \underline{else} \ S$ 之间的左公因子 $\underline{if} \ C \ \underline{then} \ S$ 后，得到如下等价于 $G[S]$ 的文法 $G'[S]$ ：

$$\begin{aligned} S &\rightarrow \underline{if} \ C \ \underline{then} \ S \ A \mid a \\ A &\rightarrow \underline{else} \ S \mid \varepsilon \end{aligned}$$

$$C \rightarrow \underline{true} \mid \underline{false}$$

其中 A 为新增加的非终结符。

6. LL(1) 分析中的出错处理

在编译程序设计中，错误处理主要包含两个方面的任务：一是报错，发现错误时应尽可能准确指出错误位置和错误属性；二是错误恢复，尽可能进行校正，使编译工作可以继续下去，提高程序调试的效率。

关于如何设计错误处理程序，目前并没有特别值得关注的理论成果。在这一小节里，我们只简单讨论 LL(1) 分析过程中实现错误处理的最基本方法。

对于表驱动 LL(1) 分析，在以下两种情况下需要报错：

- 栈顶的终结符与当前输入符号不匹配；
- 非终结符 A 位于栈顶，面临的输入符号为 a ，但分析表 M 的表项 $M[A, a]$ 为空。

一种简单的错误恢复措施是**应急恢复**（panic-mode error recovery）。例如，在表驱动 LL(1) 分析中，我们可以专为在表项 $M[A, a]$ 为空的情形指定一些所谓的**同步符号**。在分析过程中遇到这种情形时，就跳过输入符号串中的一些符号直至遇到同步符号为止。一种简便的做法是将 $\text{First}(A)$ 或 $\text{Follow}(A)$ 中的所有符号当作 A 的同步符号，相应的处理过程可设计为：

- 跳过输入符号串中的一些符号直至遇到 $\text{Follow}(A)$ 中的符号，然后把 A 从栈中弹出，便可以使分析继续下去；
- 跳过输入符号串中的一些符号直至遇到 $\text{First}(A)$ 中的符号时，可根据 A 恢复分析。

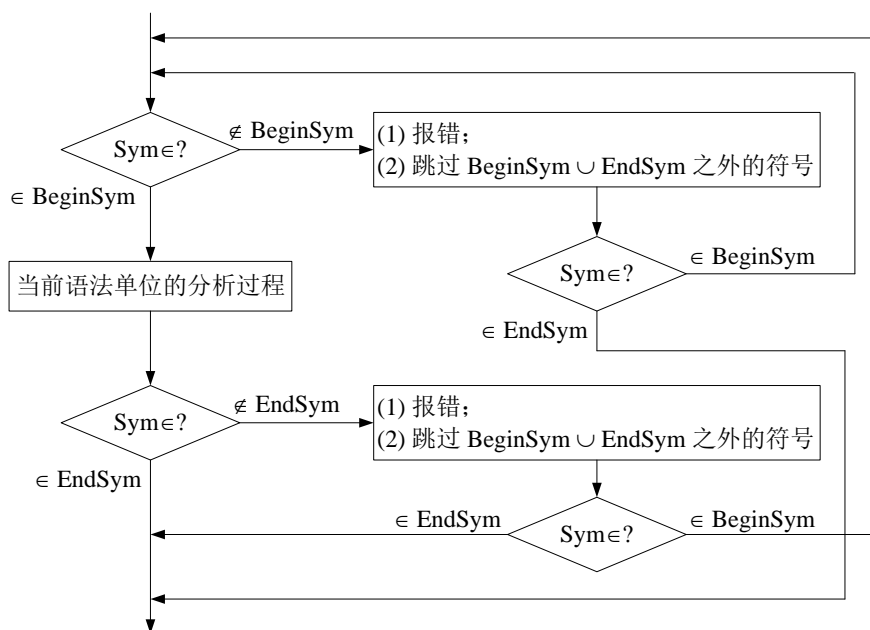


图 4 短语层恢复可采取的流程

另外一种称为**短语层恢复**（phrase-level error recovery）的措施比上述方法精确一些。原因在于，短语是与上下文相关的一个概念，而上述只考虑非终结符的同步符号的方法是与上下文无关的。有关短语的概念我们将在介绍自底向上语法分析时给出。

图 4 描述了短语层错误恢复可采取的一种流程：

- 在进入某个语法单位的分析时，检查当前符号是否属于进入该语法单位需要的符号集合 **BeginSym**。若不属于，则报错，并滤去补救的符号集合 $S = \text{BeginSym} \cup \text{EndSym}$ 外的所有符号。
- 在该语法单位分析结束时，检查当前符号是否属于离开该语法单位时需要的符号集合 **EndSym**。若不属于，则报错，并滤去补救的符号集合 $S = \text{BeginSym} \cup \text{EndSym}$ 外的所有符号。
- 无论是上述哪种情况，若遇 **BeginSym** 中的符号，则重新分析该语法单位；若遇 **EndSym** 中的符号，则退出该语法单位的分析。

我们来看在递归下降分析程序中采用短语层错误恢复的一个简单例子。设有如下文法产生式：

$$B \rightarrow [A] \mid (A)$$

$$A \rightarrow a$$

相应于非终结符 **B**，分析子程序可设计为：

```

procedure ParseB ( EndSym )
{
    if ( lookahead  $\notin$  { '[', '(' } ) {
        报错; 跳过 S 之外的单词;      /* S = { '[', '(' }  $\cup$  EndSym */
    }
    while ( lookahead  $\in$  { '[', '(' } ) {
        if ( lookahead == '[' ) {
            MatchToken( '[' );
            ParseA ( EndSym  $\cup$  { ']' } );
            MatchToken( ']' );
        }
        else {
            MatchToken( '(' );
            ParseA ( EndSym  $\cup$  { ')' } );
            MatchToken( ')' );
        }
        if ( lookahead  $\notin$  EndSym ) {
            报错; 跳过 S 之外的单词;      /* S = { '[', '(' }  $\cup$  EndSym */
        }
    }
}

```

该子程序对应于图 4 所描述的流程。**BeginSym**={ '[', '(' }，实际上是取 **First(B)**。**EndSym** 作为参数，由上一级子程序传入。值得注意的是，调用下一级分析子程序 **ParseA** 时，根据

A 不同上下文的后跟符号不同而使用了不同的参数。在方括号上下文中，使用的参数是 $\text{EndSym} \cup \{']'\}$ ；而在圆括号上下文中，使用的参数是 $\text{EndSym} \cup \{')'\}$ 。这可以体现“短语层”恢复的含义。试比较，若不区分方括号还是圆括号上下文，那么我们有可能使用 $\text{EndSym} \cup \{']', ')'\}$ 作为参数。后果如何？同学们可以思考一下。

这里，我们需要对前面（4.3.1 节）的函数 `MatchToken()` 进行一个小的修改：

```
void MatchToken(int expected)
{
    if (lookahead != expected)    //判别当前扫描的单词符号是否与期望的终结符
匹配
    {
        printf("syntax error \n"); //若不匹配，则报告出错信息，跳出
    }
    else                          //若匹配，消费掉当前单词符号
        lookahead = getToken();   //并向词法分析程序申请并读入下一个单词符号
}
```

其中，`lookahead` 为全局量，存放当前扫描的单词符号。

注意，与前面相比，修改后的 `MatchToken` 函数在报错后不退出系统。下同。

相应于非终结符 A，分析子程序可设计为：

```
procedure ParseA ( EndSym )
{
    if ( lookahead  $\notin$  { 'a' } ) {
        报错; 跳过 S 之外的单词;    /* S = { 'a' }  $\cup$  EndSym */
    }
    while ( lookahead  $\in$  { 'a' } ) {
        MatchToken ( 'a' );
        if ( lookahead  $\notin$  EndSym ) {
            报错; 跳过 S 之外的单词;    /* S = { 'a' }  $\cup$  EndSym */
        }
    }
}
```

7. LL(k) 文法（选讲）

LL(1) 分析的自顶向下预测分析策略是向前查看 1 个单词符号，然后确定应该选择哪一个产生式进行最左推导。这种方法要求文法必须是 LL(1) 文法。

如果允许向前查看任意 k ($k > 0$) 个单词符号，则进行自顶向下预测分析时对文法的要求可能会宽松一些。一般地，可以将 LL(1) 文法扩展至 LL(k) 文法。

直观来看，对于任何 CFG $G=(V, T, P, S)$ ，如果向前查看 k ($k > 0$) 个终结符，就可预测分析 $L(G[S])$ 中的所有句子，即在最左推导的每一步均能够确定应该选择哪一个产生式进行推导，这样的 CFG G 就是 LL(k) 文法。

以下是关于 $LL(k)$ 文法的几个重要结论：

- 对于整数 $k>0$ ，一个上下文无关文法是否为 $LL(k)$ 文法是可判定的。
- 对于一个上下文无关文法，是否存在一个整数 $k>0$ 使得该文法是 $LL(k)$ 文法，是不可判定的。
- 对于一个上下文无关文法，是否存在一个与之等价的 $LL(k)$ 文法 ($k>0$)，是不可判定的。
- 两个 $LL(k)$ 文法的语言是否相等是可判定的。
- $LL(k)$ 文法是无二义文法。
- $LL(k)$ 文法中不存在左递归的非终结符。
- 对于整数 $k>0$ ，不含 ϵ 产生式的 $LL(k)$ 文法的语言类真包含于不含 ϵ 产生式的 $LL(k+1)$ 文法的语言类。

这些结论在本课程的课堂讲稿 (ppt) 中有列举，目的是引起大家关注，但只是作为选讲，不属于课程的主体内容，也不会进行考察。

为满足一些同学的兴趣，以下对其中个别结论的由来及其与本讲前面部分较为相关的内容进行一些介绍，对于其余结论仅指出可进一步阅读的参考材料。

首先，需要有 $LL(k)$ 文法的严格数学定义。这里，我们引用 Aho 和 Ullman 教材[1] 中的 $LL(k)$ 文法定义。其他不同的定义方式，本质上也是一致的，如 Stearn 等人的定义[2, 3]。

先引入 **First_k 集合** 的概念。

设有 CFG $G=(V, T, P, S)$ ， $\alpha \in (V \cup T)^*$ ，以及某个固定的 $k>0$ 。定义

$\text{First}_k(\alpha) = \{ w \mid w \in T^*, \text{当 } |w| < k \text{ 时有 } \alpha \Rightarrow_{lm}^* w, \text{或者当 } |w|=k \text{ 时有 } \alpha \Rightarrow^* wx, \text{其中 } x \in (V \cup T)^* \}$ 。

CFG $G=(V, T, P, S)$ 是 **$LL(k)$ 文法**，如果对任意满足条件 $\text{First}_k(x) = \text{First}_k(y)$ 的两个最左推导

$$(1) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\beta\alpha \Rightarrow^* wx \quad \text{和}$$

$$(2) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\gamma\alpha \Rightarrow^* wy$$

都满足 $\beta = \gamma$ 。

该定义可直观解释为， G 是 $LL(k)$ 文法，如果对给定左句型 $wA\alpha$ ，以及输入串中 w 之后终结字符串的前 k 个符号构成的串 z ，满足 $z \in \text{First}_k(A\alpha)$ ，那么下一步最左推导所使用的产生式 $A \rightarrow \sigma$ 是唯一确定的。

若存在某个 $k>0$ ，文法 G 是 $LL(k)$ 文法，则称 G 是一个 **LL 文法**。

基于下列定理 (Aho 和 Ullman 教材[1]中的定理5.2)，可以帮助更好地理解 $LL(k)$ 文法。

定理: 设 CFG $G=(V, T, P, S)$ ， G 是 $LL(k)$ 文法，当且仅当下列条件成立：如果 $A \rightarrow \beta$ 和 $A \rightarrow \gamma$ 是 P 中不同的产生式，那么，对任何满足 $S \Rightarrow_{lm}^* wA\alpha$ 的左句型 $wA\alpha$ ，都有

$$\text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha) = \emptyset.$$

证明

(if) 假设 G 不是 $LL(k)$ 文法, 则存在左句型 $wA\alpha$ 以及两个不同的产生式 $A \rightarrow \beta$ 和 $A \rightarrow \gamma$ 是, 使得

$$(1) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\beta\alpha \Rightarrow^* wx \quad \text{和}$$

$$(2) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\gamma\alpha \Rightarrow^* wy$$

同时满足 $\text{First}_k(x) = \text{First}_k(y)$ 。这与 $\text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha) = \emptyset$ 矛盾。

(only if) 假设 $A \rightarrow \beta$ 和 $A \rightarrow \gamma$ 是 P 中不同的产生式, 存在满足 $S \Rightarrow_{lm}^* wA\alpha$ 的左句型 $wA\alpha$, 使得 $\text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha) \neq \emptyset$ 。则存在 $x \in \text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha)$, 使得

$$(1) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\beta\alpha \Rightarrow^* wxy \quad \text{和}$$

$$(2) S \Rightarrow_{lm}^* wA\alpha \Rightarrow_{lm} w\gamma\alpha \Rightarrow^* wxz$$

因为 G 是 $LL(k)$ 文法, 因此满足 $\beta = \gamma$, 这与 $A \rightarrow \beta$ 和 $A \rightarrow \gamma$ 是 P 中不同的产生式矛盾。

证毕。

接第3节的分析, 不难验证, 设 $k=2$, 下列文法 $G[S]$ 满足上述定理:

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aA \mid \varepsilon \\ B &\rightarrow b \mid bB \end{aligned}$$

因此, $G[S]$ 是 $LL(2)$ 文法。

以下是一个含有左递归的文法:

$$\begin{aligned} S &\rightarrow Sa \\ S &\rightarrow b \end{aligned}$$

对于该文法产生的任何终结符序列 ba^n ($n \geq 0$), 最左推导形如:

$$S \Rightarrow Sa \Rightarrow Saa \Rightarrow \dots \Rightarrow Sa^n \Rightarrow ba^n$$

对任何 $k > 0$, 取 $n \geq k-1$, 则存在左句型 Sa^{k-1} , 即 $S \Rightarrow_{lm}^* Sa^{k-1}$ 。在上述定理中, 取 w 为空串, α 为 a^{k-1} , β 为 Sa , γ 为 b , 则满足 $\text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha) \neq \emptyset$ 。因此, 该文法不是任何 $LL(k)$ 文法 ($k > 0$)。

这一结论可推广到任何含有左递归的文法 (假定不含无用符号), 无论是直接左递归还是间接左递归, 即上面所列的结论之一 “ $LL(k)$ 文法中不存在左递归的非终结符”。这一结论的证明也可参考 Rosenkrantz 和 Stearn 论文[2]中的引理5。

再看一个有关左公因子的例子。以下是一个含有左公因子的文法:

$$\begin{aligned} S &\rightarrow abA \mid abB \\ A &\rightarrow a \\ B &\rightarrow b \end{aligned}$$

根据上述定理，不难验证，该文法不是 $LL(1)$ 和 $LL(2)$ 文法，但它是 $LL(k)$ ($k \geq 3$) 文法。

另外，从上述定理可以直接得出上面所列的结论之一。

推论：任何 $LL(k)$ 文法一定是无二义文法。

设 $G=(V, T, P, S)$ 是 $LL(k)$ 文法，这一推论需要证明任何 $w \in T^*$ 均有唯一的最左推导，若不然，则一定存在如下两个不同的最左推导：

$$S \Rightarrow_{lm}^* xA\alpha \Rightarrow_{lm} x\beta\alpha \Rightarrow_{lm}^* xyz = w \text{ 和 } S \Rightarrow_{lm}^* xA\alpha \Rightarrow_{lm} x\gamma\alpha \Rightarrow_{lm}^* xyz = w$$

这里， $A \rightarrow \beta$ 和 $A \rightarrow \gamma$ 是 P 中不同的产生式，根据上述定理，应满足 $\text{First}_k(\beta\alpha) \cap \text{First}_k(\gamma\alpha) = \emptyset$ 。但从上面的两个推导，一定有 $\text{First}_k(yz) \in \text{First}_k(\beta\alpha)$ 以及 $\text{First}_k(yz) \in \text{First}_k(\gamma\alpha)$ 。矛盾。

对于上面所列的前 3 个结论，可参考 Aho 和 Ullman 教材[1]中第 5.1.6 节的讨论，其中涵盖了关于第 1 个结论中关于 $LL(k)$ 文法的一个判定算法。关于第 2、3 个结论的不可判定问题，该教材作为练习给出，参见其中的习题 5.1.11 和 5.1.12。这两个不可判定问题的证明也可参考 Rosenkrantz 和 Stearn 论文[2]中的定理 11 和定理 13，同时其定理 12 对应的是第 1 个结论。

关于“两个 $LL(k)$ 文法的语言是否相等”的可判定性证明，可参考 Aho 和 Ullman 教材[1]中的定理 8.6，或者 Rosenkrantz 和 Stearn 论文[2]中的定理 8。

对于上面所列的最后一个结论，其证明可参考 Rosenkrantz 和 Stearn 论文[2]中的定理 7。此外，由 Aho 和 Ullman 教材[1]中的定理 8.8，我们也有结论： $LL(k)$ 文法的语言类真包含于 $LL(k+1)$ 文法的语言类。

可将 4.3.2 节所介绍的表驱动 $LL(1)$ 分析算法推广至 k -预测分析算法（推广 $LL(1)$ 分析表至 $LL(k)$ 分析表）。参考 Aho 和 Ullman 教材[1]中的定理 5.1，如果在输入串后面加一个结束标志符， k -预测分析过程可以对应于一个 DPDA。我们以例 5 的下列 $LL(1)$ 文法 $G[S]$ 为例：

$$\begin{aligned} S &\rightarrow AaS \mid BbS \mid d \\ A &\rightarrow a \\ B &\rightarrow \epsilon \mid c \end{aligned}$$

下面将针对 $LL(1)$ 文法 $G[S]$ 定义一个空栈接受的 DPDA E ，使得其语言 $N(E) = \{ w\# \mid w \in L(G[S]) \}$ ，以及 $\#$ 是输入结束标志符。显然， $N(E)$ 满足前缀性质。

令 $G[S] = (V, T, P, S)$ ，其中 $V = \{S, A, B\}$ ， $T = \{a, b, c, d\}$ ，以及 P 由上面的 6 个产生式构成。我们定义 $E = (Q, T, V \cup T, \delta, q_0, \#)$ ，其中

$$Q = \{q_0\} \cup \{ \langle t, \epsilon \rangle \mid t \in T \cup \{\epsilon, \#\} \} \cup \{ \langle t, p \rangle \mid t \in T \cup \{\#\}, p \in P, \text{ 同时 } t \in PS(p) \}$$

其中 $PS(p)$ 为产生式 p 的预测集合。

δ 定义如下：

$$\delta(q_0, \epsilon, \#) = \{ \langle \epsilon, \epsilon \rangle, S \}, \delta(\langle \epsilon, \epsilon \rangle, a, S) = \{ \langle a, S \rightarrow AaS \rangle, S \}, \delta(\langle \epsilon, \epsilon \rangle, b, S) = \{ \langle b, S \rightarrow BbS \rangle, S \},$$

$$\delta(\langle \epsilon, \epsilon \rangle, c, S) = \{ \langle c, S \rightarrow BbS \rangle, S \}, \delta(\langle \epsilon, \epsilon \rangle, d, S) = \{ \langle d, S \rightarrow d \rangle, S \},$$

$$\delta(<\epsilon, \epsilon>, a, A) = \{ (<a, A \rightarrow a>, A) \}, \delta(<\epsilon, \epsilon>, b, B) = \{ (<b, B \rightarrow \epsilon>, B) \}, \delta(<\epsilon, \epsilon>, c, B) = \{ (<c, B \rightarrow c>, B) \},$$

$$\delta(<a, S \rightarrow AaS>, \epsilon, S) = \{ (<a, A \rightarrow a>, AaS) \}, \delta(<b, S \rightarrow BbS>, \epsilon, S) = \{ (<b, B \rightarrow \epsilon>, BbS) \},$$

$$\delta(<c, S \rightarrow BbS>, \epsilon, S) = \{ (<c, B \rightarrow c>, BbS) \}, \delta(<d, S \rightarrow d>, \epsilon, S) = \{ (<d, \epsilon>, d) \},$$

$$\delta(<a, A \rightarrow a>, \epsilon, A) = \{ (<a, \epsilon>, a) \} \delta(<b, B \rightarrow \epsilon>, \epsilon, B) = \{ (<b, \epsilon>, \epsilon) \}, \delta(<c, B \rightarrow c>, \epsilon, B) = \{ (<c, \epsilon>, c) \},$$

$$\delta(<a, \epsilon>, \epsilon, a) = \delta(<b, \epsilon>, \epsilon, b) = \delta(<c, \epsilon>, \epsilon, c) = \delta(<d, \epsilon>, \epsilon, d) = \delta(<\#, \epsilon>, \epsilon, \#) = \{ (<\epsilon, \epsilon>, \epsilon) \}$$

$$\delta(<\epsilon, \epsilon>, a, a) = \delta(<\epsilon, \epsilon>, b, b) = \delta(<\epsilon, \epsilon>, c, c) = \delta(<\epsilon, \epsilon>, d, d) = \delta(<\epsilon, \epsilon>, \#, \#) = \{ (<\epsilon, \epsilon>, \epsilon) \}$$

DPDA E 接受输入串 $aabd\#$ 的步骤为

$$\begin{aligned} (q_0, aabd\#, \#) &\vdash (<\epsilon, \epsilon>, aabd\#, S\#) \vdash (<a, S \rightarrow AaS>, abd\#, S\#) \\ &\vdash (<a, A \rightarrow a>, abd\#, AaS\#) \vdash (<a, \epsilon>, abd\#, aaS\#) \vdash (<\epsilon, \epsilon>, abd\#, aS\#) \\ &\vdash (<\epsilon, \epsilon>, bd\#, S\#) \vdash (<b, S \rightarrow BbS>, bd\#, S\#) \vdash (<b, B \rightarrow \epsilon>, d\#, BbS\#) \\ &\vdash (<b, \epsilon>, d\#, bS\#) \vdash (<\epsilon, \epsilon>, d\#, S\#) \vdash (<d, S \rightarrow d>, \#, S\#) \vdash (<d, \epsilon>, \#, d\#) \\ &\vdash (<\epsilon, \epsilon>, \#, \#) \vdash (<\epsilon, \epsilon>, \epsilon, \epsilon) \end{aligned}$$

这一步骤相当于图 3 所示的针对输入串 $aabd$ 的自顶向下表驱动 LL (1) 分析过程。

参考文献

- [1] Alfred V. Aho, Jefferey D. Ullman. The Theory of Parsing, Translation, and Compiling, Volume 1 & Volume 2, Prentice-Hall Series in Automatic Computation, 1972.
- [2] D. J. Rosenkrantz and R. E. Stearn, Properties of Deterministic Top-Down Grammars, Information and Control 17, 226-256, 1970.
- [3] P. M. II, Lewis, and R. E. Stearn, Syntax-directed transductions, J. ACM 15, 465-488, 1968.

课后作业

1. 设有文法 $G[S]$:

$$S \rightarrow aSb \mid aab$$

若针对该文法设计一个自顶向下预测分析过程,则需要向前察看多少个输入符号?

2. 设有文法 $G[E]$:

$$\begin{aligned}
S &\rightarrow i S A \\
A &\rightarrow e S \mid \varepsilon \\
S &\rightarrow a
\end{aligned}$$

针对文法 $G'[S]$ ，试给出各产生式右部文法符号串的 First 集合，各产生式左部非终结符的 Follow 集合，以及各产生式的预测集合 PS，即完成下表：

G 中的规则 r	$First(rhs(r))$	$Follow(lhs(r))$	$PS(r)$
$S \rightarrow i S A$			
$A \rightarrow e S$			
$A \rightarrow \varepsilon$		此处不填	
$S \rightarrow a$		此处不填	

表中的 $rhs(r)$ 表示产生式 r 右部的文法符号串， $lhs(r)$ 表示产生式 r 左部的非终结符。

3. 设有文法 $G[E]$ ：

$$\begin{aligned}
E &\rightarrow TA \\
A &\rightarrow \vee TA \mid \varepsilon \\
T &\rightarrow FB \\
B &\rightarrow \wedge FB \mid \varepsilon \\
F &\rightarrow (E) \mid i
\end{aligned}$$

试针对该文法填写下表

G 中的规则 r	$first(rhs(r))$	$follow(lhs(r))$	$PS(r)$
$E \rightarrow TA$			
$A \rightarrow \vee TA$			
$A \rightarrow \varepsilon$		此处不填	
$T \rightarrow FB$			
$B \rightarrow \wedge FB$			
$B \rightarrow \varepsilon$		此处不填	
$F \rightarrow (E)$			
$F \rightarrow i$		此处不填	

其中， $rhs(r)$ 表示产生式 r 的右部， $lhs(r)$ 表示产生式 r 的左部

4. 试验证下列文法 $G[S]$ 是 LL(1) 文法：

$$\begin{aligned}
S &\rightarrow P \mid \varepsilon \\
P &\rightarrow (P) P \mid a
\end{aligned}$$

其中 (,), 以及 a 为终结符

5. 计算下列文法中每个非终结符的 First 集和 Follow 集, 以及每个产生式的预测集合, 并判断该文法是否 LL(1)文法 (说明原因):

(1) 文法 $G_1[S]$:

$$\begin{aligned} S &\rightarrow (S) \mid aT \\ T &\rightarrow +ST \mid \varepsilon \end{aligned}$$

(2) 文法 $G_2[S]$:

$$\begin{aligned} S &\rightarrow TP \\ T &\rightarrow +PT \mid \varepsilon \\ P &\rightarrow (S) \mid a \end{aligned}$$

(3) 文法 $G_3[S]$:

$$\begin{aligned} S &\rightarrow A1B \\ A &\rightarrow 0A \mid \varepsilon \\ B &\rightarrow 0B \mid 1B \mid \varepsilon \end{aligned}$$

6. 验证如下文法是 LL(1)文法, 并基于该文法构造递归下降分析程序:

(1) 文法 $G[S]$:

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aA \\ A &\rightarrow \varepsilon \\ B &\rightarrow bB \\ B &\rightarrow \varepsilon \end{aligned}$$

(2) 文法 $G[E]$:

$$\begin{aligned} E &\rightarrow [F] A \\ A &\rightarrow E \mid \varepsilon \\ F &\rightarrow aB \\ B &\rightarrow aB \mid \varepsilon \end{aligned}$$

其中 [,], 以及 a 为终结符

7. 给出习题 5 中所有文法的预测分析表, 并根据分析表指出相应文法是否 LL(1)的, 同时验证习题 5 的结果。
8. 给出习题 6 中所有文法的 LL(1)分析表。
9. 基于图 1 中的 LL(1)分析表, 根据图 2 描述的方法, 给出输入符号串 $baacbd$ 的表驱动 LL(1)分析过程。
10. 通过变换求出与下列文法 $G[S]$ 等价的一个文法, 使其不含直接左递归:

$$\begin{aligned} S &\rightarrow AbB \\ A &\rightarrow Aa \mid a \\ B &\rightarrow Ba \mid Bb \mid b \end{aligned}$$

11. 按照本讲介绍的消除一般左递归算法消除下面文法 $G[S]$ 中的左递归（要求依非终结符的排序 S 、 Q 、 P 执行该算法）：

$$\begin{aligned} S &\rightarrow PQ \mid a \\ P &\rightarrow QS \mid b \\ Q &\rightarrow SP \mid c \end{aligned}$$

12. 按照本讲介绍的消除一般左递归算法消除下面文法 $G[P]$ 中的左递归（要求依非终结符的不同排序分别执行该算法）：

$$\begin{aligned} P &\rightarrow Qa \mid a \\ Q &\rightarrow Pb \mid b \end{aligned}$$

13. 变换下列文法，求出与其等价的一个文法，使变换后的文法不含左递归和左公因子：

(1) 文法 $G[P]$ ：

$$P \rightarrow Pa \mid Pb \mid c$$

(2) 文法 $G'[S]$ ：

$$\begin{aligned} S &\rightarrow aSAc \mid a \\ A &\rightarrow Ab \mid d \end{aligned}$$

14. 给定某类表达式文法 $G[E]$ ：

$$\begin{aligned} E &\rightarrow +ER \mid -ER \mid \underline{positive} R \\ R &\rightarrow *ER \mid \varepsilon \end{aligned}$$

其中， $+$ 和 $-$ 分别代表一元正和一元负运算， $*$ 代表普通的二元乘法运算，positive 为代表正整数（非0）的单词。

- (1) 针对文法 $G[E]$ ，下表给出各产生式右部文法符号串的 *First* 集合，各产生式左部非终结符的 *Follow* 集合，以及各产生式的预测集合 *PS*。试填充其中空白表项（共3处）的内容：

$G[E]$ 的规则 r	$First(rhs(r))$	$Follow(lhs(r))$	$PS(r)$
$E \rightarrow +ER$	+		+
$E \rightarrow -ER$	-	此处不填	-
$E \rightarrow \underline{positive} R$	<u>positive</u>	此处不填	<u>positive</u>
$R \rightarrow *ER$	*		*
$R \rightarrow \varepsilon$	ε	此处不填	

表中， $rhs(r)$ 为产生式 r 右部的文法符号串， $lhs(r)$ 为产生式 r 左部的非终结符。

- (2) $G[E]$ 不是 LL(1) 文法，试解释为什么？

- (3) 虽然 $G[E]$ 不是 LL(1) 文法，但可以采用一种强制措施，使得常规的 LL(1)

分析算法仍然可用。针对含5个单词的输入串 $+ - 20 * 18$ ，以下基于这一措施以及上述各产生式的预测集合（或预测分析表）的一个表驱动 LL(1) 分析过程：

步骤	下推栈	余留符号串	下一步动作
1	# E	$+ - 20 * \underline{18} \#$	应用产生式 $E \rightarrow + E R$
2	# $R E +$	$+ - \underline{20} * \underline{18} \#$	匹配栈顶和当前输入符号
3	# $R E$	$- \underline{20} * \underline{18} \#$	应用产生式 $E \rightarrow - E R$
4	# $R R E -$	$- \underline{20} * \underline{18} \#$	匹配栈顶和当前输入符号
5	# $R R E$	$\underline{20} * \underline{18} \#$	应用产生式 $E \rightarrow \underline{positive} R$
6	# $R R R \underline{positive}$	$\underline{20} * \underline{18} \#$	匹配栈顶和当前输入符号
7	# $R R R$	$* \underline{18} \#$	
8			
9			
10			
11			
12	# $R R R$	#	应用产生式 $R \rightarrow \epsilon$
13	# $R R$	#	应用产生式 $R \rightarrow \epsilon$
14	# R	#	应用产生式 $R \rightarrow \epsilon$
15	#	#	结束

试填写上述分析过程中第7步时使用的产生式，以及第 8~11 步的分析过程，共计13处空白；并指出采用了什么样的强制措施。