

Глава 3. Нисходящий синтаксический анализ

3.4. $LL(1)$ -грамматики

3.4.4. Основные приемы преобразования КС-грамматик в $LL(1)$ -форму

Не существует полностью универсального автоматического процесса преобразования грамматик в $LL(1)$ -форму. Это связано с тем, что данная проблема алгоритмически неразрешима. Отсутствие общего решения проблемы не означает невозможности ее решения для частных случаев.

Во многих случаях КС-грамматики, которые не обладают $LL(1)$ -свойствами, можно преобразовать в $LL(1)$ -грамматики с помощью левой факторизации (см. разд. 1.6.5). Рассмотрим это преобразование на примере следующей грамматики:

$$S \rightarrow aSb \mid aSc \mid \varepsilon$$

Определим множества направляющих символов (учитывая, что грамматика пополнена маркером конца строки \perp , т. е. в предположении, что в грамматике имеется продукция $S' \rightarrow S\perp$):

$$DS(S \rightarrow aSb) = \{a\},$$

$$DS(S \rightarrow aSc) = \{a\},$$

$$DS(S \rightarrow \varepsilon) = Follow(S) = \{b, c, \perp\}.$$

Направляющий символ a является общим для двух первых продукций, т. е. это не $LL(1)$ -грамматика. Применим для них правило левой факторизации (общим префиксом является aS). В результате получится $LL(1)$ -грамматика

$$S \rightarrow aSX \mid \varepsilon$$

$$X \rightarrow b \mid c$$

Таким образом, факторизация как бы выносит за скобки направляющие символы.

Многие языки программирования содержат такие часто используемые конструкции, как списки. Поэтому рассмотрим описание различных типов списков с помощью $LL(1)$ -грамматик.

Грамматика $L \rightarrow a ; L \mid a$ для порождения списка символов a , разделенных символом-разделителем (для определенности в качестве символа-разделителя используем точку с запятой), не обладает $LL(1)$ -свойствами. Применив левую факторизацию, получим $LL(1)$ -грамматику

$$L \rightarrow aR$$

$$R \rightarrow ; aR \mid \varepsilon$$

Аналогично, применив левую факторизацию к продукциям $L \rightarrow aL \mid a$, порождающих список символов a без символов-разделителей, получим $LL(1)$ -грамматику

$$L \rightarrow aR$$

$$R \rightarrow aR \mid \varepsilon$$

Рассмотренные грамматики не допускают пустой список. Если это необходимо, то достаточно добавить в грамматики продукцию $L \rightarrow \varepsilon$. Для списка из нуля или более символов a возможно применение грамматики $L \rightarrow aL \mid \varepsilon$, но иногда она может оказаться непригодной для трансляции.

Следует заметить, что грамматика, содержащая левую рекурсию, не является $LL(1)$ -грамматикой. Рассмотрим продукции

$$A \rightarrow A\alpha \text{ (леворекурсивная продукция по } A\text{),}$$

$$A \rightarrow a$$

$DS(A \rightarrow A\alpha) = \{a\}$ и $DS(A \rightarrow a) = \{a\}$, т. е. множества направляющих символов пересекаются.

По тем же причинам грамматика, содержащая леворекурсивный цикл, не может быть $LL(1)$ -грамматикой.

Как уже отмечалось, левую рекурсию всегда можно исключить из грамматики, преобразовав ее в правую рекурсию, которая не вызывает никаких проблем для нисходящего разбора (см. разд. 1.6.3 и 1.6.4).

После устранения левой рекурсии в соответствии с преобразованием из разд. 1.6.3 для получения $LL(1)$ -грамматики всегда требуется последующая левая факторизация. Можно сформулировать правило, объединяющее эти преобразования.

Пусть множество продукций содержит леворекурсивные продукции $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m$ и все остальные A -продукции $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$, не являющиеся леворекурсивными. Результатом устранения левой рекурсии будут продукции

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

$$A' \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_m | \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A'$$

Как видим, имеются продукции с общими префиксами, поэтому выполним левую факторизацию:

$$A \rightarrow \beta_1 X | \beta_2 X | \dots | \beta_n X$$

$$X \rightarrow \varepsilon | A'$$

$$A' \rightarrow \alpha_1 X | \alpha_2 X | \dots | \alpha_m X$$

Выполнив замену вхождений нетерминала A' , получим окончательный результат:

$$A \rightarrow \beta_1 X | \beta_2 X | \dots | \beta_n X$$

$$X \rightarrow \varepsilon | \alpha_1 X | \alpha_2 X | \dots | \alpha_m X$$

Таким образом, модифицированное правило устранения левой рекурсии можно сформулировать следующим образом.

Пусть множество продукций содержит леворекурсивные продукции $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m$ и остальные A -продукции $A \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$, не являющиеся леворекурсивными. Тогда новая эквивалентная грамматика будет иметь продукции вида:

$$A \rightarrow \beta_1 X | \beta_2 X | \dots | \beta_n X$$

$$X \rightarrow \varepsilon | \alpha_1 X | \alpha_2 X | \dots | \alpha_m X$$

Для иллюстрации преобразуем в $LL(1)$ -форму следующую грамматику (предполагая наличие продукции $E' \rightarrow E\perp$):

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

$$T \rightarrow T \times F \mid F$$

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

Она не является $LL(1)$ -грамматикой, поскольку содержит леворекурсивные productions. Применим рассмотренное выше модифицированное правило устранения левой рекурсии:

$$E \rightarrow TX$$

$$X \rightarrow + TX \mid \varepsilon$$

$$T \rightarrow FY$$

$$Y \rightarrow \times FY \mid \varepsilon$$

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

Это $LL(1)$ -грамматика. Покажем это, определив множества направляющих символов для productions с одинаковыми левыми частями:

$$DS(X \rightarrow + TX) = \{+\},$$

$$DS(X \rightarrow \varepsilon) = \{(), \perp\},$$

$$DS(Y \rightarrow \times FY) = \{\times\},$$

$$DS(Y \rightarrow \varepsilon) = \{+,), \perp\},$$

$$DS(F \rightarrow (E)) = \{(),$$

$$DS(F \rightarrow i) = \{i\}.$$

Множества направляющих символов всех пар альтернативных productions не пересекаются, что позволяет при разборе строки детерминированно выбирать нужную production.

$$E \rightarrow TX$$

$$X \rightarrow + TX \mid \varepsilon$$

$$T \rightarrow FY$$

$$Y \rightarrow \times FY \mid \varepsilon$$

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

$$DS(X \rightarrow + TX) = \{+\},$$

$$DS(X \rightarrow \varepsilon) = \{\}, \perp\},$$

$$DS(Y \rightarrow \times FY) = \{\times\},$$

$$DS(Y \rightarrow \varepsilon) = \{+, \), \perp\},$$

$$DS(F \rightarrow (E)) = \{\{\},$$

$$DS(F \rightarrow i) = \{i\}.$$

Процесс вывода строки $i \times (i+i) \perp$ соответствует следующей левосторонней схеме вывода:

$$\begin{aligned} E' &\Rightarrow E \perp \Rightarrow TX \perp \Rightarrow FYX \perp \Rightarrow iYX \perp \Rightarrow i \times FYX \perp \Rightarrow i \times (E) YX \perp \\ &\Rightarrow i \times (TX) YX \perp \Rightarrow i \times (FYX) YX \perp \Rightarrow i \times (iYX) YX \perp \Rightarrow i \times (iX) YX \perp \\ &\Rightarrow i \times (i+TX) YX \perp \Rightarrow i \times (i+FYX) YX \perp \Rightarrow i \times (i+iYX) YX \perp \\ &\Rightarrow i \times (i+iX) YX \perp \Rightarrow i \times (i+i) YX \perp \Rightarrow i \times (i+i) X \perp \Rightarrow i \times (i+i) \perp. \end{aligned}$$

К сожалению, процесс факторизации нельзя распространить на общий случай. Следующий пример показывает, что может произойти. Рассмотрим грамматику

$$S \rightarrow Ac|Bd$$

$$A \rightarrow eAf|a$$

$$B \rightarrow eBg|b$$

Первые две S -продукции в своих множествах направляющих символов содержат общий символ e . Для проведения факторизации предварительно выполним замену вхождений нетерминалов A и B , чтобы направляющие символы явно присутствовали в этих продукциях:

$$S \rightarrow eAfc|ac|eBgd|bd$$

Выполняя факторизацию, эти продукции можно заменить следующими:

$$S \rightarrow ac|bd|eS_1$$

$$S_1 \rightarrow Afc|Bgd$$

Продукции для S_1 аналогичны первоначальным продукциям для S и имеют пересекающиеся множества направляющих символов. Можно повторить преобразование этих продукций, как это было сделано с продукциями для S :

$$S_1 \rightarrow eAffc|afc|eBggd|bgd$$

В результате факторизации получим

$$S_1 \rightarrow afc|bgd|eS_2$$

$$S_2 \rightarrow Affc|Bggd$$

Продукции для S_2 аналогичны продукциям для S_1 и S , но длиннее их, и теперь очевидно, что этот процесс бесконечный.

Это означает, что все попытки преобразовать грамматику в $LL(1)$ -форму не всегда приводят к результату. Это является прямым следствием того, что $LL(1)$ -языки являются только подклассом более широкого класса языков, допускающих детерминированный разбор.