

Практическое занятие №5. Синтаксический анализ с использованием магазинных автоматов

Д.Ю. Чалый

24 декабря 2010 г.

1 Построение эквивалентного магазинного автомата для КС-грамматики

Рассмотрим следующую грамматику:

$$\begin{array}{lcl} S & \rightarrow & ABa \mid DE \\ A & \rightarrow & EC \mid AC \\ B & \rightarrow & b \mid \epsilon \\ C & \rightarrow & c \\ D & \rightarrow & d \mid \epsilon \\ E & \rightarrow & BAB \mid c \end{array}$$

Тогда эквивалентный этой грамматике МП-автомат будет таким: $M = (\{q\}, \{a, b, c, d\}, \{S, A, B, C, D, E\} \cup \{a, b, c, d\}, \delta, q, S, \{q\})$, где:

- $\delta(q, \epsilon, S) = \{(q, ABa), (q, DE)\};$
- $\delta(q, \epsilon, A) = \{(q, EC), (q, AC)\};$
- $\delta(q, \epsilon, B) = \{(q, b), (q, \epsilon)\};$
- $\delta(q, \epsilon, C) = \{(q, c)\};$
- $\delta(q, \epsilon, D) = \{(q, d), (q, \epsilon)\};$
- $\delta(q, \epsilon, E) = \{(q, BAB), (q, c)\};$
- $\delta(q, a, a) = \{(q, \epsilon)\};$
- $\delta(q, b, b) = \{(q, \epsilon)\};$
- $\delta(q, c, c) = \{(q, \epsilon)\};$
- $\delta(q, d, d) = \{(q, \epsilon)\};$

В общем-то построение автомата является простой задачей, а в динамике он ведет себя следующим образом. Если в стеке наверху находится нетерминал, то он заменяется на одну из правых частей продукции для этого нетерминала (так как МП-автомат является недетерминированным вычислительным устройством, то нельзя сказать на какую), а если сверху стека находится терминал, то происходит его сравнение с очередным символом входной строки. Если сравнение успешно, то символ изымается из стека, в противном случае автомат останавливается.

Проблема состоит в том, что МП-автомат является недетерминированным по своей природе и мы должны как-то промоделировать его детерминированными средствами (все алгоритмы в своей основе построены на положении, что на основе текущего состояния мы однозначным, детерминированным образом определяем как изменить это состояние).

Для МП-автоматов есть два алгоритма — поиск в ширину и поиск в глубину. Однако они применимы только когда исходная грамматика не является леворекурсивной. Данная грамматика имеет левую рекурсию. После преобразования (проделайте это сами и не забудьте про скрытую левую рекурсию) получится следующая грамматика¹:

$$\begin{aligned}
 S &\rightarrow ABa \mid DE \\
 A &\rightarrow EC \mid ECA' \\
 A' &\rightarrow C \mid CA' \\
 B &\rightarrow b \mid \epsilon \\
 C &\rightarrow c \\
 D &\rightarrow d \mid \epsilon \\
 E &\rightarrow bAB \mid c \mid bABE' \mid cE' \\
 E' &\rightarrow CA'B \mid CB \mid CBE' \mid CA'BE'
 \end{aligned}$$

Перед рассмотрением алгоритмов необходимо поименовать продукции. Будем их именовать символом, который находится в голове продукции, а индекс будет задавать номер продукции по порядку вышеупомянутой грамматики. Например, продукция $S \rightarrow ABa$ будем иметь имя S_1 , а продукция $A' \rightarrow CA' - A'_2$.

Пусть на вход подана строка $\alpha = dccb$. Необходимо найти вывод этой строки в данной грамматике. Кстати, вообще говоря, нет необходимости явно строить МП-автомат для грамматики, так как он имеет простую структуру, и работать непосредственно с грамматикой.

2 Алгоритм поиска в ширину

Алгоритм поиска ширину начинает свою работу со нижеследующего состояния, которое соответствует начальной конфигурации магазинного автомата:

¹Обратите внимание что приведенная грамматика имеет ϵ -продукции. В общем случае перед применением алгоритма удаления левой рекурсии от них необходимо избавляться, но для данной конкретной грамматики алгоритм применим напрямую. Наличие ϵ -продукций позволяет сократить объем грамматики.

(1)

	$dc cb\sharp$
	$S\sharp$

Состояние удобно задавать в виде таблицы, где в правой части записывается сверху остаток входной строки, а снизу стек МП-автомата. Алгоритм моделирования в ширину состоит из применения двух видов шагов: раскрытие нетерминала и сравнение терминалов. Сначала совершаем раскрытие согласно тех продуктов, которые имеются для нетерминала сверху стека (верхний символ располагает слева, в данном случае это символ S , второй символ' символ \sharp , находится под ним):

	$dc cb\sharp$
S_1	$A \bar{B} a \sharp$
S_2	$D E \sharp$

Так как для S имеются две продукции, то мы исходную строчку заменили на две, а слева записали имена продуктов, которые мы применили. Шаг раскрытия нетерминала повторяется до тех пор, пока в правой части таблицы стеки начинаются с нетерминалов. Следующий шаг преобразует состояние следующим образом:

	$dc cb\sharp$
$S_1 A_1$	$E C B a \sharp$
$S_1 A_2$	$E C A' B a \sharp$
$S_2 D_1$	$d E \sharp$
$S_2 D_2$	$E \sharp$

Как видно, слева формируются строки, которые хранят историю применяемых продуктов. Применяем шаг раскрытия продуктов до тех пор, пока есть хоть одна строка, в левой части которой стек начинается с нетерминала. В итоге получим:

	$dc cb\sharp$
$S_1 A_1 E_1$	$b A B C B a \sharp$
$S_1 A_1 E_2$	$c C B a \sharp$
$S_1 A_1 E_3$	$b A B E' C B a \sharp$
$S_1 A_1 E_4$	$c E' C B a \sharp$
$S_1 A_2 E_1$	$b A B C A' B a \sharp$
$S_1 A_2 E_2$	$c C A' B a \sharp$
$S_1 A_2 E_3$	$b A B E C A' B a \sharp$
$S_1 A_2 E_4$	$c E' C A' B a \sharp$
$S_2 D_1$	$d E \sharp$
$S_2 D_2 E_1$	$b A B \sharp$
$S_2 D_2 E_2$	$c \sharp$
$S_2 D_2 E_3$	$b A B E' \sharp$
$S_2 D_2 E_4$	$c E' \sharp$

Теперь все строки в правой части таблицы начинаются с терминального символа. Производим сравнение этого символа с очередным символом на входе, символом d. Если символ сверху стека не совпадает с ним, то вычеркиваем строку. Если же символ совпадает, то моделируем его успешное сравнение, перенося из правой части таблицы в левую. Непосредственно после этого действия получим:

d	$ccb\sharp$
S_2D_1d	$E\sharp$

У нас снова в правой части таблицы появились строчки, которые начинаются с нетерминалов. Раскрываем нетерминалы так же, как и ранее:

d	$ccb\sharp$
$S_2D_1dE_1$	$bAB\sharp$
$S_2D_1dE_2$	$c\sharp$
$S_2D_1dE_3$	$bAB'E'\sharp$
$S_2D_1dE_4$	$cE'\sharp$

Снова все строки начинаются с терминальных символов, производим сравнение и получаем:

dc	$cb\sharp$
$S_2D_1dE_2c$	\sharp
$S_2D_1dE_4c$	$E'\sharp$

В таблице появилась строка, начинающаяся с нетерминала, раскрываем его:

dc	$cb\sharp$
$S_2D_1dE_2c$	\sharp
$S_2D_1dE_4cE'_1$	$CA'B\sharp$
$S_2D_1dE_4cE'_2$	$CB\sharp$
$S_2D_1dE_4cE'_3$	$CBE'\sharp$
$S_2D_1dE_4cE'_4$	$CA'BE'\sharp$

Следующий шаг:

dc	$cb\sharp$
$S_2D_1dE_2c$	\sharp
$S_2D_1dE_4cE'_1C_1$	$CA'B\sharp$
$S_2D_1dE_4cE'_2C_1$	$CB\sharp$
$S_2D_1dE_4cE'_3C_1$	$CBE'\sharp$
$S_2D_1dE_4cE'_4C_1$	$CA'BE'\sharp$

Производим сравнение:

dc	$b\sharp$
$S_2 D_1 d E_4 c E'_1 C_1 c$	$A' B\sharp$
$S_2 D_1 d E_4 c E'_2 C_1 c$	$B\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c$	$B E'\sharp$
$S_2 D_1 d E_4 c E'_4 C_1 c$	$A' B E'\sharp$

После раскрытия нетерминалов получим:

dc	$b\sharp$
$S_2 D_1 d E_4 c E'_1 C_1 c A'_1 C_1$	$c B\sharp$
$S_2 D_1 d E_4 c E'_1 C_1 c A'_2 C_1$	$c A' B\sharp$
$S_2 D_1 d E_4 c E'_2 C_1 c B_1$	$b\sharp$
$S_2 D_1 d E_4 c E'_2 C_1 c B_2$	\sharp
$S_2 D_1 d E_4 c E'_3 C_1 c B_1$	$b E'\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_2 E'_1 C_1$	$c A' B\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_2 E'_1 C_1$	$c B\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_2 E'_1 C_1$	$c B E'\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_2 E'_1 C_1$	$c A' B E'\sharp$
$S_2 D_1 d E_4 c E'_4 C_1 c A'_1 C_1$	$c B E'\sharp$
$S_2 D_1 d E_4 c E'_4 C_1 c A'_2 C_1$	$c A' B E'\sharp$
$S_2 D_1 d E_4 c E'_4 C_1 c A'_1 C_1$	$c B E'\sharp$
$S_2 D_1 d E_4 c E'_4 C_1 c A'_2 C_1$	$c A' B E'\sharp$

Производим сравнение:

dc	\sharp
$S_2 D_1 d E_4 c E'_2 C_1 c B_1$	\sharp
$S_2 D_1 d E_4 c E'_3 C_1 c B_1$	$E'\sharp$

Снова раскрываем нетерминалы:

dc	\sharp
$S_2 D_1 d E_4 c E'_2 C_1 c B_1$	\sharp
$S_2 D_1 d E_4 c E'_3 C_1 c B_1 E'_1 C_1$	$c A' B\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_1 E'_2 C_1$	$c B\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_1 E'_3 C_1$	$c B E'\sharp$
$S_2 D_1 d E_4 c E'_3 C_1 c B_1 E'_4 C_1$	$c A' B E'\sharp$

И, наконец, производим заключительное сравнение:

dc
$S_2 D_1 d E_4 c E'_2 C_1 c B_1 \sharp$

В итоге у нас опустошился стек и одновременно мы прочитали входную строку целиком. Следовательно, эта строка является выводимой в данной грамматике. В левой части таблицы у нас образовалась цепочка символов $S_2 D_1 d E_4 c E'_2 C_1 c B_1 \#$, называемая текущим выводом. При помощи нее можно восстановить вывод входной строки. Уберем из текущего вывода терминальные символы, получим $S_2 D_1 E_4 E'_2 C_1 B_1$. Методы моделирования магазинных автоматов строят левосторонний вывод исходной строки. Вывод начинается с начального символа грамматики S . Последовательно применяем продукцию $S_2, D_1, E_4, E'_2, C_1, B_1$ в левостороннем выводе: $S \rightarrow DE \rightarrow dE \rightarrow dcE' \rightarrow dcCB \rightarrow dcCb$.

3 Алгоритм поиска в глубину

Как было видно, при поиске в ширину таблица может значительно расти в объеме. Поиск в глубину позволяет избавиться от этого, однако алгоритм поиска содержит большее количество шагов.

Начинается поиск в глубину с того же самого состояния, что и поиск в ширину:

(1)		<u><u>dcCb#</u></u>
		<u><u>S#</u></u>

При поиске в глубину мы проходим два этапа. Первый этап состоит в раскрытии нетерминалов, можно назвать его «прямым» ходом. Выбираем первую продукцию для нетерминала S (продукции выбираются по порядку, алгоритм не делает никаких предположений о том, какую продукцию выбрать):

	<u><u>dcCb#</u></u>
<u><u>S₁</u></u>	<u><u>ABa#</u></u>

Прямой ход алгоритма продолжается до тех пор, пока в правой части таблицы стек начинается с нетерминала:

	<u><u>dcCb#</u></u>		<u><u>dcCb#</u></u>
<u><u>S₁A₁</u></u>	<u><u>ECBa#</u></u>	→	<u><u>S₁A₁E₁</u></u>

Как справа мы встретили терминал, в данном случае b , необходимо привести сравнение очередного символа на входе (символ d). Сравниваемые символы неравны и алгоритм переходит режим возврата. Начинается режим с попытки перебора продукции для последнего рассмотренного нетерминала. Последним мы рассматривали нетерминал E и продукцию E_1 , которая не подошла. Значит необходимо попробовать следующую продукцию, E_2 :

	$dccb\sharp$
$S_1 A_1 E_1$	$bABC Ba\sharp$

 \rightarrow

	$dccb\sharp$
$S_1 A_1 E_2$	$cCBa\sharp$

Снова неудача при сравнении и берем следующую продукцию:

	$dccb\sharp$
$S_1 A_1 E_2$	$cCBa\sharp$

 \rightarrow

	$dccb\sharp$
$S_1 A_1 E_3$	$bABE'CBa\sharp$

и опять...:

Снова неудача при сравнении и берем следующую продукцию:

	$dccb\sharp$
$S_1 A_1 E_3$	$bABE'CBa\sharp$

 \rightarrow

	$dccb\sharp$
$S_1 A_1 E_4$	$cE'CBa\sharp$

Все, продукции для E закончились, значит необходимо осуществить возврат по нетерминалу. Мы убираем из стека тело последней рассмотренной продукции для E и на ее место перемещаем сам нетерминал из левой части таблицы в правую:

	$dccb\sharp$
$S_1 A_1 E_4$	$cE'CBa\sharp$

 \rightarrow

	$dccb\sharp$
$S_1 A_1$	$ECBa\sharp$

Теперь у нас в левой части таблицы находится продукция A_1 , которая, как показала работа алгоритма, не подошла. Следовательно, надо попробовать следующую продукцию, A_2 :

	$dccb\sharp$
$S_1 A_1$	$ECBa\sharp$

 \rightarrow

	$dccb\sharp$
$S_1 A_2$	$ECA'Ba\sharp$

Чтобы не загромождать изложение, можно сразу увидеть, что далее рассмотрение продукции для E не приведет к успеху, мы так же безуспешно просмотрим все продукции для E и совершим возврат. Следовательно, необходимо совершить возврат по A :

	$dccb\sharp$
$S_1 A_1$	$ECA'Ba\sharp$

 \rightarrow

	$dccb\sharp$
S_1	$ABa\sharp$

Снова идем вперед, и берем следующую продукцию для S :

	$dccb\sharp$
S_1	$ABa\sharp$

 \rightarrow

	$dccb\sharp$
S_2	$DE\sharp$

Раскрываем нетерминал D :

	$dccb\sharp$
S_2	$DE\sharp$

 \rightarrow

	$dccb\sharp$
$S_2 D_1$	$dE\sharp$

Производим успешное сравнение:

	$dccb\sharp$
$S_2 D_1$	$dE\sharp$

 \rightarrow

	$dccb\sharp$
$S_2 D_1 d$	$E\sharp$

Далее без комментариев дальнейшая работа алгоритма будет такой: