

Теория формальных языков и компиляторов

Часть 1. Порождающие грамматики и языки

Лекция 4. Языки порождающих грамматик

Вернемся к рассмотрению примера 3.7:

Пример 3.7. Рассмотрим грамматику $G_4[Z = \langle \text{целое без знака} \rangle]$ с правилами:

- 1) $\langle \text{целое без знака} \rangle \rightarrow \langle \text{цифра} \rangle$
- 2) $\langle \text{целое без знака} \rangle \rightarrow \langle \text{цифра} \rangle \langle \text{целое без знака} \rangle$
- 3) $\langle \text{цифра} \rangle \rightarrow 0|1|2|3|\dots|9$

Теперь поставим задачу анализа строк: принадлежит ли цепочка «123» множеству $\langle \text{целых без знака} \rangle$?

Анализ проведём, используя определение *выводимости*:

$$\begin{aligned} \langle \text{целое без знака} \rangle &\Rightarrow^{(2)} \langle \text{цифра} \rangle \langle \text{целое без знака} \rangle \Rightarrow^{(2)} \\ &\Rightarrow^{(2)} \langle \text{цифра} \rangle \langle \text{цифра} \rangle \langle \text{целое без знака} \rangle \Rightarrow^{(1)} \\ &\Rightarrow^{(1)} \langle \text{цифра} \rangle \langle \text{цифра} \rangle \langle \text{цифра} \rangle \Rightarrow^{(3)} \\ &\Rightarrow^{(3)} \langle \underline{123} \rangle \end{aligned}$$

Используя правила грамматики, мы вывели из начального нетерминального символа $\langle \text{целое без знака} \rangle$ строку из терминальных символов. Следует отметить, что 123 в данном случае рассматривается не как целое число, а как синтаксическая конструкция, которая состоит из терминальных символов – цифр.

Определение 4.1. Языком $L(G[Z])$ называется множество *терминальных* строк (цепочек) x , порождаемых грамматикой $G[Z]$, таких что x итерационно выводится из начального символа Z .

Цепочки x называются *конечной сентенциальной формой*.

Строгое математическое определение языка имеет вид:

$$L(G[Z]) = \{x, Z \Rightarrow^* x, x \in V_T^*\}.$$

Обратите внимание, что здесь выполняется *итерация* над терминальным словарем. Это приводит к возможности существования пустых цепочек в языке.

4.1 Операции над языками

Согласно определению 4.1, язык представляет собой множество. Определим операции над этим множеством.

Определение 4.2. Объединением языков L_1 и L_2 называется новый язык, который включает строки как языка L_1 , так и языка L_2 .

$$L_1 \cup L_2 = \{x \mid x \in L_1, x \in L_2\}.$$

Определение 4.3. Конкатенацией языков L_1 и L_2 называется новый язык, строки которого образуются конкатенацией строк из L_1 и L_2 .

$$L_1 L_2 = \{xy \mid x \in L_1, y \in L_2\}$$

Длина строк нового языка равна сумме длин строк x и y .

Определение 4.4. Итерация языка L представляет новый язык L^* , который включает строки из всевозможных комбинаций строк из L любой длины (включая пустую цепочку).

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots \cup L^n \cup \dots$$

Определение 4.5. Усеченной итерацией языка L является новый язык L^+ , который включает строки из всевозможных комбинаций строк из L любой длины, начиная с 1 (исключая пустую цепочку).

$$L^+ = \bigcup_{i=1}^{\infty} L^i.$$

Содержательная интерпретация нового языка L^+ та же, что и для L^* , только без пустой цепочки.

Для произвольного формального языка L справедливо утверждение:

$$L = \{\varepsilon\} L = L \{\varepsilon\}.$$

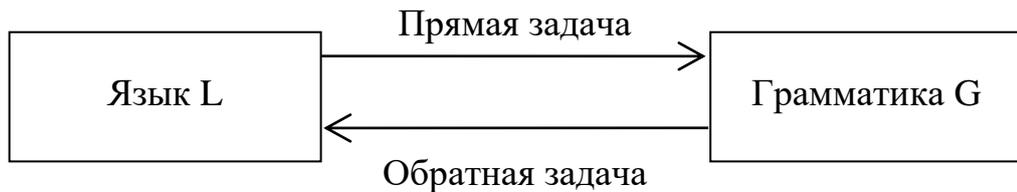
Обратите внимание, что множество, содержащее пустую строку (ε), существенно отличается от пустого множества (\emptyset). Действительно, в операциях над множествами можно записать:

$$L\emptyset = \emptyset L = \emptyset.$$

Таким образом, новый язык может породиться не только грамматикой, но и операциями над исходным языком.

4.2 Прямая и обратная задачи формальных языков и грамматик

Задача построения грамматик по языку является прямой задачей. Соответственно выбор языка по грамматике является обратной задачей.



4.2.1 Обратная задача

Рассмотрим примеры обратных задач.

Пример 4.1: Пусть грамматика $G_1[A]$ задана набором:

$$V_T = \{0,1\},$$

$$V_N = \{A\},$$

$$P: A \rightarrow 01$$

$$A \rightarrow 0A1$$

или через альтернативную конструкцию: $A \rightarrow 01 \mid 0A1$.

Требуется определить язык, порождаемый грамматикой $G_1[A]$, т.е. $L(G_1[A]) - ?$

Определим множество цепочек, порождаемых грамматикой $G_1[A]$. Для этого продемонстрируем процедуру выводимости в соответствии с правилами вывода P:

$$\begin{array}{ccccccc}
 A & \Rightarrow & 0A1 & \Rightarrow & 00A11 & \Rightarrow & 000A111 & \dots \\
 \Downarrow & & \Downarrow & & \Downarrow & & \Downarrow & \\
 01 & & 0011 & & 000111 & & \dots &
 \end{array}$$

Очевидно, что это множество включает вполне упорядоченный набор цепочек: $\{01, 0011, 000111, \dots\}$. Отсюда легко видеть, что[^]

$$L(G_1[A]) = \{0^n 1^n \mid n \geq 1\}.$$

Определение 4.6. Грамматики G и G_1 являются *эквивалентными*, если они порождают один и тот же язык, т.е. имеет место тождество множеств $L(G) = L(G_1)$.

Определение 4.7. Грамматики G и G_1 являются *почти эквивалентными*, если порождаемые ими языки отличаются не более чем на пустую цепочку ε , т.е. имеет место $L(G) = \{L(G_1) \cup \{\varepsilon\}\}$.

Пример 4.2: Рассмотрим грамматику $G_1[A]$ из примера 4.1 и грамматику $G_2[A]$, которая имеет свои правила P :

$$A \rightarrow 0B1$$

$$0B \rightarrow 00B1$$

$$B \rightarrow \varepsilon.$$

Тем же способом выводимости можно показать эквивалентность грамматик $G_1[A]$ и $G_2[A]$. Действительно,

$$\begin{array}{ccccccc} A & \Rightarrow & 0B1 & \Rightarrow & 00B11 & \Rightarrow & 000B111 \dots \\ & & \downarrow & & \downarrow & & \downarrow \\ & & 01 & & 0011 & & 000111 \dots \end{array}$$

Грамматика $G_1[A]$ обладает определенными преимуществами перед $G_2[A]$:

Во-первых, в грамматике $G_1[A]$ имеется меньшее количество нетерминальных символов и меньшее количество продукций. А это связано с меньшим количеством шагов при выводимости определенных цепочек языка.

Во-вторых, отсутствие правил вида $B \rightarrow \varepsilon$ делает грамматику более определенной.

Наконец, при разработке грамматик следует стремиться, чтобы в левой части продукций стоял нетерминальный символ, а не смешанная цепочка из терминалов и нетерминалов, как это приведено во втором правиле $G_2[A]$.

Пример 4.3: Пусть $G_3[Z]$ определена на множестве правил вывода P :

$$Z \rightarrow aA$$

$$A \rightarrow bB$$

$$B \rightarrow c$$

$$A \rightarrow Bb.$$

По аналогии с предыдущим примером воспользуемся выводимостью и проследим возможные цепочки языка:

$$\begin{array}{c} Z \Rightarrow aA \Rightarrow abB \Rightarrow abc \\ \Downarrow \\ aBb \Rightarrow acb \end{array}$$

Поэтому легко видеть, что язык $L(G_3[Z])$ включает только 2 терминальные цепочки, а именно: $L(G_3[Z]) = \{abc, acb\}$.

Пример 4.4: Определим язык, порождаемый грамматикой $G_4[Z]$:

$$V_T = \{a, b\},$$

$$V_N = \{A, Z\},$$

$$P = \{ Z \rightarrow aA; A \rightarrow Ab \}.$$

Вновь воспользуемся выводимостью, так что:

$$Z \Rightarrow aA \Rightarrow aAb \Rightarrow aAbb \Rightarrow \dots$$

Таким образом, никогда не удастся избавиться от нетерминала A , поэтому грамматика $L(G_4[Z])$ не порождает терминальных цепочек, а язык $L(G_4[Z])$ является пустым $L(G_4[Z]) = \emptyset$.

Из приведенных примеров следует, что грамматики могут порождать языки, включающие бесконечное множество строк (G_1 и G_2), ограниченное множество цепочек (G_3), и вообще не иметь терминальных цепочек (G_4).

Упражнения

1. Сравните операции над языками и операции со словарями (лекция 2).

2. Определите языки, порождаемые следующими грамматиками:

$$2.1. P: \quad 1) Z \rightarrow 11XY0 \quad 2) X \rightarrow 1X \mid 1 \mid \varepsilon \quad 3) Y \rightarrow 1Y0 \mid \varepsilon$$

$$2.2. P: \quad 1) Z \rightarrow XY - X \quad 2) X \rightarrow 1X \mid 1 \mid \varepsilon \quad 3) Y \rightarrow 1Y0 \mid \varepsilon$$

$$2.3. P: \quad 1) Z \rightarrow AB \quad 2) A \rightarrow 1A0 \mid 10 \quad 3) B \rightarrow 1B0 \mid 1$$

$$2.4. P: \quad 1) Z \rightarrow A + B \quad 2) A \rightarrow aA \mid \varepsilon \quad 3) B \rightarrow bB \mid \varepsilon$$

$$2.5. P: \quad 1) Z \rightarrow A - B \quad 2) A \rightarrow 1A00 \mid 10 \quad 3) B \rightarrow aB \mid \varepsilon$$

- 2.6. P: 1) $I \rightarrow AA$ 2) $A \rightarrow a$ 3) $A \rightarrow aa$
- 2.7. P: 1) $I \rightarrow aABc$ 2) $I \rightarrow \$$ 3) $A \rightarrow Ab$
4) $A \rightarrow cIB$ 5) $B \rightarrow bB$ 6) $B \rightarrow a$
- 2.8. P: 1) $I \rightarrow aM$ 2) $M \rightarrow A$ 3) $A \rightarrow aA$
4) $A \rightarrow B$ 5) $B \rightarrow bB$ 6) $B \rightarrow b$
- 2.9. P: 1) $I \rightarrow aA$ 2) $I \rightarrow Ic$ 3) $I \rightarrow Ab$ 4) $A \rightarrow d$

Список использованных источников

1. Шорников Ю.В. Теория и практика языковых процессоров.