

1.3. Порождающие грамматики

[Гин, 1.1], [Сал, 2.1], [АхоУль, 2.1.2], [Гла, 1.2], [Лал, с. 159–161], [Бра, с. 32–36], [ГлаМел, с. 34–48], [ГорМол, с. 354–355, 367–370], [СокКушБад, с. 12–13], [ТраБар, 1.12], [LewPar2, 4.6], [Рей, с. 28–30], [КукБей, с. 264–268]

Определение 1.41. Порождающей грамматикой (грамматикой типа 0) (generative grammar, rewrite grammar) называется четвёрка $G = \langle N, \Sigma, P, S \rangle$, где N и Σ — конечные алфавиты, $N \cap \Sigma = \emptyset$, $P \subset (N \cup \Sigma)^+ \times (N \cup \Sigma)^*$, P конечно и $S \in N$. Здесь Σ — основной алфавит (терминальный алфавит), его элементы называются терминальными символами или терминалами (terminal), N — вспомогательный алфавит (нетерминальный алфавит), его элементы называются нетерминальными символами, нетерминалами или переменными (nonterminal, variable), S — начальный символ (аксиома) (start symbol). Пары $(\alpha, \beta) \in P$ называются правилами подстановки, просто правилами или продукциями (rewriting rule, production) и записываются в виде $\alpha \rightarrow \beta$.

Пример 1.42. Пусть даны множества $N = \{S\}$, $\Sigma = \{a, b, c\}$, $P = \{S \rightarrow acSbcS, cS \rightarrow \varepsilon\}$. Тогда $\langle N, \Sigma, P, S \rangle$ является порождающей грамматикой.

Замечание 1.43. Будем обозначать элементы множества Σ строчными буквами из начала латинского алфавита, а элементы множества N — заглавными латинскими буквами. Обычно в примерах мы будем задавать грамматику в виде списка правил, подразумевая, что алфавит N составляют все заглавные буквы, встречающиеся в правилах, а алфавит Σ — все строчные буквы, встречающиеся в правилах. При этом правила порождающей грамматики записывают в таком порядке, что левая часть первого правила есть начальный символ S .

Замечание 1.44. Для обозначения n правил с одинаковыми левыми частями $\alpha \rightarrow \beta_1, \dots, \alpha \rightarrow \beta_n$ часто используют сокращённую запись $\alpha \rightarrow \beta_1 \mid \dots \mid \beta_n$.

Определение 1.45. Пусть дана грамматика G . Пишем $\phi \Rightarrow_G \psi$, если $\phi = \eta\alpha\theta$, $\psi = \eta\beta\theta$ и $(\alpha \rightarrow \beta) \in P$ для некоторых слов $\alpha, \beta, \eta, \theta$ в алфавите $N \cup \Sigma$.

Замечание 1.46. Когда из контекста ясно, о какой грамматике идёт речь, вместо \Rightarrow можно писать просто \Rightarrow .

Пример 1.47. Пусть

$$G = \langle \{S\}, \{a, b, c\}, \{S \rightarrow acSbcS, cS \rightarrow \varepsilon\}, S \rangle.$$

Тогда $cSacS \Rightarrow_G cSa$.

Определение 1.48. Если $\omega_0 \Rightarrow_G \omega_1 \Rightarrow_G \dots \Rightarrow_G \omega_n$, где $n \geq 0$, то пишем $\omega_0 \xRightarrow{*}_G \omega_n$ (другими словами, бинарное отношение $\xRightarrow{*}_G$ является рефлексивным, транзитивным замыканием бинарного отношения \Rightarrow_G , определённого на множестве $(N \cup \Sigma)^*$). При этом последовательность слов $\omega_0, \omega_1, \dots, \omega_n$ называется выводом (derivation) слова ω_n из слова ω_0 в грамматике G . Число n называется длиной (количеством шагов) этого вывода.

Замечание 1.49. В частности, для всякого слова $\omega \in (N \cup \Sigma)^*$ имеет место $\omega \xRightarrow{*}_G \omega$ (так как возможен вывод длины 0).

Пример 1.50. Пусть $G = \langle \{S\}, \{a, b\}, \{S \rightarrow aSa, S \rightarrow b\}, S \rangle$. Тогда $aSa \xRightarrow{*}_G aaaaSaaaa$. Длина этого вывода — 3.

Определение 1.51. Язык, порождаемый грамматикой G , — это множество $L(G) = \{\omega \in \Sigma^* \mid S \xRightarrow{*}_G \omega\}$. Будем также говорить, что грамматика G порождает (generates) язык $L(G)$.

Замечание 1.52. Существенно, что в определении порождающей грамматики включены два алфавита — Σ и N . Это позволило нам в определении 1.51 “отсеять” часть слов, получаемых из начального символа. А именно, отбрасывается каждое слово, содержащее хотя бы один символ, не принадлежащий алфавиту Σ .

Пример 1.53. Если $G = \langle \{S\}, \{a, b\}, \{S \rightarrow aSa, S \rightarrow bb\}, S \rangle$, то $L(G) = \{a^n bba^n \mid n \geq 0\}$.

Определение 1.54. Две грамматики эквивалентны, если они порождают один и тот же язык.

Пример 1.55. Грамматика $S \rightarrow abS, S \rightarrow a$ и грамматика $T \rightarrow aU, U \rightarrow baU, U \rightarrow \varepsilon$ эквивалентны.

1.4. Классы грамматик

[Гин, с. 23–24, 78–79], [АхоУль, 2.1.3, с. 191], [Сал, 2.1, с. 94], [Гла, 1.2, 1.3], [Бра, с. 39–45], [ГлаМел, с. 54, 63, 69–70], [ГорМол, с. 361–367], [ТраБар, 1.12], [КукБей, с. 268–271], [ЛПИИ, 5.2.1]

Определение 1.56. Контекстной грамматикой (контекстно-зависимой грамматикой, грамматикой непосредственно составляющих, НС-грамматикой, грамматикой типа 1) (context-sensitive grammar, phrase-structure grammar) называется порождающая грамматика, каждое правило которой имеет вид $\eta A \theta \rightarrow \eta \alpha \theta$, где $A \in N$, $\eta \in (N \cup \Sigma)^*$, $\theta \in (N \cup \Sigma)^*$, $\alpha \in (N \cup \Sigma)^+$.

Пример 1.57. Грамматика $S \rightarrow TS$, $S \rightarrow US$, $S \rightarrow b$, $Tb \rightarrow Ab$, $A \rightarrow a$, $TA \rightarrow AAT$, $UAb \rightarrow b$, $UAAA \rightarrow AAU$ не является контекстной (последние три правила не имеют требуемого вида).

Определение 1.58. Контекстно-свободной грамматикой (КС-грамматикой, бесконтекстной грамматикой, грамматикой типа 2) (context-free grammar) называется порождающая грамматика, каждое правило которой имеет вид $A \rightarrow \alpha$, где $A \in N$, $\alpha \in (N \cup \Sigma)^*$.

Пример 1.59. Грамматика $S \rightarrow ASTA$, $S \rightarrow AbA$, $A \rightarrow a$, $bT \rightarrow bb$, $AT \rightarrow UT$, $UT \rightarrow UV$, $UV \rightarrow TV$, $TV \rightarrow TA$ является контекстной, но не контекстно-свободной (последние пять правил не имеют требуемого вида).

Определение 1.60. Линейной грамматикой (linear grammar) называется порождающая грамматика, каждое правило которой имеет вид $A \rightarrow u$ или $A \rightarrow uBv$, где $A \in N$, $u \in \Sigma^*$, $v \in \Sigma^*$, $B \in N$.

Пример 1.61. Грамматика $S \rightarrow TT$, $T \rightarrow eTT$, $T \rightarrow bT$, $T \rightarrow a$ является контекстно-свободной, но не линейной (первые два правила не имеют требуемого вида).

Определение 1.62. Праволинейной грамматикой (рациональной грамматикой, грамматикой типа 3) (right-linear grammar) называется порождающая грамматика, каждое правило которой имеет вид $A \rightarrow u$ или $A \rightarrow uB$, где $A \in N$, $u \in \Sigma^*$, $B \in N$.

Пример 1.63. Грамматика $S \rightarrow aSa$, $S \rightarrow T$, $T \rightarrow bT$, $T \rightarrow \varepsilon$ является линейной, но не праволинейной (первое правило не имеет требуемого вида).

Пример 1.64. Грамматика $S \rightarrow T$, $U \rightarrow abba$ праволинейная.

Пример 1.65. Грамматика $S \rightarrow aS$, $S \rightarrow bS$, $S \rightarrow aaaT$, $S \rightarrow aabaT$, $S \rightarrow abaaT$, $S \rightarrow aabbaT$, $S \rightarrow ababaT$, $S \rightarrow abbbaT$, $T \rightarrow aT$, $T \rightarrow bT$, $T \rightarrow \varepsilon$ праволинейная.

Пример 1.66. Грамматика $S \rightarrow \varepsilon$, $S \rightarrow aaaS$, $S \rightarrow abbS$, $S \rightarrow babS$, $S \rightarrow aabT$, $T \rightarrow abaT$, $T \rightarrow baaT$, $T \rightarrow bbbT$, $T \rightarrow bbaS$ праволинейная. Обобщенный вариант языка, порождаемого этой грамматикой, используется в доказательстве разрешимости арифметики Пресбургера [Sip, с. 207–208].

Определение 1.67. Правила вида $\alpha \rightarrow \varepsilon$ называются ε -правилами.

Лемма 1.68. Каждая праволинейная грамматика является линейной. Каждая линейная грамматика является контекстно-свободной. Каждая контекстно-свободная грамматика без ε -правил является контекстной грамматикой.

Определение 1.69. Классы грамматик типа 0, 1, 2 и 3 образуют иерархию Хомского (Chomsky hierarchy).

Определение 1.70. Язык называется контекстным языком (контекстно-свободным языком, линейным языком, праволинейным языком), если он порождается некоторой контекстной грамматикой (соответственно контекстно-свободной грамматикой, линейной грамматикой, праволинейной грамматикой). Контекстно-свободные языки называются также алгебраическими языками.

Пример 1.71. Пусть дан произвольный алфавит $\Sigma = \{a_1, \dots, a_n\}$. Тогда язык Σ^* является праволинейным, так как он порождается грамматикой $S \rightarrow \varepsilon$, $S \rightarrow a_1S$, \dots , $S \rightarrow a_nS$.

Правила обычно записывают в виде $\alpha \rightarrow \beta$ или $\alpha ::= \beta$ что означает α порождает (состоит из) β .

Язык, порождаемый грамматикой. - это множество цепочек, которые состоят только из терминалов и выводятся, начиная с одного, особо выделенного, нетерминала S , называемого начальным символом или аксиомой грамматики. Среди множества правил грамматики должно присутствовать хотя бы одно правило $S \rightarrow \beta$.

Рассмотрим ряд грамматик и обсудим алгоритм порождения, применяемый для вывода цепочек языка.

Пример 1.1.

$\langle \text{предложение} \rangle ::= \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$
 $\langle \text{подлежащее} \rangle ::= \text{мать} \mid \text{отец}$
 $\langle \text{группа сказуемого} \rangle ::= \langle \text{сказуемое} \rangle \langle \text{дополнение} \rangle$
 $\langle \text{сказуемое} \rangle ::= \text{любит} \mid \text{обожает} \mid \text{боготворит}$
 $\langle \text{дополнение} \rangle ::= \text{сына} \mid \text{дочь} \mid \square$

Если имеется множество правил, то ими можно воспользоваться для того, чтобы **вывести** или **породить** цепочку (предложение) по следующей схеме. Начнем с начального символа грамматики - $\langle \text{предложение} \rangle$, найдем правило, в котором $\langle \text{предложение} \rangle$ слева от $::=$, и подставим вместо $\langle \text{предложение} \rangle$ цепочку, которая расположена справа от $::=$, то есть

$\langle \text{предложение} \rangle \Rightarrow \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$.

Таким образом, мы заменяем синтаксическое понятие на одну из цепочек, из которых оно может состоять. Повторим процесс. Возьмем один из нетерминалов в цепочке $\langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$, например $\langle \text{подлежащее} \rangle$; найдем правило, где $\langle \text{подлежащее} \rangle$ находится слева от $::=$, и заменим $\langle \text{подлежащее} \rangle$ в исходной цепочке на соответствующую цепочку, которая находится справа от $::=$. Это даст

$\langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle \Rightarrow \text{мать} \langle \text{группа сказуемого} \rangle$.

Символ " \Rightarrow " означает, что один символ слева от \Rightarrow в соответствии с правилом грамматики заменяется цепочкой, находящейся справа от \Rightarrow . Полный вывод одного предложения будет таким:

$\langle \text{предложение} \rangle \Rightarrow \langle \text{подлежащее} \rangle \langle \text{группа сказуемого} \rangle$
 $\Rightarrow \text{мать} \langle \text{группа сказуемого} \rangle$
 $\Rightarrow \text{мать} \langle \text{сказуемое} \rangle \langle \text{дополнение} \rangle$
 $\Rightarrow \text{мать любит} \langle \text{дополнение} \rangle$
 $\Rightarrow \text{мать любит сына}.$

Этот вывод предложения запишем сокращенно, используя новый символ \Rightarrow^+ : $\langle \text{предложение} \rangle \Rightarrow^+ \text{мать любит сына}.$

На каждом шаге можно заменить **любой** нетерминал. В приведенном выше выводе всегда заменялся самый левый из них.

Вывод, на каждом шаге которого заменяется самый левый нетерминал сентенциальной формы, называется **левым (левосторонним) выводом**. Существует и часто используется также **правый (правосторонний) вывод**, который получается, если в сентенциальной форме заменять всегда самый правый нетерминал.

Обратите внимание на то, что предложенная грамматика используется для описания многих предложений. Девять правил грамматики, если считать каждую альтернативу за отдельное правило, а так оно и есть, определяют двенадцать предложений (цепочек) языка:

мать любит сына	мать обожает сына	мать боготворит сына
мать любит дочь	мать обожает дочь	мать боготворит дочь
отец любит сына	отец обожает сына	отец боготворит сына
отец любит дочь	отец обожает дочь	отец боготворит дочь

Одно из назначений грамматики как раз и состоит в том, чтобы описывать **все** цепочки языка с помощью приемлемого числа правил. Это особенно важно, если учесть, что количество предложений в языке, чаще всего, бесконечно.

Рассмотрим еще один пример полезной грамматики

Пример 1.2. Грамматика целого числа без знака содержит следующие 13 правил

(1) $\langle \text{число} \rangle ::= \langle \text{чс} \rangle$	$S \rightarrow A$
(2) $\langle \text{чс} \rangle ::= \langle \text{цифра} \rangle \langle \text{чс} \rangle$	$A \rightarrow AB$
(3) $\langle \text{чс} \rangle ::= \langle \text{цифра} \rangle$	$A \rightarrow B$
(4) $\langle \text{цифра} \rangle ::= 0$	$B \rightarrow 0$
(5) $\langle \text{цифра} \rangle ::= 1$	$B \rightarrow 1$
(6) $\langle \text{цифра} \rangle ::= 2$	$B \rightarrow 2$
(7) $\langle \text{цифра} \rangle ::= 3$	$B \rightarrow 3$
(8) $\langle \text{цифра} \rangle ::= 4$	$B \rightarrow 4$
(9) $\langle \text{цифра} \rangle ::= 5$	$B \rightarrow 5$
(10) $\langle \text{цифра} \rangle ::= 6$	$B \rightarrow 6$
(11) $\langle \text{цифра} \rangle ::= 7$	$B \rightarrow 7$
(12) $\langle \text{цифра} \rangle ::= 8$	$B \rightarrow 8$
(13) $\langle \text{цифра} \rangle ::= 9$	$B \rightarrow 9$

Пусть G - грамматика. Будем говорить, что цепочка α непосредственно порождает цепочку β , и обозначим $\alpha \Rightarrow \beta$, если для некоторых цепочек φ и ψ можно написать $\alpha = \varphi U \psi$, $\beta = \varphi U' \psi$, где $U ::= \gamma$ правило грамматики G . Будем также говорить, что β непосредственно выводима из α или что β непосредственно приводится (редуцируется, сворачивается) к α .

Цепочки φ и ψ конечно, могут быть пустыми. Следовательно, для любого правила $A \rightarrow \alpha$ грамматики G имеет место $A \Rightarrow \alpha$. На рис. 1.1 даны некоторые примеры непосредственных выводов для грамматики $G(\langle \text{число} \rangle)$ из примера 1.2 и обозначений предыдущего определения.

Будем говорить, что α порождает β или β приводится к α и записывать $\alpha \Rightarrow^+ \beta$, если существует последовательность непосредственных выводов $\alpha = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_n = \beta$, где $n > 0$. Эта последовательность называется **выводом длины n**. Будем писать $\alpha \Rightarrow^* \beta$, если $\alpha \Rightarrow^+ \beta$ или $\alpha = \beta$.

α	β	Номера правил	φ	ψ
$\langle \text{число} \rangle$	$\Rightarrow \langle \text{чс} \rangle$	1	ϵ	ϵ
$\langle \text{чс} \rangle$	$\Rightarrow \langle \text{цифра} \rangle \langle \text{чс} \rangle$	2	ϵ	ϵ
$\langle \text{цифра} \rangle \langle \text{чс} \rangle$	$\Rightarrow 2 \langle \text{чс} \rangle$	6	ϵ	$\langle \text{чс} \rangle$
$2 \langle \text{чс} \rangle$	$\Rightarrow 2 \langle \text{цифра} \rangle$	3	2	ϵ
$2 \langle \text{цифра} \rangle$	$\Rightarrow 25$	9	2	ϵ

Рисунок 1.1. Вывод для грамматики целых чисел без знака

Если просмотреть все строки на рисунке 1.1, то мы получим $\langle \text{число} \rangle \Rightarrow \langle \text{чс} \rangle \Rightarrow \langle \text{цифра} \rangle \langle \text{чс} \rangle \Rightarrow 2 \langle \text{чс} \rangle \Rightarrow 2 \langle \text{цифра} \rangle \Rightarrow 25$.

Таким образом, $\langle \text{число} \rangle \Rightarrow^+ 25$ и длина вывода равна 5. (Если длина вывода известна можно, записывать в явном виде $\langle \text{число} \rangle \Rightarrow^5 25$.)

Заметим, что пока в цепочке есть хотя бы один нетерминал, из нее можно вывести новую цепочку, однако если нетерминальные символы отсутствуют, то вывод завершен. Неслучайно "терминалом" (*terminal* - заключительный, конечный) называют символ, который не встречается в левой части ни одного из