

Лекція 7-8.

ТЕОРІЯ

ФОРМАЛЬНИХ ГРАМАТИК

4. ФОРМАЛЬНЫЕ ГРАММАТИКИ И ЯЗЫКИ

4.1. Общие сведения

В последние три десятилетия появилось большое количество работ по общей теории языков и грамматик. Можно выделить четыре научных направления, которые удалось объединить по методам их исследования в одну общую задачу теории языков.

Первое из этих направлений связано с построением формальной, или математической, лингвистики, которая начала особенно быстро развиваться в тот период, когда были сформулированы вопросы машинного перевода. Эта проблема требовала формализации понятий словарь, грамматика, язык, их классификации и умения относить конкретные словари, грамматики, языки к определенному классу.

Интерес к задачам такого рода еще более усилился с появлением искусственных языков программирования. С появлением трансляторов проблема перевода во многом определила построение общей теории вычислительных машин, а сами языки программирования стали все более приближаться к формально построенным математическим конструкциям.

Независимо от указанных двух направлений развивалось построение формальных моделей динамических систем. Для создания продуктивной теории эти модели должны были быть, с одной стороны, достаточно узкими, а с другой - достаточно широкими, чтобы охватить некоторый общий класс прикладных задач. Типичным примером такого рода является модель конечного автомата. Эта модель позволяет описать многие процессы, заданные на конечных множествах и развивающиеся в счетном времени, что позволило создать для нее продуктивную теорию. Если в эту модель ввести бесконечность по какому-либо параметру, то это приводит к общему классу систем, например, к машинам Тьюринга.

Независимо и параллельно развивалась общая теория алгоритмов как ветвь современной математики. Развитие вычислительной техники поставило перед математической теорией алгоритмов новую задачу: стало необходимым классифицировать алгоритмы по различным признакам. Эквивалентность понятий "алгоритм" и "машина Тьюринга" позволила предположить, что поиски классификации алгоритмов окажутся связанными с поисками промежуточных моделей между моделями конечного автомата и машиной Тьюринга.

Таким образом, перечисленные четыре направления оказались тесно связанными. Теория языков, порожденная чисто лингвистическими задачами, оказалась в центре интересов математиков, занимающихся теорией алгоритмов, и ученых, разрабатывающих абстрактные модели динамических систем и теоретические основы автоматизации.

Теория формальных языков и грамматик является разделом математической лингвистики - специфической математической дисциплины, ориентированной на изучение структуры естественных и искусственных языков. По характеру используемого математического аппарата теория формальных грамматик и языков близка к теории алгоритмов и теории автоматов.

На интуитивном уровне язык можно определить как некоторое множество предложений с заданной структурой, имеющих определенное значение. Правила, определяющие допустимые конструкции языка, составляют синтаксис языка. Значение, или смысл фразы, определяется семантикой языка.

Если бы все языки состояли из конечного числа предложений, то не было бы проблемы синтаксиса. Но, так как язык содержит неограниченное (или довольно большое) число правильно построенных фраз, то возникает проблема описания бесконечных языков с помощью каких-либо конечных средств.

Различают следующие виды описания языков:

- 1) порождающее, которое предполагает наличие алгоритма, последовательно порождающего все правильные предложения языка;
- 2) распознающее, которое предполагает наличие алгоритма, определяющего принадлежность любой фразы данному языку.

4.2. Основные понятия порождающих грамматик

Алфавит - это непустое конечное множество. Элементы алфавита называются символами.

Цепочка над алфавитом $\Sigma = \{a_1, a_2, \dots, a_n\}$ есть конечная последовательность элементов a_i . Множество всех цепочек над алфавитом Σ обозначается Σ^* .

Длина цепочки x равна числу ее элементов и обозначается $|x|$. Цепочка нулевой длины называется пустой и обозначается символом ε . Соответственно, непустая цепочка определяется как цепочка ненулевой длины. Пусть имеется алфавит $\Sigma = \{a, b\}$. Тогда множество всех цепочек определяется следующим образом: $\Sigma^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$.

Цепочки x и y равны, если они одинаковой длины и совпадают с точностью до порядка символов, из которых они состоят.

Над цепочками x и y определена операция сцепления (конкатенации, произведения), которая получается дописыванием символов цепочки y за символами цепочки x .

Язык L над алфавитом Σ представляет собой множество цепочек над Σ . Необходимо различать пустой язык $L=\emptyset$ и язык, содержащий только пустую цепочку: $L=\{\varepsilon\} \neq \emptyset$.

Формальный язык L над алфавитом Σ - это язык, выделенный с помощью конечного множества некоторых формальных правил.

Пусть M и L - языки над алфавитами. Тогда конкатенация $LM = \{xy \mid x \in L, y \in M\}$. В частности, $\{\varepsilon\}L = L\{\varepsilon\} = L$. Используя понятие произведения, определим итерацию L^* и усеченную итерацию L^+ множества L :

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

$$L^* = \bigcup_{i=0}^{\infty} L^i,$$

где i - степень языка, L определяется рекурсивно следующим образом:

$$L^0 = \{\varepsilon\};$$

$$L^1 = L;$$

$$L^{n+1} = L^n L = L L^n;$$

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

Например, если задан язык $L=\{a\}$, тогда $L^* = \{\varepsilon, a, aa, aaa, \dots\}$, $L^+ = \{a, aa, aaa, \dots\}$.

Порождающая грамматика - это упорядоченная четверка $G = (V_T, V_N, P, S)$, где V_T - конечный алфавит, определяющий множество терминальных символов; V_N - конечный алфавит, определяющий множество нетерминальных символов; P - конечное множество правил вывода, т.е. множество пар следующего вида:

$$u \rightarrow v, \quad \text{где } u, v \in (V_T \cup V_N)^*;$$

S - начальный символ (аксиома грамматики), $S \in V_N$.

Из терминальных символов состоят цепочки языка, порожденного грамматикой. Аксиомой называется символ в левой части первого правила вывода грамматики.

Для того чтобы различать терминальные и нетерминальные символы, принято обозначать терминальные символы строчными, а нетерминальные символы заглавными буквами латинского алфавита.

В грамматике G цепочка x непосредственно порождает цепочку y , если $x = \alpha u \beta$, $y = \alpha v \beta$ и $u \rightarrow v \in P$, т.е. цепочка y непосредственно выводится из x , что обозначается $x \Rightarrow y$.

Языком, порождаемым грамматикой $G = (V_T, V_N, P, S)$, называется множество терминальных цепочек, выводимых в грамматике G из аксиомы:

$L(G) = \{x \mid x \in V_T^*; S \Rightarrow^* x\}$, где \Rightarrow^* - выводимость.

Пример. Дана грамматика $G = (V_T, V_N, P, S)$, в которой $V_T = \{a, b\}$, $V_N = \{S\}$, $P = \{S \rightarrow aSb, S \rightarrow ab\}$. Определить язык, порождаемый этой грамматикой.

Решение. Используя рекурсию, выведем несколько цепочек языка, порождаемого данной грамматикой:

$S \rightarrow ab$;

$S \rightarrow aSb \Rightarrow aabb$;

$S \rightarrow aSb \Rightarrow aaSbb \Rightarrow aaabbb$; . . .

Определим язык, порожденный данной грамматикой:

$L(G) = \{a^n b^n \mid n > 0\}$.

Говоря о представлении грамматик, нужно отметить, что множество правил вывода грамматики может приводиться без специального указания множеств терминалов и нетерминалов. В таком случае обычно предполагается, что грамматика содержит в точности те терминальные и нетерминальные символы, которые встречаются в правилах вывода.

Также предполагается, что правые части правил, для которых совпадают левые части, можно записать в одну строку с использованием разделителя. Так, правила вывода из приведенного примера можно записать следующим образом: $S \rightarrow aSb \mid ab$.

4.3. Классификация грамматик

Правила порождающих грамматик позволяют осуществлять преобразования строк. Ограничения же на виды правил позволяют выделить классы грамматик. Рассмотрим классификацию, которую предложил Н. Хомский.

Грамматики типа 0 - это грамматики, на правила вывода которых нет ограничений.

Грамматики типа 1, или контекстные грамматики - это грамматики, все правила которых имеют следующий вид: $xAu \rightarrow x\varphi u$, где $A \in V_N$, $x, u, \varphi \in (V_N \cup V_T)^+$.

Грамматики типа 2 - это бесконтекстные, или контекстно-свободные грамматики (КС-грамматики). Правила вывода для этих грамматик имеют следующий вид: $A \rightarrow \varphi$, где $A \in V_N$, $\varphi \in (V_N \cup V_T)^*$.

Грамматики типа 3 - это автоматные грамматики, которые делятся на два типа:

а) левосторонние (леворекурсивные), правила вывода для которых имеют следующий вид: $A \rightarrow Aa \mid a$, где $A \in V_N$;

б) правосторонние (праворекурсивные), правила вывода для которых имеют следующий вид: $A \rightarrow Aa \mid a$.

Язык L называется языком типа i , если существует грамматика типа i , порождающая язык L .

4.3.1. Методика решения задач

Пример 1. Дана порождающая грамматика $G = (V_T, V_N, P, S)$, в которой $V_T = \{a, d, e\}$, $V_N = \{B, C, S\}$, $P = \{S \rightarrow aB, B \rightarrow Cd \mid dC, C \rightarrow e\}$. Выписать терминальные цепочки, порожденные данной грамматикой, и определить длину их вывода.

Решение. Получим следующие терминальные цепочки:

$S \rightarrow aB \rightarrow aCd \rightarrow aed$,

$S \rightarrow aB \rightarrow adC \rightarrow ade$,

длина вывода которых равна 3.

Пример 2. Дана грамматика $G = (V_T, V_N, P, C)$, в которой $V_T = \{a, b, c, d, e\}$, $V_N = \{A, B, C, D, E\}$, $P = \{A \rightarrow ed, B \rightarrow Ab, C \rightarrow Bc, C \rightarrow dD, D \rightarrow Ae, E \rightarrow bc\}$. Определить, принадлежит ли языку $L(G)$ цепочка $eadabcb$.

Решение. Выведем возможные терминальные цепочки из аксиомы с помощью заданных правил вывода:

$$C \rightarrow Bc \rightarrow Abc \rightarrow edbc,$$

$$C \rightarrow dD \rightarrow dAe \rightarrow dede.$$

Так как первые же терминальные символы полученных цепочек не совпадают с заданной, можно сделать вывод о том, что цепочка $eadabcb$ не принадлежит языку $L(G)$.

Пример 3. Построить КС-грамматику (грамматику типа 2), порождающую цепочки из 0 и 1 с одинаковым числом тех и других символов.

Решение. Определим множества, задающие грамматику: $V_T = \{0, 1\}$; $V_N = \{S\}$; $P = \{S \rightarrow 0S1, S \rightarrow 1S0, S \rightarrow \epsilon, S \rightarrow S01, S \rightarrow S10\}$.

Пример 4. Описать язык, порождаемый следующими правилами: $S \rightarrow bSS \mid a$.

Решение. Построим несколько терминальных цепочек, применяя заданные правила вывода:

$S \rightarrow a$;

$S \rightarrow bSS \rightarrow baa$;

$S \rightarrow bSS \rightarrow bbSSS \rightarrow bbaaaa$;

$S \rightarrow bSS \rightarrow bbSSbSS \rightarrow bbaabaa$;

$S \rightarrow bSS \rightarrow bbSSbSS \rightarrow bbbSSbSSbbSSbSS \rightarrow \dots$

Из полученных цепочек видно, что:

а) цепочки всегда начинаются с терминала b , кроме аксиомы, и заканчиваются терминалом a ;

б) количество терминалов a в любой цепочке всегда на 1 больше, чем b .

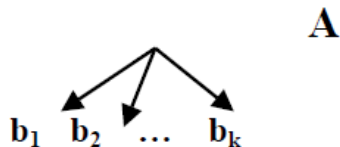
Исходя из этих выводов, определим язык, порождаемый заданными правилами, следующим образом:

$L(G) = \{\alpha \mid \alpha \in \{a, b\}^*, |\alpha| = |b| + 1, \text{ причем цепочки начинаются с терминала } b \text{ и заканчиваются терминалом } a.\}$

4.4. Грамматический разбор

В КС-грамматике может быть несколько выводов, эквивалентных в том смысле, что в них применяются одни и те же правила к одинаковым в промежуточном выводе цепочкам, но в различном порядке. Например, для грамматики G с правилами вывода $S \rightarrow ScS|b|a$ возможны следующие эквивалентные выводы терминальной цепочки acb : $S \rightarrow ScS \rightarrow Scb \rightarrow acb$ и $S \rightarrow ScS \rightarrow acS \rightarrow acb$.

Деревом вывода цепочки x в КС-грамматике $G = (V_T, V_N, P, S)$ называется упорядоченное дерево, каждая вершина которого помечена символом из множества $V \cup V_N \cup \{\varepsilon\}$ так, что каждому правилу $A \rightarrow b_1b_2b_3 \dots b_k$, используемому при выводе цепочки x , в дереве вывода соответствует поддерево с корнем A и прямыми потомками $b_1, b_2, b_3, \dots, b_k$, которое приведено на рисунке:



В силу того, что цепочка $x \in L(G)$ выводится из аксиомы S , корнем вывода всегда является аксиома. Внутренние узлы вывода соответствуют нетерминальным символам. Концевые вершины дерева вывода - листья - это вершины, не имеющие потомков. Такие вершины соответствуют либо терминалам, либо пустым символам (ϵ) при условии, что среди правил грамматики имеются правила с пустой правой частью. При чтении слева направо концевые вершины дерева вывода образуют цепочку x .

Дерево вывода часто называют деревом грамматического разбора, или синтаксическим деревом, а процесс построения дерева вывода - грамматическим разбором (синтаксическим анализом). Одной цепочке языка может соответствовать больше одного дерева, так как эта цепочка может иметь разные выводы, порождающие разные деревья.

КС-грамматика $G = (V_T, V_N, P, S)$ называется неоднозначной (неопределенной), если существует цепочка $x \in L(G)$, имеющая два или более дерева вывода.

Рассмотрим пример.

Пусть даны две КС-грамматики:

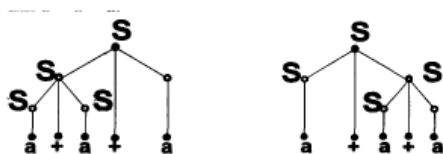
$G_1 = (V_T, V_N, P_1, S), V_N = \{S\},$

$P_1 = \{S \rightarrow S+S \mid S \mid S \cdot S \mid (S) \mid a\};$

$G_2 = (V_T, V_N, P_2, S), V_N = \{S, A, B\},$

$P_2 = \{S \rightarrow S + A \mid A \mid A, A \rightarrow A \cdot B \mid A / B \mid B, B \rightarrow a \mid (S)\},$ содержащие в множестве терминальных символов знаки операций, круглые скобки и символ a . Определить, являются ли грамматики однозначными. Если какая-либо из них неоднозначна, привести пример цепочки, для которой существует два различных дерева вывода.

Решение: Грамматика G_1 является неоднозначной, так как она имеет правила с правой частью, содержащей два вхождения нетерминала S и знак арифметической операции. Построим два дерева вывода для простейшей цепочки $a + a + a$:



Грамматика G_2 является однозначной, так как не содержит правил с двойным вхождением нетерминального символа. Так, для цепочки $a + a + a$ она имеет только одно дерево вывода.

В некоторых случаях неоднозначность в грамматиках может устраняться путем эквивалентных преобразований. Однако в общем случае проблема устранения неоднозначности неразрешима. На практике от неоднозначности избавляются путем задания словесных ограничений, называемых контекстными условиями языка, которые включаются в его семантику.

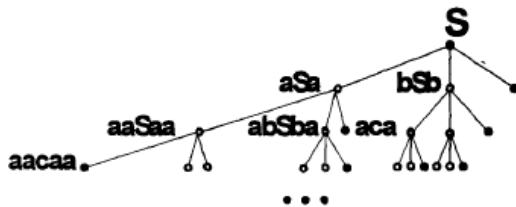
Различают две стратегии грамматического разбора: восходящую и нисходящую, которые соответствуют способу построения синтаксических деревьев. При нисходящей стратегии разбора дерево строится от корня аксиомы вниз, к терминальным вершинам. Главной задачей при нисходящем разборе является выбор того правила, которое следует применить на рассматриваемом шаге. При восходящем разборе дерево строится от терминальных вершин к корню дерева (аксиоме).

Преобразование цепочки, обратное порождению, называется редукцией.

4.4.1. Представление грамматики в виде графа

Дерево грамматического разбора не следует путать с представлением грамматики в виде графа. Граф грамматики в качестве вершин содержит сентенциальные формы (любые цепочки, выводимые из аксиомы).

Рассмотрим представление грамматики G в виде графа: $G = (V_T, V_N, P, S)$, в которой $V_T = \{a, b, c\}$, $V_N = \{S\}$, $P = \{S \rightarrow aSa \mid bSb \mid c\}$.



4.5. Преобразования КС-грамматик

Часто требуется изменить грамматику таким образом, чтобы она удовлетворяла определенным требованиям, не изменяя при этом порождаемый грамматикой язык. Для этого используются эквивалентные преобразования КС-грамматик, некоторые из которых рассмотрены ниже.

4.5.1. Удаление правил вида $A \rightarrow B$

Преобразование первого типа состоит в удалении правил $A \rightarrow B$, или $\langle \text{нетерминал} \rangle \rightarrow \langle \text{нетерминал} \rangle$.

Покажем, что для любой КС-грамматики можно построить эквивалентную грамматику, не содержащую правил вида: $A \rightarrow B$, где A и B - нетерминальные символы.

Пусть имеется КС-грамматика $G=(V_T, V_N, P, S)$, где множество нетерминалов $V_N=\{A_1, A_2, \dots, A_n\}$. Разобьем P на два непересекающихся множества: $P = P_1 \cup P_2$. В P_1 включены все правила вида $A_i \rightarrow A_k$, в P_2 включены все остальные правила, т.е. $P_2 = P \setminus P_1$. Затем для каждого A_i определим множество правил $P(A_i)$, включив в него все такие правила $A_i \rightarrow \varphi$, что $A_i \rightarrow^* A_j$ и $A_j \rightarrow \varphi$, где $A_j \rightarrow \varphi \in P_2$. Построим эквивалентную КС-грамматику $G_3 = (V_T, V_N, P_3, S)$, в которой множества терминальных и нетерминальных символов, а также аксиома совпадают с теми же объектами исходной грамматики G , а множество правил P_3 получено объединением правил множества P_2 и правил $P(A_i)$ для всех $1 \leq i \leq n$:

$$P_3 = \bigcup_{i=1}^n P(A_i) \cup P_2.$$

Пример. Пусть задана грамматика G со следующими правилами вывода $S \rightarrow aFb \mid A$; $A \rightarrow aA \mid B$; $B \rightarrow aSb \mid S$; $F \rightarrow bc \mid bFc$.

Построим множества правил P_2 , $P(S)$, $P(A)$, $P(B)$, $P(F)$.

Определим правила для P_2 : $P_2 = \{S \rightarrow aFb; A \rightarrow aA; B \rightarrow aSb; F \rightarrow bc \mid bFc\}$.

Определим правила для $P(S)$: $S \Rightarrow A \Rightarrow B$ или $S \Rightarrow^* A$; $S \Rightarrow B$, где \Rightarrow^* обозначает непосредственную выводимость. $P(S) = \{S \rightarrow aA; S \rightarrow aSb\}$.

Определим правила для $P(A)$: $A \Rightarrow B \Rightarrow S$ или $A \Rightarrow^* B$; $A \Rightarrow S$. $P(A) = \{A \rightarrow aSb; A \rightarrow aFb\}$.

Определим правила для $P(B)$: $B \Rightarrow S \Rightarrow A$ или $B \Rightarrow^* S$; $B \Rightarrow^* A$. $P(B) = \{B \rightarrow aFb; B \rightarrow aA\}$.

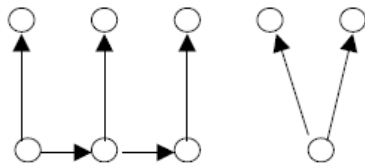
Определим правила для $P(F)$: так как непосредственно выводимых нетерминалов не существует, то $P(F) = \emptyset$.

Объединив полученные правила, можно записать грамматику G_2 , эквивалентную исходной:

$S \rightarrow aFb \mid aSb \mid aA;$	$A \rightarrow aA \mid aSb \mid aFb;$
$B \rightarrow aA \mid aSb \mid aFb;$	$F \rightarrow bc \mid bFc.$

Графическая модификация метода

Аналитическое преобразование по рассмотренному алгоритму оказывается довольно сложным. При автоматизированном преобразовании грамматик проще применить графическую модификацию этого метода. С этой целью каждому нетерминальному символу и каждой правой части правил из множества P_2 поставлена в соответствие вершина графа. Из вершины с меткой U в вершину с меткой V направлено ребро, если в грамматике существует правило $U \rightarrow V$.



В эквивалентную грамматику будут включены правила вида $A \rightarrow w$, $A \in V_N$; $w \in (V_T \cup V_N)^*$, если из вершины с меткой A существует путь в вершину с меткой w .

$S \rightarrow aFb \mid aSb \mid aA$;

$A \rightarrow aA \mid aSb \mid aFb$;

$B \rightarrow aA \mid aSb \mid aFb$;

$F \rightarrow bc \mid bFc$.

Получено то же множество правил P , что и аналитическим методом.

4.5.2. Построение неукорачивающей грамматики

Грамматика, не содержащая правил с пустой правой частью, называется неукорачивающей грамматикой.

В грамматике с правилами вида $A \rightarrow \epsilon$ длина выводимой цепочки при переходе от k -го шага к $(k+1)$ -му уменьшается. Поэтому грамматики с правилами вида $A \rightarrow \epsilon$ называются укорачивающими. Восходящий синтаксический разбор в укорачивающих грамматиках сложнее по сравнению с разбором в неукорачивающих грамматиках, т.к. при редукции необходимо отыскать такой фрагмент входной цепочки, в которую можно вставить пустой символ.

Покажем, что для произвольной КС-грамматики, порождающей язык без пустой цепочки, можно построить эквивалентную неукорачивающую КС-грамматику.

Построим множество всех нетерминальных символов грамматики $G=(V_T, V_N, P, S)$, из которых выводится пустая цепочка, выделив следующие множества:

$$W_1 = \{A \mid A \rightarrow \varepsilon \in P\},$$

$$W_{n+1} = W_n \cup \{B \mid B \rightarrow \varphi \in P, \varphi \in W_n^*\}.$$

Затем найдем множество правил эквивалентной грамматики в два этапа:

а) удалив из множества P исходной грамматики правила с пустой правой частью $P_1' = P \setminus \{A \rightarrow \varepsilon \mid A \rightarrow \varepsilon \in P\}$;

б) получив новые правила $A \rightarrow \varphi'$ после удаления из каждого правила исходной грамматики $A \rightarrow \varphi \in P$ те нетерминалы, которые вошли в множество W_n по правилу:

$$P_1'' = \{A \rightarrow \varphi' \mid A \rightarrow \varphi \in P'; \varphi = \varphi_1 B \varphi_2 \mid B \in W_n; \varphi' = \varphi_1 B \varphi_2\}.$$

Повторив п.б) для каждого нетерминала, принадлежащего множеству W_n , получим эквивалентную грамматику $G = (V_T, V_N, P_3, S)$.

Пример. Пусть задана грамматика G со следующими правилами вывода: $S \rightarrow AbA \mid cAb \mid Bb$; $A \rightarrow aAb \mid \varepsilon$; $B \rightarrow AA \mid a$. Необходимо:

- 1) построить множество нетерминалов, из которых выводится ε ;
- 2) построить неукорачивающую грамматику, эквивалентную исходной.

Для того чтобы построить множество всех нетерминалов грамматики, из которых выводится пустая цепочка, выделим следующие множества:

$$W_1 = \{A \mid A \rightarrow \varepsilon \in P\};$$

$$W_{m+1} = W_m \cup \{B \mid B \rightarrow \varphi \in P, \varphi \in W_m^*\}.$$

Так как мы имеем правило $A \rightarrow \varepsilon \in P$, то можно построить множество $W_1 = \{A\}$, включающее нетерминал A .

Построим множество W_2 . С нетерминалом A связан нетерминал B , т.е. существует правило $B \rightarrow AA$ и $A \in W_1$. Следовательно, $W_2 = \{A, B\}$.

$W_3 = W_2$, т. к. множество $W_3 = \{B \mid B \rightarrow \varphi \in P, \varphi \in W_m^*\}$ является пустым.

Исключив правило, содержащее пустую цепочку в правой части, получим неукорачивающую грамматику $G1$ следующего вида:

$$S \rightarrow AbA \mid cAb \mid Bb \mid bA \mid Ab \mid cb \mid b;$$

$$A \rightarrow aAb \mid ab;$$

$$B \rightarrow AA \mid A \mid a.$$

4.5.3. Построение грамматики с продуктивными нетерминалами

Нетерминальный символ A называется непродуктивным (непроизводящим), если он не порождает ни одной терминальной цепочки, т.е. не существует вывода $A \rightarrow^* x$, где $x \in V_T^*$. Поэтому представляет интерес удаление из грамматики всех непродуктивных нетерминальных символов. Рассмотрим, как для произвольной КС-грамматики можно построить эквивалентную КС-грамматику, все нетерминальные символы которой продуктивны. С этой целью выделяется множество $W_1 = \{A \mid A \rightarrow \varphi \in P, \varphi \in V_T^*\}$. Затем строится множество W_1, W_2, \dots, W_{n+1} по следующим правилам:

$$W_{n+1} = W_n \cup \{B \mid B \rightarrow x \in P, x \in (V_T \cup W_n)^*\}.$$

Пусть задана грамматика G со следующими правилами вывода: $S \rightarrow SA \mid BSb \mid bAb; A \rightarrow aSa \mid bb; B \rightarrow bBb \mid BaA$. Построим множество продуктивных нетерминалов:

$$W_1 = \{A\}, \text{ т.к. в множестве } P \text{ есть правило } A \rightarrow bb;$$

$$W_2 = \{A, S\}, \text{ т.к. имеется правило } S \rightarrow bAb \text{ и } A \in W_1;$$

$$W_3 = W_2.$$

Все символы множества $V_N \setminus W_n$ являются непродуктивными, не используются в выводе никакой терминальной цепочки и их можно удалить из грамматики вместе со всеми правилами, в которые они входят. Грамматика G_1 , эквивалентная исходной грамматике, будет включать следующее множество правил:

$$S \rightarrow SA \mid bAb; A \rightarrow aSa \mid bb.$$

В грамматике G_1 все нетерминалы продуктивны.

4.5.4. Построение грамматики, аксиома которой зависит от всех нетерминалов

Существует еще один тип нетерминальных символов, которые можно удалять из грамматики вместе с правилами, в которые они входят. Например, в грамматике G , заданной множеством правил $P: S \rightarrow aSb \mid ba; A \rightarrow aAa \mid bba$, нетерминал A не участвует ни в каком выводе, т.к. из аксиомы нельзя вывести цепочку, содержащую этот нетерминал. Поэтому из заданной грамматики можно удалить нетерминал A . Рассмотрим, как для произвольной КС-грамматики можно построить эквивалентную КС-грамматику, аксиома которой зависит от всех нетерминальных символов.

Для этого построим множество нетерминалов, от которых зависит аксиома. С этой целью выделим множества W_1, W_2, \dots, W_{n+1} по следующим правилам:

$$W_1 = \{S\},$$

$$W_{n+1} = W_n \cup \{B \mid A \rightarrow xBy \in P, A \in W_n\}.$$

Нетерминал $B \in W_n$ тогда и только тогда, когда аксиома S зависит от B . Все нетерминалы, не содержащиеся в W_n , можно удалить из грамматики вместе с правилами, в которые они входят.

Пример. Пусть дана КС-грамматика G :

$S \rightarrow abS \mid ASa \mid ab$; $A \rightarrow abAa \mid ab$; $B \rightarrow bAab \mid bB$.

Найдем нетерминалы, от которых зависит аксиома:

$W_1 = \{S\}$,

$W_2 = \{S, A\}$, т.к. имеется правило $S \rightarrow ASa$ и $S \in W_1$;

$W_3 = W_2$.

Эквивалентная грамматика G_1 , аксиома которой зависит от всех нетерминальных символов:

$S \rightarrow abS \mid ASa \mid ab$; $A \rightarrow abAa \mid ab$.

4.5.5. Удаление правил с терминальной правой частью

Пусть в грамматике G имеется с терминальной правой частью $A \rightarrow \beta$, где $\beta \in V_T^*$. Тогда любой вывод с использованием этого правила имеет вид: $S \xrightarrow{*} xAy \rightarrow x\beta y$. Здесь нетерминал A в сентенциальной форме xAy появился на предыдущем шаге вывода $B \rightarrow uAv$. Если в это правило вместо A подставить β , получим правило $B \rightarrow u\beta v$. Тогда длина вывода некоторой цепочки сократится на один шаг. Следовательно, для того чтобы удалить терминальное правило грамматики $A \rightarrow \beta$, необходимо учесть следующие правила:

- а) если для нетерминала A больше нет правил, тогда во всех правых частях A заменяется на β ;
- б) если для нетерминала A есть другие правила, тогда добавляются новые правила, в которых A заменяется на β .

Пример. Пусть дана КС-грамматика G:

$S \rightarrow aSb \mid bAa \mid B$; $A \rightarrow ABA \mid b \mid \epsilon$; $B \rightarrow ab \mid ba$.

Удалим правила для нетерминала B. Тогда эквивалентная грамматика G_1 будет включать следующие правила:

$S \rightarrow aSb \mid bAa \mid ab \mid ba$; $A \rightarrow Aaba \mid Abaa \mid b \mid \epsilon$.

4.5.6. Построение эквивалентной праворекурсивной КС-грамматики

Некоторые специальные методы грамматического разбора неприменимы к леворекурсивным или праворекурсивным грамматикам, поэтому рассмотрим устранение левой или правой рекурсии. В общем случае избавиться от рекурсии в правилах грамматики невозможно, т.к. бесконечные языки порождаются грамматиками с конечным числом правил только благодаря рекурсии. Поэтому можно говорить лишь о преобразовании одного вида рекурсии в другой. Рассмотрим, как для любой леворекурсивной КС-грамматики построить эквивалентную праворекурсивную КС-грамматику.

Пусть нетерминал A - леворекурсивен, т.е. для него имеются правила следующего вида:

$A \rightarrow Ax_1 \mid Ax_2 \mid \dots \mid Ax_p \mid w_1 \mid \dots \mid w_k$, где x_i и w_j - цепочки над множеством $V_T \cup V_N$.

Введем дополнительные нетерминалы B и D и указанные правила заменим на эквивалентные им:

$A \rightarrow AB \mid D$;

$B \rightarrow x_1 \mid x_2 \mid \dots \mid x_p$;

$D \rightarrow w_1 \mid w_2 \mid \dots \mid w_k$.

В результате замены получим вывод, который имеет вид:

$A \rightarrow AB \rightarrow AB^2 \rightarrow AB^3 \rightarrow \dots \rightarrow AB^* \rightarrow DB^*$.

Таким образом, для нетерминала A можно определить эквивалентные правила:

$A \rightarrow DK$;

$K \rightarrow BK \mid \varepsilon$.

Выполняя подстановку D и B в эти правила, получим следующие праворекурсивные правила:

$A \rightarrow w_1K \mid w_2K \mid \dots \mid w_kK$;

$K \rightarrow x_1K \mid x_2K \mid \dots \mid x_pK \mid \varepsilon$.

Пример. Пусть задана грамматика G со следующими правилами вывода: $S \rightarrow Sa \mid ba$. Требуется построить эквивалентную ей праворекурсивную грамматику.

Для решения задачи воспользуемся рассмотренным выше алгоритмом. В исходной грамматике один леворекурсивный нетерминал S . Построим для него вывод:

$$S \rightarrow SB \mid D;$$

$$B \rightarrow a; \quad D \rightarrow ba. \quad \text{Вывод из } S \text{ имеет вид: } S \rightarrow SB \rightarrow SBB \rightarrow \dots \rightarrow DB^*.$$

Запишем эквивалентные правила для S :

$$S \rightarrow DK; \quad K \rightarrow BK \mid \varepsilon.$$

Подставим в эти правила B и D и получим эквивалентную праворекурсивную грамматику G_1 :

$$S \rightarrow baK; \quad K \rightarrow aK \mid \varepsilon.$$

Рассмотрим вывод цепочки $baaaa$ в исходной леворекурсивной грамматике G :

$$S \rightarrow Sa \rightarrow Saa \rightarrow Saaa \rightarrow baaaa.$$

Эту же цепочку можно вывести в эквивалентной праворекурсивной грамматике G_1 :

$$S \rightarrow baK \rightarrow baaK \rightarrow baaaK \rightarrow baaaaK \rightarrow baaaa.$$