

## Лекция №3 Грамматики.

### Формальное определение грамматики. Типы грамматик и их свойства.

Для нас наибольший интерес представляет одна из систем генерации языков - грамматики. Понятие грамматики сначала было формализовано лингвистами при изучении естественных языков. Предполагалось, что это может помочь при их автоматической трансляции. Однако, лучшие результаты в этом направлении достигнуто при описании не естественных языков, а языков программирования. примером может служить способ описания синтаксиса языков программирования с помощью БНФ - формы Бекуса-Наура.

Определение. Грамматика - это четверка  $G = (N, T, P, S)$ , где

$N$  - алфавит нетерминальных символов;

$T$  - алфавит терминальных символов,  $NT = \emptyset$ ;

$P$  - конечное множество правил вида, где  $(NT)^* N (NT)^* (NT)^*$ ;

$S$  - начальный символ (или аксиома) грамматики.

Мы будем использовать большие латинские буквы для обозначения нетерминальных символов, малые латинские буквы с начала алфавита для обозначения терминальных символов, малые латинские буквы с конца алфавита для обозначения цепочек с  $T^*$  и, наконец, малые греческие буквы для обозначения цепочек с  $(NT)^*$ .

Будем использовать также сокращенную запись  $A_1 | A_2 | \dots | A_n$  для обозначения группы правил  $A_1, A_2, \dots, A_n$ .

Определим на множестве  $(NT)^*$  бинарное отношение выводимости таким образом: если  $P$ , то для всех,  $(NT)^*$ . Если  $1, 2$ , то говорят, что цепочка  $2$  непосредственно выведена из  $1$ .

Мы будем использовать также рефлексивно-транзитивное и транзитивное замыкания отношения, а также его степень  $k \geq 0$  (обозначаются соответственно  $^*$  и  $^k$ ). Если  $1^* 2$  ( $1 \rightarrow^* 2$ ,  $1 \xrightarrow{k} 2$ ), то говорят, что цепочка  $2$  выведена (нетривиально выведена, выведена за  $k$  шагов) с  $1$ .

Если  $k \geq 0$ , то существует последовательность шагов

где  $k = 0$  и  $k$ . Последовательность цепочек  $0, 1, 2, \dots, k$  в этом случае называют выводом из.

Сентенциальной форме грамматики  $G$  называется цепочка, выводимая из ее начального символа.

Языку, порождаемой грамматикой  $G$  (обозначается  $L(G)$ ), называется множество всех ее терминальных сентенциальных форм, то есть

Грамматики  $G_1$  и  $G_2$  называются эквивалентными, если они порождают одну и ту же речь, то есть  $L(G_1) = L(G_2)$ .

Пример 2.5. Грамматика  $G = (\{S, B, C\}, \{a, b, c\}, P, S)$ , где

$P = \{S \rightarrow aSBC, S \rightarrow aBC, CB \rightarrow BC, aB \rightarrow ab, bB \rightarrow bb, bC \rightarrow bc, cC \rightarrow cc\}$  порождает язык  $L(G) = \{a^n b^n c^n \mid n > 0\}$ .

Действительно, применяем  $n-1$  раз правило 1 и получаем  $a^{n-1}S(BC)^{n-1}$ , затем один раз правило 2 и получаем  $a^n(BC)^n$ , затем  $n(n-1)/2$  раз правило 3 и получаем  $a^n B^n C^n$ .

Затем используем правило 4 и получаем  $a^n b B^{n-1} C^n$ . затем применяем  $n-1$  раз правило 5 и получаем  $a^n b^n C^n$ . Затем применяем правило 6 и  $n-1$  раз правило 7 и получаем  $a^n b^n c^n$ . Можно показать, что язык  $L(G)$  состоит из цепочек только такого вида.

Пример 2.6. Рассмотрим грамматику  $G = (\{S\}, \{0, 1\}, \{S \rightarrow 0S1, S \rightarrow 01\}, S)$ . легко видеть, что цепочка  $000111 \in L(G)$ , так как существует заключение

Нетрудно показать, что грамматика порождает язык  $L(G) = \{0^n 1^n \mid n > 0\}$ .

Пример 2.7. Рассмотрим грамматику  $G = (\{S, A\}, \{0, 1\}, \{S \rightarrow 0S, S \rightarrow 0A, A \rightarrow 1A, A \rightarrow 1\}, S)$ . Нетрудно показать, что грамматика порождает язык  $L(G) = \{0^n 1^m \mid n, m > 0\}$ .

### Типы грамматик и их свойства

Рассмотрим классификацию грамматик (предложенную Н. Хомского), основанную на виду их правил.

Определение. Пусть дана грамматика  $G = (N, T, P, S)$ . тогда

если правила грамматики не удовлетворяют никаким ограничениям, то ее называют грамматикой типа 0, или грамматикой без ограничений.

если

каждое правило грамматики, кроме  $S \rightarrow \epsilon$ , имеет вид, где  $| \alpha | \geq 1$ , и

в том случае, когда  $S \in P$ , символ  $S$  не встречается в правых частях правил,

то грамматику называют грамматикой типа 1 или несокращения.

если каждое правило грамматики имеет вид  $A \rightarrow \alpha$ , где  $A \in N$ ,  $(\alpha \in T^+)$ , то ее называют грамматикой типа 2, или контекстно-свободной (КС-грамматикой).

если каждое правило грамматики имеет вид  $A \rightarrow xB$  или  $A \rightarrow x$ , где  $A, B \in N$ ,  $x \in T^+$  то ее называют грамматикой типа 3, или праволинейной.

Легко видеть, что грамматика в примере 2.5 - несокращения, в примере 2.6 - контекстно-свободная, в примере 2.7 - праволинейная.

Язык, порожденная грамматикой типа  $i$ , называют языком типа  $i$ . Язык типа 0 называют также языком без ограничений, язык типа 1 - контекстно-зависимым (КЗ), язык типа 2 - контекстно-свободным (КС), язык типа 3 - праволинейным.

**Теорема 2.1.** Каждая контекстно-свободная речь может быть порождена неукорачивающей контекстно-свободной грамматикой.

**Доказательство.** Пусть  $L$  - контекстно-свободная речь. Тогда существует контекстно-свободная грамматика  $G = (N, T, P, S)$ , порождает  $L$ .

Построим новую грамматику  $G' = (N', T, P', S')$  следующим образом:

1. Если в  $P$  есть правило вида  $A \rightarrow B_1 B_2 \dots B_k$ , где  $k \geq 1$ ,  $B_i \in T \cup N$ , и ни с одной цепочки  $j$  ( $0 \leq j < k$ ) не выводится  $\epsilon$ , то включить в  $P'$  все правила (кроме  $A \rightarrow \epsilon$ ) вида

где  $X_i$  - это либо  $B_i$ , или  $\epsilon$ .

2. Если  $S \rightarrow \epsilon$ , то включить в  $P'$  правила  $S' \rightarrow S$ ,  $S' \rightarrow \epsilon$  и положить  $N' = N \cup \{S'\}$ . В противном случае положить  $N' = N$  и  $S' = S$ .

Порождает грамматика пустую цепочку можно установить следующим простым алгоритмом:

Шаг 1. Строим множество  $N_0 = N \mid N \rightarrow e$

Шаг 2. Строим множество  $N_i = N \mid N \rightarrow \alpha \cdot N_{i-1}^*$

Шаг 3. Если  $N_i = N_{i-1}$ , перейти к шагу 4, иначе шаг 2.

Шаг 4. Если  $S \rightarrow N_i$  значит  $S \rightarrow e$  /

Легко видеть, что  $G'$  - неукорачивающая грамматика. Можно показать по индукции,  $L(G') = L(G)$ . \_\_

Пусть  $K_i$  - класс всех языков типа  $i$ . Доказано, что справедливо следующее (Строгое) включения:

$K_3 \supset K_2 \supset K_1 \supset K_0$ .

Заметим, что если речь порождается некоторой грамматикой, это не значит, что он не может быть порожден грамматикой с более сильными ограничениями на правила. Приведенный ниже пример иллюстрирует этот факт.

Пример 2.8. Рассмотрим грамматику  $G = (\{S, A, B\}, \{0, 1\}, \{S \rightarrow AB, A \rightarrow 0A, A \rightarrow 0, B \rightarrow 1B, B \rightarrow 1\}, S)$ . Эта грамматика является контекстно-свободной. Легко показать, что  $L(G) = \{0^n 1^m \mid n, m > 0\}$ .

Однако, в примере 2.7 приведен праволинейная порождающая грамматика ту же язык.

Покажем что существует алгоритм, позволяющий для произвольного КЗ-языка  $L$  в алфавите  $T$ , и произвольной цепочки  $w \in T^*$  определить, принадлежит ли  $w$  языку  $L$ .

Теорема 2.2. Каждый контекстно-зависимый язык является рекурсивным языком.

Доказательство. Пусть  $L$  - контекстно-зависимый язык. тогда существует некоторая неукорачивающая грамматика  $G = (N, T, P, S)$ , порождает  $L$ .

Пусть  $w \in T^*$  и  $|w| = n$ . Если  $n = 0$ , то есть  $w = e$ , то принадлежность  $w \in L$  проверяется тривиальным образом. Так что будем предполагать, что  $n > 0$ .

Определим множество  $T_m$  как множество строк  $u$  (NT) + длины не более  $n$  таких, что вывод  $S \Rightarrow^* u$  имеет не более  $m$  шагов. Ясно, что  $T_0 = \{S\}$ .

Легко показать, что  $T_m$  можно получить из  $T_{m-1}$  просматривая, какие строки с длиной, меньшей или равной  $n$  можно вывести из строк с  $T_{m-1}$  применением одного правила, то есть

$$T_m = T_{m-1} \cup \{u \mid v \Rightarrow u \text{ для некоторого } v \in T_{m-1}, \text{ где } |u| \leq n\}.$$

если  $S \Rightarrow^* u$  и  $|u| \leq n$ , то  $u \in T_m$  для некоторого  $m$ . Если с  $S$  не выводится  $u$  или  $|u| > n$ , то  $u$  не принадлежит  $T_m$  ни для какого  $m$ .

Очевидно, что  $T_m \supseteq T_{m-1}$  для всех  $m \geq 1$ . Поскольку  $T_m$  зависит только от  $m$ , если  $T_m = T_{m-1}$ , то  $T_m = T_{m+1} = T_{m+2} = \dots$  процедура будет вычислять  $T_1, T_2, T_3, \dots$  пока для некоторого  $m$  не обнаружится  $T_m = T_{m-1}$ . если  $w$  не принадлежит  $T_m$ , то не принадлежит и  $L(G)$ , поскольку для  $j > m$  выполнено  $T_j = T_m$ . если  $w \in T_m$ , то  $S \Rightarrow^* w$ .

Покажем, что существует такое  $m$ , что  $T_m = T_{m-1}$ . Поскольку для каждого  $i \geq 1$  справедливо  $T_i \supseteq T_{i-1}$ , то если  $T_i \neq T_{i-1}$ , то число элементов в  $T_i$  крайней мере на 1 больше, чем в  $T_{i-1}$ . пусть  $|NT| = K$ . Тогда число строк в (NT) + длины меньшей или равной  $n$  равно  $k + k_2 + \dots + k_n \leq nk_n$ . Только эти строки могут быть в любом  $T_i$ . Значит,  $T_m = T_{m-1}$  для некоторого  $m \leq nk_n$ . Таким образом, процедура, вычисляемый  $T_i$  для всех  $i \geq 1$  до тех пор, пока не будут найдены два равных множеств, гарантированно заканчивается, значит, это алгоритм.