

Семейство лемм о накачке.
Нормальная форма Грейбах.
Теорема Хомского–Шутценберже

Теория формальных языков
2022 г.



Связь КС-грамматик и АПГ

Рассмотрим язык сентенциальных форм, но только без учёта терминальных символов. Получим алфавитную префиксную грамматику (АПГ, см. лекцию 4). При этом каждое применение правила переписывания такой грамматики выбрасывает не более чем один терминальный символ слева, а стирающее правило — ровно один.

Рассмотрим КС-грамматику, соответствующую ей АПГ и путь вывода слова, получаемый с помощью АПГ.

$$S \rightarrow aS \mid bS \mid aN_a \mid bN_b$$

$$N_a \rightarrow N_a E \mid bN_a E \mid cVB$$

$$N_b \rightarrow aN_b E \mid bN_b E \mid VA$$

$$V \rightarrow aV \mid bV \mid a \mid b$$

$$A \rightarrow a \quad B \rightarrow b \quad E \rightarrow a \mid b$$

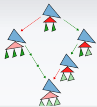
$$S \rightarrow S \mid N_a \mid N_b$$

$$N_a \rightarrow N_a E \mid VB$$

$$N_b \rightarrow N_b E \mid VA$$

$$V \rightarrow V \mid \varepsilon$$

$$A \rightarrow \varepsilon \quad B \rightarrow \varepsilon \quad E \rightarrow \varepsilon$$



Связь КС-грамматик и АПГ

$$S \rightarrow aS \mid bS \mid aN_a \mid bN_b$$

$$N_a \rightarrow N_a E \mid bN_a E \mid cVB$$

$$N_b \rightarrow aN_b E \mid bN_b E \mid VA$$

$$V \rightarrow aV \mid bV \mid a \mid b$$

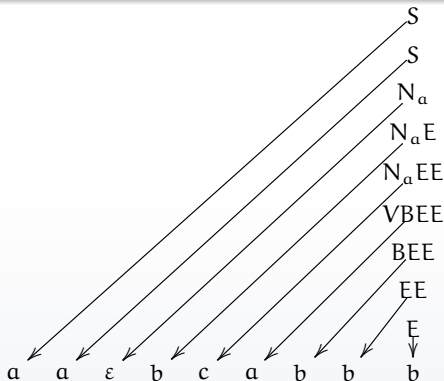
$$A \rightarrow a \quad B \rightarrow b \quad E \rightarrow a \mid b$$

$$S \rightarrow S \mid N_a \mid N_b$$

$$N_a \rightarrow N_a E \mid VB$$

$$N_b \rightarrow N_b E \mid VA$$

$$V \rightarrow V \mid \varepsilon$$

$$A \rightarrow \varepsilon \quad B \rightarrow \varepsilon \quad E \rightarrow \varepsilon$$




Основная теорема

Следствие теоремы Турчина

Пусть G — алфавитная префиксная грамматика, в которой N правил с непустой правой частью, и максимальная длина правой части правила равна M . Тогда любая последовательность порождаемых ею слов

$a_1 \dots a_n$

\dots

ε

длиной не менее $N^M \cdot (n + 1)$ содержит пару вида

$\tau_1 = \Phi\Theta_0$, $\tau_2 = \Phi\Psi\Theta_0$, такую, что $|\Phi| \leq M$ и на отрезке $[\tau_1, \tau_2]$ нет слов длины меньше $|\Theta_0| + 1$.

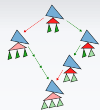


Первая лемма о накачке

Классическая лемма о накачке

Пусть G — КС-язык. Тогда существует $p \in \mathbb{N}$ такое, что любое слово $w \in L(G)$ длины не меньше p имеет представление вида $x_1 y_1 z y_2 x_2$, где $|y_1 y_2| \geq 1$, $|y_1 z y_2| \leq p$, и все слова вида $x_1 y_1^k z y_2^k x_2$ также принадлежат $L(G)$.

Доказательство: при левостороннем разборе выбираем самую последнюю пару сентенциальных форм.



Вторая лемма о накачке

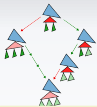
Пусть G — КС-язык. Тогда существует $p \in \mathbb{N}$ такое, что любое слово $w \in L(G)$ длины не меньше p имеет представление вида $x_1 y_1 z y_2 x_2$, где $|y_2| \geq 1$, $|x_1 y_1| \leq p$, $|y_2| \leq p$ либо y_2 накачивается отдельно, $|x_2| \leq p$ либо x_2 накачивается отдельно, и все слова вида $x_1 y_1^k z y_2^k x_2$ также принадлежат $L(G)$.



Варианты лемм о накачке

- Хотелось бы сдвигать начало отрезка накачки вперёд на любое константное количество букв, аналогично — конец отрезка накачки (используя свойство реверса).
- Понятие накачки может быть применено рекурсивно к некоторым достаточно длинным подсловам выбранного слова. Т.е. можно выкидывать из слова подслова, накачиваемые отдельно, без риска выйти из языка.

Для первого нужна гарантия того, что если порождается k букв слова, то длина сентенциальных форм увеличивается не более чем в $f(k)$ раз (где f — хотя бы полином).



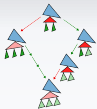
Бонус: лемма Огдена

Пусть L — КС-язык. Тогда существует такое число n , что в любом слове w , $|w| \geq n$, при отметке n или более букв, w представляется в виде $x_1y_1zy_2x_2$, причём либо во всех трех из x_1 , y_1 , z есть отмеченные буквы, либо они есть во всех трех из z , y_2 , x_2 , в слове y_1zy_2 отмечено не более n букв, и $\forall k(x_1y_1^kzy_2^kx_2 \in L)$.

Исследуем «плохой» язык $\{a^mb^nc^nd^n \mid m > 0\} \cup \{b^ic^jd^k\}$ с помощью леммы Огдена. Абельяр (антагонист) выбирает n . Элоиза (т.е. мы) строит слово $ab^{2n}c^{2n}d^{2n}$ и отмечает n последних букв d . Абельяр может разбить слово $ab^{2n}c^{2n}d^{2n}$ на $x_1y_1zy_2x_2$ двумя способами:

- отмечены x_1 , y_1 , z , накачиваться может только d^{2n} .
- отмечены x_2 , y_2 , z , накачивается либо d^{2n} , либо d^{2n} совместно с c^{2n} , b^{2n} или a .

Оба типа накачки выводят из языка, поскольку при любой положительной накачке число вхождений букв b или c расходится с числом вхождений d в слово.



Н.ф. Хомского и левосторонний вывод

- Могут быть непродуктивные левосторонние цепочки:
 $A \rightarrow AB \rightarrow \dots AB^n \rightarrow \dots$
- Есть гарантия роста слова при развёртке, но нет определённости, по какому префиксу.



Нормальная форма Грейбах

Определение

Грамматика G ($\varepsilon \notin L(G)$) находится в GNF (н.ф. Грейбах) \Leftrightarrow каждое её правило имеет вид $A_i \rightarrow a_j \alpha$, где $A_i \in N$, $\alpha \in N^*$, $a_j \in \Sigma$.

- Левосторонний разбор по грамматике в GNF на каждом шагу переписывания порождает терминальный символ.



Нормальная форма Грейбах

Определение

Грамматика G ($\varepsilon \notin L(G)$) находится в GNF (н.ф. Грейбах) \Leftrightarrow каждое её правило имеет вид $A_i \rightarrow a_j \alpha$, где $A_i \in N$, $\alpha \in N^*$, $a_j \in \Sigma$.

- Левосторонний разбор по грамматике в GNF на каждом шагу переписывания порождает терминальный символ.
- Для приведения к GNF нужно «вытащить из рекурсии» возможные first-терминалы, порождаемые нетерминалами грамматики.



Нормальная форма Грейбах

Определение

Грамматика G ($\varepsilon \notin L(G)$) находится в GNF (н.ф. Грейбах) \Leftrightarrow каждое её правило имеет вид $A_i \rightarrow a_j \alpha$, где $A_i \in N$, $\alpha \in N^*$, $a_j \in \Sigma$.

- Левосторонний разбор по грамматике в GNF на каждом шагу переписывания порождает терминальный символ.
- Для приведения к GNF нужно «вытащить из рекурсии» возможные first-терминалы, порождаемые нетерминалами грамматики.
 - явно найти все завершающиеся цепочки вывода;

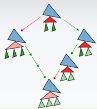


Нормальная форма Грейбах

Определение

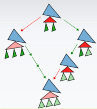
Грамматика G ($\varepsilon \notin L(G)$) находится в GNF (н.ф. Грейбах) \Leftrightarrow каждое её правило имеет вид $A_i \rightarrow a_j \alpha$, где $A_i \in N$, $\alpha \in N^*$, $a_j \in \Sigma$.

- Левосторонний разбор по грамматике в GNF на каждом шагу переписывания порождает терминальный символ.
- Для приведения к GNF нужно «вытащить из рекурсии» возможные first-терминалы, порождаемые нетерминалами грамматики.
 - явно найти все завершающиеся цепочки вывода;
 - рассмотреть язык-реверс сентенциальных форм.
- По умолчанию считаем, что к GNF приводится CNF (н.ф. Хомского).



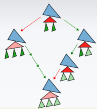
Первый способ приведения к GNF

- 1 Нумеруем нетерминалы в правых частях правил в порядке их вхождения;
- 2 (по исчерпанию по i , начиная с $i = 1$) Если имеется правило вида $A_i \rightarrow B_j \beta$, где $j < i$, тогда подставляем вместо B_j все правые части α_k правил вида $B_j \rightarrow \alpha_k$.
- 3 Если после этого все правила имеют вид либо $A_i \rightarrow a\alpha$, $a \in \Sigma$, либо $A_i \rightarrow B_j \beta$, причём $i < j$, тогда GNF получается последовательной развёрткой B_j .



Первый способ приведения к GNF

- 1 Нумеруем нетерминалы в правых частях правил в порядке их вхождения;
- 2 (по исчерпанию по i , начиная с $i = 1$) Если имеется правило вида $A_i \rightarrow B_j \beta$, где $j < i$, тогда подставляем вместо B_j все правые части α_k правил вида $B_j \rightarrow \alpha_k$.
- 3 Если после этого все правила имеют вид либо $A_i \rightarrow a\alpha$, $a \in \Sigma$, либо $A_i \rightarrow B_j \beta$, причём $i < j$, тогда GNF получается последовательной развёрткой B_j . Существует лексикографический порядок на функциональных символах из N .



Первый способ приведения к GNF

- 1 Нумеруем нетерминалы в правых частях правил в порядке их вхождения;
- 2 (по исчерпанию по i , начиная с $i = 1$) Если имеется правило вида $A_i \rightarrow V_j \beta$, где $j < i$, тогда подставляем вместо V_j все правые части α_k правил вида $V_j \rightarrow \alpha_k$.
- 3 Если после этого все правила имеют вид либо $A_i \rightarrow a\alpha$, $a \in \Sigma$, либо $A_i \rightarrow V_j \beta$, причём $i < j$, тогда GNF получается последовательной развёрткой V_j . Существует лексикографический порядок на функциональных символах из N .
- 4 Если есть правила вида $A_i \rightarrow A_i \alpha$, тогда устраняем левую рекурсию.
- 5 Увеличиваем i , если ещё остались нетерминалы, не приведённые к ГНФ.



Устранение левой рекурсии

- ❶ Предположим, для A_i нашлось n леворекурсивных правил и m упорядоченных лексикографически:

$$A_i \rightarrow A_i \alpha_1$$

$$A_i \rightarrow \beta_1$$

...

...

$$A_i \rightarrow A_i \alpha_n$$

$$A_i \rightarrow \beta_m$$

- ❷ Вводим новый нетерминал A'_i такой, что его вес меньше всех прочих, и заменяем правила на:

$$A'_i \rightarrow \alpha_1 A'_i \mid \alpha_1$$

$$A_i \rightarrow \beta_1 \mid \beta_1 A'_i$$

...

...

$$A'_i \rightarrow \alpha_n A'_i \mid \alpha_n$$

$$A_i \rightarrow \beta_m \mid \beta_m A'_i$$

- ❸ После всех таких замен грамматика лексикографически упорядочена по левому разбору, и GNF получается последовательной левой развёрткой.



Второй способ приведения к GNF

Алгоритм Блюма–Коха (1999).

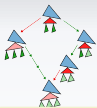
Неформальное описание

- Рассмотрим язык sentенциальных форм с переписыванием только по левому разбору. Он регулярен, и в конечное состояние его НКА ведут стрелки, помеченные терминалами.
- Для такого языка легко построить инверсный \Rightarrow множество терминалов-префиксов, которые может породить данный нетерминал.



Второй способ: порождение НКА

- ❶ По каждому нетерминалу V строим автомат $M_V = \langle N_V \cup \{S_V\}, \Sigma \cup N, V_V, \{S_V\}, \delta \rangle$ (S_V — новое состояние, N_V — множество нетерминалов CFG, индексированное нетерминалом V). Правила перехода δ :
 - $\langle C_V, E, M \rangle \Leftrightarrow M = \{D_V \mid C \rightarrow DE \in P\};$
 - $\langle C_V, a, \{S_V\} \rangle \Leftrightarrow C \rightarrow a \in P.$
- ❷ Строим реверс к M_V , получаем НКА M_V^R .
- ❸ Строим грамматику $G'_V = \langle N_V \cup \{S_V\}, \Sigma \cup N, R'_V, S_V \rangle$ для M_V^R с правилами переписывания:
 - $S_V \rightarrow aC_V \Leftrightarrow \langle S_V, a, C_V \rangle \in \delta^R$ и $C_V \neq V_V$ либо из V_V есть стрелки в M_V^R ;
 - $S_V \rightarrow a \Leftrightarrow \langle S_V, a, V_V \rangle \in \delta^R$;
 - $D_V \rightarrow EC_V \Leftrightarrow \langle D_V, E, C_V \rangle \in \delta^R$ и $C_V \neq V_V$ либо из V_V есть стрелки в M_V^R ;
 - $C_V \rightarrow E \Leftrightarrow \langle C_V, E, V_V \rangle \in \delta^R.$



Окончание конструкции

Теперь по всем G'_i строим окончательный вариант грамматики $G_B = \langle N_B \cup \{S_B\}, \Sigma, R_B, S_B \rangle$ с правилами:

- $S_B \rightarrow aC_B, S_B \rightarrow aC_B \in R'_B$;
- $S_B \rightarrow a S_B \rightarrow a \in R'_B$;
- $D_B \rightarrow \alpha C_B \Leftrightarrow D_B \rightarrow EC_B \in R'_B \ \& \ S_E \rightarrow \alpha$ (по всем таким α и E);
- $D_B \rightarrow \alpha \Leftrightarrow D_B \rightarrow E \in R'_B \ \& \ S_E \rightarrow \alpha$ (по всем таким α и E).

Грамматика $\bigcup_{i \in \mathbb{N}} G_i$ со стартовым символом S_S — это искомая GNF для исходной грамматики G .

Последние два правила разворачивают неразмеченные нетерминалы в стартовые правила грамматик их сентенциальных форм, поэтому автоматы M_B имеет смысл строить только для тех нетерминалов, которые встречаются в исходной грамматике не только первыми в правых частях правил.



Пример преобразования грамматики по Блему–Коху

Привести к GNF грамматику некорректных сумм двоичных чисел (почему некорректных?)

$$S \rightarrow S + S \mid D \quad D \rightarrow D0 \mid D1 \mid 1 \mid (S)$$



Пример преобразования грамматики по Блему–Коху

Привести к GNF грамматику некорректных сумм двоичных чисел (почему некорректных?)

$$S \rightarrow S + S \mid D \quad D \rightarrow D0 \mid D1 \mid 1 \mid (S)$$

Сначала избавляемся от цепного правила $S \rightarrow D$. Потом строим порождающую структуру A_V сентенциальных форм по левостороннему разбору с финальным состоянием N_V и стартовым V_V . Каждому нетерминалу V соответствует своя структура.

$$\begin{array}{llllll} \text{Для } A_S : & S_S \xrightarrow{+S} S_S & S_S \xrightarrow{0} D_S & S_S \xrightarrow{1} D_S & S_S \xrightarrow{1} N_S \\ & S_S \xrightarrow{(S)} N_S & D_S \xrightarrow{0} D_S & D_S \xrightarrow{1} D_S & D_S \xrightarrow{1} N_S & D_S \xrightarrow{(S)} N_S \end{array}$$



Пример преобразования грамматики по Блему–Коху

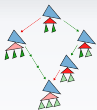
Привести к GNF грамматику некорректных сумм двоичных чисел (почему некорректных?)

$$S \rightarrow S + S \mid D \quad D \rightarrow D0 \mid D1 \mid 1 \mid (S)$$

Сначала избавляемся от цепного правила $S \rightarrow D$. Потом строим порождающую структуру A_V сентенциальных форм по левостороннему разбору с финальным состоянием N_V и стартовым V_V . Каждому нетерминалу V соответствует своя структура.

$$\begin{array}{llllll} \text{Для } A_S : & S_S \xrightarrow{+S} S_S & S_S \xrightarrow{0} D_S & S_S \xrightarrow{1} D_S & S_S \xrightarrow{1} N_S \\ & S_S \xrightarrow{(S)} N_S & D_S \xrightarrow{0} D_S & D_S \xrightarrow{1} D_S & D_S \xrightarrow{1} N_S & D_S \xrightarrow{(S)} N_S \end{array}$$

$$\text{Для } A_D : \quad D_D \xrightarrow{0} D_D \quad D_D \xrightarrow{1} D_D \quad D_D \xrightarrow{1} N_D \quad D_D \xrightarrow{(S)} N_D$$



Пример преобразования грамматики по Блуму–Коху

Превращаем структуры в праволинейные (меняя местами нетерминалы левых и правых частей правил и стартовые состояния с финальными):

$$\begin{aligned}
 \text{Для } A_S : \quad & S_S \xrightarrow{+S} S_S \quad D_S \xrightarrow{0} S_S \quad D_S \xrightarrow{1} D_S \quad N_S \xrightarrow{1} S_S \\
 & N_S \xrightarrow{(S)} S_S \quad D_S \xrightarrow{0} D_S \quad D_S \xrightarrow{1} D_S \quad N_S \xrightarrow{1} D_S \quad N_S \xrightarrow{(S)} D_S \\
 \\
 \text{Для } A_D : \quad & D_D \xrightarrow{0} D_D \quad D_D \xrightarrow{1} D_D \quad N_D \xrightarrow{1} D_D \quad N_D \xrightarrow{(S)} D_D
 \end{aligned}$$

Извлекаем праволинейные (почти) грамматики. В правой части правила может не быть нетерминалов, если там стоял нетерминал V_V .

$$\begin{array}{lllll}
 \text{q-RLG } G_S : & S_S \rightarrow +SS_S & D_S \rightarrow 0S_S & D_S \rightarrow 1D_S & N_S \rightarrow 1S_S \\
 & N_S \rightarrow (S)S_S & D_S \rightarrow 0D_S & D_S \rightarrow 1D_S & N_S \rightarrow (S)D_S \\
 & S_S \rightarrow +S & D_S \rightarrow 0 & D_S \rightarrow 1 & N_S \rightarrow (S)
 \end{array}$$

Извлекаем праволинейные (почти) грамматики. В правой части правила может не быть нетерминалов, если там стоял нетерминал V_V .

$$\begin{array}{lllll} \text{q-RLG } G_S : & S_S \rightarrow +SS_S & D_S \rightarrow 0S_S & D_S \rightarrow 1D_S & N_S \rightarrow 1S_S \\ N_S \rightarrow (S)S_S & D_S \rightarrow 0D_S & D_S \rightarrow 1D_S & N_S \rightarrow 1D_S & N_S \rightarrow (S)D_S \\ S_S \rightarrow +S & D_S \rightarrow 0 & D_S \rightarrow 1 & N_S \rightarrow 1 & N_S \rightarrow (S) \end{array}$$

$$\begin{array}{lllll} \text{q-RLG } G_D : & D_D \rightarrow 0D_D & D_D \rightarrow 1D_D & N_D \rightarrow 1D_D & N_D \rightarrow (S)D_D \\ D_D \rightarrow 0 & D_D \rightarrow 1 & N_D \rightarrow 1 & N_D \rightarrow (S) & \end{array}$$

Извлекаем праволинейные (почти) грамматики. В правой части правила может не быть нетерминалов, если там стоял нетерминал V_V .

$$\begin{array}{lllll} \text{q-RLG } G_S : & S_S \rightarrow +SS_S & D_S \rightarrow 0S_S & D_S \rightarrow 1D_S & N_S \rightarrow 1S_S \\ & N_S \rightarrow (S)S_S & D_S \rightarrow 0D_S & D_S \rightarrow 1D_S & N_S \rightarrow 1D_S & N_S \rightarrow (S)D_S \\ & S_S \rightarrow +S & D_S \rightarrow 0 & D_S \rightarrow 1 & N_S \rightarrow 1 & N_S \rightarrow (S) \end{array}$$

$$\begin{array}{lllll} \text{q-RLG } G_D : & D_D \rightarrow 0D_D & D_D \rightarrow 1D_D & N_D \rightarrow 1D_D & N_D \rightarrow (S)D_D \\ & D_D \rightarrow 0 & D_D \rightarrow 1 & N_D \rightarrow 1 & N_D \rightarrow (S) \end{array}$$

Заменяем неразмеченные нетерминальные символы V исходной грамматики на N_V . В данном случае нет правил, в которых неразмеченные нетерминалы стояли бы первыми в правых частях, поэтому достаточно просто заменить их на N_V . Иначе пришлось бы заменять их на все возможные правые части α правил вида $N_V \rightarrow \alpha$. Стартовый символ — N_S . GNF почти построена!

Извлекаем праволинейные (почти) грамматики. В правой части правила может не быть нетерминалов, если там стоял нетерминал V_V .

Заменяем неразмеченные нетерминальные символы V исходной грамматики на N_V . В данном случае нет правил, в которых неразмеченные нетерминалы стояли бы первыми в правых частях, поэтому достаточно просто заменить их на N_V . Иначе пришлось бы заменять их на все возможные правые части α правил вида $N_V \rightarrow \alpha$. Стартовый символ — N_S . GNF почти построена!

q-GNF для G :

$S_S \rightarrow +N_S S_S$	$D_S \rightarrow 0S_S$	$D_S \rightarrow 1D_S$	$N_S \rightarrow 1S_S$	
$N_S \rightarrow (N_S)S_S$	$D_S \rightarrow 0D_S$	$D_S \rightarrow 1D_S$	$N_S \rightarrow 1D_S$	$N_S \rightarrow (N_S)D_S$
$S_S \rightarrow +N_S$	$D_S \rightarrow 0$	$D_S \rightarrow 1$	$N_S \rightarrow 1$	$N_S \rightarrow (N_S)$

Извлекаем праволинейные (почти) грамматики. В правой части правила может не быть нетерминалов, если там стоял нетерминал V_V .

Заменяем неразмеченные нетерминальные символы V исходной грамматики на N_V . В данном случае нет правил, в которых неразмеченные нетерминалы стояли бы первыми в правых частях, поэтому достаточно просто заменить их на N_V . Иначе пришлось бы заменять их на все возможные правые части α правил вида $N_V \rightarrow \alpha$. Стартовый символ — N_S . GNF почти построена!

q-GNF для G :

$S_S \rightarrow +N_S S_S$	$D_S \rightarrow 0S_S$	$D_S \rightarrow 1D_S$	$N_S \rightarrow 1S_S$	
$N_S \rightarrow (N_S)S_S$	$D_S \rightarrow 0D_S$	$D_S \rightarrow 1D_S$	$N_S \rightarrow 1D_S$	$N_S \rightarrow (N_S)D_S$
$S_S \rightarrow +N_S$	$D_S \rightarrow 0$	$D_S \rightarrow 1$	$N_S \rightarrow 1$	$N_S \rightarrow (N_S)$

Осталось обернуть в delay-нетерминалы терминальные символы правых частей правил, кроме первого. Здесь это символ $)$.



Теорема Хомского–Шутценберже

Пусть PAREN_n — язык из $4 * n$ элементов
 $\{[1,]_1, \dots, [n,]_n, (1,)_1, \dots, (n,)_n\}$.

Теорема

Любой CF-язык получается гомоморфизмом из языка
 $L' = \text{PAREN}_n \cap R$, где R — регулярный.

Пусть G — грамматика L в нормальной форме Хомского.
Пронумеруем правила G и поставим им в соответствие следующие.



Теорема Хомского–Шутценберже

Пусть PAREN_n — язык из $4 * n$ элементов
 $\{[1,]_1, \dots, [n,]_n, (1,)_1, \dots, (n,)_n\}$.

Теорема

Любой CF-язык получается гомоморфизмом из языка
 $L' = \text{PAREN}_n \cap R$, где R — регулярный.

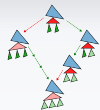
Пусть G — грамматика L в нормальной форме Хомского.
Пронумеруем правила G и поставим им в соответствие следующие.

- ❶ Если правило n имеет вид $A \rightarrow BC$, тогда порождаем правило $A \rightarrow [{}_n B]_n ({}_n C)_n$.
- ❷ Если правило n имеет вид $A \rightarrow a$, тогда порождаем правило $A \rightarrow [{}_n]_n ({}_n)_n$.



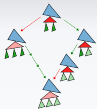
Свойства языка $L(G')$

- Все $]_n$ строго предшествуют $(_n$.



Свойства языка $L(G')$

- Все $]_n$ строго предшествуют $(_n$.
- Ни одна $)_n$ не предшествует непосредственно левой скобке.



Свойства языка $L(G')$

- Все $]_n$ строго предшествуют $(_n$.
- Ни одна $)_n$ не предшествует непосредственно левой скобке.
- Если правило n — это $A \rightarrow BC$, тогда $[_n$ непосредственно предшествует некоторой $[_p$, так же как и $(_n$.



Язык R

$R = \{x \in \{[j,]_j, (j,)_j\}^* \mid x \text{ начинается с } [{}_n \text{ для некоторого правила } n : A \rightarrow \dots \& \text{ все }]_n \text{ предшествуют } ({}_n\}.$



Язык R

$R = \{x \in \{[j,]_j, (j,)_j\}^* \mid x \text{ начинается с } [{}_n \text{ для некоторого правила } n : A \rightarrow \dots \& \text{ все }]_n \text{ предшествуют } ({}_n\}.$

Можно убедиться, что $L(G') = R \cap \text{PAREN}_n.$



Язык R

$R = \{x \in \{[j,]_j, (j,)_j\}^* \mid x \text{ начинается с } [{}_n \text{ для некоторого правила } n : A \rightarrow \dots \& \text{ все }]_n \text{ предшествуют } ({}_n\}.$

Можно убедиться, что $L(G') = R \cap \text{PAREN}_n.$

Осталось определить h . Если n — нефинальное правило, то $h([{}_n) = h(]_n) = h(({}_n) = h()_n) = \varepsilon$. Иначе $h([{}_n) = a$, для остальных скобок так же.



Значение теоремы Х.-Ш.

Возможно разделить парсинг любого КС-языка на две стадии: лексический анализ (проверка условия R) и разбор правильных скобочных структур.

Замечание: поскольку гомоморфизм h не обязан быть инъективным, разбор ПСП не всегда можно определить однозначно. Пример: $\{a^n b^n\} \cup \{a^n b^{2n}\}$ (полностью неоднозначность устранить нельзя, т.к. этот язык не является детерминированным). Однако Т.Х.Ш. даёт подсказку, как строить КС-грамматики: надо найти в языке все скрытые «скобочные структуры».



Построение грамматики по Х.-Ш.

Построить КС-грамматику для языка $\{a^n b^m c^k \mid n = 2 * m - k\}$.

Ищем возможную скобочную структуру. Для этого сначала избавимся от вычитания: $n + k = 2 * m$. Значит, буквы a должны балансироваться буквами b справа (т.е. буквы b являются «закрывающими скобками» для a), а буквы c — буквами b слева (т.е. буквы b являются «открывающими» для c). Возможны два случая: n и k оба чётны либо оба нечётны. Построим соответствующие им разбиения: $\{a^{2*n'} b^{n'} b^{k'} c^{2*k'}\}$ и $\{a a^{2*n'} b^{n'} b b^{k'} c^{2*k'} c\}$. Дальнейшее построение грамматики уже очевидно. Заметим, что гомоморфизм подразумевает минимум четыре вида скобок: пара $(_2a,)_b$, пара $(_b,)_{2c}$, внешняя пара $[_a,]_c$ (для нечётного варианта) и $[_b,]_\varepsilon$ для него же, чтобы породить внутреннюю букву b .



Построение грамматики по Х.-Ш.

Построить КС-грамматику для языка $\{a^n b^m c^k \mid n = 2 * m - k\}$.

Как итог, получаем язык, гомоморфно порождаемый языком Дика над $\{(2a,)_b, (b,)_{2c}, [b,]_\varepsilon, [a,]_c\}$ со следующим лексером:

- ❶ До $(2a$ может идти лишь единственная $[a$.
- ❷ После $)_b$ распознаётся одна $[b$, если распозналась $[a$.
- ❸ После $)_b$ или $]_\varepsilon$ не может идти ничего другого, кроме $(b$ или $]_c$ (последняя — только после $]_\varepsilon$).
- ❹ После $)_{2c}$ не может быть ничего, кроме $)_{2c}$ или $]_c$.

Дополнительное условие на существование $[a$ уже не требуется — оно следует из сбалансированности ПСП.

Конструкция выше отличается от используемой в доказательстве теоремы — в целях экономии, в ней почти нет скобок, гомоморфно отображаемых в пустое слово.



Дополнительный пример

Построить КС-грамматику для $L_{\neq} = \{w_1cw_2 \mid w_i \in \{a, b\}^+ \text{ \& } w_1 \neq w_2\}$.

Классический пример грамматики с не-КС дополнением. Чтобы расшифровать неравенство, раскроем его в дизъюнкцию: «слово w_1 короче, чем w_2 ; либо w_2 короче, чем w_1 ; либо существует такое i , что w_1 и w_2 различаются в i -й позиции». Здесь условия не взаимоисключающие: достаточно одного из них, чтобы слово принадлежало L_{\neq} , но могут выполняться и два сразу.

Перепишем первое условие: $w_1cw'_1w_2$, где $|w_2| > 0$ и $|w_1| = |w'_1|$.

Очевидно, что «открывающими скобками» будут буквы из w_1 , «закрывающими» — из w'_1 , а «скобки» для w_2 замкнуты на самом w_2 .

Чтобы обеспечить четыре вида соответствий букв по счёту, придётся ввести четыре пары скобок для w_1 и w'_1 : $\{(a,)_a, (b,)_b, [a,]_b, [b,]_a\}$. И две пары скобок для w_2 : $\{\{a, \}_\varepsilon, \{b, \}'_\varepsilon\}$ и пара скобок для порождения c : $(c \text{ и })_\varepsilon$ (в нижних индексах — гомоморфные образы скобок).



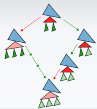
Дополнительный пример

Построить КС-грамматику для $L_{\neq} = \{w_1cw_2 \mid w_i \in \{a, b\}^+ \text{ \& } w_1 \neq w_2\}$.

Осталось построить регулярные условия. Для языка $w_1cw'_1w_2$, где $|w_2| > 0$ и $|w_1| = |w'_1|$, их можно описать следующим образом:

- ❶ После $(_a, (_b, [_a, [_b$ всегда идёт либо опять одна из таких скобок, либо $(_c$. Скобка $(_c$ единственна.
- ❷ После $)_c$, а также скобок $)_a,)_b,]_b,]_a$, могут идти либо $)_a,)_b,]_b,]_a$, либо фигурные скобки.
- ❸ В каждом слове есть хотя бы одна фигурная скобка. После открывающей фигурной скобки обязательно сразу идёт закрывающая, и после первой встреченной фигурной скобки все остальные скобки — тоже фигурные.

Представление Хомского–Шутценбергера для языка $w_1w_2cw'_1$, где $|w_2| > 0$ и $|w_1| = |w'_1|$, строится симметрично.



Дополнительный пример

Построить КС-грамматику для $L_{\neq} = \{w_1cw_2 \mid w_i \in \{a, b\}^+ \text{ \& } w_1 \neq w_2\}$.

Осталось разобрать $\{w_1t_1w_2cw_3t_2w_4 \mid |w_1| = |w_3| \text{ \& } t_1 \neq t_2\}$. Очевидно, что в нём «открывающими» будут элементы w_1 , закрывающими — элементы w_3 , t_1 и t_2 — уникальные скобки, отображающиеся в разные элементы алфавита, а w_2 и w_4 закрываются сами собой (как w_2 в предыдущем языке). Для $t_1 + t_2$ -скобок назначим пары $[_a^t,]_b^t$ и $[_b^t,]_a^t$, остальные обозначения сохраним те же. Лексер языка:

- ❶ После $(_a, (_b, [_a, [_b$ идёт либо опять одна из таких скобок, либо $[_a^t$ -скобка. $[_a^t$ единственна, за ней следует либо $\{_a$, либо $\{'_b$, либо $(_c$.
- ❷ За открывающей фигурной скобкой следует закрывающая, и после первой встреченной фигурной скобки все остальные скобки — тоже фигурные, до конца строки либо до чтения единственной $(_c$.
- ❸ После $)_\varepsilon$, а также скобок $)_a,)_b,]_b,]_a$, могут идти либо $)_a,)_b,]_b,]_a$, либо скобка $]_a^t$. За $]_a^t$ следует EOL или $\{_a$ или $\{'_b$.

Чтобы получить прообраз языка L_{\neq} , объединим все три лексера.