

Задание 6

Грамматика

Ключевые слова¹: язык, регулярный язык, ДКА, НКА, алгебра регулярных выражений, грамматики, уравнения с регулярными коэффициентами.

1 Грамматика

Одна из больших проблем науки, которую мы с вами изучаем – определения. Их слишком много и они отличаются друг от друга, хотя в итоге конечно описывают одни и те же классы языков. Я призываю на экзамене пользоваться определениями из книги Серебрякова, хотя при выполнении задания вы можете пользоваться эквивалентными определениями из другой литературы.

Определение 1. Грамматика Γ определяется через

- N – множество нетерминальных символов
- T – множество терминальных символов
- P – множество правил вывода, $P \subseteq (N \cup T)^* \times (N \cup T)^*$.
- S – аксиома, $S \in N$.

Все множества из описания грамматики конечные. При этом, $N \cap T = \emptyset$. Принято обозначение $\Gamma = G(N, T, P, S)$. При описании грамматики приняты следующие соглашения. Нетерминалы обозначают заглавными буквами A, B, C, \dots терминалы обозначают строчными буквами, смешанные цепочки из $(N \cup T)^*$ обозначают греческими буквами α, β, γ . Слово $w \in T^*$ порождается грамматикой Γ , если существует последовательность правил вывода, начинающаяся с правила вида $S \rightarrow \alpha$, в результате применения которых порождается слово w . Под применением правила $\alpha \rightarrow \beta$, понимается, что подслово α заменяется на подслово β

¹минимальный необходимый объем понятий и навыков по этому разделу)

В зависимости от ограничений, налагаемых на правила вывода, получаются разные классы языков. В рамках этого задания нас пока интересует только два последних типа.

- Если на множество правил P не накладывается ограничений, то есть правила имеют вид $\alpha \rightarrow \beta$, то грамматика называется грамматикой типа 0 по Хомскому
- Грамматика, в которых правила имеют вид $\alpha A \beta \rightarrow \alpha \gamma \beta$, $|\gamma| > 0$ называются грамматиками типа 1 или Контекстно-зависимыми. В качестве исключения грамматике может принадлежать правило $S \rightarrow \varepsilon$, но тогда нетерминал S не может встречаться в правых частях.
- Грамматика, в которых правила имеют вид $A \rightarrow \alpha$, называются грамматиками типа 2 или Контекстно-Свободными грамматиками.
- Грамматика, в которых правила имеют вид $A \rightarrow xB$ или $A \rightarrow x$, $x \in T^*$, называются грамматиками типа 3 или праволинейными грамматиками.

В определении КЗ-грамматики существенно, что она является *неукорачивающей*, т.е. правая часть правил всегда длиннее левой. Эквивалентное определение из книги Серебрякова гласит, что в КЗ-грамматике все правила, кроме быть может $S \rightarrow \varepsilon$, имеют вид $\alpha \rightarrow \beta$, $|\alpha| < |\beta|$. Опять-таки, если есть правило $S \rightarrow \varepsilon$, то нетерминал S в правых частях правил встречаться не может.

Очень часто грамматиками типа 3 называют грамматика, в которых правила вывода имеют вид $A \rightarrow xB$ или $A \rightarrow x$, $x \in T$, также допускается правило $S \rightarrow \varepsilon$ с всё той же оговоркой, что аксиома не может встречаться в правой части. Такие грамматика называются *праволинейными регулярными* грамматиками.

Упражнение 1. Доказать, что праволинейные грамматика и праволинейные регулярные грамматика эквивалентны, т.е. порождают один и тот же тип языков.

Определение 2. Грамматика типа 3 является *неоднозначной*, если существует более одного способа вывести хотя бы одно слово из языка, порождённого грамматикой.

Для КС-грамматик это определение неприемлемо. Вдумчивый читатель может подумать почему, прежде чем переходить к следующему разделу.

Левoliniейные грамматики определяются аналогично праволинейным: в них правила имеют вид $A \rightarrow Bx$ или $A \rightarrow x$.

Приведём пример КС-грамматики. При записи правил используют вспомогательное обозначение $A \rightarrow \alpha|\beta$, которое означает, что в грамматике есть два правила: $A \rightarrow \alpha$ и $A \rightarrow \beta$.

Пример 1. Грамматика G задана правилами:

$$\begin{aligned} S &\rightarrow aAB \\ A &\rightarrow aA|a \\ B &\rightarrow bB|b \end{aligned}$$

Слово $aabb$ выводится грамматикой. Последовательность применений правил вывода такая:

$$\begin{aligned} S &\rightarrow aAB \\ a\underline{A}B & \\ A &\rightarrow a \\ aa\underline{B} & \\ B &\rightarrow bB \\ aab\underline{B} & \\ B &\rightarrow b \\ aabb & \end{aligned}$$

1.1 Вывод, левый вывод, дерево разбора

Выводом цепочки α называется такая последовательность применений правил с указанием раскрываемого нетерминала, что применяя правила из неё начиная с аксиомы получается цепочка α . Если цепочка α не содержит нетерминалов, то α принадлежит языку, порождаемому КС-грамматикой. Нам будет удобно пользоваться такими понятиями как левый вывод и правый вывод. *Левым выводом* называют такой вывод, что

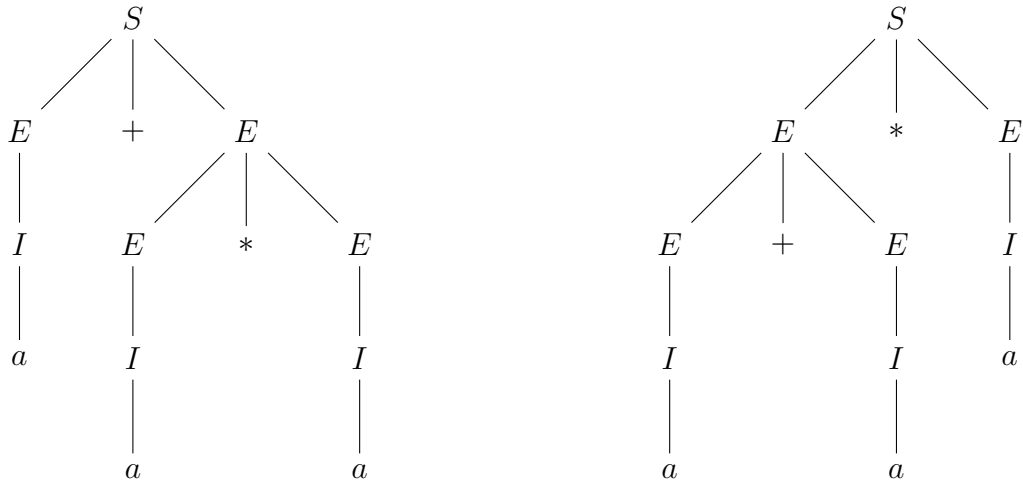
на каждом его шаге раскрывается самый левый нетерминал в промежуточной цепочке. Вывод в примере 1 является левым. Правый вывод определяется аналогично.

Также мы будем использовать деревья вывода. С формальным определением дерева вывода вы можете познакомиться, например, в книге Хопкрофта, Мотвани и Ульмана, а мы рассмотрим пример деревьев вывода и затем дадим неформальное описание этого понятия.

Пример 2. Грамматика G , $T = \{a, +, *\}$ задана правилами:

$$\begin{aligned} S &\rightarrow E \\ E &\rightarrow E + E \\ E &\rightarrow E * E \\ E &\rightarrow I \\ I &\rightarrow a \end{aligned}$$

Построим деревья вывода для слова $a + a * a$:



Теперь опишем понятие дерева вывода, которое также называется деревом разбора. Зафиксируем грамматику G . *Деревом вывода* для слова w называется упорядоченное дерево, в корне которого находится аксиома S , каждая вершина помечена нетерминалом, терминалом или пустым словом, если вершина помечена терминалом или ε , то эта вершина

является листом, если же вершина помечена нетерминалом A , то для некоторого правила $A \rightarrow X_1 X_2 \dots X_n \in P$ ($X_i \in N \cup T$) вершины-дети A помечены символами $X_1, X_2 \dots X_n$ слева направо. Листья дерева вывода образуют слово w .

Как мы видим, синтаксически деревья вывода принципиально разные: в первом случае выражение интерпретируется как $a + (a * a)$, а во втором как $(a + a) * a$, что приводит к непредсказуемому результату при выполнении стандартных операций!

Эта проблема приводит нас к новому важному понятию – неоднозначности. Грамматика G является *неоднозначной*, если хотя бы для одного слова существует два различных дерева вывода.

Левый и правый вывод помимо удобства важны тем, что они фактически задают порядок обхода дерева вывода, поэтому каждому левому выводу соответствует ровно одно дерево разбора, а каждому дереву разбора соответствует ровно один левый вывод.

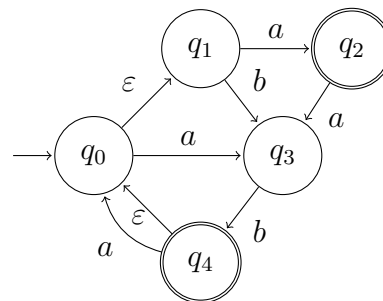
Упражнение 2. Доказать, что грамматика G является однозначной тогда и только тогда, когда каждое слово, порождаемое G имеет ровно один левый(правый) вывод.

2 Задачи

Внимание, все задачи на построение автоматов должны быть снабжены диаграммами!

Задача 1.

Предложите алгоритм построения праволинейной грамматики по автомату \mathcal{A} и докажете его корректность. Постройте по автомату \mathcal{A} регулярную праволинейную грамматику G по предложенному алгоритму. Если вместо построения своего алгоритма, Вы возьмёте его из книжки, укажите это и не списывайте страницами, пожалуйста.



Задача 2. Постройте автомат по грамматике G :

$$S \rightarrow abaA|abB|\varepsilon, A \rightarrow aB|aa, B \rightarrow bA|aS$$

Задача 3. Является ли грамматика G из предыдущей задачи однозначной?

Задача 4. Верно ли, что праволинейная грамматика G однозначна тогда и только тогда, когда построенный по ней автомат является детерминированным?

Задача 5. Назовём грамматику линейной, если в правой части её правил может быть не более одного нетерминала. Верно ли, что для любой линейной грамматики G , $L(G) \in \text{REG}$?

Ещё раз напоминаю, что задачи, помеченные \dagger являются дополнительными, поэтому списывать их из книжек – бессмысленное увеличение энтропии.

Определение 3. Для языка $L \subseteq \{\sigma_1, \sigma_2, \dots, \sigma_n\}^* = \Sigma_n^*$ и языков $L_{\sigma_1}, L_{\sigma_2}, \dots, L_{\sigma_n} \subseteq \Sigma_n^*$, подстановкой в L языков $L_{\sigma_1}, \dots, L_{\sigma_n}$ назовём язык L' , такой что для всех слов $w = w[1] \dots w[n]$ из языка L справедливо $L_{w[1]}L_{w[2]} \dots L_{w[n]} \subseteq L'$

Задача 6[†]. Доказать, что регулярные языки замкнуты относительно операции подстановки.

Определение 4. Даны алфавиты Σ и Δ . Для языка $L \subseteq \Sigma \times \Delta$ определены операции проекции на Σ^* и Δ^* . Проекцией L на Σ^* называется язык $L_\Sigma = \{w \in \Sigma^* \mid \exists v \in \Delta^* : (w, v) \in L\}$. Проекция L на Δ^* определяется аналогичным образом.

Задача 7[†]. Доказать, что регулярные языки замкнуты относительно операции проекции.

Определение 5. Для языка $L_\Sigma \subseteq \Sigma^*$, Δ -цилиндром называется язык L , такой что $L = \{w \mid w = (u, v), u \in L_\Sigma, v \in \Delta^*\}$

Задача 8[†]. Показать, что Σ -проекция Δ -цилиндра L есть L_Σ . Доказать, что регулярные языки замкнуты относительно операции цилиндра.

Задача 9. На семинаре я «доказал», что грамматика $G : S \rightarrow aSb \mid SS \mid \varepsilon$ порождает язык правильных скобочных выражений с одним типом скобок – язык Дика D_1 . На самом деле, я дал доказательство только в одну сторону: что любое слово, выведенное из этой грамматики будет правильным скобочным выражением. После чего порадовавшись, что меня никто за руку не поймал, в назидание оставляю доказательство в другую сторону в качестве домашней задачи. Напомню, что мы договорились считать, что слово w является правильным скобочным выражением, если его скобочный итог $d(w)$ равен нулю, и при этом скобочный итог любого префикса w неотрицательный. Скобочным итогом называется разница между числом открывающих и закрывающих скобок: $d(u) = |u|_a - |u|_b$.

Задача 10. Является ли грамматика G для языка D_1 из предыдущей задачи однозначной?