

Автомат Глушкова



Линеаризация

Определение

Если регулярное выражение $r \in \mathcal{RE}$ содержит n вхождений букв алфавита Σ , тогда линеаризованное регулярное выражение $\text{Linearize}(r)$ получается из r приписыванием i -ой по счёту букве, входящей в r , индекса i .

Пример

Рассмотрим регулярное выражение:

$$(ba \mid b)aa(a \mid ab)^*$$

Его линеаризованная версия:

$$(b_1a_2 \mid b_3)a_4a_5(a_6 \mid a_7b_8)^*$$



Множества First, Last, Follow

Определение

Пусть $r \in \mathcal{RE}$, тогда:

- множество First — это множество букв, с которых может начинаться слово из $\mathcal{L}(r)$ (если $\varepsilon \in \mathcal{L}(r)$, то оно формально добавляется в First);
- множество Last — это множество букв, которыми может заканчиваться слово из $\mathcal{L}(r)$;
- множество Follow(c) — это множество букв, которым может предшествовать c. Т.е.
 $\{d \in \Sigma \mid \exists w_1, w_2 (w_1 c w_2 \in \mathcal{L}(r))\}.$

Множество Follow в теории компиляции обычно определяется иначе — это множество символов, которые могут идти за выводом из определённого нетерминального символа. Два этих определения можно унифицировать, если рассматривать каждую букву в r как «обёрнутую»



First, Last, Follow — пример

Построим указанные множества для регулярного выражения
 $r = (ba \mid b)aa(a \mid ab)^*$.

Начнём с исходного регулярного выражения.

Исходное регулярное выражение

- $\text{First}(r) = \{b\}$.
- $\text{Last}(r) = \{a, b\}$.
- $\text{Follow}_r(a) = \{a, b\}$; $\text{Follow}_r(b) = \{a\}$.

Хотя данные множества описывают, как устроены слова из $\mathcal{L}(r)$ локально, однако они не исчерпывают всей информации о языке, поскольку разные вхождения букв в регулярное выражения никак не различаются.

Например, по множествам First и Last можно предположить, что $b \in \mathcal{L}(r)$, хотя это не так.



First, Last, Follow — пример

Построим указанные множества для регулярного выражения
 $r = (ba \mid b)aa(a \mid ab)^*$.

Вспомним, что $r_{\text{Lin}} = (b_1a_2 \mid b_3)a_4a_5(a_6 \mid a_7b_8)^*$.

Линеаризованное выражение

- $\text{First}(r_{\text{Lin}}) = \{b_1, b_3\}$.
- $\text{Last}(r_{\text{Lin}}) = \{a_5, a_6, b_8\}$.
- $\text{Follow}_{r_{\text{Lin}}}(b_1) = \{a_2\}$; $\text{Follow}_{r_{\text{Lin}}}(a_2) = \{a_4\}$;
 $\text{Follow}_{r_{\text{Lin}}}(b_3) = \{a_4\}$; $\text{Follow}_{r_{\text{Lin}}}(a_4) = \{a_5\}$;
 $\text{Follow}_{r_{\text{Lin}}}(a_5) = \{a_6, a_7\}$; $\text{Follow}_{r_{\text{Lin}}}(a_6) = \{a_6, a_7\}$;
 $\text{Follow}_{r_{\text{Lin}}}(a_7) = \{b_8\}$; $\text{Follow}_{r_{\text{Lin}}}(b_8) = \{a_6, a_7\}$.

В описании данных множеств содержится исчерпывающая информация о языке $\mathcal{L}(r_{\text{Lin}})$.



Конструкция автомата Глушкова

Алгоритм построения $\text{Glushkov}(r)$

- Строим линеаризованную версию r : $r_{\text{Lin}} = \text{Linearize}(r)$.
- Ищем $\text{First}(r_{\text{Lin}})$, $\text{Last}(r_{\text{Lin}})$ и $\text{Follow}_{r_{\text{Lin}}}(c)$ для всех $c \in \Sigma_{r_{\text{Lin}}}$.
- Все состояния автомата, кроме начального (назовём его S), соответствуют буквам $c \in \Sigma_{r_{\text{Lin}}}$.
- Из начального состояния строим переходы в те состояния, для которых $c \in \text{First}(r_{\text{Lin}})$. Переходы имеют вид $S \xrightarrow{c}$.
- Переходы из состояния c соответствуют элементам d множества $\text{Follow}_{r_{\text{Lin}}}(c)$ и имеют вид $c \xrightarrow{d}$.
- Конечные состояния — такие, что $c \in \text{Last}(r_{\text{Lin}})$, а также S , если $\varepsilon \in \mathcal{L}(R)$.
- Теперь стираем разметку, построенную линеаризацией, на переходах автомата. Конструкция завершена.



Пример автомата Глушкова

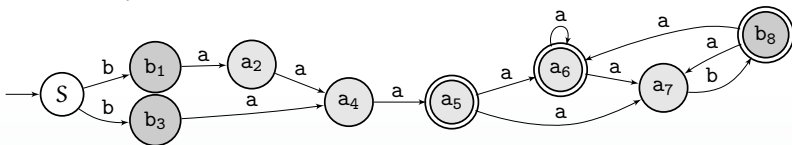
Исходное регулярное выражение:

$$(ba \mid b)aa(a \mid ab)^*$$

Линеаризованное регулярное выражение:

$$(b_1a_2 \mid b_3)a_4a_5(a_6 \mid a_7b_8)^*$$

Автомат Глушкова:



Подграфы, распознающие регулярные выражения, являющиеся подструктурами исходного, не имеют общих вершин. Это свойство автомата Глушкова используется в реализациях `match`-функций некоторых библиотек регулярных выражений.



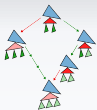
Свойства автомата Глушкова

- Не содержит ε -переходов.
- Число состояний равно длине регулярного выражения (без учёта регулярных операций), плюс один (стартовое состояние).
- В общем случае недетерминированный.

Примечание

Для 1-однозначных регулярных выражений r автомат $\text{Glushkov}(r)$ является детерминированным. Эту его особенность активно используют в современных библиотеках регулярных выражений, например, в RE2. Выигрыш может получиться колоссальным: например, $\text{Thompson}((a^*)^*)$ является экспоненциально неоднозначным, а $\text{Glushkov}((a^*)^*)$ однозначен и детерминирован!

Бисимуляция



Labelled Transition Systems

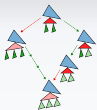
Понятие бисимуляции возникло в контексте систем размеченных переходов (LTS).

Определение

Labelled Transition System — тройка $\langle S, \Sigma, Q \rangle$, где S — множество состояний, Σ — множество меток, Q — множество переходов (троек из $S \times \Sigma \times S$).

LTS похожи на конечные автоматы, но допускают бесконечные множества S и Q . Кроме того, в LTS нет начальных и финальных состояний.

Трансформационный моноид также строится в контексте LTS, то есть без учёта финальности состояний. Поэтому из ДКА, по которому строится трансформационный моноид, предварительно удаляются все ловушки, иначе в нём могут появиться правила переписывания, не имеющие никакого отношения к языку ДКА.



Симуляция и бисимуляция

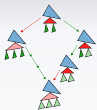
Определение

Если \lesssim — симуляция для $LTS = \langle S, \Sigma, Q \rangle$, то

$$\forall p, q \in S (p \lesssim q \Rightarrow (\exists p', a((p \xrightarrow{a} p') \Rightarrow \exists q'(q \xrightarrow{a} q' \ \& \ p' \lesssim q'))))$$

Если одновременно выполняются условия $p \lesssim q$ и $q \lesssim p$, то говорят, что p и q находятся в отношении бисимуляции (обозначается $p \sim q$).

Можно считать, что если $p \lesssim q$, то множество путей в LTS, стартующих в p , вкладывается в множество путей с началом в q . Бисимуляция состояний в единственной LTS легко обобщается и на бисимуляцию между двумя разными LTS. Поскольку в них нет начальных состояний, и они не обязаны быть связными, можно рассматривать несколько LTS как одну LTS с несколькими компонентами и искать бисимуляцию между элементами этих компонент.



Бисимилярность НКА

Чтобы определить отношение бисимуляции на конечных автоматах, к отношению бисимуляции на LTS нужно добавить ограничения на бисимуляцию начальных и конечных состояний. Более точно, для бисимуляции НКА \mathcal{A}_1 и \mathcal{A}_2 необходимы следующие условия:

- 1 каждому состоянию \mathcal{A}_1 бисимилярно состояние \mathcal{A}_2 , и наоборот;
- 2 стартовому состоянию \mathcal{A}_1 бисимилярно стартовое состояние \mathcal{A}_2 ;
- 3 каждому финальному состоянию \mathcal{A}_1 бисимилярно финальное состояние \mathcal{A}_2 , и наоборот.

Лемма

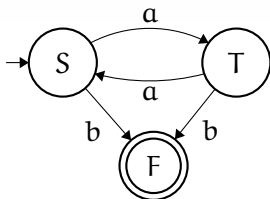
Бисимилярные НКА распознают равные языки.



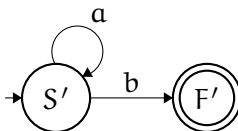
Пример бисимилярных НКА

Рассмотрим следующие два автомата, распознающие язык a^*b .

Автомат \mathcal{A}_1 :



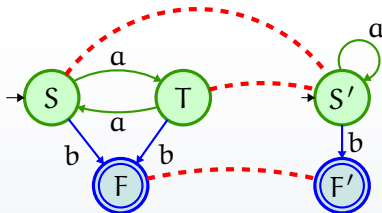
Автомат \mathcal{A}_2 :



Их бисимуляция:

$$\{\langle S, S' \rangle, \langle T, S' \rangle, \langle F, F' \rangle\}$$

Состояния S и T бисимилярны одному и тому же состоянию S' .





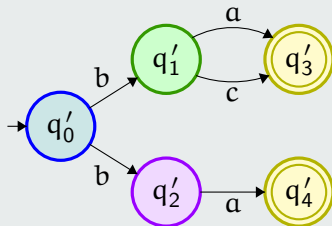
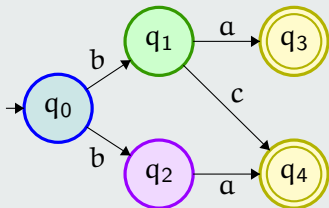
Бисимуляция и равенство

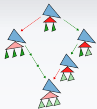
В равных НКА состояния бисимилярны, однако только условия существования бисимуляции и биекции бисимилярных состояний недостаточно, чтобы гарантировать равенство.

Пример неравных бисимилярных НКА

(автор примера: А. Д. Дельман)

Следующие два автомата бисимилярны и имеют одинаковое число состояний, однако не равны:



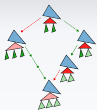


Бисимилярность состояний в НКА

Определение

Состояния q_i, q_j в НКА \mathcal{A} бисимилярны ($q_i \sim_{\mathcal{A}} q_j$), если они связаны LTS-бисимуляцией и имеют одинаковую финальность в \mathcal{A} .

- С учётом определения выше, бисимуляцию НКА можно переформулировать как отношение бисимуляции состояний НКА такое, что стартовые состояния бисимилярны.
- Отношение $\sim_{\mathcal{A}}$ имеет важное свойство: бисимилярные состояния в автомате можно объединить без изменения его семантики. Это преобразование часто позволяет существенно упростить НКА.



Слияние по бисимуляции

Бисимилярность состояний в НКА

Определение

Состояния q_i, q_j в НКА \mathcal{A} бисимилярны ($q_i \sim_{\mathcal{A}} q_j$), если они связаны LTS-бисимуляцией и имеют одинаковую финальность в \mathcal{A} .

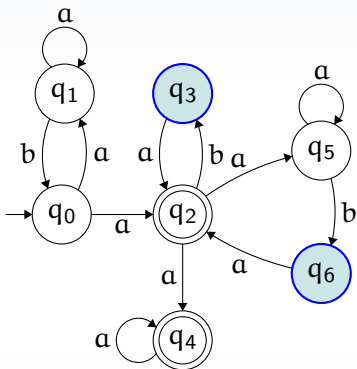
Пример

Все состояния-ловушки в любом полном автомате (т.е. с явно присутствующими переходами по всем буквам алфавита) бисимилярны друг другу. Все финальные состояния без переходов из них (кроме как в ловушки) также бисимилярны.

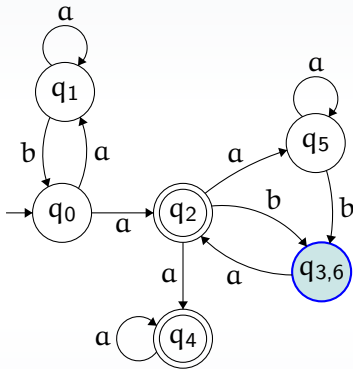


Пример слияния по бисимуляции

Исходный автомат:



Итоговый автомат:



Бисимуляция: $\{\{q_3, q_6\}, \bigcup_{i \neq 3 \wedge i \neq 6} \{q_i\}\}$

Кроме q_3 и q_6 , все состояния не бисимилярны никаким другим.

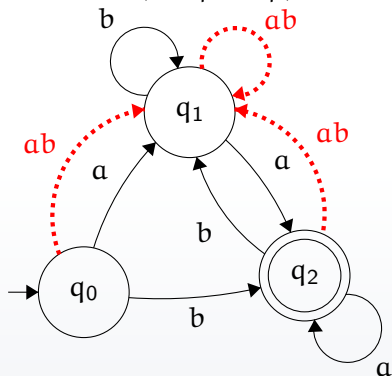
Например, $q_1 \not\sim q_5$, поскольку $q_1 \xrightarrow{b} q_0$, $q_5 \xrightarrow{b} q_6$, но $q_0 \xrightarrow{a} q_1$ (и q_1 — не финальное), а из q_6 есть переход только в финальное состояние q_2 .

Трансформационный МОНОИД



Функции переходов по слову в ДКА

Правила перехода в ДКА \mathcal{A} над алфавитом Σ и множеством состояний Q определяются функцией $\Sigma \times Q \rightarrow Q$. Если специализировать её по первому аргументу, получится функция $F_\xi : Q \rightarrow Q$ ($\xi \in \Sigma$). Эту функцию можно продолжить на строки, положив $F_\xi \circ F_\eta = F_{\eta\xi}$.



Пусть мы строим функцию переходов по слову ab в автомате \mathcal{A} . Сначала определим функции F_a , F_b , определяющие его поведение на буквах a и b . Тогда поведение переходов на слове ab получится композицией F_a и F_b .

	q_0	q_1	q_2
a	q_1	q_2	q_2
b	q_2	q_1	q_1
ab	q_1	q_1	q_1



Функции переходов по слову в ДКА

Свойства множества функций переходов ДКА \mathcal{A}

- Существует единичная функция F_ε такая, что $F_\varepsilon \circ F_\xi = F_\xi \circ F_\varepsilon = F_\xi$.
- Композиция \circ ассоциативна.

Таким образом, функции переходов по словам из Σ^* в ДКА \mathcal{A} образуют моноид относительно композиции.

Если ДКА представлен в краткой (trim) форме, некоторые переходы могут вести «в никуда». На самом деле они ведут в (единственное!) состояние-ловушку, существование которого неявно подразумевается. Однако наличие нескольких ловушек в ДКА повлечёт ошибки при построении функции переходов.



Определение и свойства

Определение

Трансформационный моноид $\mathcal{M}_{\mathcal{A}}$ для ДКА \mathcal{A} — это моноид функций F_{ξ} таких, что $F_{\xi}(q_i) = q_j \Leftrightarrow (q_i \xrightarrow{\xi} q_j \text{ в } \mathcal{A})$. Иначе можно сказать, что трансформационный моноид $\mathcal{M}_{\mathcal{A}}$ определяется множеством классов эквивалентности $\{w \mid w \in \Sigma^+\}$ таким, что $w_i = w_j \Leftrightarrow F_{w_i} = F_{w_j}$.

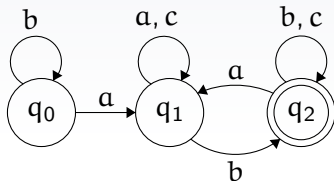


Определение и свойства

- $\mathcal{M}_{\mathcal{A}}$ определяется фактормножеством классов эквивалентности и правилами переписывания, задающими эквивалентность. ε обычно не включается в множество w_i .
- Поскольку множество функций F_{w_i} в случае ДКА конечно, то $\mathcal{M}_{\mathcal{A}}$ содержит конечное число классов эквивалентности (верно и обратное: каждый такой моноид определяет ДКА).
- Трансмоноид строится для ДКА без ловушек; переход в ловушку обозначается в таблице переходов просто прочерком.
- Для единообразия записи трансформаций и перестановок в алгебре, в таблице переходов пишут только номера состояний \mathcal{A} .



Построение трансф. моноида

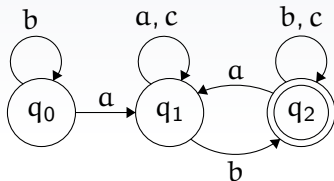


Определим соответствие между буквами и множествами переходов по ним и будем расширять этот список новыми словами в лексикографическом порядке. Если очередное слово задаёт такую же трансформацию, как и уже рассмотренное, порождаем соответствующее правило переписывания.

Классы эквивалентности				Правила переписывания	
	0	1	2		
a	1	1	1		
b	0	2	2		
c	—	1	2		



Построение трансф. моноида



Классы эквивалентности				Правила переписывания	
	0	1	2		
a	1	1	1	$aa \rightarrow a$	$ac \rightarrow a$
b	0	2	2	$ba \rightarrow a$	$bb \rightarrow b$
c	—	1	2	$cb \rightarrow bc$	$cc \rightarrow c$
ab	2	2	2	$abc \rightarrow ab$	$bca \rightarrow ca$
bc	—	2	2	$cab \rightarrow bc$	
ca	—	1	1		

Всего классов эквивалентности: 6



Синтаксический моноид

Определим отношение синтаксической конгруэнтности слов:

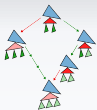
$$w_i \sim_{\mathcal{L}} w_j \Leftrightarrow \forall x, y (x w_i y \in \mathcal{L} \Leftrightarrow x w_j y \in \mathcal{L})$$

Синтаксический моноид $\mathcal{M}(\mathcal{L})$ — это множество его классов эквивалентности относительно $\sim_{\mathcal{L}}$. То есть такая полугруппа с единицей над $w \in \Sigma^*$, что $w_i = w_j \Leftrightarrow w_i \sim_{\mathcal{L}} w_j$ (равенство здесь понимается в алгебраическом смысле: как возможность преобразовать w_i и w_j к одному и тому же слову).

Лемма

Синтаксический моноид регулярного языка \mathcal{L} совпадает с трансф. моноидом минимального ДКА, его распознающего.

Синтаксический моноид (так же, как и минимальный ДКА) — атрибут *языка*, а трансф. моноид — атрибут *конкретного ДКА*.



Суффиксная конгруэнтность

Предшествующие понятия рассматривали структуру переходов автомата без учёта начальных и конечных состояний, хотя неявно они использовались, чтобы удалить недостижимые состояния и состояния-ловушки при подготовке к построению моноида.

Однако если чуть-чуть специализировать отношение $\sim_{\mathcal{L}}$, положив возможные префиксы пустыми, получится отношение эквивалентности, напрямую зависящее от положения стартовых и финальных состояний в минимальном ДКА.

Определим отношение эквивалентности по Нероуду как:

$$w_i \equiv_{\mathcal{L}} w_j \Leftrightarrow \forall y (w_i y \in \mathcal{L} \Leftrightarrow w_j y \in \mathcal{L})$$

Обозначение $\sim_{\mathcal{L}}$ может использоваться в литературе как в смысле синтаксической конгруэнции, так и в смысле эквивалентности по Нероуду. Лучше дополнительно уточнить.



Критерий регулярности языка

Теорема Майхилла–Неруда

Язык \mathcal{L} регулярен тогда и только тогда, когда множество классов эквивалентности по $\equiv_{\mathcal{L}}$ конечно.

\Rightarrow : Пусть \mathcal{L} регулярен. Тогда он порождается некоторым DFA \mathcal{A} с конечным числом состояний N . Значит, множество $\{q_i \mid q_0 \xrightarrow{w} q_i\}$ конечно, а для любых двух w_1, w_2 таких, что $q_0 \xrightarrow{w_1} q_i$ и $q_0 \xrightarrow{w_2} q_i$, выполняется $w_1 \equiv_{\mathcal{L}} w_2$.

\Leftarrow : Пусть все слова в Σ^* принадлежат N классам эквивалентности A_1, \dots, A_n по $\equiv_{\mathcal{L}}$. Построим по ним DFA \mathcal{A} , распознающий \mathcal{L} . Классы A_i сопоставим состояниям:

- Начальным объявим класс эквивалентности A_0 такой, что $\varepsilon \in A_0$.
- Конечными объявим такие A_j , что $\forall w \in A_j (w \in \mathcal{L})$.
- Если $w \in A_i$, $w a_k \in A_j$, тогда добавляем в δ правило $\langle A_i, a_k, A_j \rangle$.

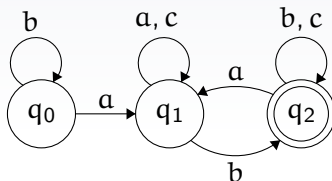


Трансформационный моноид и $\equiv_{\mathcal{L}}$

- Классы эквивалентности по Майхиллу–Нероуду можно извлечь из трансформационного моноида минимального ДКА для языка \mathcal{L} : они являются подмножеством факторслов, которые переводят стартовое состояние в какое-то другое (непустое) состояние.
- Для каждой пары таких факторслов w_i и w_j , переводящих стартовое состояние в разные состояния q_i и q_j , в синтаксическом моноиде обязательно найдётся различающий суффикс (т.е. класс эквивалентности u такой, что $w_i u \in \mathcal{L}$ & $w_j u \notin \mathcal{L}$, либо наоборот).
- Если в ДКА существует ловушка (возможно, неявная), то в трансформационном моноиде найдётся класс эквивалентности, переводящий в неё стартовое состояние. Таких классов может быть несколько, но с точки зрения эквивалентности по Майхиллу–Нероуду, они не различаются.



Пример



Включим в число факторслов ε и выделим по одному факторслову для каждого состояния q_i , переводящему стартовое слово в q_i . Для каждого из них определим множество суффиксов, которые оставляют слова этих классов в языке. После этого достаточно собрать вместе все суффиксы и префиксы и выкинуть из полученной таблицы дубли столбцов.

Факторслова–префиксы			Таблица классов эквивалентности			
префикс	$0 \rightarrow ?$	суффиксы		ε	b	ab
ε	0	ab	ε	—	—	+
a	1	b, ab	a	—	+	+
ab	2	ε, b, ab	ab	+	+	+