

Lecture 4 – Regular expressions, Kleene's theorem, string substitution

NTIN071 Automata and Grammars

Jakub Bulín (KTIML MFF UK)

Spring 2024

** Adapted from the Czech-lecture slides by Marta Vomlelová with gratitude.
The translation, some modifications, and all errors are mine.*

Recap of Lecture 3

- Nondeterministic finite automata (NFA): can 'guess' the right path to accepting, computation described by a state tree.
- ϵ -transitions: allow to change states without reading any input
- Subset construction: every NFA and ϵ -NFA is equivalent to a DFA (but can be easier to design, much smaller).
- Regular languages are closed under set operations (union, intersection, complement, difference)
- And under string operations (concatenation, iteration and positive iteration, reverse, left and right quotient)

1.8 Regular expressions

Regular expressions (RE)

- an algebraic description of languages
- declarative: express the form of the words we want to accept
- can describe all, and only, regular languages
- can be viewed as a programming language, a user/friendly description of a finite automaton

Example

- grep command in UNIX.
- Python module re
- lexical analysis, e.g. Flex (description via 'tokens' \leftrightarrow RE)

Note: syntax analysis needs a stronger tool, context-free grammars

The definition

A **regular expression** α over (finite, nonempty) Σ , $\alpha \in \text{RegE}(\Sigma)$ and the matching language $L(\alpha)$, are defined inductively:

expression	language	note
\emptyset	$L(\emptyset) = \emptyset$	empty expression
ϵ	$L(\epsilon) = \{\epsilon\}$	empty string
a	$L(\mathbf{a}) = \{a\}$	for all $a \in \Sigma$
$(\alpha + \beta)$	$L((\alpha + \beta)) = L(\alpha) \cup L(\beta)$	union (grep, re use ' ')
$(\alpha\beta)$	$L((\alpha\beta)) = L(\alpha)L(\beta)$	concatenation
α^*	$L(\alpha^*) = L(\alpha)^*$	iteration (Kleene star)

Examples, notation

Example

- The language of alternating 0s and 1s can be expressed as:
 - $(01)^* + (10)^* + 1(01)^* + 0(10)^*$
 - $(\epsilon + 1)(01)^*(\epsilon + 0)$
- $L((0^*10^*10^*1)^*0^*) = \{w \in \{0, 1\}^* \mid |w|_1 \equiv 0 \pmod{3}\}$

We often omit parentheses:

- priority of operators: iteration $*$ $>$ concatenation $>$ union $+$
- associativity of concatenation, union $+$
- outer parentheses

We could define, and will sometimes use, positive iteration α^+

Theorem (Kleene's theorem)

A language is regular, iff it is matched by some regular expression.

We will prove it by giving two constructions:

1. from RE to ϵ -NFA (which can be converted to a DFA)
2. from a DFA to a RE (but we could start from a ϵ -NFA)

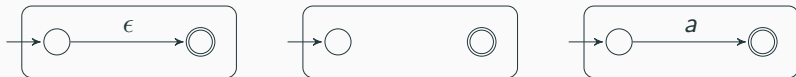
For 2. we also mention a better algorithm: **state elimination**

RE to ϵ -NFA

By induction on the structure of α , construct a ϵ -NFA E s.t.
 $L(\alpha) = L(E)$ with three additional properties:

1. Exactly one accepting state.
2. No incoming edges into the initial state.
3. No outgoing edges from the accepting state.

Induction base: α is the empty string ϵ , empty set \emptyset , or a letter **a**



Induction step: $\alpha + \beta$, $\alpha\beta$, α^* (next slide)



RE to ϵ -NFA: Induction step

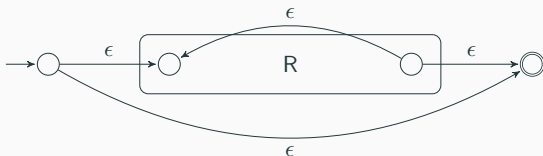
Addition $\alpha + \beta$



Concatenation $\alpha\beta$



Iteration α^*



Assume the states are $Q = \{1, \dots, n\}$ and the start state is $q_0 = 1$.

Construct a RE $R_{ij}^{(k)}$ matching words that transition from state i into state j and all intermediate states (if any) have index $\leq k$.

Then we set $\alpha = \sum_{j \in F_A} R_{1j}^{(n)}$ (from start to some accepting state)

Iteratively construct $R_{ij}^{(k)}$ for $k = 0, \dots, n$ (finite induction).

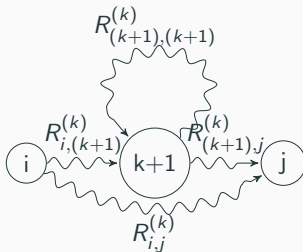
Induction base: $k = 0$

- If $i \neq j$, set $R_{ij}^{(0)} = \mathbf{a_1} + \dots + \mathbf{a_m}$ where a_1, \dots, a_m are symbols on edges from i into j ($R_{ij}^{(0)} = \emptyset$ or $R_{ij}^{(0)} = \mathbf{a}$ for $m = 0, 1$).
- If $i = j$, $R_{ii}^{(0)} = \epsilon + \mathbf{a_1} + \dots + \mathbf{a_m}$ where a_i 's are on loops on i .

DFA to RE: Induction step

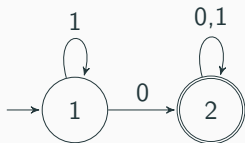
Once we have $R_{ij}^{(k)}$ for all $i, j \in Q$, we can construct $R_{ij}^{(k+1)}$:

$$R_{ij}^{(k+1)} = R_{ij}^{(k)} + R_{i,(k+1)}^{(k)} (R_{(k+1),(k+1)}^{(k)})^* R_{(k+1),j}^{(k)}$$



- paths $i \rightsquigarrow j$ not going through $k+1$: already in $R_{ij}^{(k)}$
- paths $i \rightsquigarrow j$ going through $k+1$ one or more times: $i \rightsquigarrow k+1$ (first visit), loop on $k+1$, finally (last visit) $k+1 \rightsquigarrow j$ □

Example



Apply the construction, simplify:

$$\alpha = R_{12}^{(2)} = \mathbf{1^*0(0+1)^*}$$

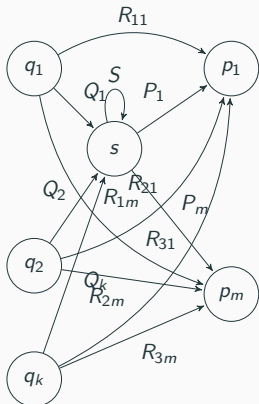
$R_{11}^{(0)}$	$\epsilon + \mathbf{1}$	$=$
$R_{12}^{(0)}$	$\mathbf{0}$	$=$
$R_{21}^{(0)}$	\emptyset	$=$
$R_{22}^{(0)}$	$(\epsilon + \mathbf{0} + \mathbf{1})$	$=$
$R_{11}^{(1)}$	$\epsilon + \mathbf{1} + (\epsilon + \mathbf{1})(\epsilon + \mathbf{1})^*(\epsilon + \mathbf{1})$	$= \mathbf{1^*}$
$R_{12}^{(1)}$	$\mathbf{0} + (\epsilon + \mathbf{1})(\epsilon + \mathbf{1})^*\mathbf{0}$	$= \mathbf{1^*0}$
$R_{21}^{(1)}$	$\emptyset + \emptyset(\epsilon + \mathbf{1})^*(\epsilon + \mathbf{1})$	$= \emptyset$
$R_{22}^{(1)}$	$\epsilon + \mathbf{0} + \mathbf{1} + \emptyset(\epsilon + \mathbf{1})^*\mathbf{0}$	$= \epsilon + \mathbf{0} + \mathbf{1}$
$R_{11}^{(2)}$	$\mathbf{1^*} + \mathbf{1^*0}(\epsilon + \mathbf{0} + \mathbf{1})^*\emptyset$	$= \mathbf{1^*}$
$R_{12}^{(2)}$	$\mathbf{1^*0} + \mathbf{1^*0}(\epsilon + \mathbf{0} + \mathbf{1})^*(\epsilon + \mathbf{0} + \mathbf{1})$	$= \mathbf{1^*0(0+1)^*}$
$R_{21}^{(2)}$	$\emptyset + (\epsilon + \mathbf{0} + \mathbf{1})(\epsilon + \mathbf{0} + \mathbf{1})^*\emptyset$	$= \emptyset$
$R_{22}^{(2)}$	$\epsilon + \mathbf{0} + \mathbf{1} + (\epsilon + \mathbf{0} + \mathbf{1})(\epsilon + \mathbf{0} + \mathbf{1})^*(\epsilon + \mathbf{0} + \mathbf{1})$	$= (\mathbf{0} + \mathbf{1})^*$

State elimination algorithm

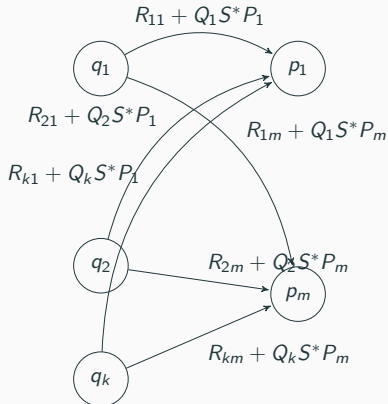
State elimination: the idea

Idea: Allow edges labelled by RE, iteratively remove nodes. (More efficient, avoids duplicity.)

State s selected for elimination



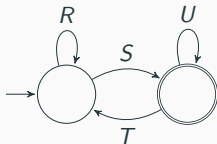
After s is eliminated.



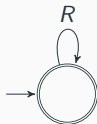
State elimination: the algorithm

For every accepting $q \in F$ eliminate all states $p \in Q \setminus \{q, q_0\}$.

- for $q \neq q_0$: $\text{RegE}(q) = (R + SU^*T)^*SU^*$



- for $q = q_0$: $\text{RegE}(q) = R^*$

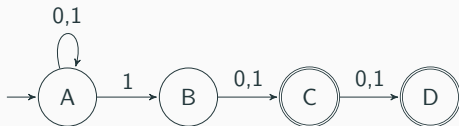


Finally, union over all accepting states: $\text{RegE}(A) = \sum_{q \in F} \text{RegE}(q)$

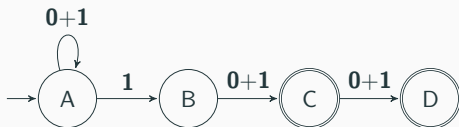
(Elimination order: first nonaccepting and noninitial states.)

State elimination: an example

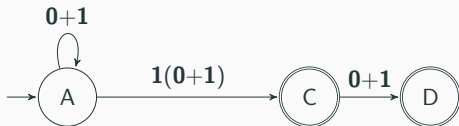
The original automaton:



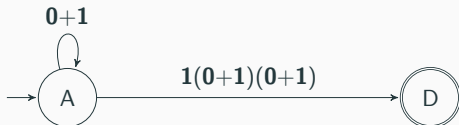
Replace letters by RE:



Eliminate B:



Eliminate C:



$$(0+1)^*1(0+1) + (0+1)^*1(0+1)(0+1)$$

Algebraic description of regular languages

Let $RL(\Sigma)$ denote the smallest set of languages over Σ that:

- contains \emptyset and $\{x\}$ for any letter $x \in \Sigma$, and
- is closed under union, concatenation, and iteration.

That is, for $A, B \in RL(\Sigma)$ also $A \cup B, A.B, A^* \in RL(\Sigma)$. Note that:

- $\{\epsilon\} \in RL(\Sigma)$ since $\{\epsilon\} = \emptyset^*$
- $\Sigma \in RL(\Sigma)$ since $\Sigma = \bigcup_{x \in \Sigma} \{x\}$ (a finite union)
- $\Sigma^* \in RL(\Sigma)$
- any finite language over Σ is in $RL(\Sigma)$.

Theorem (A restatement of Kleene's Theorem)

A language over Σ is regular, iff it is in $RL(\Sigma)$.

Some properties to simplify RE (will not be tested)

$$L.\emptyset = \emptyset.L = \emptyset$$

$$\{\epsilon\}.L = L.\{\epsilon\} = L$$

$$(L^*)^* = L^*$$

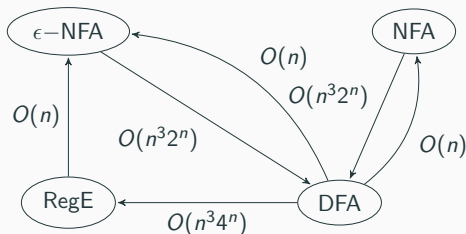
$$(L_1 \cup L_2)^* = L_1^*(L_2.L_1^*)^* = L_2^*(L_1.L_2^*)^*$$

$$(L_1.L_2)^R = L_2^R.L_1^R$$

$$\partial_w(L_1 \cup L_2) = \partial_w(L_1) \cup \partial_w(L_2)$$

$$\partial_w(\Sigma^* - L) = \Sigma^* - \partial_w L.$$

Converting between representations



- NFA to DFA
 - ϵ -closure in $O(n^3)$ (search n states $\times n^2$ arcs)
 - subset construction, DFA with up to 2^n states; for each state need $O(n^3)$ time to compute transitions.
- DFA to NFA or ϵ -NFA: simple modification of the table
- DFA to RE: $O(n^3 4^n)$
- RE to ϵ -NFA: $O(n)$

String substitution

This slide is under construction