

Regular expressions

3.1

are a formalism for describing a certain class of languages:
declarative rather than computational view

Definition: given an alphabet Σ , regular expressions are strings over the alphabet $\Sigma \cup \{+, *, (,), \cdot, \epsilon, \emptyset\}$ defined inductively as follows:

- basis: ϵ, \emptyset , and each $a \in \Sigma$ is a R.E.
- inductive step: if E and F are R.E., then so are:
 - $E + F$ (union)
 - $E \cdot F$ (concatenation)
 - E^* (closure)
 - (E) (parentheses)

Example: $a \cdot (a + b)^* \cdot b^* \cdot a$

Definition: language $\mathcal{L}(E)$ defined by a R.E. E

is also defined inductively:

- $\mathcal{L}(\epsilon) = \{\epsilon\}$ empty word
- $\mathcal{L}(\emptyset) = \emptyset$ empty language
- $\mathcal{L}(E + F) = \mathcal{L}(E) \cup \mathcal{L}(F)$
- $\mathcal{L}(E \cdot F) = \mathcal{L}(E) \cdot \mathcal{L}(F)$ concatenation
- for each $a \in \Sigma$, $\mathcal{L}(a) = \{a\}$
- $\mathcal{L}(E^*) = \mathcal{L}(E)^*$

concatenation of two languages L_1 and L_2 :

$$L_1 \cdot L_2 = \{w \mid w = x \cdot y, x \in L_1, y \in L_2\}$$

Example: $E = \epsilon + 1 \Rightarrow \mathcal{L}(E) = \{\epsilon, 1\}$

$$F = \epsilon + 0 + 1 \Rightarrow \mathcal{L}(F) = \{\epsilon, 0, 1\}$$

$$G = E \cdot F \Rightarrow \mathcal{L}(G) = \{\epsilon, 0, 1, 10, 11\} \\ = (\epsilon + 1) \cdot (\epsilon + 0 + 1)$$

• $\mathcal{L}(E^*) = (\mathcal{L}(E))^*$... closure

closure of a language L ?

we first define the powers of a language L :

- $L^0 = \{\epsilon\}$
- $L^k = L^{k-1} \cdot L$

hence $L^k = \{w \mid w = x_1 \dots x_k, \text{ with } \forall i, x_i \in L\}$

closure of L : $L^* = L^0 \cup L^1 \cup L^2 \cup \dots = \bigcup_{k \geq 0} L^k$

Example: $\begin{cases} E = 0+1 & \Rightarrow \mathcal{L}(E) = \{0, 1\} \\ F = E^* & \Rightarrow \mathcal{L}(F) = \text{set of all binary strings} \end{cases}$

$\begin{cases} E = 0 \cdot 0 & \Rightarrow \mathcal{L}(E^*) = \{\epsilon, 00, 0000, 000000, \dots\} \\ & = \text{all even-length strings of 0's} \end{cases}$

Positive closure of a language L

$L^+ = L^1 \cup L^2 \cup \dots$

We can introduce a positive closure operator in R.E.

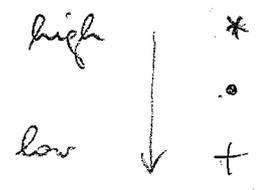
$\mathcal{L}(E^+) = (\mathcal{L}(E))^+$

Note: we have to distinguish between an expression E and the language $\mathcal{L}(E)$ defined by E

When we write $E = F$, we usually mean not syntactic equality but equality of the corresponding languages, i.e. $\mathcal{L}(E) = \mathcal{L}(F)$.

In other words, equality is in the algebra of R.E.

Precedence of operators:



example: $E + F \cdot G^* = E + (F \cdot (G^*))$

Algebraic laws for R.E.

similar to the laws for arithmetic expressions, we can express laws for R.E. : treat + as sum

• as product

(note! the laws are not the same as for arithmetic expressions)

- associativity of \cdot and $+$

$$(E \cdot F) \cdot G = E \cdot (F \cdot G) = E \cdot F \cdot G$$

$$(E + F) + G = E + (F + G) = E + F + G$$

- commutativity of $+$

$$E + F = F + E$$

Note: \cdot is not commutative: $E \cdot F \neq F \cdot E$

- distributivity:

1) left distributive law of \cdot over $+$: $E \cdot (F + G) = E \cdot F + E \cdot G$

2) right " " : $(F + G) \cdot E = F \cdot E + G \cdot E$

Proof of (1): the law actually holds for arbitrary languages, (OPTIONAL) and does not require E, F, G to be R.E.

Hence, we prove: for arbitrary languages L, M, N :

$$L \cdot (M \cup N) = L \cdot M \cup L \cdot N$$

We show that for a string w we have $w \in L \cdot (M \cup N)$ iff $w \in L \cdot M \cup L \cdot N$

"Only-if": $w \in L \cdot (M \cup N) \Rightarrow w = x \cdot y$ with $x \in L, y \in M \cup N$
Since $y \in M \cup N$, either $y \in M$ or $y \in N$ (or both).

If $y \in M$, then $w = x \cdot y \in L \cdot M$, hence $w \in L \cdot M \cup L \cdot N$
(similarly for $y \in N$)

"If": $w \in L \cdot M \cup L \cdot N$, hence either $w \in L \cdot M$ or $w \in L \cdot N$.

If $w \in L \cdot M$, then $w = x \cdot y$, with $x \in L, y \in M$. (Similarly for $w \in L \cdot N$)
Hence $y \in M \cup N$, and $w = x \cdot y \in L \cdot (M \cup N)$.

Exercise 3.1.1 Write R.E.s for the following languages

- a) $\{w \in \{a, b, c\}^* \mid w \text{ contains at least one } a \text{ and at least one } b\}$
- b) $\{w \in \{0, 1\}^* \mid w\text{'s tenth symbol from the right is a } 1\}$
- c) $\{w \in \{0, 1\}^* \mid w \text{ contains at most one pair of consecutive } 1\text{'s}\}$

Exercise 3.1.2 Write R.E.s for the following languages

- a) The set of all strings over $\{0, 1\}$ s.t. every pair of adjacent 0's appears before any pair of adjacent 1's
- b) The set of strings of 0's and 1's whose number of 0's is divisible by 5

Solutions:

$$3.1.1 \text{ a) } (c^* \cdot a \cdot (a+c)^* \cdot b \cdot (a+b+c)^* + c^* \cdot b \cdot (b+c)^* \cdot a \cdot (a+b+c)^*)^*$$

$$b) (0+1)^* \underbrace{1 \cdot (0+1) \cdot \dots \cdot (0+1)}_{3 \text{ times}}$$

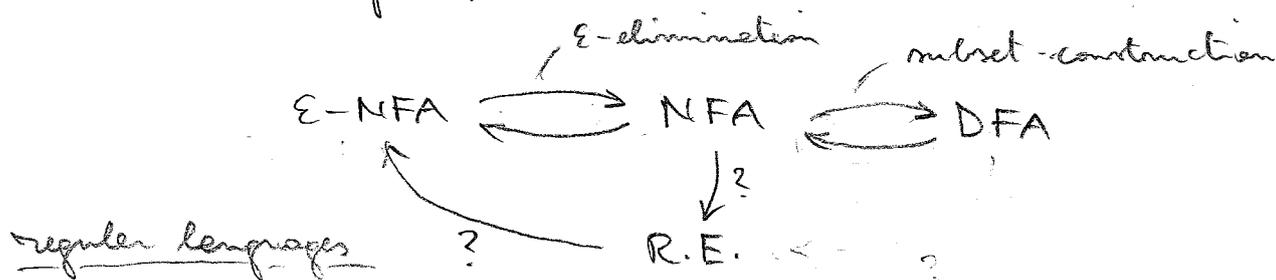
$$c) 0^* \cdot (1 \cdot 0^+)^* \cdot 1 \cdot 1 \cdot (0^+ \cdot 1) \cdot 0^* + 0^* \cdot (1 \cdot 0^+)^* \cdot (\epsilon+1)$$

or, simpler: $(0+10)^* (\epsilon+11) \cdot (0+01)^*$

$$3.1.2 \text{ a) } \underbrace{((0+10)^*)^*}_{\text{no pair of adjacent } 1\text{'s}} \cdot \underbrace{(1+10)^*}_{\text{no pair of adjacent } 0\text{'s}}$$

$$b) (1^* \cdot 0 \cdot 1^* \cdot 0 \cdot 1^* \cdot 0 \cdot 1^* \cdot 0 \cdot 1^* \cdot 0 \cdot 1^*)^*$$

What is the relationship between the classes of languages studied so far?



Theorem: (R.E. \rightarrow E-NFA)

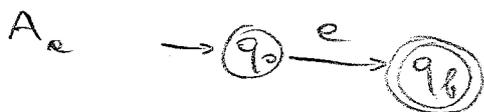
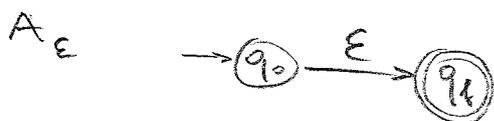
For every R.E. E there is an E-NFA A_E s.t. $L(A_E) = L(E)$.

Proof: let us call an E-NFA simple if

- it has only one final state
- the initial state has no incoming arcs
- the final state has no outgoing arcs

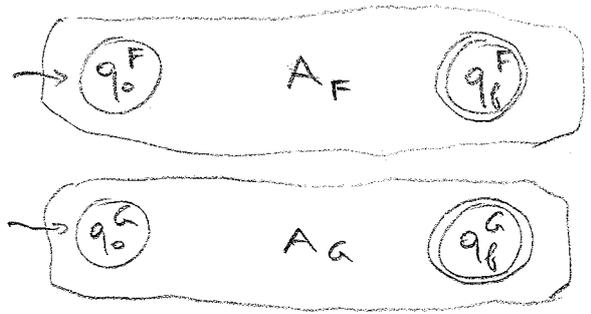
We show by structural induction that for each R.E. E there is a simple E-NFA A_E s.t. $L(E) = L(A_E)$

Basis: $E = \epsilon$, $E = \emptyset$, $E = a$ for some $a \in \Sigma$

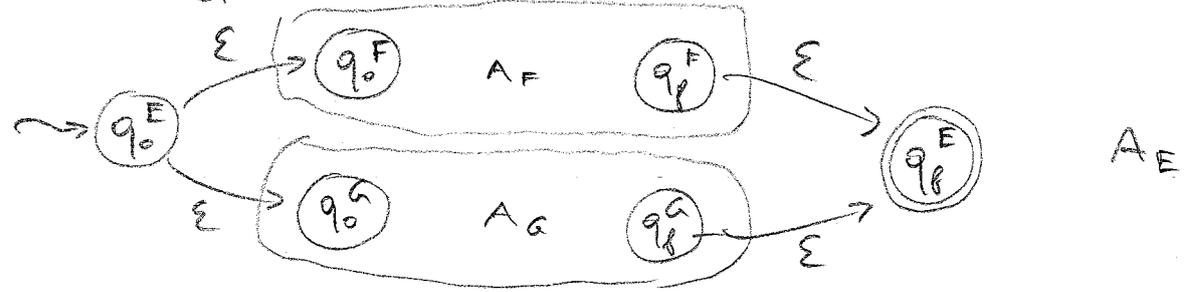


- Inductive case:
- 1) $E = F + G$
 - 2) $E = F \cdot G$
 - 3) $E = F^*$
 - 4) $E = (F)$

By I.H., there are simple ϵ -NFA's for F and G



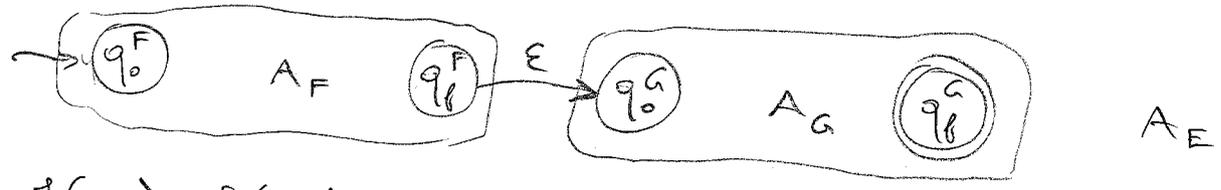
1) $E = F \cup G$



$$L(A_E) = L(A_F) \cup L(A_G) = L(F) \cup L(G) = L(F + G) = L(E)$$

by I.H.

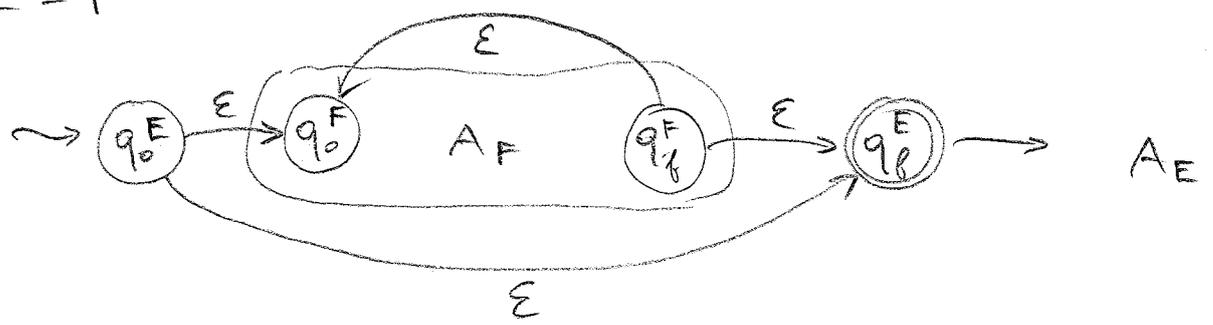
2) $E = F \cdot G$



$$L(A_E) = L(A_F) \cdot L(A_G) = L(F) \cdot L(G) = L(F \cdot G) = L(E)$$

by I.H.

3) $E = F^*$

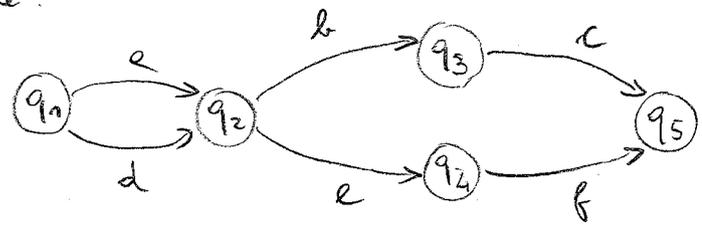


3) $E = (F)$

$$A_E = A_F$$

q.e.d.

Example:



$abc \in L_{15}^3$

$def \notin L_{15}^3$

but $def \in L_{15}^4$

$def \in L_{15}$

$L_{12}^0 = \{e, d\}$

$L_{15}^3 = \{abc, dbc\}$

$L_{15}^4 = L_{15}^5 = L_{15} = \{abc, dbc, def, def\}$

Note: $L_{ij}^n = L_{ij}$

Hence, we are done if we can construct R.E.s E_{ij}^k for L_{ij}^k .

We can simply take $E_{ij} = E_{ij}^n$, and hence $E_A = \sum_{q_j \in F} E_{1j}^n$.

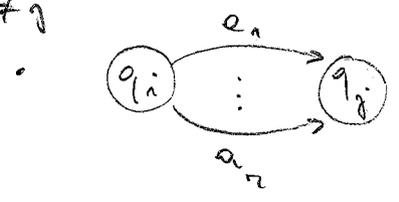
We construct E_{ij}^k by induction on k :

Basis: we construct E_{ij}^0 , for all $i, j \in \{1, \dots, n\}$

Since $k=0$, we cannot go through any intermediate state.

2 cases, each with 2 subcases:

$i \neq j$



$E_{ij}^0 = a_1 + \dots + a_n$



$E_{ij}^0 = \emptyset$

$i = j$



$E_{ii}^0 = \epsilon + a_1 + \dots + a_n$



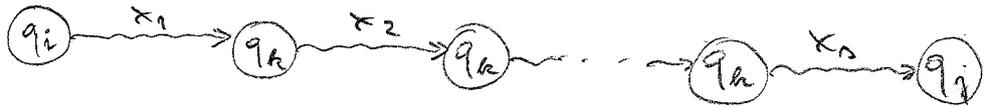
$E_{ii}^0 = \epsilon$

Induction: assume we have constructed $E_{ij}^{k-1} \quad \forall i, j \in \{1, \dots, n\}$ (3.10)

we show how to construct E_{ij}^k

Observe:

- L_{ij}^k will include L_{ij}^{k-1}
- it additionally will contain those words that lead through q_k at least once, when going from q_i to q_j



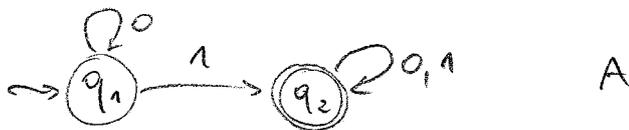
$w = x_1 x_2 \dots x_n$ where: \rightarrow represents transitions going at most through $\{x_1, \dots, x_{k-1}\}$

then $x_1 \in L_{ik}^{k-1}$
 $x_2, \dots, x_{n-1} \in L_{kk}^{k-1}$
 $x_n \in L_{kj}^{k-1}$

$\Rightarrow w \in L_{ik}^{k-1} \cdot (L_{kk}^{k-1})^* \cdot L_{kj}^{k-1}$

$\Rightarrow E_{ij}^k = E_{ij}^{k-1} + E_{ik}^{k-1} \cdot (E_{kk}^{k-1})^* \cdot E_{kj}^{k-1}$

Example:



k	E_{11}^k	E_{12}^k	E_{21}^k	E_{22}^k
0	$\epsilon + 0$	1	\emptyset	$\epsilon + 0 + 1$
1	$(\epsilon + 0) + (\epsilon + 0)(\epsilon + 0)^* = 0^*$	$1 + (\epsilon + 0) \cdot (\epsilon + 0)^* \cdot 1 = 0^* \cdot 1$	\emptyset	$\epsilon + 0 + 1$
2	not needed	$0^* \cdot 1 + (0^* \cdot 1) \cdot (\epsilon + 0 + 1)^* \cdot (\epsilon + 0 + 1) = 0^* \cdot 1 \cdot (0 + 1)^*$	not needed	not needed

$\mathcal{L}(A) =$
 $\mathcal{L}(E_{12}^2) =$
 $\mathcal{L}(E_{12}^2) =$
 $E_A^1 = E_{12}^2 = E_{12}^2 = 0^* \cdot 1 \cdot (0 + 1)^*$

↑ OPTIONAL

Theorem (NFA \rightarrow R.E.)

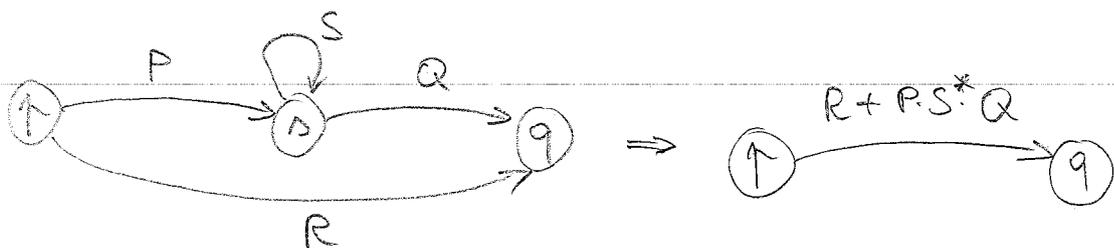
For every NFA A there is a R.E. E_A s.t. $\mathcal{L}(E_A) = \mathcal{L}(A)$

Proof sketch: We show how to construct E_A by eliminating states of A .

Consider the elimination of a state s :

- If there was a path from state p to state q over s ,
- after eliminating s the path does no longer exist
- \Rightarrow we have to compensate for that

We add a regular expression "connecting" p and q and capturing the missing path.



We can eliminate in this way all states except initial and final states:

Strategy: Assume A has k final states q_1, \dots, q_k .

For each final state q_i , eliminate all states except q_i, q_0

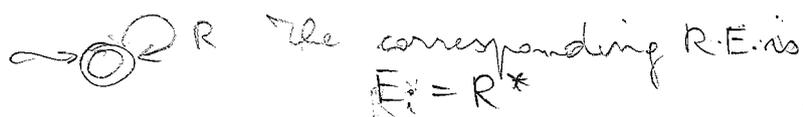
a) If $q_i \neq q_0$, we are left with:



The corresponding R.E. is $E_i = (R + S \cdot U^* \cdot T)^* \cdot S \cdot U^*$

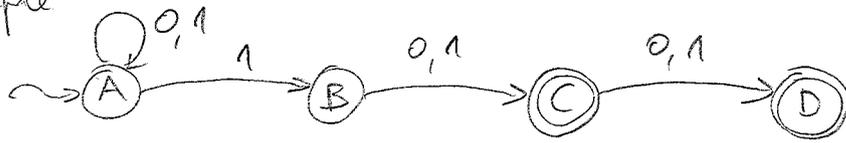
b) If $q_i = q_0 \in F$, we must eliminate all states except q_0 .

We are left with

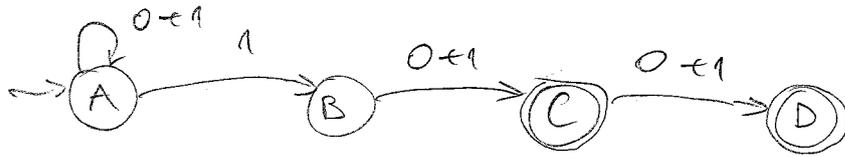


Then E_A is the sum of all E_i : $E_A = E_1 + E_2 + \dots + E_k$

Example



We view all edge labels as R.E.s (missing labels mean \emptyset)

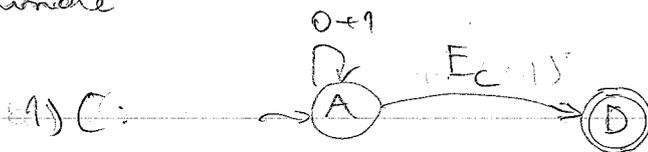


Eliminate B:



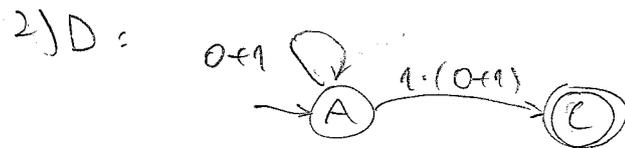
$$E_B = \emptyset + 1 \cdot \emptyset^* \cdot (0+1) = 1 \cdot \emptyset^* \cdot (0+1) = 1 \cdot (0+1)$$

Eliminate



$$E_C = \emptyset + 1 \cdot (0+1) \cdot \emptyset^* \cdot (0+1) = 1 \cdot (0+1) \cdot (0+1)$$

$$E_1 = (0+1)^* \cdot E_C = (0+1)^* \cdot 1 \cdot (0+1) \cdot (0+1)$$



$$E_2 = (0+1)^* \cdot 1 \cdot (0+1)$$

$$E = E_1 + E_2 = (0+1)^* \cdot 1 \cdot (0+1) \cdot (0+1) + (0+1)^* \cdot 1 \cdot (0+1)$$

$$= (0+1)^* \cdot 1 \cdot (0+1) (E + 0+1)$$