

CSC 355
2/4/2026

Regular Expressions (regexes) ①

Time permitting: DFA min example

Def. Let A & B be languages over Σ .

$$AB = \{xy : x \in A \text{ \& } y \in B\} \quad (\text{concat of } A \text{ \& } B)$$

$$A^* := \{\epsilon\} \cup A \cup AA \cup AAA \cup \dots \cup A^n \cup \dots$$

$$\begin{array}{cccc} \parallel & \parallel & \parallel & \parallel \\ A^0 & A^1 & A^2 & A^3 \end{array}$$

$$= \bigcup_{n \geq 0} A^n$$

Concat, $*$ -operator ("Kleene closure"), & union ($A \cup B$)

are the regular operators. A regular expression (regex)

is an expression that represents (denotes) a language

built from simple languages using regular operators.

Given a regex r , we let $L(r)$ be the lang denoted by r .

Define regex syntax & semantics recursively as follows

	regex r	$L(r)$
atomic regexes	\emptyset	\emptyset (empty lang)
	$\forall a \in \Sigma$ a	$\{a\}$

non-atomic regexes

for regexes s, t	
(s,t regexes)	s + t
	st
	s*

- $L(s) \cup L(t)$
- $L(s)L(t)$
- $L(s)^*$

Parens can be used for grouping to override precedence rules:

- + — lowest prec.
- concat — next higher
- * — highest

Ex: $\Sigma = \{a, b, c\}$

ab $L(ab) = L(a)L(b) = \{a\}\{b\} = \{ab\}$

abacca $L(abacca) = \dots = \{abaccaca\}$

$(a+b)(c+a) \quad L((a+b)(c+a)) = L(aa + ac + ba + bc) = \{aa, ac, ba, bc\}$

+ is commutative: $r + s \equiv s + r$
equivalent; denoting the same lang

assoc $(r + s) + t \equiv r + (s + t)$

concat distributes over +:
 $r(s+t) \equiv rs + rt$
 $(s+t)r \equiv sr + tr$

concat is not commutative (order matters), but (3)

" is associative (grouping does not matter)

$$\begin{aligned} a^* \quad L(a^*) &= \{a\}^* = \{\varepsilon\} \cup \{a\} \cup \{a\}\{a\} \cup \dots \\ &= \{\varepsilon, a, a^2, a^3, \dots, a^n, \dots\} \\ &= \{a^n : n \geq 0\} \end{aligned}$$

$$ab^* \quad L(ab^*) = \{a\}\{b\}^* = \dots = \{a, ab, abb, \dots, ab^n, \dots\}$$

$$(ab)^* \quad L((ab)^*) = \{ab\}^* = \{\varepsilon, ab, abab, ababab, \dots, (ab)^n, \dots\}$$

$$\begin{aligned} (a+b)^* \quad L((a+b)^*) &= \{\varepsilon\} \cup \{a, b\} \cup \{a, b\}\{a, b\} \cup \dots \\ &= \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots\} \\ &\quad \text{all strings over } \{a, b\} \end{aligned}$$

(Σ^* is just $*$ -operator on Σ)

$$a + b^* \quad L(a + b^*) = \{a, \varepsilon, b, b^2, b^3, \dots\}$$

Theorem: $L(r)$ is regular for every regex r .

Proof: by construction. Given regex r we construct a equivalent ϵ -NFA N such that $L(N) = L(r)$.

($\therefore \exists$ DFA equiv to r by previous results).

Construction is recursive on syntax of r :

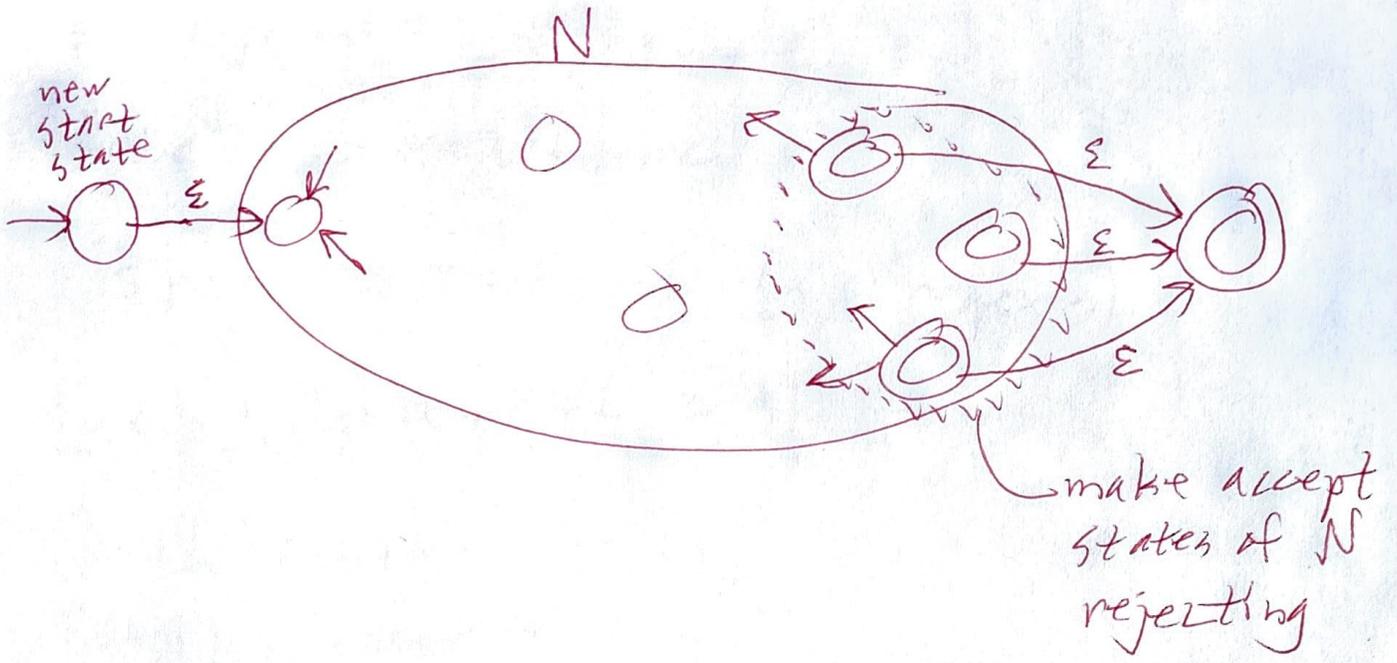
r	equivalent ϵ -NFA (clean)
\emptyset	
$(a \in \Sigma)$ a	
$(s, t \text{ rejects})$ $s + t$	
st	

Def: An ϵ -NFA is clean if (a) It has exactly one accepting state, (b) No transitions out of the accepting state; (c) no transitions into the start state.

start state
different from the

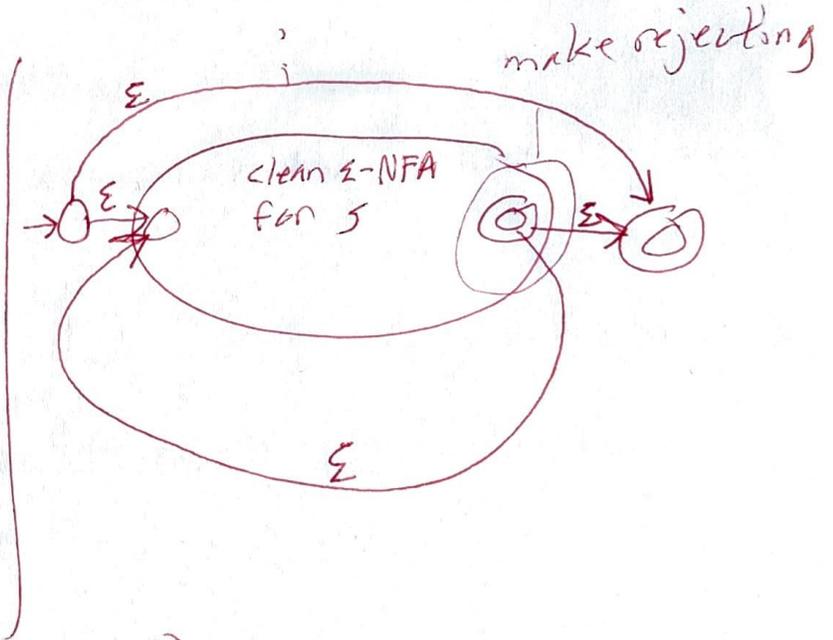
Prop: Every ϵ -NFA has an equivalent clean NFA.

Proof: Given an ϵ -NFA N :



Construction (continued)

⋮
S*



Regex shorthand (in popular use)

r regex

$$r^+ := r r^* \quad (\text{one or more } r\text{'s})$$

$$\epsilon := \emptyset^*$$

$$L(\epsilon) = L(\emptyset^*) = \{\epsilon\} \cup \emptyset \cup \emptyset\emptyset \cup \emptyset\emptyset\emptyset \cup \dots = \{\epsilon\}$$

$r?$ (optional r) $:= r + \epsilon$

"abc" $:=$ ~~abc~~ abc

$s|t$ $:= s + t$ (s, t regexes)

$[a < b]$ $:= a + c + b \equiv a|b|c$

$[a-z]$ $:= a + b + c + \dots + z$

$[0-9]$ — any digit

$[A-Za-z0-9]$ — any alphanumeric character

\cdot — any char (except $\backslash n$)

$[^ \dots]$ — any single char other than what is \dots

Unsigned integer constants $[0-9]^+$ (one or more digits)

Identifiers in C, C++, Java $[A-Za-z_][A-Za-z0-9_]^*$

Unsigned Floating point constants (in Pascal)

$[0-9]^+ \cdot [0-9]^+ ([Ee][+-]? [0-9]^+)?$

Let r be a regex and w a string.

" r matches w "
" w matches r " } just mean $w \in L(r)$