

CSCE 355

Regular Expressions

①

2/5/2025

Patterns that match strings

Fix an alphabet Σ .

Def. A regular expression (regex) over Σ is built from the following syntax:

Atomic: \emptyset (for $a \in \Sigma$)

\emptyset — does not match anything

a — matches the string "a" and nothing else

Nonatomic: Let s and t be regexes over Σ .

Then

(union) $s + t$ — string can match s or t (or both)

(concat.) $s \cdot t$ — matches concatenations of strings matching s followed by t .

(Kleene closure) s^*
(star operator)
are regexes

— matches the concatenation of any number (zero or more) of strings matching s .

$+$ & concat are binary ops (infix)

$*$ is unary postfix

precedences (lowest to highest): $+$, concat, $*$

②

Can use parens to control grouping.

Def: Every regex r over Σ^1 defines a language $L(r) \subseteq \Sigma^{1*}$. [The strings that match r , or ~~the~~ equiv, that r matches.]

Formally, $L(r)$ is defined inductively on the syntax of r :

$$L(\emptyset) = \emptyset$$

$$(a \in \Sigma^1) \quad L(a) = \{a\}$$

$$L(s+t) = L(s) \cup L(t) \quad s, t \text{ regexes}$$

$$L(st) = L(s)L(t)$$

$$[\text{Def: } L_1, L_2 \subseteq \Sigma^{1*} : L_1 L_2 := \{xy : x \in L_1, \& y \in L_2\}]$$

$$L(s^*) = \underbrace{\{\epsilon\}}_{L(s)^0} \cup \underbrace{L(s)}_{L(s)^1} \cup \underbrace{L(s)L(s)}_{L(s)^2} \cup \dots \cup \dots \cup L(s)^n \cup \dots$$

Ex: $\Sigma^1 = \{a, b\}$

$$a+bb = a+(bb)$$
$$L(a+bb) = L(a+(bb)) = L(a) \cup L(bb)$$

$$= \{a\} \cup \{b\}\{b\} = \{a\} \cup \{bb\}$$

$$= \boxed{\{a, bb\}}$$

$$L((a+b)b) = L(a+b)L(b) = (L(a) \cup L(b))L(b)$$

$$= (\{a\} \cup \{b\})\{b\} = \{a, b\}\{b\}$$

$$= \{ab, bb\} \neq L(a+bb)$$

$$L(a^*) = \{\epsilon, a, aa, aaa, \dots, a^n, \dots\}$$

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

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

$$L((a+b)^*) = \{\epsilon, a, b, aa, ab, ba, bb, \\ aaa, aab, aba, abb, baa, \dots\}$$

all possible strings over $\Sigma^1 = \{a, b\}$

$$\left[\begin{array}{l} \Sigma^1 \\ \text{"} \\ \{a, b\} \end{array} \right]^* = \text{set of all strings over } \Sigma^1$$

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

3

Let ε ~~be the~~ be the regex \emptyset^* (4)

Then $L(\varepsilon) = \{\varepsilon\}$. ("syntactic sugar")

More syntactic sugar:

Let r, s be regexes.

$r|s$ means $r \vee s$

r^+ means $rr^* = r(r^*)$

"one or more occurrences of r "

$r^?$ means $r \cup \varepsilon$

"optional r , or zero or one occurrence of r "

~~"aba"~~

"abc" means abc

"a+b" means "a" + "b" ~~AB~~

(+ treated as an alphabet symbol)

- (period) means Σ (matches any single char from the alphabet)

If $\Sigma = \{a, b, c\}$, then $L(\cdot) = \{a, b, c\}$

Character classes

$[abc] = a + b + c$, lang is $\{a, b, c\}$
 ↑ ↑
 square brackets

~~[abc]~~ $[abc] = [bca] = \dots$

$\Sigma = \text{ASCII char set}$

$[a-z] = [abcd \dots z]$ lowercase letters

$[A-Z] = [AB \dots Z]$ capital letters

$[\wedge a-z]$ — matches anything except $a \dots z$
 ↑
 complement

Regexes for common prog. lang. constructs:

— unsigned integer constants

$[0-9]^+$

— identifiers

$[A-Za-z_][A-Za-z0-9_]^*$

— optionally signed int constants

(6)

$[+-]?[0-9]^+$

— (unsigned float constants) a la Pascal:

one or more digits

then decimal point

then one or more digits

then option exponent following E or e

$[0-9]^+ \cdot ([0-9]^+)^{\text{exponent}} ([Ee][+-]?[0-9]^+)?$

Theorem: For every regex r over Σ ,
there is an ϵ -NFA N over Σ such that
 $L(r) = L(N)$ (r & N are equivalent)

Proof is by construction.

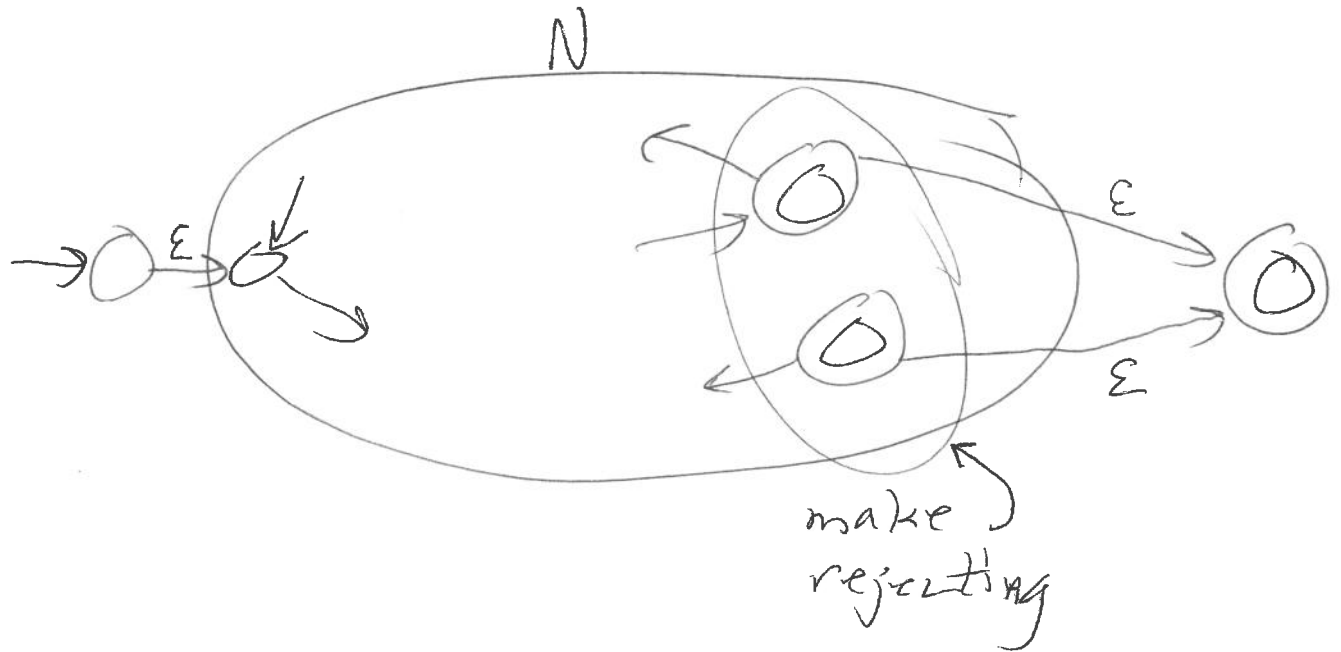
Def: An ϵ -NFA is clean if

— it has ~~at~~ exactly one accept state,
and this state is distinct from the
start state

- there are no transitions into the start state
- there are no transitions out of the accepting state.

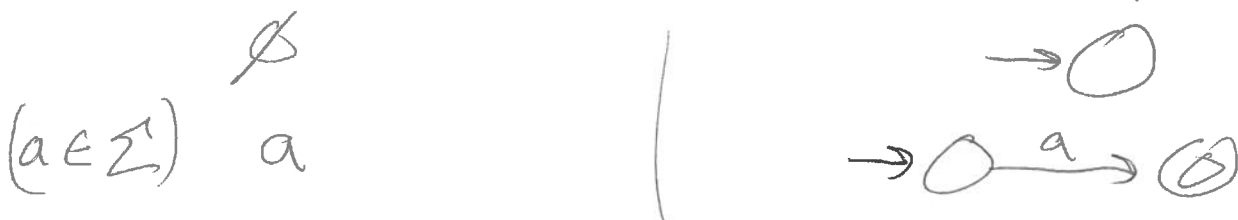
Prop: For every ϵ -NFA there is an equiv clean NFA.

Proof: Let N be any ϵ -NFA

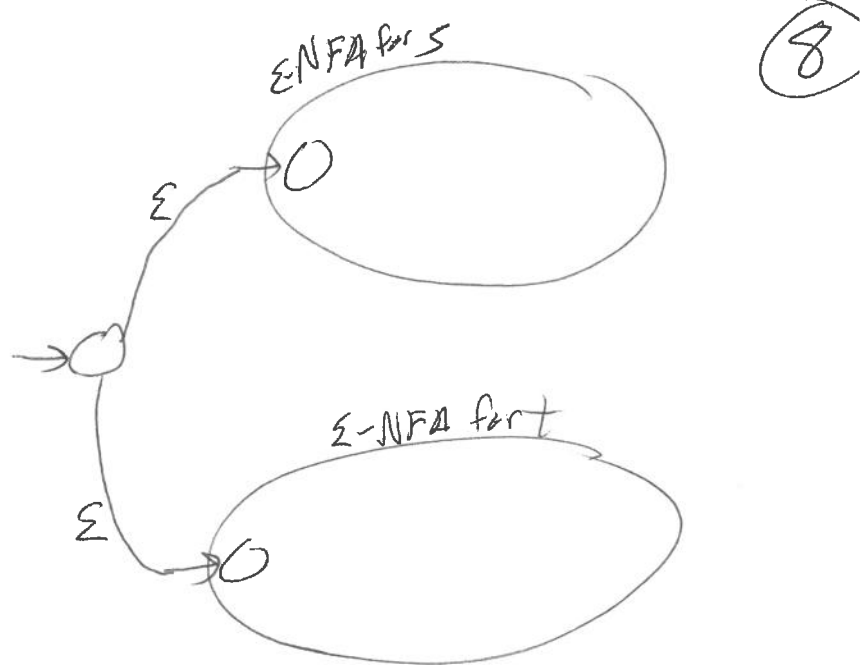


Correctness proof omitted

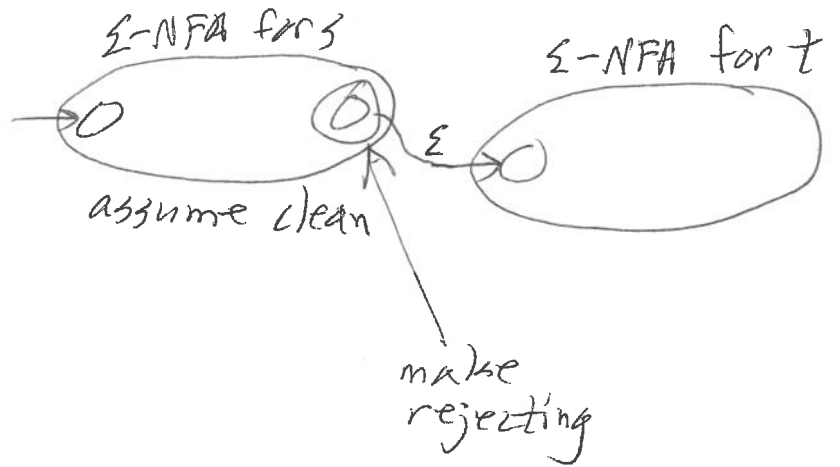
Proof of the thm (the construction) by induction on the regex syntax:



$s+t$



st



s^*

