

Regular Expressions

CS 2800: Discrete Structures, Fall 2014

Sid Chaudhuri

Regular expressions (“regex”-es) are
defined *inductively*

(start with simple base expressions, construct more
complicated ones recursively)

Empty set

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

Empty string

$$L(\boldsymbol{\varepsilon}) = \{\boldsymbol{\varepsilon}\}$$

Literal character

$$L(\mathbf{x}) = \{x\}$$

e.g. $L(\mathbf{1}) = \{1\}$, $L(\mathbf{2}) = \{2\}$, $L(\mathbf{a}) = \{a\}$

Concatenation

$$L(\mathbf{AB}) = \{ab \mid a \in L(A), b \in L(B)\}$$

e.g. $L(\mathbf{12}) = \{12\}$, $L(\mathbf{aabb}) = \{aabb\}$

Concatenation

$$L(AB) = \{ab \mid a \in L(A), b \in L(B)\}$$

e.g. $L(\mathbf{12}) = \{12\}$, $L(\mathbf{aabb}) = \{aabb\}$,
 $L(\mathbf{a}\epsilon) = ?$

Concatenation

$$L(\mathbf{AB}) = \{ab \mid a \in L(A), b \in L(B)\}$$

e.g. $L(\mathbf{12}) = \{12\}$, $L(\mathbf{aabb}) = \{aabb\}$,
 $L(\mathbf{a\epsilon}) = \{a\}$

Concatenation

$$L(\mathbf{AB}) = \{ab \mid a \in L(A), b \in L(B)\}$$

e.g. $L(\mathbf{12}) = \{12\}$, $L(\mathbf{aabb}) = \{aabb\}$,
 $L(\mathbf{a\varepsilon}) = \{a\}$, $L(\mathbf{a\emptyset}) = ?$

Concatenation

$$L(\mathbf{AB}) = \{ab \mid a \in L(A), b \in L(B)\}$$

e.g. $L(\mathbf{12}) = \{12\}$, $L(\mathbf{aabb}) = \{aabb\}$,
 $L(\mathbf{a\varepsilon}) = \{a\}$, $L(\mathbf{a\emptyset}) = \emptyset$

Alternation

$$L(A|B) = L(A) \cup L(B)$$

e.g. $L(\mathbf{1|2}) = \{1, 2\}$, $L(\mathbf{aa|bb}) = \{aa, bb\}$

Alternation

$$L(A|B) = L(A) \cup L(B)$$

e.g. $L(1|2) = \{1, 2\}$, $L(\mathbf{aa|bb}) = \{aa, bb\}$,
 $L(\mathbf{a|\epsilon}) = ?$

Alternation

$$L(A|B) = L(A) \cup L(B)$$

e.g. $L(1|2) = \{1, 2\}$, $L(\mathbf{aa|bb}) = \{aa, bb\}$,
 $L(\mathbf{a|\epsilon}) = \{a, \epsilon\}$

Alternation

$$L(A|B) = L(A) \cup L(B)$$

e.g. $L(1|2) = \{1, 2\}$, $L(\mathbf{aa}|\mathbf{bb}) = \{aa, bb\}$,
 $L(\mathbf{a}|\boldsymbol{\varepsilon}) = \{a, \varepsilon\}$, $L(\mathbf{a}|\emptyset) = ?$

Alternation

$$L(A|B) = L(A) \cup L(B)$$

e.g. $L(1|2) = \{1, 2\}$, $L(\mathbf{aa}|\mathbf{bb}) = \{aa, bb\}$,
 $L(\mathbf{a}|\boldsymbol{\varepsilon}) = \{a, \varepsilon\}$, $L(\mathbf{a}|\emptyset) = \{a\}$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(a^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,
 $L((\mathbf{ab})^*) = ?$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,
 $L((\mathbf{ab})^*) = \{\varepsilon, ab, abab, ababab, \dots\}$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = ?$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$,

$L(\boldsymbol{\varepsilon}^*) = ?$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$,

$L(\boldsymbol{\varepsilon}^*) = \{\varepsilon\}$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$,

$L(\boldsymbol{\varepsilon}^*) = \{\varepsilon\}$, $L(\emptyset^*) = ?$

Kleene star

$$L(A^*) = \{\varepsilon\} \cup \{x_1x_2\dots x_n \mid n \in \mathbf{N}, x_i \in L(A)\}$$

e.g. $L(\mathbf{a}^*) = \{\varepsilon, a, aa, aaa, aaaa, \dots\}$,

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

$L((\mathbf{a|b})^*) = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, aab, aba, \dots\}$,

$L(\boldsymbol{\varepsilon}^*) = \{\varepsilon\}$, $L(\emptyset^*) = \{\varepsilon\}$

Playing with regexes

- <http://regex101.com/>
- <http://rubular.com/>
- <http://www.google.com/search?q=online+regex+tester>

Kleene's Theorem

- A language is recognized by a regular expression (that is, it is a “regular language”) if and only if it is recognized by a finite automaton

Kleene's Theorem

- A language is recognized by a regular expression (that is, it is a “regular language”) if and only if it is recognized by a finite automaton
 - Regex has FA
 - Relatively simple construction

Kleene's Theorem

- A language is recognized by a regular expression (that is, it is a “regular language”) if and only if it is recognized by a finite automaton
 - Regex has FA
 - Relatively simple construction
 - FA has regex
 - Tricky to prove

Regex \rightarrow FA

- For every regular expression, there is a finite automaton that recognizes the same language

Regex \rightarrow FA

- For every regular expression, there is a finite automaton that recognizes the same language
- We will construct an ε -NFA

Regex \rightarrow FA

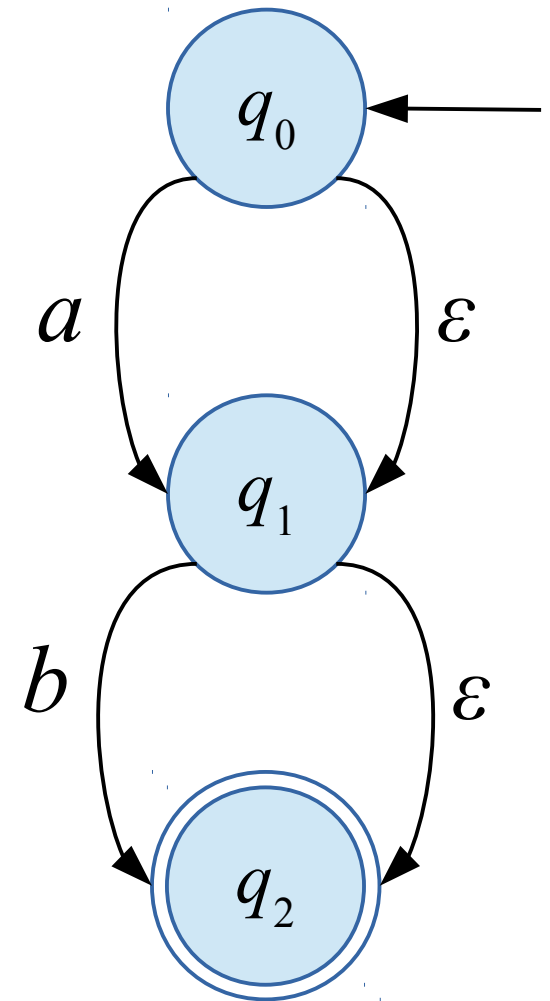
- For every regular expression, there is a finite automaton that recognizes the same language
- We will construct an ε -NFA
 - ... which can be converted to an NFA

Regex \rightarrow FA

- For every regular expression, there is a finite automaton that recognizes the same language
- We will construct an ε -NFA
 - ... which can be converted to an NFA
 - ... which can be converted to a DFA

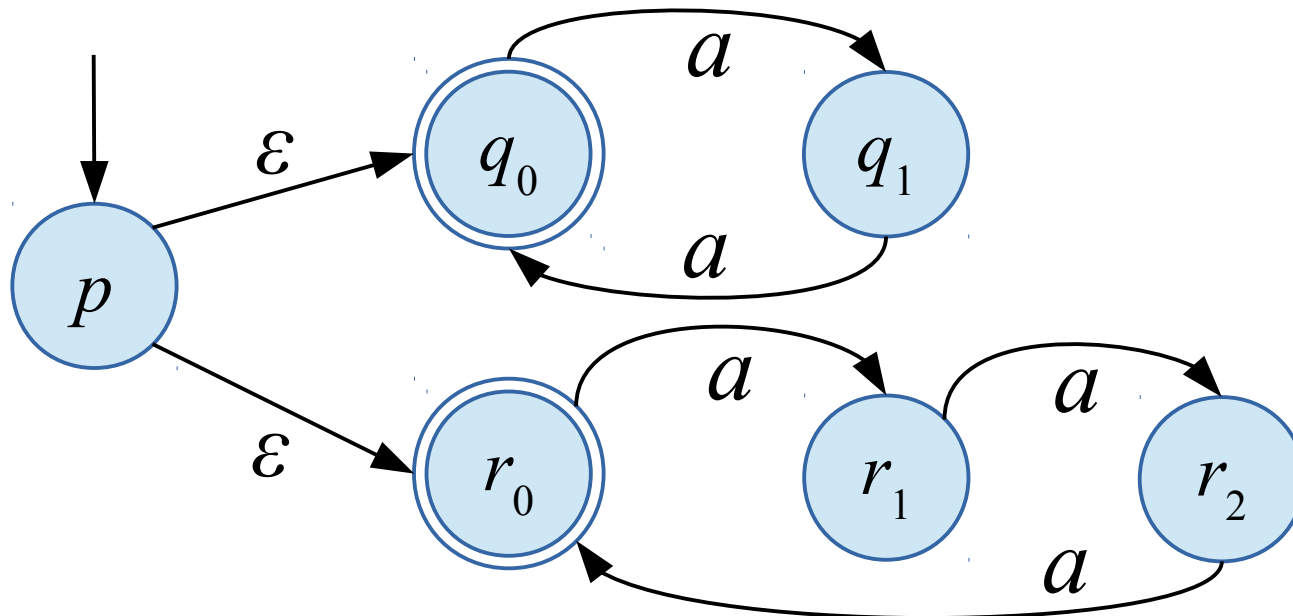
Recap: NFAs with epsilon transitions

- Just like ordinary NFAs, but...
 - Can “instantaneously” change state *without* reading an input symbol
 - Valid transitions of this type are shown by arcs labeled ' ϵ '
 - Note that ϵ does not suddenly become a member of the alphabet. Instead, we assume ϵ does not belong to *any* alphabet – it's a special symbol.



Why ε -NFAs?

- Suitable for representing “or” relations
- E.g. $L = \{ a^n \mid n \in \mathbf{N} \text{ is divisible by 2 or 3} \}$



- ... but they're equivalent to NFAs and DFAs