CS304: Automata and Formal Languages

Lec 8

Regular Languages, Regular Expressions, and Kleene's Theorem

Rachit Nimavat

August 14, 2025

Outline

- Recap
- Regular Expressions
- Regular Languages
- 4 Proving Kleene's Theorem

Alphabet Σ : A finite, non-empty set of symbols

String w: A finite sequence of symbols taken from an alphabet

Language L: A set of strings over a particular alphabet

Consider $\Sigma = \left\{0,1\right\}$ and $L = \left\{0\right\}^* \left\{1\right\} \equiv 0^*1.$

Deterministic Finite Automata (DFA)

Formal Definition

DFA is a 5-tuple $A = (Q, \Sigma, \delta, q_0, F)$

'symbol'	description
Q	finite set of states
Σ	underlying finite alphabet
δ	transition function between states
$q_0 \in Q$	start state
$F \subseteq Q$	accepting states

in our code

 $\begin{array}{l} \texttt{state} \in \{0,1,2\} \\ \Sigma = \{ \text{'0', '1'} \} \\ \texttt{the 'core logic' in C code} \\ \texttt{state} = 0 \\ \texttt{state} = 1 \end{array}$

Deterministic Finite Automata (DFA)

Visualizing DFA

State Diagram is the intuitive way we visualize and work with DFAs

Conventions:

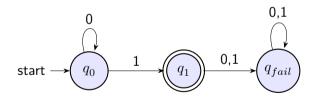
States Q are circles

Start state (q_0) has an incoming arrow labeled start

Accepting states (F) are double circles

Transitions (δ) are arrows between states, labeled with input symbols from Σ

Our DFA for $L=0^*1$:



Non-Deterministic Finite Automata (NFA)

Visualizing NFA

State Diagram is the intuitive way we visualize and work with NFAs

Conventions:

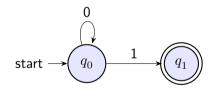
States Q are circles

Start state (q_0) has an incoming arrow labeled start

Accepting states (F) are double circles

Transitions (δ) are arrows between states, labeled with input symbols from Σ

Our NFA for $L=0^*1$:



NFA vs DFA

NFA

$$N = (Q, \Sigma, \delta, q_0, F)$$

$$\delta: Q \times (\Sigma \cup \{\varepsilon\}) \mapsto 2^Q$$



Multiple Transitions. Possibly multiple outgoing arrows for same symbol

Zero Transitions. Possibly no outgoing arrow for a symbol (that path "dies" has license to kill) ε -**Transitions.** Can change state without consuming an input symbol

DFA

$$D = (Q, \Sigma, \delta, q_0, F)$$

$$\delta:Q\times\Sigma\mapsto Q$$

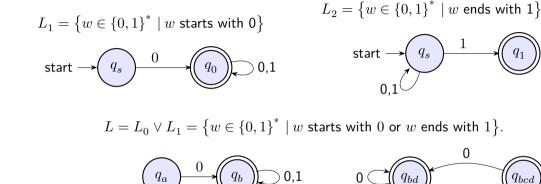
No such powers!



Theorem 1.1

A language is recognized by an NFA if and only if it is recognized by a DFA.

Example



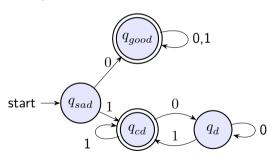
start $\longrightarrow (q_{sad})$

 q_d

 q_{cd}

Example - Simplified DFA

 $L = L_0 \vee L_1 = \big\{w \in \{0,1\}^* \mid w \text{ starts with } 0 \text{ or } w \text{ ends with } 1\big\}.$



What Next?

- we've learned so far to describe languages using machines: DFAs and NFAs
- they are good at recognizing strings, but are cumbersome to write
- consider find function in a text editor
 - we provide it a pattern and ask to find all matching strings in the text we don't give it a DFA
 - that 'pattern-language' is built on Regular Expressions (regex or RE)

Regular Expressions (RE)

Recursive Defn

A regular expression (RE) over an alphabet Σ is defined as follows:

Primitive Regular Expressions (Base case):

```
\emptyset : \text{ is an RE denoting empty language } L(\emptyset) = \{\}
```

$$\varepsilon$$
: is a RE the language containing only the empty string $L(\varepsilon) = \{\varepsilon\}$ a : is an RE for a symbol $a \in \Sigma$ denoting $L(a) = \{a\}$

Union (+ or |): (R+S) is an RE, where $L(R+S)=L(R)\cup L(S)$

Operations (Inductive Step): If R and S are RE, then the following are also REs:

```
e.g. a+b describes the language \{a,b\}
Concat (·): (RS) is an RE, where L(RS) = L(R)L(S)
e.g. (a+b)c describes the language \{ac,bc\}
```

Kleene Star (*):
$$(R^*)$$
 is an RE, where $L(R^*) = (L(R))^*$

Regular Expressions

Examples

RE 0*1 is any number of 0's followed by a 1

Language described by $0^*1: \{1, 01, 001, 0001, \dots, \}$

RE $(a+b)^*$ is any string composed of a's and b's

Language described by $(a+b)^* \colon \left\{ \varepsilon, a, b, ab, ba, aab, aba, baa, \dots, \right\}$

(0+1)*00(0+1)* is any binary string with 00 as substring

RE $(a + \varepsilon)b$ is an optional a followed by b

Language described by $(a + \varepsilon)b$: $\{b, ab\}$

Regular Expressions

Operator Precedence

To reduce parentheses, there is a standard order of operations. The operators are evaluated in the following order, from highest to lowest precedence:

- 1 Kleene Star * think: power
- 2 Concat · think: multiplication
- 3 Union + or | think: addition

$$\mathsf{RE}\ a + bc^* \ \mathsf{is} \ \mathsf{parsed} \ \mathsf{as:} \ a + (b(c^*)).$$

Language described by $a + bc^*$: $\{a, b, bc, bcc, bccc, \dots, \}$

RE
$$(a+b)c^*$$
 is parsed as: $(a+b)(c^*)$.

Language described by $(a + b)c^*$: $\{a, b, ac, bc, acc, bcc, \dots, \}$

Regular Language

Defn: A language L is called a <u>regular language</u> if and only if there exists a regular expression R that describes it: L=L(R)

But wait, didn't we already define regular languages earlier? A language L is regular iff 'if and only if':

- \exists an NFA A such that L = L(A)
- or equivalently, \exists a DFA D such that L=L(D)

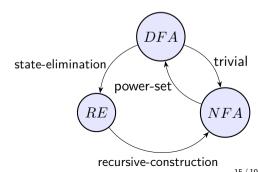
Regular Language



Theorem 3.1 (Kleene's Theorem)

All these definitions are equivalent:

- \exists a DFA D such that L = L(D)
- \exists an NFA A such that L = L(A)
- \exists an RE R such that L = L(R)



Kleene's Theorem: RE ⇒ NFA

Step-by-step algorithm to convert an RE to an NFA $\,$

Primitive Regular Expressions (Base case):

- a (where $a \in \Sigma$)
- 6
- (

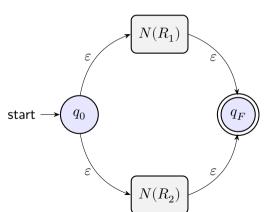
RE	NFA
a (where $a \in \Sigma$)	start $\longrightarrow q_0$ $\xrightarrow{a} q_1$
a (where $a \in \Delta$)	
ϵ	$start \longrightarrow q_0$
Ø	$start \longrightarrow q_0$

Kleene's Theorem: RE \implies NFA

Step-by-step algorithm to convert an RE to an NFA

Combining Machines (Inductive step):

- $R_1 + R_2$ (union)

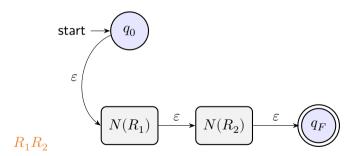


Kleene's Theorem: RE ⇒ NFA

Step-by-step algorithm to convert an RE to an NFA

Combining Machines (Inductive step):

- R_1R_2 (concatenation)



Kleene's Theorem: RE ⇒ NFA

Step-by-step algorithm to convert an RE to an NFA

Combining Machines (Inductive step):

- R_1^* (Kleene star)

