

CMPE 322/327 - Theory of Computation

Week 4: Pattern Matching & Regular Expressions

Burak Ekici

March 14-18, 2022

Outline

- 1

A Quick Recap
- 2

Pattern Matching
- 3

Regular Expressions
- 4

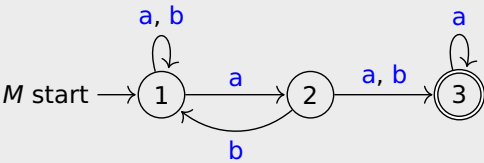
Homomorphisms

Definitions

- nondeterministic finite automaton (NFA) is quintuple $N = (Q, \Sigma, \Delta, s, F)$ with
 - Q : finite set of states
 - Σ : input alphabet
 - $\Delta : Q \times \Sigma \rightarrow 2^Q$: transition function
 - $S \subseteq Q$: start state
 - $F \subseteq Q$: final (accept) states

Example

$N = (Q, \Sigma, \Delta, S, F)$



- $Q := \{1, 2, 3\}$
- $\Sigma := \{a, b\}$
- $\Delta : Q \times \Sigma \rightarrow 2^Q$
- $S := \{1\}$
- $F := \{3\}$

Δ	a	b
1	{1, 2}	{1}
2	{3}	{1, 3}
3	{3}	\emptyset

Definitions

- nondeterministic finite automaton (NFA)** is quintuple $N = (Q, \Sigma, \Delta, s, F)$ with
 - 1 Q : finite set of states
 - 2 Σ : input alphabet
 - 3 $\Delta: Q \times \Sigma \rightarrow 2^Q$: transition function
 - 4 $S \subseteq Q$: set of start states
 - 5 $F \subseteq Q$: final (accept) states
- $\widehat{\Delta}: 2^Q \times \Sigma^* \rightarrow 2^Q$ is inductively defined by

$$\widehat{\Delta}(A, \varepsilon) = A \qquad \widehat{\Delta}(A, xa) = \bigcup_{q \in \widehat{\Delta}(A, x)} \Delta(q, a)$$
- string $x \in \Sigma^*$ is **accepted** by N if $\widehat{\Delta}(S, x) \cap F \neq \emptyset$

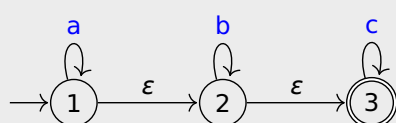
5/32

Definitions

- NFA with ε -transitions (NFA_ε) is sextuple $N = (Q, \Sigma, \varepsilon, \Delta, S, F)$ such that
 - 1 $\varepsilon \notin \Sigma$
 - 2 $N_\varepsilon = (Q, \Sigma \cup \{\varepsilon\}, \Delta, S, F)$ is NFA over alphabet $\Sigma \cup \{\varepsilon\}$
- ε -closure** of set $A \subseteq Q$ is defined as $C_\varepsilon(A) = \bigcup \{\widehat{\Delta}_{N_\varepsilon}(A, x) \mid x \in \{\varepsilon\}^*\}$
- $\widehat{\Delta}_N: 2^Q \times \Sigma^* \rightarrow 2^Q$ is inductively defined by

$$\widehat{\Delta}_N(A, \varepsilon) = C_\varepsilon(A) \qquad \widehat{\Delta}_N(A, xa) = \bigcup \{C_\varepsilon(\Delta(q, a)) \mid q \in \widehat{\Delta}_N(A, x)\}$$

Example



$$C_\varepsilon(\{1\}) = \{1, 2, 3\}$$

$$\widehat{\Delta}(\{1\}, b) = \{2, 3\}$$

6/32

- Theorem

every set accepted by NFA is regular
- Theorem

every set accepted by NFA_ϵ is regular
- Theorem

regular sets are effectively closed under concatenation
- Theorem

regular sets are effectively closed under asterate

Outline

- 1 A Quick Recap
- 2 Pattern Matching
- 3 Regular Expressions
- 4 Homomorphisms

Pattern matching is important for

- lexical analysis of programs
 - scripting languages (Perl, Ruby)
- search engines (Google Code Search)
 - DNA analysis

Applications of Regular expressions: grep

- grep foo file returns lines in file containing pattern foo
- basis for more powerful tools like awk, sed, perl

Some Patterns

^	matches beginning of line	.	matches any character
\$	matches end of line	[abc]	matches a or b or c
c	matches character c	[a-zA-Z]	matches any letter

Example

grep "0" file returns lines containing 0

grep "0\$" file returns lines ending with 0

grep "b.g" file returns lines containing e.g. bag, big, bug, buggy

Pattern matching is important for

• lexical analysis of programs

• scripting languages (Perl, Ruby)

• search engines (Google Code Search)

• DNA analysis

Definitions

• pattern is string α that represents set of strings $L(\alpha) \subseteq \Sigma^*$

atomic pattern α	$L(\alpha)$	compound pattern α	$L(\alpha)$
$a \in \Sigma$	$\{a\}$	$\beta + \gamma$	$L(\beta) \cup L(\gamma)$
ϵ	$\{\epsilon\}$	$\beta \cap \gamma$	$L(\beta) \cap L(\gamma)$
• \emptyset	\emptyset	$\beta \gamma$	$L(\beta)L(\gamma)$
$\#$	Σ	β^*	$L(\beta)^*$
$@$	Σ^*	β^+	$L(\beta)^+$
		$\sim \beta$	$\sim L(\beta) = \Sigma^* - L(\beta)$

• string $x \in \Sigma^*$ matches pattern α if $x \in L(\alpha)$

Example

pattern	matched string
$@a@a@a@$	strings containing at least 3 occurrences of a
$@a@b@$	strings containing a followed later by b
$\#n \sim a$	single letters except a
$(\#n \sim a)^*$	strings without a

Questions

- how difficult is pattern matching?
- is pattern equivalence ($L(\alpha) = L(\beta)$) decidable?
- which operators are **redundant**?

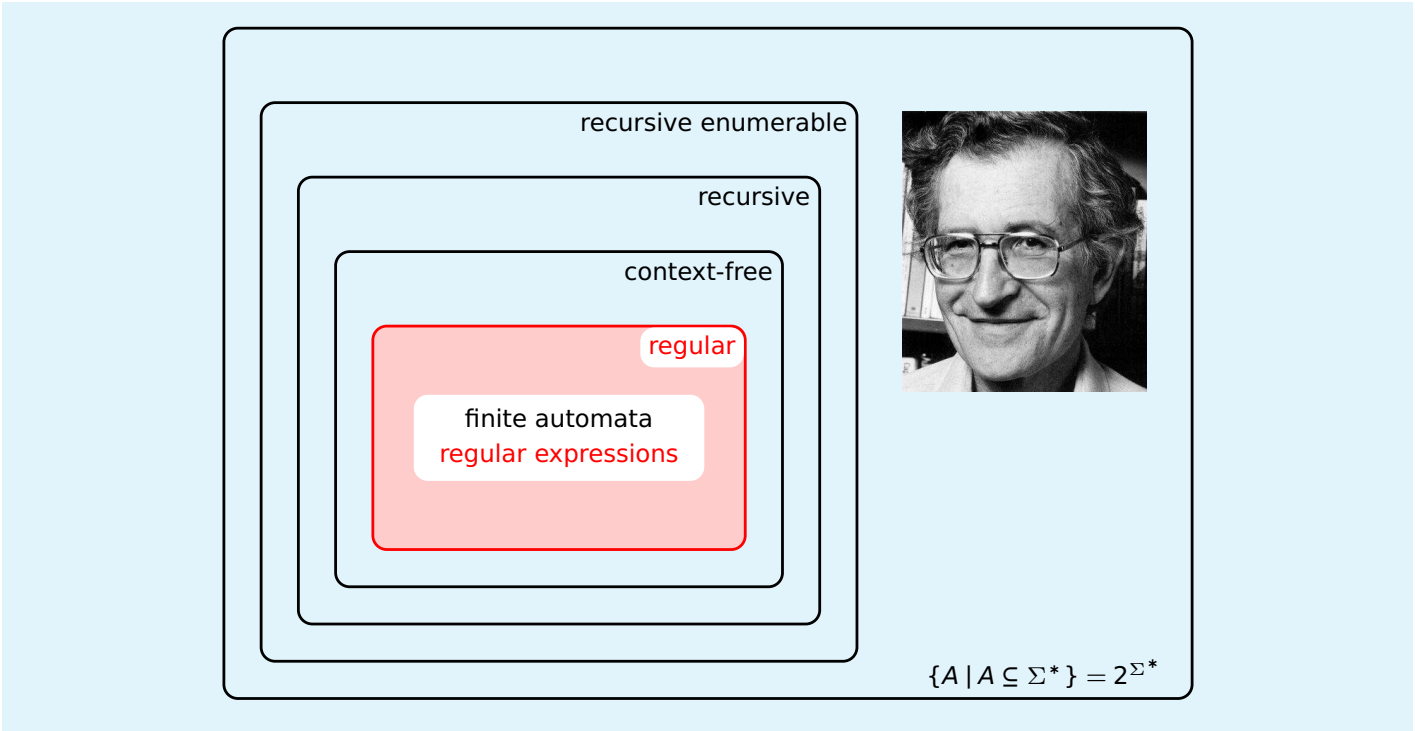
ϵ	\equiv	$\sim (\# @) \equiv \emptyset^*$	
$@$	\equiv	$\#^*$	
α^+	\equiv	$\alpha\alpha^*$	
$\#$	\equiv	$a_1 \dots a_n$	if $\Sigma = \{a_1 \dots a_n\}$
$\alpha \cap \beta$	\equiv	$\sim (\sim \alpha + \sim \beta)$	
$\sim \alpha$	\equiv	?	

Notation

$\alpha \equiv \beta \quad \text{if } L(\alpha) = L(\beta)$

Outline

- 1 A Quick Recap
- 2 Pattern Matching
- 3 Regular Expressions
- 4 Homomorphisms



Definition

regular expressions are restricted patterns which use only

$a \in \Sigma \quad \epsilon \quad \emptyset \quad \alpha + \beta \quad \alpha^* \quad \alpha\beta$

Theorem

finite automata, patterns, and regular expressions are **equivalent**:

for all $A \subseteq \Sigma^*$

① A is regular

\iff ② $A = L(\alpha)$ for some pattern α

\iff ③ $A = L(\alpha)$ for some regular expression α

Proof.

③ \implies ②

② \implies ①

① \implies ③

trivial (every regular expression is a pattern)

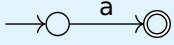
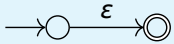
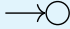
induction on α (see slides #17 – 18)

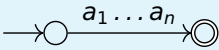

(see slide #19)

Proof. (2 ⇒ 1)

statement: for any pattern α , $L(\alpha)$ is regular
induction on pattern α

1 atomic pattern $\Sigma = \{a_1, \dots, a_n\}$

α	$L(\alpha)$	finite automaton
$a \in \Sigma$	$\{a\}$	
ϵ	$\{\epsilon\}$	
\emptyset	\emptyset	

α	$L(\alpha)$	finite automaton
$\#$	Σ	
$@$	Σ^*	

Proof. (2 ⇒ 1)

statement: for any pattern α , $L(\alpha)$ is regular
induction on pattern α

2 compound patterns

α	$L(\alpha)$
$\beta + \gamma$	$L(\beta) \cup L(\gamma)$
$\beta \cap \gamma$	$L(\beta) \cap L(\gamma)$
$\beta \gamma$	$L(\beta)L(\gamma)$

α	$L(\alpha)$
β^*	$L(\beta)^*$
β^+	$L(\beta)^+$
$\sim \beta$	$\sim L(\beta)$

$L(\beta)$ and $L(\gamma)$ are regular according to induction hypothesis
hence $L(\alpha)$ is regular according to closure properties of regular sets

Proof. (1 \Rightarrow 3 – An idea)

given $\text{NFA}_\varepsilon N_\varepsilon = (Q, \Sigma, \varepsilon, \Delta, S, F)$

$\forall Y \subseteq Q \quad \forall u, v \in Q$ construct regular expression α_{uv}^Y such that

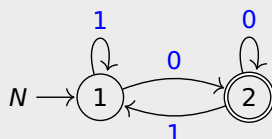
$$x \in L(\alpha_{uv}^Y) \iff \boxed{\exists \text{ a path from } u \text{ to } v \text{ labeled } x \ (v \in \widehat{\Delta}(\{u\}, x)) \text{ such that all intermediate states belong to } Y}$$

Definitions

- $\alpha_{uv}^\emptyset := \begin{cases} a_1 + \dots + a_k & \text{if } u \neq v \text{ and } k > 0 \\ \emptyset & \text{if } u \neq v \text{ and } k = 0 \\ a_1 + \dots + a_k + \varepsilon & \text{if } u = v \text{ and } k > 0 \\ \varepsilon & \text{if } u = v \text{ and } k = 0 \end{cases} \quad \{a_1, \dots, a_k\} := \{a \in \Sigma \cup \{\varepsilon\} \mid v \in \Delta(u, a)\}$
- $\alpha_{uv}^Y := \alpha_{uv}^{Y-\{q\}} + \alpha_{uq}^{Y-\{q\}} (\alpha_{qq}^{Y-\{q\}})^* \alpha_{qv}^{Y-\{q\}}$ for some fixed $q \in Y$

19/32

Example



$L(N) = L(\alpha)$ with

$$\alpha = \alpha_{12}^{\{1,2\}} = \alpha_{12}^{\{1\}} + \alpha_{12}^{\{1\}} (\alpha_{22}^{\{1\}})^* \alpha_{22}^{\{1\}} \quad (q = 2)$$

$$\alpha_{12}^{\{1\}} = \alpha_{12}^\emptyset + \alpha_{11}^\emptyset (\alpha_{11}^\emptyset)^* \alpha_{12}^\emptyset = 0 + (1 + \varepsilon)(1 + \varepsilon)^* 0$$

$$\alpha_{22}^{\{1\}} = \alpha_{22}^\emptyset + \alpha_{21}^\emptyset (\alpha_{11}^\emptyset)^* \alpha_{22}^\emptyset = (0 + \varepsilon) + 1(1 + \varepsilon)^* 0$$

$$\alpha_{12}^\emptyset = 0 \quad \alpha_{11}^\emptyset = 1 + \varepsilon \quad \alpha_{22}^\emptyset = 0 + \varepsilon \quad \alpha_{21}^\emptyset = 1$$

$$\begin{aligned} \alpha &= (0 + (1 + \varepsilon)(1 + \varepsilon)^* 0) + (0 + (1 + \varepsilon)(1 + \varepsilon)^* 0)((0 + \varepsilon) + 1(1 + \varepsilon)^* 0)^*((0 + \varepsilon) + 1(1 + \varepsilon)^* 0) \\ &\equiv (0 + 1)^* 0 \end{aligned}$$

20/32

Outline

- 1 A Quick Recap
- 2 Pattern Matching
- 3 Regular Expressions
- 4 Homomorphisms

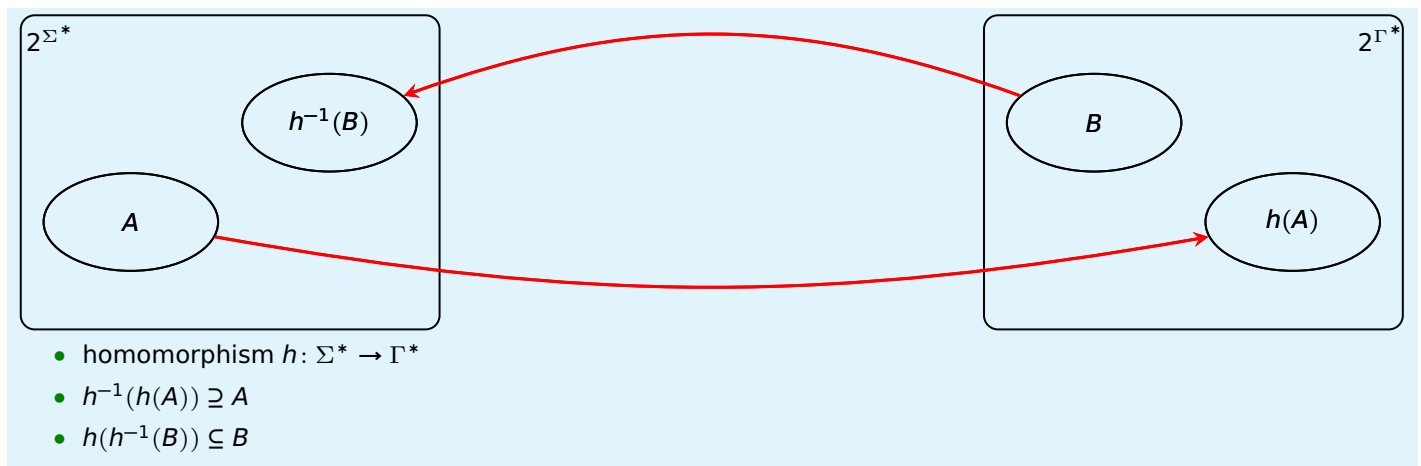
Theorem

regular sets are effectively closed under homomorphic image and preimage

Definitions

- homomorphism is mapping $h: \Sigma^* \rightarrow \Gamma^*$ such that
$$h(\varepsilon) = \varepsilon \qquad h(xy) = h(x)h(y)$$
so homomorphism is completely determined by its effect on Σ

if $A \subseteq \Sigma^*$ then	$h(A)$	$=$	$\{h(x) \mid x \in A\} \subseteq \Gamma^*$	“image of A under h ”
if $B \subseteq \Gamma^*$ then	$h^{-1}(B)$	$=$	$\{x \mid h(x) \in B\} \subseteq \Sigma^*$	“preimage of B under h ”



Example

$\Sigma = \Gamma = \{0, 1\}$ $h(0) = 11$ $h(1) = 1$ $A = \{0\}$ $\Sigma = \Gamma = \{0, 1\}$
 $h(0) = 11$ $h(1) = 1$ $A = B = \{0\}$

- $h^{-1}(h(A)) = h^{-1}(\{11\}) = \{0, 11\} \supset A$
- $h(h^{-1}(B)) = h(\emptyset) = \emptyset \subset B$

23/32

Lemma

$A \subseteq \{0, 1\}^*$ is regular $\implies \{xy \mid x1y \in A\}$ is regular

Proof.

- $\Sigma = \{0, 1\}$ and $\Gamma = \{0, 1, 2\}$
- define homomorphisms $h, i: \Gamma^* \rightarrow \Sigma^*$ by

$$h(0) = 0 \quad h(1) = h(2) = 1 \quad i(0) = 0 \quad i(1) = 1 \quad i(2) = \varepsilon$$
- $h^{-1}(A) = \{x \mid h(x) \in A\}$
- $h^{-1}(A) \cap L((0+1)^*2(0+1)^*) = \{x2y \mid x1y \in A\}$
- $\{xy \mid x1y \in A\} = i(h^{-1}(A) \cap L((0+1)^*2(0+1)^*))$ is regular

24/32

Theorem

regular sets are effectively closed under homomorphic image and **preimage**

Proof.

- DFA $M = (Q, \Gamma, \delta, s, F)$
- homomorphism $h: \Sigma^* \rightarrow \Gamma^*$
- $h^{-1}(L(M)) = L(M')$ for DFA $M' = (Q, \Sigma, \delta', s, F)$ with $\delta'(q, a) := \widehat{\delta}(q, h(a))$
- claim: $\widehat{\delta}'(q, x) = \widehat{\delta}(q, h(x)) \quad \forall x \in \Sigma^* \quad \forall q \in Q$
- proof of claim: induction on $|x|$ (see next slide)

25/32

proof of the claim

claim: $\widehat{\delta}'(q, x) = \widehat{\delta}(q, h(x)) \quad \forall x \in \Sigma^* \quad \forall q \in Q$

- base case: $|x| = 0$ thus $x = \varepsilon$

$$\widehat{\delta}'(q, \varepsilon) = q = \widehat{\delta}(q, h(\varepsilon))$$

- step case: $|x| > 0$ thus $x = ya$ s.t. $|y| = |x| - 1$ with IH: $\widehat{\delta}'(q, y) = \widehat{\delta}(q, h(y))$

$$\begin{aligned} \widehat{\delta}'(q, ya) &= \delta'(\widehat{\delta}'(q, y), a) && \text{(by definition of } \widehat{\delta}') \\ &= \delta'(\widehat{\delta}(q, h(y)), a) && \text{(by induction hypothesis IH)} \\ &= \widehat{\delta}(\widehat{\delta}(q, h(y)), h(a)) && \text{(by definition of } \delta') \\ &= \widehat{\delta}(q, h(y)h(a)) && \text{(by distributivity of } \widehat{\delta} - \text{w3.pdf, slide 10)} \\ &= \widehat{\delta}(q, h(ya)) && \text{(by definition of homomorphism)} \\ &= \widehat{\delta}(q, h(x)) \end{aligned}$$

□

26/32

Proof. (closedness under complement homomorphic preimage)

statement: $L(M') = h^{-1}(L(M))$

$$\begin{aligned} \forall x \in \Sigma^*, x \in L(M') &\iff \widehat{\delta'}(s, x) \in F && \text{(by definition of acceptance)} \\ &\iff \widehat{\delta}(s, h(x)) \in F && \text{(by claim proven in slide 26)} \\ &\iff h(x) \in L(M) && \text{(by definition of acceptance)} \\ &\iff x \in h^{-1}(L(M)) && \text{(by definition of homomorphic preimage)} \end{aligned}$$

□

Example

This page has too many overlays. Please refer to the original slides (w4.pdf) to monitor the whole content.

Theorem

regular sets are effectively closed under **homomorphic image** and preimage

Proof.

- regular expression α over Σ
- homomorphism $h: \Sigma^* \rightarrow \Gamma^*$
- $h(L(\alpha)) = L(\alpha')$ for regular expression α' defined inductively:

$$\begin{aligned} \mathbf{a}' &= h(\mathbf{a}) && \text{for } \mathbf{a} \in \Sigma \\ \epsilon' &= \epsilon \\ \emptyset' &= \emptyset \end{aligned}$$

$$\begin{aligned} (\beta + \gamma)' &= \beta' + \gamma' \\ (\beta\gamma)' &= \beta'\gamma' \\ (\beta^*)' &= (\beta')^* \end{aligned}$$

Definitions

- Hamming distance** $H(x, y)$ is number of places where bit strings x and y differ (if $|x| \neq |y|$ then $H(x, y) = \infty$)
- $N_k(A) := \{x \in \{0, 1\}^* \mid H(x, y) \leq k \text{ for some } y \in A\}$

Lemma

$A \subseteq \{0, 1\}^*$ is regular $\implies \forall k \in \mathbb{N}, N_k(A)$ is regular

Proof.

$$\begin{aligned} D_k &= \{x \in (\{0, 1\} \times \{0, 1\})^* \mid x \text{ contains at most } k \text{ pairs } (0, 1) \text{ or } (1, 0)\} && \text{is regular} \\ &= \{x \in (\{0, 1\} \times \{0, 1\})^* \mid H(\text{fst}(x), \text{snd}(x)) \leq k\} \\ N_k(A) &= \text{fst}(\text{snd}^{-1}(A) \cap D_k) \end{aligned}$$

Example

This page has too many overlays. Please refer to the original slides (w4.pdf) to monitor the whole content.

Thanks! & Questions?