CM0081 Formal Languages and Automata § 3.1 Regular Expressions

Andrés Sicard-Ramírez

Universidad EAFIT

Semester 2024-1

Preliminaries

Conventions

- The number and page numbers assigned to chapters, examples, exercises, figures, quotes, sections and theorems on these slides correspond to the numbers assigned in the textbook [Hopcroft, Motwani and Ullman (1979) 2007].
- The natural numbers include the zero, that is, $\mathbb{N} = \{0, 1, 2, ... \}$.

The power set of a set A, that is, the set of its subsets, is denoted by $\mathcal{P}A$.

Introduction: Description of regular languages

 $(\mathbf{01})(\mathbf{01})^* + (\mathbf{010})(\mathbf{010})^*$



Algebraic description

Machine-like description

Features

Algebraic description of regular languages

Features

Algebraic description of regular languages

Declarative ('user-friendly') way to express the strings that belong to the language

Features

- Algebraic description of regular languages
- Declarative ('user-friendly') way to express the strings that belong to the language

Uses

Search commands (e.g. GREP)

Features

- Algebraic description of regular languages
- Declarative ('user-friendly') way to express the strings that belong to the language

Uses

- Search commands (e.g. GREP)
- Lexical-analyzer generators (e.g. LEX and ALEX)

Features

- Algebraic description of regular languages
- Declarative ('user-friendly') way to express the strings that belong to the language

Uses

- Search commands (e.g. GREP)
- Lexical-analyzer generators (e.g. LEX and ALEX)
- Domain specific languages (DSLs)

Notation

The power set of a set A is denoted $\mathcal{P}A$.

Definition

Let L, L_1 and L_2 be languages on an alphabet Σ .

(i) Union of languages:

$$\label{eq:posterior} \begin{split} & \cup: \mathcal{P}\Sigma^* \times \mathcal{P}\Sigma^* \to \mathcal{P}\Sigma^* \\ L_1 \cup L_2 \coloneqq \{ \, x \mid x \in L_1 \text{ or } x \in L_2 \, \}. \end{split}$$

(iii) **Powers** of a language:

 $(-)^{(-)} : \mathcal{P}\Sigma^* \times \mathbb{N} \to \mathcal{P}\Sigma^*$ $L^0 := \{\varepsilon\},$ $L^{n+1} := L \cdot L^n.$

(ii) Concatenation of languages:

 $\begin{array}{l} \cdot:\mathcal{P}\Sigma^*\times\mathcal{P}\Sigma^*\to\mathcal{P}\Sigma^* \\ L_1\cdot L_2:=\{\,x\cdot y\mid x\in L_1 \text{ and } y\in L_2\,\}. \\ L^*:=\bigcup L^n. \end{array} \tag{iv} \begin{array}{l} \text{Kleene closure of a language:} \\ (-)^*:\mathcal{P}\Sigma^*\to\mathcal{P}\Sigma^* \\ L^*:=\bigcup L^n. \end{array}$

Operations on Languages

Examples

▶ If $L = \{0, 1\}$, then L^* consists of all strings of 0's and 1's and the empty word.

Examples

▶ If $L = \{0, 1\}$, then L^* consists of all strings of 0's and 1's and the empty word.

$$\blacktriangleright \text{ If } L = \{ 0^n \mid n \ge 1 \}, \text{ then } L^* = L \cup \{ \varepsilon \}.$$

Examples

- If $L = \{0, 1\}$, then L^* consists of all strings of 0's and 1's and the empty word.
- If $L = \{ 0^n \mid n \ge 1 \}$, then $L^* = L \cup \{ \varepsilon \}$.

• If $L = \{0, 11\}$, then L^* consists of the empty word and those strings of 0's and 1's such that the 1's come in pairs.

Examples

- If $L = \{0, 1\}$, then L^* consists of all strings of 0's and 1's and the empty word.
- If $L = \{ 0^n \mid n \ge 1 \}$, then $L^* = L \cup \{ \varepsilon \}$.
- If $L = \{0, 11\}$, then L^* consists of the empty word and those strings of 0's and 1's such that the 1's come in pairs.
- 🕨 Powers on 🖉

$$\begin{split} & \emptyset^0 = \{\varepsilon\}, \\ & \emptyset^i = \emptyset, \quad \text{for } i \geq 1, \\ & \emptyset^* = \{\varepsilon\}. \end{split}$$

Let Σ be an alphabet. The **regular expressions** (regex's) on Σ are inductively defined by:

Let Σ be an alphabet. The **regular expressions** (regex's) on Σ are inductively defined by:



(i) ε is a regex, (ii) \emptyset is regex and (iii) If $a \in \Sigma$ then a is a regex.

Let Σ be an alphabet. The **regular expressions** (regex's) on Σ are inductively defined by:

Basis step

(i) ε is a regex,
(ii) Ø is regex and

(iii) If $a \in \Sigma$ then a is a regex.

Inductive step If E and F are regex's then (i) E + F is a regex, (ii) $E \cdot F$ is a regex, (iii) E^* is a regex and (iv) (E) is a regex.

Precedence of Operators

Order of precedence and associative

Precedence from highest to lowest: (), *, \cdot and +.

Associative: The operators \cdot and + are left-associative.

Precedence of Operators

Order of precedence and associative

Precedence from highest to lowest: (), *, \cdot and +.

Associative: The operators \cdot and + are left-associative.

Example

 $\begin{aligned} \mathbf{01}^* + \mathbf{1} &= (\mathbf{0}(\mathbf{1}^*)) + \mathbf{1} \\ &\neq (\mathbf{01})^* + \mathbf{1} \\ &\neq \mathbf{0}(\mathbf{1}^* + \mathbf{1}) \end{aligned}$

Definition

Let E be a regular expression. The **language denoted** by E, denoted by L(E), is inductively defined by:

Definition

Let E be a regular expression. The **language denoted** by E, denoted by L(E), is inductively defined by:



```
\begin{split} \mathbf{L}(\varepsilon) &\coloneqq \{\varepsilon\},\\ \mathbf{L}(\emptyset) &\coloneqq \emptyset,\\ \mathbf{L}(\pmb{a}) &\coloneqq \{a\}. \end{split}
```

Let E be a regular expression. The **language denoted** by E, denoted by L(E), is inductively defined by:



```
\begin{split} \mathbf{L}(\varepsilon) &:= \{\varepsilon\},\\ \mathbf{L}(\emptyset) &:= \emptyset,\\ \mathbf{L}(\pmb{a}) &:= \{a\}. \end{split}
```

Inductive step

Let ${\rm L}(E)$ and ${\rm L}(F)$ be the languages denoted by the regular expressions E and F, then

$$\begin{split} \mathbf{L}(E+F) &\coloneqq \mathbf{L}(E) \cup \mathbf{L}(F), \\ \mathbf{L}(E \cdot F) &\coloneqq \mathbf{L}(E) \cdot \mathbf{L}(F), \\ \mathbf{L}(E^*) &\coloneqq (\mathbf{L}(E))^*, \\ \mathbf{L}((E)) &\coloneqq \mathbf{L}(E). \end{split}$$

E	L(E)
a + b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$

E	L(E)
a + b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$
a^*	$\{arepsilon, a, aa, aaa,\}$

E	L(E)
a + b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$
a^*	$\{\varepsilon, a, aa, aaa, \ldots\}$
$(\boldsymbol{a} + \boldsymbol{b})(\boldsymbol{a} + \boldsymbol{b})$	$\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})\cdot\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})=\{a,b\}\cdot\{a,b\}=\{aa,ab,ba,bb\}$

E	L(E)
a + b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$
a^*	$\{\varepsilon, a, aa, aaa, \ldots\}$
$(\boldsymbol{a} + \boldsymbol{b})(\boldsymbol{a} + \boldsymbol{b})$	$\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})\cdot\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})=\{a,b\}\cdot\{a,b\}=\{aa,ab,ba,bb\}$
$oldsymbol{a} + (oldsymbol{a}oldsymbol{b})^*$	$\{a,\varepsilon,ab,abab,ababab,\ldots\}$

E	$\mathrm{L}(E)$
a+b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$
a^*	$\{\varepsilon, a, aa, aaa, \ldots\}$
$(\boldsymbol{a} + \boldsymbol{b})(\boldsymbol{a} + \boldsymbol{b})$	$\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})\cdot\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})=\{a,b\}\cdot\{a,b\}=\{aa,ab,ba,bb\}$
$oldsymbol{a} + (oldsymbol{a}oldsymbol{b})^*$	$\{a,\varepsilon,ab,abab,ababab,\ldots\}$
$(0+1)^*01(0+1)^*$	$\{ x01y \mid x, y \in \{0,1\}^* \}$

E	L(E)
a+b	$\mathcal{L}(\boldsymbol{a}) \cup \mathcal{L}(\boldsymbol{b}) = \{a\} \cup \{b\} = \{a, b\}$
a^*	$\{\varepsilon, a, aa, aaa, \ldots\}$
$(\boldsymbol{a} + \boldsymbol{b})(\boldsymbol{a} + \boldsymbol{b})$	$\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})\cdot\mathcal{L}(\boldsymbol{a}+\boldsymbol{b})=\{a,b\}\cdot\{a,b\}=\{aa,ab,ba,bb\}$
$oldsymbol{a} + (oldsymbol{a}oldsymbol{b})^*$	$\{a,\varepsilon,ab,abab,ababab,\ldots\}$
$(0+1)^*01(0+1)^*$	$\{x01y \mid x, y \in \{0,1\}^*\}$
$oldsymbol{a}_{oldsymbol{i}}(oldsymbol{a}_1+oldsymbol{a}_2+\cdots+oldsymbol{a}_{oldsymbol{n}})^*$	$\{w\in\Sigma^*\mid w \text{ starts by } a_i\}$

Example

Write a regular expression for the language L defined by

```
L = \{ w \in \{0,1\}^* \mid 0 \text{ and } 1 \text{ alternate in } w \}.
```

Example

Write a regular expression for the language L defined by

```
L = \{ w \in \{0, 1\}^* \mid 0 \text{ and } 1 \text{ alternate in } w \}.
```

Solution.

 $(\mathbf{01})^* + (\mathbf{10})^* + \mathbf{0}(\mathbf{10})^* + \mathbf{1}(\mathbf{01})^*$

Example

Write a regular expression for the language L defined by

```
L = \{ w \in \{0, 1\}^* \mid 0 \text{ and } 1 \text{ alternate in } w \}.
```

Solution.

 $(\mathbf{01})^* + (\mathbf{10})^* + \mathbf{0}(\mathbf{10})^* + \mathbf{1}(\mathbf{01})^*$

Other solution.

 $(\varepsilon+\mathbf{1})(\mathbf{01})^*(\varepsilon+\mathbf{0})$

Example

The regular expression

$(\mathbf{10}+\mathbf{0})^*(\varepsilon+\mathbf{1})$

denotes the set of strings of 0's and 1's that have no two adjacent 1's.

Example

Write a regular expression for denoting the set of strings over $\Sigma = \{0, 1\}$ not ending in 01.

Example

Write a regular expression for denoting the set of strings over $\Sigma = \{0, 1\}$ not ending in 01.

Solution.

 $\varepsilon + \mathbf{0} + \mathbf{1} + (\mathbf{0} + \mathbf{1})^* (\mathbf{00} + \mathbf{10} + \mathbf{11})$

Derivatives of Regular Expressions

Observation

The material on derivatives of regular expressions is from [Brzozowski 1964].

Observation

The material on derivatives of regular expressions is from [Brzozowski 1964].

Definition

Let $L \subseteq \Sigma^*$ be a language and $a \in \Sigma$ a symbol. We define the **derivative** of L by a, denoted by $\partial_a L$, by

 $\begin{array}{l} \partial_a:\mathcal{P}\Sigma^*\to\mathcal{P}\Sigma^*\\ \partial_aL=\{\,x\in\Sigma^*\mid ax\in L\,\}. \end{array}$

Observation

The material on derivatives of regular expressions is from [Brzozowski 1964].

Definition

Let $L \subseteq \Sigma^*$ be a language and $a \in \Sigma$ a symbol. We define the **derivative** of L by a, denoted by $\partial_a L$, by

 $\begin{array}{l} \partial_a:\mathcal{P}\Sigma^*\to\mathcal{P}\Sigma^*\\ \partial_aL=\{\,x\in\Sigma^*\mid ax\in L\,\}. \end{array}$

Example

$$\begin{split} \partial_a \{abab, abba\} &= \{bab, bba\}, \\ \partial_a \mathcal{L}(\boldsymbol{ab^*}) &= \mathcal{L}(\boldsymbol{b^*}), \\ \partial_b \mathcal{L}(\boldsymbol{ab^*}) &= \emptyset. \end{split}$$

Let E be a regular expression on Σ and let $a \in \Sigma$ be a symbol. We define recursively the **derivative** of E by a, denoted $\partial_a E$, by

 $\partial_a: \mathrm{RegEx} \to \mathrm{RegEx}$

$$\begin{array}{ll} \partial_a \emptyset = \emptyset, & \partial_a (E+F) = \partial_a E + \partial_a F, \\ \partial_a \varepsilon = \emptyset, & \partial_a (EF) = \begin{cases} (\partial_a E)F + \partial_a F, & \text{if } \varepsilon \in \mathcal{L}(E), \\ (\partial_a E)F, & \text{otherwise}, \end{cases} \\ \partial_a (E^*) = (\partial_a E)E^*. \end{array}$$

Let E be a regular expression on Σ and let $w \in \Sigma^*$ be a string. We define recursively the **derivative** of E by w, denoted $\partial_w E$, by

$$\begin{split} \partial_w: \mathrm{RegEx} &\to \mathrm{RegEx} \\ \partial_\varepsilon E = E, \\ \partial_{ax} E &= \partial_a (\partial_x E). \end{split}$$

Derivatives of Regular Expressions

Theorem (Brzozowski [1964], Theorem 4.2)

Let E be a regular expression on Σ and let $w\in \Sigma^*$ be a string. Then

 $w\in \mathcal{L}(E) \qquad \Leftrightarrow \qquad \varepsilon\in \mathcal{L}(\partial_w E).$

Libraries

Observation

Theoretical regular expressions \neq practical regular expressions.

Libraries

Observation

Theoretical regular expressions \neq practical regular expressions.

Some programming languages with support to regular expressions .NET, C, HASKELL, JAVA, MATHEMATICA, MATLAB and PERL.

Algorithms

Algorithms

See the ${\rm HASKELL}$ implementation of some algorithms on regular expressions in the course homepage.

Applications

Some programs that use regular expressions

GREP: Print lines matching a pattern AWK: Pattern scanning and processing language SED: Stream editor for filtering and transforming text ALEX, FLEX and LEX: Lexical-analyser generators EMACS and VIM: Test editors MYSQL and ORACLE: Databases

Applications

Reading

§ 3.3. Applications of Regular Expressions.

Applications

Reading

§ 3.3. Applications of Regular Expressions.

In the above section are defined:

 $E^+ := EE^*$ (one or many times operator) $E? := \varepsilon + E$ (zero or one time operator)



'Rob's implementation itself is a superb example of beautiful code: compact, elegant, efficient, and useful. It's one of the best examples of recursion that I have ever seen.'

Brian Kernighan, p. 3.

References

Brzozowski, J. A. (1964). Derivates of Regular Expressions. Journal of the ACM 11.4, pp. 481–494. DOI: 10.1145/321239.321249 (cit. on pp. 35–37, 40).

