

# Chapter 1

## Regular Languages

### 1.1 Languages

**Definition 1.1.** An **alphabet** is a finite nonempty set of symbols.

**Definition 1.2.** Let  $\Sigma$  be an alphabet.

- A **string** over  $\Sigma$  is a finite sequence of symbols from  $\Sigma$ . The collection of all strings over  $\Sigma$  is denoted by  $\Sigma^*$ .
- The **length** of a string  $w$ , denoted by  $|w|$ , is the number of symbols it contains.
- The string containing no symbols is called the **empty string**, denoted by  $\epsilon$ .

**Definition 1.3.** A subset of  $\Sigma^*$  is called a **language** over  $\Sigma$ .

## 1.2 Deterministic Finite State Automata

**Definition 1.4.** A **deterministic finite state automaton** (DFA) is a tuple

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

where each component is as follows.

- $Q$  is a finite set of **states**.
- $\Sigma$  is a finite set of input symbols.
- $\delta : Q \times \Sigma \rightarrow Q$  is a function, called the **transition function**.
- $q_0 \in Q$  is called the **start state**.
- $F \subseteq Q$  is called the **accepting states**.

**Definition 1.5.** Let  $A = (Q, \Sigma, \delta, q_0, F)$  be a DFA.

- (a) For each symbol  $a \in \Sigma$ , we define  $\delta_a : Q \rightarrow Q$  to be the function such that  $\delta_a(p) = \delta(p, a)$  for any states  $p, q \in Q$ .
- (b) For each string  $w \in \Sigma^*$ , we define  $\delta_w : Q \rightarrow Q$  as follows.
  - $\delta_\epsilon$  is the identity function.
  - For any strings  $x \in \Sigma^*$  and any symbol  $a \in \Sigma$ , the function  $\delta_{xa}$  satisfies  $\delta_{xa}(p) = \delta_a(\delta_x(p))$  for any  $p \in Q$ .

**Definition 1.6.** Let  $A = (Q, \Sigma, \delta, q_0, F)$  be a DFA.

- We say that  $A$  accepts a string  $w \in \Sigma^*$  if  $\delta_w(q_0) \in F$ .
- The **language** of  $A$ , denoted  $L(A)$ , is defined as the set of strings that are accepted by  $A$ .

**Definition 1.7.** A language  $L$  is **regular** if there exists a DFA  $A$  such that  $L(A) = L$ .

## 1.3 Nondeterministic Finite State Automata

**Definition 1.8.** A **nondeterministic finite state automaton** (NFA) is a tuple

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

where each component is as follows.

- $Q$  is a finite set of **states**.
- $\Sigma$  is a finite set of input symbols.
- $\delta : Q \times \Sigma \times Q$  is a relation, called the **transition relation**.
- $q_0 \in Q$  is called the **start state**.
- $F \subseteq Q$  is called the **accepting states**.

**Definition 1.9.** Let  $A = (Q, \Sigma, \delta, q_0, F)$  be an NFA.

- (a) For each symbol  $a \in \Sigma$ , we define  $\delta_a \subseteq Q \times Q$  to be the relation such that  $(p, q) \in \delta_a$  if and only if  $(p, a, q) \in \delta$  for any states  $p, q \in Q$ .
- (b) For each string  $w \in \Sigma^*$ , we define  $\delta_w \subseteq Q \times Q$  as follows.
  - $\delta_\epsilon$  is the identity relation.
  - For any strings  $x \in \Sigma^*$ , any symbol  $a \in \Sigma$  and any states  $p, q \in Q$ ,

$$(p, q) \in \delta_{xa}$$

if and only if there exists a state  $r \in Q$  such that

$$(p, r) \in \delta_x \quad \text{and} \quad (r, q) \in \delta_a.$$

**Definition 1.10.** Let  $A = (Q, \Sigma, \delta, q_0, F)$  be an NFA.

- We say that  $A$  accepts a string  $w \in \Sigma^*$  if there exists  $q \in F$  such that  $(q_0, q) \in \delta_w$ .
- The **language** of  $A$ , denoted  $L(A)$ , is defined as the set of strings that are accepted by  $A$ .

**Theorem 1.11.** For every NFA  $A$ , there is a DFA  $A'$  with  $L(A') = L(A)$ .

*Proof.* Let  $A = (Q, \Sigma, \delta, q_0, F)$ . We construct  $A' = (\mathcal{P}(Q), \Sigma, \Delta, \{q_0\}, \Phi)$  as follows.

- $\Delta : \mathcal{P}(Q) \times \Sigma \rightarrow \mathcal{P}(Q)$  is the function with

$$\Delta_a(P) = \bigcup_{p \in P} \{q \in Q : (p, q) \in \delta_a\}$$

for any  $P \subseteq Q$  and  $a \in \Sigma$ .

- $\Phi = \{P \subseteq Q : P \cap F \neq \emptyset\}$ .

Now we prove that for any  $w \in \Sigma^*$ , for any  $q \in Q$  and for any  $P \subseteq Q$ , we have  $q \in \Delta_w(P)$  if and only if  $(p, q) \in \delta_w$  for some  $p \in P$ . For the induction basis, let  $w = \epsilon$ , and we have

$$q \in \Delta_\epsilon(P) \Leftrightarrow q \in P \Leftrightarrow (p, q) \in \delta_\epsilon \text{ for some } p \in P.$$

For the induction step, let  $w = xa$ , where  $a$  is the last symbol of  $w$ . Note that by the construction of  $\Delta$ , we have  $q \in \Delta_a(P)$  if and only if  $(p, q) \in \delta_a$  for some  $p \in P$ . Thus, we can conclude that

$$\begin{aligned} q \in \Delta_{xa}(P) &\Leftrightarrow q \in \Delta_a(\Delta_x(P)) \\ &\Leftrightarrow (r, q) \in \delta_a \text{ for some } r \in \Delta_x(P) \\ &\quad \text{and } (p, r) \in \delta_x \text{ for some } p \in P \\ &\Leftrightarrow (p, q) \in \delta_{xa} \text{ for some } p \in P. \end{aligned}$$

Finally we prove that  $L(A') = L(A)$ , which is given by

$$\begin{aligned} w \in L(A') &\Leftrightarrow \Delta_w(\{q_0\}) \in \Phi \\ &\Leftrightarrow \Delta_w(\{q_0\}) \cap F \neq \emptyset \\ &\Leftrightarrow q \in \Delta_w(\{q_0\}) \text{ for some } q \in F \\ &\Leftrightarrow (p, q) \in \delta_w \text{ for some } q \in F \text{ and } p \in \{q_0\} \\ &\Leftrightarrow (q_0, q) \in \delta_w \text{ for some } q \in F \\ &\Leftrightarrow w \in L(A). \end{aligned}$$

□

## 1.4 Regular Expressions

**Definition 1.12.** Let  $\Sigma$  be an alphabet. A **regular expression** over  $\Sigma$  is a string in the minimal language over  $\Sigma \cup \{\emptyset, \epsilon, *, +, (, )\}$  that satisfies the following conditions.

1.  $\emptyset$  is a regular expression.
2.  $\epsilon$  is a regular expression.
3. If  $a \in \Sigma$ , then  $a$  is a regular expression.
4. If  $e_1$  and  $e_2$  are regular expressions, then so is  $(e_1 e_2)$ .
5. If  $e_1$  and  $e_2$  are regular expressions, then so is  $(e_1 + e_2)$ .
6. If  $e$  is a regular expression, then so is  $(e)^*$ .

**Definition 1.13.** A regular expression  $e$  over an alphabet  $\Sigma$  defines a language  $L(e)$  as follows.

1.  $L(\emptyset) = \emptyset$ .
2.  $L(\epsilon) = \{\epsilon\}$ .
3.  $L(a) = \{a\}$  for each  $a \in \Sigma$ .
4.  $L((e_1 e_2)) = L(e_1)L(e_2)$  for each regular expressions  $e_1$  and  $e_2$ .
5.  $L((e_1 + e_2)) = L(e_1) \cup L(e_2)$  for each regular expressions  $e_1$  and  $e_2$ .
6.  $L((e)^*) = L(e)^*$  for each regular expression  $e$ .

**Remark.** From now on, we may omit parentheses if there is no ambiguity.

**Lemma 1.14.** If  $L$  is a regular language over an alphabet  $\Sigma$ , then there is a regular expression  $e$  over  $\Sigma$  such that  $L(e) = L$ .

*Proof.* Since  $L$  is regular, there exists a DFA  $A = (Q, \Sigma, \delta, q_0, F)$  with  $L(A) = L$ . Suppose that  $Q = \{p_1, p_2, \dots, p_n\}$  with  $p_1 = q_0$ . For any  $i, j \in \{1, \dots, n\}$  and for any  $k \in \{0, \dots, n\}$ , let  $L_{ij}^{(k)}$  denote the language of strings  $w$  such that

- $\delta_w(p_i) = p_j$ , and
- for each string  $x$  with  $\epsilon \sqsubset x \sqsubset w$ , we have  $\delta_x(p_i) = p_\ell$  for some  $\ell \in \{1, \dots, k\}$ .

We are going to prove that for all  $i, j \in \{1, \dots, n\}$  and  $k \in \{0, \dots, n\}$ , there exists a regular expression  $e_{ij}^{(k)}$  such that

$$L(e_{ij}^{(k)}) = L_{ij}^{(k)}.$$

The proof is by induction on  $k$ . For the induction basis, let  $k = 0$ . Let  $\Pi_{ij} \subseteq \Sigma$  denote the set of symbols  $a$  with  $\delta_a(p_i) = p_j$ . If  $i \neq j$ , we have

$$L_{ij}^{(0)} = \bigcup_{a \in \Pi_{ij}} \{a\},$$

and thus we can construct  $e_{ij}^{(0)}$  by

$$e_{ij}^{(0)} = \sum_{a \in \Pi_{ij}} a.$$

(If  $\Pi_{ij} = \emptyset$ , then the summation is defined as  $\emptyset$ .) If  $i = j$ , we have

$$L_{ii}^{(0)} = \{\epsilon\} \cup \bigcup_{a \in \Pi_{ii}} \{a\},$$

and thus we can construct  $e_{ii}^{(0)}$  by

$$e_{ii}^{(0)} = \epsilon + \sum_{a \in \Pi_{ii}} a.$$

Now for the induction step, let  $k \geq 1$ . Suppose that  $w \in L_{ij}^{(k)}$ . If there is no string  $x$  with  $\epsilon \sqsubset x \sqsubset w$  such that  $\delta_x(p_i) = p_k$ , then we have

$$w \in L_{ij}^{(k-1)}.$$

Otherwise, let  $x_0, x_1, \dots, x_\ell$  be all strings with  $\epsilon \sqsubset x_0 \sqsubset x_1 \sqsubset \dots \sqsubset x_\ell \sqsubset w$  such that

$$\delta_{x_0}(p_i) = \delta_{x_1}(p_i) = \dots = \delta_{x_\ell}(p_i) = p_k.$$

Let  $u_0, u_1, \dots, u_{\ell+1}$  be strings such that

$$w = u_0 u_1 \dots u_{\ell+1},$$

and  $x_h = u_0 u_1 \dots u_h$  for each  $h \in \{0, \dots, \ell\}$ . Since  $u_0 \in L_{ik}^{(k-1)}$ ,  $u_{\ell+1} \in L_{kj}^{(k-1)}$ , and  $u_h \in L_{kk}^{(k-1)}$  for each  $h \in \{1, \dots, \ell\}$ , it follows that

$$w \in L_{ik}^{(k-1)} \left( L_{kk}^{(k-1)} \right)^* L_{kj}^{(k-1)}.$$

As a result, we have

$$L_{ij}^{(k)} \subseteq L_{ij}^{(k-1)} \cup L_{ik}^{(k-1)} \left( L_{kk}^{(k-1)} \right)^* L_{kj}^{(k-1)},$$

implying

$$L_{ij}^{(k)} = L_{ij}^{(k-1)} \cup L_{ik}^{(k-1)} \left( L_{kk}^{(k-1)} \right)^* L_{kj}^{(k-1)}.$$

Therefore, we can construct  $e_{ij}^{(k)}$  by

$$e_{ij}^{(k)} = e_{ij}^{(k-1)} + e_{ik}^{(k-1)} \left( e_{kk}^{(k-1)} \right)^* e_{kj}^{(k-1)}.$$

Now we can construct the regular expression  $e$  with  $L(e) = L$  as follows. Let  $\Phi$  be the set of integers  $j \in \{1, \dots, n\}$  such that  $p_j \in F$ . Since

$$L = \bigcup_{j \in \Phi} L_{1j}^{(n)},$$

we can construct  $e$  by

$$e = \sum_{j \in \Phi} e_{1j}^{(n)},$$

which completes the proof. □