

Chapter 5

Regular Languages and Regular Expressions

5.1 Directed Graphs and Paths

It is often useful to view DFA's and NFA's as labeled directed graphs.

Definition 5.1. A *directed graph* is a quadruple $G = (V, E, s, t)$, where V is a set of *vertices, or nodes*, E is a set of *edges, or arcs*, and $s, t: E \rightarrow V$ are two functions, s being called the *source* function, and t the *target* function. Given an edge $e \in E$, we also call $s(e)$ the *origin* (or *source*) of e , and $t(e)$ the *endpoint* (or *target*) of e .

Remark: the functions s, t need not be injective or surjective. Thus, we allow “isolated vertices.”

Example: Let G be the directed graph defined such that

$$E = \{e_1, e_2, e_3, e_4, e_5, e_6, e_7, e_8\},$$

$$V = \{v_1, v_2, v_3, v_4, v_5, v_6\}, \text{ and}$$

$$s(e_1) = v_1, \quad s(e_2) = v_2, \quad s(e_3) = v_3, \quad s(e_4) = v_4,$$

$$s(e_5) = v_2, \quad s(e_6) = v_5, \quad s(e_7) = v_5, \quad s(e_8) = v_5,$$

$$t(e_1) = v_2, \quad t(e_2) = v_3, \quad t(e_3) = v_4, \quad t(e_4) = v_2,$$

$$t(e_5) = v_5, \quad t(e_6) = v_5, \quad t(e_7) = v_6, \quad t(e_8) = v_6.$$

Such a graph can be represented by the following diagram:

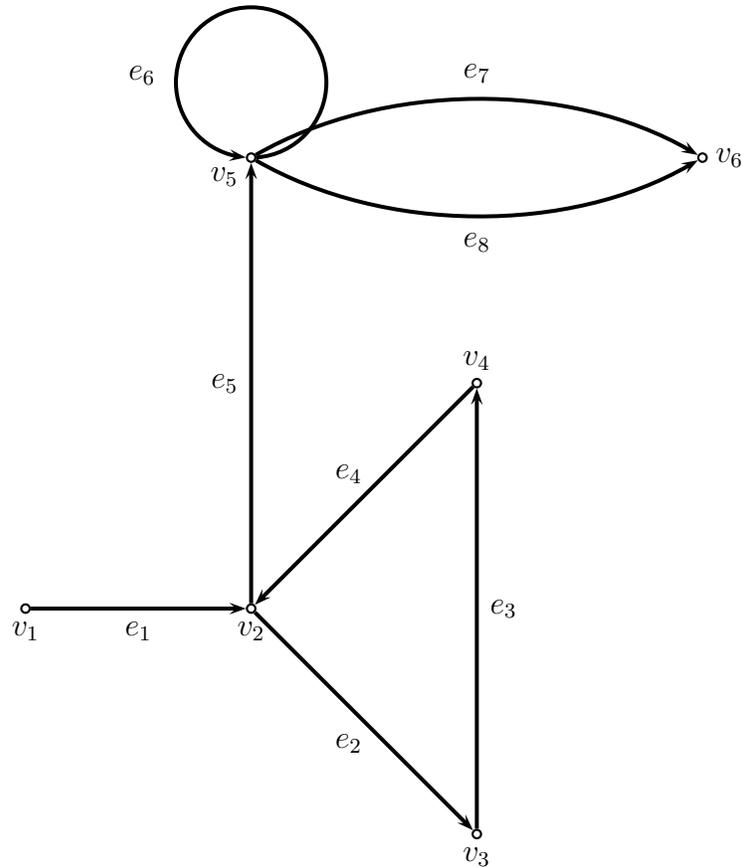


Figure 5.1: A directed graph.

In drawing directed graphs, we will usually omit edge names (the e_i), and sometimes even the node names (the v_j).

We now define paths in a directed graph.

Definition 5.2. Given a directed graph $G = (V, E, s, t)$, for any two nodes $u, v \in V$, a *path from u to v* is a triple $\pi = (u, e_1 \dots e_n, v)$, where $e_1 \dots e_n$ is a string (sequence) of edges in E such that, $s(e_1) = u$, $t(e_n) = v$, and $t(e_i) = s(e_{i+1})$, for all i such that $1 \leq i \leq n - 1$. When $n = 0$, we must have $u = v$, and the path (u, ϵ, u) is called the *null path from u to u* . The number n is the *length* of the path. We also call u the *source* (or *origin*) of the path, and v the *target* (or *endpoint*) of the path. When there is a nonnull path π from u to v , we say that *u and v are connected*.

Remark: In a path $\pi = (u, e_1 \dots e_n, v)$, the expression $e_1 \dots e_n$ is a **sequence**, and thus, the e_i are **not** necessarily distinct.

For example, the following are paths:

$$\pi_1 = (v_1, e_1 e_5 e_7, v_6),$$

$$\pi_2 = (v_2, e_2 e_3 e_4 e_2 e_3 e_4 e_2 e_3 e_4, v_2),$$

and

$$\pi_3 = (v_1, e_1 e_2 e_3 e_4 e_2 e_3 e_4 e_5 e_6 e_6 e_8, v_6).$$

Clearly, π_2 and π_3 are of a different nature from π_1 . Indeed, they contain cycles. This is formalized as follows.

Definition 5.3. Given a directed graph $G = (V, E, s, t)$, for any node $u \in V$ a *cycle (or loop) through u* is a nonnull path of the form $\pi = (u, e_1 \dots e_n, u)$ (equivalently, $t(e_n) = s(e_1)$). More generally, a nonnull path $\pi = (u, e_1 \dots e_n, v)$ *contains a cycle* iff for some i, j , with $1 \leq i \leq j \leq n$, $t(e_j) = s(e_i)$. In this case, letting $w = t(e_j) = s(e_i)$, the path $(w, e_i \dots e_j, w)$ is a cycle through w . A path π is *acyclic* iff it does not contain any cycle. Note that each null path (u, ϵ, u) is acyclic.

Obviously, a cycle $\pi = (u, e_1 \dots e_n, u)$ through u is also a cycle through every node $t(e_i)$. Also, a path π may contain several different cycles.

Paths can be concatenated as follows.

Definition 5.4. Given a directed graph $G = (V, E, s, t)$, two paths $\pi_1 = (u, e_1 \dots e_m, v)$ and $\pi_2 = (u', e'_1 \dots e'_n, v')$ can be *concatenated* provided that $v = u'$, in which case their *concatenation* is the path

$$\pi_1\pi_2 = (u, e_1 \dots e_m e'_1 \dots e'_n, v').$$

It is immediately verified that the concatenation of paths is associative, and that the concatenation of the path $\pi = (u, e_1 \dots e_m, v)$ with the null path (u, ϵ, u) or with the null path (v, ϵ, v) is the path π itself.

The following fact, although almost trivial, is used all the time, and is worth stating in detail.

Lemma 5.1. *Given a directed graph $G = (V, E, s, t)$, if the set of nodes V contains $m \geq 1$ nodes, then every path π of length at least m contains some cycle.*

A consequence of lemma 5.1 is that in a finite graph with m nodes, given any two nodes $u, v \in V$, in order to find out whether there is a path from u to v , it is enough to consider paths of length $\leq m - 1$.

Indeed, if there is path between u and v , then there is some path π of minimal length (not necessarily unique, but this doesn't matter).

If this minimal path has length at least m , then by the lemma, it contains a cycle.

However, by deleting this cycle from the path π , we get an even shorter path from u to v , contradicting the minimality of π .

We now turn to labeled graphs.

5.2 Labeled Graphs and Automata

In fact, we only need edge-labeled graphs.

Definition 5.5. A *labeled directed graph* is a tuple $G = (V, E, L, s, t, \lambda)$, where V is a set of *vertices, or nodes*, E is a set of *edges, or arcs*, L is a set of *labels*, $s, t: E \rightarrow V$ are two functions, s being called the *source* function, and t the *target* function, and $\lambda: E \rightarrow L$ is the *labeling function*. Given an edge $e \in E$, we also call $s(e)$ the *origin* (or *source*) of e , $t(e)$ the *endpoint* (or *target*) of e , and $\lambda(e)$ the *label* of e .

Note that the function λ need not be injective or surjective. Thus, distinct edges may have the same label.

Example: Let G be the directed graph defined such that

$$E = \{e_1, e_2, e_3, e_4, e_5, e_6, e_7, e_8\},$$

$$V = \{v_1, v_2, v_3, v_4, v_5, v_6\}, L = \{a, b\},$$

and

$$s(e_1) = v_1, \quad s(e_2) = v_2, \quad s(e_3) = v_3, \quad s(e_4) = v_4,$$

$$s(e_5) = v_2, \quad s(e_6) = v_5, \quad s(e_7) = v_5, \quad s(e_8) = v_5,$$

$$t(e_1) = v_2, \quad t(e_2) = v_3, \quad t(e_3) = v_4, \quad t(e_4) = v_2,$$

$$t(e_5) = v_5, \quad t(e_6) = v_5, \quad t(e_7) = v_6, \quad t(e_8) = v_6.$$

$$\lambda(e_1) = a, \quad \lambda(e_2) = b, \quad \lambda(e_3) = a, \quad \lambda(e_4) = a,$$

$$\lambda(e_5) = b, \quad \lambda(e_6) = a, \quad \lambda(e_7) = a, \quad \lambda(e_8) = b.$$

Such a labeled graph can be represented by the following diagram:

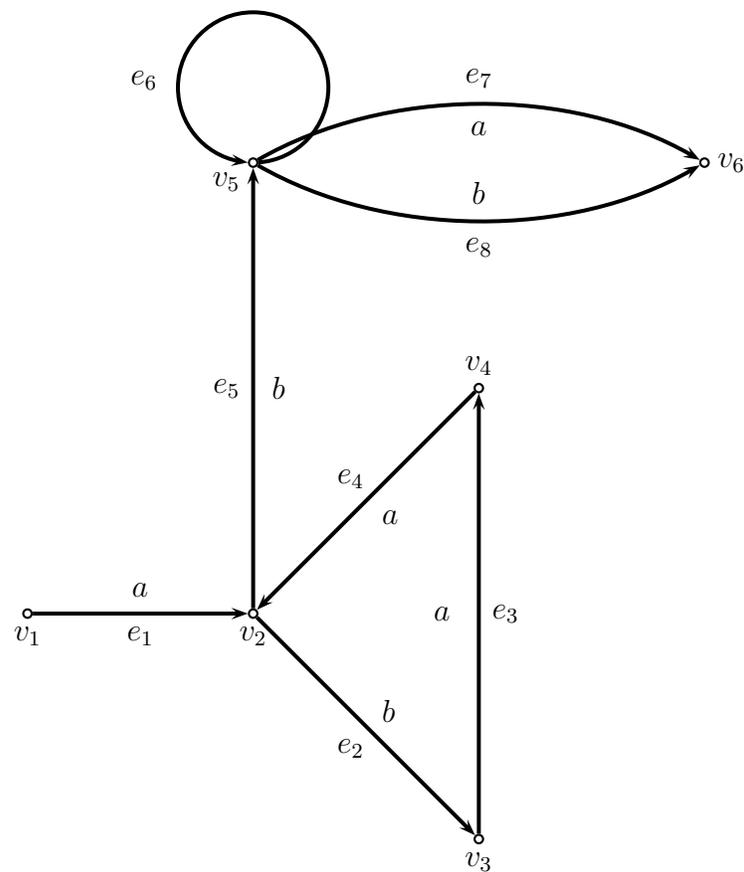


Figure 5.2: A labeled directed graph.

In drawing labeled graphs, we will usually omit edge names (the e_i), and sometimes even the node names (the v_j).

Paths, cycles, and concatenation of paths are defined just as before (that is, we ignore the labels). However, we can now define the *spelling* of a path.

Definition 5.6. Given a labeled directed graph $G = (V, E, L, s, t, \lambda)$ for any two nodes $u, v \in V$, for any path $\pi = (u, e_1 \dots e_n, v)$, the *spelling of the path π* is the string of labels

$$\lambda(e_1) \dots \lambda(e_n).$$

When $n = 0$, the spelling of the null path (u, ϵ, u) is the null string ϵ .

For example, the spelling of the path

$$\pi_3 = (v_1, e_1 e_2 e_3 e_4 e_2 e_3 e_4 e_5 e_6 e_6 e_8, v_6)$$

is

$$abaabaabaab.$$

Every DFA and every NFA can be viewed as a labeled graph, in such a way that the set of spellings of paths from the start state to some final state is the language accepted by the automaton in question.

Given a DFA $D = (Q, \Sigma, \delta, q_0, F)$, where $\delta: Q \times \Sigma \rightarrow Q$, we associate the labeled directed graph $G_D = (V, E, L, s, t, \lambda)$ defined as follows:

$$V = Q, \quad E = \{(p, a, q) \mid q = \delta(p, a), p, q \in Q, a \in \Sigma\},$$

$$L = \Sigma, \quad s((p, a, q)) = p, \quad t((p, a, q)) = q,$$

$$\text{and } \lambda((p, a, q)) = a.$$

Such labeled graphs have a special structure that can easily be characterized.

It is easily shown that a string $w \in \Sigma^*$ is in the language $L(D) = \{w \in \Sigma^* \mid \delta^*(q_0, w) \in F\}$ iff w is the spelling of some path in G_D from q_0 to some final state.

Similarly, given an NFA $N = (Q, \Sigma, \delta, q_0, F)$, where $\delta: Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$, we associate the labeled directed graph $G_N = (V, E, L, s, t, \lambda)$ defined as follows:

$$V = Q$$

$$E = \{(p, a, q) \mid q \in \delta(p, a), p, q \in Q, a \in \Sigma \cup \{\epsilon\}\},$$

$$L = \Sigma \cup \{\epsilon\}, \quad s((p, a, q)) = p, \quad t((p, a, q)) = q,$$

$$\lambda((p, a, q)) = a.$$

Remark: When N has no ϵ -transitions, we can let $L = \Sigma$.

Such labeled graphs have also a special structure that can easily be characterized.

Again, a string $w \in \Sigma^*$ is in the language

$L(N) = \{w \in \Sigma^* \mid \delta^*(q_0, w) \cap F \neq \emptyset\}$ iff w is the spelling of some path in G_N from q_0 to some final state.

5.3 The Closure Definition of the Regular Languages

Let $\Sigma = \{a_1, \dots, a_m\}$ be an alphabet.

Informally, we define the family of languages $R(\Sigma)$ using the following rules:

- (1) The languages $\{a_1\}, \dots, \{a_m\}$, the empty language, and the trivial language $\{\epsilon\}$, called *base languages*, belong to $R(\Sigma)$.
- (2a) If L_1 and L_2 belong to $R(\Sigma)$, then $L_1 \cup L_2$ also belongs to $R(\Sigma)$.
- (2b) If L_1 and L_2 belong to $R(\Sigma)$, then $L_1 L_2$ also belongs to $R(\Sigma)$.
- (2c) If L belongs to $R(\Sigma)$, then L^* also belongs to $R(\Sigma)$.

The issue is to show that the above rules define a family of languages which is the smallest family containing the base languages and closed under union, concatenation, and Kleene $*$.

We define the family $(R(\Sigma)_n)$ of sets of languages as follows:

$$\begin{aligned} R(\Sigma)_0 &= \{\{a_1\}, \dots, \{a_m\}, \emptyset, \{\epsilon\}\}, \\ R(\Sigma)_{n+1} &= R(\Sigma)_n \cup \{L_1 \cup L_2, L_1L_2, L^* \mid \\ &\quad L_1, L_2, L \in R(\Sigma)_n\}. \end{aligned}$$

Then, we define $R(\Sigma)$ as

$$R(\Sigma) = \bigcup_{n \geq 0} R(\Sigma)_n.$$

For example, if $\Sigma = \{a, b\}$, we have

$$\begin{aligned} R(\Sigma)_1 &= \{\{a\}, \{b\}, \emptyset, \{\epsilon\}, \\ &\quad \{a, b\}, \{a, \epsilon\}, \{b, \epsilon\}, \\ &\quad \{aa\}, \{ab\}, \{ba\}, \{bb\}, \{a\}^*, \{b\}^*\}. \end{aligned}$$

Some of the languages that will appear in $R(\Sigma)_2$ are:

$$\{a, bb\}, \{ab, ba\}, \{abb\}, \{aabb\}, \{a\}\{a\}^*, \{aa\}\{b\}^*, \{bb\}^*.$$

Observe that each family $R(\Sigma)_n$ contains a *finite number* of languages.

Definition 5.7. *Regular languages, Version 2 = $R(\Sigma)$.*

Consider the following properties of a family of languages, $\mathcal{L} \subseteq 2^{\Sigma^*}$:

- (1) $\{a_1\}, \dots, \{a_m\}, \emptyset, \{\epsilon\} \in \mathcal{L}$
- 2(a) If $L_1 \in \mathcal{L}$ and $L_2 \in \mathcal{L}$, then $L_1 \cup L_2 \in \mathcal{L}$
- 2(b) If $L_1 \in \mathcal{L}$ and $L_2 \in \mathcal{L}$, then $L_1L_2 \in \mathcal{L}$
- 2(c) If $L \in \mathcal{L}$, then $L^* \in \mathcal{L}$.

If properties 2(a), 2(b) and 2(c) hold, we say that the family, \mathcal{L} , is *closed under union, concatenation and Kleene **.

Proposition 5.2. *The family $R(\Sigma)$ is the smallest family of languages which contains the (atomic) languages $\{a_1\}, \dots, \{a_m\}, \emptyset, \{\epsilon\}$, and is closed under union, concatenation, and Kleene $*$.*

Proof sketch. To prove that $R(\Sigma)$ satisfies properties (1), 2(a), 2(b) and 2(c), use the fact that $R(\Sigma)_n \subseteq R(\Sigma)_{n+1}$ for all $n \geq 0$.

To prove that for any family, \mathcal{L} , if \mathcal{L} satisfies properties (1), 2(a), 2(b) and 2(c), then $R(\Sigma) \subseteq \mathcal{L}$, prove that $R(\Sigma)_n \subseteq \mathcal{L}$ by induction on n . \square

Note: a given language L may be built up in different ways. For example,

$$\{a, b\}^* = (\{a\}^* \{b\}^*)^*.$$

The definition of the regular languages that we just gave is not very convenient to manipulate them in a practical way.

A better formalism is to represent regular languages in terms of certain strings called regular expressions.

5.4 Regular Expressions

Given an alphabet $\Sigma = \{a_1, \dots, a_m\}$, consider the new alphabet

$$\Delta = \Sigma \cup \{+, \cdot, *, (,), \emptyset, \epsilon\}.$$

Informally, we define the family of regular expressions $\mathcal{R}(\Sigma)$ using the following rules:

- (1) The strings a_1, \dots, a_m , the empty string ϵ , and the empty set \emptyset , called *base regular expressions*, belong to $\mathcal{R}(\Sigma)$.
- (2a) If R_1 and R_2 are regular expressions (*i.e.*, belong to $\mathcal{R}(\Sigma)$), then $(R_1 + R_2)$ is a regular expression (*i.e.*, belongs to $\mathcal{R}(\Sigma)$).
- (2b) If R_1 and R_2 are regular expressions (*i.e.*, belong to $\mathcal{R}(\Sigma)$), then $(R_1 \cdot R_2)$ is a regular expression (*i.e.*, belongs to $\mathcal{R}(\Sigma)$).
- (2c) If R is a regular expression (*i.e.*, belongs to $\mathcal{R}(\Sigma)$), then R^* is a regular expression (*i.e.*, belongs to $\mathcal{R}(\Sigma)$).

More precisely, we define the family $(\mathcal{R}(\Sigma)_n)$ of languages over Δ as follows:

$$\begin{aligned}\mathcal{R}(\Sigma)_0 &= \{a_1, \dots, a_m, \emptyset, \epsilon\}, \\ \mathcal{R}(\Sigma)_{n+1} &= \mathcal{R}(\Sigma)_n \cup \{(R_1 + R_2), (R_1 \cdot R_2), R^* \mid \\ &\quad R_1, R_2, R \in \mathcal{R}(\Sigma)_n\}.\end{aligned}$$

Then, we define $\mathcal{R}(\Sigma)$ as

$$\mathcal{R}(\Sigma) = \bigcup_{n \geq 0} \mathcal{R}(\Sigma)_n.$$

Note that every language $\mathcal{R}(\Sigma)_n$ is finite.

For example, if $\Sigma = \{a, b\}$, we have

$$\begin{aligned} \mathcal{R}(\Sigma)_1 = \{ & a, b, \emptyset, \epsilon, \\ & (a + b), (b + a), (a + a), (b + b), (a + \epsilon), (\epsilon + a), \\ & (b + \epsilon), (\epsilon + b), (a + \emptyset), (\emptyset + a), (b + \emptyset), (\emptyset + b), \\ & (\epsilon + \epsilon), (\epsilon + \emptyset), (\emptyset + \epsilon), (\emptyset + \emptyset), \\ & (a \cdot b), (b \cdot a), (a \cdot a), (b \cdot b), (a \cdot \epsilon), (\epsilon \cdot a), \\ & (b \cdot \epsilon), (\epsilon \cdot b), (\epsilon \cdot \epsilon), (a \cdot \emptyset), (\emptyset \cdot a), \\ & (b \cdot \emptyset), (\emptyset \cdot b), (\epsilon \cdot \emptyset), (\emptyset \cdot \epsilon), (\emptyset \cdot \emptyset), \\ & a^*, b^*, \epsilon^*, \emptyset^* \}. \end{aligned}$$

Some of the regular expressions appearing in $\mathcal{R}(\Sigma)_2$ are:

$$\begin{aligned} & (a + (b \cdot b)), ((a \cdot b) + (b \cdot a)), ((a \cdot b) \cdot b), \\ & ((a \cdot a) \cdot (b \cdot b)), (a \cdot a^*), ((a \cdot a) \cdot b^*), (b \cdot b)^*. \end{aligned}$$

Definition 5.8. $\mathcal{R}(\Sigma)$ is the set of *regular expressions* (over Σ).

Proposition 5.3. *The language $\mathcal{R}(\Sigma)$ is the smallest language which contains the symbols $a_1, \dots, a_m, \emptyset, \epsilon$, from Δ , and such that $(R_1 + R_2)$, $(R_1 \cdot R_2)$, and R^* , also belong to $\mathcal{R}(\Sigma)$, when $R_1, R_2, R \in \mathcal{R}(\Sigma)$.*

For simplicity of notation, write

$$(R_1 R_2)$$

instead of

$$(R_1 \cdot R_2).$$

Examples: $R = (a + b)^*$, $S = (a^* b^*)^*$.

$$T = (((a + b)^* a) \underbrace{((a + b) \cdots (a + b))}_n).$$

5.5 Regular Expressions and Regular Languages

Every regular expression $R \in \mathcal{R}(\Sigma)$ can be viewed as the *name*, or *denotation*, of some language $L \in \mathcal{R}(\Sigma)$. Similarly, every language $L \in \mathcal{R}(\Sigma)$ is the *interpretation* (or *meaning*) of some regular expression $R \in \mathcal{R}(\Sigma)$.

Think of a regular expression R as a *program*, and of $\mathcal{L}(R)$ as the result of the *execution* or *evaluation*, of R by \mathcal{L} .

This can be made rigorous by defining a function

$$\mathcal{L}: \mathcal{R}(\Sigma) \rightarrow \mathcal{R}(\Sigma).$$

This function is defined recursively:

$$\begin{aligned}
 \mathcal{L}[a_i] &= \{a_i\}, \\
 \mathcal{L}[\emptyset] &= \emptyset, \\
 \mathcal{L}[\epsilon] &= \{\epsilon\}, \\
 \mathcal{L}[(R_1 + R_2)] &= \mathcal{L}[R_1] \cup \mathcal{L}[R_2], \\
 \mathcal{L}[(R_1R_2)] &= \mathcal{L}[R_1]\mathcal{L}[R_2], \\
 \mathcal{L}[R^*] &= \mathcal{L}[R]^*.
 \end{aligned}$$

Proposition 5.4. *For every regular expression $R \in \mathcal{R}(\Sigma)$, the language $\mathcal{L}[R]$ is regular (version 2), i.e. $\mathcal{L}[R] \in R(\Sigma)$. Conversely, for every regular (version 2) language $L \in R(\Sigma)$, there is some regular expression $R \in \mathcal{R}(\Sigma)$ such that $L = \mathcal{L}[R]$.*

Note: the function \mathcal{L} is **not** injective.

Example: If $R = (a + b)^*$, $S = (a^*b^*)^*$, then

$$\mathcal{L}[R] = \mathcal{L}[S] = \{a, b\}^*.$$

For simplicity, we often denote $\mathcal{L}[R]$ as L_R . As examples, we have

$$\begin{aligned}\mathcal{L}[(((ab)b) + a)] &= \{a, abb\} \\ \mathcal{L}[((((a^*b)a^*)b)a^*)] &= \{w \in \{a, b\}^* \mid w \text{ has} \\ &\qquad\qquad\qquad \text{two } b\text{'s}\} \\ \mathcal{L}[(((((((a^*b)a^*)b)a^*)^*a^*)] &= \{w \in \{a, b\}^* \mid w \text{ has an} \\ &\qquad\qquad\qquad \text{even } \# \text{ of } b\text{'s}\} \\ \mathcal{L}[(((((((((((a^*b)a^*)b)a^*)^*a^*)b)a^*)] &= \{w \in \{a, b\}^* \mid w \text{ has an} \\ &\qquad\qquad\qquad \text{odd } \# \text{ of } b\text{'s}\}\end{aligned}$$

Remark. If

$$R = (((a + b)^*a)\underbrace{((a + b) \cdots (a + b))}_n),$$

it can be shown that any minimal DFA accepting L_R has 2^{n+1} states.

Yet, both $((a + b)^*a)$ and $\underbrace{((a + b) \cdots (a + b))}_n$ denote languages that can be accepted by “small” DFA’s (of size 2 and $n + 2$).

Definition 5.9. Two regular expressions $R, S \in \mathcal{R}(\Sigma)$ are *equivalent*, denoted as $R \cong S$, iff $\mathcal{L}[R] = \mathcal{L}[S]$.

It is immediate that \cong is an equivalence relation.

The relation \cong satisfies some (nice) identities. For example:

$$\begin{aligned} (((aa) + b) + c) &\cong ((aa) + (b + c)) \\ ((aa)(b(cc))) &\cong (((aa)b)(cc)) \\ (a^*a^*) &\cong a^*, \end{aligned}$$

and more generally

$$\begin{aligned} ((R_1 + R_2) + R_3) &\cong (R_1 + (R_2 + R_3)), \\ ((R_1R_2)R_3) &\cong (R_1(R_2R_3)), \\ (R_1 + R_2) &\cong (R_2 + R_1), \\ (R^*R^*) &\cong R^*, \\ R^{**} &\cong R^*. \end{aligned}$$

There are algorithms to test equivalence of regular expressions, but their complexity is exponential.

It is an *open problem* to prove that the problem cannot be decided in polynomial time.

5.6 Regular Expressions and NFA's

Proposition 5.5. *There is an algorithm, which, given any regular expression $R \in \mathcal{R}(\Sigma)$, constructs an NFA N_R accepting L_R , i.e., such that $L_R = L(N_R)$.*

In order to ensure the correctness of the construction as well as to simplify the description of the algorithm it is convenient to assume that our NFA's satisfy the following conditions:

1. Each NFA has a *single* final state, t , distinct from the start state, s .
2. There are *no incoming transitions* into the the start state, s , and *no outgoing transitions* from the final state, t .
3. Every state has at most two incoming and two outgoing transitions.

Here is the algorithm, sometimes called the *sombrero construction*.

For the base case, either

(a) $R = a_i$, in which case, N_R is the following NFA:

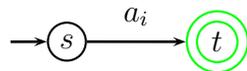


Figure 5.3: NFA for a_i .

(b) $R = \epsilon$, in which case, N_R is the following NFA:

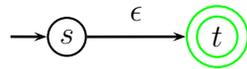


Figure 5.4: NFA for ϵ .

(c) $R = \emptyset$, in which case, N_R is the following NFA:



Figure 5.5: NFA for \emptyset .

The recursive clauses are as follows:

(i) If our expression is $(R + S)$, the algorithm is applied recursively to R and S , generating NFA's N_R and N_S , and then these two NFA's are combined in parallel as shown in Figure 5.6:

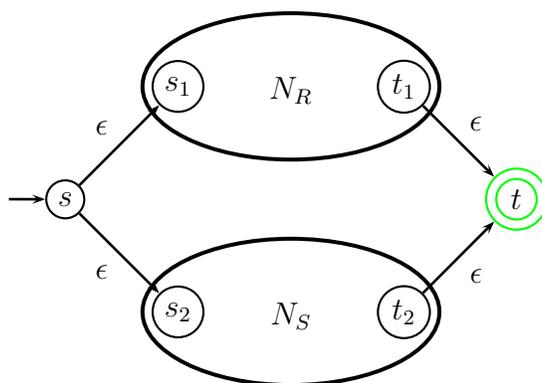


Figure 5.6: NFA for $(R + S)$.

(ii) If our expression is $(R \cdot S)$, the algorithm is applied recursively to R and S , generating NFA's N_R and N_S , and then these NFA's are combined sequentially as shown in Figure 5.7 by merging the “old” final state, t_1 , of N_R , with the “old” start state, s_2 , of N_S :

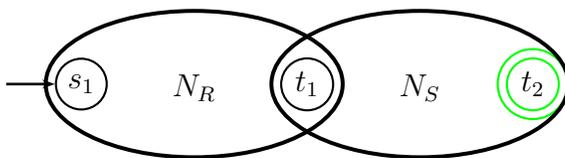


Figure 5.7: NFA for $(R \cdot S)$.

Note that since there are no incoming transitions into s_2 in N_S , once we enter N_S , there is no way of reentering N_R , and so the construction is correct (it yields the concatenation $L_R L_S$).

(iii) If our expression is R^* , the algorithm is applied recursively to R , generating the NFA N_R . Then we construct the NFA shown in Figure 5.8 by adding an ϵ -transition from the “old” final state, t_1 , of N_R to the “old” start state, s_1 , of N_R and, as ϵ is not necessarily accepted by N_R , we add an ϵ -transition from s to t :

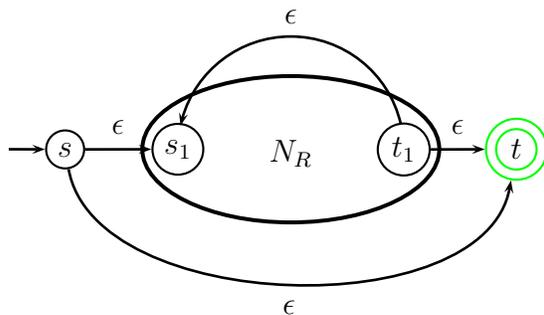


Figure 5.8: NFA for R^* .

Since there are no outgoing transitions from t_1 in N_R , we can only loop back to s_1 from t_1 using the new ϵ -transition from t_1 to s_1 and so the NFA of Figure 5.8 does accept L_R^* .

As a corollary of this construction, we get

Reg. languages version 2 \subseteq Reg. languages, version 1.

The reader should check that if one constructs the NFA corresponding to the regular expression $(a + b)^*abb$ we obtain the NFA shown in Figure 5.9.

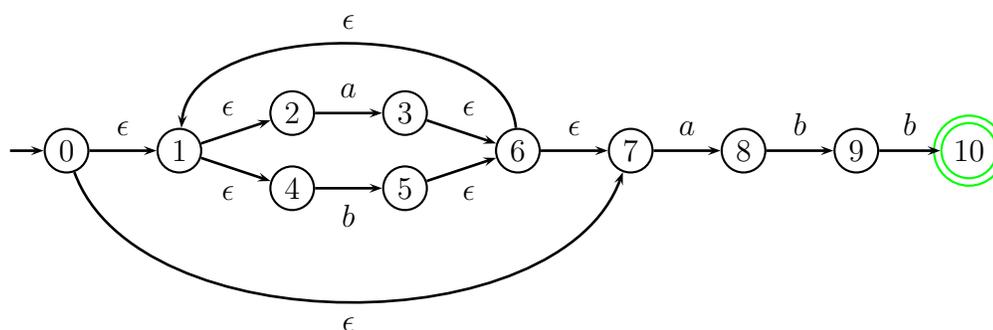


Figure 5.9: An NFA for $R = (a + b)^*abb$.

If we apply the subset construction to the above NFA, we get the following DFA:

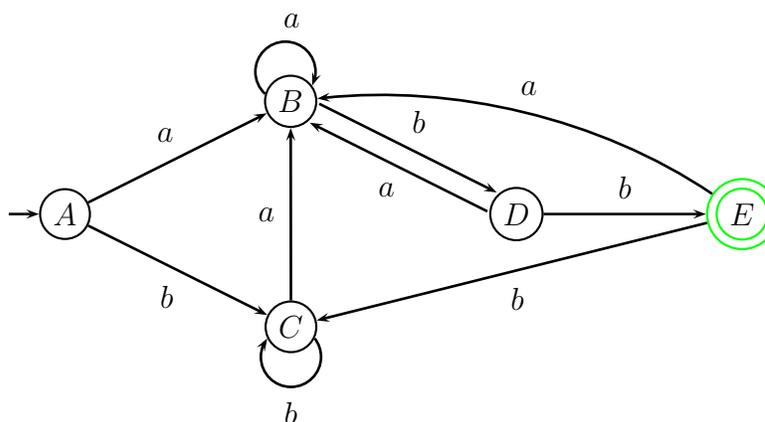


Figure 5.10: A non-minimal DFA for $\{a, b\}^*\{abb\}$.

Proposition 5.6. *There is an algorithm, which, given any NFA N , constructs a regular expression $R \in \mathcal{R}(\Sigma)$, denoting $L(N)$, i.e., such that $L_R = L(N)$.*

As a corollary,

Reg. languages version 1 \subseteq Reg. languages, version 2.

This is the *node elimination algorithm*.

The general idea is to allow more general labels on the edges of an NFA, namely, regular expressions. Then, such generalized NFA's are simplified by eliminating nodes one at a time, and readjusting labels.

Preprocessing, phase 1:

If necessary, we need to add a new start state with an ϵ -transition to the old start state, if there are incoming edges into the old start state.

If necessary, we need to add a new (unique) final state with ϵ -transitions from each of the old final states to the new final state, if there is more than one final state or some outgoing edge from any of the old final states.

At the end of this phase, the start state, say s , is a source (no incoming edges), and the final state, say t , is a sink (no outgoing edges).

Preprocessing, phase 2:

We need to “flatten” parallel edges. For any pair of states (p, q) ($p = q$ is possible), if there are k edges from p to q labeled u_1, \dots, u_k , then create a single edge labeled with the regular expression

$$u_1 + \dots + u_k.$$

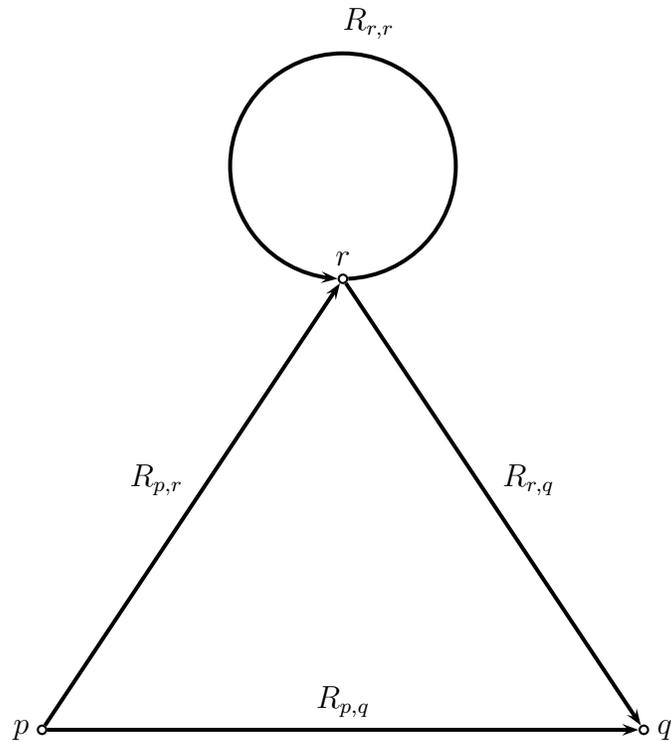
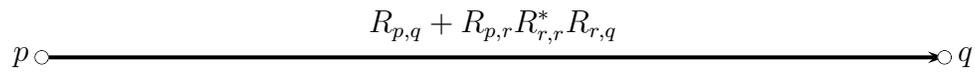
For any pair of states (p, q) ($p = q$ is possible) such that there is **no** edge from p to q , we put an edge labeled \emptyset .

At the end of this phase, the resulting “*generalized NFA*” is such that for any pair of states (p, q) (where $p = q$ is possible), there is a unique edge labeled with some regular expression denoted as $R_{p,q}$. When $R_{p,q} = \emptyset$, this really means that there is no edge from p to q in the original NFA N .

By interpreting each $R_{p,q}$ as a function call (really, a macro) to the NFA $N_{p,q}$ accepting $\mathcal{L}[R_{p,q}]$ (constructed using the previous algorithm), we can verify that the original language $L(N)$ is accepted by this new generalized NFA.

Node elimination only applies if the generalized NFA has at least one node distinct from s and t .

Pick any node r distinct from s and t . For every pair (p, q) where $p \neq r$ and $q \neq r$, replace the label of the edge from p to q as indicated below:

Figure 5.11: Before Eliminating node r .Figure 5.12: After Eliminating node r .

At the end of this step, delete the node r and all edges adjacent to r .

Note that $p = q$ is possible, in which case the triangle is “flat”. It is also possible that $p = s$ or $q = t$. Also, this step is performed for all **pairs** (p, q) , which means that both (p, q) and (q, p) are considered (when $p \neq q$).

Note that this step only has an effect if there are edges from p to r and from r to q in the original NFA N . Otherwise, r can simply be deleted, as well as the edges adjacent to r .

Other simplifications can be made. For example, when $R_{r,r} = \emptyset$, we can simplify $R_{p,r}R_{r,r}^*R_{r,q}$ to $R_{p,r}R_{r,q}$. When $R_{p,q} = \emptyset$, we have $R_{p,r}R_{r,r}^*R_{r,q}$.

The order in which the nodes are eliminated is irrelevant, although it affects the size of the final expression.

The algorithm stops when the only remaining nodes are s and t . Then, the label R of the edge from s to t is a regular expression denoting $L(N)$.

For example, let

$$L = \{w \in \Sigma^* \mid w \text{ contains an odd number of } a\text{'s} \\ \text{or an odd number of } b\text{'s}\}.$$

An NFA for L after the preprocessing phase is:

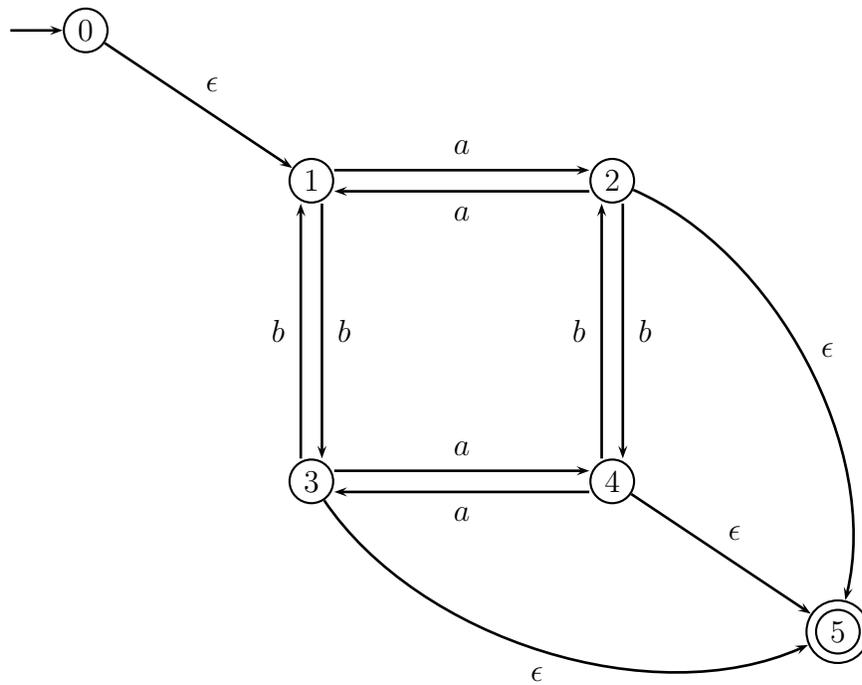


Figure 5.13: NFA for L (after preprocessing phase).

After eliminating node 2:

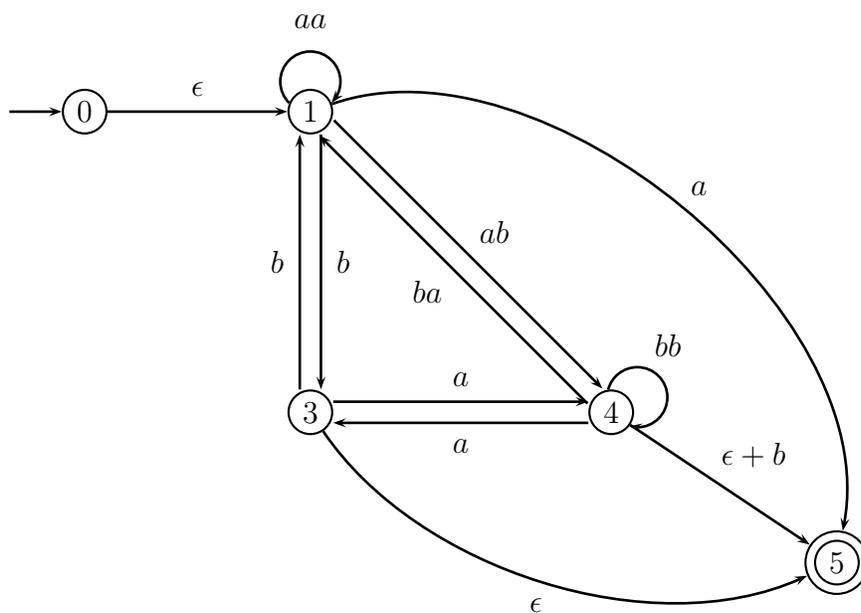


Figure 5.14: NFA for L (after eliminating node 2).

After eliminating node 3:

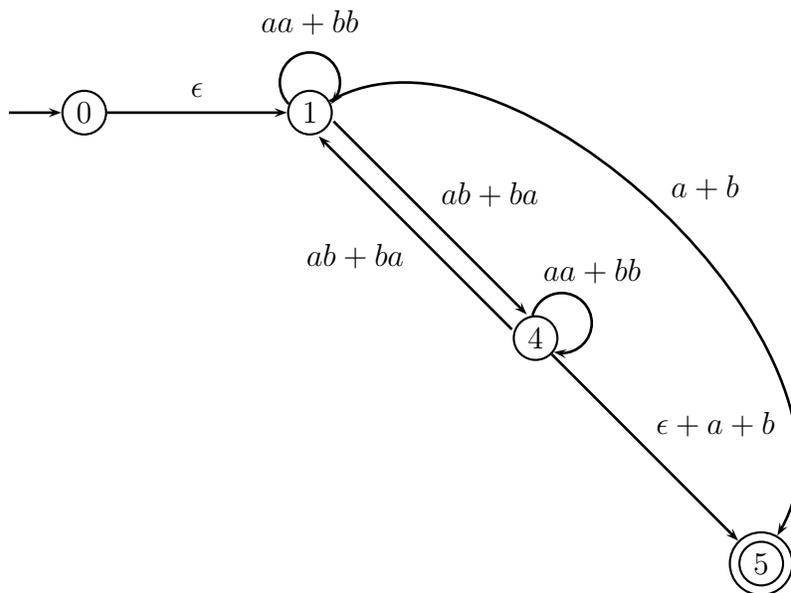


Figure 5.15: NFA for L (after eliminating node 3).

After eliminating node 4:

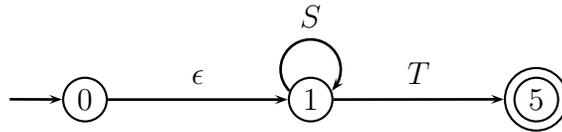


Figure 5.16: NFA for L (after eliminating node 4).

where

$$T = a + b + (ab + ba)(aa + bb)^*(\epsilon + a + b)$$

and

$$S = aa + bb + (ab + ba)(aa + bb)^*(ab + ba).$$

Finally, after eliminating node 1, we get:

$$R = (aa + bb + (ab + ba)(aa + bb)^*(ab + ba))^* \\ (a + b + (ab + ba)(aa + bb)^*(\epsilon + a + b)).$$

5.7 Applications of Regular Expressions: Lexical analysis, Finding patterns in text

Regular expressions have several practical applications. The first important application is to *lexical analysis*.

A *lexical analyzer* is the first component of a *compiler*.

The purpose of a lexical analyzer is to scan the source program and break it into atomic components, known as *tokens*, i.e., substrings of consecutive characters that belong together logically.

Examples of tokens are: identifiers, keywords, numbers (in fixed point notation or floating point notation, etc.), arithmetic operators (+, ·, −, ^), comparison operators (<, >, =, <>), assignment operator (:=), etc.

Tokens can be described by regular expressions. For this purpose, it is useful to enrich the syntax of regular expressions, as in UNIX.

For example, the 26 upper case letters of the (roman) alphabet, A, \dots, Z , can be specified by the expression

$$[A-Z]$$

Similarly, the ten digits, $0, 1, \dots, 9$, can be specified by the expression

$$[0-9]$$

The regular expression

$$R_1 + R_2 + \dots + R_k$$

is denoted

$$[R_1 R_2 \dots R_k]$$

So, the expression

$$[A-Za-z0-9]$$

denotes any letter (upper case or lower case) or digit. This is called an *alphanumeric*.

If we define an identifier as a string beginning with a letter (upper case or lower case) followed by any number of alphanumerics (including none), then we can use the following expression to specify identifiers:

$$[A-Za-z][A-Za-z0-9]^*$$

There are systems, such as **lex** or **flex** that accept as input a list of regular expressions describing the tokens of a programming language and construct a lexical analyzer for these tokens.

Such systems are called *lexical analyzer generators*. Basically, they build a DFA from the set of regular expressions using the algorithms that have been described earlier.

Usually, it is possible associate with every expression some action to be taken when the corresponding token is recognized

Another application of regular expressions is finding patterns in text.

Using a regular expression, we can specify a “vaguely defined” class of patterns.

Take the example of a street address. Most street addresses end with “Street”, or “Avenue”, or “Road” or “St.”, or “Ave.”, or “Rd.”.

We can design a regular expression that captures the shape of most street addresses and then convert it to a DFA that can be used to search for street addresses in text.

For more on this, see Hopcroft-Motwani and Ullman.

5.8 Summary of Closure Properties of the Regular Languages

The family of regular languages is closed under many operations. In particular, it is closed under the following operations listed below. Some of the closure properties are left as a homework problem.

- (1) Union, intersection, relative complement.
- (2) Concatenation, Kleene $*$, Kleene $+$.
- (3) Homomorphisms and inverse homomorphisms.
- (4) gsm and inverse gsm mappings, a -transductions and inverse a -transductions.

Another useful operation is substitution.

Given any two alphabets Σ, Δ , a *substitution* is a function, $\tau: \Sigma \rightarrow 2^{\Delta^*}$, assigning some language, $\tau(a) \subseteq \Delta^*$, to every symbol $a \in \Sigma$.

A substitution $\tau: \Sigma \rightarrow 2^{\Delta^*}$ is extended to a map $\tau: 2^{\Sigma^*} \rightarrow 2^{\Delta^*}$ by first extending τ to strings using the following definition

$$\begin{aligned}\tau(\epsilon) &= \{\epsilon\}, \\ \tau(ua) &= \tau(u)\tau(a),\end{aligned}$$

where $u \in \Sigma^*$ and $a \in \Sigma$, and then to languages by letting

$$\tau(L) = \bigcup_{w \in L} \tau(w),$$

for every language $L \subseteq \Sigma^*$.

Observe that a homomorphism is a special kind of substitution.

A substitution is a *regular* substitution iff $\tau(a)$ is a regular language for every $a \in \Sigma$. The proof of the next proposition is left as a homework problem.

Proposition 5.7. *If L is a regular language and τ is a regular substitution, then $\tau(L)$ is also regular. Thus, the family of regular languages is closed under regular substitutions.*

