

CS 181 (Introduction to Formal Languages and Automata Theory)

May 9, 2022

Contents

1	Deterministic finite automata (DFAs)	3
1.1	Basic notions	3
1.2	Deterministic finite automata	4
1.3	Designing DFAs	6
1.4	Formal definitions	8
2	Nondeterminism	10
2.1	Basic notions	10
2.2	Using shortcuts	11
2.3	Pattern matching	12
2.4	Alternatives	13
2.5	Formal definitions	15
3	Equivalence of DFAs and NFAs	17
3.1	Example	17
3.2	General theorem	19
3.3	Blow-up in size	19
4	Closure	21
4.1	Complement	21
4.2	Union and intersection	23
4.3	Difference	26
4.4	Concatenation	27
4.5	Kleene star	28
4.6	Reverse	29
4.7	Prefix and suffix	30
4.8	Substring	31
4.9	Shuffle	32
5	Regular expressions	33
5.1	Basic notions	33
5.2	Equivalence with automata	35
5.3	Conversion example	36
5.4	Conversion example #2	38
6	Pumping lemma	40
6.1	First nonregular language	40

6.2	Pumping lemma	41
6.3	Examples	42
6.4	True or false	44
7	Myhill-Nerode theorem	45
7.1	When the pumping lemma fails	45
7.2	Indistinguishability	46
7.3	Myhill-Nerode theorem	48
7.4	Proving nonregularity	49
7.5	Proving optimality	50
8	Context-free grammars	52
8.1	Context-free grammars defined	52
8.2	Reading a CFG	54
8.3	Constructing CFGs	55
8.4	Ambiguity	58
8.5	Ambiguity in English	63
9	Pushdown Automata	64
9.1	Definitions	64

1 Deterministic finite automata (DFAs)

1.1 Basic notions

Definition 1.1.1

An **alphabet** is any finite set of symbols.

Example 1.1.2. Binary alphabet: $\{0, 1\}$

Example 1.1.3. English alphabet: $\{a, b, \dots, c\}$

Definition 1.1.4

A **string** is any finite sequence of symbols from a given alphabet.

Example 1.1.5. 001010110101

Example 1.1.6. abracadabra

Example 1.1.7. ε (empty string)

Definition 1.1.8

A **language** is a set of strings over a given alphabet.

Example 1.1.9. \emptyset (empty language)

Example 1.1.10. $\{\varepsilon\}$

Example 1.1.11. $\{\text{acclaim, aim, brim}, \dots\}$

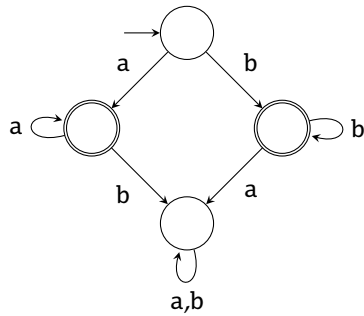
Example 1.1.12. $\{0, 1, 00, 11, \dots\}$

Definition 1.1.13

A **computational device** is a mechanism that inputs a string and either accepts or rejects it.

1.2 Deterministic finite automata

- Choose an alphabet: $\{a, b\}$.
- Draw states.
- Choose start state and accept states.
- Draw transitions (out of every state on every symbol).



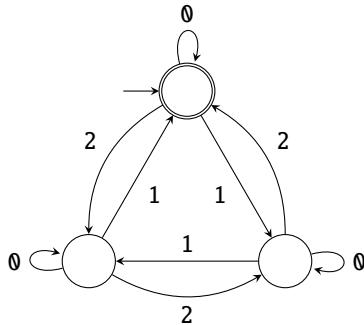
Input	Output
ϵ	reject
ab	reject
aaa	accept
bb	accept

In words, this machine accepts nonempty strings of all a's or all b's.

Definition 1.2.1

The **language** of a DFA is the set of all strings it accepts.

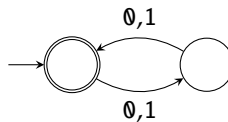
Example 1.2.2.



Input	Output
$00\dots 0$	accept
12	accept
111	accept
20	reject
1	reject

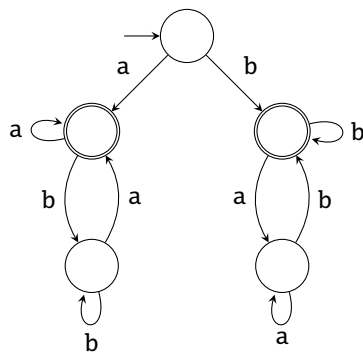
$$\Sigma = \{0, 1, 2\}, \mathcal{L} = \{w : 3 \mid \sum w_i\}$$

Example 1.2.3.



$$\Sigma = \{0, 1\}, \mathcal{L} = \{w : 2 \mid |w|\}$$

Example 1.2.4.



$$\Sigma = \{a, b\}, \mathcal{L} = \{w : w \neq \varepsilon \wedge w_1 = w_{|w|}\}$$

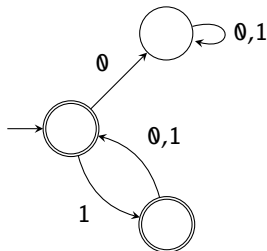
1.3 Designing DFAs

We will be using the binary alphabet $\{0, 1\}$.

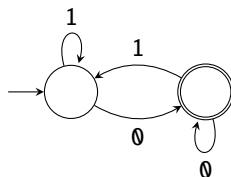
Example 1.3.1. A DFA for $\mathcal{L} = \emptyset$ is given by:



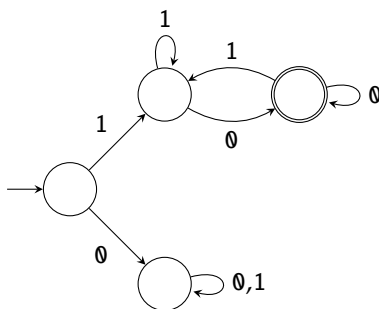
Example 1.3.2. A DFA for $\mathcal{L} = \{w : \text{every odd position is a 1}\}$ is given by:



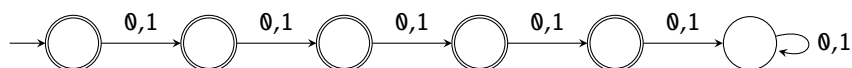
Example 1.3.3. A DFA for $\mathcal{L} = \{w : w \text{ ends in } 0\}$ is given by:



Example 1.3.4. A DFA for $\mathcal{L} = \{w : w \text{ begins with 1, ends with } 0\}$ is given by:



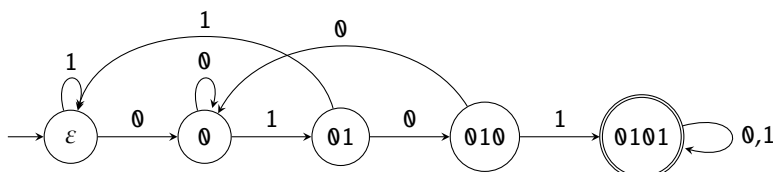
Example 1.3.5. A DFA for $\mathcal{L} = \{w : |w| \leq 4\}$ is given by:



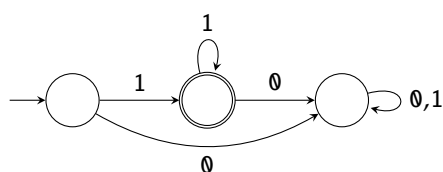
Example 1.3.6. A DFA for $\mathcal{L} = \{w : 1000 \mid |w|\}$ is given by:

In words, each state represents a remainder modulo 1000, and only the 0 state is accepting.

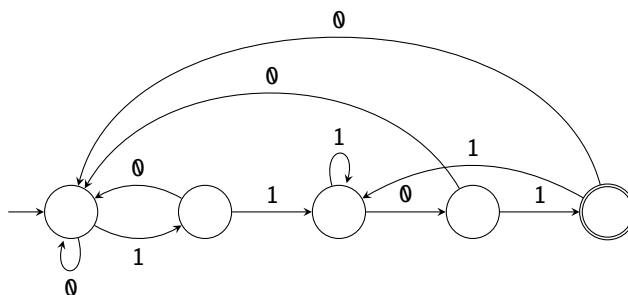
Example 1.3.7. A DFA for $\mathcal{L} = \{w : w \text{ contains } 0101 \text{ as a substring}\}$ is given by:



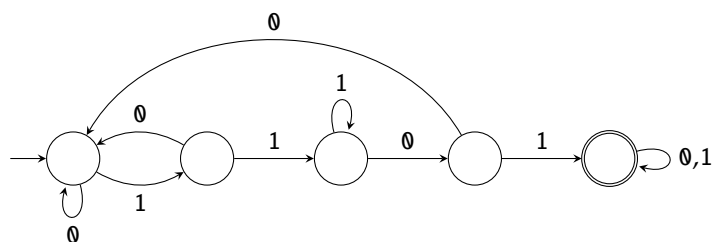
Example 1.3.8 (Week 1 Discussion). A DFA for $\mathcal{L} = \{w : |w| > 0 \wedge w \text{ contains only 1's}\}$ is given by:



Example 1.3.9 (Week 1 Discussion). A DFA for $\mathcal{L} = \{w : w \text{ ends in } 1101\}$ is given by:



Example 1.3.10 (Week 1 Discussion). A DFA for $\mathcal{L} = \{w : w \text{ contains } 1101\}$ is given by:



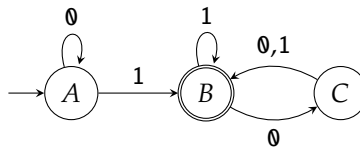
1.4 Formal definitions

Definition 1.4.1

A DFA is a tuple $(Q, \Sigma, \delta, q_0, F)$ where

- Q = set of states,
- Σ = alphabet,
- δ = transition function ($\delta: Q \times \Sigma \rightarrow Q$),
- q_0 = start state ($q_0 \in Q$), and
- F = set of accept states ("favorable"? states, $F \subseteq Q$).

Example 1.4.2. A formal description of the DFA



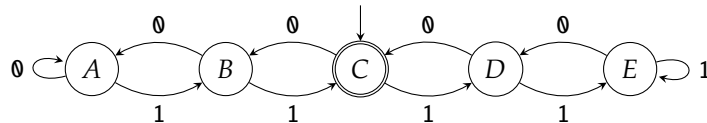
is given by $(\{A, B, C\}, \{0, 1\}, \delta, A, \{B\})$ where δ is defined by the table

	0	1
A	A	B
B	C	B
C	B	B

Example 1.4.3. The graph for the DFA $(\{A, B, C, D, E\}, \{0, 1\}, \delta, C, \{C\})$ where δ is defined by the table

	0	1
A	A	B
B	A	C
C	B	D
D	C	E
E	D	E

is given by:



Example 1.4.4. A formal description for Example 1.3.6 is given by $(\{0, 1, 2, \dots, 999\}, \{0, 1\}, \delta, 0, \{0\})$ where $\delta(q, \sigma) = (q + 1) \bmod 1000$.

Definition 1.4.5

DFA $(Q, \Sigma, \delta, q_0, F)$ **accepts** a string $w = w_1 w_2 \dots w_n$ iff

$$\delta(\dots \delta(\delta(q_0, w_1), w_2) \dots, w_n) \in F.$$

Definition 1.4.6

DFA D **recognizes** the language \mathcal{L} iff

$$\mathcal{L} = \{w : D \text{ accepts } w\}.$$

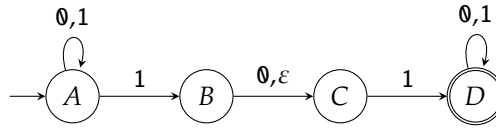
Note.

- Every DFA recognizes exactly 1 language.
- A language has either 0 or ∞ DFAs recognizing it.

2 Nondeterminism

2.1 Basic notions

Example 2.1.1.



- Choose an alphabet: $\{0, 1\}$.
- Draw states.
- Choose start state and accept states. The steps so far are the same as those of a DFA.
- Draw transitions. A state may have any number of transitions on a given symbol. A state may also have transitions on ϵ .

Definition 2.1.2

An NFA **accepts** w iff there is *at least* one path to an accept state.

Example 2.1.3. The output table for Example 2.1.1 is given by

Input	Accepting path	Output
ϵ	-	reject
0	-	reject
1	-	reject
010110	AABCDDDD	accept
010	-	reject
11	ABCD	accept.

Then the language is the set of all strings containing 101 or 11.

2.2 Using shortcuts

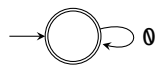
Example 2.2.1. An NFA for $\mathcal{L} = \emptyset$ is given by:



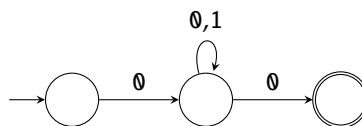
Example 2.2.2. An NFA for $\mathcal{L} = \{\varepsilon\}$ is given by:



Example 2.2.3. An NFA for $\mathcal{L} = \{w : w \text{ doesn't contain } 1\}$ is given by:

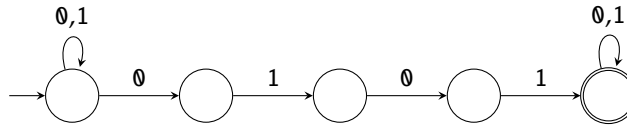


Example 2.2.4. An NFA for $\mathcal{L} = \{w : |w| \geq 2 \text{ and } w \text{ starts and ends with } 0\}$ is given by:

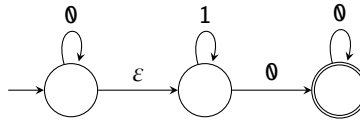


2.3 Pattern matching

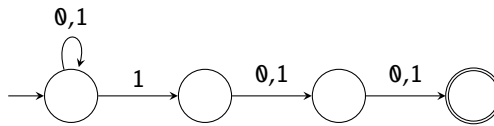
Example 2.3.1. An NFA for $\mathcal{L} = \{w : \text{contains } 0101\}$ is given by:



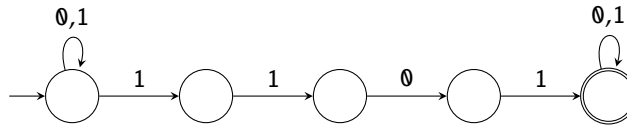
Example 2.3.2. An NFA for $\mathcal{L} = \{w : w = \underbrace{00\dots0}_{\geq 0} \underbrace{11\dots1}_{\geq 0} \underbrace{00\dots0}_{\geq 1}\}$ is given by:



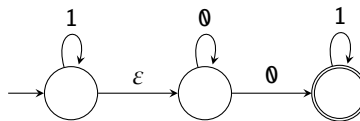
Example 2.3.3. An NFA for $\mathcal{L} = \{w : w \text{ has a 1 in the 3rd position from the end}\}$ is given by:



Example 2.3.4 (Week 1 Discussion). An NFA for $\mathcal{L} = \{w : w \text{ contains } 1101\}$ is given by:

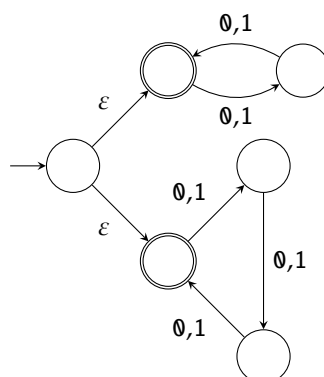


Example 2.3.5 (Week 1 Discussion). An NFA for $\mathcal{L} = \{w : w = \underbrace{11\dots1}_{\geq 0} \underbrace{00\dots0}_{\geq 1} \underbrace{11\dots1}_{\geq 0}\}$ is given by:

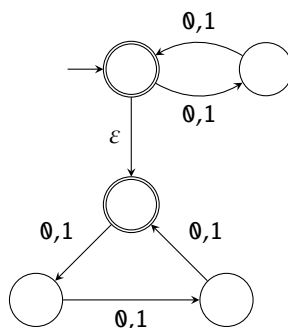


2.4 Alternatives

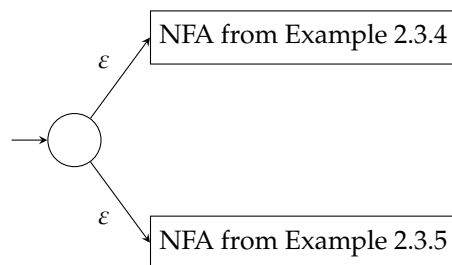
Example 2.4.1. An NFA for $\mathcal{L} = \{w : 2 \mid |w| \vee 3 \mid |w|\}$ is given by:



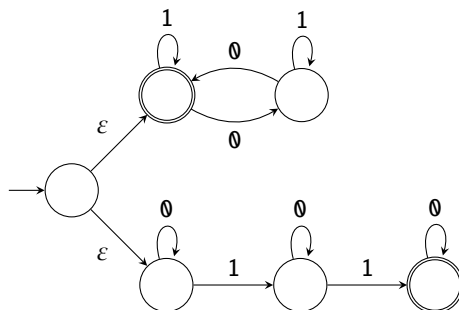
Note that the following is not valid due to side effects:



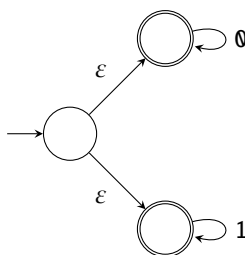
Example 2.4.2 (Week 1 Discussion). An NFA for $\mathcal{L} = \{w : w \text{ contains } 1101 \vee w = \underbrace{11\dots1}_{\geq 0} \underbrace{00\dots0}_{\geq 1} \underbrace{11\dots1}_{\geq 0}\}$ is given by:



Example 2.4.3. An NFA for $\mathcal{L} = \{w : w \text{ contains an even number of } 0\text{'s, or exactly two } 1\text{'s}\}$ is given by:



Example 2.4.4. An NFA for $\mathcal{L} = \{w : w \text{ does not contain both } 0 \text{ and } 1\}$ is given by:



2.5 Formal definitions

Definition 2.5.1

An NFA is a tuple $(Q, \Sigma, \delta, q_0, F)$ where

- Q = set of states,
- Σ = alphabet,
- δ = transition function ($\delta: Q \times (\Sigma \cup \{\varepsilon\}) \rightarrow \mathcal{P}(Q)$),
- q_0 = start state ($q_0 \in Q$), and
- F = set of accept states ($F \subseteq Q$).

Example 2.5.2. A formal description for NFA from Example 2.1.1 is given by

$$(\{A, B, C, D\}, \{\emptyset, 1\}, \delta, A, \{D\})$$

where δ is defined by the table

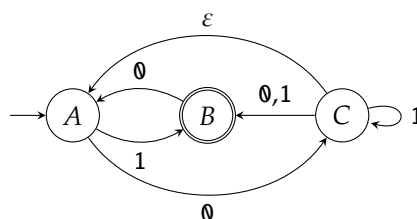
	\emptyset	1	ε
A	$\{A\}$	$\{A, B\}$	\emptyset
B	$\{C\}$	\emptyset	$\{C\}$
C	\emptyset	$\{D\}$	\emptyset
D	$\{D\}$	$\{D\}$	\emptyset .

Note that as pictures are not precise, we had to make an assumption on the alphabet of the NFA. Also note that although adding transition from A to A on ε does not change the language, it does not match the drawing and therefore represents a different NFA. Make sure to transcribe the given NFA.

Example 2.5.3. The graph for the DFA $(\{A, B, C\}, \{\emptyset, 1\}, \delta, A, \{B\})$ where δ is defined by the table

	\emptyset	1	ε
A	$\{C\}$	$\{B\}$	\emptyset
B	$\{A\}$	\emptyset	\emptyset
C	$\{B\}$	$\{B, C\}$	$\{A\}$

is given by



Definition 2.5.4

NFA $(Q, \Sigma, \delta, q_0, F)$ **accepts** a string w iff

$$\begin{aligned} \exists q_1, q_2, \dots, q_m \in Q \exists \sigma_0, \sigma_1, \dots, \sigma_{m-1} \in \Sigma \cup \{\varepsilon\} : \\ q_1 \in \delta(q_0, \sigma_0) \wedge q_2 \in \delta(q_1, \sigma_1) \wedge \dots \\ \wedge q_m \in \delta(q_{m-1}, \sigma_{m-1}) \wedge q_m \in F \\ \wedge \sigma_0 \sigma_1 \dots \sigma_{m-1} = w, \end{aligned}$$

or, in words, an accept state is reachable from q_0 via some path on input w .

Definition 2.5.5

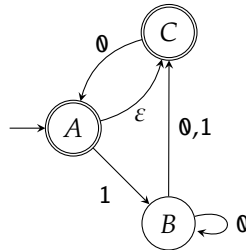
NFA N **recognizes** the language \mathcal{L} iff

$$\mathcal{L} = \{w : N \text{ accepts } w\}.$$

3 Equivalence of DFAs and NFAs

3.1 Example

Does this NFA accept the string 0110?



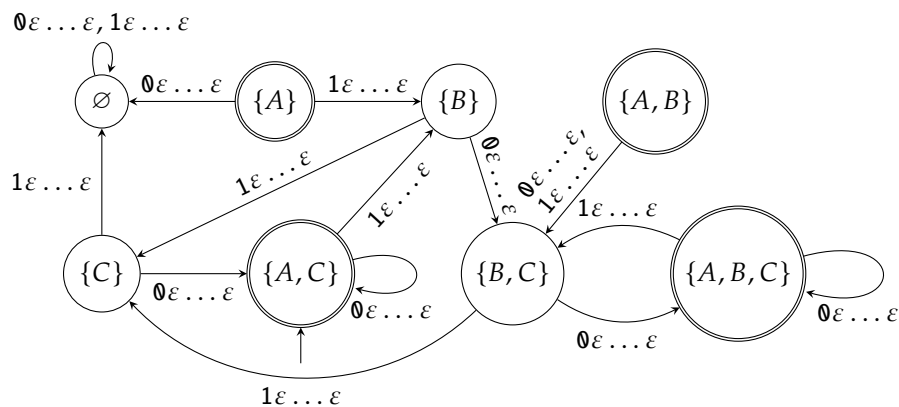
Recall that the NFA accepts 0110 iff an accept state is reachable from the start state via some path $\underbrace{\varepsilon \dots \varepsilon}_{\geq 0} \underbrace{0 \varepsilon \dots \varepsilon}_{\geq 0} \underbrace{1 \varepsilon \dots \varepsilon}_{\geq 0} \underbrace{1 \varepsilon \dots \varepsilon}_{\geq 0} \underbrace{0 \varepsilon \dots \varepsilon}_{\geq 0}$.

Path	Possible end states
$\varepsilon \dots \varepsilon$	A, C
$\varepsilon \dots \varepsilon 0 \varepsilon \dots \varepsilon$	A, C
$\varepsilon \dots \varepsilon 0 \varepsilon \dots \varepsilon 1 \varepsilon \dots \varepsilon$	B
$\varepsilon \dots \varepsilon 0 \varepsilon \dots \varepsilon 1 \varepsilon \dots \varepsilon 1 \varepsilon \dots \varepsilon$	C
$\varepsilon \dots \varepsilon 0 \varepsilon \dots \varepsilon 1 \varepsilon \dots \varepsilon 1 \varepsilon \dots \varepsilon 0 \varepsilon \dots \varepsilon$	A, C

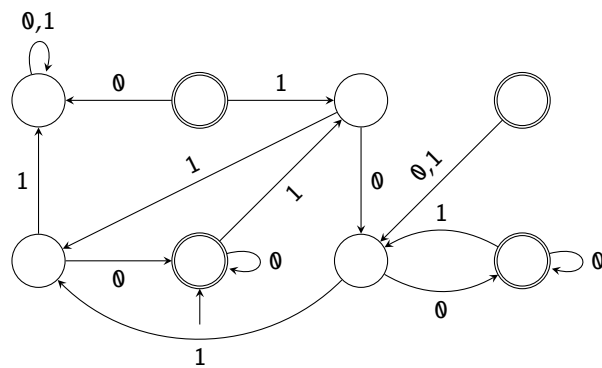
Since C is an accept state, the NFA accepts the string 0110.

Can we convert this to a DFA? Yes.

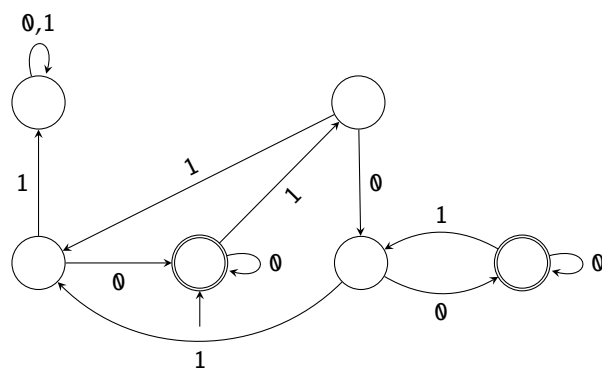
Step 1 Use subsets of the states of the NFA as the states of the DFA. Accept states are subsets that contain accept states of the NFA. The start state is the subset of all reachable states from the start state of the NFA via ε . For subsets containing more than one element, the transition is the union of all transitions for each individual element.



Step 2 Clean up.



Step 3 Optionally drop unreachable states.



3.2 General theorem

Theorem 3.2.1

Every NFA N can be converted to a DFA D that recognizes the same language.

Proof. Given NFA $N = (Q, \Sigma, \delta, q_0, F)$ define DFA $D = (\mathcal{P}(Q), \Sigma, \Delta, S_0, \mathcal{F})$ where

- $S_0 = \{q \in Q : q \text{ is reachable from } q_0 \text{ via a path } \underbrace{\varepsilon \dots \varepsilon}_{\geq 0}\}$,
- $\Delta(S, \sigma) = \{q \in Q : q \text{ is reachable from a state in } S \text{ via a path } \delta \underbrace{\varepsilon \dots \varepsilon}_{\geq 0}\}$, and
- $\mathcal{F} = \{S \subseteq Q : S \text{ contains a state in } F\}$.

Then N accepts a string $w = w_1 w_2 \dots w_n \Leftrightarrow$ a state in F is reachable via a path $\underbrace{\varepsilon \dots \varepsilon}_{\geq 0} w_1 \underbrace{\varepsilon \dots \varepsilon}_{\geq 0} w_2 \dots w_n \underbrace{\varepsilon \dots \varepsilon}_{\geq 0}$
 \Leftrightarrow a state in \mathcal{F} is reachable from S_0 via the path $w_1 w_2 \dots w_n \Leftrightarrow D$ accepts w . \square

3.3 Blow-up in size

Definition 3.3.1

A language \mathcal{L} is called **regular** iff it is recognized by some NFA or, equivalently, some DFA.

Last time: If \mathcal{L} has an NFA with k states then \mathcal{L} has a DFA with 2^k states.

We will show that this conclusion is the best possible.

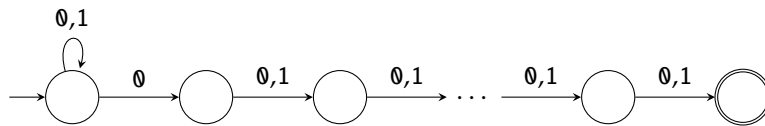
Theorem 3.3.2

Define $\Sigma = \{0, 1\}$, $k = \text{any positive integer}$, $\mathcal{L}_k = \{w : w \text{ has a } 0 \text{ in the } k\text{-th position from the end}\}$.

- \mathcal{L}_k has a NFA of size $k + 1$, and
- \mathcal{L}_k does not have a DFA of size $< 2^k$.

Proof.

(a)



(b) Let D be any DFA with $< 2^k$ states.

Input	End state of D
$000\dots 00$	q_i
$000\dots 01$	q_j
\vdots	\vdots
$111\dots 11$	q_k

If the left hand side of the table enumerates all 2^k strings of length k , then the right hand side cannot be all distinct due to the pigeonhole principle. Then there exists strings $u \neq v$ of length k such that

$$\text{end state on } u = \text{end state on } v.$$

Say $u_i = 0, v_i = 1$. Then

$$\text{end state on } u \underbrace{00\dots 0}_{i-1} = \text{end state on } v \underbrace{00\dots 0}_{i-1}.$$

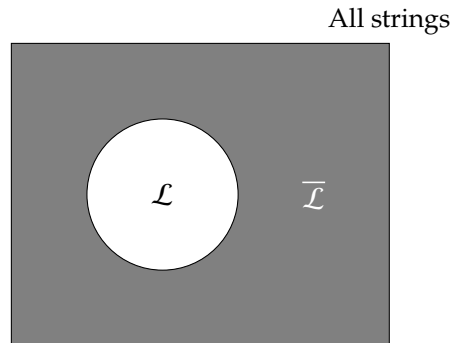
However, note that the left hand side has a 0 in the k -th position from the end while the right hand side has a 1 in the k -th position from the end.

Then we conclude that D does not recognize \mathcal{L}_k .

□

4 Closure

4.1 Complement



Definition 4.1.1

The **complement** of the language \mathcal{L} is

$$\overline{\mathcal{L}} = \{w : w \notin \mathcal{L}\}.$$

Theorem 4.1.2

If \mathcal{L} is regular then $\overline{\mathcal{L}}$ is regular. In other words, regular languages are closed under complement.

Proof. Swap the accept and reject states of the DFA for \mathcal{L} to get a DFA for $\overline{\mathcal{L}}$.

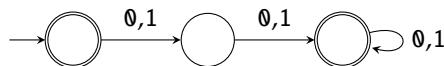
Formally, if DFA $(Q, \Sigma, \delta, q_0, F)$ recognizes \mathcal{L} then DFA $(Q, \Sigma, \delta, q_0, Q \setminus F)$ recognizes $\overline{\mathcal{L}}$. □

Example 4.1.3. $\mathcal{L} = \{w : w \neq 0, 1\}$.

Note that $\overline{\mathcal{L}} = \{w : w = 0 \vee w = 1\}$. The DFA for $\overline{\mathcal{L}}$ is given by:

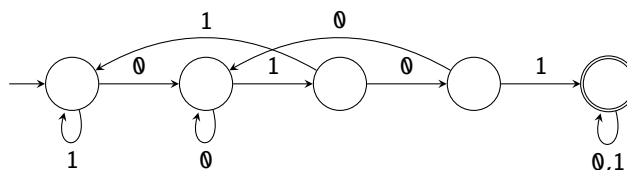


Then the DFA for \mathcal{L} is given by

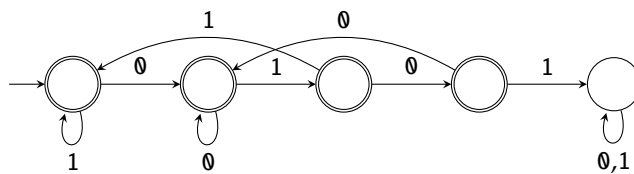


Example 4.1.4. $\mathcal{L} = \{w : w \text{ does not contain } 0101\}$.

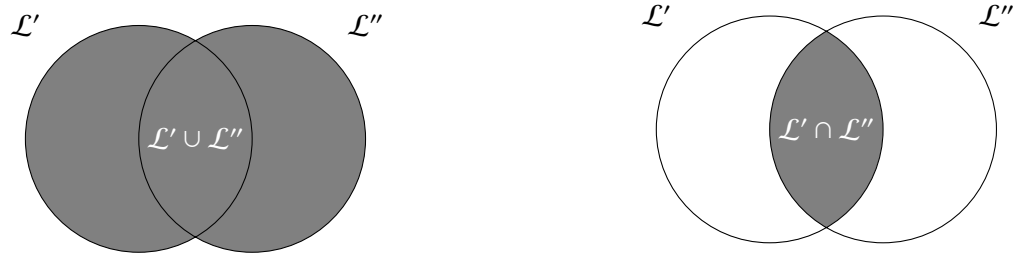
Note that $\overline{\mathcal{L}} = \{w : w \text{ contains } 0101\}$. The DFA for $\overline{\mathcal{L}}$ is given by:



Then the DFA for \mathcal{L} is given by



4.2 Union and intersection



Definition 4.2.1

The **union** of languages \mathcal{L}' and \mathcal{L}'' is

$$\mathcal{L}' \cup \mathcal{L}'' = \{w : w \in \mathcal{L}' \vee w \in \mathcal{L}''\}.$$

Definition 4.2.2

The **intersection** of languages \mathcal{L}' and \mathcal{L}'' is

$$\mathcal{L}' \cap \mathcal{L}'' = \{w : w \in \mathcal{L}' \wedge w \in \mathcal{L}''\}.$$

Theorem 4.2.3

If \mathcal{L}' , \mathcal{L}'' are regular, then so are

- (a) $\mathcal{L}' \cup \mathcal{L}''$ and
- (b) $\mathcal{L}' \cap \mathcal{L}''$.

In other words, regular languages are closed under union and intersection.

Proof.

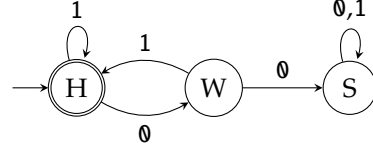
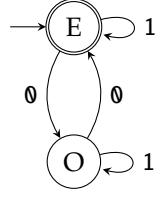
- (a) Connect an initial state to the initial states for the NFAs for \mathcal{L}' and \mathcal{L}'' via ε .
- (b) If \mathcal{L}' and \mathcal{L}'' are regular then $\overline{\mathcal{L}'}$ and $\overline{\mathcal{L}''}$ are regular. Then $\overline{\mathcal{L}' \cap \mathcal{L}''}$ is regular. Then, by De Morgan's law, $\mathcal{L}' \cap \mathcal{L}'' = \overline{\overline{\mathcal{L}'} \cup \overline{\mathcal{L}''}}$ is regular.

□

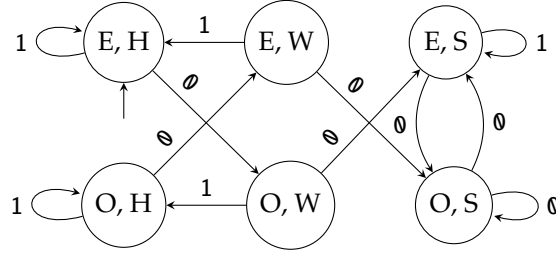
Proof. **Idea:** ¹

$\mathcal{L}' = \{w : w \text{ has an even number of } 0\text{'s}\}.$

$\mathcal{L}'' = \{w : \text{every } 0 \text{ is immediately followed by a } 1\}.$

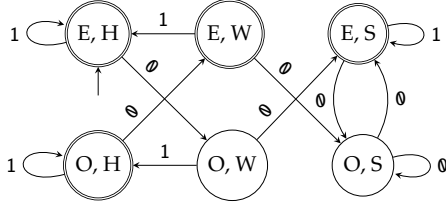


Consider the construction:

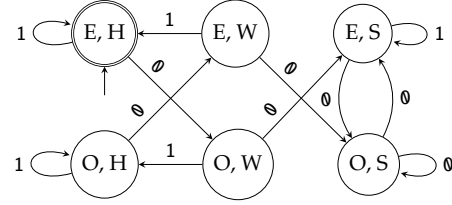


Note that we can create DFAs for either union or intersection by assigning appropriate accept states:

$\mathcal{L} = \mathcal{L}' \cup \mathcal{L}''.$



$\mathcal{L} = \mathcal{L}' \cap \mathcal{L}''.$



Note that this construction is better than that in the first proof since the number of nodes in this construction is $O(n^2)$ while that in the first proof is $O(4^n)$.

Formally, write the DFA for \mathcal{L}' as $(Q', \Sigma, \delta', q'_0, F')$ and that for \mathcal{L}'' as $(Q'', \Sigma, \delta'', q''_0, F'')$.

Then the DFA for $\mathcal{L}' \cup \mathcal{L}''$ is given by $(Q' \times Q'', \Sigma, \delta, (q'_0, q''_0), (F' \times Q'') \cup (Q' \times F''))$ where $\delta((q', q''), \sigma) := (\delta'(q', \sigma), \delta''(q'', \sigma))$ for all $q' \in Q'$, $q'' \in Q''$, and $\sigma \in \Sigma$.

Similarly, the DFA for $\mathcal{L}' \cap \mathcal{L}''$ is given by $(Q' \times Q'', \Sigma, \delta, (q'_0, q''_0), F' \times F'')$ where $\delta((q', q''), \sigma) := (\delta'(q', \sigma), \delta''(q'', \sigma))$ for all $q' \in Q'$, $q'' \in Q''$, and $\sigma \in \Sigma$. □

Corollary 4.2.4

If $\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_n$ are regular, so are

- $\mathcal{L}_1 \cup \mathcal{L}_2 \cup \dots \cup \mathcal{L}_n$ and
- $\mathcal{L}_1 \cap \mathcal{L}_2 \cap \dots \cap \mathcal{L}_n$.

Proof. Induction. □

¹Here “E” stands for “EVEN”, “O” stands for “ODD”, “H” stands for “HAPPY”, “W” stands for “WORRIED”, and “S” stands for “SAD”.

Corollary 4.2.5

Every finite language is regular.

Proof. Let \mathcal{L} be finite. Then we can write

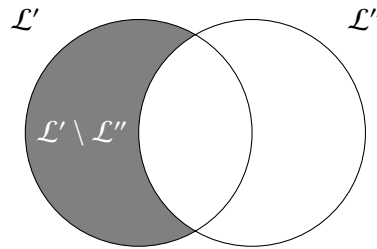
$$\mathcal{L} = \bigcup_{w \in \mathcal{L}} \{w\},$$

that is, as a union of singleton languages.

Note that an NFA can be easily found for any singleton language, that is, singleton languages are regular. Then \mathcal{L} is regular by Corollary 4.2.4. \square

4.3 Difference

$$\mathcal{L}' \setminus \mathcal{L}'' = \{w \in \mathcal{L}' : w \notin \mathcal{L}''\}$$



Theorem 4.3.1

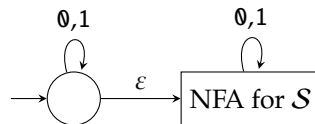
If \mathcal{L}' and \mathcal{L}'' are regular then so is $\mathcal{L}' \setminus \mathcal{L}''$.

Proof. $\mathcal{L}' \setminus \mathcal{L}'' = \mathcal{L}' \cap \overline{\mathcal{L}''}$. □

Corollary 4.3.2

If \mathcal{A}, \mathcal{S} are regular, then so is $\mathcal{L} = \{w \in \mathcal{A} : \text{no substring of } w \text{ is in } \mathcal{S}\}$.

Proof. Let $\mathcal{S}' = \{w : \text{some substring of } w \text{ is in } \mathcal{S}\}$. Note that \mathcal{S}' is a superset of \mathcal{S} . An NFA for \mathcal{S}' is given by:



where the self loop on NFA for \mathcal{S} refers to self loops on the accept states of the NFA.

Since \mathcal{S}' is regular by construction, $\mathcal{L} = \mathcal{A} \setminus \mathcal{S}'$ is regular. □

4.4 Concatenation

Definition 4.4.1

The **concatenation** of two languages \mathcal{L}' and \mathcal{L}'' , is given by

$$\mathcal{L}'\mathcal{L}'' = \{w'w'' : w' \in \mathcal{L}' \quad \wedge \quad w'' \in \mathcal{L}''\}.$$

This is sometimes denoted as $\mathcal{L}' \circ \mathcal{L}''$.

Example 4.4.2. Given $\mathcal{L}' = \{000, 111\}$ and $\mathcal{L}'' = \{10, 010\}$, we have

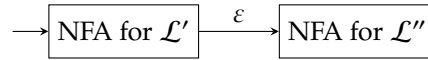
$$\mathcal{L}'\mathcal{L}'' = \{00010, 000010, 11110, 111010\}.$$

Observation: $\mathcal{L}\emptyset = \emptyset\mathcal{L} = \emptyset$ and $\mathcal{L}\{\varepsilon\} = \{\varepsilon\}\mathcal{L} = \mathcal{L}$ for any language \mathcal{L} .

Theorem 4.4.3

If \mathcal{L}' and \mathcal{L}'' are regular then so is $\mathcal{L}'\mathcal{L}''$.

Proof. An NFA for $\mathcal{L}'\mathcal{L}''$ is given by:



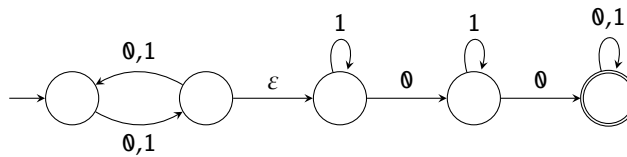
where the ε transition from the first NFA to the second NFA refer to those from the accept states in the first NFA to the start state of the second and all accept states in the first NFA are converted to reject states.

Formally, write the first NFA as $(Q', \Sigma, \delta', q'_0, F')$ and the second NFA as $(Q'', \Sigma, \delta'', q''_0, F'')$. Then we

- make every state in F' a reject state and
- add ε transitions from every state in F' to q''_0 .

□

Example 4.4.4. An NFA for $\{uw : |u| \text{ is odd, } v \text{ contains at least two } 0\text{'s}\}$ is given by:



4.5 Kleene star

Definition 4.5.1

The **Kleene star** of a language \mathcal{L} is given by

$$\mathcal{L}^* = \{w : w \text{ is the concatenation of zero or more strings in } \mathcal{L}\}.$$

Example 4.5.2. The Kleene star of $\mathcal{L} = \{a, b\}$ is given by

$$\mathcal{L}^* = \{\varepsilon, a, b, aa, ab, ba, bb, \dots\}.$$

In words, $\{a, b\}^*$ = the set of all strings over the alphabet $\{a, b\}$.

Observation: $\emptyset^* = \{\varepsilon\}^* = \{\varepsilon\}$.

Observation: \mathcal{L}^* is finite if and only if $\mathcal{L} = \emptyset$ or $\mathcal{L} = \{\varepsilon\}$.

Theorem 4.5.3

If \mathcal{L} is regular then so is \mathcal{L}^* .

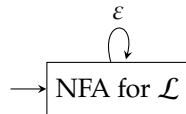
Proof. **(Incorrect)** Note that

$$\mathcal{L}^* = \{\varepsilon\} \cup \mathcal{L} \cup \mathcal{L}\mathcal{L} \cup \dots \cup \underbrace{\mathcal{L} \dots \mathcal{L}}_n \cup \dots,$$

and since the union of regular languages is regular, \mathcal{L}^* is regular.

However, note that the theorem that the union of regular languages is regular does not work for an infinite number of languages and therefore this proof fails. \square

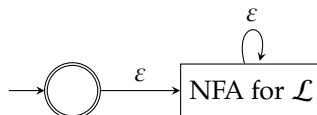
Proof. **(Incorrect)** An NFA for \mathcal{L}^* is given by



where the self loop on the NFA refers to ε transitions from the accept states of the NFA to the start state and the start state is converted to an accept state.

However, this construction fails if q_0 is originally rejecting *and* part of a cycle. \square

Proof. **(Correct)** An NFA for \mathcal{L}^* is given by



where the self loop on the NFA refers to ε transitions from the accept states of the NFA to the start state.

In words, add ε -transitions from every state in F to q_0 and add a new start state q'_0 and link it with an ε -arrow to q_0 . \square

4.6 Reverse

Definition 4.6.1

For a string $w = w_1w_2 \dots w_n$, its **reverse** is given by

$$w^R = w_nw_{n-1} \dots w_2w_1.$$

Definition 4.6.2

The **reverse** of a language \mathcal{L} is given by

$$\mathcal{L}^R = \{w^R : w \in \mathcal{L}\}.$$

Example 4.6.3. $\{011, abc\}^R = \{110, cba\}.$

Theorem 4.6.4

If \mathcal{L} is regular then so is \mathcal{L}^R .

Proof. Consider an NFA $(Q, \Sigma, \delta, q_0, F)$ for \mathcal{L} . Reverse all the transitions, create a new start state, and link it with ε -transitions to states in F . Make q_0 the only accept state. \square

4.7 Prefix and suffix

Definition 4.7.1

The **prefix language** of a language \mathcal{L} is given by

$$\text{prefix}(\mathcal{L}) = \{u : uv \in \mathcal{L} \text{ for some string } v, \text{ possibly empty}\}.$$

Definition 4.7.2

The **suffix language** of a language \mathcal{L} is given by

$$\text{suffix}(\mathcal{L}) = \{v : uv \in \mathcal{L} \text{ for some string } u, \text{ possibly empty}\}.$$

Theorem 4.7.3

If \mathcal{L} is regular then so is $\text{prefix}(\mathcal{L})$.

Proof. Given a DFA $(Q, \Sigma, \delta, q_0, F)$ for \mathcal{L} , a DFA for $\text{prefix}(\mathcal{L})$ is given by $(Q, \Sigma, \delta, q_0, F')$ where

$$F' = \{q \in Q : \text{there exists a path, possibly empty, from } q \text{ to a state in } F\}$$

□

Corollary 4.7.4

If \mathcal{L} is regular then so is $\text{suffix}(\mathcal{L})$.

Proof. $\text{suffix}(\mathcal{L}) = (\text{prefix}(\mathcal{L}^R))^R$.

□

4.8 Substring

Definition 4.8.1

The **substring language** of a language \mathcal{L} is given by

$$\text{substring}(\mathcal{L}) = \{v : uvw \in \mathcal{L} \text{ for some strings } u, w\}.$$

Theorem 4.8.2

If \mathcal{L} is regular then so is $\text{substring}(\mathcal{L})$.

Proof. $\text{substring}(\mathcal{L}) = \text{suffix}(\text{prefix}(\mathcal{L}))$. □

4.9 Shuffle

Definition 4.9.1

The **shuffle language** of two languages \mathcal{L}' and \mathcal{L}'' is given by

$$\{a_1 b_1 a_2 b_2 \dots a_k b_k : a_1, \dots, a_k \text{ strings}, b_1, \dots, b_k \text{ strings}, a_1 \dots a_k \in \mathcal{L}', b_1 \dots b_k \in \mathcal{L}''\}.$$

Proof. Given an NFA $(Q', \Sigma, \delta', q'_0, F')$ for \mathcal{L}' and an NFA $(Q'', \Sigma, \delta'', q''_0, F'')$ for \mathcal{L}'' , an NFA for $\mathcal{L}' \diamond \mathcal{L}''$ is given by $(Q' \times Q'', \Sigma, \Delta, (q'_0, q''_0), F' \times F'')$ where

$$\Delta((q', q''), \sigma) = (\delta'(q', \sigma) \times \{q''\}) \cup (\{q'\} \times \delta''(q'', \sigma)).$$

□

5 Regular expressions

5.1 Basic notions

Definition 5.1.1

A regular expression is

- \emptyset ,
- $\{\varepsilon\}$,
- $\{\sigma\}$ for $\sigma \in \Sigma$,
- (R^*) ,
- $(R_1 \cup R_2)$, and
- $(R_1 \circ R_2)$.

Example 5.1.2. $((\{\emptyset\} \cup \{1\})^* \cup \{1\})$.

Example 5.1.3. $((\{\emptyset\} \circ \{1\}) \circ \{1\}) \circ \{1\}$.

Example 5.1.4. $((\{\emptyset\} \circ \{1\})^* \circ \{1\}) \circ ((\{\emptyset\} \cup \{1\})^*)$

Shorthand

- $R^k = \underbrace{(R \circ R \circ \dots \circ R)}_k$
- $R^+ = (R \circ (R^*))$
- $\Sigma = (\{\emptyset\} \cup \{1\})$ assuming alpha bet is $\{\emptyset, 1\}$
- may omit $\{\}, (), \circ$

Precedence: $*$ then \circ then \cup .

Note. Sometimes $()$ are necessary.

Example 5.1.5. Example 5.1.2 can be simplified to $\Sigma^* \cup \{1\} = \Sigma^*$.

Example 5.1.6. Example 5.1.3 can be simplified to $\emptyset 111$.

Example 5.1.7. Example 5.1.4 can be simplified to $\Sigma^* 1 \Sigma^*$.

Example 5.1.8. A regular expression for $\{w : w \text{ contains exactly one } 1\}$ is $\emptyset^* 1 \emptyset^*$.

Example 5.1.9. A regular expression for $\{w : w \text{ contains at least one } 1\}$ is $\Sigma^* 1 \Sigma^*$.

Example 5.1.10. A regular expression for $\{w : w \text{ contains } \emptyset \emptyset 1 \text{ as a substring}\}$ is $\Sigma^* \emptyset \emptyset 1 \Sigma^*$.

Example 5.1.11. A regular expression for $\{w : w \text{ does not contain } \emptyset \emptyset \text{ as a substring}\}$ is $1^*(\emptyset 1^+)^*(\varepsilon \cup \emptyset)$.

Example 5.1.12. A regular expression for $\{w : |w| \text{ is even}\}$ is $(\Sigma^2)^*$.

Example 5.1.13. A regular expression for $\{w : 5 \mid |w|\}$ is $(\Sigma^5)^*$.

Example 5.1.14. A regular expression for $\{w : w \neq \varepsilon \wedge w_1 = w_{|w|}\}$ is $\Sigma \cup \emptyset \Sigma^* \emptyset \cup 1 \Sigma^* 1$.

Example 5.1.15. A regular expression for $\{\emptyset 1, 1 \emptyset\}$ is $\emptyset 1 \cup 1 \emptyset$.

Example 5.1.16. The regular expression $(\emptyset \cup \varepsilon) 1^*$ can be arguably simplified to $\emptyset 1^* \cup 1^*$.

Example 5.1.17. The regular expression $1^*\emptyset$ can be simplified to \emptyset .

Example 5.1.18. The regular expression \emptyset^* can be simplified to ε .

Example 5.1.19. The regular expression $(\emptyset \cup 1^+\emptyset)^*$ can be simplified to $((\varepsilon \cup 1^+)\emptyset)^*$ and then to $(1^*\emptyset)^*$. In words, all strings that do not end in a 1.

Example 5.1.20.

Regular expression	Shortest accepted string	Shortest rejected string
$\emptyset^*\emptyset\emptyset 11^*$	$\emptyset\emptyset 1$	ε
$\emptyset^*\emptyset\emptyset(\emptyset \cup 1)^*$	$\emptyset\emptyset$	ε
$(\emptyset\emptyset)^*(11\emptyset)^*(11)^*$	ε	$\emptyset, 1$
$\emptyset^*(\emptyset 1 \cup 1\emptyset)^*1^*$	ε	$1\emptyset\emptyset, 11\emptyset$

5.2 Equivalence with automata

Theorem 5.2.1

The language of every regular expression is regular.

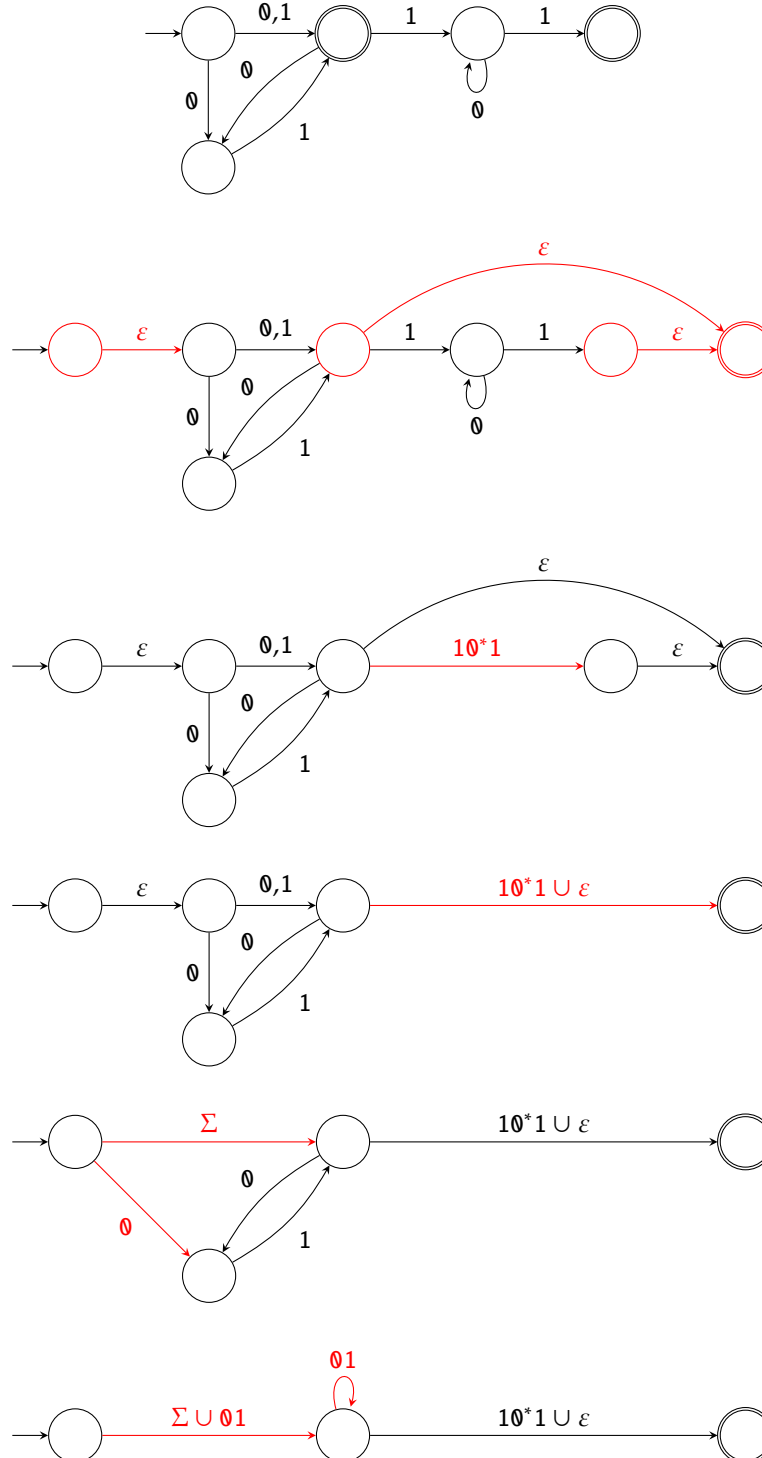
Proof. Induction. □

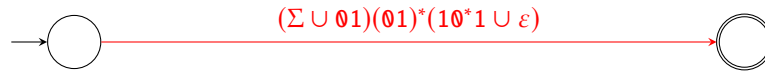
Theorem 5.2.2: Kleene's theorem

Then language of every NFA can be represented by a regular expression.

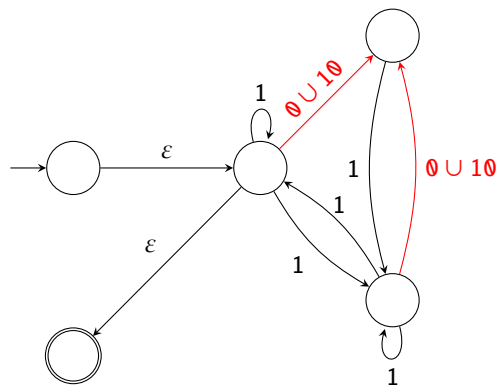
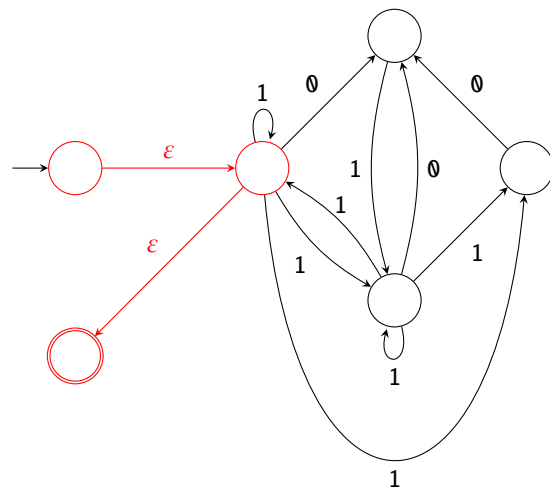
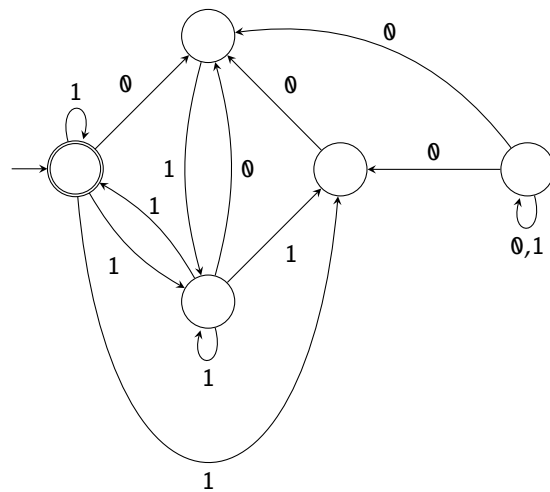
Proof. Given an NFA, create a new state and a new and unique accept state (making every state in the original NFA rejecting). Choose a state other than the start state and the accept state. Call the states from which the incoming transitions originate donors and those at which the outgoing transitions arrive receivers. Note that a state can be both a donor and a receiver. Then for each donor-receiver pair, replace the transitions with a single transition $\alpha\beta^*\gamma$ where α is the incoming transition, β is the looping transition, and γ is the outgoing transition. Remove the chosen state. Collapse any parallel edges created in the process by unioning them. If there are states other than the start state or the accept state, repeat the process. If there are no transitions from the start state to the accept state then the regular expression is \emptyset . Otherwise there is exactly one transition α , the desired regular expression. □

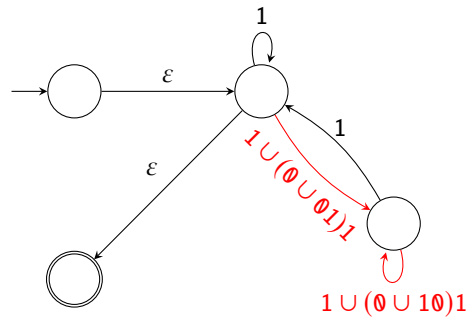
5.3 Conversion example



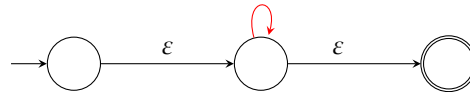


5.4 Conversion example #2

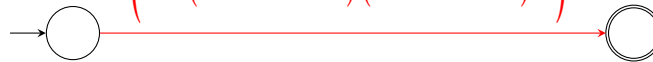




$$1 \cup (1 \cup (\emptyset \cup 1\emptyset)1)(1 \cup (\emptyset \cup 1\emptyset)1)^* 1$$



$$(1 \cup (1 \cup (\emptyset \cup 1\emptyset)1)(1 \cup (\emptyset \cup 1\emptyset)1)^* 1)^*$$



6 Pumping lemma

6.1 First nonregular language

Theorem 6.1.1

The language $\mathcal{L} = \{0^n 1^n : n \geq 1\}$ is not regular.

Proof. By contradiction. Suppose we have a DFA D for \mathcal{L} . Let p be the number of states in D .

- Run the DFA on $w = \underbrace{00\dots0}_p \underbrace{11\dots1}_p$.
- Note that w must end in an accept state.
- By the pigeonhole principle, some state q must have been visited at least twice while reading $\underbrace{00\dots0}_p$.
- Write $w = xyz$ where x is the portion of w before the first visit of q , y is the portion of w between the first and the second visit of q , and z is the portion of w after the second visit of q .
- Since D accepts $w = xyz$, D also accepts $xyyz$, which contains more than p 0's and exactly p 1's. Then D does not recognize \mathcal{L} .

□

6.2 Pumping lemma

Lemma 6.2.1

Let \mathcal{L} be a regular language. Then there exists $p \in \mathbb{N}$ such that every string $w \in \mathcal{L}$ of length $\geq p$ can be written as

$$w = xyz,$$

where

- x, y, z are strings,
- $y \neq \varepsilon$,
- $|xy| \leq p$, and
- $xy^iz \in \mathcal{L}$ for $i = 0, 1, 2, \dots$

To put in a slightly irresponsible wording, all long enough strings in \mathcal{L} can be pumped.

Proof. Let D be a DFA for \mathcal{L} . Define p to be the number of states in D .

- Run the DFA on w .
- Note that w must end in an accept state.
- By the pigeonhole principle, some state q must have been visited at twice while reading $w_1w_2 \dots w_p$.
- Define x to be the portion of w before the first visit of q , y to be the portion of w between the first and the second visit of q , and z to be the portion of w after the second visit of q .
- Note that x, y , and z are strings, $y \neq \varepsilon$, $|xy| \leq |w_1w_2 \dots w_p| = p$, and $xy^iz \in \mathcal{L}$ for all i .

□

Original pumping lemma ("pumping property")

$$\mathcal{L} \text{ reg} \Rightarrow (\exists p \forall w \in \mathcal{L} : |w| \geq p \Rightarrow \exists x, y, z : w = xyz \wedge y \neq \varepsilon \wedge |xy| \leq p \wedge \forall i : xy^iz \in \mathcal{L}).$$

Contrapositive form of pumping lemma

$$(\forall p \exists w \in \mathcal{L} : |w| \geq p \Rightarrow \forall x, y, z : w = xyz \wedge y \neq \varepsilon \wedge |xy| \leq p \Rightarrow \exists i : xy^iz \notin \mathcal{L}) \Rightarrow \mathcal{L} \text{ not regular}$$

6.3 Examples

Example 6.3.1. Consider the language $\mathcal{L} = \{0^n 1^n : n \geq 0\}$.

Let p be arbitrary. Define $w = 0^p 1^p \in \mathcal{L}$. Consider arbitrary x, y, z such that $w = xyz$, $y \neq \varepsilon$, and $|xy| \leq p$. Note that $xy^2z \notin \mathcal{L}$.

Then the pumping property fails and \mathcal{L} is not regular.

Example 6.3.2. Consider the language $\mathcal{L} = \{w : w \text{ contains as many } 0\text{'s as } 1\text{'s}\}$.

Note that

$$\mathcal{L} \cap \underbrace{0^* 1^*}_{\text{regular}} = \underbrace{\{0^n 1^n : n \geq 0\}}_{\text{nonregular}}.$$

Then by the closure of regular languages under intersection, \mathcal{L} is not regular.

Example 6.3.3. Consider the language $\mathcal{L} = \{0^n 1^m : n \neq m\}$.

Note that

$$\underbrace{0^* 1^*}_{\text{regular}} \setminus \mathcal{L} = \underbrace{\{0^n 1^n : n \geq 0\}}_{\text{nonregular}}.$$

Then by the closure of regular languages under set difference, \mathcal{L} is not regular.

Example 6.3.4. Consider the language $\mathcal{L} = \{0^n 1^m : n > m\}$ and the string $\overbrace{00 \dots 0}^{p+1} \overbrace{11 \dots 1}^p$. Consider any $y \neq \varepsilon$ such that $|xy| \leq p$. Note that $xz \notin \mathcal{L}$. Then \mathcal{L} is not regular.

Example 6.3.5. Consider the language $\mathcal{L} = \{10^n 1^n : n \geq 1\}$ and the string $1 \overbrace{00 \dots 0}^p \overbrace{11 \dots 1}^p$. Then $y \neq \varepsilon$ either contains a leading 1, in which case xz does not have a leading 1 and therefore $xz \notin \mathcal{L}$, or does not contain a leading 1, in which case xz has a different number of 0's as trailing 1's and therefore $xz \notin \mathcal{L}$. Then \mathcal{L} is not regular.

Example 6.3.6. Consider the language $\mathcal{L} = \{1^{n^2} : n \geq 0\}$ and the string $\overbrace{11 \dots 1}^{p^2}$. Consider any $y \neq \varepsilon$ such that $|xy| \leq p$. Then

$$\begin{aligned} p^2 &< |xy^2z| = |xyz| + |y| \\ &\leq p^2 + p < (p+1)^2. \end{aligned}$$

Then $xy^2z \notin \mathcal{L}$ and \mathcal{L} is not regular.

Example 6.3.7. Consider the language $\mathcal{L} = \{ww : w \in \{0, 1\}^*\}$ and the string $\overbrace{00 \dots 0}^p 1 \overbrace{00 \dots 0}^p 1$. Consider any $y \neq \varepsilon$ such that $|xy| \leq p$. Note that $xz \notin \mathcal{L}$ since the first half of xz ends with 0 while the second half ends with 1. Then \mathcal{L} is not regular.

Example 6.3.8. Consider the language $\mathcal{L} = \{www : w \in \{0, 1\}^*\}$. It can be shown that this is not regular.

Example 6.3.9. Consider the language $\mathcal{L} = \{w : w = w^R\}$ and the string $\underbrace{00 \dots 0}_p 1 \underbrace{00 \dots 0}_p$. Consider any $y \neq \varepsilon$ such that $|xy| \leq p$. Note that $xz \notin \mathcal{L}$. Then \mathcal{L} is not regular.

Example 6.3.10. Consider the language $\mathcal{L} = \{w : |w| \text{ prime}\}$. Let q be a prime greater than p . Consider the string $\underbrace{00 \dots 0}_q$. Consider any $y \neq \varepsilon$ such that $|xy| \leq p$. Then $xy^{q+1}z \notin \mathcal{L}$ since $|xy^{q+1}z| = \underbrace{q}_{\geq 2} \underbrace{(1 + |y|)}_{\geq 2}$ is not prime. Then \mathcal{L} is not regular.

6.4 True or false

Example 6.4.1. The statement “regular languages are closed under arbitrary unions, i.e., even infinite unions” is false. Note that every nonregular language is an infinite union of singleton languages, which are regular. For a concrete example, note that $\{0^n 1^n : n \geq 0\} = \bigcup_{n=0}^{\infty} \{0^n 1^n\}$.

Example 6.4.2. The statement “subsets of regular languages are regular” is false. Note that Σ^* is regular yet $\{0^n 1^n : n \geq 0\} \subseteq \Sigma^*$.

Example 6.4.3. The statement “subsets of nonregular languages are nonregular” is false. Note that \emptyset is regular yet $\emptyset \subseteq \{0^n 1^n : n \geq 0\}$.

Example 6.4.4. The statement “nonregular languages are closed under complement” is true since \mathcal{L} regular $\Leftrightarrow \overline{\mathcal{L}}$ regular, which implies \mathcal{L} nonregular $\Leftrightarrow \overline{\mathcal{L}}$ nonregular.

Example 6.4.5. The statement “nonregular languages are closed under union” is false. Note that $\mathcal{L} = \{0^n 1^n : n \geq 0\}$ is nonregular and so is $\overline{\mathcal{L}}$, yet $\mathcal{L} \cup \overline{\mathcal{L}} = \Sigma^*$ is regular.

Example 6.4.6. The statement “nonregular languages are closed under intersection” is false. Note that $\mathcal{L} = \{0^n 1^n : n \geq 0\}$ is nonregular and so is $\overline{\mathcal{L}}$, yet $\mathcal{L} \cap \overline{\mathcal{L}} = \emptyset$ is regular.

Example 6.4.7. The statement “nonregular languages are always infinite” is true since all finite languages are regular.

7 Myhill-Nerode theorem

7.1 When the pumping lemma fails

Example 7.1.1. Consider the language $\mathcal{L} = 0^*1^* \cup \{1^m 0^n 1^n : m \geq 1 \wedge n \geq 0\}$. We can consider strings of the forms 0^+ , 1^+ , 0^+1^+ , $1^+0^n1^n$, etc. We can check that all of the listed forms satisfy the pumping property by taking $x = \varepsilon$ and $y = w_1$. Then in this case the pumping lemma is useless. However \mathcal{L} is nonregular since $\mathcal{L} \cap 10^+1^+ = \{10^n1^n : n \geq 1\}$ is nonregular.

7.2 Indistinguishability

Definition 7.2.1

Let \mathcal{L} be a given language. Strings x and y are called \mathcal{L} -indistinguishable if for all strings w ,

$$xw \in \mathcal{L} \Leftrightarrow yw \in \mathcal{L}.$$

Notation: $x \equiv_{\mathcal{L}} y$

Example 7.2.2. Consider the language $\mathcal{L} = \{w : 3 \mid |w|\}$. Note that

- $1 \equiv_{\mathcal{L}} 0$,
- $\varepsilon \not\equiv_{\mathcal{L}} 0$,
- $\varepsilon \equiv_{\mathcal{L}} 000$, and
- $1 \equiv_{\mathcal{L}} 1010 \equiv_{\mathcal{L}} 1110001$.

Example 7.2.3. Consider the language $\mathcal{L} = \{w : 3 \mid \sum w_i\}$. Note that

- $1 \not\equiv_{\mathcal{L}} 0$,
- $\varepsilon \equiv_{\mathcal{L}} 0 \equiv_{\mathcal{L}} 000$,
- $1 \equiv_{\mathcal{L}} 101011$

Proposition 7.2.4

For any language \mathcal{L} , $\equiv_{\mathcal{L}}$ is an equivalence relation.

Proof.

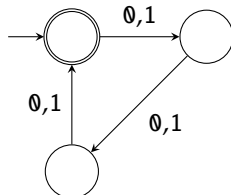
- Reflexive: $x \equiv_{\mathcal{L}} x$.
- Symmetric: $x \equiv_{\mathcal{L}} y \Leftrightarrow y \equiv_{\mathcal{L}} x$.
- Transitive: $xw \in \mathcal{L} \Leftrightarrow yw \in \mathcal{L} \Leftrightarrow zw \in \mathcal{L}$.

□

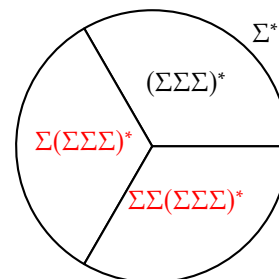
Note. Every equivalence class is either entirely in \mathcal{L} or entirely in $\overline{\mathcal{L}}$.

Example 7.2.5. Consider the language $\mathcal{L} = \{w : 3 \mid |w|\}$.

A DFA for \mathcal{L} is given by:

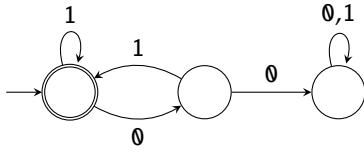


Then the language can be partitioned into:

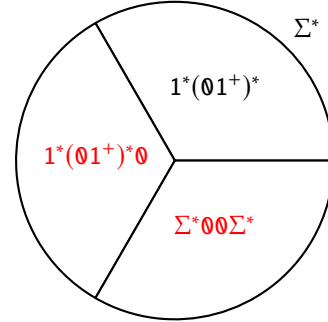


Example 7.2.6. Consider the language $\mathcal{L} = \{w : \text{every } 0 \text{ is immediately followed by a } 1\}$.

A DFA for \mathcal{L} is given by:

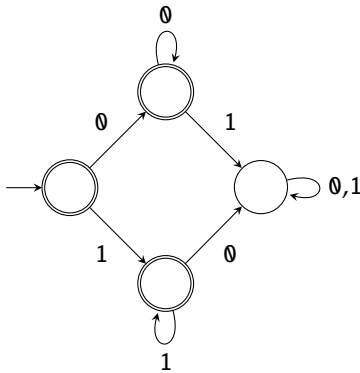


Then the language can be partitioned into:

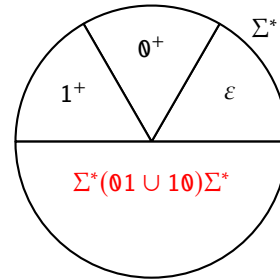


Example 7.2.7. Consider the language $\mathcal{L} = \{w : w \text{ does not contain both a } 0 \text{ and a } 1\}$.

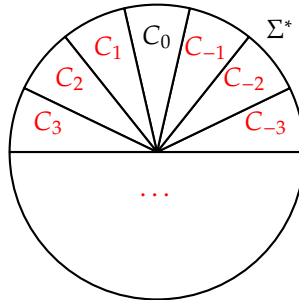
A DFA for \mathcal{L} is given by:



Then the language can be partitioned into:



Example 7.2.8. Consider the language $\mathcal{L} = \{w : w \text{ contains as many } 0\text{'s as } 1\text{'s}\}$. Note that this language is nonregular. The language can be partitioned into:



where

$$C_i = \{w : \text{number of } 0\text{'s in } w - \text{number of } 1\text{'s in } w = i\}.$$

Example 7.2.9. Consider the language $\mathcal{L} = \{w : w^R = w\}$. Then

$$w \neq w' \Rightarrow w \mathbf{0^n 1 0^n} w^R \in \mathcal{L} \wedge w' \mathbf{0^n 1 0^n} w^R \notin \mathcal{L}$$

for large enough n . Then every $w \in \Sigma^*$ is in its own equivalence class.

7.3 Myhill-Nerode theorem

Theorem 7.3.1

For every language \mathcal{L} , the minimum number of states in a DFA for \mathcal{L} is equal to the number of equivalence classes under $\equiv_{\mathcal{L}}$.

If this number is infinite, we interpret it as “no such DFA exists,” that is, \mathcal{L} is nonregular.

Proof. \geq : By contradiction.

Fix a DFA for \mathcal{L} with k states. Suppose $\exists x_1, x_2, \dots, x_{k+1}$, each in a different equivalence class of $\equiv_{\mathcal{L}}$.

Input string	End state
x_1	q_{27}
x_2	q_6
x_3	q_5
\vdots	\vdots
x_{k+1}	q_{12}

Note that the left hand side of the table enumerates $k + 1$ strings while the right hand side only has k possibilities. Then for some x_i, x_j where $i \neq j$, the DFA ends in the same state on reading x_i as x_j . Then for all strings w we have

$$x_i w \in \mathcal{L} \Leftrightarrow x_j w \in \mathcal{L},$$

that is,

$$x_i \equiv_{\mathcal{L}} x_j,$$

a contradiction.

\leq : Suppose \mathcal{L} has k equivalence classes. Pick a representative from each class: $S_1 = \varepsilon, S_2, \dots, S_k$. Then a DFA for \mathcal{L} is given by

$$(\{S_1, S_2, \dots, S_k\}, \Sigma, S_1, \delta, \{S_1, S_2, \dots, S_k\} \cap \mathcal{L})$$

where

$$\delta(S_i, \sigma) = \text{representative for the equivalence class of } S_i \sigma.$$

□

7.4 Proving nonregularity

There are 3 methods to prove nonregularity:

1. Pumping lemma,
2. Pumping lemma with closure properties, and
3. Myhill-Nerode theorem.

Example 7.4.1. Consider the languages $\mathcal{L}_1 = \{w : w^R\}$ and $\mathcal{L}_2 = \{w : w \text{ has equally many 0's and 1's}\}$. we have shown that $\equiv_{\mathcal{L}_1}$ and $\equiv_{\mathcal{L}_2}$ have infinitely many equivalence classes. Then \mathcal{L}_1 and \mathcal{L}_2 are nonregular.

Example 7.4.2. Consider the language $\mathcal{L} = \{ww : w \in \{0, 1\}^*\}$. Note that

$$i \neq j \Rightarrow 0^i 1 0^j 1 \in \mathcal{L} \wedge 0^j 1 0^i 1 \notin \mathcal{L}.$$

Then $\varepsilon, 0, 00, 000, \dots$ are in different classes. Then $\equiv_{\mathcal{L}}$ has infinitely many classes. Then \mathcal{L} is nonregular.

Example 7.4.3. Consider the language $\mathcal{L} = \{0^n 1^m : n \neq m\}$. Note that

$$i \neq j \Rightarrow 0^i 1^j \in \mathcal{L} \wedge 0^j 1^i \notin \mathcal{L} \Rightarrow 0^i \not\equiv_{\mathcal{L}} 0^j.$$

Then $\varepsilon, 0, 00, 000, \dots$ are in different classes. Then $\equiv_{\mathcal{L}}$ has infinitely many classes. Then \mathcal{L} is nonregular.

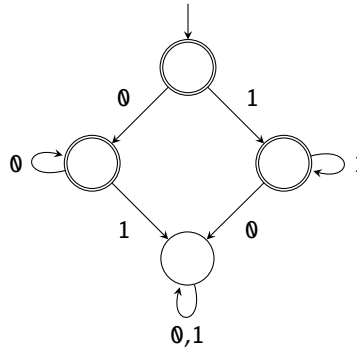
Example 7.4.4. Consider the language $\mathcal{L} = 0^* 1^* \cup \bigcup_{n=0}^{\infty} 1^* 0^n 1^n$. The nonregularity of this language was not provable by using the pumping lemma. However, note that

$$i \neq j \Rightarrow 10^i 1^i \in \mathcal{L} \wedge 10^j 1^i \notin \mathcal{L} \Rightarrow 10^i \not\equiv_{\mathcal{L}} 10^j.$$

Then $10, 100, 1000, \dots$ are in different classes. Then $\equiv_{\mathcal{L}}$ has infinitely many classes. Then \mathcal{L} is nonregular.

7.5 Proving optimality

Example 7.5.1. We claim that the smallest possible DFA for the language $0^* \cup 1^*$ is given by:

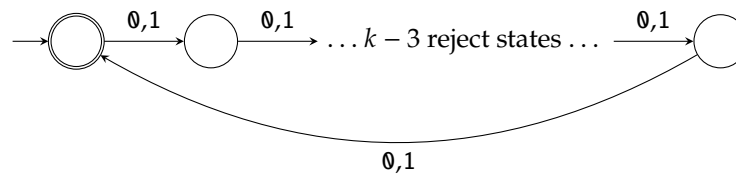


Consider the following table of distinguishing suffixes:

	ε	0	1	01
ε				
0	1			
1	0	0		
01	ε	ε	ε	

Then $\varepsilon, 0, 1, 01$ are in different classes, that is, there are at least 4 equivalence classes. Then the minimum size of any DFA for \mathcal{L} is 4. Then our DFA is optimal.

Example 7.5.2. Find the smallest possible DFA for $\mathcal{L}_k = \{w : k \mid |w|\}$.



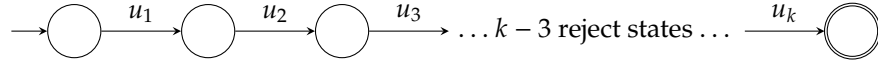
To show that no smaller DFA exists, note that

$$i \neq j \pmod{k} \Rightarrow 0^i 0^{k-i} \in \mathcal{L}_k \wedge 0^j 0^{k-i} \notin \mathcal{L}_k.$$

Corollary 7.5.3

For each k , \mathcal{L}_k is a language that can be recognized by a DFA with k states but not $k - 1$ states.

Example 7.5.4. Given $u = u_1u_2 \cdots u_k$, define $\mathcal{L} = \Sigma^*u_1u_2 \cdots u_k\Sigma^*$. Find the smallest possible DFA for \mathcal{L} .



Note that one still needs to figure out the transitions for $\neg u_1, \neg u_2, \dots, \neg u_k$ from the respective nodes.

However, to show that there exists no smaller DFA for \mathcal{L} , let $i, h \in \{0, 1, \dots, k\}$ where $i > j$. Note that

$$u_1u_2 \cdots u_i \underbrace{u_{i+1} \cdots u_j}_{< k} u_{i+1} \cdots u_k \in \mathcal{L} \wedge u_1u_2 \cdots u_j \underbrace{u_{i+1} \cdots u_k}_{< k} \notin \mathcal{L}.$$

Then $\varepsilon, u_1, u_1u_2, \dots, u_1u_2 \cdots u_k$ are in different classes of $\equiv_{\mathcal{L}}$. Then \mathcal{L} has at least $k + 1$ classes.

8 Context-free grammars

8.1 Context-free grammars defined

Example 8.1.1. A context-free grammar has an alphabet, for brevity we use $\Sigma = \{a, b\}$ here. These are called “terminals” (alphabet symbols).

The second component of a grammar is a set of variables $V = \{S, A, B\}$. These are called “nonterminals.”

Then there are what are known as substitution rules, for instance,

$$\begin{aligned} S &\rightarrow A \\ S &\rightarrow B \\ A &\rightarrow aAb \\ A &\rightarrow \varepsilon \\ B &\rightarrow bBa \\ B &\rightarrow \varepsilon. \end{aligned}$$

Here S is denoted the start variable.

How is this a computational device? What strings does this accept? What strings does this reject?

It all comes down to *derivation*:

$S \rightarrow A$	rule #1
$\rightarrow aAb$	rule #3
$\rightarrow aaAbb$	rule #3
$\rightarrow aaaAbbb$	rule #3
$\rightarrow aaabbb$	rule #4

If a string can be derived by some application of the rule, then the string is in the language, and vice versa.

Definition 8.1.2

The **language** of a context-free grammar (CFG) is the set of all strings that can be derived.

Shorthand:

$$\begin{aligned} S &\rightarrow A \mid B \\ A &\rightarrow aAb \mid \varepsilon \\ B &\rightarrow bBa \mid \varepsilon \end{aligned}$$

The language is $\mathcal{L} = \{a^n b^n \cup b^n a^n : n \geq 0\}$.

Definition 8.1.3

A **context-free grammar** is a tuple (V, Σ, R, S) where

- V = a set of variables,
- Σ = the alphabet,
- R = the set of rules, and
- S = the start variable.

Then the CFG we discussed above can be written as

$$(\{S, A, B\}, \{a, b\}, \{S \rightarrow A \mid B, A \rightarrow aAb \mid \varepsilon, B \rightarrow bBa \mid \varepsilon\}, S).$$

8.2 Reading a CFG

Example 8.2.1. Consider the CFG with rules

$$\begin{aligned} S &\rightarrow AA \\ A &\rightarrow A. \end{aligned}$$

The language is \emptyset .

Example 8.2.2. Consider the CFG with a single rule

$$S \rightarrow abS \mid a.$$

The language is $(ab^*)a$, or equivalently $a(ba)^*$.

Example 8.2.3. Consider the CFG with a single rule

$$S \rightarrow aSb \mid \varepsilon.$$

The language is $\{a^n b^n : n \geq 0\}$.

Example 8.2.4. Consider the CFG with rules

$$\begin{aligned} S &\rightarrow aSa \mid B \\ B &\rightarrow bB \mid \varepsilon. \end{aligned}$$

The language is $\{a^n b^m a^n : n \geq 0 \wedge m \geq 0\}$ or $a^n b^* a^n$ where $n \geq 0$.

Example 8.2.5. Consider the CFG with a single rule

$$S \rightarrow SS \mid aS \mid b.$$

The language is Σ^*b .

Example 8.2.6. Consider the CFG with rules

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow aAa \mid \varepsilon \\ B &\rightarrow bBb \mid \varepsilon. \end{aligned}$$

The language is $(aa)^*(bb)^*$.

Example 8.2.7. Consider the CFG with rules

$$\begin{aligned} S &\rightarrow \varepsilon \mid aT \mid bT \\ T &\rightarrow aS \mid bS. \end{aligned}$$

The language is $(\Sigma\Sigma)^*$.

8.3 Constructing CFGs

Example 8.3.1. Consider the language $\mathcal{L} = \{a^n(a \cup b)^n : n \geq 0\}$.

A CFG for this language is given by

$$S \rightarrow aSa \mid aSb \mid \varepsilon.$$

Example 8.3.2. Consider the language $\mathcal{L} = \{(a \cup b)^n a(a \cup b)^n : n \geq 0\}$.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow TST \mid a \\ T &\rightarrow a \mid b. \end{aligned}$$

Example 8.3.3. Consider the language $\mathcal{L} = \{a^n b^m : 0 \leq m \leq 3n\}$.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow aSBBB \mid \varepsilon \\ B &\rightarrow b \mid \varepsilon. \end{aligned}$$

Example 8.3.4. Consider the language $\mathcal{L} = \{w : w = w^R\}$.

A CFG for this language is given by

$$S \rightarrow aSa \mid bSb \mid a \mid b \mid \varepsilon.$$

Example 8.3.5. Consider the language $\mathcal{L} = \{a^n b^m : n \geq m\}$.

A CFG for this language is given by

$$S \rightarrow aS \mid aSb \mid \varepsilon.$$

Example 8.3.6. Consider the language $\mathcal{L} = \{a^n b^m : n > m\}$.

A CFG for this language is given by

$$S \rightarrow aS \mid aSb \mid a.$$

Example 8.3.7. Consider the language $\mathcal{L} = \{w : w \text{ contains at least three a's}\}$.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow GaGaGaG \\ G &\rightarrow aG \mid bG \mid \varepsilon. \end{aligned}$$

Example 8.3.8. Consider the language $\mathcal{L} = \{w : w \text{ contains abb as a substring}\}$.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow GabbG \\ G &\rightarrow aG \mid bG \mid \varepsilon. \end{aligned}$$

Example 8.3.9. Consider the language $\mathcal{L} = \{a^n b^{m+n} a^m : n, m \geq 0\}$.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow LR \\ L &\rightarrow aLb \mid \varepsilon \\ R &\rightarrow bRa \mid \varepsilon. \end{aligned}$$

Example 8.3.10. Consider the alphabet $\Sigma = \{a, b, \varepsilon, \emptyset, (,), \cup, *, \}$ and the language $\mathcal{L} =$ all regular expressions.

A CFG for this language is given by

$$S \rightarrow \varepsilon \mid \emptyset \mid a \mid b \mid (SS) \mid (S \cup S) \mid (S^*).$$

Example 8.3.11. Consider a programming language:

- identifiers = $(a \cup \dots \cup z)(a \cup \dots \cup z \cup \emptyset \cup \dots \cup 9)^*$
- numbers = positive integers
- statements of 2 types:
 - assignment
 - conditional
- program = a single statement

Now consider the alphabet $\Sigma = \{a, b, \dots, z, \emptyset, 1, \dots, 9, >, =, <, \sqcup, \P\}$ and the language $\mathcal{L} =$ all valid programs.

A CFG for this language is given by

$$\begin{aligned} S &\rightarrow A \mid C \\ D^+ &\rightarrow 1 \mid 2 \mid \dots \mid 9 && \text{nonzero digit} \\ D &\rightarrow \emptyset \mid D^+ && \text{digit} \\ L &\rightarrow a \mid b \mid \dots \mid z && \text{letter} \\ I &\rightarrow ID \mid IL \mid L && \text{identifier} \\ N &\rightarrow ND \mid D^+ && \text{number} \\ Q &\rightarrow I \mid N && \text{numeric quantity} \end{aligned}$$

$W \rightarrow \sqcup \mid \P \mid W$	whitespace
$W' \rightarrow \varepsilon \mid W$	optional whitespace
$T \rightarrow QW' = W'Q \mid QW' < W'Q \mid QW' > W'Q$	test
$A \rightarrow IW' = W'Q$	assignment
$C \rightarrow \text{if } WT \text{ then } WS \mid \text{if } WTW \text{ then } WSW \text{ then } WS$	conditional.

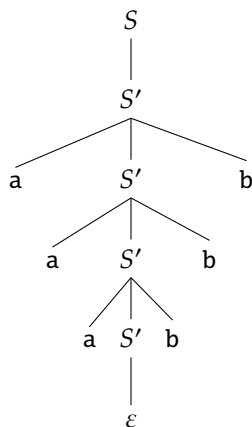
8.4 Ambiguity

Example 8.4.1. Consider the CFG with rules

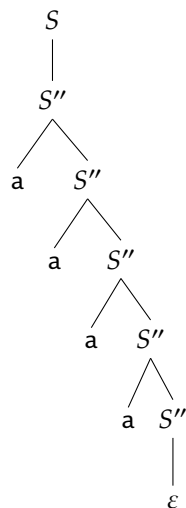
$$\begin{aligned} S &\rightarrow S' \mid S'' \\ S' &\rightarrow aS'b \mid \varepsilon \\ S'' &\rightarrow aS'' \mid \varepsilon. \end{aligned}$$

Note that this grammar generates the language $\{a^n b^n : n \geq 0\} \cup a^*$.

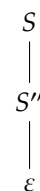
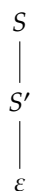
A parse tree for the string aaabbb is given by:



A parse tree for the string aaaa is given by:



We can find two parse trees for the string ε :



Definition 8.4.2

A grammar is called **ambiguous** if a string has more than one parse tree. Otherwise the grammar is called “unambiguous.”

Example 8.4.3. Note that the CFG with rules

$$\begin{aligned} S &\rightarrow S' \mid S'' \\ S' &\rightarrow aS'b \mid \varepsilon \\ S'' &\rightarrow aS'' \mid \varepsilon \end{aligned}$$

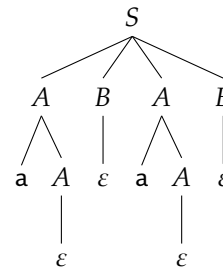
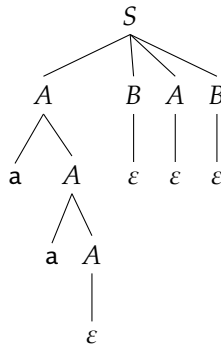
is ambiguous since the string ε has two parse trees. The ambiguity is fixed in the CFG with rules

$$\begin{aligned} S &\rightarrow S' \mid S'' \\ S' &\rightarrow aS'b \mid \varepsilon \\ S'' &\rightarrow aS'' \mid a. \end{aligned}$$

Example 8.4.4. Note that the CFG with rules

$$\begin{aligned} S &\rightarrow ABAB \\ A &\rightarrow aA \mid \varepsilon \\ B &\rightarrow bB \mid \varepsilon \end{aligned}$$

generates language $a^*b^*a^*b^*$ and ambiguously generates the string aa :



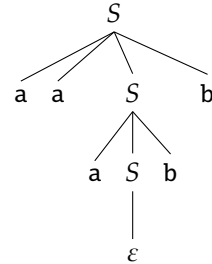
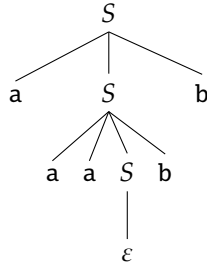
The ambiguity is fixed in the CFG with rules

$$\begin{aligned} S &\rightarrow \varepsilon \mid \overbrace{A \mid B \mid AB \mid BA \mid ABA \mid BAB \mid ABAB}^{8 \text{ disjoint languages}} \\ A &\rightarrow aA \mid a \\ B &\rightarrow bB \mid b. \end{aligned}$$

Example 8.4.5. Note that the CFG with a single rule

$$S \rightarrow aSb \mid aaSb \mid \varepsilon$$

generates the language $\{a^n b^m : m \leq n \leq 2m\}$ and ambiguously generates the string aaabb:



The ambiguity is fixed in the CFG with rules

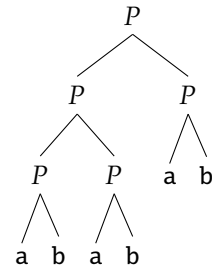
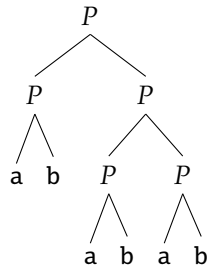
$$S \rightarrow aSb \mid T \mid \varepsilon$$

$$T \rightarrow aaSb \mid aab.$$

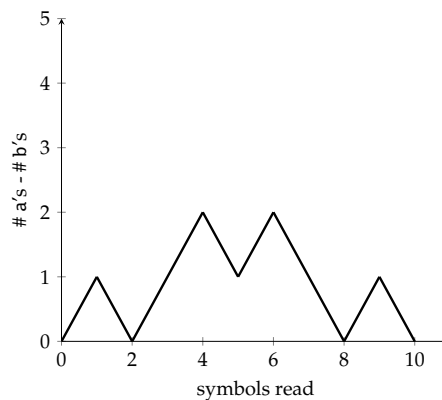
Example 8.4.6. Note that the CFG with a single rule

$$P \rightarrow PP \mid aPb \mid ab$$

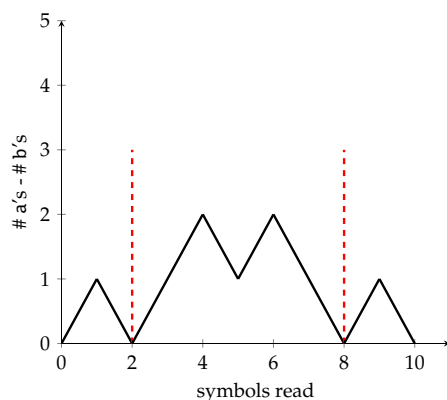
generates the language of properly nested parentheses, where a represents (and b represents), and ambiguously generates the string ababab:



Consider the following graph:



Note that we can unambiguously split this graph into properly parenthesized expressions that have an outer pair of parentheses:

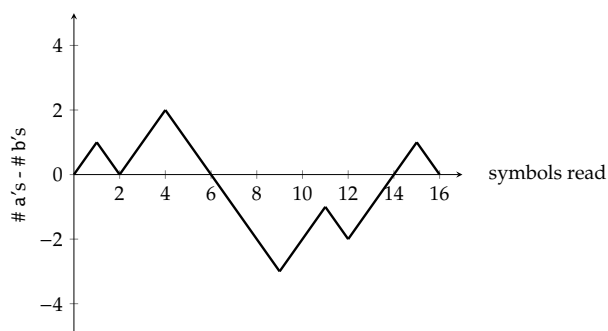


Then the ambiguity is fixed in the CFG with the single rule

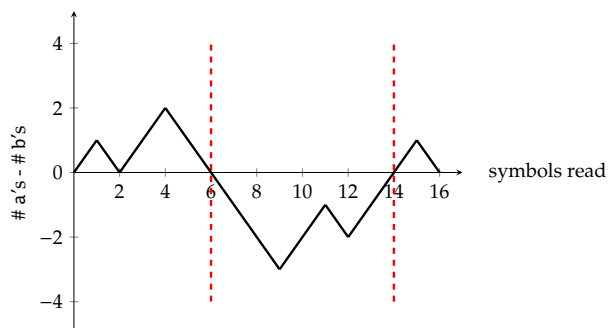
$$P \rightarrow \underbrace{ab \mid aPb}_{\text{1 pair of outer parentheses}} \mid \underbrace{abP \mid aPbP}_{\geq 2 \text{ pairs of outer parentheses}}$$

Example 8.4.7. Consider the language $\mathcal{L} = \{w : w \text{ has equally many a's and b's}\}$.

Consider the following graph:



Note that we can unambiguously split this graph into maximal contiguous sections that lie completely above or below the x -axis:



Then the ambiguity is fixed in the CFG with rules

$$S \rightarrow S' \mid S'' \mid \varepsilon$$

$$S' \rightarrow P \mid PS''$$

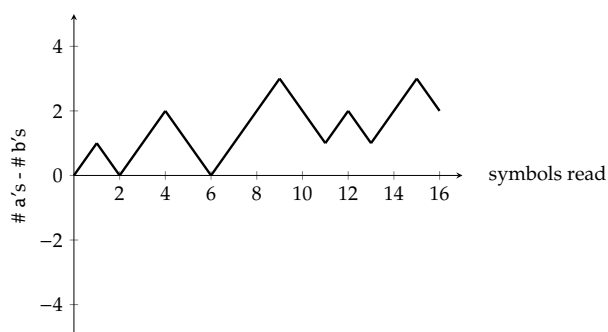
$$S'' \rightarrow Q \mid QS'$$

$$P \rightarrow ab \mid aPb \mid abP \mid aPbP$$

$$Q \rightarrow ba \mid bQa \mid baQ \mid bQaQ$$

Example 8.4.8. Consider the language $\mathcal{L} = \{w : \text{every prefix of } w \text{ has at least as many a's as b's}\}$.

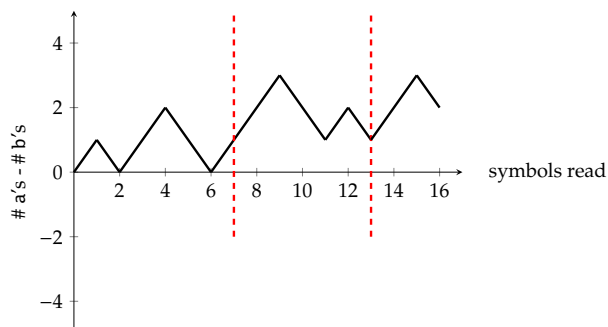
Consider the following graph:



The two cases we can have are where

1. the balance starts at 0 and ends at 0, and where
2. the balance starts at 0 and ends up positive.

These two cases are disjoint:



Then the ambiguity is fixed in the CFG with a single rule

$$S \rightarrow \underbrace{\varepsilon \mid P}_{\text{case 1}} \mid \underbrace{aP \mid PaS}_{\text{case 2}}.$$

8.5 Ambiguity in English

Consider the following sentences:

The hunter attacked the bear with a stick.

The complex houses married students.

The bird the cat the dog the boy saw barked at chased flew away.

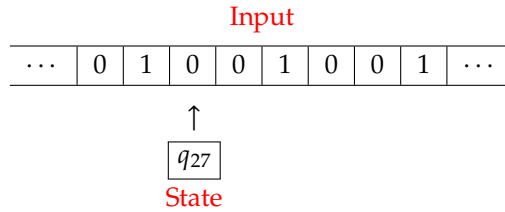
Ship-shipping ships ship ship-shipping ships.

These examples drive home the point that English is very subtle and that you cannot represent it with a context-free grammar.

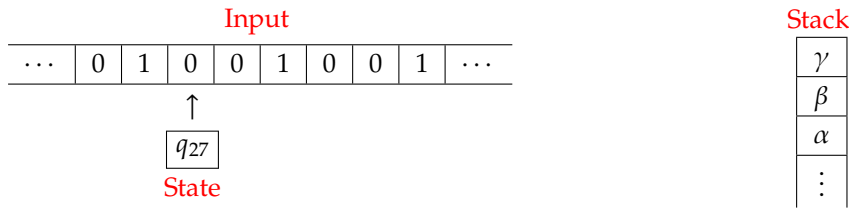
9 Pushdown Automata

9.1 Definitions

We may represent an NFA in the following way:

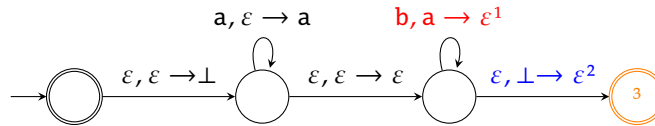


Then a PDA may be represented in the following way:



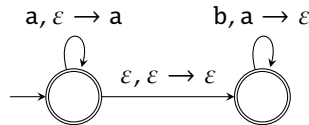
Transitions in PDAs are in the form of $\alpha, \beta \rightarrow \gamma$ where α is the input symbol to read (read nothing if $\alpha = \epsilon$), β is the symbol to pop from the top of the stack (do not pop if $\beta = \epsilon$), and γ is the symbol to push onto the stack (do not push if $\gamma = \epsilon$).

Example 9.1.1. A PDA for $\mathcal{L} = \{a^n b^n : n \geq 0\}$ is given by:

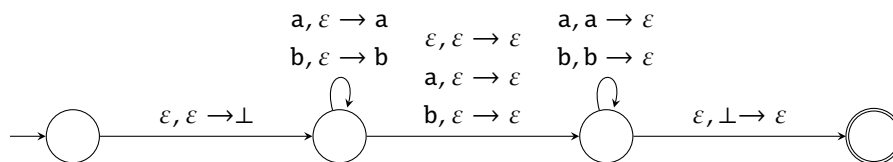


Note that the input alphabet is $\Sigma = \{a, b\}$ while the stack alphabet is $\Gamma = \{\perp, a\}$ where \perp is the bottom marker used to the bottom of the stack.

Example 9.1.2. A PDA for $\mathcal{L} = \{a^n b^m : n \geq m \geq 0\}$ is given by:



Example 9.1.3. A PDA for $\mathcal{L} = \{w : w = w^R\}$ is given by:

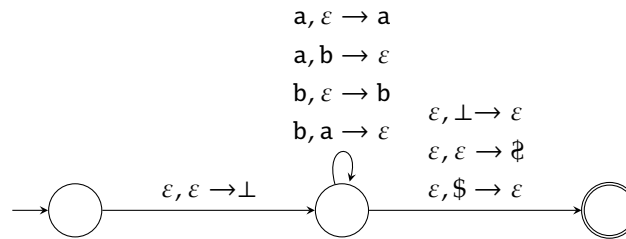


¹Usable only if next symbol is b and a is on top of the stack.

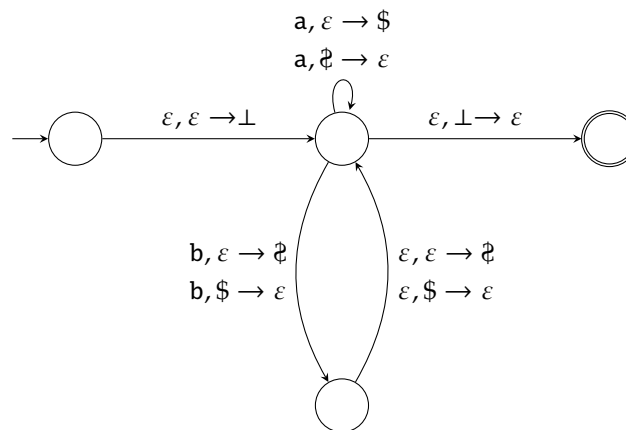
²Usable only if \perp is on top of the stack.

³The PDA blows up if input is not finished.

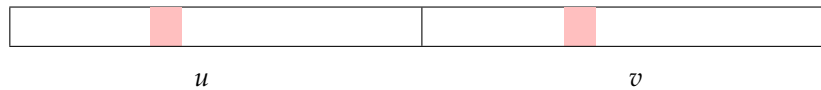
Example 9.1.4. A PDA for $\mathcal{L} = \{w : w \text{ has equally many a's and b's}\}$ is given by:



Example 9.1.5. A PDA for $\mathcal{L} = \{w : w \text{ has twice as many a's as b's}\}$ is given by:

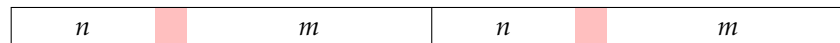


Example 9.1.6. Consider the language $\mathcal{L} = \{uv : |u| = |v| \wedge u \neq v\}$. A string $w \in \mathcal{L}$ may be visually represented by



where the pink cells represent disagreement.

One may intuitively divide the string up like so:



However, it is very hard to find a PDA for this language from this. Dividing it up in a different way makes it a lot easier:



Then a PDA for \mathcal{L} is given by:

