# 1

# Formal Languages

## *Mathematical Preliminaries*

A set is a collection of elements. Those elements can be anything, including other sets.

Sets can be described by listing their elements, such as $\{a, b, c, \ldots\}$.

The empty set is denoted $\varnothing$, or $\{\,\}$.

A set $A$ is *finite* if it has a finite number of elements, that is, if there is a natural number $n \in \mathbb{N}$ such that $A$ has $n$ elements. If no such $n$ exists, then $A$ is *infinite*.

The main relation on sets is *set membership*, written $a \in A$: $a$ is an element of set $A$.

Sets are *extensional*: two sets are equal if they have exactly the same elements.

Another relation on sets is *set inclusion*, written $A \subseteq B$: $A$ is a subset of $B$, meaning that every element of $A$ is an element of $B$).

Some properties of $\subseteq$:

$\varnothing \subseteq A$ for every $A$

$A \subseteq A$ for every $A$

If $A \subseteq B$ and $B \subseteq C$, then $A \subseteq C$

$A = B$ if and only if $A \subseteq B$ and $B \subseteq A$.

If $P$ is a property, then $\{x \mid P(x)\}$ is the set of all elements satisfying the

property. (Technically speaking, there are some restrictions on what makes up an acceptable property in this context — but they won't impact us.)

Common operations on sets:

$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$

$A \cap B = \{x \mid x \in A \text{ and } x \in B\}$

$\overline{A} = \{x \in \mathcal{U} \mid x \notin A\}$ where $\mathcal{U}$ is a universe of elements with $A \subseteq \mathcal{U}$ — the definition of ¯ therefore depends on the universe under consideration, which will often be clear from context

$A \backslash B = \{x \mid x \in A \text{ but } y \notin B\}$

$A \times B = \{\langle x, y \rangle \mid x \in A, y \in B\}$, the set of all pairs of elements from $A$ and $B$. This generalizes in the obvious way to products $A_1 \times A_2 \times \cdots \times A_k$.

A function $f : A \longrightarrow B$ associates (or maps) every element of $A$ to an element of $B$. Set $A$ is the domain of the function, and $B$ is the codomain. The image of $A$ under $f$ is the subset of $B$ defined by $\{b \in B \mid f(a) = b \text{ for some } a \in A\}$.

If $f : A \longrightarrow B$ and $g : B \longrightarrow C$, then the *composition* $g \circ f : A \longrightarrow C$ defined by $(g \circ f)(x) = g(f(x))$.

A function $f : A \longrightarrow B$ is *one-to-one* if it maps distinct elements of $A$ into distinct elements of $B$ (that is, if $a \neq b$, then $f(a) \neq f(b)$). A function $f : A \longrightarrow B$ is *onto* if every element of $B$ is in the image of $A$ under $f$ (that is, if for every element $b \in B$ there is an element $a \in A$ with $f(a) = b$). A function $f : A \longrightarrow B$ is a *one-to-one correspondance* if it is both one-to-one and onto.

## Decisions Problems

Intuitively, a computation is a way to "implement" a mathematical function $f : A \longrightarrow B$. A big part of the course is to build an understanding of what that intuition means.

Arbitrary functions between arbitrary sets $A$ and $B$ is too broad a class of functions to work with. Historically, researchers have looked at two classes of functions to study computation:

1. **Natural number functions** of the form

$$f : \mathbb{N} \times \cdots \times \mathbb{N} \longrightarrow \mathbb{N}$$

2. **Decision problems** of the form

$$d : A \longrightarrow \{1, 0\}$$

where 1 can be interpreted as true and 0 as false. Decision problems are predicates on the domain $A$.

For the first half of this course, we will consider study computability for decision problems. Example of decision problems include determining if a graph is planar, or determining if two natural numbers are coprime. Restricting to decision problems looks like a limitation, in that many functions of interest we may be interested in computing are not in that form, such as addition, but we can capture addition using a decision problems by considering the decision problem that takes three numbers and determines if the sum of the first two is equal to the third.

One reason why decision problems form an interesting class of problems for which to study computability is that they can be "summarized" using simple sets. You can "summarize" a decision problem

$$d : A \longrightarrow \{1, 0\}$$

using the set

$$S = \{a \in A \mid d(a) = 1\}$$

Given such a set $S$, you can *reconstruct d* readily, by taking

$$d(a) = \begin{cases} 1 & \text{if } a \in S \\ 0 & \text{otherwise} \end{cases}$$

It is easy to check that you get the same decision problem back. Thus, you can study a decision problem by studying the set $S$ that summarizes it.

We will make a further assumption, that the domain $A$ for decision problems will be a set of *strings*.

Let $\Sigma$ be a non-empty finite set we will call the *alphabet*. A *string over $\Sigma$* is a (possibly empty) finite sequence of elements of $\Sigma$, usually written $a_1 \cdots a_k$, where $a_i \in \Sigma$. For example, aaccabc is a string over alphabet $\{a, b, c\}$. It

is also a string over alphabet $\{\mathsf{a}, \mathsf{b}, \mathsf{c}, \mathsf{d}\}$. We will use $u, v, w$ to range over strings.

The length of $u = a_1 \ldots a_k$, written $|u|$, is $k$.

The empty string is written $\epsilon$. It has length 0.

The set of all strings over $\Sigma$ is denoted $\Sigma^*$. Note that this is an infinite set. (Why?)

If $u = a_1 \ldots a_k$ and $v = b_1 \cdots b_m$ are strings over $\Sigma$, then the *concatenation uv* is the string $a_1 \cdots a_k b_1 \cdots b_m$. Note that $\epsilon u = u\epsilon = u$ for every string $u$.

We define $u^0 = \epsilon$, $u^1 = u$, $u^2 = uu$, $u^3 = uuu$, etc.

We will transform the problem of of determining whether a decision problem is computable into the problem of determining whether a set of strings is computable. But of course, what we will mean w hen we say that a set of strings $A$ is computable is that the corresponding decision problem constructed from $A$ as above is computable. Why do we do this? Because sets of strings are much easier to work with: they are sets, and we can manipulate sets in many different ways.

## *Languages*

Because sets of strings are so important, let's give them a name. A *formal language* (usually called only a language) over alphabet $\Sigma$ is a set of strings over $\Sigma$.

Since languages are sets, we inherit the usual set operations $A \cup B$, $A \cap B$, $\overline{A}$ (where the universe of $A$ is taken to be $\Sigma^*$).

Because a language $A$ is a set of strings specifically, we can also define more specific operations.

$A \cdot B = \{uv \mid u \in A \text{ and } v \in B\}$, that is, the set of all strings obtained by concatenating a string of $A$ and a string of $B$.

$A^0 = \{\epsilon\}$

$A^1 = A$

$A^2 = A \cdot A$

$A^3 = A \cdot A \cdot A$

4

etc...

$$A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \cdots = \bigcup_{k \geq 0} A^k$$

The $^*$ operation is called the *Kleene star*.

Some properties that are easy to verify:

$$\emptyset \cdot A = A \cdot \emptyset = \emptyset$$

$$\{\epsilon\} \cdot A = A \cdot \{\epsilon\} = A$$

**Example:** $\Sigma = \{a, b\}$, and $A = \{aa, bb\}$. Then $A^* = \{\epsilon, aaaa, aabb, bbaa,$ bbbb, aaaaaa, aaaabb, aabbaa, aabbbb, bbaaaa, bbaabb, bbbbaa, bbbbbb, $\dots\}$.

This explains why I wrote $\Sigma^*$ for the set of all strings: if we consider $\Sigma$ as a set of strings, each of length 1, then $\Sigma^*$ according to the above definition indeed gives the set of all strings over alphabet $\Sigma$.

The operations $\cup$, $\cdot$, and $^*$ are called the *regular operations*.

### *Regular Languages*

A language over alphabet $\Sigma = \{a_1, \dots, a_k\}$ is *regular* if it can be obtained from the sets $\emptyset, \{\epsilon\}, \{a_1\}, \dots, \{a_k\}$ and finitely many applications of the regular operations.

**Example:** consider the language of all even-length strings over alphabet $\{a, b\}$. It can be obtained as:

$$(( \{a\} \cup \{b\} \cup \{c\}) \cdot (\{a\} \cup \{b\} \cup \{c\}))^*$$

Therefore, that language is regular.

**Example:** consider the language of all strings over $\{a, b, c\}$ that start and end with an a:
$$(\{a\} \cdot (\{a\} \cup \{b\} \cup \{c\})^* \cdot \{a\}) \cup \{a\}$$

Therefore, that language is also regular.

Regular languages form a very natural class of languages, and we will see soon that they arise out of an equally natural model of computation.

Some natural ways to form regular languages:

- The empty language is regular, pretty much by definition.

- Every singleton language (language with a single string) is regular. This is easy to see: language $\{a_1 \ldots a_k\}$ is obtained by taking $\{a_1\} \cup \cdots \cup \{a_k\}$.

- Every finite language (i.e., language with a finite number of strings) is regular. Again, this is easy to see. Say $A = \{u_1, \ldots, u_k\}$. By the previous statement, each of $\{u_1\}$, $\ldots$, $\{u_k\}$ are regular. Therefore, their union is regular.

- if $A$ and $B$ are regular languages, then $A \cup B$ and $A \cdot B$ are regular languages. That follows directly from the definition.

- If $A$ is regular, then $A^*$ is regular. Again, this follows directly from the definition. This means, in particular, that $\Sigma^*$, the set of all strings over a given alphabet, is regular.

- If $A$ and $B$ are regular, then $A \cap B$ is regular. We'll show this is the case later, because we don't have enough techniques to show that just yet.

- If $A$ is regular, then $\overline{A} = \Sigma^* - A$ is regular. Again, we don't have enough techniques to show that just yet.

- If $A$ is regular, then $rev(A)$, the set of all strings from $A$ but in reverse (where the reverse of $a_1 \ldots a_k$ is just $a_k \ldots a_1$) is regular. We'll be able to show that next section.

You might be led to believe at this point that every language is regular. That's not true. The most natural non-regular language is probably the language of all *palindrome* strings over an alphabet $\Sigma$. A palindrome string is a string with the property that it and its reverse are equal, such as aaa, or abba. It's not obvious that this is not regular, and we don't have enough techniques to show that yet.

### *Regular Epressions*

Regular expressions are a convenient notation for regular languages.

A regular expression over alphabet $\Sigma$ is defined by the following syntax:

$$r ::= \quad 1$$

$$\mathbb{0}$$

$$a \quad \text{for every } a \in \Sigma$$

$$(r_1{+}r_2)$$

$$(r_1 r_2)$$

$$(r_1^*)$$

We usually drop parentheses, under the assumption that $r_1^*$ binds tighter than concatenation $r_1 r_2$ which binds tighter than $r_1{+}r_2$. For example, ab+ac is a regular expression, as is a(b+c) and a*(b+c)*.

A regular expression $r$ denotes a language $[\![ r ]\!]$ over $\Sigma$ in the following way:

$$[\![ 1 ]\!] = \{\epsilon\}$$

$$[\![ \mathbb{0} ]\!] = \varnothing$$

$$[\![ a ]\!] = \{a\}$$

$$[\![ r_1{+}r_2 ]\!] = [\![ r_1 ]\!] \cup [\![ r_2 ]\!]$$

$$[\![ r_1 r_2 ]\!] = [\![ r_1 ]\!] \cdot [\![ r_2 ]\!]$$

$$[\![ r_1^* ]\!] = [\![ r_1 ]\!]^*$$

For example:

$$
\begin{aligned}
[\![ ab{+}ac ]\!] &= [\![ ab ]\!] \cup [\![ ac ]\!] \\
&= ([\![ a ]\!] \cdot [\![ b ]\!]) \cup ([\![ a ]\!] \cdot [\![ c ]\!]) \\
&= (\{a\} \cdot \{b\}) \cup (\{a\} \cdot \{c\}) \\
&= \{ab\} \cup \{ac\} \\
&= \{ab, ac\}
\end{aligned}
$$

$$
\begin{aligned}
[\![ a(b{+}c) ]\!] &= [\![ a ]\!] \cdot [\![ b{+}c ]\!] \\
&= [\![ a ]\!] \cdot ([\![ b ]\!] \cup [\![ c ]\!]) \\
&= \{a\} \cdot (\{b\} \cup \{c\}) \\
&= \{a\} \cdot (\{b, c\}) \\
&= \{ab, ac\}
\end{aligned}
$$

$$[\![ a^*(b{+}c)^* ]\!] = [\![ a^* ]\!] \cdot [\![ (b{+}c)^* ]\!]$$

$$= [\![a]\!]^* \cdot [\![b+c]\!]^*$$
$$= \{a\}^* \cdot ([\![b]\!] \cup [\![c]\!])^*$$
$$= \{a\}^* \cdot (\{b\} \cup \{c\})^*$$
$$= \{a\}^* \cdot \{b, c\}^*$$

And thinking about this last set (which is difficult to write down), it is basically the set of all strings obtained by concatenating a sequence of as (including none) to a sequence of bs and cs (in any order, including none). So aaaaaa is in this set, as is aaaab, aaaabbbb, aaaabbbcbcbc, etc.

**Theorem:** A language $A$ is regular exactly if there is a regular expression $r$ such that $[\![r]\!] = A$.

We can use regular expressions to show that if $A$ is regular, then $rev(A) = \{rev(u) \mid u \in A\}$, where $rev(u)$ is the reverse of string $u$, is regular: since $A$ is regular, there is a regular expression $r$ with $[\![r]\!] = A$. We take this regular expression and transform it into regular expression $\hat{r}$ as follow:

$$\hat{1} = 1$$
$$\hat{0} = 0$$
$$\hat{a} = a$$
$$\widehat{r_1+r_2} = \hat{r_2}+\hat{r_1}$$
$$\widehat{r_1 r_2} = \hat{r_2} \cdot \hat{r_1}$$
$$\widehat{r_1^*} = (\hat{r_1})^*$$

Now, if $[\![r]\!] = A$, then $[\![\hat{r}]\!] = rev(A)$, and therefore $rev(A)$ is regular.