

FORMAL LANGUAGES

Alphabets and Strings

- An **alphabet** Σ is a finite set of **characters** (or **symbols**).
- A **word**, or **sequence**, or **string over** Σ is any group of 0 or more consecutive characters of Σ .
- The **length** of a word is the number of characters in the word.
- The **null string** is the string of length 0. It is denoted ε or λ .
- A string of length n is really an ordered n -tuple of characters written without parentheses or commas.
- Given two strings x and y over Σ , the **concatenation** of x and y is the string xy obtained by putting all the characters of y right after x .

Languages over an alphabet

Let Σ be an alphabet. A **formal language over** Σ is a set of strings over Σ .

- \emptyset is the **empty language** (over Σ)
- $\Sigma^n = \{\text{all strings over } \Sigma \text{ that have length } n\}$ where $n \in \mathbb{N}$
- $\Sigma^+ =$ the **positive closure** of $\Sigma = \{\text{all strings over } \Sigma \text{ that have length } \geq 1\}$
- $\Sigma^* =$ the **Kleene closure** of $\Sigma = \{\text{all strings over } \Sigma\}$

Operations on Languages

Let Σ be an alphabet. Let L and L' be two languages defined over Σ .

The following operations define new languages over Σ :

- The **concatenation of** L and L' , denoted LL' , is $LL' = \{xy \mid x \in L \wedge y \in L'\}$
- The **union of** L and L' , denoted $L \cup L'$, is $L \cup L' = \{x \mid x \in L \vee x \in L'\}$
- The **Kleene closure of** L , denoted L^* , is $L^* = \{x \mid x \text{ is a concatenation of any finite number of strings in } L\}$. Note that $\varepsilon \in L^*$.

REGULAR EXPRESSIONS

Definition

Let Σ be an alphabet. The following are **regular expressions (r.e.)** over Σ :

- I. BASE: ε and each individual symbol of Σ are regular expressions.
- II. RECURSION: if r and s are regular expressions over Σ , then the following are also regular expressions over Σ :
 - (rs) the concatenation of r and s
 - $(r | s)$ r or s
 - (r^*) the Kleene closure of r
- III. RESTRICTION: The only regular expressions over Σ are the ones defined by I and II above.

Order of Precedence of Regular Expression Operations

- The order of precedence of r.e. operators are, from highest to lowest:
- Highest: $()$ * concatenation | : lowest

Languages Defined by Regular Expressions

Let Σ be an alphabet. Define a function L as follows:

$$L: \begin{cases} \{\text{all r. e.'s over } \Sigma\} \rightarrow \{\text{all languages over } \Sigma\} \\ r \mapsto L(r) = \text{the language defined by } r \end{cases}$$

- I. BASE: $L(\varepsilon) = \{\varepsilon\}, \forall a \in \Sigma L(a) = \{a\}$
- II. RECURSION: If $L(r)$ and $L(s)$ are the languages defined by the regular expressions r and s over Σ , then
 - $L(rs) = L(r)L(s)$
 - $L(r|s) = L(r) \cup L(s)$
 - $L(r^*) = (L(r))^*$

Variations

Some definitions of regular expressions and regular languages define \emptyset to be a r.e. with $L(\emptyset) = \emptyset$

Shorthand:

- $[a-c] = a|b|c$
- $[^a-c] = \text{any letter other than } a, b, c$
- $r^+ = rr^*$
- $r\{n\} = r$ is concatenated n times
- $[a-c x-z] = [a-c, x-z] = a|b|c|x|y|z$
- $.$ means any character
- $r? = (r|\varepsilon)$
- $r\{n,m\}$ r concatenated n to m times