

Formal languages, grammars, and automata

Herman Geuvers

(All info on

<http://www.cs.ru.nl/~herman/onderwijs/FLGA/>

or google me)

Overview

Topics

Languages:	regular	context-free	[natural languages]	[enumerable]
Automata:	finite	push-down	[bounded Turing machine]	[Turing machine]
Grammars:	regular	context-free	[context-sensitive]	[NL grammars?]

Automata: **accept words of a language**
given a word, compute if it is in the language

Grammars: **generate words of a language**
produce all correct words in the language

Languages

Alphabet

An *alphabet* Σ is a (finite) set of symbols

Examples

$$\Sigma_1 = \{a\}$$

$$\Sigma_2 = \{0, 1\}$$

$$\Sigma_3 = \{A, C, G, T\}$$

$$\Sigma_4 = \{a, b, c, d, \dots, x, y, z\}$$

$$\Sigma_5 = \{s \mid s \text{ is an ascii symbol}\}$$

$$\Sigma_6 = \{\}$$

Japanese alphabet: 2×52 signs

$$\Sigma_7 = \{, \dots\}$$

Chinese alphabet: ± 40.000 signs

$$\Sigma_8 = \{0, 1, +, \times, x_0, x_1, x_2, \dots\}$$

mathematical alphabet, number of signs: countably infinite

$$\Sigma_9 = \{0, 1, +, \times, x_0, x_1, x_2, \dots\} \cup \{c_r \mid r \in \mathbb{R}\}$$

mathematical alphabet, number of signs uncountably infinite

Words

A *word* (string) over Σ is a finite sequence of elements from Σ

The set Σ^* consists of all *words over* Σ

Inductive generation of words

$$\begin{array}{l} \lambda \in \Sigma^*, \quad \lambda \text{ denotes the empty word} \\ w \in \Sigma^* \Rightarrow ws \in \Sigma^*, \quad \text{for all } s \in \Sigma \end{array}$$

Note that λw is just w

Note the difference between $a \in \Sigma$ and $a \in \Sigma^*$

Think of a word as a chain of letters on a necklace:

$$\begin{array}{l} \lambda = \text{---} \\ Eva = \text{---}E\text{-}v\text{-}a\text{---} \end{array}$$

The difference between a and $\text{---}a\text{---}$ is clear

Operation on words; Language

Operations on words

$$\begin{aligned} u \in \Sigma^*, v \in \Sigma^* &\Rightarrow u \cdot v \in \Sigma^* && \text{concatenation} \\ u \in \Sigma^*, n \in \mathbb{N} &\Rightarrow u^n \in \Sigma^* && \text{repeating words} \end{aligned}$$

Inductive definitions

$$\begin{aligned} u \cdot \lambda &= u \\ u \cdot (vs) &= (u \cdot v)s \end{aligned}$$

$$\begin{aligned} u^0 &= \lambda \\ u^{k+1} &= u^k \cdot u \end{aligned}$$

usually we write concatenation $u \cdot v$ as uv

A *language over* Σ is a subset of Σ^* , notation $L \subseteq \Sigma^*$

Examples

$$\begin{aligned} L_1 &= \{w \in \{a, b\}^* \mid abba \text{ is a substring of } w\} \\ L_2 &= \{w \in \{a, b\}^* \mid abba \text{ is not a substring of } w\} \end{aligned}$$

Examples of languages

Let $\Sigma = \{a, b, c\}$.

1. $L_1 = \{a^n \mid n \in \mathbb{N} \text{ is even}\}$
2. $L_2 = \{a^n b^n \mid n \in \mathbb{N}\}$
3. $L_3 = \{a^n b^n c^n \mid n \in \mathbb{N}\}$
4. $L_4 = \{a^n \mid n \in \mathbb{N} \text{ is prime}\}$
5. $L_5 = \{n \mid n \text{ denotes an integer number}\}$
6. $L_6 = \{e \mid e \text{ is a well-formed arithmetical expression}\}$
7. $L_7 = \{P \mid P \text{ is a syntactically correct Java program}\}$
8. $L_8 = \{S \mid S \text{ is a grammatically correct English sentence}\}$

Operations on languages

Given languages $L_1, L_2, L \subseteq \Sigma^*$ we can define

$$L_1 \cup L_2$$

$$L_1 L_2$$

$$L^*$$

again languages over Σ

$$L_1 \cup L_2 = \{w \mid w \in L_1 \text{ or } w \in L_2\}$$

$$L_1 L_2 = \{w_1 w_2 \mid w_1 \in L_1 \text{ \& } w_2 \in L_2\}$$

$$L^0 = \{\lambda\}$$

$$L^{n+1} = L^n L$$

$$L^* = \bigcup_{n \in \mathbb{N}} L^n = L^0 \cup L^1 \cup L^2 \cup \dots$$

$$\neq \{w^n \mid w \in L, n \in \mathbb{N}\}$$

Regular expressions & languages over Σ

Let $\Sigma = \{a, b\}$. Then $a(ba)^*bb$ is a *regular expression* denoting

$$\begin{aligned} L &= \{a(ba)^nbb \mid n \in \mathbb{N}\} \\ &= \{abb, ababb, abababb, ababababb, \dots, a(ba)^nbb, \dots\} \end{aligned}$$

For general Σ the regular expressions over Σ are generated by

$$\text{rexp}_\Sigma ::= \emptyset \mid \lambda \mid s \mid \text{rexp}_\Sigma \text{ rexp}_\Sigma \mid \text{rexp}_\Sigma \cup \text{rexp}_\Sigma \mid \text{rexp}_\Sigma^*$$

with $s \in \Sigma$

This means $\emptyset \in \text{rexp}_\Sigma$, $\lambda \in \text{rexp}_\Sigma$, and $s \in \text{rexp}_\Sigma$ for $s \in \Sigma$ and

$$e_1, e_2 \in \text{rexp}_\Sigma \Rightarrow (e_1 \cup e_2) \in \text{rexp}_\Sigma$$

$$e_1, e_2 \in \text{rexp}_\Sigma \Rightarrow (e_1 e_2) \in \text{rexp}_\Sigma$$

$$e \in \text{rexp}_\Sigma \Rightarrow (e^*) \in \text{rexp}_\Sigma$$

For example $(abb)^*(a \cup \lambda)$ is a regular expression

Regular languages

For a regular expression e over Σ we define the language $L(e)$:

$$L(\emptyset) = \emptyset$$

$$L(\lambda) = \{\lambda\}$$

$$L(s) = \{s\}$$

$$L(e_1e_2) = L(e_1)L(e_2)$$

$$L(e_1 \cup e_2) = L(e_1) \cup L(e_2)$$

$$L(e^*) = L(e)^*$$

A language L is called *regular* if $L = L(e)$ for some $e \in \text{rexp}$

Examples

$L = \{w \mid bb \text{ occurs in } w\}$ over $\Sigma = \{a, b\}$ is regular:

$$L = L((a \cup b)^*bb(a \cup b)^*)$$

Also $L' = \Sigma^* - L = \{w \mid bb \text{ does not occur in } w\}$:

$$L' = L(a^*(baa^*)^* \cup a^*(baa^*)^*b)$$

To come

Regular Languages	Context-free Languages	Context-sensitive Languages	Enumerable Languages
Finite Automata	Push-down Automata	Bounded Turing Machines	Turing Machines
Regular Grammars	Context-free Grammars	Context-sensitive Grammars	General Grammars

Topics described in the last two columns are treated in other courses