

Talen en Automaten

Lecture 1: Regular Languages

Herman Geuvers



Outline

Organisation

Regular Languages



About this course I

Lectures

- Teachers: Jurriaan Rot (except today: Herman Geuvers)
- Weekly, 2 hours, on Tuesdays 8:45 – 10:30
- Presence not compulsory . . .
 - but active, polite attitude expected, when present
- The lectures follow:
 - these slides, available via the web
 - *Languages and Automata* by Alexandra Silva (LnA)
- Course URL for **all info (slides, exercises, schedule etc.)**:

<http://cs.ru.nl/~jrot/TnA2017/>

(Link exists in blackboard, under “Announcements”).

Check there first, before you dare to ask/mail a question!



About this course II

Exercises

- There are weekly exercises; the ones marked with **points** are to be handed in.
- Handing in is **compulsory**: To receive a grade for the course, you have to hand in **every** week.
- Exercises must be done individually
- Weekly exercise classes, on Fridays, 8:45 – 10:30 (for most of you), 10:45 - 12:30 (for Science students), and one class 15:45 – 17:30.
 - Presence not compulsory
 - Answers (for old exercises) & Questions (for new ones)
- Schedule:
 - New exercises on the web: Tuesday afternoon
 - Next exercise meeting (Friday) you can ask questions
 - Hand-in: **Tuesday before 8:45**, handwritten or typed, in the delivery boxes or by e-mail to the assistant of your group.



About this course III

Exercise Classes

Michiel de Bondt	8:45 – 10:30
Demian Janssen	8:45 – 10:30
Alexis Linard	8:45 – 10:30
David Venhoek	8:45 – 10:30
Tom van Bussel	8:45 – 10:30
• Leon Gondelmans	8:45 – 10:30
Ties Robroek	10:45 – 10:30 (Science)
Bas Steeg	10:45 – 10:30 (Science)
Jan Martens	15:45 – 17:30
Nienke Wessel	15:45 – 17:30

- see <http://cs.ru.nl/~jrot/TnA2017/exercises.html> for locations
- Please register for an exercise group via blackboard by **Tuesday November 14, 17:00**
- You will be assigned to an exercise class by **me**
- ~~Each assistant has a blue delivery box on the ground floor of the~~
Mercator 1 building



About this course IV

Examination

- There is a half-way test (Tuesday December 12, 8:30-10:30) and a final test (Wednesday January 22, 8:30-11:30).
- The final grade is composed of
 - the grade of your half-way test, **h**,
 - the grade of your final test, **f**,
 - the average grade of your exercises, **a**,
- Your final grade is $\min(10, \frac{f+h}{2} + \frac{a}{10})$
 - Additional requirements (everything at least 5) in study guide won't be applied
 - The re-exam is a full 3hrs exam about the whole course. You keep the (average) grade of the exercises.
- If you fail again, you must start all over next year (including re-doing new exercises, and additional requirements)



Languages and Automata

Let's start!



Overview

Topics

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

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

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

Examples

$$A_1 = \{a\}$$

$$A_2 = \{0, 1\}$$

$$A_3 = \{A, C, G, T\}$$

$$A_4 = \{a, b, c, d, \dots, x, y, z\}$$

$$A_5 = \{s \mid s \text{ is an ascii symbol}\}$$

$$A_6 = \{\text{あ、い、う、え、お、か、き、く、け、こ、...}\}$$

Japanese alphabet: 2×52 signs

$$A_7 = \{\text{山 川 日 雨 水 火 田, ...}\}$$

Chinese alphabet: 40.000 signs

$$A_8 = \{0, 1, +, \times, x_0, x_1, x_2, \dots\}$$

mathematical alphabet, countably infinite

$$A_9 = \{0, 1, +, \times, x_0, x_1, x_2, \dots\} \cup \{c_r \mid r \in \mathbb{R}\}$$

mathematical alphabet, uncountably infinite



Words

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

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

Inductive definition of the set of words, A^*

1. $\lambda \in A^*$ (λ denotes the empty word).
2. If $a \in A$ and $v \in A^*$, then $av \in A^*$.

Note that $a\lambda$ is just a

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

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

$$\begin{aligned}\lambda &= \text{---} \\ Eva &= \text{---}E\text{---}v\text{---}a\text{---}\end{aligned}$$

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



Operation on words

Inductive definition of the set of words, A^*

1. $\lambda \in A^*$ (λ denotes the empty word).
2. If $a \in A$ and $v \in A^*$, then $av \in A^*$.

Operations on words

$v \in A^*, u \in A^* \Rightarrow v \cdot u \in A^*$, concatenation

$v \in A^*, n \in \mathbb{N} \Rightarrow v^n \in A^*$, repetition

$v \in A^* \Rightarrow v^R \in A^*$, reverse

Inductive definitions of concatenation, repetition and reverse

$$\begin{array}{l} \lambda \cdot u = u \\ (av) \cdot u = a(v \cdot u) \end{array}$$

$$\begin{array}{l} v^0 = \lambda \\ v^{k+1} = v \cdot v^k \end{array}$$

$$\begin{array}{l} \lambda^R = \lambda \\ (av)^R = (v^R) \cdot a \end{array}$$

We write concatenation $v \cdot u$ as vu



Operation on words; Language

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

Examples (with $A = \{a, b\}$)

- $L_1 = \{w \in \{a, b\}^* \mid abba \text{ is a substring of } w\}$
- $L_2 = \{w \in \{a, b\}^* \mid w = w^R\}$



Examples of languages

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

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. $L_3 = \{a^n b^n c^n \mid n \geq 2\}$
4. $L_4 = \{a^n \mid n \in \mathbb{N} \text{ is prime}\}$

Over other alphabets:

1. $L_5 = \{n \mid n \text{ denotes an integer number}\}$
2. $L_6 = \{e \mid e \text{ is a well-formed arithmetical expression}\}$
3. $L_7 = \{P \mid P \text{ is a syntactically correct Java program}\}$
4. $L_8 = \{S \mid S \text{ is a grammatically correct English sentence}\}$



Operations on languages

Given languages $L_1, L_2, L \subseteq A^*$ we can define new languages:

$$L_1 \cup L_2 \quad L_1 \cap L_2 \quad \bar{L} \quad L_1 L_2 \quad L^*$$

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

$$L_1 \cap L_2 = \{w \mid w \in L_1 \text{ and } w \in L_2\}$$

$$\bar{L} = \{w \in A^* \mid w \notin L\}$$

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

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

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

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

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



Regular expressions

- Regular expressions are a way to **describe languages**.
- Really important concept in (theoretical) computer science
- Used a lot in text processing: search (efficiently!) for specific patterns



Example

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

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



Regular expressions and languages over an alphabet A

For general A the regular expressions over alphabet A are generated by

$$\text{rexp}_A ::= 0 \mid 1 \mid s \mid (\text{rexp}_A \text{ rexp}_A) \mid (\text{rexp}_A + \text{rexp}_A) \mid (\text{rexp}_A)^*$$

with $s \in A$

This means $0 \in \text{rexp}_A$, $1 \in \text{rexp}_A$, and $s \in \text{rexp}_A$ for $s \in A$ and

$$e_1, e_2 \in \text{rexp}_A \Rightarrow (e_1 + e_2) \in \text{rexp}_A$$

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

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

For example $(abb)^*(a + 1)$ is a regular expression



We economize on brackets

$$\text{rexp}_A ::= 0 \mid 1 \mid s \mid (\text{rexp}_A \text{ rexp}_A) \mid (\text{rexp}_A + \text{rexp}_A) \mid (\text{rexp}_A)^*$$

- We omit the outermost brackets,
- $*$ binds strongest,
- $+$ binds weakest.

So $a + ba^*$ denotes $((a + (b(a)^*)))$.

This denotes the language of either just a or b followed by a finite (possibly 0) number of a 's.



Regular languages

For a regular expression e over alphabet A we define the language $\mathcal{L}(e)$:

$$\mathcal{L}(0) = \emptyset$$

$$\mathcal{L}(1) = \{\lambda\}$$

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

$$\mathcal{L}(e_1 e_2) = \mathcal{L}(e_1) \mathcal{L}(e_2)$$

$$\mathcal{L}(e_1 + e_2) = \mathcal{L}(e_1) \cup \mathcal{L}(e_2)$$

$$\mathcal{L}(e^*) = (\mathcal{L}(e))^*$$

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



Examples

Let $A = \{a, b\}$.

- Also $L = \{w \mid w \text{ begins with } bb\}$ is regular

$$L = \mathcal{L}(bb(a + b)^*)$$

- $L = \{w \mid bb \text{ occurs in } w\}$ is regular

$$L = \mathcal{L}((a + b)^* bb(a + b)^*)$$

- $L = \{w \mid |w|_b \leq 2\}$ is regular

NB. $|w|$ denotes the **length of w** ,

$|w|_b$ denotes the **number of b 's in w**

$$L = \mathcal{L}(a^*(ba^*b + b + 1)a^*)$$

