

# Blatt 01: Reguläre Sprachen

## A1.1 Sprache von $a + a(a + b)^*a$

$$L = \{a\} \cup \{awa \mid w \in \{a,b\}^*\}$$

Es beschreibt alle Wörter, die entweder nur aus einem einzigen **a** bestehen oder mit **a** beginnen und auch enden. Dazwischen könnte auch beliebige Kombinationen aus **a** und **b** stehen.

**Beispiele:** a, aa, aba, abba

**Nicht enthalten:** b, ba, abb

## A1.2 Bezeichner in Programmiersprachen

A) Regex:

$(V|v|P|p|[A-Za-z])[A-Za-z0-9_]*[A-Za-z0-9]$

- Erstes Zeichen: V, v, P, p (ist eigentlich auch in a-z enthalten) oder ein Buchstabe (a-z, A-Z)
- Danach: Buchstaben, Ziffern oder Unterstriche erlaubt
- Letztes Zeichen: kein Unterstrich
- Länge: mindestens 2 Zeichen

**Beispiele:** va1, Vx\_2, Pa, name1

**Nicht erlaubt:** v\_, \_a, v

b) DFA

Alphabet:  $\Sigma = \{a \dots z, A \dots Z, 0 \dots 9, _\}$ . Zustände:  $Q = \{q_0, q_1, q_A, q_U, q_{\text{dead}}\}$ . Start:  $q_0$ . Endzustände:  $F = \{q_A\}$ .

**Intuition**  $q_1$ : genau ein Anfangsbuchstabe.  $q_A$ : gültiges Ende (mind 2 Zeichen, letztes Zeichen Buchstabe/Ziffer).  $q_U$ : letztes Zeichen unterstrich (kein gültiges Ende).  $q_{\text{dead}}$ : Fehler.

Übergangstabelle

	Letter	Digit	_	sonst
$q_0$	$q_1$	$q_{\text{dead}}$	$q_{\text{dead}}$	$q_{\text{dead}}$
$q_1$	$q_A$	$q_A$	$q_U$	$q_{\text{dead}}$
$q_A$	$q_A$	$q_A$	$q_U$	$q_{\text{dead}}$
$q_U$	$q_A$	$q_A$	$q_U$	$q_{\text{dead}}$
$q_{\text{dead}}$	$q_{\text{dead}}$	$q_{\text{dead}}$	$q_{\text{dead}}$	$q_{\text{dead}}$

## Beispielläufe

$\text{vCount1} : q_0 \xrightarrow{v} q_1 \xrightarrow{C} q_A \xrightarrow{o} q_A \xrightarrow{u} q_A \xrightarrow{n} q_A \xrightarrow{t} q_A \xrightarrow{1} q_A \text{ (akzeptiert)}$   
 $P_- : q_0 \xrightarrow{P} q_1 \rightarrow q_U \text{ (nicht akzeptiert)}$

## B) Reguläre Grammatik:

$$\begin{aligned} S &\rightarrow VA \mid vA \mid PA \mid pA \mid LA \\ A &\rightarrow XA \mid Y \\ X &\rightarrow a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z \mid 0 \mid 1 \mid \dots \mid 9 \mid \\ Y &\rightarrow a \mid b \mid \dots \mid z \mid A \mid B \mid \dots \mid Z \mid 0 \mid 1 \mid \dots \mid 9 \end{aligned}$$

Ableitung für  $va1$ :

$$S \Rightarrow vA \Rightarrow vXA \Rightarrow vaA \Rightarrow vaY \Rightarrow va1$$

## A1.3 Gleitkommazahlen in Python und Java

### Regex

Python:  $\wedge[0-9]+(\backslash.[0-9]+)?(e[+-]?[0-9]+)?\$$   
Java:  $\wedge[0-9]+(\backslash.[0-9]+)?([eE][+-]?[0-9]+)?\$$

### DFA

$z_0 \xrightarrow{\text{--digit--}} z_1$   
 $z_1 \xrightarrow{\text{--digit--}} z_1$   
 $z_1 \xrightarrow{\text{--.--}} z_2$   
 $z_2 \xrightarrow{\text{--digit--}} z_3$   
 $z_3 \xrightarrow{\text{--digit--}} z_3$   
 $z_1/z_3 \xrightarrow{\text{--e|E--}} z_4$   
 $z_4 \xrightarrow{\text{--+|--}} z_5$   
 $z_4/z_5 \xrightarrow{\text{--digit--}} z_6$   
 $z_6 \xrightarrow{\text{--digit--}} z_6$   
Endzustände:  $z_1$  (Ganzzahl),  $z_3$  (mit Nachkommastellen),  $z_6$  (mit Exponent)

**Reguläre Grammatik** Terminals: Ziffern 0..9, Punkt ., e, E, +, -. Nichtterminale: S, I, F, X, Y, Z.

$$\begin{aligned} S &\rightarrow dI \quad (d \text{ steht für eine Ziffer}) \\ I &\rightarrow dI \mid .F \mid eX \mid EX \mid \epsilon \\ F &\rightarrow dF \mid eX \mid EX \mid \epsilon \\ X &\rightarrow +Y \mid -Y \mid Y \\ Y &\rightarrow dZ \\ Z &\rightarrow dZ \mid \epsilon \end{aligned}$$

**bestatigung** 3.14, 2e10, 1.0E-5 werden akzeptiert.

## A1.4 Mailadressen

**Gegebener Regex:**  $(a-z)^+@(a-z).(a-z)$

### Problem:

- . steht für beliebiges Zeichen, nicht für Punkt.
- Vor und nach @ nur je ein Buchstabe, keine Wiederholung, keine Subdomains.
- Zeichenvorrat zu klein: nur a bis z, keine Ziffern und erlaubten Sonderzeichen.
- Keine Anker ^ und \$.

### Warum a + b + ... + z nicht reicht

- Beschreibt nur genau ein Zeichen; wiederholung fehlt.
- Ziffern und erlaubte Sonderzeichen fehlt.

### Verbesserte Regex

`^[A-Za-z0-9._%+~]+@[A-Za-z0-9.-]+\.[A-Za-z]{2,}$`

Beispiele: hallo@mail.de, moin.du@uni.edu

## A1.5 Der zweitletzte Buchstabe

$$\Sigma = \{1, 2, 3\}, \text{ Start } q_0, Q = \{q_0, r_1, r_2, r_3\} \cup \{s_{a,b,f}\}, F = \{s_{a,b,1}\}.$$

*Kurz:*  $r_x$  nach erstem Zeichen  $x$ . In  $s_{a,b,f}$  ist  $a$  das zweite Zeichen,  $b$  das letzte,  $f = 1$  falls zweitletztes  $= a$ .

Beispiele Gültig: 1221 (zweites = 2, zweitletztes = 2).

Ungültig: 1213 (zweites = 2, zweitletztes = 1).

## A1.6 Sprache einer regulären Grammatik

Gegeben:

$$S \rightarrow aA, \quad A \rightarrow dB \mid bA \mid cA, \quad B \rightarrow aC \mid bC \mid cA, \quad C \rightarrow \epsilon$$

Sprache:

- Start immer mit a (wegen  $S \rightarrow aA$ ).
- In  $A$  sind nur b oder c erlaubt (aber beliebig viele), ein d führt nach  $B$  daher mindestens ein d.
- In  $B$  endet a oder b sofort, c springt zurück nach  $A$  dadurch Blöcke  $c(b|c)^*d$  möglich.
- Ende ist a oder b, nie c.

also:

$$a(b|c)^*d(c(b|c)^*d)^*(a|b)$$

Beispiele: adb, acbcdb sind drin. ab, adc sind nicht drin.