

**Theoretische Informatik**

**Teil 2**

**Reguläre Ausdrücke**

Frühlingssemester 2019

L. Di Caro

D. Flumini

O. Stern



Viele Probleme in der Informatik beinhalten die **Prüfung, ob gewisse Wörter zu einer gegebenen Sprache gehören.**

- Ist eine Benutzereingabe sinnvoll?
- Gehört ein Wort zu einem Wörterbuch (Rechtschreibprüfung)?
- Bestimmte Folgen von “Events” sollen bestimmte Reaktionen auslösen (z. B. Getränkeautomat).

## Bemerkung

Da IT-Systeme nur über **endliche Speicherressourcen** verfügen, ist es daher wichtig Sprachen **endlich repräsentieren** zu können. **Unendlich grosse Sprachen** werden so einer **maschinellen Bearbeitung** überhaupt erst zugänglich gemacht.

Reguläre Ausdrücke sind Wörter, die Sprachen beschreiben, also eine Möglichkeit (gewisse) Sprachen endlich zu repräsentieren.

- Die **Syntax der regulären Ausdrücke** befasst sich mit der Frage, welche **Form** diese Wörter haben.
- In der **Semantik der regulären Ausdrücke** wird erklärt, wie man reguläre Ausdrücke als Sprachen **interpretiert**.

## Beispiel

Ein regulärer Ausdruck, der die Sprache aller Binärwörter der Länge 4 beschreibt:

$$\underbrace{(0|1)}_{\text{0 oder 1}} \underbrace{(0|1)}_{\text{nochmals}} \underbrace{(0|1)}_{\text{dreimal}} \underbrace{(0|1)}_{\text{genug}}$$

Ein passender regulärer Ausdruck ist also

$$(0|1)(0|1)(0|1)(0|1)$$

## Beispiel

Ein regulärer Ausdruck für die Sprache der Binärwörter, die das Teilwort 00 enthalten:

$\underbrace{(0|1)^*}_{\text{0 oder 1 beliebig oft}} \quad \underbrace{00}_{\text{das Teilwort}} \quad \underbrace{(0|1)^*}_{\text{0 oder 1 beliebig oft}}$

Ein passender regulärer Ausdruck ist also

$$(0|1)^*00(0|1)^*$$

## Definition (Reguläre Ausdrücke)

Es sei  $\Sigma$  ein beliebiges Alphabet. Die Sprache  $RA_{\Sigma}$  der **regulären Ausdrücke** über  $\Sigma$  ist wie folgt definiert:

- $\emptyset, \epsilon \in RA_{\Sigma}$
- $\Sigma \subset RA_{\Sigma}$
- $R \in RA_{\Sigma} \Rightarrow (R^*) \in RA_{\Sigma}$
- $R, S \in RA_{\Sigma} \Rightarrow (RS) \in RA_{\Sigma}$
- $R, S \in RA_{\Sigma} \Rightarrow (R|S) \in RA_{\Sigma}$

## Erläuterungen zur Definition

- Die Sonderzeichen  $\epsilon$  und  $\emptyset$  sind reguläre Ausdrücke.
- Jedes Symbol aus dem Alphabet  $\Sigma$  ist auch ein regulärer Ausdruck über  $\Sigma$ .
- Ist  $R$  ein regulärer Ausdruck über  $\Sigma$ , dann ist auch  $(R^*)$  ein regulärer Ausdruck über  $\Sigma$ .
- Sind  $R$  und  $S$  reguläre Ausdrücke über  $\Sigma$ , dann auch  $(RS)$  und  $(R|S)$ .

## Bemerkung

Es gibt unzählige Erweiterungen der Sprache der regulären Ausdrücke. Einige verbreitete abkürzende Schreibweisen sind:

- Ist  $R$  ein regulärer Ausdruck, dann steht  $(R^+)$  für  $R(R^*)$ .
- Ist  $R$  ein regulärer Ausdruck, dann steht  $(R?)$  für  $(R|\epsilon)$ .
- Sind  $R_1, \dots, R_k$  reguläre Ausdrücke, dann steht  $[R_1, \dots, R_k]$  für  $R_1|R_2|\dots|R_k$ <sup>1</sup>.

Diese abkürzenden Schreibweisen ( $R^+$ ,  $R?$  und  $R_1|R_2|\dots|R_k$ ) können in der Folge verwendet werden.

---

<sup>1</sup>Streng genommen müsste man hier  $R_1|(R_2|(\dots|R_k)\dots)$  schreiben.



## Beispiele

Einige reguläre Ausdrücke über dem Alphabet  $\{a, b\}$ .

- $(((((aa)^*)(b^*))(a(ba))))$
- $((a|(ab))^*)$
- $((ab)|(ba))$
- $(a(b(ba)))$

## Eigenschaften und Konventionen:

- Die Menge  $RA_{\Sigma}$  der regulären Ausdrücke über dem Alphabet  $\Sigma$  ist eine Sprache über dem Alphabet  $\{\emptyset, \epsilon, *, (, ), |\} \cup \Sigma$ .
- Der Lesbarkeit halber werden “überflüssige” Klammern weggelassen.
- Damit reguläre Ausdrücke auch mit (teilweise) weggelassenen Klammern eindeutig lesbar bleiben, gilt folgende Rangfolge der Operatoren:
  - “\*” vor “Konkatenation” und
  - “Konkatenation” vor “|”.

Der Ausdruck  $ab^*|c$  wird beispielsweise als  $((a(b^*))|c)$  gelesen.

## Definition (Die Sprache von regulären Ausdrücken)

Es sei  $\Sigma$  ein beliebiges Alphabet. Für jeden regulären Ausdruck  $R \in \text{RA}_\Sigma$  definieren wir die **Sprache**  $L(R)$  **von**  $R$  wie folgt:

- $L(\emptyset) = \emptyset$
- $L(\epsilon) = \{\epsilon\}$
- $L(a) = \{a\}$  für  $a \in \Sigma$
- $L(R^*) = L(R)^*$
- $L(R|S) = L(R) \cup L(S)$
- $L(RS) = L(R)L(S)$

## Erläuterungen zur Definition

- $\emptyset$  beschreibt die leere Sprache.
- $\epsilon$  beschreibt die Sprache  $\{\epsilon\}$ .
- Jedes Symbol  $a \in \Sigma$  beschreibt die Sprache  $\{a\}$ .
- $(R^*)$  beschreibt alle durch Konkatenation kombinierten Wörter, die von  $R$  beschrieben werden.
- $(R|S)$  beschreibt alle Wörter, die entweder von  $R$  oder von  $S$  beschrieben werden.
- $(RS)$  beschreibt die Wörter, die durch Konkatenation aus einem von  $R$  beschriebenen Wort gefolgt von einem durch  $S$  beschriebenen Wort entstehen.

## Definition

Eine Sprache  $A$  über dem Alphabet  $\Sigma$  heisst *regulär*, falls  $A = L(R)$  für einen regulären Ausdruck  $R \in \text{RA}_\Sigma$  gilt.

## Beispiele

- Für  $R = 0|(-?)[1, \dots, 9][0, \dots, 9]^*$  gilt

$L(R) =$  Menge der ganzen Zahlen in Dezimaldarstellung

- Für  $R = (0?)(10)^*(1?)$  oder  $R = (10)^*|(01)^*|(10)^*1|(01)^*0$  gilt

$L(R) =$  Menge der Binärwörter mit  
abwechselnd Nullen und Einsen

## Lemma (Rechenregeln für reguläre Ausdrücke)

*Für jedes Alphabet und alle regulären Ausdrücke  $R, S, T \in \text{RA}_\Sigma$  gelten folgende Identitäten.*

- $L(R|S) = L(S|R)$
- $L(R(ST)) = L((RS)T)$
- $L(R|(S|T)) = L((R|S)|T)$
- $L(R(S|T)) = L(RS|RT)$
- $L((R^*)^*) = L(R^*)$
- $L(R|R) = L(R)$

## Beweis.

Elementare Mengenumformungen. □

## **Anwendungen von regulären Ausdrücken:**

- Mustersuche in Texten
- Lexikalische Analyse (in Compilern); Erkennung von Schlüsselwörtern (“Token”)
- Syntax Test (bei einer einfachen Syntax)