

# Reguläre Sprachen, Endliche Automaten, und Reguläre Ausdrücke

**Programmieren und Software-Engineering  
Homomorphismen, Formale Sprachen und Syntax-Analyse**

22. Februar 2023

# Reguläre Ausdrücke: Theorie

## Regulärer Ausdruck über einer Menge $\Sigma$ von Zeichen

- ① jedes Zeichen  $z \in \Sigma$  ist ein regulärer Ausdruck
- ② sind  $r$  und  $s$  reguläre Ausdrücke, gilt
  - *Aneinanderreihung*:  $r s$  ist ein regulärer Ausdruck
  - *Alternative*:  $r \mid s$  ist ein regulärer Ausdruck:
  - *Wiederholung*:  $r^*$  ist ein regulärer Ausdruck:
  - *Klammerung*:  $(r)$  ist ein regulärer Ausdruck:

Beispiel:  $\Sigma = \{a,b\}$

- |                                      |                    |
|--------------------------------------|--------------------|
| ● $a, b, aa, ab, ba, abba, abbbaaba$ | ● $a^*b^*$         |
| ● $a b, abb ba$                      | ● $aa^*bb^*$       |
| ● $a^*, (a b)^*$                     | ● $a^*(a b)^+$     |
| ● $a(a b)^*b$                        | ● $(aa)^*(bb)^*b$  |
| ● $(a b)^*abb$                       | ● $(aa b)^*(a bb)$ |

# Praxis: Basic Regular Expressions Standard

$c$	dass Zeichen $c$	BRE
$c_1 c_2 \dots$	Zeichenfolge $c_1 c_2 \dots$	BRE
$.$	ein beliebiges Zeichen	BRE
$\backslash$	<i>escape character</i> : das folgende Zeichen ist <i>kein</i> Steuerzeichen	BRE
$\backslash .$	Punkt-Zeichen	BRE
$[c_1 c_2 \dots]$	eines der Zeichen $c_1, c_2 \dots$	BRE
$[c_1 - c_2]$	eines der Zeichen zwischen $c_1$ und $c_2$ (inkl.)	BRE
$[\sim c_1 c_2 \dots]$	jedes andere Zeichen als $c_1, c_2 \dots$	BRE
$[\sim c_1 - c_2]$	jedes andere Zeichen als eines zwischen $c_1$ und $c_2$	BRE
$\backslash w$	Buchstabe, Ziffer, $_$ : äquivalent $[A-Za-z0-9\_]$ $\backslash W$ : nicht ...	BRE
$\backslash s$	Whitespace $\backslash S$ : nicht ...	BRE
$^$	Zeilenbeginn	BRE
$\$$	Zeilenende	BRE
$\backslash b$	Wortgrenze: Wechsel $\backslash w \leftrightarrow \backslash W$ $\backslash B$ : nicht ...	BRE
$*$	<b>Wiederholung</b> beliebig oft inkl. 0-mal (lazy: $*?$ PCRE)	BRE

Praxis: **E**xtended und **P**erl **C**ompatible **R**egular **E**xpressions

<b>+</b>	Wiederholung, beliebig oft, mindestens 1-mal ( <i>lazy: +? PCRE</i> )	ERE
<b>{n}</b>	Wiederholung genau <i>n</i> -mal	ERE
<b>{m,n}</b>	Wiederholung <i>m</i> - bis <i>n</i> -mal	ERE
<b>{m,}</b>	mindestens <i>m</i> -mal	ERE
<b>{,n}</b>	höchstens <i>n</i> -mal	ERE
<b>?</b>	optional (0-mal oder 1-mal)	ERE
<b> </b>	oder (Alternative)	ERE
<b>\d</b>	Ziffer: äquivalent <b>[0-9]</b> <b>\D</b> : <i>nicht</i> ...	PCRE
<b>( ... )</b>	<b>Gruppierung</b> <sup>1</sup> und <i>back references</i> <b>\1</b> , <b>\2</b> , ...	ERE

<sup>1</sup>hauptsächlich für Ersetzungen

## Beispiele Zahlendarstellung

Ganze Zahl: 1, 123, +123, -123, 100, 001, 0, +0, -0

$[-+]?[0-9]^+$

ohne führende Nullen, inkl. Null: 1, 123, +123, -123, 100, 0

$[-+]?[1-9][0-9]^*|0$

Kommazahl: 3.14, +3.14, -3.14, 3.

$[-+]?[0-9]^+\backslash\.[0-9]^*$

ohne Vorkommateil: .14159

$[-+]?([0-9]^+\backslash\.[0-9]^*|\backslash\.[0-9]^+)$

mit Exponent: 3.14E12, 3.14e12, 3.14e-12, 3.14e+12

$[-+]?([0-9]^+\backslash\.[0-9]^*|\backslash\.[0-9]^+)([eE] [-+]?[0-9]^+)?$

## Beispiel: Aufbau einer Textzeile

z. B. `Meier, Paul.F.: DBI2_3_14:00-16:30_C3.15_`

- 1 Zeilenbeginn
- 2 Familienname (Buchstaben, 1. groß)
- 3 Beistrich
- 4 kein oder mehrere Leerzeichen
- 5 Vorname (Buchstaben, 1. groß)
- 6 ein oder mehrere Leerzeichen
- 7 2. Vorname Abkürzung (Großbuchstabe, Punkt)
- 8 Doppelpunkt
- 9 kein oder mehrere Leerzeichen
- 10 Kürzel (3 Großbuchstaben, Ziffer zwischen 1 und 4)
- 11 ein Leerzeichen
- 12 Wochentag als Nummer (1-5)
- 13 ein Leerzeichen
- 14 Uhrzeit von (als 2 Ziffern, Doppelpunkt, 2 Ziffern)
- 15 Bindestrich
- 16 Uhrzeit bis
- 17 ein oder mehrere Leerzeichen
- 18 Raum (Buchstabe A-D, Ziffer 1-5, Punkt, Ziffer, Ziffer 1-9)
- 19 kein oder mehrere Leerzeichen
- 20 Zeilenende

# Beispiel: Aufbau einer Textzeile (2)

z. B. `Meier, Paul F.: DBI2_3_14:00-16:30_C3.15_`

Zeilenbeginn	<code>^</code>
Familienname (Buchstaben, 1. groß)	<code>[A-Z] [a-z]+</code>
Beistrich	<code>,</code>
kein oder mehrere Leerzeichen	<code>_*</code>
Vorname (Buchstaben, 1. groß)	<code>[A-Z] [a-z]+</code>
ein oder mehrere Leerzeichen	<code>_+</code>
2. Vorname Abkürzung (Großbuchstabe, Punkt)	<code>[A-Z] \.</code>
Doppelpunkt	<code>:</code>
kein oder mehrere Leerzeichen	<code>_*</code>
Kürzel (3 Großbuchstaben, Ziffer zwischen 1 und 4)	<code>[A-Z] {3} [1-4]</code>
ein Leerzeichen	<code>_</code>
Wochentag (1-5)	<code>[1-5]</code>
ein Leerzeichen	<code>_</code>
Uhrzeit von (als 2 Ziffern, Doppelpunkt, 2 Ziffern)	<code>[0-9] {2} : [0-9] {2}</code>
Bindestrich	<code>-</code>
Uhrzeit bis	<code>[0-9] {2} : [0-9] {2}</code>
ein oder mehrere Leerzeichen	<code>_+</code>
Raum (Buchstabe A-D, Ziffer 1-5, Punkt, Ziffer, Ziffer 1-9)	<code>[A-D] [1-5] \. [0-9] [1-9]</code>
kein oder mehrere Leerzeichen	<code>_*</code>
Zeilenende	<code>\$</code>

## Beispiel: Aufbau einer Textzeile (3)

z. B. Meier, Paul.F.: DBI2\_3\_14:00-16:30\_C3.15\_

Regulärer Ausdruck für die ganze Zeile

```
^[A-Z][a-z]+, *[A-Z][a-z]+[A-Z]\.:  
*[A-Z]{3}[1-4]_[1-5]_[0-9]{2}:[0-9]{2}-[0-9]{2}:[0-9]{2}_+  
[A-D][1-5]\.[0-9][1-9]*$
```

Mockaroo: Generierung statt Match, d. h.  $* \Rightarrow \{0,9\}$ ,  $+ \Rightarrow \{1,9\}$

```
[[:upper:]][[:lower:]]{1,9},_{0,9}[[:upper:]][[:lower:]]{1,9}_{1,9}[[:upper:]]::  
_{0,9}[[:upper:]]{3}(1|2|3|4)_(1|2|3|4|5)_\d{2}:\d{2}-\d{2}:\d{2}_{1,9}  
(A|B|C|D)(1|2|3|4|5)_.\d(1|2|3|4|5|6|7|8|9)_{0,9}
```

## Zeichenklassen (BRE)

`[[:upper:]]`: Großbuchstaben, `[[:lower:]]`: Kleinbuchstaben, `[[:alpha:]]`: Buchstaben,  
`[[:digit:]]`: Ziffern, `[[:alnum:]]`: Buchstaben und Ziffern, `[[:blank:]]`: Leerzeichen oder  
Tabulator, `[[:punct:]]`: Sonderzeichen, ...



# Literaturübersicht I

- [1] **Berger, Krieger, Mahr: “Grundlagen der elektronischen Datenverarbeitung”, Skriptum**
- [2] Dirk W. Hoffmann: “Theoretische Informatik”, Hanser, 3. Auflage
- [3] Gernot Salzer: “Einführung in die Theorie der Informatik”, Skriptum, TU Wien, 2001
- [4] Wikipedia (Englisch): <https://en.wikipedia.org/>
- [5] Wikipedia (Deutsch): <https://de.wikipedia.org/>