

## Verwendung:

- Strings: bearbeiten, Suchen, ersetzen
- Zeichenketten validieren / auf Zusammensetzung prüfen
- Perl (hier sind die Reg Ex etwas besonderes)
- Bei web-Anwendungen (z.B. E-Mail Eingabefeld), in Shell-Skripten und in Editoren (Suchen und ersetzen) verwendet

## Ein-elementige Regex:

$[123456]$  } 1 oder 2 oder 3 ... oder 6  
 $[1-6]$

⇒ Eckige Klammer: Auflistung von Zeichen  
(immer nur einstellige Zeichen)<sup>c</sup>

$[1-35-9]$  Zwei Bereiche: 1-3 oder 5-9

Achtung: kein Leerzeichen zwischen 3 und 5

$[1-12]$  Achtung: 1 bis 12, also nicht 1 bis 1, da Zeichen einstellig sein müssen

$[a-d]$  a oder b oder c oder d  
 $[a-zA-d]$  a oder A oder ... oder d oder D

## Mehr - elementige Regex

$[1-9][a,b]$

$\Rightarrow$  1a, 1b, 2a, 2b, ..., 9a, 9b  
eine Ziffer 1-9 gefolgt von einem Buchstaben  
a oder b

$[1-9][a-zA-Z][2-4]$

$\Rightarrow$  Kombinationen möglich, Zeichenkette muss  
der Reihe nach alle 3 Regeln für seine  
3 Zeichen erfüllen.

Valid: 1c4, 9A3  
Nicht valide: 0c4, 1e4, 9A5

Achtung: Verhalten bei Zeichenkette 1c4H  
hängt vom Interpreter und vom  
Einsatzsachverhalt des Regex ab.

## Optionen:

$[1-9][0-9]?[0-9]?$  1 bis 3-stellige Zahl

$\Rightarrow$  ? besagt: Vorhergehendes Element kann  
vorkommen, muss aber nicht  
( $[0-9]$  ist ein Element)

$a \{ 5 \}$

$\Rightarrow a a a a a$  erfüllt diesen  
Regex

$\{ n \} = \wedge$  genau  $n$  mal

$a \{ 1, 3 \} b \Rightarrow a b, a a b, a a a b$  treffen  
zu

$\{ \min, \max \}$

$\Rightarrow$  Wenn nur eine mindest- oder maximal-  
anzahl angegeben wird, muss auch nur  
diese Bedingung zutreffen, z.B.

$[1-9][0-9]\{2,\}$

mindestens 3-stellig

$[1-9][0-9]\{,5\}$

maximal 6-stellig

## Beliebige Wiederholungen

Bsp.: Telefonnummer

0 7 1 3 1 1 9 6 3 3 - 6 8

0 1 8 0. 4 5. 3 3 3. 6 8

Die Länge der Telefonnummer sei an dieser  
Stelle egal, wichtig sind nur die Zeichen  
die darin vorkommen.

$[0-9/. \backslash -]^+$   $\leftarrow$  "mindestens 1 mal"

Leeresymbol  $\uparrow$

das "-" muss maskiert werden,  
da es hier als "." und nicht  
als z.B. "1-3" verwendet  
wird

=> Das + hinter der Beschreibung des in der Zeichenkette vorkommenden Zeichens besagt, dass mindestens ein Zeichen diese Regel erfüllen muss.

=> Ein \* verhält sich ähnlich, bis auf den Unterschied, dass auch kein Vorkommen erlaubt ist ("beliebig oft")

## Platzhalter

Bsp: Suche nach allen Personen, die "Thomas" mit Vornamen und "Müller" mit Nachnamen heißen, unabhängig vom Zweitnamen

Thomas \* Müller

Leertücken ↗ ↘

=> Der Punkt . steht für "ein beliebiges Zeichen", der Stern \* sagt aus, dass es beliebig oft vorkommen darf

Achtung: Zeilenumbrüche werden vom . nicht abgedeckt.  
Der Punkt muss innerhalb eckiger Klammern [ ] nicht maskiert werden.

## Negation

$[^{\wedge} a b]$

nicht (a oder b)

## Klammern

$\Rightarrow$  fassen zeichen zusammen

$B a (n a)^* n e$

$\Rightarrow$  Bane, Banane, Bananane, ...

$B a (n a) \{2, 5\} n e \Rightarrow$  Anzahl begrenzt

Thomas (Marvin) ? Müller

$\Rightarrow$  Thomas Müller, Thomas Marvin Marvin Müller

$B a (n | a) n e \Rightarrow$  Banne, Baane

$\hookrightarrow n$  oder  $a$

$\Rightarrow$  Der senkrechte Strich ist ein logisches Oder

Das Wetter ist (gut | schlecht)

## Modifikatoren

⇒ Das Verhalten des Regex kann durch Modifikatoren beeinflusst werden. Diese sind Dialekt - abhängig aber häufig ähnlich.

Syntax: [Begrenzung] [Regex] [Begrenzung] [Modifikation]

Bsp.: `/ (Ich | du) /i` (JavaScript)

- ▶ `i` Case - insensitive, Groß- und Kleinschreibung ignorieren.
- ▶ `s` Punkt wird multiline - fähig, das Wildcard - Zeichen `.` trifft auch auf Zeilenumbrüche zu.
- ▶ `m` `^` und `$` matchen auch auf Zeilenanfänge und -enden. Ohne diese Modifikation passen sie nur auf die gesamte Zeichenkette.

## Verschachtelung

`(VV (Golf | Polo) | Fiat (Panda | Panda))`

⇒ akzeptiert nur 4 verschiedene Kombinationen

$(10101) +$

=> Folge aus Nullen und Einsen, in der maximal 2 Nullen oder Einsen aufeinander folgen

## Geringe Ausdrücke

Wir wollen HTML-Tags erkennen, die links sind, also folgende Form haben:

$\langle a \text{ href} = "https: // \dots" \rangle$

Vorschlag:  $\langle a \text{ href} = ".*".* \rangle$

Problem: Der zweite Punkt frisst die Klammer-Zu  $\rangle$  ebenfalls und dazu noch die Zeichen dahinter.

Lösung 1: Punkt ersetzen

$\langle a \text{ href} = "[^"]*" [^>]* \rangle$

Lösung 2: Punkt und Stern „bändigen“

$\langle a \text{ href} = ". * ? ", * ? \rangle$

=> Das Fragezeichen sagt an dieser Stelle aus, dass nur so viele Zeichen wie nötig vom Punkt abgedeckt werden sollen.

## Gruppen

Wir wollen aus dem vorherigen Beispiel die URL extrahieren.

$\langle a \text{ href} = "(.*)".*? \rangle$   
~~~~~  
↳ Gruppe

Je nach verwendetem RegEx Dialekt kann der Inhalt der Gruppe ausgelesen werden.

Achtung: Die Nummerierung der Gruppen erfolgt der Reihenfolge der öffnenden Klammern nach

Referenzen:      Verweise auf Gruppen

$([0-9])[0-9]^*\backslash1$

~~~~~  
↑ Verweis auf Inhalt der 1. Klammer

⇒ Das Zeichen, das in der 1. Klammer gefunden wurde, muss sich an der Stelle  $\backslash 1$  auch befinden

$\backslash 0$  ist der match des gesamten RegEx



## Wortgrenzen

Bsp.: Suche nach "kai" in:

"kai fährt nach Kaiserslautern"

kai  $\Rightarrow$  2 Matches

\b kai \b  $\Rightarrow$  1 Match



Eine Wortgrenze tritt zwischen einem Wort-Zeichen und einem Nicht-Wort-Zeichen auf

## Spezial - Zeichen

Immer ein Backslash und ein Buchstabe.

Großbuchstaben stehen immer für das Gegenteil des Kleinbuchstaben

\w Wort-Zeichen, entspricht [a-zA-Z0-9]

\W Nicht-Wort-Zeichen, [^a-zA-Z0-9]

\d Ziffer („digit“) [0-9]

\D [^0-9]

\b Wortende

\B kein Wortende

\s whitespace-Zeichen: LF, CR, Space, Tab

\S kein whitespace-Zeichen

\\ Backslash-Zeichen („\“)

Der Backslash wird zum Markieren von Zeichen verwendet:

\. \+ \\* \[ \> \[ \]

\- \& \|

## Anfangs- und End-Zeichen

$\wedge$  1 \. 3 \. 20 04 \$  
~  
↑ Match nur, wenn der String hier anfängt      ~  
↑ Match nur, wenn der String hier endet

1. 3. 2004  $\Rightarrow$  Match

11. 3. 2004  $\Rightarrow$  kein Match

\d\$  $\Rightarrow$  Match Strings die auf Ziffern enden

^\d  $\Rightarrow$  Match Strings die mit " " anfangen

## Positive Lookaheads und Lookbehinds

Bsp.: Lange Ziffern-Folge

0056344087500033205791

$\Rightarrow$  Suche alle Folgen, die aus drei Ziffern ungleich 0 bestehen und von einer 0 an jeder Seite begrenzt werden.

Lösung:  $(? <= 0) [1-9] \{3\} (? = 0)$

Erklärung:

► Mitte:  $[1-9] \{3\}$

Ziffernfolge von drei Ziffern ungleich Null

► Links:  $(? <= 0)$

$? <=$   $\Rightarrow$  Sonderfunktion für „Vornedran“

Dies ist ein sog. „positive lookahead“  
(entgegen dem normalen Lese-Fluss)

$? <= 0$   $\Rightarrow$  Vornedran muss eine 0 sein

► Rechts:  $(? = 0)$

$? =$   $\Rightarrow$  Sonderfunktion für Hindendran

Dies ist ein sog. „positive lookahead“

$? = 0$   $\Rightarrow$  Hindendran muss eine 0 sein

$\setminus b$   $(\setminus d)$   $\setminus d^*$   $\setminus 1$   $\setminus b$

kann auch ausgedrückt werden als

$(\{ \geq \setminus w \}) (\setminus d) \setminus d^* \setminus 1 (\{ \geq \setminus w \})$   
↑

Achtung: lookbehinds und lookahead  
sind keine Gruppen, fallen  
also aus der Nummerierung heraus

Achtung: lookbehinds dürfen keine unbekannte  
Anzahl an Zeichen haben, also kein:

$^*$ ,  $+$ ,  $\{0, 4\}$

Negative lookahead und lookbehinds

$(\{ \geq ! 0 \}) \Rightarrow$  negativer lookahead

hier darf keine 0 folgen

Bsp.:  $\setminus b F (\{ \geq ! e \} a \setminus b) \cdot^* \setminus b$

$\Rightarrow$  alle Wörter, die mit F anfangen,  
aber nicht „feta“ sind

$(? < ! 0)$

$\Rightarrow$  negativer Lookbehind

Hier darf keine Null davor stehen

Bsp.:  $^*(? < ! \text{Müll})$  Eimer

$\Rightarrow$  alle Eimer, die keine Müll-eimer sind