

Reguläre Ausdrücke in Java Bestandteile

Zeichen

Quantoren (Quantifizierer)

Grenzen

Alternativen

Gruppen

Ungespeicherte Gruppen

Rückwärtsreferenzen

Pattern und Matcher

Matcher-Suchmethoden

Matcher-Ersetzungsmethoden

String-Methoden und reguläre Ausdrücke

Reguläre Ausdrücke in Java

Mit einem einzigen **regulären Ausdruck** kann man eine Menge unterschiedlicher Texte beschreiben, die **ähnlich aussehen**, z.B.

- ◆ Telefonnummern
- ◆ Mail-Adressen
- ◆ Datumsangaben
- ◆ Uhrzeiten

Reguläre Ausdrücke dienen dazu, solche Texte

- ◆ zu untersuchen,
- ◆ in ihre Einzelteile zu zerlegen oder
- ◆ zu verändern (suchen und ersetzen)

Es gibt verschiedene **Dialekte** für reguläre Ausdrücke, die in Einzelheiten voneinander abweichen, wie etwa

- ◆ Portable Operating System Interface-kompatibel (POSIX-kompatibel)
- ◆ Perl Compatible Regular Expressions (PCRE)

Der Begriff „Regulärer Ausdruck“ wird auch zu **Regex** oder **Regexp** abgekürzt (englisch: regular expression).

Bestandteile

Zeichen

Zeichen(folge)	passt zu
Zeichenliterale (Auswahl)	
x	dem selben Zeichen ' x '
Den folgenden Metazeichen: <code>[] () { } ? + - * ^ \$ \ .</code> muss man jedoch das Zeichen ' <code>\</code> ' voranstellen, um sie als Literale zu verwenden.	
<code>\n</code>	dem Zeilenvorschubzeichen ' <code>\n</code> '
<code>\t</code>	dem Tabulatorzeichen ' <code>\t</code> '
Zeichenklassen (Auswahl)	
<code>.</code> (ein Punkt)	einem beliebigen Zeichen (meist außer ' <code>\n</code> ')
<code>[abc]</code>	einem der Zeichen ' <code>a</code> ', ' <code>b</code> ' oder ' <code>c</code> '
<code>[^abc]</code>	jedem Zeichen außer ' <code>a</code> ', ' <code>b</code> ' oder ' <code>c</code> '
<code>[a-zA-Z]</code>	' <code>a</code> ' bis ' <code>z</code> ', ' <code>A</code> ' bis ' <code>Z</code> ' (meist ASCII-Zeichensatz)
<code>\d</code>	einer Ziffer: <code>[0-9]</code>
<code>\D</code>	allem, was keine Ziffer ist: <code>[^0-9]</code>
<code>\s</code>	Leerraum: <code>[\t\n\x0B\f\r]</code>
<code>\S</code>	allem, was kein Leerraum ist: <code>[^\s]</code>

Quantoren (Quantifizierer)

Quantoren geben an, wie oft das davor stehende Zeichen/die davor stehende [Gruppe](#) vorkommen kann/muss (Default: einmal).

Quantor	Anzahl der Wiederholungen
Gierig (wiederholt möglichst oft, steckt zurück wenn nötig)	
?	einmal oder überhaupt nicht („optional“)
*	beliebig oft oder auch überhaupt nicht
+	beliebig oft, aber mindestens einmal
{ <i>n</i> }	genau <i>n</i> Mal
{ <i>n</i> , }	mindestens <i>n</i> Mal
{ <i>n</i> , <i>m</i> }	mindestens <i>n</i> bis höchstens <i>m</i> Mal
Zurückhaltend (wiederholt möglichst selten, weitet aus wenn nötig)	
??	wie ?
*?	wie *
USW.	
Besitzergreifend (wiederholt möglichst oft und steckt nicht mehr zurück)	
?+	wie ?
*+	wie *
USW.	

Grenzen

Diese prüfen, ob eine bestimmte **Position** (kein Zeichen!) innerhalb des Textes erreicht ist.

Grenze	Position
<code>^</code>	Text-Anfang (unterscheide von: <code>[^...]</code>)
<code>\$</code>	Text-Ende bzw. Zeilenende (vor <code>'\n'</code>)
<code>\b</code>	Wortgrenze (Anfang oder Ende eines Wortes)
<code>\B</code>	jede Position außer einer Wortgrenze

Alternativen

Um zu überprüfen ob **einer von mehreren unterschiedlichen Ausdrücken** im Text vorliegt, verbindet man diese Ausdrücke mit dem Zeichen '| '.

Beispiel:

Regulärer Ausdruck	passt zu
Montag Dienstag Mittwoch	"Montag", "Dienstag" oder "Mittwoch"

Gruppen

Gruppen werden mit runden Klammern () notiert und dienen dazu,

1. mehrere **Teile eines regulären Ausdrucks** zu einer Einheit **zusammenzufassen**, auf die man dann z.B. einen [Quantor](#) anwenden kann, oder/und
2. die **zur Gruppe passenden Teile des Textes zu speichern**, damit man sie an anderer Stelle wieder verwenden kann.

Gruppen werden – beginnend bei 1 – in der Reihenfolge der öffnenden Klammern **nummeriert**. Unter dieser Nummer findet man später den gespeicherten Inhalt.

Die Gruppe 0 entspricht demjenigen Teil des Textes, der zum gesamten regulären Ausdruck passt.

Beispiele:

Regulärer Ausdruck	Text	Gruppe		
		0	1	2
<code>((ab)*)c+</code>	<code>abc</code>	<code>abc</code>	<code>ab</code>	<code>ab</code>
	<code>ababccc</code>	<code>ababccc</code>	<code>abab</code>	<code>ab</code>
	<code>c</code>	<code>c</code>	<code>--leer--</code>	<code>--leer--</code>
	<code>xyz</code>	<code>--leer--</code>	<code>--leer--</code>	<code>--leer--</code>

Ungespeicherte Gruppen

Zweck: mehrere Teile eines regulären Ausdrucks zu einer Einheit zusammenfassen, **ohne dass** die dazupassenden Teile des Textes **gespeichert werden** (wie es bei gewöhnlichen [Gruppen](#) der Fall ist).

Ungespeicherte Gruppen notiert man in der Form `(?:)`

Beispiele:

Regulärer Ausdruck	Text	Gruppe	
		0	1
<code>((?:ab)*)c+</code>	abc	abc	ab
	ababccc	ababccc	abab
	c	c	--leer--
	xyz	--leer--	--leer--

Rückwärtsreferenzen

Diese erlauben es, den zuvor ermittelten **Inhalt einer Gruppe**

♦ später im regulären Ausdruck **wieder zu verwenden**

♦ in einem **Ersetzungsausdruck** zu verwenden

Rückwärtsreferenzen werden als $\backslash n$ notiert, wobei n die Nummer der Gruppe ist.

Beispiele:

Regulärer Ausdruck	passt zu
<code>(' ") . + \1</code>	(fast) jedem Stringliteral in einfachen oder doppelten Hochkommas, z.B. <code>"abc"</code> oder <code>'abc'</code> , aber nicht zu <code>'abc"</code>
<code><(\w+)>(.*<\/\1></code>	einem zusammengehörenden Paar von öffnendem und schließendem XML-Tag, z.B. <code><h1>Titel</h1></code> oder <code>fett</code> . Der Tag-Inhalt wird in Gruppe 2 gespeichert.

Pattern und Matcher

Um in Java mit einem regulären Ausdruck zu arbeiten, geht man meist in **drei Schritten** vor:

1. Man läßt den regulären Ausdruck (der ja eine Art von „Programm“ darstellt) in ein **Pattern**-Objekt **kompilieren**. Dies ist zeitaufwändig, daher sollte man jeden regulären Ausdruck **nur einmal kompilieren**.

Dazu dient die statische Methode **Pattern.compile()**, der man den regulären Ausdruck als **String**-Argument übergibt. Ungültige reguläre Ausdrücke bewirken eine **PatternSyntaxException**.

ACHTUNG: wie in jedem Java-**String** muss auch in regulären Ausdrücken das Zeichen '****' verdoppelt werden!

2. Für den zu untersuchenden Text erzeugt man aus diesem **Pattern** einen **Matcher**. Dazu dient die **Pattern**-Methode **matcher()**.
3. Man verwendet geeignete **Matcher-Methoden**, um den Text zu untersuchen oder zu bearbeiten.

Matcher-Suchmethoden

1. [matches\(\)](#): ganzen Text vergleichen

Vergleicht den Text mit dem `Pattern` und liefert `true` bei **vollständiger** Übereinstimmung.

2. [lookingAt\(\)](#): Beginn des Textes vergleichen

Vergleicht den Text mit dem `Pattern` und liefert `true`, wenn der Text **vom Anfang an** (aber nicht notwendigerweise bis zum Ende) mit dem `Pattern` übereinstimmt.

3. [find\(\)](#): nächsten passenden Teil des Textes suchen

Sucht den **nächsten Teil** des Textes, der mit dem `Pattern` übereinstimmt und liefert `true`, wenn gefunden. Mit Hilfe der Methoden [start\(\)](#) und [end\(\)](#) kann man die folgenden passenden Textteile finden.

4. [group\(\)](#): gespeicherte Gruppeninhalte

Liefert den Inhalt einer gespeicherten [Gruppe](#), nachdem `matches()`, `lookingAt()` oder `find()` erfolgreich ausgeführt wurde.

1. [replaceAll\(\)](#)

Liefert einen neuen **String**, in dem alle Teile des Textes, die mit dem **Pattern** übereinstimmen, durch den angegebenen **Ersetzungsausdruck** ersetzt sind. Im Ersetzungsausdruck dürfen auch [Rückwärtsreferenzen](#) vorkommen.

2. [replaceFirst\(\)](#)

Wie **replaceAll()**, jedoch wird nur der allererste zum **Pattern** passende Teil ersetzt.

String-Methoden und reguläre Ausdrücke

Reguläre Ausdrücke können auch in **String**-Methoden verwendet werden. Der reguläre Ausdruck muss aber bei jedem Aufruf neu kompiliert werden (⇒ langsamer).

String-Methode	entspricht
<u>matches</u> (<i>regex</i>)	Pattern.compile(<i>regex</i>).matcher(<i>string</i>).matches()
<u>replaceFirst</u> (<i>regex</i> , <i>ersetzung</i>)	Pattern.compile(<i>regex</i>).matcher(<i>string</i>).replaceFirst(<i>ersetzung</i>)
<u>replaceAll</u> (<i>regex</i> , <i>ersetzung</i>)	Pattern.compile(<i>regex</i>).matcher(<i>string</i>).replaceAll(<i>ersetzung</i>)
<u>split</u> (<i>regex</i> , <i>limit</i>)	Pattern.compile(<i>regex</i>).split(<i>string</i>, <i>limit</i>)