

Inhaltsverzeichnis

1	Einführung	2
2	Beispiele für reguläre Ausdrücke und ihre Anwendung	2
2.1	Überprüfen einer MAC-Adresse auf Korrektheit	2
2.2	Finden von MAC-Adressen in einer Zeichenkette	3
3	Einführung in die Syntax regulärer Ausdrücke (RegexSyntax)	4
3.1	String Literale (wortwörtlicher Text) in regulären Ausdrücken	4
3.2	Zeichenklassen in regulären Ausdrücken	5
3.3	Wiederholungen von Zeichen(klassen) in regulären Ausdrücken (Quantitätsangaben)	6
3.4	Aneinanderreihung und Gruppenbildung in regulären Ausdrücken	6

1 Einführung

Das Durchsuchen von Texten nach bestimmten Wörtern oder Zeichenmustern (engl. Pattern) ist eine Standardaufgabe in der Programmierung. Möchte man z.B. alle E-Mail Adressen suchen, die in einer Webseite vorkommen, dann benötigt man eine Möglichkeit, den generellen Aufbau (das Zeichenmuster) einer E-Mail Adressen zu beschreiben. Zur Definition von Zeichenmustern gibt es eine eigene formale Sprache. Diese Metasprache nennt sich „reguläre Ausdrücke“ und wurde von dem Mathematiker Stephen Cole Kleene entwickelt. Eine Besonderheit dieser Metasprache besteht darin, dass es möglich ist, aus einem definierten Zeichenmuster (dem regulären Ausdruck) automatisch einen Algorithmus zu konstruieren. Die Java Bibliothek besitzt, wie viele andere moderne Programmiersprachen, entsprechende Methoden, um reguläre Ausdrücke nutzen zu können.

2 Beispiele für reguläre Ausdrücke und ihre Anwendung

2.1 Überprüfen einer MAC-Adresse auf Korrektheit

```

public static void main(String[] args)
{
    String sMacAdr; // Eingabe-Zeichenkette
    // Regulärer Ausdruck (Regex) für MAC-Adr:
1   String sRegexMac = "([0-9A-Fa-f]{2}\\-){5}[0-9A-Fa-f]{2}";
    // Generiere aus dem "regulären Ausdruck" den "Matcher"-Algorithmus
2   Pattern matchAlgorithm = Pattern.compile(sRegexMac);

    sMacAdr = Eingabe.getString("Enter MAC Adr: ");

    // Verbinde den "Matcher"-Algorithmus mit der Eingabe Zeichenkette
3   Matcher matcher = matchAlgorithm.matcher(sMacAdr);
    // Entspricht die Zeichenkette dem Zeichenmuster?
4   if (matcher.matches()) // Wenn Zeichenmuster -> true
        System.out.println("Adr. o.k.");
    else
        System.out.println("Adr. nicht o.k.");
}

```

1	Das Zeichenmuster („Pattern“) einer MAC-Adresse kann durch einen „regulären Ausdruck exakt beschrieben werden.
2	Aus dem regulären Ausdruck (Regex = „regular expression“) kann der Algorithmus zum Auffinden von Zeichenketten, die dem Zeichenmuster entsprechen, automatisch generiert werden („Matcher“-Algorithmus)
3	Die Eingabe-Zeichenkette wird mit Matcher-Algorithmus verknüpft
4	<code>match()</code> sucht in der Eingabe-Zeichenfolge nach dem nächsten Vorkommen des Zeichenmusters. Falls die Zeichenkette als ganzes dem Zeichenmuster genügt, liefert <code>match()</code> <code>true</code> zurück, sonst <code>false</code>

Wird, wie im Beispiel der MAC Adressüberprüfung, der Matcher-Algorithmus nur einmalig angewendet, so können mehrere Teilschritte zu einem zusammengefasst werden; zu diesem Zweck kann die Klassenmethode `Pattern.matches(...)` verwendet werden.

```
public static void main(String[] args)
{
    String sMacAdr; // Eingabe-Zeichenkette
    // Regulärer Ausdruck (Regex) für MAC-Adr:
1   String sRegexMac = "([0-9A-Fa-f]{2}\\-){5}[0-9A-Fa-f]{2}";

2   sMacAdr = Eingabe.getString("Enter MAC Adr: ");
    // Generiere aus dem "regulären Ausdruck" den "Matcher"-Algorithmus
    // Verbinde den "Matcher"-Algorithmus mit der Eingabe Zeichenkette
    // und prüfe, ob die Zeichenkette dem Zeichenmuster entspricht?

3   if (Pattern.matches(sRegexMac, sMacAdr))
        System.out.println("Adr. o.k.");
4   else
        System.out.println("Adr. nicht o.k.");
}
```

2.2 Finden von MAC-Adressen in einer Zeichenkette

```
public static void main(String[] args)
{
    boolean found = false;
    String sMacAdr; // Eingabe-Zeichenkette
1   String sRegexMac = "([0-9A-Fa-f]{2}\\-){5}[0-9A-Fa-f]{2}";
    // Generiere aus dem "regulären Ausdruck" den "Matcher"-Algorithmus
2   Pattern matchAlgorithm = Pattern.compile(sRegexMac);
    // Zeichenkette aus einer DHCP Konfigurationsdatei
    sMacAdr = "A205A;00-56-71-36-AF-15;192.168.75.110;6;7.12.2019"
              + "A205B;00-56-71-36-B3-B2;192.168.75.111;6;7.12.2019"
              + "A205C;00-56-89-41-00-00;192.168.75.112;6;7.12.2019";

3   Matcher matcher = matchAlgorithm.matcher(sMacAdr);
    // Suchen nach MAC Adressen in der Zeichenkette
4   while (matcher.find()) {
        System.out.printf("\nI found the text \"%s\" starting at " +
                          "index %3d and ending at index %3d.%n",
5        matcher.group(), matcher.start(), matcher.end());
        found = true;
    }
    if(!found){
        System.out.printf("No match found.%n");
    }
}
```

Ausgabe:

```
I found the text "00-56-71-36-AF-15" starting at index   6 and ending at index  23.
I found the text "00-56-71-36-B3-B2" starting at index  56 and ending at index  73.
I found the text "00-56-89-41-00-00" starting at index 106 and ending at index 123.
```

1	Das Zeichenmuster („Pattern“) einer MAC-Adresse kann durch einen „regulären Ausdruck exakt beschrieben werden.
2	Aus dem regulären Ausdruck (Regex = „regular expression“) kann der Algorithmus zum Auffinden von Zeichenketten, die dem Zeichenmuster entsprechen, automatisch generiert werden („Matcher“-Algorithmus)
3	Die Eingabe-Zeichenkette wird mit Matcher-Algorithmus verknüpft
4	<code>find()</code> sucht in der Eingabe-Zeichenfolge nach dem nächsten Vorkommen des Zeichenmusters. Falls in der Zeichenkette ein Teilstring vorkommt, der dem Zeichenmuster genügt, liefert <code>find()</code> <code>true</code> zurück, sonst <code>false</code>
5	<code>group()</code> liefert den durch <code>find()</code> gefundenen Teilstring. <code>start()</code> liefert den Startindex des gefundenen Teilstrings und <code>end()</code> entsprechenden Endindex.

3 Einführung in die Syntax regulärer Ausdrücke (RegexSyntax)

Um alle möglichen Fälle für Zeichenmuster abzudecken, benötigt man auch eine entsprechend umfangreiche Syntax für die Formulierung von regulären Ausdrücken. Diese komplett hier zu behandeln, würde den Rahmen sprengen. Wenn man die Syntax nicht häufig anwendet, vergisst man sie auch relativ schnell wieder, deshalb ist es sinnvoller sich bei konkretem Bedarf intensiver damit zu beschäftigen. Da die Anwendung von regulären Ausdrücken aber in der Praxis von großer Bedeutung ist, sollen hier die Grundlagen gelegt werden.

3.1 String Literale (wortwörtlicher Text) in regulären Ausdrücken

Ein Zeichenmuster kann auch wortwörtlichen Text (Literale) enthalten.

Bsp.:

```
String sText = "Hallo Welt!"; // Eingabe-Zeichenkette
// Regulärer Ausdruck bestehend aus Literal:
String sRegex = "Welt";

if (Pattern.compile(sRegex).matcher(sText).find())
    System.out.println("gefunden"); // "Welt" kommt in "Hallo Welt!" vor
else
    System.out.println("nicht gefunden");
}
```

Achtung: Kommen in dem gesuchten wortwörtlichen Text sogenannte Metazeichen vor, so müssen diese mit einem `\` maskiert werden. Metazeichen sind solche Zeichen, die in der Regex Syntax eine besondere Bedeutung haben.

Die Metazeichen, die in Java Bibliothek unterstützt werden sind folgende:

`<([{\^-= $! |]})? * + . >`

Bsp.:

```
Pattern.compile("Welt.").matcher("Hallo Welt!").find() // → true
// "Welt." kommt in "Hallo Welt!" nicht vor, wird aber trotzdem gefunden
// weil der Punkt ein Metazeichen ist, dass für ein beliebiges Zeichen steht
Pattern.compile("Welt\\.").matcher("Hallo Welt!").find() // → false
```

Möchte man in einem regulären Ausdruck innerhalb eines Literals auch Steuerzeichen verwenden, so werden diese gleich dargestellt, wie wir das bereits von normalen Stringliteralen kennen:

<code>\\</code>	The backslash character
<code>\t</code>	The tab character (' <code>\u0009</code> ')
<code>\n</code>	The newline (line feed) character (' <code>\u000A</code> ')
<code>\r</code>	The carriage-return character (' <code>\u000D</code> ')
<code>\f</code>	The form-feed character (' <code>\u000C</code> ')
<code>\a</code>	The alert (bell) character (' <code>\u0007</code> ')
<code>\e</code>	The escape character (' <code>\u001B</code> ')

3.2 Zeichenklassen in regulären Ausdrücken

In einem regulären Ausdruck kann man festlegen, dass an einer bestimmten Stelle im Zeichenmuster nur bestimmte Zeichen stehen dürfen. Dazu verwendet man Zeichenklassen.

Vordefinierte Zeichenklassen

<code>.</code>	Ein beliebiges Zeichen
<code>\d</code>	Beliebiges Ziffernzeichen; identisch mit <code>[0-9]</code>
<code>\D</code>	Alle Zeichen außer Ziffern; identisch mit <code>[^0-9]</code>
<code>\s</code>	Leerraumzeichen (White Space); identisch mit <code>[\t\n\r]</code>
<code>\S</code>	Alle Zeichen außer Leerraum; identisch mit <code>[^\s]</code>
<code>\w</code>	Jedes alphanummerisches Zeichen; identisch mit <code>[a-zA-Z_0-9]</code>
<code>\W</code>	Alle Zeichen außer Jedes alphanummerische; identisch mit <code>[^\w]</code>

Bsp.: (zur Klassenmethode `Pattern.matches` s.o. 2.1)

```
if (Pattern.matches("\\d\\d\\d", "120")) // → true
    System.out.println("Zeichenkette besteht aus 3 Ziffern");
else
    System.out.println("falsch");
```

Bsp.:

```
Pattern.matches("\\d\\d\\d", "20") // → false, es werden 3 Dezimalziffern erwartet
```

Selbstdefinierte Zeichenklassen

<code>[abc]</code>	a, b oder c, (Einfache Klasse)
<code>[^abc]</code>	Alles außer a, b oder c (Negierung)
<code>[a-zA-Z]</code>	a bis z oder A bis Z (einfacher Zeichenbereich)
<code>[a-m[n-p]]</code>	a bis m oder n bis p; identisch mit <code>[a-mm-p]</code> bzw. <code>[a-p]</code> (Vereinigung)
<code>[a-z&&[^bc]]</code>	a bis z, ohne b und ohne c (Zeichenbereich mit Ausschluss)

3.3 Wiederholungen von Zeichen(klassen) in regulären Ausdrücken (Quantitätsangaben)

In einem regulären Ausdruck kann man festlegen, dass an einer bestimmten Stelle im Zeichenmuster ein bestimmtes Zeichen, ein Zeichen einer Zeichenklasse oder ein Teilmuster mehrfach stehen darf. (X steht für ein beliebiges Teilmuster (Teilausdruck))

Quantitäten

$X?$	X genau einmal oder überhaupt nicht
X^*	X 0 Mal bis viele
X^+	X 1 Mal bis viele
$X\{n\}$	X genau n Mal
$X\{n, \}$	X mindestens n Mal bis viele
$X\{n, m\}$	X n bis m Mal

Bsp.: (Alternative Lösung zum vorletzten Bsp.)

```
if (Pattern.matches("\\d{3}", "120")) // → true
    System.out.println("Zeichenkette besteht aus 3 Ziffern");
else
    System.out.println("falsch");
```

Bsp. split-Methode:

```
sZeilenTeile = sZeile.split("[ \t]+");
```

Der Parameter der `split`-Methode ist ebenfalls ein regulärer Ausdruck und im Bsp. folgendermaßen zu interpretieren:

`[\t]` → selbstdefinierte Zeichenklasse: Leerzeichen oder Tab-Zeichen

`[\t]^+` → Quantität `+`: Zeichenmuster beschreibt eine beliebige Folge von Leer- und Tab-Zeichen, mindestens aber ein Zeichen

3.4 Aneinanderreihung und Gruppenbildung in regulären Ausdrücken

Aneinanderreihung und Gruppierung

XY	X gefolgt von Y
$X Y$	Entweder X oder Y
(X)	X als Gruppe mit späterer Bezugsmöglichkeit
$(?:X)$	X als Gruppe ohne späterer Bezugsmöglichkeit

- Eine Zeichenfolge erfüllt den regulären Ausdruck XY wenn sein Anfang X erfüllt und sein hinterer Teil Y
- Eine Zeichenfolge erfüllt den regulären Ausdruck X/Y wenn er X erfüllt oder Y

Bsp.:

```
Pattern.matches("Hallo (Welt )+", "Hallo Welt Welt ") // → true
Pattern.matches("Hallo (Welt )+", " Hallo Welt Hallo Welt ") // → false
```

- Soll ein Wiederholungszeichen nur auf einen Teilausdruck angewendet werden, so kann man die Gruppenbildung verwenden: z.B. $X(Y)^+$ → `+` bezieht sich nur auf Y