

4 Zeichenfolgen

1. Reguläre Ausdrücke
2. Backus Naur Form

Literatur, Quellen

- Friedl, J.E.F.: Reguläre Ausdrücke - 3. Aufl., Beijing ; Köln [u.a.] : O'Reilly, 2008.
- Online-Test für reguläre Ausdrücke zum einfachen Überprüfen der Übungsaufgaben
 - www.regexr.com
 - regexpal.com

Reguläre Ausdrücke

- Regeln zur Beschreibung von Zeichenketten
- Unterstützung durch zahlreiche Programmiersprachen, Tools
- Bezeichnung:
 - Englisch: regular expression
 - Abgekürzt: Regex
- Einführung
 - Ein regulärer Ausdruck beschreibt eine Menge von Worten durch Regeln
 - Typ-3 in Chomsky-Hierarchie
 - Wird durch endlichen Automat akzeptiert

Reguläre Ausdrücke

- Definition

- a definiert eine Menge, die nur das Zeichen a als Wort enthält
- $F G$ definiert Worte, die jeweils aus einem Wort aus $f \in F$ und $g \in G$ zusammengesetzt sind
- $F|G$ definiert Worte, die aus F oder G stammen
- Häufige Erweiterungen:

Y^+	eine nicht-leere Folge von Y
Y^*	eine evtl. leere Folge von Y
$Y?$	optional Y
$[abc]$	eines der Zeichen a , b oder c
$Y\{m, n\}$	eine Folge von Y mit mind. m und höchstens n Elementen
$Y\{m, \}$	eine Folge von Y mit mind. m Elementen
$Y\{n\}$	eine Folge von Y mit genau n Elementen

Reguläre Ausdrücke

- Beispiele
 - Ausdruck besteht aus beliebig vielen a oder b
 - $(a|b)^+$
 - $(a|b)^+$
 - $[ab]^+$
 - Ausdruck fängt mit a an und endet mit b
 - $a(a|b)^*b$
 - Beispiel: Reguläre Ausdrücke für Zahlen
 - $D = (0|1|2\dots9)$ oder $D=[0-9]$
 - D^+
 - $D^*.D^+$
 - $D+(.D^+)?$

Reguläre Ausdrücke in der Programmierung

- Einige vordefinierte Zeichenklassen
 - Nicht von allen Implementierungen in gleicher Weise unterstützt

.	Beliebiges Zeichen
\.	Punkt
\\	Backslash
\n	Newline
\d	Ziffer, identisch mit $[0-9]$
\D	Keine Ziffer, identisch mit $[\wedge 0-9]$
\s	Leerraumzeichen (z.B. Space, Tabulator)
\S	Kein Leerraumzeichen, identisch mit $[\wedge \backslash s]$
\w	Alphanumerisches Zeichen, identisch mit $[a-zA-Z_0-9]$
\W	Kein alphanumerische Zeichen, identisch mit $[\wedge \backslash w]$
^	Beginn des Textes
\$	Ende des Textes

Reguläre Ausdrücke

- Aufgaben
 - Definieren Sie die Anrede einer Person als regulären Ausdruck
 - Definieren Sie eine Mailadresse als regulären Ausdruck
 - Definieren Sie eine Adresse inkl. Strasse Hausnummer, Postleitzahl, Stadt als regulären Ausdruck
- Verwendung regulärer Ausdrücke
 - C, Java, PHP
 - Unix-Skripte wie egrep, sed
 - Webseiten wie z.B. <http://www.regexe.de>
 - Java siehe nächste Folie

Reguläre Ausdrücke in der Programmierung

- Verwendung regulärer Ausdrücke in Java

- Klasse String

```
public String[] split(String regex)
public boolean matches(String regex)
public String replaceAll(String regex,
                        String replacement)
```

- Klasse Pattern

```
static boolean matches(String regex,
                      CharSequence input)
```

- Beispiel

```
Pattern.matches(
    "[a-zA-Z]{2,20}\\.? [0-9]{1,3} [a-zA-Z]?",
    "Brauneggerstr. 55" )
```


Backus Naur Form (BNF)

- Beschreibung
 - Backus Naur Form (BNF) zur Beschreibung von kontextfreien Grammatiken
 - Ersetzungssysteme
 - Definition der Sprache als Menge von Sätzen, die mit Regeln erzeugt werden können
- Anwendung für
 - Programmiersprachen, wie z.B. Java, C
 - Sprachen als Schnittstelle, z.B. HTML, XML
 - Strukturen von Protokollen

HTML-Beispiel

```
<table>
  <tr>
    <td>Mo</td>
    <td>11-13</td>
    <td>AM</td>
  </tr>
  <tr>
    <td>Fr</td>
    <td>11-13</td>
    <td>AM</td>
  </tr>
</table>
```

Definition BNF

- Kontextfreie Grammatik $G = (T, N, P, S)$ besteht aus
 - T Menge der Terminalsymbole (Terminale)
entspricht der Elemente der Sprache
 - N Menge der Nichtterminalsymbole (Nichtterminale)
wird in Regeln definiert
 - $S \in N$ Startsymbol
 - $P \subseteq N \times V^*$ Menge der Produktionen, $V = T \cup N$
 $(A, x) \in P$ mit $A \in N, x \in V^*$
statt $(A, x) \in P$ schreibt man $A ::= x$
- Unterscheidung Terminalsymbole und Nichtterminalsymbole
 - Verwendung von spitzen Klammern für Nichtterminalsymbole
 - Oder: Anführungsstriche für Terminalsymbole

Definition BNF

- Beispiele

- $\langle \text{DeutscherSatz} \rangle ::= \langle \text{Subjekt} \rangle \langle \text{Prädikat} \rangle \langle \text{Objekt} \rangle$
 $\langle \text{Subjekt} \rangle ::= \langle \text{Artikel} \rangle \langle \text{Substandiv} \rangle$
 $\langle \text{Objekt} \rangle ::= \langle \text{Artikel} \rangle \langle \text{Substantiv} \rangle$
 $\langle \text{Substantiv} \rangle ::= \text{"Mann"} \mid \text{"Buch"}$
 $\langle \text{Artikel} \rangle ::= \text{"der"} \mid \text{"die"} \mid \text{"das"}$
 $\langle \text{Prädikat} \rangle ::= \text{"liest"}$
- $\text{ZifferAußerNull} ::= \text{"1"} \mid \text{"2"} \mid \text{"3"} \mid \text{"4"} \mid \text{"5"} \mid \text{"6"} \mid \text{"7"} \mid \text{"8"} \mid \text{"9"}$
- $\text{Ziffer} ::= \text{"0"} \mid \text{ZifferAußerNull}$
- $\text{Klammerung} ::= \text{"(" Liste "}"$
 $\text{Liste} ::= \text{Klammerung Liste}$
 $\text{Liste} ::=$

Ableitungsregel

- Definition Ableitung

- Sei $G=(T, N, P, S)$ eine kontextfreie Grammatik. Dann kann man mit der Produktion $A::=x \in P$ das Nichtterminal $A \in N$ durch die rechte Seite x ersetzen. Das ist eine **Ableitungsregel**.
- Mehrere Ableitungsregeln nacheinander angewandt heißen **Ableitung**

- Beispiel

- $\text{Klammerung} ::= "(" \text{ Liste } ")"$
 $\text{Liste} ::= \text{Klammerung Liste}$
 $\text{Liste} ::=$

Klammerung
 $\Rightarrow (\text{Liste})$
 $\Rightarrow (\text{Klammerung Liste})$
 $\Rightarrow (\text{Klammerung Klammerung Liste})$
 $\Rightarrow (\text{Klammerung (Liste) Liste})$
 $\Rightarrow ((\text{Liste}) (\text{Liste}) \text{Liste})$
 $\Rightarrow (() (\text{Liste}) \text{Liste})$
 $\Rightarrow (() () \text{Liste})$
 $\Rightarrow (() ())$

Beispiele

- Boolesche Term-Bäume
 - $T = \{ \text{true}, \text{false}, (,), \wedge, \vee, \neg, \text{Variable} \}$
 - $N = \text{BOOL}$
 - $S = \text{BOOL}$
 - $P = \{ \begin{array}{l} \text{BOOL} ::= \text{Variable} \\ \text{BOOL} ::= \text{true} \\ \text{BOOL} ::= \text{false} \\ \text{BOOL} ::= (\text{BOOL} \wedge \text{BOOL}) \\ \text{BOOL} ::= (\text{BOOL} \vee \text{BOOL}) \\ \text{BOOL} ::= \neg \text{BOOL} \end{array} \}$

Erweiterte BNF

- EBNF - Erweiterung der Backus Naur Form
 - Erweiterung zur besseren Lesbarkeit
- Notation
 - Alternative |
 - Gruppierung ()
 - Optionaler Inhalt []
 - Beliebige Wiederholung { }
- BNF ohne Erweiterung
 - Zahl ::= PositiveZahl | "-" Positive Zahl
 - PositiveZahl ::= (Ziffer PositiveZahl) | Ziffer
- BNF mit Erweiterung
 - Zahl ::= ["-"] { Ziffer }

Beispiele

- HTML-Tabellen

- Table ::= <table> Rows </table>
Rows ::= { Row }
Row ::= <tr> Cells </tr>
Cells ::= { Cell }
Cell ::= <td> Text </td>
Cell ::= <td> Table </td>

- Satz der Tabellen-Grammatik

- <table>
 <tr>
 <td>Mo</td>
 <td>11-13</td>
 <td>AM</td>
 </tr>
 <tr>
 <td>Fr</td>
 <td>11-13</td>
 <td>AM</td>
 </tr>
</table>

Erweiterte BNF

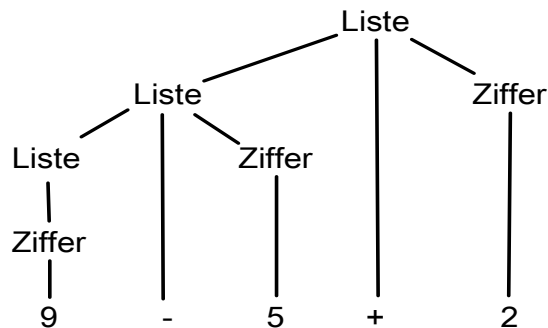
- Beispiel

```
start → expr  
expr → expr + term | term  
term → term * fact | fact  
fact → ID | INT | (expr)
```


Syntaxanalyse und Syntaxbäume

- Parsing: Zeichenkette \rightarrow Struktur
- Beispiel-Produktionen
 - Liste \rightarrow Liste + Ziffer
 - Liste \rightarrow Liste – Ziffer
 - Liste \rightarrow Ziffer
 - Ziffer \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

Parsebaum
(konkreter Syntaxbaum)



Syntaxbaum
(abstrakter Syntaxbaum)

