
Compilerbau

<http://proglang.informatik.uni-freiburg.de/teaching/compilerbau/2004/>

Übungsblatt 2

Abgabe: 10.11.2004

Aufgabe 1 (Linksquotient von M nach L):

Für zwei beliebige Sprachen L, M definiert $L \setminus M = \{w \mid \exists v \in L \text{ so dass } vw \in M\}$ den *Linksquotienten von M nach L*.

Zeige, dass für zwei reguläre Sprachen R_1 und R_2 die Sprache $R_1 \setminus R_2$ ebenso regulär ist.

Aufgabe 2 (Erweiterte reguläre Ausdrücke):

Scannergeneratoren, wie etwa Lex, benutzen als Beschreibungssprache für Tokens in der Regel nicht nur einfache reguläre Ausdrücke, sondern erweitern diese um zusätzliche Operatoren. Die folgende Tabelle zeigt die erweiterte Beschreibungssprache von Lex (wobei c für ein beliebiges Zeichen, r für einen regulären Ausdruck und s für eine Zeichenkette steht):

Ausdruck	Bedeutung	Beispiel
c	bel. (Nichtsonder-)Zeichen c	a
$\setminus c$	Zeichen c	$\setminus *$
" s "	Zeichenkette s	"**"
$.$	bel. Zeichen außer <i>newline</i>	$a.*b$
\wedge	Zeilenanfang	$\wedge abc$
$\$$	Zeilenende	$abc\$$
$[s]$	bel. Zeichen aus der Klasse s	$[abc]$
$[\wedge s]$	bel. Zeichen nicht aus der Klasse s	$[\wedge abc]$
r^*	null oder mehr rs	a^*
r^+	ein oder mehr rs	a^+
$r^?$	null oder ein r	$a^?$
$r\{m, n\}$	m bis n Wiederholungen von r	$a\{1, 5\}$
$r_1 r_2$	erst r_1 dann r_2	ab
$r_1 \mid r_2$	r_1 oder r_2	$a \mid b$
(r)	r	$(a \mid b)$
r_1 / r_2	r_1 falls gefolgt von r_2	$abc / 123$

- (i) Die Spezialbedeutung der Operatorsymbole (\setminus " $.$ \wedge $[$] $*$ $+$ $?$ $\{$ $\}$ $|$ $/$) muss ausgeschaltet werden, falls diese selber übereinstimmen sollen. Es gibt zwei Möglichkeiten, dies zu bewerkstelligen. Falls kein "-"Zeichen in einer Zeichenkette s vorkommt, kann diese mit " s " exakt auf Übereinstimmung überprüft werden (etwa in "**"). Alternativ können die Zeichen separat mit Hilfe von \setminus angegeben werden (etwa in $\setminus *$).

Gib einen regulären Ausdruck für Lex an, der die Zeichenkette " \setminus " beschreibt.

- (ii) Eine Zeichenklasse, die mit dem Symbol \wedge beginnt, beschreibt eine *Komplementär-Zeichenklasse*. Mit einer Komplementär-Zeichenklasse erkennt man alle Zeichen, die

nicht in einer bestimmten Zeichenklasse enthalten sind ($[\hat{a}]$ etwa beschreibt alle Zeichen außer a).

Zeige, dass zu jedem regulären Ausdruck mit Komplementär-Zeichenklassen ein äquivalenter regulärer Ausdruck ohne Komplement-Zeichenklassen existiert.

- (iii) Der reguläre Ausdruck $r\{m,n\}$ erkennt m bis n Wiederholungen des Musters r . Beispielsweise erkennt $a\{1,5\}$ Zeichenketten mit ein bis fünf a s.

Zeige, dass zu jedem regulären Ausdruck mit Wiederholungsoperator ein äquivalenter regulärer Ausdruck ohne Wiederholungsoperator existiert.

Aufgabe 3 (Ableitungen regulärer Ausdrücke):

In der Vorlesung wurde die Ableitung eines regulären Ausdrucks r bezüglich eines Symbols a , $D(r, a)$, zusammen mit einer Hilfsfunktion E definiert.

Beweise die folgenden Hilfslemmas über Ableitungen regulärer Ausdrücke:

- (i) $L(E(r)) = L(r) \cap \{\epsilon\}$
- (ii) Für alle $a \in \Sigma$ gilt: $w \in L(D(r, a))$ iff $aw \in L(r)$.

Beweise mit Hilfe der beiden Hilfslemmas das Repräsentationslemma für reguläre Sprachen:

- (iii) $L(r) = L(E(r)) \cup \bigcup_{a \in \Sigma} aL(D(r, a))$

Aufgabe 4 (SAX-Parser):

SAX¹ ist eine einheitliche Schnittstelle für XML-Parser, die über eine Event-Schnittstelle mit der Außenwelt kommunizieren.

Implementiere in Java² einen (vereinfachten) SAX-konformen Parser für XML-Dokumenten mit Hilfe des Scannergenerators JLex³ oder JFlex⁴.

Gehe etwa nach folgendem “Kochrezept” vor:

- (i) Untersuche die SAX-Parser-API-Klasse `DefaultHandler`⁵ und finde heraus, welche Art von Tokens Dein Scanner/Parser erkennen muss. (Dein Parser sollte später zumindest sinnvolle Events durch Methodenaufrufe von `characters`, `endDocument`, `endElement`, `fatalError`, `ignorableWhitespace`, `startDocument` und `startElement` generieren können, wobei Attribute eines Elements ignoriert oder als gemeinsamer String behandelt werden können.)
- (ii) Untersuche die Spezifikation von XML⁶ und definiere für jede Klasse von Tokens einen passenden regulären Ausdruck.
- (iii) Implementiere Deinen SAX-Parser, indem Du eine passende JLex-Spezifikation schreibst, wobei die zu den Regeln assoziierten *Actions* passende SAX-Events erzeugen.

¹<http://sax.sourceforge.net/>

²Die Aufgabe kann auch in jeder anderen Programmiersprache bearbeitet werden, für die es einen Scannergenerator gibt.

³<http://www.cs.princeton.edu/~appel/modern/java/JLex/>

⁴<http://jflex.de/>

⁵<http://sax.sourceforge.net/apidoc/org/xml/sax/helpers/DefaultHandler.html>

⁶<http://www.w3.org/TR/REC-xml>