

Anwendungen kontextfreier Grammatiken

Vortrag im Rahmen des Proseminars Theoretische Informatik WS 10/11

von Jan Lottermoser

Ursprünglich waren kontextfreie Grammatiken als Mittel zur Beschreibung natürlicher Sprachen gedacht. Diese Erwartung hat sich jedoch nicht erfüllt. Erst als rekursive Konzepte in der Informatik stark zu nahmen, fanden kontextfreie Grammatiken erste Anwendungen.

Im Folgenden werden eine alte und eine neuere Anwendung kontextfreier Grammatiken vorgestellt:

1. Kontextfreie Grammatiken zur Beschreibung von Programmiersprachen (Parser)
2. Extensible Markup Language (XML) und Dokumenttypdefinitionen (DTD)

1. Beschreibung von Programmiersprachen (Parser)

Es gibt wichtige Aspekte von Programmiersprachen, die nicht durch reguläre Ausdrücke beschrieben werden können. Es folgen zwei Beispiele:

Beispiel a) $L = \text{Dyck-Sprache (Sprache der korrekten Klammerausdrücke)}$

$G = (\{B\}, \{(\,)\}, P, B)$ wobei P besteht aus:

$B \rightarrow BB \mid (B) \mid \varepsilon$

Behauptung: $L(G) = L$

Beweis: (1) $L(G) \subseteq L$

Offensichtlich erzeugt G nur Wörter aus L , da zwei Wörter aus L konkateniert oder ein Wort aus L umklammert wieder ein Wort aus L ergeben und ε auch in L liegt.

(2) $L \subseteq L(G)$

Induktion über Länge n eines Wortes w aus L . (n ist Anzahl der Klammerpaare)

IA: $n=0$: $w = \varepsilon$: $B \xrightarrow{G} \varepsilon$

$n=1$: $w = ()$: $B \xrightarrow{G} (B) \xrightarrow{G} ()$

IS: $n \geq 2$: $w = (u)$, $|u| = n$, $u \in \{(\,)\}^*$

1. Fall: $u \in L$

Nach Ind. Ann. gilt $B \xrightarrow{G} u$ Somit ist Ableitung für w : $B \xrightarrow{G} (B) \xrightarrow{G}^* (u) = w$

2. Fall: $u \notin L \Rightarrow u = x)y(z, \Rightarrow x, y, z \in L$

Nach Ind. Ann. gilt: $B \xrightarrow{G}^* x$ und $B \xrightarrow{G}^* y$ und $B \xrightarrow{G}^* z$ Somit ist die Ableitung für w:

$$B \xrightarrow{G} BB \xrightarrow{G} BBB \xrightarrow{G} (B)BB \xrightarrow{G} (B)B(B) \xrightarrow{G}^* (x)y(z) = (x)y(z) = (u) = w \quad \square$$

Beweis, dass L nicht regulär ist: (Pumping Lemma)

Wäre L regulär, so gäbe es ein $k \in \mathbb{N}$ für das sich jedes w aus L mit $|w| \geq k$ zerlegen lässt in

$w = uvw$ mit $|v| \geq 1$ und $|uv| \leq k$, sodass gilt: $uv^i w \in L$ für $i = 0, 1, 2, \dots$

Sei $w = ({}^k)^k$. Da $|uv| \leq k$ gilt, besteht v nur aus einer oder mehr öffnender Klammern. uw enthält somit mehr schließende als öffnende Klammern. \square

Es gibt jedoch auch Strukturen in Programmiersprachen, die nicht immer ausgewogene Zeichenreihen enthalten. Ein Beispiel dafür ist die if-else-Struktur in C. Eine if-Anweisung kann mit oder ohne else-Klausel vorkommen.

Beispiel b) L = Sprache der korrekten if-else-Klauseln über dem Alphabet $\{i, e\}$

$G = (\{S\}, \{i, e\}, P, \{S\})$ wobei P aus den folgenden Produktionen besteht:

$$S \rightarrow \varepsilon \mid SS \mid iS \mid iSe$$

Behauptung: $L(G) = L$

Beweis: (1) $L(G) \subseteq L$

Offensichtlich erzeugt G nur Wörter aus L: ε ist korrekte if-else-Klausel. Zwei konkatenierte korrekte if-else-Klauseln sind wieder korrekt. Vor eine korrekte if-else-Klausel kann man ein if stellen und sie bleibt korrekt und man kann eine korrekte if-else-Klausel mit i und e umranden, da das e sich entweder auf das erste i bezieht oder auf eines danach.

(2) $L \subseteq L(G)$

Induktion über die Länge n eines Wortes w aus L.

$$\text{IA: } n=0: w=\varepsilon: S \xrightarrow{G} \varepsilon = w \quad n=1: w=i: S \xrightarrow{G} iS \xrightarrow{G} i = w$$

$$\text{IS: } n \geq 2: w = ixex$$

1. Fall: e bezieht sich auf i. Dann muss gelten: $x \in L$ und nach Ind. Ann. $S \xrightarrow{G}^* x$

$$\text{Also: } S \xrightarrow{G} iSe \xrightarrow{G}^* ixex = w$$

2. Fall: e bezieht sich auf ein i innerhalb von x. Somit hat x die Form $x = uiv$. Da das letzte e in w sich auf das i zwischen u und v bezieht, ist v aus L und nach Ind. Ann.

$$S \xrightarrow{G}^* v \quad \text{Da } w = uiuve \text{ aus L und } ive \text{ aus L muss auch } iu \text{ aus L sein. Somit gilt:}$$

$$S \xrightarrow{G} SS \xrightarrow{G} iSS \xrightarrow{G}^* iuS \xrightarrow{G} iuiSe \xrightarrow{G}^* iuive = w \quad \square$$

Der **Beweis**, dass L nicht regulär ist, ist analog zu dem aus Beispiel a) für das Wort $w = i^k e^k$. \square

1.1 Der YACC-Parsergenerator

Ein Parser ist eine Funktion, die die Struktur eines Quelltextprogramms erkennt und sie durch einen Ableitungsbaum (Parsebaum) repräsentiert. Der YACC Befehl ist auf allen UNIX Systemen vorhanden und erzeugt aus gegebener kfg einen Parser.

2. Markupsprachen und XML

Markup Sprachen oder Auszeichnungssprachen sind Sprachen, die bestimmte Markierungen, genannt Tags, enthalten. Wörter dieser Sprachen werden Dokumente genannt und die Tags sagen etwas über die Semantik der verschiedenen Zeichenreihen innerhalb der Dokumente aus. Die bekannteste Auszeichnungssprache ist HTML und wie jede Programmiersprache kann sie durch eine kontextfreie Grammatik beschrieben werden. Einen Teil einer HTML Grammatik zeigt folgende Abbildung:

1. $Char \rightarrow a \mid A \mid \dots$
2. $Text \rightarrow \varepsilon \mid Char Text$
3. $Doc \rightarrow \varepsilon \mid Element Doc$
4. $Element \rightarrow Text \mid \langle EM \rangle Doc \langle /EM \rangle \mid \langle P \rangle Doc \mid \langle OL \rangle List \langle /OL \rangle$
5. $ListItem \rightarrow \langle LI \rangle Doc$
6. $List \rightarrow \varepsilon \mid ListItem List$

2.1 Extensible Markup Language (XML)

XML ist eine erweiterbare Auszeichnungssprache. Dies bedeutet, dass man mithilfe von sogenannten Dokumenttypdefinitionen (DTD) die Sprache um Tags erweitern kann und dabei die zulässigen Strukturen innerhalb der Tags definiert.

Eine DTD ist im Grunde eine kontextfreie Grammatik mit spezieller Syntax. Eine DTD hat folgende Form:

```
<!DOCTYPE Name-der-DTD [ Liste der Elementdefinitionen ]>
```

Eine Elementdefinition sieht wiederum so aus:

```
<!ELEMENT Elementname (Beschreibung des Elements)>
```

Hier lässt sich schon erahnen, dass die Sprache zur Beschreibung von DTDs wieder durch eine kfg beschrieben werden kann.

Die Beschreibung eines Elements ist im Grunde ein regulärer Ausdruck. Grundlage dieser reg. Ausdrücke sind wieder Elemente und der spezielle Term #PCDATA, der für Text steht, der keine Tags beinhaltet. Die Operatoren sind Folgende:

1. | steht für die Vereinigung
2. Ein Komma ist die Konkatenation
3. * ist der gewöhnliche Sternoperator, + steht für ein oder mehr Vorkommen und ? steht für null oder ein Vorkommen

Die untere Abbildung zeigt ein Beispiel einer DTD für die Beschreibung von PCs.

```

<!DOCTYPE PcSpecs [
<!ELEMENT PCS          (PC*)>
<!ELEMENT PC          (MODEL, PRICE, PROCESSOR, RAM, DISK+)>
<!ELEMENT MODEL       (#PCDATA)>
<!ELEMENT PRICE       (#PCDATA)>
<!ELEMENT PROCESSOR   (MANF, MODEL, SPEED)>
<!ELEMENT MANF        (#PCDATA)>
<!ELEMENT MODEL       (#PCDATA)>
<!ELEMENT SPEED       (#PCDATA)>
<!ELEMENT RAM         (#PCDATA)>
<!ELEMENT DISK        (HARDDISK | CD | DVD)>
<!ELEMENT HARDDISK    (MANF, MODEL, SIZE)>
<!ELEMENT SIZE        (#PCDATA)>
<!ELEMENT CD          (SPEED)>
<!ELEMENT DVD         (SPEED)>
]>

```

Diese Notation für die Definition von DTDs entspricht nicht ganz der einer kfG, es gibt jedoch Regeln zur Umwandlung einer solchen kfG in eine gewöhnliche. Wir zeigen dies durch Induktion über die Größe des Ausdrucks um Produktionsrumpf.

Induktionsanfang:

Ist der Rumpf eine Verkettung von Elementen, so ist die Produktion schon in der gewöhnlichen kfG Notation:

`<!ELEMENT HARDDISK (MANF, MODEL, SIZE)>` wird zu
 Harddisk → Manf Model Size

Induktionsschritt:

Andernfalls sind abhängig vom letzten Operator im Produktionsrumpf fünf Fälle zu unterscheiden:

1. Die Produktion hat die Form $A \rightarrow E_1 E_2$, wobei E_1 und E_2 zulässige Ausdrücke sind. Wir führen zwei neue Variablen B und C ein, die sonst nirgends in der Grammatik vorkommen, und ersetzen die Produktion durch $A \rightarrow BC$, $B \rightarrow E_1$, $C \rightarrow E_2$. Die erste Produktion ist in gewöhnlicher Notation und da die Rümpfe der anderen beiden Produktionen kürzer sind als der Rumpf der ursprünglichen Produktion, kann man sie nach Induktionsannahme auch in die gewöhnliche Notation umwandeln. Dies gilt auch für alle folgenden Fälle.
2. Die Produktion hat die Form $A \rightarrow E_1 | E_2$. Sie wird ersetzt durch $A \rightarrow E_1$, $A \rightarrow E_2$.
3. Die Produktion hat die Form $A \rightarrow (E_1)^*$ und wird ersetzt durch $A \rightarrow BA$, $A \rightarrow \varepsilon$, $B \rightarrow E_1$.
4. Die Produktion hat die Form $A \rightarrow (E_1)^+$ und wird ersetzt durch $A \rightarrow BA$, $A \rightarrow B$, $B \rightarrow E_1$.
5. Die Produktion hat die Form $A \rightarrow (E_1)^?$ und wird ersetzt durch $A \rightarrow \varepsilon$, $A \rightarrow E_1$.