

Grundlagen Theoretischer Informatik 2

WiSe 2011/12 in Trier

Henning Fernau

Universität Trier

fernau@uni-trier.de

Grundlagen Theoretischer Informatik 2 Gesamtübersicht

- Organisatorisches; Einführung
- Ersetzungsverfahren: [Grammatiken und Automaten](#)
- Rekursionstheorie-Einführung
- Komplexitätstheorie-Einführung

Organisatorisches

Vorlesungen FR 12.15-13.45 im F59

Vorschlag: 12.25-13.55 ab 2. Semesterwoche

Übungsbetrieb in Form von einer Übungsgruppe

BEGINN: als Sprechstunde in der ersten Semesterwoche

FR 14.00-15.30, F 59

Dozentensprechstunde DO 13-14 in meinem Büro H 410 (4. Stock)

Mitarbeitersprechstunde (Stefan Gulan) DO 13-14 H 413

Tutorensprechstunde (Martin Puppe) DO 13-14 H 407

Chomsky-Normalform für kontextfreie Grammatiken

Eine kontextfreie Grammatik $G = (N, T, P, S)$ ist in Chomsky-Normalform, falls alle ihre Regeln die folgende Form haben:

$$A \rightarrow BC \text{ oder } A \rightarrow a$$

für $A, B, C \in N$ und $a \in T$.

Kurz: $P \subset (N \times N^2) \cup (N \times T)$

Satz: Zu jeder kontextfreien Grammatik G gibt es eine Grammatik G' in Chomsky-Normalform mit $L(G') = L(G) \setminus \{\lambda\}$.

Hinweis: Mit dem Verfahren aus der vorigen Vorlesung können wir zunächst eine zu G äquivalente λ -freie kfG \tilde{G} konstruieren mit $L(\tilde{G}) = L(G) \setminus \{\lambda\}$.

1. Lange rechte Seiten mit Terminalzeichen

Wir nennen eine Satzform *gemischt*, wenn sie sowohl Terminal- als auch Nicht-terminalzeichen enthält.

Lemma: Zu jeder kfG G lässt sich eine äquivalente kfG G' konstruieren, bei welcher rechte Regelseiten mit Terminalzeichen die Länge höchstens Eins haben. Mit G ist auch G' λ -frei.

Damit enthält G' auch keine gemischten rechten Regelseiten.

1. Führe für jedes Terminalzeichen a ein spezielles Nichtterminal X_a ein.
2. Führe neue Regeln $X_a \rightarrow a$ ein.
3. Ersetze a in rechter Regelseite w durch X_a , falls $w \neq a$.

\implies Nur noch Regeln der Formen

$$X \rightarrow a \text{ mit } a \in T \text{ oder } X \rightarrow w \text{ mit } w \in N^*$$

2. Lange rechte Seiten mit Nichtterminalzeichen

Entfernung aller Regeln $X \rightarrow B_1 \dots B_k$ mit $k \geq 3$

Dazu neue Nichtterminale C_1, C_2, \dots, C_{k-2} und neue Regeln:

$$\begin{aligned} X &\rightarrow B_1 C_1 \\ C_1 &\rightarrow B_2 C_2 \\ &\dots \\ C_{k-2} &\rightarrow B_{k-1} B_k \end{aligned}$$

\implies Nur noch Regeln der Formen

$$X \rightarrow a \text{ mit } a \in T \text{ oder } X \rightarrow B_1 B_2 \text{ mit } B_1, B_2 \in N$$

oder

$$X \rightarrow B \text{ mit } B \in N$$

Die letztgenannten Kettenregeln zu entfernen können wir bereits (letzte VL)!

Der Weg zur Chomsky-Normalform Zusammenfassung:

0. Eliminiere λ -Regeln.
1. Handle zu lange rechte Seiten mit Terminalzeichen.
2. Handle zu lange rechte Seiten mit Nichtterminalzeichen.
3. Entferne Kettenregeln (durch Abkürzungen).

Beachte:

Schritt i zerstört nicht die durch Schritt $j < i$ gewonnenen Eigenschaften!

Ein Beispiel: Gegeben sei $G = (\{S, T\}, \{a, b, c\}, P, S)$ mit

$$P = \{S \rightarrow aT \mid TbT \mid T, T \rightarrow S \mid c\}$$

1) neue Nichtterminale X_a, X_b , neue Regeln $X_a \rightarrow a, X_b \rightarrow b$

Ersetzung: $S \rightarrow aT$ durch $S \rightarrow X_aT$; $S \rightarrow TbT$ durch $S \rightarrow TX_bT$

neue Regelmengemenge:

$$\{S \rightarrow X_aT \mid TX_bT \mid T, T \rightarrow S \mid c, X_a \rightarrow a, X_b \rightarrow b\}$$

Ein Beispiel (Forts.):

2) Neues Nichtterminal C

Ersetzung: $S \rightarrow TX_bT$ durch $S \rightarrow TC$ und $C \rightarrow X_bT$

neue Regelmenge

$$\{S \rightarrow X_aT | TC | T, T \rightarrow S | c, X_a \rightarrow a, X_b \rightarrow b, C \rightarrow X_bT\}$$

Ein Beispiel (Forts.):

Schritt 3: Eliminiere die beiden Kettenregeln $S \rightarrow T$ und $T \rightarrow S$.

Zuerst $S \rightarrow T$: füge die beiden Regeln $S \rightarrow S$ und $S \rightarrow c$ hinzu.

Bei $T \rightarrow S$ kommen die Regeln $T \rightarrow X_a T$, $T \rightarrow T$ und $T \rightarrow TC$ hinzu.

Wir haben jetzt die Regelmenge

$$\{S \rightarrow X_a T | TC | T | S | c, T \rightarrow S | c | X_a T | TC | T, X_a \rightarrow a, X_b \rightarrow b, C \rightarrow X_b T\}$$

Neue Regeln entstehen nicht mehr,

\rightsquigarrow Entferne die vier Regeln $S \rightarrow S | T$ und $T \rightarrow S | T$.

\rightsquigarrow Zu G äquivalente Grammatik $G' = (\{S, T, X_a, X_b, C\}, \{a, b, c\}, P', S)$ in Chomsky-Normalform, wobei:

$$P' = \{S \rightarrow X_a T | TC | c, T \rightarrow c | X_a T | TC, X_a \rightarrow a, X_b \rightarrow b, C \rightarrow X_b T\}$$

Satz: Es gibt einen Algorithmus, der zu jeder vorgelegten kfG $G = (N, T, P, S)$ und jedem $w \in T^*$ entscheidet, ob $w \in L(G)$ gilt.

Beweis: Wird $w = \lambda$ angefragt, so berechnen wir N_1 wie im Beweis I in der vorigen VL. Es gilt:
 $\lambda \in L(G)$ gdw. $S \in N_1$.

Wird $w \neq \lambda$ angefragt, so überführen wir G in eine Chomky-Normalform kfG G' mit der Eigenschaft $w \in L(G)$ gdw. $w \in L(G')$.

Da G' keine λ -Regeln enthält, kann man " $w \in L(G)$?" durch Berechnen aller aus G' ableitbaren Satzformen der Länge höchstens $\ell(w)$ beantworten.

Cocke, Younger und Kasami haben gezeigt, wie man die Komplexität des beschriebenen Verfahrens durch **dynamisches Programmieren** erheblich vermindern kann.

Das Verfahren von Cocke, Younger und Kasami (*CYK-Algorithmus*)

Satz: Ist eine kfG G in Chomsky-Normalform fixiert, so lässt sich die Frage “ $w \in L(G)$?” in einer Zeit beantworten, die sich durch ein kubisches Polynom in $\ell(w)$ abschätzen lässt.

Beweis: Betrachte $w = a_1 \dots a_n$.

Lege 2–dimensionale dreiecksförmige Tabelle T an mit Nonterminalmengen als Einträgen, so dass $T[i, j] = \{A \in N \mid A \xrightarrow{*} a_i \dots a_j\}$. T kann man wie folgt “von oben nach unten” berechnen:

$$T[i, i] = \{A \in N \mid A \rightarrow a_i \in P\}$$

$$T[i, j] = \{A \in N \mid A \rightarrow BC \in P \wedge \exists i \leq k < j : B \in T[i, k], C \in T[k + 1, j]\}$$

Die CYK-Tabelle

a_1	a_2	\dots	a_n
$T[1, 1]$	$T[2, 2]$	\dots	$T[n, n]$
$T[1, 2]$	$T[2, 3]$	\dots	
\vdots	\vdots		
$T[1, n-1]$	$T[2, n-1]$		
$T[1, n]$			

$a_1 \dots a_n \in L(G)$ gdw. das Startsymbol erscheint in $T[1, n]$.

Der CYK-Algorithmus

FOR $j = 1, \dots, n$

$T[i, i] = \{A \in N \mid A \rightarrow \alpha_i \in P\}$

FOR $j = 1, \dots, n - 1$

FOR $i = 1, \dots, n - j$

$T[i, j + i] = \emptyset$

FOR $k = i, \dots, j + i - 1$

IF $\exists A \rightarrow BC \in P : B \in T[i, k], C \in T[k + 1, j + i]$

THEN $T[i, j + i] = T[i, j + i] \cup \{A\}$

Ein Beispiel

$S \rightarrow AX, X \rightarrow SB, S \rightarrow c, A \rightarrow a, B \rightarrow b$ erzeugt $L = \{a^n cb^n \mid n \geq 0\}$.

a	a	c	b	b
A	A	S	B	B
	S	X		
	X			
S				

Ein **binärer Wurzelbaum** ist gegeben durch ein Tripel $B = (V, \phi, r)$ mit ausgezeichneter *Wurzel* $r \in V$ und einer *Vater-Abbildung* $\phi : V \setminus \{r\} \rightarrow V$ mit der Eigenschaft

$$\forall v \in V : \underbrace{\#\{u \in V \mid \phi(u) = v\}}_{\kappa(v) :=} \leq 2$$

$\kappa(v)$ liefert also die *Kinder* von v .

Knoten v mit $\kappa(v) = \emptyset$ heißen *Blätter*.

Lemma: Der Ableitungsbaum eines jeden von einer kontextfreien Grammatik in Chomsky-Normalform akzeptierten Wortes ist als binärer Wurzelbaum auffassbar.

Die **Höhe** eines binären Wurzelbaumes $B = (V, \phi, r)$ ist gegeben durch

$$h(B) = \max_{v \in V \setminus \{r\}} \{k \in \mathbb{N} \mid \phi^k(v) = r\}.$$

Lemma: Hat ein binärer Wurzelbaum B mehr als 2^h Blätter, so gilt $h(B) > h$.

Beweis: Wir zeigen die **Kontraposition** per Induktion:

Gilt $h(B) \leq h$, so hat $B = (V, \phi, r)$ höchstens 2^h viele Blätter.

$h = 0$ ist trivial.

Es gelte $h > 0$. Daher gilt $\kappa(r) \neq \emptyset$.

Betrachte $r' \in \kappa(r)$. $V' := \{v \in V \mid \exists k \in \mathbb{N} : \phi^k(v) = r'\}$.

ϕ' sei die Einschränkung von ϕ auf $V' \setminus \{r'\}$.

$B' = (V', \phi', r')$ ist ein binärer Wurzelbaum der Höhe höchstens $h - 1$.

Auf B' ist die Induktionshypothese anwendbar: B' hat höchstens 2^{h-1} viele Blätter.

\rightsquigarrow B hat höchstens $\#\kappa(r) \cdot 2^{h-1} \leq 2^h$ viele Blätter.

Lemma: Ist $G = (N, T, P, S)$ eine kfG in Chomsky-Normalform und ist $w \in L(G)$ mit $\ell(w) > 2^{\#N}$, so gilt für jeden Ableitungsbaum von w bzgl. G , dass es einen Weg von S zu einem Blatt gibt, auf dem mehr als $\#N$ viele Nichtterminalzeichen ersetzt werden.

Beweis: Nach dem vorigen Lemma hat der unterliegende binäre Wurzelbaum die Höhe größer als $\#N$. Es gibt also einen Weg in besagtem Ableitungsbaum von der mit S beschrifteten Wurzel zu einem Blatt, dem mehr als $\#N$ Regelanwendungen entspricht \leadsto Behauptung.

Folg.: Auf besagtem Weg von der Wurzel zum Blatt im Ableitungsbaum von w finden wir also nach dem Schubfachprinzip zwei Regelanwendungen $A \rightarrow v$ und $A \rightarrow u$ mit gleicher linker Seite.

Ein Pumping-Lemma für Typ-2

Satz: Zu jeder kfS L gibt es eine Konstante $n > 0$, sodass jedes Wort $w \in L$ mit $\ell(w) \geq n$ als Konkatination $w = uvxyz$ dargestellt werden kann mit geeigneten u, v, x, y, z mit folgenden Eigenschaften:

1. $\ell(v) > 0$ oder $\ell(y) > 0$;

2. $\ell(vxy) \leq n$;

3. $\forall i \geq 0 : uv^i xy^i z \in L$.

Beweis: Mit L ist auch $L \setminus \{\lambda\}$ vom Typ 2. Sei also o.E. L λ -frei.

L werde durch kfG $G = (N, T, P, S)$ in Chomsky-Normalform beschrieben.

Es sei $n = 2^{\#N}$.

Gibt es kein $w \in L$ mit $\ell(w) \geq n$, so stimmt die Aussage leer.

Sonst wähle ein $w \in L$ mit $\ell(w) \geq n$.

Nach obiger Folgerung gibt es zwei Regelanwendungen $A \rightarrow u'$ und $A \rightarrow v'$ auf einem Weg in einem Ableitungsbaum von w .

\rightsquigarrow Es gibt Linksableitung

$$S \xRightarrow{*} uA\zeta \xRightarrow{*} uvA\eta\zeta \xRightarrow{*} uvx\eta\zeta \xRightarrow{*} uvxy\zeta \xRightarrow{*} uvxyz.$$

Wegen Chomsky-NF (keine Kettenregeln oder λ -Regeln) gilt $\ell(v) > 0$ oder $\ell(y) > 0$.

$\ell(vxy) \leq n$ ergibt sich aus dem Beweis der Folgerung sowie wegen Chomsky-NF.

Daher gibt es auch Ableitungen

$$S \xRightarrow{*} uA\zeta \xRightarrow{*} ux\zeta \xRightarrow{*} uxz$$

und allgemeiner

$$S \rightarrow uA\zeta \xRightarrow{*} uvA\eta\zeta \xRightarrow{*} uvvA\eta\eta\zeta \xRightarrow{*} uv^iA\eta^i\zeta \xRightarrow{*} uv^ixy^iz.$$

Ein Beispiel zur Anwendung des Pumping-Lemmas:

$$L = \{a^k b^k c^k \mid k \in \mathbb{N}\} \notin \mathcal{L}_2$$

Wäre L kontext-frei, so gäbe es Pumping-Konstante n .

Wähle $w = a^n b^n c^n$ als genügend langes Wort.

Diskutiere Zerlegungen $w = uvxyz$.

$\ell(vxy) \leq n \rightsquigarrow vxy$ enthält höchstens a 's gefolgt von b 's oder b 's gefolgt von c 's.

Beide Fälle sind analog; diskutiere also den ersten.

“Nullpumpen” liefert ein Wort, in dem alle a - und b -Vorkommen zusammen weniger als die doppelte Anzahl von c 's ausmachen: \nexists Wortstruktur.

Eine Anwendung von KfG:

Beschreibung der Syntax (grammatisches Gerüst) natürlicher Sprachen

Satz \rightarrow NP VP

NP \rightarrow Artikel Nomen

VP \rightarrow Verb

Artikel \rightarrow "der"

Artikel \rightarrow "die"

Nomen \rightarrow "Hund"

Nomen \rightarrow "Hunde"

Nomen \rightarrow "Katze"

Verb \rightarrow "beißen"

Satz \Rightarrow NP VP \Rightarrow^2 Artikel Nomen Verb \Rightarrow^3 "die" "Hunde" "beißen"

Satz \Rightarrow NP VP \Rightarrow^2 Artikel Nomen Verb \Rightarrow^3 "die" "Katze" "beißen"

Satz \Rightarrow NP VP \Rightarrow^2 Artikel Nomen Verb \Rightarrow^3 "der" "Katze" "beißen"

Eine Anwendung von KfG: Beschreibung arithmetischer Ausdrücke

$\Sigma = \{v, -, +, *, /, (,)\}$, $N = \{E\}$.

Die Regeln seien die folgenden:

$E \rightarrow (E)$

$E \rightarrow -(E)$

$E \rightarrow (E + E)$

$E \rightarrow (E - E)$

$E \rightarrow (E * E)$

$E \rightarrow (E/E)$

$E \rightarrow v$

Beispielableitung:

$E \Rightarrow -(E)$

$\Rightarrow -(E + E)$

$\Rightarrow -((E * E) + E)$

$\Rightarrow -((-E) * E + E)$

$\Rightarrow -((-E) * E + (E - E))$

$\stackrel{*}{\Rightarrow} -((-v) * v + (v - v))$

Hinweis: Der *Scanner*, ein endlicher Automat, produziert idealerweise als Ausgabe die Eingabe für den *Parser* zwecks *syntaktischer Analyse* eines Programmtextes.

“Nebensächlichkeiten” wie Zahlen oder Zahlenvariablen werden in einen Statthalter v bersetzt.

Ein Wort hat viele Ableitungen

Linksableitung:

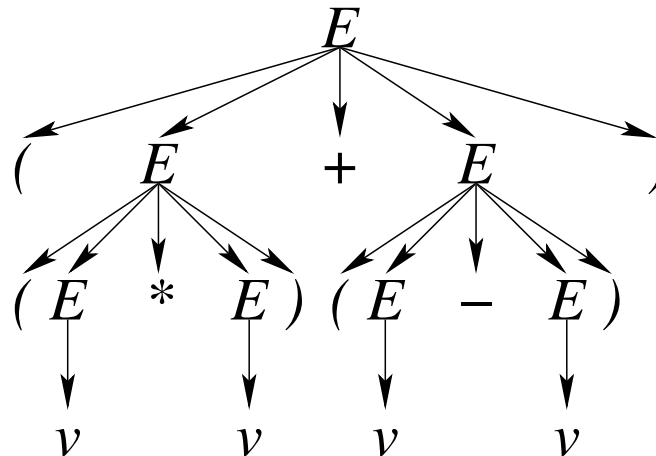
$$\begin{aligned} E &\Rightarrow -(E) \\ &\Rightarrow -(E + E) \\ &\Rightarrow -((E * E) + E) \\ &\Rightarrow -((-E) * E) + E) \\ &\Rightarrow -((-v) * E) + E) \\ &\Rightarrow -((-v) * v) + E) \\ &\Rightarrow -((-v) * v) + (E - E)) \\ &\Rightarrow -((-v) * v) + (v - E)) \\ &\Rightarrow -((-v) * v) + (v - v)) \end{aligned}$$

Rechtsableitung:

$$\begin{aligned} E &\Rightarrow -(E) \\ &\Rightarrow -(E + E) \\ &\Rightarrow -(E + (E - E)) \\ &\Rightarrow -(E + (E - v)) \\ &\Rightarrow -(E + (v - v)) \\ &\Rightarrow -((E * E) + (v - v)) \\ &\Rightarrow -((E * v) + (v - v)) \\ &\Rightarrow -((-E) * v) + (v - v)) \\ &\Rightarrow -((-v) * v) + (v - v)) \end{aligned}$$

Ein Syntaxbaum (oder *Ableitungsbaum*) für

$$E \Rightarrow (E + E) \Rightarrow ((E * E) + E) \Rightarrow ((E * E) + (E - E)) \xRightarrow{*} ((v * v) + (v - v))$$



Linksableitung: Tiefensuche mit Linksabstieg durch Syntaxbaum

Rechtsableitung: Tiefensuche mit Rechtsabstieg durch Syntaxbaum

Operatorbaum: ein kompakter Syntaxbaum für das ursprüngliche Beispiel

