

Einführung in die Computerlinguistik

Kontextfreie Grammatiken

Hinrich Schütze

Center for Information and Language Processing

2019-01-14

Die Grundfassung dieses Foliensatzes wurde von Prof. Dr. Stefan Evert erstellt. Fehler und Mängel sind ausschließlich meine Verantwortung.

- 1 Kontextfreie Grammatiken
- 2 Top-down parsing
- 3 CYK

1 Kontextfreie Grammatiken

2 Top-down parsing

3 CYK

- Mächtiger als reguläre Sprachen / Automaten
- Rekursive Strukturen
- Konstituentengrammatik

Definition kontextfreie Grammatik (CFG)

Eine kontextfreie Grammatik G über dem Alphabet Σ ist ein Quadrupel $G = (V, \Sigma, P, S)$. Die Elemente von V heißen **Variablen** oder **Nichtterminalsymbole**, entsprechend werden die Zeichen aus Σ auch als **Terminalsymbole** bezeichnet. Wir nehmen stets $V \cap \Sigma = \emptyset$ an. üblicherweise verwenden wir für Terminalsymbole Kleinbuchstaben $a, b, c, \dots \in \Sigma$ und für Variablen Großbuchstaben $A, B, C, \dots \in V$. Zur Unterscheidung von Wörtern $u, v, w, \dots \in \Sigma^*$ bezeichnen wir Zeichenketten, die sowohl Variablen als auch Terminalsymbole enthalten können, als **Terme** und verwenden dafür griechische Kleinbuchstaben $\alpha, \beta, \gamma, \dots \in (V \cup \Sigma)^*$. $S \in V$ ist eine spezielle Variable, die **Startsymbol** genannt wird. $P \subseteq V \times (V \cup \Sigma)^*$ schließlich ist die Menge der **Produktionen**: jede Produktion ist von der Form $A \rightarrow \alpha$, wobei A eine Variable und α ein beliebiger Term ist.

Eine Produktion $A \rightarrow \alpha \in P$ wird auch als **Regel** bezeichnet, A als **linke Seite** und α als **rechte Seite** der Regel.

Zur Vereinfachung der Notation dürfen Regeln $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$ mit identischer linken Seite zusammengefasst werden:
 $A \rightarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$.

Ein **Ableitungsschritt** $\delta \Rightarrow_G \delta'$ überführt einen Term δ durch Ersetzung genau einer Variable in einen Term δ' . In der formalen Darstellung schreiben wir $\delta = \beta A \gamma$ und $\delta' = \beta \alpha \gamma$, wobei $A \in V$ die genannte Variable ist, die durch einen Term α ersetzt wird. Der Ableitungsschritt $\beta A \gamma \Rightarrow_G \beta \alpha \gamma$ ist **zulässig**, wenn es eine Produktion $A \rightarrow \alpha \in P$ gibt.

Eine **Ableitung** ist eine beliebige Folge von zulässigen Ableitungsschritten: $\alpha_1 \Rightarrow_G \alpha_2 \Rightarrow_G \dots \Rightarrow_G \alpha_n$. Wir schreiben kurz $\alpha_1 \Rightarrow_G^* \alpha_n$ und sagen, dass α_n aus α_1 **ableitbar** ist. Der Index G kann dabei ausgelassen werden, sofern klar ist, bezüglich welcher Grammatik G die Ableitung durchgeführt wird.

Definition kontextfreie Sprache

Die von G beschriebene **formale Sprache** $\mathcal{L}[G]$ ist die Menge aller Wörter w , die aus dem Startsymbol ableitbar sind:

$$\mathcal{L}[G] := \{w \in \Sigma^* \mid S \Rightarrow^* w\}.$$

Eine Sprache $L \subseteq \Sigma^*$ heißt **kontextfrei**, wenn sie durch eine kontextfreie Grammatik G beschrieben werden kann, d.h. wenn $L = \mathcal{L}[G]$ gilt.

Als Beispiel betrachten wir eine kontextfreie Grammatik G_1 für einfache arithmetische Ausdrücke über dem Alphabet $\Sigma_1 = \{0, 1, \dots, 9, +, *\}$. $G_1 = (V_1, \Sigma_1, P_1, S)$ ist folgendermaßen definiert:

$$V_1 := \{S, T\}$$

$$P_1 := \{S \rightarrow T,$$

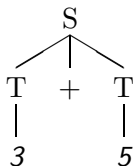
$$T \rightarrow T + T \mid T * T$$

$$T \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9\}.$$

Die Variable T repräsentiert dabei jeweils einen arithmetischen Term. Das Wort $w = 3 + 5$ gehört zu $\mathcal{L}[G_1]$, da es aus S ableitbar ist: $S \Rightarrow T \Rightarrow T + T \Rightarrow 3 + T \Rightarrow 3 + 5$, also kurz $S \Rightarrow^* 3 + 5$.

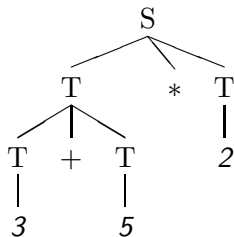
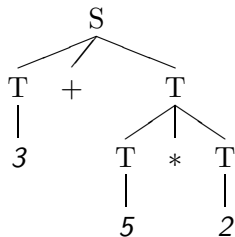
Ableitungsbaum für

$$S \Rightarrow T \Rightarrow T + T \Rightarrow 3 + T \Rightarrow 3 + 5$$



Die Ableitung von w bezüglich der Grammatik G_1 ist nicht eindeutig: eine andere mögliche Ableitung ist $S \Rightarrow T \Rightarrow T + T \Rightarrow T + 5 \Rightarrow 3 + 5$. Beide Varianten führen jedoch auf denselben Ableitungsbaum.

Ein Wort w kann bezüglich einer CFG G auch mehrere verschiedene Ableitungsbäume besitzen. Z.B. hat $w = 3 + 5 * 2 \in \mathcal{L}[G_1]$ die folgenden beiden Ableitungsbäume:



Die beiden Ableitungsbäume weisen w unterschiedliche **Struktur** zu, was für Anwendungen von großer Bedeutung ist. Man bezeichnet den Ableitungsbaum daher auch als **Analyse** von w durch die Grammatik. (Man denke z.B. an ein Taschenrechnerprogramm, das arithmetische Ausdrücke anhand ihres Ableitungsbaums auswertet. In diesem Fall wäre die linke Analyse die gewünschte, da $*$ stärker bindet als $+$). Eine Grammatik, in der es ein Wort w mit mehreren verschiedenen Ableitungsbäumen (bzw. Linksableitungen) gibt, heißt **mehrdeutig** / **ambig**.

Im obigen Beispiel wäre es wünschenswert, G_1 so abzuändern, dass $w = 3 + 5 * 2$ nur noch eine Analyse besitzt (nämlich die durch den linken Baum dargestellte). Eine solche Grammatik ist $G_2 = (V_2, \Sigma_1, P_1, S)$ mit

$$V_2 := \{S, P\}$$

$$P_2 := \{S \rightarrow S + S | P,$$

$$S \rightarrow 0|1|2|3|4|5|6|7|8|9,$$

$$P \rightarrow P * P,$$

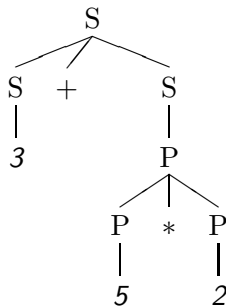
$$P \rightarrow 0|1|2|3|4|5|6|7|8|9\}.$$

Dabei stehen die Variablen S und P anschaulich für *Summe* und *Produkt*.

Verbesserte Grammatik:

Beispiel $w = 3 + 5 * 2$

Bezüglich G_2 besitzt w nur noch die folgende eindeutige Analyse:



1 Kontextfreie Grammatiken

2 Top-down parsing

3 CYK

cfgtopdown.odp

- 1 Kontextfreie Grammatiken
- 2 Top-down parsing
- 3 CYK

cyk.odp

cyk,2009.odp

time,flies,cyk.pdf

- Formale Definition CFG:
Terminale, Variablen, Startsymbol, Produktionen
- CYK
- Ableitungsbäume