

Kapitel: Die Chomsky Hierarchie

Definition

Eine **Grammatik** $G = (\Sigma, V, S, P)$ besteht aus:

- einem endlichen Alphabet Σ ,
- einer endlichen Menge V von Variablen (oder Nichtterminalen) mit $\Sigma \cap V = \emptyset$,
- dem Startsymbol $S \in V$ und
- einer endlichen Menge P von Produktionen der Form $u \rightarrow v$ mit

$$u \in (\Sigma \cup V)^* V (\Sigma \cup V)^* \quad \text{und} \quad v \in (\Sigma \cup V)^*$$

Eine Grammatik G erzeugt eine Sprache $L(G) \subseteq \Sigma^*$ wie folgt:

Beginne mit dem Startsymbol S und wende dann Produktionen an:

Eine Produktion $u \rightarrow v$ ersetzt ein Vorkommen von u durch ein Vorkommen von v .

Details: nächste Folie

Die von einer Grammatik erzeugte Sprache

Sei $G = (\Sigma, V, S, P)$ eine Grammatik.

- Für eine Produktion $u \rightarrow v$ und ein Wort $w_1 = xuy$ wird das Wort $w_2 = xvy$ abgeleitet. Wir schreiben

$$xuy \Rightarrow xvy.$$

- Für Worte $r, s \in (\Sigma \cup V)^*$ schreiben wir

$$r \xRightarrow{*} s \quad :\Leftrightarrow \quad \text{Es gibt Worte } w_1 = r, w_2, \dots, w_k = s, \text{ für } k \geq 1, \text{ so dass} \\ w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_k.$$

s kann in null, einem oder mehreren Schritten aus r abgeleitet werden.

- $L(G) = \{w \in \Sigma^* : S \xRightarrow{*} w\}$ ist die von der Grammatik G erzeugte Sprache.

Beispiel: Eine Grammatik, die die Sprache $L := \{a^n b^n c^n : n \in \mathbb{N}_{>0}\}$ erzeugt

Wir wissen bereits, dass diese Sprache nicht kontextfrei ist — somit haben wir keine Chance, eine kontextfreie Grammatik zu finden, die L erzeugt.

Eine nicht-kontextfreie Grammatik $G = (\Sigma, V, S, P)$ mit $L(G) = L$:

- ▶ $\Sigma = \{a, b, c\}$
- ▶ $V = \{S, B, C\}$
- ▶ Startsymbol S
- ▶ P besteht aus folgenden Produktionen:

S	\rightarrow	$aSBC \mid aBC$
CB	\rightarrow	BC
aB	\rightarrow	ab
bB	\rightarrow	bb
bC	\rightarrow	bc
cC	\rightarrow	cc

Behauptung: $L(G) = \{a^n b^n c^n : n \in \mathbb{N}_{>0}\}$.

Beweis: siehe Tafel.

- (0) Allgemeine Grammatiken, wie sie am Anfang dieses Kapitels definiert wurden, heißen **Typ 0-Grammatiken**.
- (1) Eine Grammatik $G = (\Sigma, V, S, P)$ heißt **monoton** oder **Typ 1-Grammatik**, falls für jede Produktion $(u \rightarrow v) \in P$ gilt: $|u| \leq |v|$.

Beispiel: Die Grammatik für $\{a^n b^n c^n : n \geq 1\}$ der vorherigen Folie ist monoton.

Eine Grammatik $G = (\Sigma, V, S, P)$ heißt **kontextsensitiv**, falls jede Produktion in P von der Form $uXv \rightarrow uxv$ ist, mit $X \in V$, $u, x, v \in (\Sigma \cup V)^*$ und $x \neq \varepsilon$.

Beachte: Per Definition ist jede kontextsensitive Grammatik ist monoton. Umgekehrt kann man zeigen (hier ohne Beweis), dass es für jede monotone Grammatik eine kontextsensitive Grammatik gibt, die dieselbe Sprache erzeugt. In der Literatur wird daher der Begriff "**kontextsensitive Grammatik**" oft als **Synonym** für "**monotone Grammatik**" verwendet.

Sonderfall: Erzeugen des leeren Worts

Per Definition gilt für monotone Grammatiken G , dass $\varepsilon \notin L(G)$.

Um auch Sprachen erzeugen zu können, die das leere Wort enthalten, erlauben wir als Typ 1-Grammatiken auch Grammatiken der Form

$$(\Sigma, V \cup \{S'\}, S', P \cup \{S' \rightarrow \varepsilon \mid S\}),$$

wobei (Σ, V, S, P) eine monotone Grammatik ist.

- (2) Eine Grammatik $G = (\Sigma, V, S, P)$ heißt **kontextfrei** oder **Typ 2-Grammatik**, falls jede Produktion in P von der Form $X \rightarrow x$ ist, mit $X \in V$ und $x \in (\Sigma \cup V)^*$.

Wir wissen bereits, dass die Typ 2-Grammatiken genau die kontextfreien Sprachen erzeugen.

Beachte: Für jede Typ 2-Grammatik G kann man leicht eine Typ 1-Grammatik konstruieren, die dieselbe Sprache erzeugt:

1. Sei $V_\varepsilon = \{X \in V : X \xrightarrow{*} \varepsilon\}$.
 2. Für jedes $A \in V_\varepsilon$ tue Folgendes:
Entferne aus P alle Produktionen der Form $A \rightarrow \varepsilon$ und füge für jede Produktion der Form $(B \rightarrow uAv) \in P$ mit $uv \neq \varepsilon$ eine zusätzliche Produktion der Form $B \rightarrow uv$ ein.
 3. Falls $S \in V_\varepsilon$, so wähle ein neues Startsymbol S' und füge die zusätzlichen Produktionen $S' \rightarrow S$ und $S' \rightarrow \varepsilon$ ein.
- (3) Eine Grammatik $G = (V, \Sigma, S, P)$ heißt **regulär** oder **Typ 3-Grammatik**, falls jede Produktion in P von der Form $X \rightarrow \varepsilon$ oder $X \rightarrow aY$ mit $X, Y \in V$ und $a \in \Sigma$ ist.

Wir wissen bereits, dass die Typ 3-Grammatiken genau die regulären Sprachen erzeugen.

Die Sprachen der Chomsky Hierarchie

Für jedes $i \in \{0, 1, 2, 3\}$ sei

\mathcal{L}_i die Klasse aller Sprachen, die von Typ i -Grammatiken erzeugt werden.

Es gilt:

- ▶ \mathcal{L}_3 ist die Klasse aller **regulären Sprachen** (kurz auch: **REG** oder **Typ 3-Sprachen**).
- ▶ \mathcal{L}_2 ist die Klasse aller **kontextfreien Sprachen** (kurz auch: **CFL**, für “context-free language”, oder **Typ 2-Sprachen**).
- ▶ \mathcal{L}_1 wird die Klasse aller **kontextsensitiven Sprachen** genannt (kurz auch: **CSL**, für “context-sensitive language”, oder **Typ 1-Sprachen**).
- ▶ \mathcal{L}_0 wird die Klasse aller **Typ 0-Sprachen** genannt.

Es gilt:

$$\mathcal{L}_0 \quad \not\supseteq \quad \mathcal{L}_1 \quad \not\supseteq \quad \mathcal{L}_2 \quad \not\supseteq \quad \mathcal{L}_3.$$

$\underbrace{\hspace{10em}}_{H \in \mathcal{L}_0 \setminus \mathcal{L}_1} \quad \underbrace{\hspace{10em}}_{\{a^n b^n c^n : n \geq 1\} \in \mathcal{L}_1 \setminus \mathcal{L}_2} \quad \underbrace{\hspace{10em}}_{\{a^n b^n : n \geq 1\} \in \mathcal{L}_2 \setminus \mathcal{L}_3}$

Frage:

- ▶ Was genau sind die Typ 0-Sprachen?
- ▶ Was genau sind die Typ 1-Sprachen?

Antwort:

Satz:

- (a) Die **Typ 0**-Sprachen sind genau die **semi-entscheidbaren** Sprachen $L \subseteq \Sigma^*$.
- (b) Die **Typ 1**-Sprachen sind genau die Sprachen $L \subseteq \Sigma^*$, die von einer nichtdeterministischen Turingmaschine mit **linear beschränktem Platz** entschieden werden können. Solche Turingmaschinen heißen auch **linear beschränkte Automaten** (kurz: **LBA**).

Beweis von (a):

- ▶ Jede Typ 0-Sprache ist semi-entscheidbar:

Dies erhält man leicht durch einen Algorithmus, der bei Eingabe von $w \in \Sigma^*$ nach und nach sämtliche möglichen Ableitungen der Typ 0-Grammatik G durchprobiert und anhält, falls er eine Ableitung für w gefunden hat.

- ▶ Jede semi-entscheidbare Sprache wird von einer Typ 0-Grammatik erzeugt:

Sei T eine Turingmaschine, die $L \subseteq \Sigma^*$ semi-entscheidet.

Wir konstruieren eine Grammatik G mit $L(G) = L$.

- Berechnungen von T beginnen natürlich stets mit der Eingabe w , aber Ableitungen von w enden mit w .
- Also sollten wir die Grammatik G so konstruieren, dass die Berechnungen von T „rückwärts“ simuliert werden.

O.B.d.A gibt es nur einen Zustand q_h , in dem T anhält, und wenn T hält, ist das Band leer.

Angenommen, wir unterbrechen die Berechnung von T auf Eingabe w zu einem beliebigen Zeitpunkt.

Wir beschreiben die aktuelle Konfiguration der TM genauso wie im Beweis der Unentscheidbarkeit des Postschen Korrespondenzproblems:

Die Situation, in der T im Zustand q ist, die Bandbeschriftung $\alpha_1 \cdots \alpha_{i-1} \alpha_i \cdots \alpha_N$ hat und der Kopf auf Position i steht, beschreiben wir durch das Wort

$$\alpha_1 \cdots \alpha_{i-1} q \alpha_i \cdots \alpha_N \in (\Gamma \cup Q)^*.$$

- Unsere Grammatik G erzeugt diese Beschreibungen als Zwischenschritte.
- G erzeugt zuerst die Endkonfiguration, in der T mit leerem Band anhält, indem sie die folgende Produktion nutzt:

$$S \rightarrow \square q_h \square$$

- Für jeden Befehl $\delta(q, a) = (q', b, \leftarrow)$ nehmen wir die Produktion

$$q'cb \rightarrow cqa$$

für alle $c \in \Gamma$ auf:

Wenn die Konfiguration $* \dots * q'cb * \dots *$ schon erzeugt wurde, können wir damit die mögliche Vorgänger-Konfiguration $* \dots * cqa * \dots *$ erzeugen.

- Für jeden Befehl $\delta(q, a) = (q', b, \rightarrow)$ nehmen wir die Produktion $bq' \rightarrow qa$ auf.
- Für jeden Befehl $\delta(q, a) = (q', b, \downarrow)$ nehmen wir die Produktion $q'b \rightarrow qa$ auf.

- Am Ende der Ableitung haben wir ein Wort der Form $\square^* q_0 w \square^*$ erzeugt.
- Jetzt lösche q_0 und alle \square -Symbole mit weiteren Produktionen.

Die dadurch entstehende Grammatik erzeugt genau die Worte, die von T akzeptiert werden. □

Typ 1-Sprachen und linearer Speicherplatz

Satz:

$L \subseteq \Sigma^*$ ist genau dann kontextsensitiv, wenn $L \subseteq \Sigma^*$ von einer nichtdeterministischen Turingmaschinen auf linearem Speicherplatz entschieden wird.

Insbesondere ist jede kontextsensitive Sprache entscheidbar.

Beweis:

“ \implies ”: Es gelte $L = L(G)$, für eine Typ-1 Grammatik G .

Warum kann L von einer nichtdeterministischen Turingmaschine T auf linearem Platz akzeptiert werden?

- T “rät” eine Ableitung $S \xRightarrow{*} w$ von w .
- Falls $S \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_\ell$ mit $u_\ell = w$ eine Ableitung von w ist, so ist $|u_i| \leq |u_{i+1}| \leq |w|$, für alle $i < \ell$, da G monoton ist.
- Daher kann T zu jedem Zeitpunkt i die Worte u_i, u_{i+1}, w auf seinem Band speichern.
- Insgesamt reicht linearer Speicherplatz $O(|w|)$ aus!

“ \impliedby ”: Übung (ähnlich wie der entsprechende Beweis für Typ 0-Sprachen). □

Die Chomsky-Hierarchie besteht aus folgenden Klassen von Sprachen:

- ▶ **Typ 0-Sprachen**, d.h. Sprachen, die von **allgemeinen Grammatiken** erzeugt werden.
Dies sind genau die **semi-entscheidbaren Sprachen**.
- ▶ **Typ 1-Sprachen**, d.h. Sprachen, die von **monotonen (oder kontextsensitiven) Grammatiken** erzeugt werden.
Dies sind genau die Sprachen, die von **linear beschränkten Automaten** (d.h. nichtdeterministische Turingmaschinen mit linear beschränktem Platz) entschieden werden.
- ▶ **Typ 2-Sprachen**, d.h. Sprachen, die von **kontextfreien Grammatiken** erzeugt werden.
Dies sind genau die Sprachen, die von **nichtdeterministischen Kellerautomaten** akzeptiert werden.
- ▶ **Typ 3-Sprachen**, d.h. Sprachen, die von **regulären Grammatiken** erzeugt werden.
Dies sind genau die Sprachen, die von **endlichen Automaten** akzeptiert werden.

Trennende Beispiele:

- ▶ Eine Typ 2-Sprache (kontextfrei), die keine Typ 3-Sprache (regulär) ist:

$$\{a^n b^n : n \geq 1\}.$$

- ▶ Eine Typ 1-Sprache (kontextsensitiv), die keine Typ 2-Sprache (kontextfrei) ist:

$$\{a^n b^n c^n : n \geq 1\}.$$

- ▶ Eine Typ 0-Sprache (semi-entscheidbar), die keine Typ 1-Sprache ist:

$$\text{Halteproblem } H := \{ \langle T \rangle w : T \text{ ist eine TM, die bei Eingabe } w \text{ hält} \} \subseteq \{0, 1\}^*$$

- ▶ Eine Sprache, die keine Typ 0-Sprache ist:

$$\bar{H} := \{0, 1\}^* \setminus H.$$