

FORMALE SYSTEME

2. Vorlesung: Grammatiken und die Chomsky-Hierarchie

Hannes Straß

Folien: © Markus Krötzsch, <https://iccl.inf.tu-dresden.de/web/FS2020>, CC BY 3.0 DE

TU Dresden, 14. Oktober 2021

Sprachen beschreiben

Wiederholung

- Formale Sprachen sind in Praxis und Theorie sehr wichtig.
- Ein **Alphabet** ist eine nichtleere, endliche Menge von Symbolen.
- Ein **Wort** ist eine endliche Sequenz von Symbolen.
- Eine **Sprache** ist eine Menge von Wörtern.
- Es gibt viele **Operationen auf Sprachen** ($\cap, \cup, \bar{}, \circ, *, +, \backslash$).
- Man kann Sprachen auf viele Arten beschreiben (von denen wir noch einige genauer kennen lernen werden).

Wie kann man Sprachen beschreiben?

Die Operationen aus der vorigen Vorlesung können das schon ganz gut:

Beispiel (Wiederholung): Die Menge aller gültigen Schreibweisen für Dezimalzahlen kann wie folgt definiert werden:

$$\text{Dezimalzahl} = (\{\epsilon\} \cup \{+, -\}) \circ (\{\mathbf{0}\} \cup (\mathbf{Z} \setminus \{\mathbf{0}\}) \circ \mathbf{Z}^*) \circ (\{\epsilon\} \cup \{.\} \mathbf{Z}^+)$$

mit $\mathbf{Z} = \{\mathbf{0}, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$.

Wir haben hier eine unendliche Sprache mit einfachen Grundoperationen aus endlichen Sprachen konstruiert.

Funktioniert das für alle Sprachen?

Wie viele Wörter gibt es?

Satz: Es gibt **abzählbar viele** Wörter über jedem Alphabet Σ . Jede Sprache ist also entweder endlich oder abzählbar unendlich.

Wdh.: Eine Menge M ist **abzählbar** gdw. es eine injektive Funktion $f : M \rightarrow \mathbb{N}$ gibt.

Idee: Zum „Abzählen“ reihen wir die Wörter in Σ^* der Länge nach auf, beginnend mit den kurzen Wörtern.

Zum Beispiel: $\epsilon, a, b, aa, ab, ba, bb, aaa, aab, \dots$

Beweis (die gleiche Idee, aber etwas formaler):

- Laut Definition ist $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots = \bigcup_{i \geq 0} \Sigma^i$.
- Jede der Mengen Σ^i ist endlich (es gibt nur endlich viele Wörter der Länge i).
- Es gibt abzählbar viele Mengen Σ^i (eine für jede natürliche Zahl)
- Eine abzählbare Vereinigung endlicher Mengen ist abzählbar.
- Also ist Σ^* abzählbar. □

Wie viele Sprachen gibt es?

Satz: Es gibt **überabzählbar viele** Sprachen über jedem beliebigen Alphabet Σ .

Beweis:

- Die Menge aller Wörter Σ^* ist **abzählbar unendlich**, selbst wenn Σ nur ein Symbol enthält.
- Eine Sprache ist eine beliebige Teilmenge von Σ^* .
- Die Menge aller Sprachen über Σ ist die Menge aller Teilmengen (**Potenzmenge**) von Σ^* .
- Die Kardinalität der Potenzmenge ist immer größer als die der Menge selbst (**Satz von Cantor**).
- Also gibt es überabzählbar viele Sprachen. □

Aus Kardinalitäten lernen

Zusammenfassung: Über jedem beliebigen Alphabet Σ gibt es

- **abzählbar** viele Wörter
- **überabzählbar** viele Sprachen

Beobachtung: Wir beschreiben Sprachen mit „Wörtern“:

- Natürlichsprachliche Texte sind Wörter
- Computerprogramme sind Wörter
- Mathematische Formeln sind Wörter

Es gibt **viel mehr** Sprachen als Texte, Programme oder Formeln!

Satz: Fast alle Sprachen können in keiner Weise endlich beschrieben werden. Fast alle Sprachen werden durch kein Computerprogramm korrekt erkannt.

Ein starker Satz, aber er gibt uns keine Hinweise, welche Sprachen man wie beschreiben kann. Zum Beispiel gibt es eindeutig beschreibbare Sprachen, für deren Erkennung es kein Computerprogramm gibt.

Grammatiken

Sprachen mit Grammatiken beschreiben

Grammatiken sind die wichtigste Methode zur Beschreibung von Sprachen.

Sie verwenden **Ersetzungsregeln** um Wörter zu **erzeugen**.

Beispiel: Die folgende Grammatik erzeugt Dezimalzahlen.

Dezimalzahl	\rightarrow Vorzeichen DZahl DZahl	(1)
Vorzeichen	\rightarrow + -	(2)
DZahl	\rightarrow NatZahl NatZahl . Ziffernreihe	(3)
NatZahl	\rightarrow Ziffer ZifferAb1 Ziffernreihe	(4)
Ziffernreihe	\rightarrow Ziffer Ziffer Ziffernreihe	(5)
ZifferAb1	\rightarrow 1 2 ... 9	(6)
Ziffer	\rightarrow 0 1 2 ... 9	(7)

Grammatiken anwenden

Dezimalzahl	\rightarrow Vorzeichen DZahl DZahl	(1)
Vorzeichen	\rightarrow + -	(2)
DZahl	\rightarrow NatZahl NatZahl . Ziffernreihe	(3)
NatZahl	\rightarrow Ziffer ZifferAb1 Ziffernreihe	(4)
Ziffernreihe	\rightarrow Ziffer Ziffer Ziffernreihe	(5)
ZifferAb1	\rightarrow 1 2 ... 9	(6)
Ziffer	\rightarrow 0 1 2 ... 9	(7)

Wir können folgende Ersetzungsschritte ausführen:

$$\begin{aligned}
 \text{Dezimalzahl} &\stackrel{(1)}{\Rightarrow} \text{Vorzeichen DZahl} \stackrel{(2)}{\Rightarrow} -\text{DZahl} \stackrel{(3)}{\Rightarrow} -\text{NatZahl.Ziffernreihe} \\
 &\stackrel{(4)}{\Rightarrow} -\text{Ziffer.Ziffernreihe} \stackrel{(7)}{\Rightarrow} -3.\text{Ziffernreihe} \\
 &\stackrel{(5)}{\Rightarrow} -3.\text{Ziffer Ziffernreihe} \stackrel{(5)}{\Rightarrow} -3.\text{Ziffer Ziffer} \\
 &\stackrel{(7)}{\Rightarrow} -3.1 \text{ Ziffer} \stackrel{(7)}{\Rightarrow} -3.14
 \end{aligned}$$

Definition: Grammatik

Eine **Grammatik** G ist ein 4-Tupel $G = \langle V, \Sigma, P, S \rangle$ bestehend aus:

- V eine Menge von **Variablenamen**,
- Σ ein Alphabet, disjunkt zu V (d.h., $\Sigma \cap V = \emptyset$),
- P eine Menge von **Produktionsregeln** der Form $w \rightarrow v$ für beliebige Wörter w und v über $\Sigma \cup V$, wobei w mindestens eine Variable enthält (d.h. $w, v \in (\Sigma \cup V)^*$ und $w \notin \Sigma^*$),
- S eine **Startvariable** aus V (d.h. $S \in V$).

Die Elemente von Σ nennt man auch **Terminalsymbole**, die Elemente von V analog **Nichtterminalsymbole**.

Die zuvor verwendete Schreibweise $w \rightarrow v_1 | v_2 | \dots | v_n$ ist eine Abkürzung für die Produktionsregeln $w \rightarrow v_1, w \rightarrow v_2, \dots, w \rightarrow v_n$.

Beispiele: Grammatiken

Beispiel: Eine Grammatik $\langle V, \Sigma, P, \text{Start} \rangle$ sei gegeben durch $V = \{\text{Start, Summe, Var}\}$, $\Sigma = \{(\text{,}, +, *, \text{x}, \text{y}, \text{z})\}$ und Regeln P :

$$\begin{aligned}
 \text{Start} &\rightarrow \text{Summe} | \text{Summe} * \text{Start} \\
 \text{Summe} &\rightarrow (\text{Var} + \text{Var}) \\
 \text{Var} &\rightarrow \text{x} | \text{y} | \text{z}.
 \end{aligned}$$

Diese Grammatik erzeugt zum Beispiel das Wort $(\text{x} + \text{y}) * (\text{x} + \text{z})$.

Beispiel: Die selbe Sprache wird auch durch die Grammatik $\langle \{\text{Start, Var}\}, \Sigma, P', \text{Start} \rangle$ definiert, mit Σ wie oben und Regeln P' :

$$\begin{aligned}
 \text{Start} &\rightarrow (\text{Var} + \text{Var}) \\
 (\text{Var} + \text{Var}) &\rightarrow (\text{Var} + \text{Var}) * (\text{Var} + \text{Var}) \\
 \text{Var} &\rightarrow \text{x} | \text{y} | \text{z}.
 \end{aligned}$$

Es ist also erlaubt, (1) Terminalsymbole durch Regeln zu verändern und (2) mehr als ein Symbol auf einmal zu ersetzen.

Beispiele: Diskussion

Die Beispiele werfen **wichtige Fragen** auf:

- Wie kann man herausfinden, ob zwei unterschiedliche Grammatiken die selbe Sprache beschreiben?
- Muss man Grammatiken so allgemein definieren oder genügen auch vereinfachte Formen?
- Kann man Grammatiken (automatisch) vereinfachen, ohne die Sprache zu verändern?

Antworten in dieser Vorlesung ...

Erst einmal sollten wir die „**Sprache einer Grammatik**“ ordentlich definieren.

Definition: Sprache einer Grammatik

Die **von einer Grammatik** $G = \langle V, \Sigma, P, S \rangle$ erzeugte Sprache $L(G)$ besteht aus allen Wörtern über Σ , die man von S ableiten kann:

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$$

Vereinfachung: Für einfache Beispiele nehmen wir ab jetzt an, dass alle Großbuchstaben Variablen und alle Kleinbuchstaben Terminalsymbole sind, wobei S die Startvariable ist.

Beispiel: Wir betrachten die Grammatik G_{id} mit den Regeln:

$$S \rightarrow BA \quad A \rightarrow BA \mid ZA \mid \epsilon \quad B \rightarrow b \quad Z \rightarrow z$$

(Vereinfachte Version der „Bezeichner“, wobei alle Buchstaben als b und alle Ziffern als z abgekürzt werden.)

Dann gilt $L(G_{id}) = \{b\} \circ \{b, z\}^*$.

Definition: Ableitung

Sei $\langle V, \Sigma, P, S \rangle$ eine Grammatik. Die **1-Schritt-Ableitungsrelation** ist eine binäre Relation \Rightarrow zwischen Wörtern aus $(V \cup \Sigma)^*$, so dass $u \Rightarrow v$ genau dann, wenn:

$$u = w_1 x w_2 \text{ und } v = w_1 y w_2 \text{ und es gibt eine Regel } x \rightarrow y \in P$$

wobei $w_1, w_2, x, y \in (V \cup \Sigma)^*$ beliebige Wörter sind.

Die **Ableitungsrelation** \Rightarrow^* ist der reflexive transitive Abschluss von \Rightarrow , das heißt $u \Rightarrow^* v$ genau dann, wenn:

$$u = w_0 \Rightarrow w_1 \Rightarrow \dots \Rightarrow w_{n-1} \Rightarrow w_n = v$$

wobei $n \in \mathbb{N}$ und $w_0, \dots, w_n \in (V \cup \Sigma)^*$ beliebige Wörter sind. Insbesondere gilt $u \Rightarrow^* u$ für alle $u \in (V \cup \Sigma)^*$ (Fall $n = 0$).

Anmerkung: Der Begriff „Herleitungsrelation“ ist auch gebräuchlich. Wir verwenden „Ableitung“ und „Herleitung“ synonym.

Anmerkung 2: Manche Autor:innen schreiben \vdash statt \Rightarrow .

Beobachtung: Der Parameter n oben korrespondiert zur Anzahl der einzelnen Ableitungsschritte.

Beweis zum Beispiel (1)

Für die Grammatik G_{id} mit den Regeln

$$S \rightarrow BA \quad A \rightarrow BA \mid ZA \mid \epsilon \quad B \rightarrow b \quad Z \rightarrow z$$

gilt $L(G_{id}) = \{b\} \circ \{b, z\}^*$?

Das müsste man erst einmal beweisen ...

Wdh.: Um $M = N$ für Mengen M, N zu beweisen, zeigen wir (1) $M \subseteq N$ und (2) $N \subseteq M$.

Beweis: Wir zeigen die beiden Richtungen

(1) $L(G_{id}) \subseteq \{b\} \circ \{b, z\}^*$ und (2) $L(G_{id}) \supseteq \{b\} \circ \{b, z\}^*$ getrennt.

(1) Die Behauptung ist: „Jedes Wort in $L(G_{id})$ ist in $\{b\} \circ \{b, z\}^*$.“

- $L(G_{id}) \subseteq \{b, z\}^*$, weil $\{b, z\}$ das Alphabet von G_{id} ist
- jedes Wort in $L(G_{id})$ beginnt mit b , weil S nur durch BA ersetzt werden kann und B nur durch b
- also ist $L(G_{id}) \subseteq \{b\} \circ \{b, z\}^*$

Beweis zum Beispiel (2)

$S \rightarrow BA \quad A \rightarrow BA \mid ZA \mid \epsilon \quad B \rightarrow b \quad Z \rightarrow z$

Beweis (Fortsetzung): Es fehlt noch die Rückrichtung.

(2) Die Behauptung ist: „Jedes Wort in $\{b\} \circ \{b, z\}^*$ ist in $L(G_{id})$.“

- Sei $bs_1 \dots s_n$ mit $n \geq 0$ ein beliebiges Wort in $\{b\} \circ \{b, z\}^*$ (also $s_1, \dots, s_n \in \{b, z\}$)
- Wir zeigen $bs_1 \dots s_n \in L(G_{id})$ indem wir eine Ableitung angeben, d.h. wir zeigen $S \Rightarrow^* bs_1 \dots s_n$
- Die Ableitung beginnt immer so: $S \Rightarrow BA \Rightarrow bA$
- Jetzt müssen wir nur noch $A \Rightarrow^* s_1 \dots s_n$ zeigen.

Beweis zum Beispiel (4)

$S \rightarrow BA \quad A \rightarrow BA \mid ZA \mid \epsilon \quad B \rightarrow b \quad Z \rightarrow z$

Beweis (Fortsetzung): Der „rekursiv konstruierte Beweis“ ist natürlich eine **vollständige Induktion**:

Induktionsbehauptung: $A \Rightarrow^* s_1 \dots s_n$ für alle $n \geq 0$ und $s_1, \dots, s_n \in \{b, z\}$

Induktionsanfang: für $n = 0$ gilt $A \Rightarrow^* \epsilon$

Induktionsvoraussetzung: $A \Rightarrow^* s_1 \dots s_n$ gilt für n und alle $s_1, \dots, s_n \in \{b, z\}$

Induktionsschritt: $A \Rightarrow^* s_1 \dots s_{n+1}$ gilt für alle $s_1, \dots, s_{n+1} \in \{b, z\}$, denn:

- Wenn $s_1 = b$, dann $A \Rightarrow BA \Rightarrow bA$.
- Wenn $s_1 = z$, dann $A \Rightarrow ZA \Rightarrow zA$.
- Behauptung folgt, da $A \Rightarrow^* s_2 \dots s_{n+1}$ (Induktionsvoraussetzung).

Zusammenfassung Beweis Teil (2):

Für beliebige Wörter w gilt: Wenn $w \in \{b\} \circ \{b, z\}^*$, dann $w \in L(G_{id})$.

Anders gesagt: $\{b\} \circ \{b, z\}^* \subseteq L(G_{id})$. □

Beweis zum Beispiel (3)

$S \rightarrow BA \quad A \rightarrow BA \mid ZA \mid \epsilon \quad B \rightarrow b \quad Z \rightarrow z$

Beweis (Fortsetzung): Wir behaupten $A \Rightarrow^* s_1 \dots s_n$ für beliebige $n \geq 0$ und $s_1, \dots, s_n \in \{b, z\}$. Die Ableitung ist in jedem konkreten Fall leicht gefunden, zum Beispiel für $n = 2$ und $s_1 s_2 = zb$ ergibt sich $A \Rightarrow ZA \Rightarrow zA \Rightarrow zBA \Rightarrow zbA \Rightarrow zb$.

Aber es gibt unendlich viele Fälle!

~> Wir zeigen, wie man die Ableitung rekursiv konstruieren kann.

- Für $n = 0$ ist $s_1 \dots s_n = \epsilon$ und $A \Rightarrow \epsilon$ gilt direkt
- Für $n > 0$ und $s_1 = b$ verwenden wir $A \Rightarrow BA \Rightarrow bA$ und von da ab (rekursiv) die Ableitung $A \Rightarrow^* s_2 \dots s_n$.
- Für $n > 0$ und $s_1 = z$ verwenden wir $A \Rightarrow ZA \Rightarrow zA$ und von da ab (rekursiv) die Ableitung $A \Rightarrow^* s_2 \dots s_n$.

Die Rekursion terminiert, da wir mit jedem Rekursionsschritt kürzere Wörter ableiten.

Quiz: Sprache einer Grammatik

Die von einer Grammatik $G = \langle V, \Sigma, P, S \rangle$ erzeugte Sprache $L(G)$ besteht aus allen Wörtern über Σ , die man von S ableiten kann:

$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$$

Quiz: Wir betrachten die Grammatik $G = \langle V, \{0, 1\}, P, S \rangle$ mit den Regeln:

...

Die Chomsky-Hierarchie

Sprachen einteilen nach Grammatiken

Die meisten Sprachen können **nicht** mit Grammatiken beschrieben werden. (Es gibt überabzählbar viele Sprachen, aber nur abzählbar viele Grammatiken.)

- Die meisten **praktisch relevanten** Sprachen haben aber eine Grammatik.
- Manche Sprachen kann man sogar mit einfacheren **Sonderformen von Grammatiken** beschreiben.



Ministerio de Cultura de la Nación
2015, CC-BY-SA 2.0

Noam Chomsky, 2015

Man kann Sprachen nach Komplexität ihrer Grammatiken unterteilen.

Die Chomsky-Hierarchie

Die **Chomsky-Hierarchie** unterteilt Grammatiken in vier Stufen:

- **Typ 0:** beliebige Grammatiken
- **Typ 1^a:** **kontextsensitive Grammatiken:**
Alle Regeln $w \rightarrow v$ erfüllen die Bedingung $|w| \leq |v|$.
- **Typ 2:** **kontextfreie Grammatiken:**
Alle Regeln haben die Form $A \rightarrow v$ für eine Variable A .
- **Typ 3:** **reguläre Grammatiken:**
Alle Regeln haben eine der Formen

$$A \rightarrow cB \quad A \rightarrow c \quad A \rightarrow \epsilon$$

wobei A und B Variablen sind und c ein Terminalsymbol ist.

^aErste Version. Wir werden Typ 1 später noch leicht erweitern.

Chomsky-Hierarchie: Beispiele

Beispiel: Eine Typ-3-Grammatik (regulär):

$$S \rightarrow bA \quad A \rightarrow zA \quad A \rightarrow z$$

Beispiel: Eine Typ-2-Grammatik (kontextfrei, nicht regulär):

$$S \rightarrow bBz \quad B \rightarrow BzB \quad B \rightarrow \epsilon$$

Beispiel: Eine Typ-1-Grammatik (kontextsensitiv, nicht kontextfrei):

$$S \rightarrow bC \quad bC \rightarrow Sz \quad Sz \rightarrow bz$$

Beispiel: Eine Typ-0-Grammatik (nicht kontextsensitiv):

$$S \rightarrow bzDzDz \quad zD \rightarrow zzD \quad zzzDz \rightarrow z$$

Alle vier Grammatiken erzeugen die gleiche Sprache $\{b\} \circ \{z\}^+$.

Kontextfrei und kontextsensitiv

Der Namensgebung liegt die folgende Idee zugrunde:

Typische kontextsensitive Regel:

$$aBc \rightarrow avc$$

„B kann im Kontext von a und c durch v ersetzt werden.“

Typische kontextfreie Regel:

$$B \rightarrow v$$

„B kann in jedem Kontext durch v ersetzt werden.“

Allerdings können kontextsensitive Regeln auch andere Formen haben.

Quiz: Chomsky-Hierarchie

Die **Chomsky-Hierarchie** unterteilt Grammatiken in vier Stufen:

- **Typ 0:** beliebige Grammatiken
- **Typ 1^a:** **kontextsensitive Grammatiken:**
Alle Regeln $w \rightarrow v$ erfüllen die Bedingung $|w| \leq |v|$.
- **Typ 2:** **kontextfreie Grammatiken:**
Alle Regeln haben die Form $A \rightarrow v$ für eine Variable A.
- **Typ 3:** **reguläre Grammatiken:**
Alle Regeln haben eine der Formen $A \rightarrow cB$ $A \rightarrow c$ $A \rightarrow \epsilon$,
wobei A und B Variablen sind und c ein Terminalsymbol ist.

Quiz: Gegeben sei die Grammatik $G = \langle V, \{0, 1\}, P, S \rangle$ mit den Regeln

...

Typen von Sprachen

Die Chomsky-Hierarchie kann auf Sprachen angewendet werden:

Eine Sprache **L** ist genau dann **vom Typ i**, wenn sie durch eine Grammatik **G** von Typ **i** beschrieben wird, d.h. wenn $L(G) = L$.
Es gibt demnach **reguläre**, **kontextfreie** und **kontextsensitive** Sprachen.

Die Beispiele haben gezeigt: Man kann die gleiche Sprache oft mit Grammatiken verschiedenen Typs beschreiben.

Aber wir werden sehen: Für viele Sprachen gibt es einen maximalen Typ von Grammatik, z.B. gibt es kontextfreie Sprachen, die nicht regulär sind.

Beziehungen zwischen Chomskys Typen

Bilden die Typen wirklich eine Hierarchie

„Typ 3 \subseteq Typ 2 \subseteq Typ 1 \subseteq Typ 0“?

- Jede Grammatik von Typ 1, 2 oder 3 ist auch von Typ 0 (per Definition)
- Jede reguläre Grammatik ist kontextfrei (die Typ-3-Regeln $A \rightarrow cB$, $A \rightarrow c$ und $A \rightarrow \epsilon$ sind auch von Typ 2)
- Aber nicht jede kontextfreie Grammatik ist kontextsensitiv (Regeln der Form $A \rightarrow \epsilon$ sind nicht erlaubt)

Allgemein können Typ-1-Sprachen in unserer Definition nie das leere Wort enthalten!

Typ 1 mit ϵ -Regeln

Wir erweitern unsere Definition wie folgt:

Eine Grammatik $G = \langle V, \Sigma, P, S \rangle$ ist genau dann von Typ 1 (kontextsensitiv), wenn eine der folgenden Bedingungen gilt:

- (1) Alle Regeln $w \rightarrow v$ erfüllen die Bedingung $|w| \leq |v|$ (ursprüngliche Definition).
- (2) Es gibt eine Regel $S \rightarrow \epsilon$ und alle anderen Regeln $w \rightarrow v$ erfüllen zwei Bedingungen:
 - (a) $|w| \leq |v|$ (insbesondere ist also $v \neq \epsilon$);
 - (b) S kommt nicht in v vor.
 (Sonderfall mit einer einzelnen ϵ -Regel.)

Beispiel: Die reguläre Grammatik G mit den Regeln

$$S \rightarrow \epsilon \mid A \qquad A \rightarrow \text{nyan}A \qquad A \rightarrow \text{nyan}$$

ist kontextsensitiv nach der erweiterten Definition. $L(G) = \{\text{nyan}^*\}$.

Erzeugung ϵ -freier Grammatiken

Der folgende Algorithmus kann ϵ -Regeln eliminieren.

Eingabe: kontextfreie Grammatik (CFG) $G = \langle V, \Sigma, P, S \rangle$

Ausgabe: ϵ -freie CFG $G' = \langle V', \Sigma, P', S' \rangle$ mit $L(G') = L(G)$

- Initialisiere $P' := P$ und $V' := V$
- Berechne $V_\epsilon = \{A \in V \mid A \Rightarrow^* \epsilon\}$ (Details dazu später)
- Entferne alle ϵ -Regeln aus P'
- Solange es in P' eine Regel $B \rightarrow xAy$ gibt, mit

$$A \in V_\epsilon \qquad |x| + |y| \geq 1 \qquad B \rightarrow xy \notin P'$$

wähle eine solche Regel und setze $P' := P' \cup \{B \rightarrow xy\}$

- Falls $S \in V_\epsilon$, dann definiere ein neues Startsymbol $S' \notin V$, setze $V' := V' \cup \{S'\}$ und $P' := P' \cup \{S' \rightarrow S, S' \rightarrow \epsilon\}$.
- Falls $S \notin V_\epsilon$, dann verwenden wir einfach $S' := S$ als Startsymbol.

(CFG steht für **context-free grammar** und wird oft als Abkürzung verwendet.)

Kontextfrei vs. kontextsensitiv mit ϵ

Manche Typ-2-Grammatiken sind weiterhin nicht von Typ 1:

Beispiel: Die kontextfreie (und reguläre) Grammatik

$$S \rightarrow \text{bA} \qquad A \rightarrow \text{zA} \qquad A \rightarrow \epsilon$$

ist nicht kontextsensitiv, nicht einmal mit ϵ -Erweiterung.

Eine kontextfreie Grammatik $G = \langle V, \Sigma, P, S \rangle$ ist ϵ -frei, wenn eine der folgenden Bedingungen gilt:

- (1) Es gibt keine Regel der Form $A \rightarrow \epsilon$
- (2) Es gibt eine Regel $S \rightarrow \epsilon$ und bei allen anderen Regeln $A \rightarrow v$ ist $v \neq \epsilon$ und S kommt nicht in v vor.

Offensichtlich: Jede ϵ -freie kontextfreie Grammatik ist kontextsensitiv.

Weniger offensichtlich: ϵ -freie kontextfreie Grammatiken erzeugen die gleichen Sprachen wie kontextfreie Grammatiken allgemein.

Idee: ϵ -Regeln durch „Anwendung“ auf andere Regeln entfernen. Geht das immer?

$$S \rightarrow \text{bA} \mid \text{b} \qquad A \rightarrow \text{zA} \mid \text{z} \qquad A \rightarrow \epsilon$$

Berechnung von V_ϵ

Die Variablen $A \in V$ mit $A \Rightarrow^* \epsilon$ kann man leicht rekursiv finden:

Eingabe: kontextfreie Grammatik $G = \langle V, \Sigma, P, S \rangle$

Ausgabe: V_ϵ für diese Grammatik

- Initialisiere $V_\epsilon := \{A \mid A \rightarrow \epsilon \in P\}$
- Solange es eine Regel $B \rightarrow A_1 \dots A_n \in P$ gibt, so dass gilt

$$B \notin V_\epsilon \qquad \text{und} \qquad A_i \in V_\epsilon \text{ für alle } i \in \{1, \dots, n\}$$

wähle eine solche Regel und setze $V_\epsilon := V_\epsilon \cup \{B\}$.

Beispiel: $S \rightarrow \text{aB} \mid \text{BAA} \qquad B \rightarrow \text{a} \mid \text{AC} \qquad A \rightarrow \text{bb} \mid \epsilon$

$$C \rightarrow \text{DS} \mid \epsilon \qquad D \rightarrow \text{aAS}$$

Initialisierung: $V_\epsilon = \{A, C\}$

Schritt 1: $V_\epsilon = \{A, C, B\}$

Schritt 2: $V_\epsilon = \{A, C, B, S\}$ – Terminierung

Beispiel: Eliminierung von ϵ -Regeln

$S \rightarrow bBz$ $B \rightarrow BzB$ $B \rightarrow \epsilon$

$V_\epsilon = \{B\}$

$P' = \{$
 $S \rightarrow bBz,$
 $B \rightarrow BzB,$
 ~~$B \rightarrow \epsilon,$~~
 $S \rightarrow bz,$
 $B \rightarrow Bz,$
 $B \rightarrow zB,$
 $B \rightarrow z,$
 $\}$

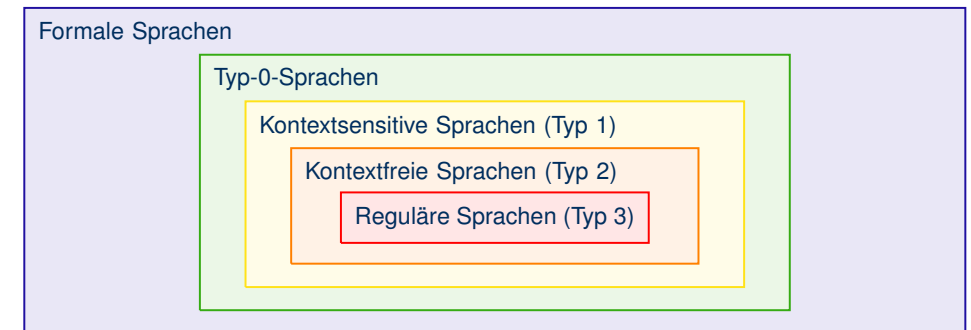
- ▶ • Initialisiere $P' := P$
- ▶ • Entferne alle ϵ -Regeln aus P'
- ▶ • Solange es in P' eine Regel $B \rightarrow xAy$ gibt, sodass $A \in V_\epsilon, |x| + |y| \geq 1, B \rightarrow xy \notin P'$
wähle eine solche Regel und setze $P' := P' \cup \{B \rightarrow xy\}$

Zusammenfassung der Ergebnisse

Wir können jede CFG in eine ϵ -freie CFG umschreiben. Daher gilt:

Satz: Jede kontextfreie Sprache ist kontextsensitiv im erweiterten Sinn (mit ϵ -Sonderfall).

Damit gelten alle Inklusionen der Hierarchie. Grafisch dargestellt:



Zusammenfassung und Ausblick

Nicht alle Sprachen sind mathematisch beschreibbar oder durch Computer entscheidbar.

Induktion (= rekursive Konstruktion von Beweisen) ist oft nötig, um Behauptungen über unendliche Sprachen zu beweisen

Grammatiken können (manche) Sprachen beschreiben. Die **Chomsky-Hierarchie** teilt Grammatiken in vier Typen ein.

Offene Fragen:

- Was gewinnt man durch die Beschränkung auf spezielle Typen?
- Sind die Typen wirklich unterschiedlich?
- Wie kommt Berechnung ins Spiel?