

Kapitel 3

Kontextfreie Sprachen

3.2

Chomsky Normalform (CNF)

Prof. Dr. Robert Preis
Fachbereich Informatik
Fachhochschule Dortmund
Robert.Preis@fh-dortmund.de

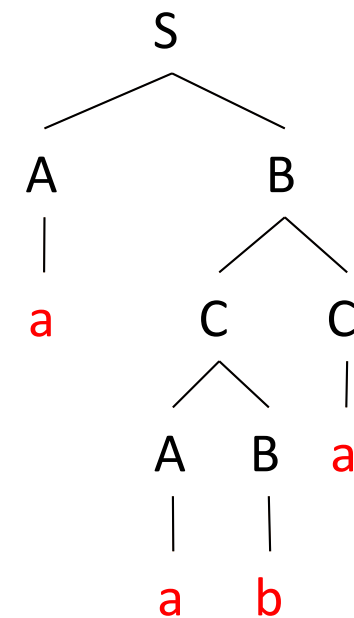
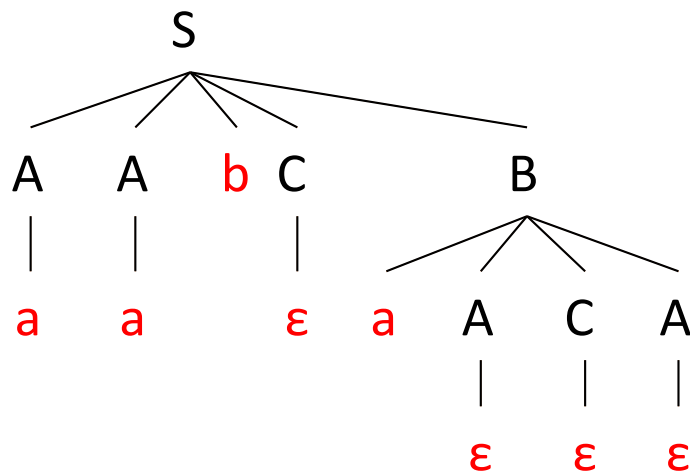
Alle Materialien (Folien, Übungsblätter, etc.) dieser Veranstaltung sind urheberrechtlich geschützt und nur von Teilnehmern dieser Veranstaltung und im Rahmen dieser zu verwenden. Eine anderweitige Verwendung oder Verbreitung ist nicht gestattet.

Ableitungsbäume von kontextfreien Grammatiken

$G=(V,T,P,S)$ mit $P = \{$
 $S \rightarrow AAbCB \mid ABBC,$
 $A \rightarrow BACS \mid a \mid \varepsilon,$
 $B \rightarrow CC \mid aACA \mid ab,$
 $C \rightarrow AB \mid a \mid \varepsilon,$
 $D \rightarrow BcCa \mid BAaC\}$

$H=(V,T,P,S)$ mit $P = \{$
 $S \rightarrow AB \mid BC,$
 $A \rightarrow BA \mid a,$
 $B \rightarrow CC \mid b,$
 $C \rightarrow AB \mid a \}$

Wie sehen die Ableitungsbäume für das Wort $w = aaba$ aus?



Bei G ist der Ableitungsbaum recht chaotisch, bei H ist es ein Binärbaum.

Grammatik in Chomsky-Normalform

Ein Grammatik $G_{\text{CNF}} = (V, T, P, S)$ ist in

Chomsky-Normalform (CNF)

wenn G **keine unnützen Variablen** enthält und jede Produktion eine der drei Formen hat:

1. $A \rightarrow BC$ mit $A, B, C \in V$
2. $A \rightarrow a$ mit $A \in V$ und $a \in T$
3. $S \rightarrow \varepsilon$ nur falls S Startvariable und S nirgends auf der rechten Seite enthalten ist. Dies ist wichtig für $\varepsilon \in L(G_{\text{CNF}})$.

*Jede kontextfreie Grammatik G kann man
in eine Grammatik G_{CNF} in CNF umwandeln
mit $L(G_{\text{CNF}}) = L(G)$.*

Umformung in CNF

Die Umformung einer allgemeinen kontextfreien Grammatik in eine kontextfreie Grammatik in Chomsky Normalform erfolgt in mehreren Schritten:

1. **Eliminierung von ϵ -Produktionen $A \rightarrow \epsilon$ falls A kein Startzustand**
(die Produktionen müssen rechts mindestens die Länge 1 haben)
2. **Eliminierung von Einheitsproduktionen $A \rightarrow B$**
(Produktionen der Länge 1 dürfen nicht aus einer Variablen bestehen)
3. **Eliminierung unnützer Symbole**
(keine unnütze Variablen erlaubt)
4. **Separieren von Terminalen und Variablen in Produktionen (X-Regel)**
(Produktionen haben rechts entweder Variablen oder Terminale)
5. **Aufspalten von Produktionen $A \rightarrow \alpha$ mit $|\alpha| > 2$ (Y-Regel)**
(Produktionen haben rechts höchstens die Länge 2)

Schritt 1: Elimination von ϵ -Übergängen

ϵ -Produktionen sind überflüssig (Ausnahme: $S \rightarrow \epsilon$).

Eine Variable $A \in V$ ist **eliminierbar**, falls $A \rightarrow^* \epsilon$.

(Sonderfall: falls S eliminierbar UND S ist irgendwo auf der rechten Seite:
erzeuge neue Startvariable S' und $S' \rightarrow \epsilon \mid S$.)

$P = \{$
 $S \rightarrow \mathbf{AB} \mid \mathbf{CD} \mid abc,$ S ist eliminierbar, weil A und B eliminierbar sind
 $A \rightarrow D \mid a\mathbf{AA}b \mid \mathbf{\epsilon},$ A ist eliminierbar
 $B \rightarrow b\mathbf{BB} \mid \mathbf{\epsilon},$ B ist eliminierbar
 $C \rightarrow \mathbf{c},$
 $D \rightarrow cD\}$

Neue Grammatik:

$P_1 = \{$
 $S \rightarrow \mathbf{AB} \mid \mathbf{CD} \mid abc \mid \mathbf{\epsilon} \mid \mathbf{A} \mid \mathbf{B},$
 $A \rightarrow D \mid a\mathbf{AA}b \mid \mathbf{aAb} \mid \mathbf{ab},$
 $B \rightarrow b\mathbf{BB} \mid \mathbf{bB} \mid \mathbf{b},$
 $C \rightarrow \mathbf{c},$
 $D \rightarrow cD\}$

*Die rechten Seiten
haben alle die
Länge ≥ 1
(außer $S \rightarrow \epsilon$) !*

Schritt 1: Elimination von ε -Übergängen: Verfahren

Bestimmung aller eliminierbaren Variablen:

1. Ist $A \rightarrow \varepsilon \in P$ dann ist A eliminierbar.
2. Ist $A \rightarrow X_1 \dots X_n \in P$ und alle X_i eliminierbar, dann ist A eliminierbar.

Verfahren zur Erzeugung einer Grammatik ohne ε -Regeln (ohne eliminierbaren Variablen):

1. Bestimme alle eliminierbaren Variablen.
2. Falls S eliminierbar ist:
 - falls S nirgendwo auf der rechten Seite vorkommt: erzeuge $S \rightarrow \varepsilon$
 - sonst erzeuge neue Startvariable S' und $S' \rightarrow \varepsilon \mid S$.
3. Für $A \rightarrow \alpha \in P$ mit eliminierbaren Symbolen X_1, \dots, X_m in α erzeuge neue Regeln (maximal $2^m - 1$): Streiche dabei jeweils alle möglichen Teilmengen von $\{X_1, \dots, X_m\}$ aus α .
4. Entferne alle Regeln der Form $A \rightarrow \varepsilon$ (nur falls A keine Startvariable).

Die neue Grammatik erzeugt die selben Wörter!

Die neue Grammatik enthält keine ε -Regeln mehr außer $S \rightarrow \varepsilon$!

Schritt 2: Elimination von Einheitsproduktionen

Einheitsproduktionen $A \rightarrow B$ dürfen nicht vorkommen und sind auch unnötig.

(A, B) ist ein **Einheitspaar**, falls $A \rightarrow^* B$.

$P_1 = \{$	
$S \rightarrow AB \mid CD \mid abc \mid \epsilon \mid A \mid B,$	Einheitspaare $(S,S), (S,A), (S,B), (S,D)$
$A \rightarrow D \mid aAAb \mid aAb \mid ab,$	Einheitspaare $(A,A), (A,D)$
$B \rightarrow bBB \mid bB \mid b,$	Einheitspaar (B,B)
$C \rightarrow c,$	Einheitspaar (C,C)
$D \rightarrow cD\}$	Einheitspaar (D,D)

Neue Grammatik:

$P_2 = \{$	
$S \rightarrow AB \mid CD \mid abc \mid \epsilon \mid aAAb \mid aAb \mid ab \mid bBB \mid bB \mid b \mid cD,$	
$A \rightarrow aAAb \mid aAb \mid ab \mid cD,$	
$B \rightarrow bBB \mid bB \mid b,$	
$C \rightarrow c,$	
$D \rightarrow cD\}$	

*Auf den rechten Seiten
steht keine Variable
mehr alleine!*

Schritt 2: Elimination von Einheitsproduktionen: Verfahren

Bestimmung aller Einheitspaare:

1. Alle Paare (A,A) für $A \in V$ sind Einheitspaare.
2. Ist (A, B) ein Einheitspaar und $B \rightarrow C \in P$ dann ist (A,C) Einheitspaar.

Verfahren zur Erzeugung einer Grammatik ohne Einheitsproduktionen:

1. Bestimme alle Einheitspaare in G .
2. Für jedes Einheitspaar (A, B) erzeuge neue Produktionen $\{ A \rightarrow \alpha \mid B \rightarrow \alpha \in P \text{ ist keine Einheitsproduktion} \}$.
3. Lösche alle Einheitsproduktionen.

Die neue Grammatik erzeugt die selben Wörter!

Die neue Grammatik enthält keine Einheitsproduktionen mehr!

Schritt 3: Elimination von unnützen Variablen

Nur nützliche Variablen machen Sinn.

Eine Symbol X ist **nützlich**, falls es eine Ableitung $S \rightarrow^* \alpha X \beta \rightarrow^* w \in T^*$ gibt, d.h. es ist **erreichbar** ($S \rightarrow^* \alpha X \beta$) und **erzeugend** ($X \rightarrow^* v \in T^*$).

$P_2 = \{$
 $S \rightarrow AB | CD | abc | \varepsilon | aAAb | aAb | ab | bBB | bB | b | cD,$ erzeugend / erreichbar
 $A \rightarrow aAAb | aAb | ab | cD,$ erzeugend / erreichbar
 $B \rightarrow bBB | bB | b,$ erzeugend / erreichbar
 $C \rightarrow c,$ erz./err., nach Löschung von $S \rightarrow CD$ **nicht erreichbar**
 $D \rightarrow cD\}$ **nicht erzeugend / erreichbar**

Neue Grammatik:

$P_3 = \{$
 $S \rightarrow AB | abc | \varepsilon | aAAb | aAb | ab | bBB | bB | b,$ *Alle Variablen*
 $A \rightarrow aAAb | aAb | ab,$ *sind erreichbar*
 $B \rightarrow bBB | bB | b\}$ *und erzeugend!*

Schritt 3: Elimination von unnützen Variablen: Verfahren

Verfahren zur Bestimmung aller erzeugender Symbole:

1. Alle Terminalsymbole $a \in T$ sind erzeugend.
2. Ist $A \rightarrow X_1 \dots X_n \in P$ und alle X_i erzeugend, dann ist A erzeugend.

Verfahren zur Bestimmung aller erreichbaren Symbole:

1. S ist erreichbar.
2. Ist $A \rightarrow X_1 \dots X_n \in P$ und A erreichbar dann sind alle X_i erreichbar.

Verfahren zur Elimination von unnützen Variablen:

1. Lösche alle nicht erzeugenden Variablen (mit Produktionen).
2. Lösche alle nicht erreichbaren Variablen (mit Produktionen).

Die neue Grammatik erzeugt die selben Wörter!

Die neue Grammatik enthält keine unnützen Variablen mehr!

Schritt 4: Separieren von Terminalen und Variablen (X-Regel)

Auf der rechten Seite der Regeln müssen entweder
NUR mind. 2 Variablen oder NUR ein Terminal stehen.

$$\begin{aligned} P_3 = \{ \\ S &\rightarrow AB \mid abc \mid \epsilon \mid aAAb \mid aAb \mid ab \mid bBB \mid bB \mid b, \\ A &\rightarrow aAAb \mid aAb \mid ab, \\ B &\rightarrow bBB \mid bB \mid b \} \end{aligned}$$

Neue Grammatik:

$$\begin{aligned} P_4 = \{ \\ S &\rightarrow AB \mid X_a X_b X_c \mid \epsilon \mid X_a A A X_b \mid X_a A X_b \mid X_a X_b \mid X_b B B \mid X_b B \mid b, \\ A &\rightarrow X_a A A X_b \mid X_a A X_b \mid X_a X_b, \\ B &\rightarrow X_b B B \mid X_b B \mid b, \\ X_a &\rightarrow a, \\ X_b &\rightarrow b, \\ X_c &\rightarrow c \} \end{aligned}$$

*Auf den rechten Seiten ist
entweder genau ein Terminal oder
sind mindestens 2 Variablen!*

Schritt 4: Separieren von Terminalen und Variablen: Verfahren

Verfahren zur Separation von Terminalsymbolen und Variablen:

1. Für jedes Terminalsymbol $a \in T$ erzeuge neue Variable X_a .
2. In jeder Produktion $A \rightarrow \alpha$ mit $|\alpha| \geq 2$, in der $a \in T$ auftaucht, ersetze a durch X_a .
3. Ergänze Produktionen $X_a \rightarrow a$ für alle $a \in T$.
4. Lösche Variable X_a , falls nicht verwendet.

Die neue Grammatik erzeugt die selben Wörter!

Die neue Grammatik separiert die Terminale und Variablen!

Schritt 5: Aufspalten von $A \rightarrow \alpha$ mit $|\alpha| > 2$ (Y-Regel)

Auf der rechten Seite müssen alle Regeln mit Variablen
aus **genau zwei** Variablen bestehen.

$$P_4 = \{$$

$$S \rightarrow AB \mid X_a X_b X_c \mid \varepsilon \mid X_a A X_b \mid X_a A X_b \mid X_a X_b \mid X_b B B \mid X_b B \mid b,$$

$$A \rightarrow X_a A X_b \mid X_a A X_b \mid X_a X_b,$$

$$B \rightarrow X_b B B \mid X_b B \mid b,$$

$$X_a \rightarrow a, \quad X_b \rightarrow b, X_c \rightarrow c\}$$

Neue Grammatik:

$$P_5 = \{$$

$$S \rightarrow AB \mid X_a Y_1 \mid \varepsilon \mid X_a Y_2 \mid X_a Y_4 \mid X_a X_b \mid X_b Y_5 \mid X_b B \mid b,$$

$$Y_1 \rightarrow X_b X_c, Y_2 \rightarrow A Y_3, Y_3 \rightarrow A X_b, Y_4 \rightarrow A X_b, Y_5 \rightarrow B B,$$

$$A \rightarrow X_a Y_6 \mid X_a Y_8 \mid X_a X_b, \quad Y_6 \rightarrow A Y_7, Y_7 \rightarrow A X_b, Y_8 \rightarrow A X_b,$$

$$B \rightarrow X_b Y_9 \mid X_b B \mid b, \quad Y_9 \rightarrow B B,$$

$$X_a \rightarrow a, X_b \rightarrow b, X_c \rightarrow c\}$$

Alle Regeln sind in CNF !

Schritt 5: Aufspalten von $A \rightarrow \alpha$ mit $|\alpha| > 2$: Verfahren

Verfahren zur Aufspaltung von Produktionen $A \rightarrow \alpha$ mit $|\alpha| > 2$:

Ersetze jede Produktion

$$A \rightarrow X_1 \dots X_k$$

durch $k-1$ Produktionen mit $k-2$ neuen Variablen Y_1, \dots, Y_{k-2}

$$A \rightarrow X_1 Y_1, \quad Y_1 \rightarrow X_2 Y_2, \quad \dots \quad Y_{k-2} \rightarrow X_{k-1} X_k$$

Die neue Grammatik erzeugt die selben Wörter!

Die neue Grammatik ist in CNF!

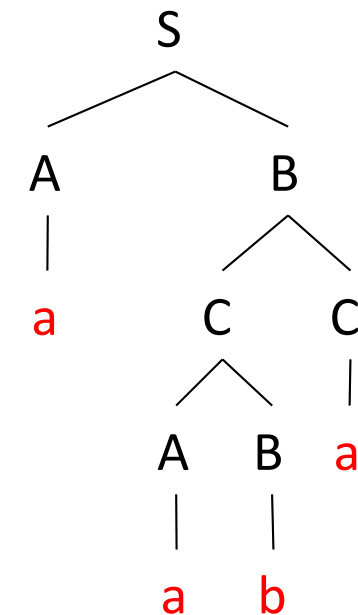
Ableitung bei CNF-Grammatiken: Binärbäume

Wie sehen die Ableitungsbäume von CNF-Grammatiken aus?

Sei $G=(V,T,P,S)$ mit

$P = \{ \begin{array}{l} S \rightarrow AB \mid BC, \\ A \rightarrow BA \mid a, \\ B \rightarrow CC \mid b, \\ C \rightarrow AB \mid a \end{array} \}$

Wie sieht ein Ableitungsbaum für $w = „aaba“$ aus?



Beobachtungen für einen Ableitungsbaum für ein Wort $w \neq \varepsilon$:

- Es gibt genau $|w|$ Blätter.
- Die Variablenknoten bilden einen Binärbaum (d.h. immer entweder 0 oder 2 Kinder als Variablen).
- Es gibt genau $2|w|-1$ Variablenknoten.

Jede Ableitung eines Wortes w hat die Länge $2|w|-1$.

Zusammenfassung

- Man kann jede kontextfreie Grammatik in eine Grammatik in **Chomsky-Normalform** umwandeln.
- CNF-Grammatiken erzeugen **binäre Ableitungsbäume** (max. 2 Kinder pro Knoten). Analysen und Beweise von Eigenschaften werden einfacher.