

7. Normal Forms for Context-free Grammars

Goal: Transform a given context-free grammar G into another grammar G' that

- has the same language but
- a simpler structure.

Motivation:

- Algorithms can then rely on the simpler structure which will make them faster (unless the blow-up in the conversion is prohibitive).
- It is easier for us to understand properties of context-free languages if the grammars have a simpler structure.

7.1 Getting Rid of ϵ -Productions, Unit Productions, and Useless Non-Terminals

Lemma:

For every CFG G , we can construct a CFG G'

- without ϵ -productions ($A \rightarrow \epsilon$) and
- without unit productions ($A \rightarrow B$)

satisfying $L(G') = L(G) \setminus \{\epsilon\}$.

Note: The result is a statement of existence of G' .
But it says more: We also know how to obtain G' .
If we not only know an algorithm exists but we can actually give the method, we say the algorithm is effective.

Proof: Let $G = (N, \Sigma, P, S)$.

We define $\hat{G} := (N, \Sigma, \hat{P}, S)$.

Here, \hat{P} is the smallest set containing P and closed under the following rules:

(a) If $A \rightarrow \alpha B \beta$ and $B \rightarrow \epsilon$ are in \hat{P} , then $A \rightarrow \alpha \beta$ is in \hat{P} .

(b) If $A \rightarrow B$ and $B \rightarrow \gamma$ are in \hat{P} , then $A \rightarrow \gamma$ is in \hat{P} .

Note that only finitely many productions are added. The reason is that the length of the right-hand sides is bounded by the length of the longest right-hand side in P .

Moreover, note that $L(\hat{P}) = L(P)$.

" \supseteq " Holds because $P \subseteq \hat{P}$.

" \subseteq " Holds because each new production can be simulated by two old productions that caused it to be added

(formally, this needs an induction).

Claim:
For every $w \in \Sigma^+$, every derivation $S \Rightarrow_{\hat{P}}^* w$ of minimal length neither uses ϵ -productions nor unit productions.

Proof:

Consider $w \neq \epsilon$ and a shortest derivation $S \Rightarrow_{\hat{P}}^* w$.

Towards a contradiction, assume that $B \rightarrow \epsilon$ is used at some point,

say $S \Rightarrow^* \alpha_1 B \alpha_2 \Rightarrow \alpha_1 \alpha_2 \Rightarrow^* w$.

• One of α_1 and α_2 is $\neq \epsilon$, otherwise $w = \epsilon$.

Thus, the B is not S but has been introduced at some point,

say with a production $A \rightarrow \beta_1 B \beta_2$:

$$S \Rightarrow^m \gamma_1 A \gamma_2 \Rightarrow \gamma_1 \beta_1 B \beta_2 \gamma_2 \Rightarrow^r \alpha_1 B \alpha_2 \Rightarrow \alpha_1 \alpha_2 \Rightarrow^k w.$$

• But by Rule (a), also $A \rightarrow \beta_1 \beta_2$ is in \hat{P} .

If we apply that production instead, we obtain

$$S \Rightarrow^m \gamma_1 A \gamma_2 \Rightarrow \gamma_1 \beta_1 \beta_2 \gamma_2 \Rightarrow^s \alpha_1 \alpha_2 \Rightarrow^k w.$$

• This contradicts the assumption that the given derivation was of minimal length. \square

• For unit productions, the reasoning is similar.

• Since we neither need ϵ -productions nor unit productions to generate $w \neq \epsilon$, we remove them from \tilde{G}

and obtain G' with $L(G') = L(G) \setminus \{\epsilon\}$. \square

Definition:

Let $G = (N, \Sigma, P, S)$. A non-terminal $X \in N$ is useful,

if $S \Rightarrow^* \alpha X \beta \Rightarrow^* w \in \Sigma^*$ for some $\alpha, \beta \in (N \cup \Sigma)^*$.

Otherwise, X is called useless.

Theorem:

For every CFG G , we can construct a CFG G' without ϵ -productions, without unit productions, and without useless non-terminals

so that $L(G') = L(G) \setminus \{\epsilon\}$.

Proof:

Consider the grammar G'' without ϵ -productions and without unit productions obtained with the previous lemma.

- With a fixed point computed backwards, we determine the non-terminals that can derive a terminal word.
- With a fixed point computed forwards, we determine the non-terminals that are reachable from S .
- Intersecting the sets and removing the remaining non-terminals together with the productions that use them yields the desired grammar. □

7.2 Chomsky Normal Form

Definition:

A CFG is in Chomsky normal form (CNF) if all productions take the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a.$$

Theorem:

For every CFG G , we can construct a CFG G' in CNF and with $L(G') = L(G) \setminus \{\epsilon\}$.

We base our construction on the previous theorem. As a result, all non-terminals will be useful, there are no ϵ -productions and no unit productions.

Proof:

- We first apply the previous theorem to turn G into G'' without useless non-terminals, ϵ -productions, and unit productions.

- We introduce a new non-terminal \bar{A}_a for every $a \in \Sigma$ that occurs in the right-hand side of a production.
We replace every occurrence of a in a right-hand side by \bar{A}_a .
We add the production $\bar{A}_a \rightarrow a$.
(Note that this step does not introduce useless non-terminals.)

- Now all productions are of the form
 $A \rightarrow a$ or $A \rightarrow B_1 \dots B_k$ with $k \geq 2$.

- For any production $\bar{A} \rightarrow B_1 \dots B_k$ with $k \geq 3$,
we introduce a new non-terminal C
and replace the production by

$$\bar{A} \rightarrow B_1 C \quad \text{and} \quad C \rightarrow B_2 \dots B_k.$$

We repeat this until all productions have length 2. \square

Example:

We derive a grammar in CNF for

$$\{a^n b^n \mid n \in \mathbb{N}\} \setminus \{\epsilon\} = \{a^n b^n \mid n \geq 1\}.$$

Starting from

$$S \rightarrow a S b \mid \epsilon \quad \text{for } \{a^n b^n \mid n \in \mathbb{N}\}.$$

We remove ϵ -productions as described in the first lemma.

This yields

$$S \rightarrow a S b \mid a b \quad \text{generating } \{a^n b^n \mid n \geq 1\}.$$

We add non-terminals \bar{A}, \bar{B} and replace the production by

$$S \rightarrow \bar{A} S \bar{B} \mid \bar{A} \bar{B} \quad \bar{A} \rightarrow a \quad \bar{B} \rightarrow b.$$

Finally, we add a non-terminal C and replace $S \rightarrow \bar{A} S \bar{B}$ by

$$S \rightarrow \bar{A} C \quad C \rightarrow S \bar{B}.$$

7.3 Greibach Normal Form

Definition:

A CFG G is in Greibach normal form (GNF).

if every production is of the form

$$A \rightarrow a B_1 \dots B_k, \quad k \geq 0.$$

Theorem:

For every CFG G , we can construct a CFG G' in GNF with $L(G') = L(G) \setminus \{\epsilon\}$.

We will base our construction on the Chomsky normal form. Therefore, we will obtain $k \leq 2$, no useless non-terminals, no ϵ -productions, no unit productions.

Proof:

Let $G = (N, \Sigma, P, S)$ be in CNF, $A \in N$, $a \in \Sigma$.

We define

$$R_{A,a} := \{ \beta \in N^* \mid A \Rightarrow_{SL}^* a \beta \}.$$

The strong left derivation relation $\Rightarrow_{SL}^* \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$ is defined by

$\alpha_1 \Rightarrow_{SL}^* \alpha_2$ if α_2 can be obtained from α_1

by a sequence of derivation steps in which productions are only applied to the leftmost symbol

• The set $R_{A,a}$ is a regular language over the alphabet N .
Indeed, take the left-linear grammar

$$\begin{aligned} & (\{ B' \mid B \in N \}, N, \{ B' \rightarrow C'D \mid B \rightarrow CDE \in P \}, A') \\ & \cup \{ B' \rightarrow \epsilon \mid B' \rightarrow a \} \end{aligned}$$

to generate it.

It contains a non-terminal B' for every non-terminal B in G .
 The former non-terminals N now form the terminal symbols.
 The productions always replace the leftmost non-terminal,
 thus simulating a sequence of left derivations.
 The process stops when the terminal symbol of interest a
 has been found.

// Note that it is a substantial conceptual twist
 to think about words over non-terminals
 as a regular language.

You should really think this through.

• It is not difficult to show that for every left-linear grammar $G_{L,a}$
 there is a strongly right-linear grammar $G_{R,a}$
 with the same language.

Strong here means the productions are of the form

$$X \rightarrow B.Y \quad \text{or} \quad X \rightarrow \epsilon.$$

where X, Y are non-terminals and B is a terminal symbol.

Let $T_{R,a}$ be the start symbol of $G_{R,a}$.

• We can assume that the non-terminals in all $G_{R,a}$ and in G
 are pairwise disjoint.

We form the grammar G_1 by adding to G
 all non-terminals and productions of all $G_{R,a}$.

The start symbol of G_1 is S .

The productions of G_1 are of the form

$$X \rightarrow \epsilon \quad X \rightarrow b \quad X \rightarrow B.Y, \quad X, B, Y \text{ non-terminals.}$$

Moreover, $L(G_1) = L(G)$ since none of the new non-terminals is reachable.

- We define the grammar G_2 from G_1 by replacing $X \rightarrow B.Y$ with the rules $X \rightarrow a.T_{B,a}.Y$ for all $a \in \Sigma$.

What is going on here?

The idea is to guess the leading terminal symbol a that will be generated in a derivation sequence from B .

The possible sequences of non-terminals following this a (that are generated from B)

are captured by $G_{B,a}$ with start symbol $T_{B,a}$.

It may seem restrictive to only consider sequences of non-terminals following a .

But we do the same replacement also for the rules in $G_{B,a}$.

Therefore, we capture all sentential forms derivable from B .

The productions in G_2 are of the form

$X \rightarrow \epsilon$, $X \rightarrow b$, and $X \rightarrow a.T_{B,a}.Y$.

- We get rid of ϵ -productions using the construction in the first lemma.

This construction does not introduce unit productions.

To see this, note that every non- ϵ -production has a terminal on the right-hand side.

- The resulting grammar G_3

- is in CNF and

- satisfies $L(G_3) = L(G) \setminus \{\epsilon\}$.

- It has no ϵ -productions and no unit productions.
- It has no useless non-terminals if we make sure to only consider $R_{A,a} \neq \emptyset$ (and modify $G_{A,a}$ accordingly).

□

Example:

Consider the grammar in CNF for balanced parentheses:

$$G: \begin{array}{l} S \rightarrow AB \mid AC \mid SS \\ C \rightarrow SB \end{array} \quad \begin{array}{l} A \rightarrow [\\ B \rightarrow] \end{array}$$

We compute the regular languages

$$(1) R_{S,\epsilon} = (B \cup C).S^*$$

$$(2) R_{C,\epsilon} = (B + C).S^*.B$$

$$(3) R_{A,\epsilon} = \{\epsilon\} = R_{B,]}.$$

All other languages are \emptyset .

Corresponding strongly right-linear grammars are

$$(1) T_{S,\epsilon} \rightarrow B.X \mid C.X \quad X \rightarrow S.X \mid \epsilon$$

$$(2) T_{C,\epsilon} \rightarrow B.Y \mid C.Y \quad Y \rightarrow S.Y \mid B.Z \quad Z \rightarrow \epsilon$$

$$(3) T_{A,\epsilon} \rightarrow \epsilon$$

$$T_{B,]} \rightarrow \epsilon.$$

Combining these grammars with G , we obtain:

$$S \rightarrow [\cdot T_{A,E} \cdot B \mid [\cdot T_{A,E} \cdot C \mid [\cdot T_{S,E} \cdot S$$

$$T_{S,E} \rightarrow] \cdot T_{B,] } \cdot X \mid [\cdot T_{C,E} \cdot X$$

$$T_{C,E} \rightarrow] \cdot T_{B,] } \cdot Y \mid [\cdot T_{C,E} \cdot Y$$

$$T_{A,E} \rightarrow \epsilon$$

$$T_{B,] } \rightarrow \epsilon$$

$$C \rightarrow [\cdot T_{S,E} \cdot B$$

$$X \rightarrow [\cdot T_{S,E} \cdot X \mid \epsilon$$

$$Y \rightarrow [\cdot T_{S,E} \cdot Y \mid] \cdot T_{B,] } \cdot Z$$

$$A \rightarrow [$$

$$B \rightarrow]$$

$$Z \rightarrow \epsilon.$$

Removing ϵ -productions yields:

$$S \rightarrow [\cdot B \mid [\cdot C \mid [\cdot T_{S,E} \cdot S$$

$$T_{S,E} \rightarrow] \mid] \cdot X \mid [\cdot T_{C,E} \mid [\cdot T_{C,E} \cdot X$$

$$T_{C,E} \rightarrow] \cdot Y \mid [\cdot T_{C,E} \cdot Y$$

$$C \rightarrow [\cdot T_{S,E} \cdot B$$

$$X \rightarrow [\cdot T_{S,E} \cdot X \mid [\cdot T_{S,E}$$

$$Y \rightarrow [\cdot T_{S,E} \cdot Y \mid]$$

$$A \rightarrow [$$

$$B \rightarrow].$$