

Chomsky Normal Form

- A standardized form for a CFG.
- One of the main goals is to be able to check if string w is generated by grammar G .
- What we want: Simple rules that are easy to check.
- What we don't want: loops, Empty derivations
Useless rules.

Consider $S \rightarrow OS1 \mid SOS1 \mid T$
 $T \rightarrow S \mid \epsilon$

1010000
.
10S10T00
10S1000
101000

Def 2.8: A context-free grammar is in Chomsky Normal Form if every rule is of the form

$$A \rightarrow BC$$

$$A \rightarrow a$$

$$\begin{aligned} A \rightarrow BC &\rightarrow CDC \\ &\rightarrow cDC \\ &\rightarrow c \hat{a} C \\ &\rightarrow cac \end{aligned}$$

where a is a terminal, A, B, C are variables and B and C are not the start variable. Also, we allow $S \rightarrow \epsilon$.

One of the advantages of the Chomsky Normal Form is the "predictability" in the derivation of a string. Any string of length n requires exactly $2n-1$ steps in its derivation.

↓
EXERCISE!

Also Chomsky NF results in an efficient algorithm for checking if w is generated by G .

Natural Question at this stage. What if there is no equivalent CFG that is in the Chomsky Normal Form?

Theorem 2.9: Any CFL is generated by a CFG in Chomsky Normal Form.

Proof: We provide a process that demonstrates how to convert any CFG into Chomsky NF.

Optional Step: Remove useless symbols and productions

$S \rightarrow AB \mid a$
 $A \rightarrow a$

→ B is a useless var.

1. Add new start variable S_0 .

$$S_0 \rightarrow S$$

This guarantees that S_0 does not appear in the R.H.S of any rule.

2. Remove ϵ rules.

Say there is the rule $A \rightarrow \epsilon$. Then modify rules with A in R.H.S.

If $R \rightarrow uAv$ was a rule, then replace it with

$$R \rightarrow uAv \mid uv$$

(here u, v are string of terminals and variables)

If $R \rightarrow A$ is a rule, then replace it with

$$R \rightarrow A | \epsilon$$

unless $R \rightarrow \epsilon$ was already removed.

3. Remove unit rules like $A \rightarrow B$.

If there is a rule $B \rightarrow u$, add $A \rightarrow u$ unless we removed that rule.

4. Restructure rules with long RHS.

Example:

$$S \rightarrow TST | aB$$

$$T \rightarrow B | S$$

$$B \rightarrow b | \epsilon$$

1. Add S_0 as new start variable.

$$S_0 \rightarrow S$$

$$S \rightarrow TST \mid aB$$

$$T \rightarrow B \mid S$$

$$B \rightarrow b \mid \epsilon$$

2. Remove $B \rightarrow \epsilon$

$$S_0 \rightarrow S$$

$$S \rightarrow TST \mid aB \mid a$$

$$T \rightarrow B \mid \underline{\epsilon} \mid S$$

$$B \rightarrow b$$

$T \rightarrow \epsilon$

$$S_0 \rightarrow S$$

$$S \rightarrow TST \mid ST \mid TS \mid \cancel{S} \mid aB \mid a$$

$$T \rightarrow B \mid S$$

$$B \rightarrow b$$

redundant rule
can be removed.

3. Unit Rules

$$S_0 \rightarrow S$$

$$S_0 \rightarrow TST \mid TS \mid ST \mid aB \mid a$$

$$S \rightarrow TST \mid TS \mid ST \mid aB \mid a$$

$$T \rightarrow B \mid S$$

$$B \rightarrow b$$

$$T \rightarrow B$$

$$T \rightarrow b \mid S$$

$$T \rightarrow S$$

$$S_0 \rightarrow TST \mid TS \mid ST \mid aB \mid a$$

$$S \rightarrow TST \mid TS \mid ST \mid aB \mid a$$

$$T \rightarrow TST \mid TS \mid ST \mid aB \mid a \mid b$$

$$B \rightarrow b$$

4. Restructure rules with more than one symbol in the RHS.

let $TS = \underline{Z}$, $a = \underline{A}$.

$$S_0 \rightarrow ZT \mid TS \mid ST \mid AB \mid a$$

$$S \rightarrow ZT \mid TS \mid ST \mid AB \mid a$$

$$T \rightarrow ZT \mid TS \mid ST \mid AB \mid a \mid b$$

$$Z \rightarrow TS$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Now the above CFG is in Chomsky Normal Form and by construction, we have ensured that it is equivalent to the original grammar.