

## Chomsky Normal Form (CNF)

A CFG  $G$  is said to be in CNF if it has no useless symbols and every production is of the form  $A \rightarrow BC$  or  $A \rightarrow a$ .

Given any CFL  $L$  s.t.  $L - \{\epsilon\}$  is nonempty we can construct a CFG  $G$  in CNF such that  $L(G) = L - \{\epsilon\}$ .

Method of constructing such a  $G$  :

- 1) Start with any CFG  $G_1$  such that  $L(G_1) = L$ .
- 2) Simplify  $G_1$  to get  $G_2$  such that  $L(G_2) = L - \{\epsilon\}$  and every production of  $G_2$  is of the form  $A \rightarrow X_1 X_2 \dots X_k$   $k > 1$  or  $A \rightarrow a$ . If any  $X_i = b$  then introduce a variable  $B_i$  and a production  $B_i \rightarrow b$ . Then we get an equivalent CFG  $G_3$  where every production is of the form  $A \rightarrow B_1 B_2 \dots B_k$   $k > 1$  or  $A \rightarrow a$ .
- 3) If  $k > 2$  for any production  $A \rightarrow B_1 B_2 \dots B_k$  replace this by  $A \rightarrow B_1 B_1'$ ,  $B_1' \rightarrow B_2 B_2'$ ,  $\dots$ ,  $B_{k-2}' \rightarrow B_{k-1} B_k$ . This way we get a Grammar  $G_4$  which is in CNF and  $L(G_4) = L(G_1) - \{\epsilon\} = L - \{\epsilon\}$ .

Example 1 :  $S \rightarrow a A a \mid b B b \mid \epsilon$

$A \rightarrow C \mid a$

$B \rightarrow C \mid b$

$C \rightarrow C D \mid \epsilon$

$D \rightarrow A \mid B \mid \epsilon$

First we have to simplify this which was given as a Homework in the last lesson.

The result is :  $S \rightarrow a A a \mid b B b \mid a a \mid b b$

$A \rightarrow a \mid C D \mid a b \mid b$

$B \rightarrow b \mid C D \mid a b \mid a$

$C \rightarrow C D \mid a b \mid a \mid b$

$$D \rightarrow CD \mid ab \mid a \mid b$$

Converting this to CNF is now easy. For the terminals a and b in the body of productions with more than one symbol, we introduce variables U and V and productions  $U \rightarrow a$  and  $V \rightarrow b$  getting the Grammar

$$S \rightarrow UAU \mid VBV \mid UU \mid VV$$

$$A \rightarrow a \mid CD \mid UV \mid b$$

$$B \rightarrow b \mid CD \mid UV \mid a$$

$$C \rightarrow CD \mid UV \mid a \mid b$$

$$D \rightarrow CD \mid UV \mid a \mid b$$

$$U \rightarrow a$$

$$V \rightarrow b$$

After that we merely have to take care of the productions  $S \rightarrow UAU \mid VBV$ . These yield  $S \rightarrow A_1U \mid B_1V$ ,  $A_1 \rightarrow UA$  and  $B_1 \rightarrow VB$ . The final Grammar in CNF is

$$S \rightarrow A_1U \mid B_1V \mid UU \mid VV$$

$$A_1 \rightarrow UA$$

$$B_1 \rightarrow VB$$

$$A \rightarrow a \mid CD \mid UV \mid b$$

$$B \rightarrow b \mid CD \mid UV \mid a$$

$$C \rightarrow CD \mid UV \mid a \mid b$$

$$D \rightarrow CD \mid UV \mid a \mid b$$

$$U \rightarrow a, V \rightarrow b$$

Example 2 :  $S \rightarrow A B C \mid B a B$

$A \rightarrow a A \mid B a C \mid a a a$

$B \rightarrow b B b \mid a \mid D$

$C \rightarrow C A \mid A C$

$D \rightarrow \epsilon$

This has to be first simplified. This was also given as a Homework. The result is

$S \rightarrow B a B \mid B a \mid a B \mid a$

$B \rightarrow b B b \mid b b \mid a$

For the terminals a and b in the bodies of productions of length > 1 we introduce the variables U, V and the productions  $U \rightarrow a$  and  $V \rightarrow b$  getting the Grammar

$S \rightarrow B U B \mid B U \mid U B \mid a$

$B \rightarrow V B V \mid V V \mid a, U \rightarrow a, V \rightarrow b$

Proceeding in the standard way now for bodies of certain productions having > 2 variables we get the final equivalent Grammar in CNF

$S \rightarrow B B_1 \mid B U \mid U B \mid a$

$B_1 \rightarrow U B$

$B \rightarrow V B_2 \mid V V \mid a$

$B_2 \rightarrow B V$

$U \rightarrow a$

$V \rightarrow b$