# Chomsky Normal Form

*We introduce Chomsky Normal Form, which is used to answer questions about context-free languages*

# Chomsky Normal Form

**Chomsky Normal Form.** *A grammar where every production is either of the form $A \to BC$ or $A \to$ c (where $A$, $B$, $C$ are arbitrary variables and c an arbitrary symbol).*

Example:

$$S \to AS \mid \text{a}$$
$$A \to SA \mid \text{b}$$

(If language contains $\varepsilon$, then we allow $S \to \varepsilon$ where $S$ is start symbol, and forbid $S$ on RHS.)

## Why Chomsky Normal Form?

The key advantage is that in Chomsky Normal Form, every derivation of a string of $n$ letters has exactly $2n - 1$ steps.

Thus: *one can determine if a string is in the language by exhaustive search of all derivations.*

## Conversion

The conversion to Chomsky Normal Form has four main steps:

1. Get rid of all $\varepsilon$ productions.

2. Get rid of all productions where RHS is one variable.

3. Replace every production that is too long by shorter productions.

4. Move all terminals to productions where RHS is one terminal.

# 1) Eliminate $\varepsilon$ Productions

Determine the nullable variables (those that generate $\varepsilon$) (algorithm given earlier).

Go through all productions, and for each, omit every possible subset of nullable variables.

For example, if $P \rightarrow A\mathrm{x}B$ with both $A$ and $B$ nullable, add productions $P \rightarrow \mathrm{x}B \mid A\mathrm{x} \mid \mathrm{x}$.

After this, delete all productions with empty RHS.

## 2) Eliminate Variable Unit Productions

A unit production is where RHS has only one symbol.

Consider production $A \to B$. Then for every production $B \to \alpha$, add the production $A \to \alpha$. Repeat until done (but don't re-create a unit production already deleted).

# 3) Replace Long Productions by Shorter Ones

For example, if have production $A \rightarrow BCD$, then replace it with $A \rightarrow BE$ and $E \rightarrow CD$.

(In theory this introduces many new variables, but one can re-use variables if careful.)

## 4) Move Terminals to Unit Productions

For every terminal on the right of a non-unit production, add a substitute variable.

For example, replace production $A \to \textsf{b}C$ with productions $A \to BC$ and $B \to \textsf{b}$.

Consider the CFG:

$S \rightarrow \mathsf{a}X\mathsf{b}X$

$X \rightarrow \mathsf{a}Y \mid \mathsf{b}Y \mid \varepsilon$

$Y \rightarrow X \mid \mathsf{c}$

The variable $X$ is nullable; and so therefore is $Y$. After elimination of $\varepsilon$, we obtain:

$S \rightarrow \mathsf{a}X\mathsf{b}X \mid \mathsf{ab}X \mid \mathsf{a}X\mathsf{b} \mid \mathsf{ab}$

$X \rightarrow \mathsf{a}Y \mid \mathsf{b}Y \mid \mathsf{a} \mid \mathsf{b}$

$Y \rightarrow X \mid \mathsf{c}$

After elimination of the unit production $Y \to X$, we obtain:

$S \to \mathsf{a}X\mathsf{b}X \mid \mathsf{ab}X \mid \mathsf{a}X\mathsf{b} \mid \mathsf{ab}$

$X \to \mathsf{a}Y \mid \mathsf{b}Y \mid \mathsf{a} \mid \mathsf{b}$

$Y \to \mathsf{a}Y \mid \mathsf{b}Y \mid \mathsf{a} \mid \mathsf{b} \mid \mathsf{c}$

Now, break up the RHSs of $S$; and replace a by $A$, b by $B$ and c by $C$ wherever not units:

$S \rightarrow EF \mid AF \mid EB \mid AB$

$X \rightarrow AY \mid BY \mid a \mid b$

$Y \rightarrow AY \mid BY \mid a \mid b \mid c$

$E \rightarrow AX$

$F \rightarrow BX$

$A \rightarrow a$

$B \rightarrow b$

$C \rightarrow c$

Convert the following CFG into Chomsky Normal Form:

$$S \rightarrow A b A$$
$$A \rightarrow A a \mid \varepsilon$$

After the first step, one has:

$$S \rightarrow AbA \mid bA \mid Ab \mid b$$
$$A \rightarrow Aa \mid a$$

The second step does not apply. After the third step, one has:

$$S \rightarrow TA \mid bA \mid Ab \mid b$$
$$A \rightarrow Aa \mid a$$
$$T \rightarrow Ab$$

And finally, one has:

$$S \to TA \mid BA \mid AB \mid \text{b}$$
$$A \to AC \mid \text{a}$$
$$T \to AB$$
$$B \to \text{b}$$
$$C \to \text{a}$$

## Summary

There are special forms for CFGs such as Chomsky Normal Form, where every production has the form $A \rightarrow BC$ or $A \rightarrow c$. The algorithm to convert to this form involves (1) determining all nullable variables and getting rid of all $\varepsilon$-productions, (2) getting rid of all variable unit productions, (3) breaking up long productions, and (4) moving terminals to unit productions.