# Chomsky Normal Form

* Helpful to have a simplified form for $G$.
* Need to be efficient to check if $w \in L(G)$.
* Should not have loops (in derivation).
* Should not have empty derivations, or useless rules.
* Should have simple rules that are easy to check.

**Def 2.8:** A context-free grammar is in _Chomsky Normal Form_ if every rule is of the form

$$A \to BC$$
$$A \to a$$

where $a$ is a terminal, $A, B, C$ are variables and $B$ and $C$ should not be the start variable. Also, we allow $S \to \varepsilon$.

**Theorem 2.9** : Any CFL is generated by a CFG in

## Chomsky Normal Form.

**Proof:** We convert any CFG into CNF.

**Optional Step:** Remove useless symbols & productions.

1. Add new start variable $S_0$.

$$S_0 \rightarrow S$$

This guarantees that $S_0$ does not appear in the RHS of any rule.

2. Remove $\varepsilon$ rules.
Say $A \rightarrow \varepsilon$. Then modify rules with $A$ in RHS. If $R \rightarrow uAv$ was a rule, then

$$R \rightarrow uAv \mid uv.$$

$$R \rightarrow A \quad \text{becomes} \quad R \rightarrow A \mid \varepsilon$$

unless $R \rightarrow \varepsilon$ is already removed.

3. Remove unit rules like $A \rightarrow B$.
When $B \rightarrow u$ appears, add $A \rightarrow u$ unless it was already removed.

4. Restructure long RHS.

## 4. Restructure long RHS.

Example :

$$S \to TST \mid aB$$
$$T \to B \mid S$$
$$B \to b \mid \varepsilon$$

**(1) Add** $S_0 \to S$.

$$S_0 \to S$$
$$S \to TST \mid aB$$
$$T \to B \mid S$$
$$B \to b \mid \varepsilon$$

**(2) Remove** $B \to \varepsilon$

$$S_0 \to S$$
$$S \to TST \mid aB \mid a$$
$$T \to B \mid \varepsilon \mid S$$
$$B \to b$$

$T \to \varepsilon$

$$S_0 \to S$$
$$S \to TST \mid TS \mid ST \mid \cancel{S} \mid aB \mid a$$
$$T \to B \mid S$$
$$B \to b$$

**(3) Unit Rules**

$$\phantom{S} \to TST \mid TS \mid ST \mid aB \mid a$$

(3) Unit rules

$S_0 \rightarrow S$

$S_0 \rightarrow TST \mid TS \mid ST \mid aB \mid a$
$S \rightarrow TST \mid TS \mid ST \mid aB \mid a$
$T \rightarrow B \mid S$
$B \rightarrow b$

---

$T \rightarrow B$
$T \rightarrow S$

$S_0 \rightarrow TST \mid TS \mid ST \mid aB \mid a$
$S \rightarrow TST \mid TS \mid ST \mid aB \mid a$
$B \rightarrow b$
$T \rightarrow b \mid S$

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

---

(4) Restructure rules with more than one symbol in RHS.

Let $TS = Z$, $a = A$.

$S_0 \rightarrow ZT \mid TS \mid ST \mid AB \mid a$
$S \rightarrow ZT \mid TS \mid ST \mid AB \mid a$
$T \rightarrow ZT \mid TS \mid ST \mid AB \mid a$
$Z \rightarrow TS$

$$A \longrightarrow a$$
$$B \longrightarrow b.$$

The main advantage of CNF is that there is "predictability" in the derivation of a string. Any string of length n requires exactly 2n-1 steps for derivation. Why?

Chomsky Normal Form results in an algorithm for checking if w is generated by G.

## CYK algorithm

John Cocke    1970
Daniel Younger 1967
Tado  Kasami   1965

One naive approach is to try out all derivations with 2n-1 steps. This is not time efficient. The CYK algorithm is a dynamic programming algorithm.

Use subproblems to solve larger problems.

Use subproblems to solve larger problems.

$\Theta(n^3)$ time.

Let $v = a_1 a_2 \ldots a_n$ where each $a_i \in \Sigma$.

$w_{ij} = a_i a_{i+1} \ldots a_j$ for all $1 \leq i \leq j \leq n$.

CYK builds $T_{ij}$ for all $1 \leq i \leq j \leq n$, such that $\underline{T_{ij} = \{ A \mid A \overset{*}{\Rightarrow} w_{ij} \}}$, leading up to $T_{1n}$. Finally, check if $\underline{S \in T_{1n}}$. That is, $\underline{S \overset{*}{\Rightarrow} w_{1n} = \omega}$.

$$T_{11} \quad T_{22} \quad T_{33} \quad \ldots \ldots \quad T_{n,n}$$
$$T_{12} \quad T_{23} \quad T_{34} \ldots \quad T_{n-1,n}$$
$$T_{13} \quad T_{24} \ldots T_{n-2,n}$$
$$\vdots$$
$$T_{1,n-1} \quad T_{2,n}$$
$$T_{1,n}$$

Order of computing $T_{i,j}$'s.

For $k=0$ to $n-1$, compute all $T_{(i, i+k)}$

For $k = 0$ to $n-1$, compute all $'(i, i+k)$

$k = 0$. $T_{i,i} = \{A \mid A \overset{*}{\Rightarrow} w_{i,i} = a_i\}$.

All the rules $A \to a_i$

For $k > 0$, $A \in T_{i, i+k}$ iff $B \in T_{i, i+j}$ and
$C \in T_{i+j+1, i+k}$ and $A \to BC$ is a rule.

---

If $w = \varepsilon$, accept if $S \to \varepsilon$ is a rule.

For $i = 1$ to $n$
$\qquad A \in T_{i,i} \iff A \to a_i$ is a rule.

For $k = 1$ to $n-1$.
$\qquad$ For $i = 1$ to $n-k$
$\qquad\qquad$ For $j = 0$ to $k-1$

$\qquad\qquad\qquad$ Check all rules $R \to AB$.
$\qquad\qquad\qquad$ If $T_{i, i+j}$ contains $A$, and
$\qquad\qquad\qquad\qquad T_{i+j+1, i+k}$ contains $B$
$\qquad\qquad\qquad\qquad$ then $T_{i, i+k} = T_{i, i+k} \cup \{R\}$.

If $S \in T_{1,n}$, say $w \in L(G)$.

If $SET_{1,n}$, say $\omega \in L(G)$.

Else $\omega \notin L(G)$.

Running time: $O(n^3 r)$ where $r$ is no. of rules.

$r$ is considered to be a constant

Correctness is evident from the algorithm.

## Example:

$S \rightarrow AB \mid BC$

$A \rightarrow BA \mid a$

$B \rightarrow CC \mid b$

$C \rightarrow AB \mid a$

$\omega = baaba$

| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 5 | SAC | SAC | B | SA | AC |
| 4 | — | B | SC | B | |
| 3 | — | B | AC | | |
| 2 | SA | AC | | | |
| 1 | B | | | | |

Exercise: Verify this matrix.