# Chomsky Normal Form. (CNF)

$G = (N, \Sigma, P, S)$

$S \rightarrow SS \mid 0 \mid 1 \mid \epsilon$

$S \xrightarrow[G]{1} SS \xrightarrow[G]{1} SSS \xrightarrow[G]{1} SSSS \xrightarrow[G]{1} SSS \xrightarrow[G]{1} SS \xrightarrow[G]{1} 0S \xrightarrow[G]{1} 01$

$G = (N, \Sigma, P, S)$ is in CNF if all productions

$A \rightarrow BC \qquad A \rightarrow a \qquad A, B, C \in N, \ a \in \Sigma.$

Example. $S \rightarrow [S] \mid SS \mid \epsilon \qquad \color{red}{\rightarrow G_1 \text{ not in CNF}}$

$S \rightarrow AB \mid AC \mid SS \qquad C \rightarrow SB \qquad A \rightarrow [ \qquad B \rightarrow ]$

$\color{red}{\hookrightarrow G_2}$

Claim. $L(G_1) = L(G_2)$

One step progress. $\begin{cases} - \text{\# of Nonterminals increase by 1} \\ - \text{\# of terminals increase by 1.} \end{cases}$

Theorem. For any CFG $G$, there is a CFG $G'$ in CNF such that $L(G') = L(G) - \{\epsilon\}$.

**Lemma 1.** For any CFG $G = (N, \Sigma, P, S)$ there is a CFG $G'$ with no $\varepsilon$-productions or unit-productions such that $L(G') = L(G) - \{\varepsilon\}$

**Proof.** Let $\hat{P}$ be the smallest set of productions containing P and closed under the rules:

(a) if $A \to \alpha B \beta$ and $B \to \varepsilon$ are in $\hat{P}$ then $A \to \alpha \beta \in \hat{P}$

(b) if $A \to B$ and $B \to \gamma$ are in $\hat{P}$ then $A \to \gamma \in \hat{P}$.

Note: $\hat{P}$ is finite.
{ Finitely many new production rules are added }
{ Each new RHS is a substring of an old RHS. }

$\hat{G} = (N, \Sigma, \hat{P}, S)$

We have $L(G) \subseteq L(\hat{G})$ Since $P \subseteq \hat{P}$

$L(G) = L(\hat{G})$ - Each new production was included because of rule (a) or (b) - can be simulated in 2 steps by two productions that caused it to be included.

**Claim 2.** For any non-null $x \in \Sigma^*$, any derivation $S \xrightarrow{*}_{\hat{G}} x$ of minimum length does not use $\epsilon$-or unit productions.

**Proof.** Let $x \neq \epsilon$. Let $S \xrightarrow{*}_{\hat{G}} x$ be the minimum length derivation.

Suppose an $\epsilon$-production $B \rightarrow \epsilon$ is used at some point

$$S \xrightarrow{*}_{\hat{G}} \gamma B \delta \xrightarrow{1}_{\hat{G}} \gamma \delta \xrightarrow{*}_{\hat{G}} x.$$

At least one of $\gamma$ or $\delta$ is non-null $\Longrightarrow$ $B$ was introduced from a production of the form $A \rightarrow \alpha B \beta$.

$$S \xrightarrow{m}_{\hat{G}} \eta A \theta \xrightarrow{1}_{\hat{G}} \eta \alpha B \beta \theta \xrightarrow{n}_{\hat{G}} \gamma B \delta \xrightarrow{1}_{\hat{G}} \gamma \delta \xrightarrow{k}_{\hat{G}} x$$

$$\text{for } m, n, k \geq 0$$

By rule (a) $A \rightarrow \alpha \beta \in \hat{P}$.

But then we have a strictly shorter derivation of $x$

$$S \xrightarrow{m}_{\hat{G}} \eta A \theta \xrightarrow{1}_{\hat{G}} \eta \alpha \beta \theta \xrightarrow{n}_{\hat{G}} \gamma \delta \xrightarrow{k}_{\hat{G}} x.$$

This gives a contradiction.

# Unit Productions

Let $x \neq \epsilon$. Consider a derivation $S \overset{*}{\underset{\hat{G}}{\Rightarrow}} x$ of minimum length.

Suppose a unit production $A \to B$ is used at some point

$$S \overset{*}{\underset{\hat{G}}{\Rightarrow}} \alpha A \beta \overset{1}{\underset{\hat{G}}{\Rightarrow}} \alpha B \beta \overset{*}{\underset{\hat{G}}{\Rightarrow}} x.$$

$B$ must be removed later by applying a production $B \to \gamma$.

$$S \overset{m}{\underset{\hat{G}}{\Rightarrow}} \alpha A \beta \overset{1}{\underset{\hat{G}}{\Rightarrow}} \alpha B \beta \overset{n}{\underset{\hat{G}}{\Rightarrow}} \eta B \theta \overset{1}{\underset{\hat{G}}{\Rightarrow}} \eta \gamma \theta \overset{k}{\underset{\hat{G}}{\Rightarrow}} x.$$

By rule (b), $A \to \gamma \in \hat{P}$.

But then, there is a shorter derivation of $x$

$$S \overset{m}{\underset{\hat{G}}{\Rightarrow}} \alpha A \beta \overset{1}{\underset{\hat{G}}{\Rightarrow}} \alpha \gamma \beta \overset{n}{\underset{\hat{G}}{\Rightarrow}} \eta \gamma \theta \overset{k}{\underset{\hat{G}}{\Rightarrow}} x.$$

This is a contradiction.

Claim 2 implies we can remove the $\epsilon$-productions and unit productions from $\hat{P}$ without changing the language.

# Chomsky Normal Form.

By Lemma 1, $L(G) = L(\hat{G})$ and $\hat{P}$ does not have $\epsilon$-productions or unit productions.

For each terminal $a \in \Sigma$ introduce a new nonterminal $A_a$ and add the production rule $A_a \to a$.

Replace all occurences of $a$ on the RHS of old productions (except productions of the form $B \to a$) with $A_a$. Then all productions are of the form:

$$A \to a \quad \text{or} \quad A \to B_1 B_2 \cdots B_k \quad k \geq 2$$

$\underbrace{B_1 B_2 \cdots B_k}_{\text{nonterminals.}}$

For any production of the form $A \to B_1 B_2 \cdots B_k$ with $k \geq 3$, introduce a new nonterminal $C$ and replace with

$$A \to B_1 C \quad \text{and} \quad C \to B_2 \cdots B_k.$$

Repeat until all RHS of all productions are of length atmost 2.