Chomsky Normal Form

Chomsky Normal Form

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

 $A \rightarrow BC$ $A \rightarrow a$

where **a** is any terminal and **A**, **B**, and **C** are any variables — except that **B** and **C** may not be the start variable.

In addition, we permit the rule $S \to \epsilon$ where S is the start variable if the language of the grammar contains ϵ .

Chomsky Normal Form

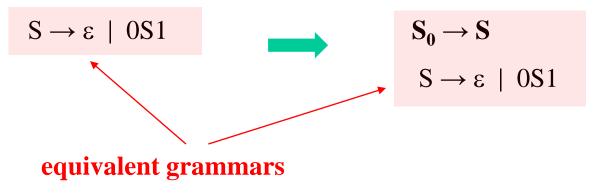
Theorem: Every (non-empty) context free language can be generated by a contextfree grammar in Chomsky normal form.

- In order to obtain **an equivalent grammar in Chomsky normal form** for any given CFG G, we will have the following conversion steps:
 - 1. Add a new start variable S_0 and a new production rule $S_0 \rightarrow S$ where S is the original start variable of G.
 - 2. Eliminate ε -productions (productions of the form $A \rightarrow \varepsilon$). After this conversion step, only one ε -production ($S_0 \rightarrow \varepsilon$) if the language of G contains ε .
 - 3. Eliminate unit productions (productions of the form $A \rightarrow B$ where A and B are variables).
 - **4.** Eliminate useless symbols. Useless symbols do not appear in any derivation of a terminal string from the start symbol.
 - **5.** Convert the remaining rules into Chomsky normal form adding new variables and rules.

Add A New Start Variable

- We add a new start variable S_0 and the rule $S_0 \rightarrow S$, where S was the original start variable.
- This change guarantees that the start variable does not occur on the right-hand side of a rule.
- The new grammar is equivalent to the original grammar (i.e *they generate same language*).

Example:



Eliminate ε-productions nullable variables

- In order to remove **ε-productions**, first we will determine **nullable variables**.
- A variable A is said to **nullable** if $A \stackrel{\circ}{\Rightarrow} \varepsilon$.
- We can compute nullable(G), the set of all nullable symbols of a CFG G=(V,T,P,S) as follows:

Basis:

nullable(G) = {A : $A \rightarrow \varepsilon \in P$ }

Induction:

If $\{C_1,...,C_k\} \subseteq nullable(G)$ and $A \rightarrow C_1,...,C_k \in P$, then $nullable(G) = nullable(G) \cup \{A\}$

Eliminate ε-productions

nullable variables: Example

• A CFG G_1 $S_0 \rightarrow S$ $S \rightarrow \varepsilon \mid 0S1$

→ nullable(G_1) = {S, S₀}

- A CFG G_2 $S_0 \rightarrow S$ $S \rightarrow AB$
 - $A \rightarrow aAA \mid \epsilon$
 - $B \rightarrow bBB \mid \epsilon$

→ nullable(G_2) = {A, B, S, S₀}

Eliminate ε-productions

Steps for ε-production elimination for CFG G:

- 1. Find nullable(G), the set of all *nullable* symbols of G.
- 2. Generate new rules from a rule R by eliminating *nullable* variables from its right-side, if *nullable* variables appears on its right-side.
 - The number of new rules depends on the number of *nullable* variables on the right-side. If there are k *nullable* variables, we have to generate 2^k-1 new rules.
 - Generated new rule is added if it is not already among the rules.
 - If $R \rightarrow \alpha A\beta$ is a rule and A is the only *nullable* variable on $\alpha A\beta$, generate and add the new rule $R \rightarrow \alpha\beta$.
 - If $R \rightarrow \alpha A \beta B \gamma$ is a rule and A and B are only *nullable* variables on $\alpha A \beta B \gamma$, generate and add the new rules $R \rightarrow \alpha \beta B \gamma$, $R \rightarrow \alpha A \beta \gamma$ and $R \rightarrow \alpha \beta \gamma$.
 - ...
- 3. Remove all ε -productions $A \rightarrow \varepsilon$ (except $S_0 \rightarrow \varepsilon$) from the rules.

The new grammar that is obtained by **eliminating** ε -productions is equivalent to the original grammar (i.e *they generate same language*).

Eliminate ε-productions: *Example*

$$S_0 \to S \qquad \qquad S \to \epsilon \ | \ 0S1$$

- nullable(G) = {S, S_0 }
- Since S is nullable,
 - generate $S_0 \rightarrow \epsilon$ from $S_0 \rightarrow S$
- Since S is nullable,
 - generate $S \rightarrow 01$ from $S \rightarrow 0S1$
- After all generations, we have the following rules:

 $S_0 \rightarrow \varepsilon \mid S$ $S \rightarrow \varepsilon \mid 01 \mid 0S1$

• Remove all ε -productions except $S_0 \rightarrow \varepsilon$, our final grammar is:

Eliminate ε-productions: *Example*

from $S_0 \rightarrow S$

from $S \rightarrow AB$

from $A \rightarrow aAA$

from $B \rightarrow bBB$

$$S_0 \to S \qquad S \to AB \qquad A \to aAA \mid \underset{X}{\epsilon} \qquad B \to bBB \mid \underset{X}{\epsilon}$$

nullable(G) = {A, B, S, S_0 }

Generate new rules:

 $S_0 \rightarrow \varepsilon$ $S \rightarrow A \qquad S \rightarrow B$ $A \rightarrow aA$ $A X_{aA}$

$$A \rightarrow aA \qquad A \stackrel{X}{\rightarrow} aA \qquad A \rightarrow a \\ B \rightarrow bB \qquad B \stackrel{X}{\rightarrow} bB \qquad B \rightarrow b$$

IJ

• Remove
$$\varepsilon$$
-productions
 $S_0 \rightarrow S \mid \varepsilon$
 $S \rightarrow AB \mid A \mid B$
 $A \rightarrow aAA \mid aA \mid a$
 $B \rightarrow bBB \mid bB \mid b$
• equivalent grammars

S ₩ε

Eliminate Unit Productions

- $A \rightarrow B$ is a unit production, whenever A and B are variables.
- Unit productions can be eliminated from a grammar to obtain a grammar without unit productions.
 - The resulting grammar that is obtained by eliminating unit productions will be equivalent to the original grammar.
- We will remove unit productions one by one from the grammar.
- Remove a unit production $A \rightarrow B$ from the grammar.
 - Then, whenever a rule $B \rightarrow u$ appears, we add the rule $A \rightarrow u$ unless this was a unit rule previously removed.
- We repeat these steps until we eliminate all unit rules.

 $S_0 \rightarrow \varepsilon \mid S$ $S \rightarrow 01 \mid 0S1$

- Unit productions: $\{ S_0 \rightarrow S \}$
- Remove $S_0 \rightarrow S$,
 - Add $S_0 \rightarrow 01$ and $S_0 \rightarrow 0S1$
- The resulting grammar after eliminating unit productions.

 $S_0 \rightarrow S \mid \varepsilon$ $S \rightarrow AB \mid A \mid B$ $A \rightarrow aAA \mid aA \mid a$ $\mathbf{B} \rightarrow \mathbf{b}\mathbf{B}\mathbf{B} \mid \mathbf{b}\mathbf{B} \mid \mathbf{b}$ - Unit productions: { $S_0 \rightarrow S, S \rightarrow A, S \rightarrow B$ } - Remove $S \rightarrow B$, add $S \rightarrow bBB \mid bB \mid b$ - Remove $S \rightarrow A$, add $S \rightarrow aAA \mid aA \mid a$ - Remove $S_0 \rightarrow S$, add $S_0 \rightarrow AB \mid aAA \mid a \mid bBB \mid bB \mid b$ The resulting grammar after eliminating unit productions. $S_0 \rightarrow \epsilon \mid AB \mid aAA \mid aA \mid a \mid bBB \mid bB \mid b$ $S \rightarrow AB \mid aAA \mid aA \mid a \mid bBB \mid bB \mid b$ equivalent grammars $A \rightarrow aAA \mid aA \mid a$ $B \rightarrow bBB \mid bB \mid b$

BBM401 Automata Theory and Formal Languages

 $E \rightarrow E+T | T$ $T \rightarrow T^*F | F$ $F \rightarrow G^{A}F | G$ $G \rightarrow id | (E)$

- Unit productions: { $E \rightarrow T, T \rightarrow F, F \rightarrow G$ }
- Remove $F \rightarrow G$, add $F \rightarrow id \mid (E)$
- Remove $T \rightarrow F$, add $T \rightarrow G^{A}F \mid id \mid (E)$
- Remove $E \rightarrow T$, add $E \rightarrow T^*F \mid G^{\wedge}F \mid id \mid (E)$
- The resulting grammar after eliminating unit productions.

```
\begin{array}{l} E \rightarrow E + T \mid T^*F \mid G^*F \mid id \mid (E) \\ T \rightarrow T^*F \mid G^*F \mid id \mid (E) \\ F \rightarrow G^*F \mid id \mid (E) \\ G \rightarrow id \mid (E) \end{array}
```

equivalent grammars

Eliminating unit productions in different order do not change the result.

 $E \rightarrow E+T \mid T$ $T \rightarrow T^*F \mid F$ $\mathbf{F} \rightarrow \mathbf{G}^{\mathbf{F}} \mid \mathbf{G}$ $G \rightarrow id \mid (E)$ - Unit productions: { $E \rightarrow T, T \rightarrow F, F \rightarrow G$ } - Remove $E \rightarrow T$, add $E \rightarrow T^*F \mid F$ - Remove $\mathbf{T} \rightarrow \mathbf{F}$, add $\mathbf{T} \rightarrow \mathbf{G}^{\mathsf{A}}\mathbf{F} \mid \mathbf{G}$ - Remove $\mathbf{F} \rightarrow \mathbf{G}$, add $\mathbf{F} \rightarrow \mathbf{id} \mid (\mathbf{E})$ - Remove newly introduced unit productions - Remove $E \rightarrow F$, add $E \rightarrow G^{A}F \mid id \mid (E)$ - Remove $T \rightarrow G$, add $T \rightarrow id \mid (E)$ The resulting grammar after eliminating unit productions. $E \rightarrow E+T \mid T^*F \mid G^{A}F \mid id \mid (E)$ $T \rightarrow T^*F \mid G^{\uparrow}F \mid id \mid (E)$ equivalent grammars $F \rightarrow G^{A}F \mid id \mid (E)$

 $G \rightarrow id \mid (E)$

Eliminate Useless Symbols

• A symbol X is useful for a grammar G=(V,T,P,S), if there is a derivation $S \stackrel{*}{\Rightarrow} \alpha X \beta \stackrel{*}{\Rightarrow} w$

for a terminal string w.

- Symbols that are not useful are called useless.
- A symbol X is generating if $X \stackrel{*}{\Rightarrow} w$ for some string $w \in T^*$.
- A symbol X is reachable if $S \stackrel{*}{\Rightarrow} \alpha X\beta$ for some $\{\alpha,\beta\} \subseteq (V \cup T)^*$.
- If we eliminate non-generating symbols first, and then non-reachable symbols, we will be left with only useful symbols.
 - The grammar that is obtained by eliminating useless symbols will be equivalent to the original grammar.

Eliminate Useless Symbols computing generating symbols

- For a grammar G = (V,T,P,S), the generating symbols generating(G) are computed by the following closure algorithm:
- **Basis:** generating(G) = T

Induction:

If $X \rightarrow \epsilon \in P$ or $X \rightarrow A_1 \dots A_n \in P$ where $\{A_1, \dots, A_n\} \subseteq generating(G)$ then generating(G) = generating(G) $\cup \{X\}$

Example:

- Let G be $S \rightarrow AB|a, A \rightarrow b$
- Initially, generating(G) = $\{a,b\}$
- A will be in generating(G) because of $A \rightarrow b$
- S will be in generating(G) because of $S \rightarrow a$
- Thus, generating(G)={a,b,A,S} and non-generating symbols are {B}

Eliminate Useless Symbols computing reachable symbols

- For a grammar G = (V,T,P,S), the reachable symbols **reachable**(G) are computed by the following closure algorithm:
- **Basis:** reachable(G) = $\{S\}$

Induction:

If $X \in reachable(G)$ and $X \rightarrow \alpha \in P$ then

add all symbols in α to reachable(G).

Example:

- Let G be $S \rightarrow AB|a, A \rightarrow b, C \rightarrow a$
- Initially, reachable(G) = $\{S\}$
- A and B will be in reachable(G) because of $S \rightarrow AB$
- a will be in reachable(G) because of $S \rightarrow a$
- b will be in reachable(G) because of $A \rightarrow b$
- Thus, reachable(G)={S,A,B,a,b} and non-reachable symbols are {C}

Eliminate Useless Symbols

Steps to eliminate useless symbols from G = (V,T,P,S) :

- 1. Compute generating(G).
- 2. Remove all productions containing at least one non-generating symbol in order to create a new grammar G_1 (a grammar without non-generating symbols).
 - Remove a production if a non-generating symbol appears in that production (on its rightside or its left-side)
- 3. Compute reachable(G_1).
- 4. Remove all productions containing at least one non-reachable symbol in order to create a new grammar G_2 without useless symbols (a grammar without non-reachable symbols and non-generating symbols).

The new grammar G_2 (a grammar without useless symbols) will be equivalent to the original grammar G.

Eliminate Useless Symbols: *Example*

- G: $S \rightarrow AB \mid a, A \rightarrow b$
- Compute generating(G):
 - generating(G)= $\{a,b,A,S\}$ and **non-generating symbols are \{B\}.**
- Remove productions containing non-generating symbols:
 - Remove $S \rightarrow AB$ because it contains B.
 - Thus, following G_1 is a grammar without non-generating symbols.
 - G_1 is $S \rightarrow a, A \rightarrow b$
- Compute reachable(G₁):
 - reachable(G_1)={S,a} and **non-reachable symbols are {A,b}.**
- Remove productions containing non-reachable symbols:
 - Remove $\mathbf{A} \rightarrow \mathbf{b}$ because it contains A (and/or b).
- Grammar G₂ without useless symbols (non-generating and non-reachable symbols):
 G₂: S→a

Convert the remaining rules into CNF

- Steps to obtain an equivalent grammar in Chomsky Normal Form:
- 1. Add a new start variable S_0^{-1}
- 2. Eliminate ε-productions.
- 3. Eliminate unit productions.
- 4. Eliminate useless symbols.
 - i. Eliminate non-generating symbols.
 - ii. Eliminate non-reachable symbols.
- **5.** Convert the remaining rules into CNF:

Now, to obtain a grammar in CNF, we want every rule to be the form

 $A \rightarrow BC$ $A \rightarrow a$

- i. Arrange that all bodies of length 2 or more consists of only variables.
- ii. Break bodies of length 3 or more into a cascade of two-variable-bodied productions.

cleanup steps

Convert the remaining rules into CNF

- i. Arrange that all bodies of length 2 or more consists of only variables.
 - For every terminal a that appears in a body of length≥2, create a new variable, say X_a, and replace a by X_a in all bodies.
 - Then add a new rule $X_a \rightarrow a$.
- ii. Break bodies of length 3 or more into a cascade of two-variable-bodied productions.
 - For each rule of the form

 $A \rightarrow B_1, \dots, B_k$

k \geq 3, introduce new variables Y_1, \dots, Y_{k-2} and replace the rule with

$$\begin{array}{l} \mathbf{A} \rightarrow \mathbf{B}_{1}\mathbf{Y}_{1} \\ \mathbf{Y}_{1} \rightarrow \mathbf{B}_{2}\mathbf{Y}_{2} \\ \cdots \\ \mathbf{Y}_{k-3} \rightarrow \mathbf{B}_{k-2}\mathbf{Y}_{k-2} \\ \mathbf{Y}_{k-2} \rightarrow \mathbf{B}_{k-1}\mathbf{B}_{k} \end{array}$$

Convert the remaining rules into CNF: Example

 $\begin{array}{lll} S_0 \rightarrow \epsilon & \mid 01 \mid 0S1 \\ S \rightarrow 01 \mid 0S1 \end{array} & already cleaned grammar \end{array}$

Arrange that all bodies of length 2 or more consists of only variables.

$$\begin{split} \mathbf{S}_0 &\to \mathbf{\epsilon} ~|~ \mathbf{X}_0 \mathbf{X}_1 ~|~ \mathbf{X}_0 \mathbf{S} ~\mathbf{X}_1 \\ \mathbf{S} &\to \mathbf{X}_0 \mathbf{X}_1 ~|~ \mathbf{X}_0 \mathbf{S} ~\mathbf{X}_1 \\ \mathbf{X}_0 &\to \mathbf{0} \\ \mathbf{X}_1 &\to \mathbf{1} \end{split}$$

Break bodies of length 3 or more into two-variable-bodied productions.

$\mathbf{S}_0 \to \mathbf{X}_0 \mathbf{S} \mathbf{X}_1 \clubsuit$	$S_0 \rightarrow X_0 Y_1$	$Y_1 \rightarrow S X_1$
$\mathbf{S} \to \mathbf{X}_0 \mathbf{S} \mathbf{X}_1 \qquad \clubsuit$	$S \to X_0 Y_2$	$\mathbf{Y}_2 \rightarrow \mathbf{S} \ \mathbf{X}_1$
Grammar in CNF:		
$\mathbf{S}_0 \rightarrow \mathbf{\epsilon} \mid \mathbf{X}_0 \mathbf{X}_1 \mid \mathbf{X}_0 \mathbf{Y}_1$	$\mathbf{Y}_1 \rightarrow \mathbf{S} \ \mathbf{X}_1$	
$\mathbf{S} \rightarrow \mathbf{X}_{0} \mathbf{X}_{1} \mid \mathbf{X}_{0} \mathbf{Y}_{2}$	$\mathbf{Y}_2 \rightarrow \mathbf{S} \ \mathbf{X}_1$	
$\mathbf{X}_{0} \rightarrow 0$		
$X_1 \rightarrow 1$		

Converting into CNF: A Full Example

 $S \rightarrow ABA$ $A \rightarrow aA \mid \varepsilon$ $B \rightarrow bBc \mid \varepsilon$

Step 1. Add a new start variable S₀

Step 2. Eliminate ε-productions.

nullable(G) = {A, B, S, S_0 }

$$\begin{array}{ll} S_0 \rightarrow S & S_0 \rightarrow S \mid \epsilon \\ S \rightarrow ABA & S \rightarrow ABA \mid BA \mid AA \mid AB \mid A \mid B \mid A \mid \epsilon \\ A \rightarrow aA \mid \epsilon & A \rightarrow aA \mid a \mid \epsilon \\ B \rightarrow bBc \mid \epsilon & B \rightarrow bBc \mid \epsilon \mid bc \end{array}$$

$$\begin{split} \mathbf{S}_0 &\to \mathbf{S} \mid \mathbf{\epsilon} \\ \mathbf{S} &\to \mathbf{A}\mathbf{B}\mathbf{A} \mid \mathbf{B}\mathbf{A} \mid \mathbf{A}\mathbf{A} \mid \mathbf{A}\mathbf{B} \mid \mathbf{A} \mid \mathbf{B} \\ \mathbf{A} &\to \mathbf{a}\mathbf{A} \mid \mathbf{a} \\ \mathbf{B} &\to \mathbf{b}\mathbf{B}\mathbf{c} \mid \mathbf{b}\mathbf{c} \end{split}$$

Converting into CNF: A Full Example

Step 3. Eliminate unit productions.

 $\begin{array}{l} S_0 \rightarrow S \mid \epsilon \\ S \rightarrow ABA \mid BA \mid AA \mid AB \mid A \mid B \\ A \rightarrow aA \mid a \\ B \rightarrow bBc \mid bc \end{array}$

$$\begin{split} S_0 &\rightarrow \epsilon \mid ABA \mid BA \mid AA \mid AB \mid aA \mid a \mid bBc \mid bc \\ S &\rightarrow ABA \mid BA \mid AA \mid AB \mid aA \mid a \mid bBc \mid bc \\ A &\rightarrow aA \mid a \\ B &\rightarrow bBc \mid bc \end{split}$$

Step 4. Eliminate useless symbols.

- i. Eliminate non-generating symbols. none
- ii. Eliminate non-reachable symbols. S

Chomsky Normal Form (CNF) Converting into CNF: A Full Example

Step 5. Convert the remaining rules into CNF:

Arrange that all bodies of length 2 or more consists of only variables.

$S_0 \rightarrow \epsilon \mid ABA \mid BA \mid AA \mid AB \mid aA \mid a \mid bBc \mid bc$	$S_0 \rightarrow \epsilon \mid ABA \mid BA \mid AA \mid AB \mid XA \mid a \mid YBZ \mid YZ$
$\mathbf{A} \rightarrow \mathbf{a} \mathbf{A} \mid \mathbf{a}$	$\mathbf{A} \rightarrow \mathbf{X}\mathbf{A} \mid \mathbf{a}$
$\mathbf{B} \rightarrow \mathbf{b}\mathbf{B}\mathbf{c} \mid \mathbf{b}\mathbf{c}$	$\mathbf{B} \rightarrow \mathbf{YBZ} \mid \mathbf{YZ}$
·	$\mathbf{X} \rightarrow \mathbf{a}$
	$\mathbf{Y} \rightarrow \mathbf{b}$
	$\mathbf{Z} \rightarrow \mathbf{c}$

Break bodies of length 3 or more into two-variable-bodied productions.

$S_0 \rightarrow \epsilon \mid ABA \mid BA \mid AA \mid AB \mid XA \mid a \mid YBZ \mid YZ$	$S_0 \rightarrow \epsilon \mid AC \mid BA \mid AA \mid AB \mid XA \mid a \mid YD \mid YZ$
$A \rightarrow XA \mid a$	$\mathbf{C} \rightarrow \mathbf{B}\mathbf{A} \qquad \mathbf{D} \rightarrow \mathbf{B}\mathbf{Z}$
$B \rightarrow YBZ YZ$	$\mathbf{A} \rightarrow \mathbf{X}\mathbf{A} \mid \mathbf{a}$
$\mathbf{X} \rightarrow \mathbf{a}$	$\mathbf{B} \rightarrow \mathbf{YE} \mid \mathbf{YZ}$
$\mathbf{Y} \rightarrow \mathbf{b}$	$\mathbf{E} \rightarrow \mathbf{BZ}$
$Z \rightarrow c$	$\mathbf{X} \rightarrow \mathbf{a}$
	$\mathbf{Y} \rightarrow \mathbf{b}$
	$\mathbf{Z} \rightarrow \mathbf{c}$

Converting into CNF: A Full Example

Grammar in CNF:

$$\begin{split} \mathbf{S}_0 &\rightarrow \epsilon \,|\, \mathbf{AC} \,|\, \mathbf{BA} \,|\, \mathbf{AB} \,|\, \mathbf{AB} \,|\, \mathbf{XA} \,|\, \mathbf{a} \,|\, \mathbf{YD} \,|\, \mathbf{YZ} \\ \mathbf{C} &\rightarrow \mathbf{BA} \\ \mathbf{D} &\rightarrow \mathbf{BZ} \\ \mathbf{A} &\rightarrow \mathbf{XA} \,|\, \mathbf{a} \\ \mathbf{B} &\rightarrow \mathbf{YE} \,|\, \mathbf{YZ} \\ \mathbf{E} &\rightarrow \mathbf{BZ} \\ \mathbf{X} &\rightarrow \mathbf{a} \\ \mathbf{Y} &\rightarrow \mathbf{b} \\ \mathbf{Z} &\rightarrow \mathbf{c} \end{split}$$