

Normal Forms for CFG's

- Eliminating Useless Variables
- Removing Epsilon
- Removing Unit Productions
- Chomsky Normal Form

Variables That Derive Nothing

- ◆ Consider a CFG:

- ◆ $S \rightarrow AB$

- ◆ $A \rightarrow aA \mid a,$

- ◆ $B \rightarrow AB$

- ◆ Although A derives all strings of a's, B derives no terminal strings.

- ◆ Thus, S derives nothing, and the language is **empty**.

Testing Whether a Variable Derives Some Terminal String

- ◆ **Basis:** If there is a production $A \rightarrow w$, where w has no variables, then A derives a terminal string.
- ◆ **Induction:** If there is a production $A \rightarrow \alpha$, where α consists only of terminals and variables known to derive a terminal string, then A derives a terminal string.

Algorithm to Eliminate Variables That Derive Nothing: Non Generating Symbols

1. Discover all variables that derive terminal strings.
2. For all other variables, remove all productions in which they appear either on the left or the right.

Example: Eliminate Non Generating Variables

$S \rightarrow AB \mid C, A \rightarrow aA \mid a, B \rightarrow bB, C \rightarrow c$

- ◆ **Basis:** A and C are identified because of $A \rightarrow a$ and $C \rightarrow c$.
- ◆ **Induction:** S is identified because of $S \rightarrow C$.
- ◆ Nothing else can be identified.
- ◆ **Result:** $S \rightarrow C, A \rightarrow aA \mid a, C \rightarrow c$

Unreachable Symbols

- ◆ Another way a terminal or variable deserves to be eliminated is if it cannot appear in any derivation from the start symbol.
- ◆ **Basis**: We can reach S (the start symbol).
- ◆ **Induction**: if we can reach A , and there is a production $A \rightarrow \alpha$, then we can reach all symbols of α .

Unreachable Symbols

- ◆ Easy inductions in both directions show that when we can discover no more symbols, then we have all and only the symbols that appear in derivations from S .
- ◆ **Algorithm**: Remove from the grammar all symbols not discovered reachable from S and all productions that involve these symbols.

Eliminating Useless Symbols

- ◆ A symbol is *useful* if it appears in some derivation of some terminal string from the start symbol.
- ◆ Otherwise, it is *useless*.
Eliminate all useless symbols by:
 1. Eliminate symbols that derive no terminal string.
 2. Eliminate unreachable symbols.

Example: Useless Symbols

$S \rightarrow AB|AC, A \rightarrow C, C \rightarrow c, B \rightarrow bB$

- ◆ If we eliminated unreachable symbols first, we would find everything is reachable. So no unreachable here.
- ◆ For testing generating symbol, we can get A , C and S are generating and B is not.
- ◆ A, C,S would never get eliminated but B is eliminated.
- ◆ So after eliminating useless symbols:

$S \rightarrow AC, A \rightarrow C, C \rightarrow c$

Epsilon Productions

- ◆ We can almost avoid using productions of the form $A \rightarrow \epsilon$ (called *ϵ -productions*).
 - ◆ The problem is that ϵ cannot be in the language of any grammar that has no ϵ -productions.
- ◆ **Theorem:** If L is a CFL, then $L - \{\epsilon\}$ has a CFG with no ϵ -productions.

Nullable Symbols

- ◆ To eliminate ϵ -productions, we first need to discover the *nullable variables* = variables A such that $A \Rightarrow^* \epsilon$.
- ◆ **Basis**: If there is a production $A \rightarrow \epsilon$, then A is nullable.
- ◆ **Induction**: If there is a production $A \rightarrow \alpha$, and all symbols of α are nullable, then A is nullable.

Example: Nullable Symbols

$S \rightarrow AB, A \rightarrow aA \mid \epsilon, B \rightarrow bB \mid A$

- ◆ **Basis:** A is nullable because of $A \rightarrow \epsilon$.
- ◆ **Induction:** B is nullable because of $B \rightarrow A$.
- ◆ Then, S is nullable because of $S \rightarrow AB$.

Eliminating ϵ -Productions

- ◆ **Key idea:** turn each production $A \rightarrow X_1 \dots X_n$ into a family of productions.
- ◆ For each subset of nullable X 's, there is one production with those eliminated from the right side "in advance."
 - ◆ Except, if all X 's are nullable, do not make a production with ϵ as the right side.

Example: Eliminating ϵ -Productions

$S \rightarrow ABC, A \rightarrow aA \mid \epsilon, B \rightarrow bB \mid \epsilon, C \rightarrow \epsilon$

◆ A, B, C, and S are all nullable.

◆ New grammar:

$S \rightarrow \cancel{ABC} \mid AB \mid \cancel{AC} \mid \cancel{BC} \mid A \mid B \mid \cancel{C}$

$A \rightarrow aA \mid a$

$B \rightarrow bB \mid b$

Note: C is now useless.
Eliminate its productions.

Removal of Unit Productions

- ◆ A *unit production* is one whose right side consists of exactly one variable.
- ◆ These productions can be eliminated.
- ◆ **Key idea:** If $A \Rightarrow^* B$ by a series of unit productions, and $B \rightarrow \alpha$ is a non-unit-production, then add production $A \rightarrow \alpha$.
- ◆ Then, drop all unit productions.

Unit Productions – (2)

- ◆ Find all pairs (A, B) such that $A \Rightarrow^* B$ by a sequence of unit productions only.
- ◆ **Basis**: Surely (A, A) .
- ◆ **Induction**: If we have found (A, B) , and $B \rightarrow C$ is a unit production, then add (A, C) .

Proof That We Find Exactly the Right Pairs

- ◆ By induction on the order in which pairs (A, B) are found, we can show $A \Rightarrow^* B$ by unit productions.
- ◆ Conversely, by induction on the number of steps in the derivation by unit productions of $A \Rightarrow^* B$, we can show that the pair (A, B) is discovered.

Simplification of CFG

- ◆ **Theorem:** if L is a CFL, then there is a CFG for $L - \{\epsilon\}$ that has:
 1. No useless symbols.
 2. No ϵ -productions.
 3. No unit productions.
- ◆ i.e., every right side is either a single terminal or has length ≥ 2 .

Simplification of CFG

- ◆ **Proof:** Start with a CFG for L.
- ◆ Perform the following steps in order:
 1. Eliminate ϵ -productions(Nullable).
 2. Eliminate unit productions.
 3. Eliminate variables that derive no terminal string (useless).
 4. Eliminate variables not reached from the start symbol (useless). **Must be first. Can create unit productions or useless variables.**

Chomsky Normal Form(CNF)

- ◆ A CFG is said to be in *Chomsky Normal Form* if every production is of one of these two forms:
 1. $A \rightarrow BC$ (right side is two variables).
 2. $A \rightarrow a$ (right side is a single terminal).
- ◆ **Theorem:** If L is a CFL, then $L - \{\epsilon\}$ has a CFG in CNF.

Proof of CNF Theorem

- ◆ **Step 1:** “Simplify” the grammar, so every production right side is either a single terminal or of length at least 2.
- ◆ **Step 2:** For each right side \neq a single terminal, make the right side all variables.
 - ◆ For each terminal a create new variable A_a and production $A_a \rightarrow a$.
 - ◆ Replace a by A_a in right sides of length > 2 .

Example: Step 2

- ◆ Consider production $A \rightarrow BcDe$.
- ◆ We need variables A_c and A_e with productions $A_c \rightarrow c$ and $A_e \rightarrow e$.
 - ◆ **Note:** you create at most one variable for each terminal, and use it everywhere it is needed.
- ◆ Replace $A \rightarrow BcDe$ by $A \rightarrow BA_cDA_e$.

CNF Proof – Continued

- ◆ **Step 3:** Break right sides longer than 2 into a chain of productions with right sides of two variables.
- ◆ **Example:** $A \rightarrow BCDE$ is replaced by $A \rightarrow BF$, $F \rightarrow CG$, and $G \rightarrow DE$.
 - ◆ F and G must be used nowhere else.

Example of Step 3 – Continued

- ◆ Recall $A \rightarrow BCDE$ is replaced by $A \rightarrow BF$, $F \rightarrow CG$, and $G \rightarrow DE$.
- ◆ In the new grammar, $A \Rightarrow BF \Rightarrow BCG \Rightarrow BCDE$.
- ◆ **More importantly**: Once we choose to replace A by BF , we must continue to BCG and $BCDE$.
 - ◆ Because F and G have only one production.

CNF: Example

◆ Given CFG, Convert to CNF

- ◆ $S \rightarrow AAC$
- ◆ $A \rightarrow aAb \mid \epsilon$
- ◆ $C \rightarrow aC \mid a$

◆ Eliminating ϵ -Production:

- ◆ Here Variable A is Nullable since $A \rightarrow \epsilon$, so eliminating ϵ -production
 - $S \rightarrow AAC \mid AC \mid C$
 - $A \rightarrow aAb \mid ab$
 - $C \rightarrow aC \mid a$

CNF: Example

◆ Eliminating unit-Production:

- ◆ Discover the unit pair first
- ◆ (S, C) are unit pairs since $S \rightarrow C$, so eliminating Unit productions we have
 - $S \rightarrow AAC \mid AC \mid aC \mid a$
 - $A \rightarrow aAb \mid ab$
 - $C \rightarrow aC \mid a$

◆ Here are no useless symbols since S , A and C all variables are generating and reachable

CNF: Example

◆ Hence Simplified CFG is:

- $S \rightarrow AAC \mid AC \mid aC \mid a$
- $A \rightarrow aAb \mid ab$
- $C \rightarrow aC \mid a$

◆ Converting to CNF

- ◆ For each production of the form

$A \rightarrow X_1X_2X_3 \dots X_n$ for which any X_i introduce new variable for that terminal, So

- ◆ $S \rightarrow AAC \mid AC \mid X_aC \mid a$
- ◆ $X_a \rightarrow a$
- ◆ $A \rightarrow X_aAX_b \mid X_aX_b$
- ◆ $X_b \rightarrow b$
- ◆ $C \rightarrow X_aC \mid a$

CNF: Example

◆ For $S \rightarrow AAC$, Replace with $S \rightarrow AT_1, T_1 \rightarrow AC$

◆ For $A \rightarrow X_aAX_b$, Replace with $A \rightarrow X_aT_2, T_2 \rightarrow AX_b$

◆ So grammar has now all the productions in any one of the form $A \rightarrow BC$ or $A \rightarrow a$. So the CNF grammar is:

- ◆ $S \rightarrow AT_1 \mid AC \mid X_aC \mid a$
- ◆ $T_1 \rightarrow AC$
- ◆ $X_a \rightarrow a$
- ◆ $A \rightarrow X_aT_2 \mid X_aX_b$
- ◆ $T_2 \rightarrow AX_b$
- ◆ $X_b \rightarrow b$
- ◆ $C \rightarrow X_aC \mid a$

Exercise: Convert to CNF

◆ $S \rightarrow AACD$

◆ $A \rightarrow aAb \mid \epsilon$

◆ $C \rightarrow aC \mid a$

◆ $D \rightarrow aDa \mid bDb \mid \epsilon$

◆ Solution:

◆ Removing ϵ productions:

- Here A and D are nullable since $A \rightarrow \epsilon$, $D \rightarrow \epsilon$. So removing ϵ -production
- Replace $S \rightarrow AACD$ by $S \rightarrow AACD \mid ACD \mid AAC \mid AC \mid CD \mid C$
- Replace $A \rightarrow aAb$ by $A \rightarrow aAb \mid ab$
- Replace $D \rightarrow aDa$ by $D \rightarrow aDa \mid aa$ and $D \rightarrow bDb$ by $D \rightarrow bDb \mid bb$

◆ Grammar after removing ϵ -Production

- ◆ $S \rightarrow AACD \mid ACD \mid AAC \mid AC \mid CD \mid C$
- ◆ $A \rightarrow aAb \mid ab$
- ◆ $C \rightarrow aC \mid a$
- ◆ $D \rightarrow aDa \mid aa \mid bDb \mid bb$

◆ **Removing Unit Production:** only (S,C) is unit pair since $S \rightarrow C$ and no others, so removing unit production

- ◆ **Replace $S \rightarrow C$ by $S \rightarrow aC \mid a$ i.e. non unit production of C**

◆ Now, Grammar will be

- ◆ $S \rightarrow AACD \mid ACD \mid AAC \mid AC \mid CD \mid aC \mid a$
- ◆ $A \rightarrow aAb \mid ab$
- ◆ $C \rightarrow aC \mid a$
- ◆ $D \rightarrow aDa \mid aa \mid bDb \mid bb$

Example....

- ◆ Removing the useless symbols: Here Variable S,A,C D all are generating symbols. Also all variables are reachable so no useless symbols. So simplified grammar is :

- ◆ $S \rightarrow AACD \mid ACD \mid AAC \mid AC \mid CD \mid aC \mid a$
- ◆ $A \rightarrow aAb \mid ab$
- ◆ $C \rightarrow aC \mid a$
- ◆ $D \rightarrow aDa \mid aa \mid bDb \mid bb$

- ◆ **Convert this grammar into CNF**

- ◆ **At first converting terminals which are not single at any productions.**

- ◆ Replace $S \rightarrow aC$ by $S \rightarrow X_a C$ and $X_a \rightarrow a$
- ◆ Replace $A \rightarrow aAb$ by $A \rightarrow X_a A X_b$ and $X_b \rightarrow b$
- ◆ Replace $C \rightarrow aC$ by $C \rightarrow X_a C$
- ◆ Replace $D \rightarrow aDa$ by $D \rightarrow X_a D X_a$ and $D \rightarrow bDb$ by $D \rightarrow X_b D X_b$

◆ **Now Grammar will be:**

- ◆ $S \rightarrow AACD \mid ACD \mid AAC \mid AC \mid CD \mid X_a C \mid a$
- ◆ $X_a \rightarrow a$
- ◆ $A \rightarrow X_a A X_b \mid X_a X_b$
- ◆ $X_b \rightarrow b$
- ◆ $C \rightarrow X_a C \mid a$
- ◆ $D \rightarrow X_a D X_a \mid X_a X_a \mid X_b D X_b \mid X_b X_b$

◆ **Now replacing production on right side having more than 2 variables**

- ◆ Replace $S \rightarrow AACD$ by $S \rightarrow AE$, $E \rightarrow AF$ and $F \rightarrow CD$
- ◆ Replace $S \rightarrow ACD$ by $S \rightarrow AF$
- ◆ Replace $S \rightarrow AAC$ by $S \rightarrow AG$, $G \rightarrow AC$
- ◆ Replace $A \rightarrow X_a A X_b$ by $A \rightarrow X_a H$ and $H \rightarrow A X_b$
- ◆ **Replace $D \rightarrow X_a D X_a$ by $D \rightarrow X_a I$ and $I \rightarrow D X_a$**
- ◆ **Replace $D \rightarrow X_b D X_b$ by $D \rightarrow X_b J$ and $J \rightarrow D X_b$**

◆ Now Grammar in CNF will be:

- ◆ $S \rightarrow AE \mid AF \mid AG \mid AC \mid CD \mid X_a C \mid a$
- ◆ $E \rightarrow AF$
- ◆ $F \rightarrow CD$
- ◆ $G \rightarrow AC$
- ◆ $X_a \rightarrow a$
- ◆ $A \rightarrow X_a H \mid X_a X_b$
- ◆ $H \rightarrow AX_b$
- ◆ $X_b \rightarrow b$
- ◆ $C \rightarrow X_a C \mid a$
- ◆ $D \rightarrow X_a I \mid X_a X_a \mid X_b J \mid X_b X_b$
- ◆ $I \rightarrow DX_a$
- ◆ $J \rightarrow DX_b$

Greibach Normal Form

◆ A CFG is in Greibach Normal Form if the Productions are in the following forms –

◆ $A \rightarrow b$

◆ $A \rightarrow bD_1 \dots D_n$

where A, D_1, \dots, D_n are non-terminals and b is a terminal.

Algorithm to Convert a CFG into GNF

- ◆ **Step 1** – If the start symbol **S** occurs on some right side, create a new start symbol **S'** and a new production **S' → S**.
- ◆ **Step 2** – Remove Null productions. (Using the Null production removal algorithm discussed earlier)
- ◆ **Step 3** – Remove unit productions. (Using the Unit production removal algorithm discussed earlier)
- ◆ **Step 4** – Remove all direct and indirect left-recursion.
- ◆ **Step 5** – Do proper substitutions of productions to convert it into the proper form of GNF.

GNF: Example

◆ Convert the following CFG into CNF

◆ $S \rightarrow XY \mid Xb \mid d$

◆ $X \rightarrow aX \mid a$

◆ $Y \rightarrow Xb \mid c$

Solution

◆ Here, **S** does not appear on the right side of any production and there are no unit or null productions in the production rule set. So, we can skip Step 1 to Step 3.

GNF:Example

Step 4

◆ Now after replacing

◆ X in $S \rightarrow XY \mid Xc \mid d$

With $aX \mid a$ we obtain

◆ $S \rightarrow aXY \mid aY \mid aXc \mid mc \mid d.$

◆ And after replacing X in $Y \rightarrow Xb \mid c$
with the right side of $X \rightarrow aX \mid a$ we
obtain $Y \rightarrow aXb \mid ab \mid c.$

GNF: Example

- ◆ Two new productions $C \rightarrow c$ and $D \rightarrow d$ are added to the production set and then we came to the final GNF as the following –

- ◆ $S \rightarrow aXY \mid aY \mid aXC \mid aC \mid d$

- ◆ $X \rightarrow aX \mid a$

- ◆ $Y \rightarrow aXD \mid aD \mid c$

- ◆ $C \rightarrow c$

- ◆ $D \rightarrow d$

- ◆ Now all productions are in the GNF

Left Recursion

- ◆ A grammar becomes left-recursive if it has any non-terminal 'A' whose derivation contains 'A' itself as the left-most symbol.
- ◆ Left-recursive grammar is considered to be a problematic situation for top-down parsers.
- ◆ It becomes hard for it to judge when to stop parsing the left non-terminal and it goes into an infinite loop.

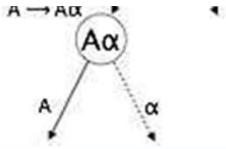
Example of Left Recursion

1. $A \Rightarrow Aa \mid \beta$ direct left recursion for A

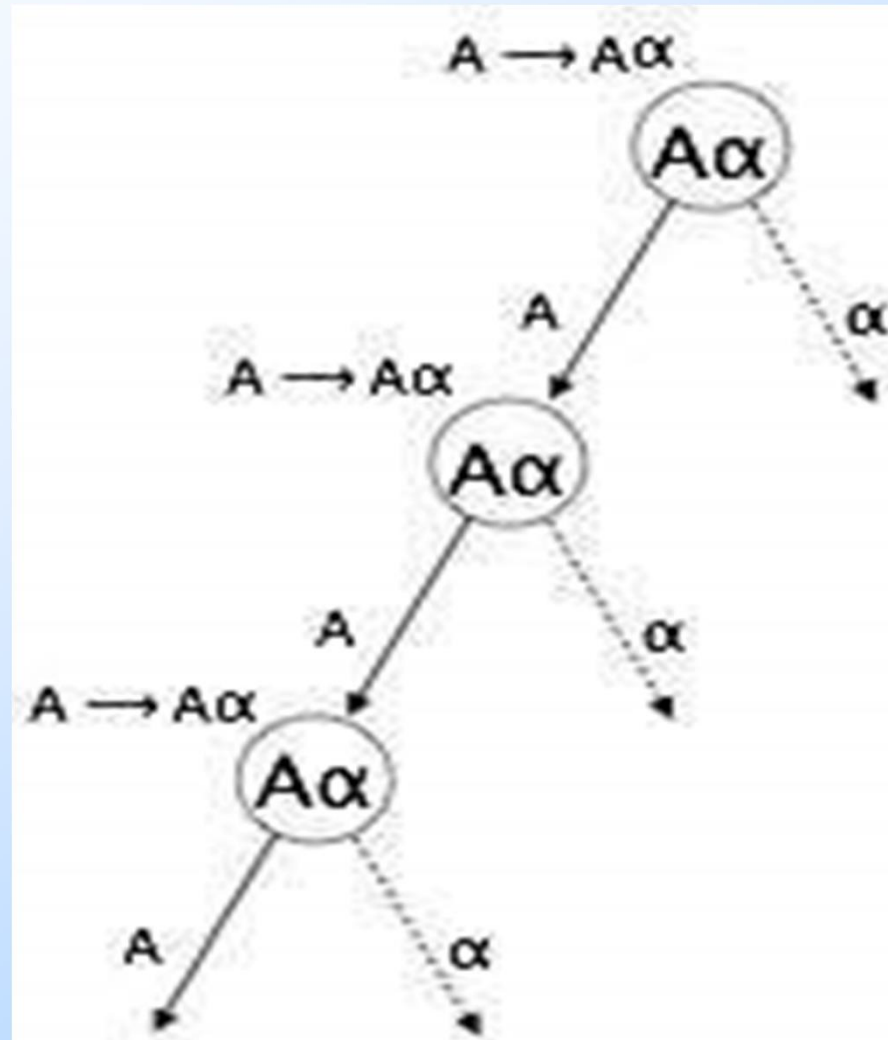
2. $S \Rightarrow Aa \mid \beta$

$A \Rightarrow Sd$ indirect left recursion for S

- ◆ (1) is an example of immediate left recursion, where A is any non-terminal symbol and a represents a string of non-terminals.
- ◆ (2) is an example of indirect-left recursion.



Pictorial View of Derivation from left –recursive grammar



Removal of Left Recursion

- ◆ One way to remove left recursion is to use the following technique:
- ◆ The production $A \rightarrow Aa \mid \beta$ is converted into following productions
 - ◆ $A \rightarrow \beta A'$
 - ◆ $A' \rightarrow aA' \mid \epsilon$

Or without ϵ we can write grammar as:

- ◆ $A \rightarrow \beta A' \mid \beta$
- ◆ $A' \rightarrow aA' \mid a$

Removal of Left Recursion

◆ In general, for grammar

◆ $A \rightarrow Aa_1 \mid Aa_2 \mid Aa_3 \mid \dots \mid Aa_m \mid \beta_1 \mid \beta_2 \dots \mid \beta_n$
is converted into following productions

◆ $A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$

◆ $A' \rightarrow a_1 A' \mid a_2 A' \mid a_3 A' \mid \dots \mid a_m A' \mid \epsilon$

Or without ϵ we can write grammar as:

◆ $A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

◆ $A' \rightarrow a_1 A' \mid a_2 A' \mid a_3 A' \mid \dots \mid a_m A' \mid \epsilon a_1 \mid a_2 \mid a_3 \mid \dots \mid a_m$

Exercise: Remove left recursion

1.

◆ $E \rightarrow E + T \mid T$

◆ $T \rightarrow T * F \mid F$

◆ $F \rightarrow (E) \mid a$

2. $S \rightarrow S + S \mid S - S \mid S * S \mid S / S \mid (S) \mid a$