# Chomsky Normal Form and Pushdown Automata

Mark Greenstreet, CpSc 421, Term 1, 2006/07

- Chomsky Normal Form

- Push Down Automata

# Chomsky Normal Form (CNF)

- A CFG is in Chomsky Normal Form iff
  - Every rule is of the form
    - $A \rightarrow x$, where $A$ is a variable and $x$ is a terminal, or
    - $A \rightarrow BC$, where $A$, $B$ and $C$ are variables.
  - There are no rules of the form $A \rightarrow \epsilon$ unless $A$ is the start variable.

- We'll show that for every CFG, the is a CFG in Chomsky Normal Form that generates the same language.

- CNF is handy at times for proofs. To prove some property of CFLs we can start by writing:

  Let $L$ be an arbitrary CFL, and let $G$ be a CNF CFG for $L$

  . . .

# Every CFL has a CNF Grammar

- Every CFL has a CNF Grammar.

# Every CFL has a CNF Grammar

- Every CFL has a CNF Grammar.

- Proof:
    - Let $L$ be an arbitrary CFL, and let $G$ be a CNF CFG for $L$.

    - QED  ☺

# Every CFL has a CNF Grammar

- Every CFL has a CNF Grammar.

- Proof:
    - Let $L$ be an arbitrary CFL, and let $G$ be a CFG for $L$.
    - We will find each rule of $G$ that violates the restrictions of CNF and replace it other rules that generate the same language but satisfy CNF.
        - First, we'll remove $\epsilon$ rules: $A \rightarrow \epsilon$, where $A$ is not the start symbol.
        - Second, we'll remove unit rules: $A \rightarrow B$, where $A$ and $B$ are variables.
        - Third, we'll convert all rules of the form $A \rightarrow u$, where $u$ has three or more variables or symbols into multiple rules of the form $A \rightarrow A_1 A_2$.
        - Fourth, we'll fix all rules of the form $A \rightarrow Bc$, $A \rightarrow bC$ or $A \rightarrow bc$ where $A$, $B$, and $C$ are variables, and $b$ and $c$ are terminals.

# Every CFL has a CNF Grammar

- Every CFL has a CNF Grammar.

- Proof:
    - Let $L$ be an arbitrary CFL, and let $G$ be a CFG for $L$.
    - We will find each rule of $G$ that violates the restrictions of CNF and replace it other rules that generate the same language but satisfy CNF.
        - First, we'll remove $\epsilon$ rules: $A \rightarrow \epsilon$, where $A$ is not the start symbol.
        - Second, we'll remove unit rules: $A \rightarrow B$, where $A$ and $B$ are variables.
        - Third, we'll convert all rules of the form $A \rightarrow u$, where $u$ has three or more variables or symbols into multiple rules of the form $A \rightarrow A_1 A_2$.
        - Fourth, we'll fix all rules of the form $A \rightarrow Bc$, $A \rightarrow bC$ or $A \rightarrow bc$ where $A$, $B$, and $C$ are variables, and $b$ and $c$ are terminals.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

- Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

- Introduce a new start variable, $S_0$, and the rule $S_0 \to S$

- Eliminate $\epsilon$ transitions.

- Remove unit rules.

- Fix rules that produce long strings.

- Replace terminals with variables.

- We have produced a CNF grammar that generates $L$.
  $\square$.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

○ Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

● Introduce a new start variable, $S_0$, and the rule $S_0 \to S$

○ Eliminate $\epsilon$ transitions.

○ Remove unit rules.

○ Fix rules that produce long strings.

○ Replace terminals with variables.

○ We have produced a CNF grammar that generates $L$.
  $\square$.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

- Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

- Introduce a new start variable, $S_0$, and the rule $S_0 \rightarrow S$

- Eliminate $\epsilon$ transitions. For each rule of the form $A \rightarrow \epsilon$:
  - eliminate the rule, $A \rightarrow \epsilon$, and
  - for every rule of the form $B \rightarrow uAv$, add a new rule $B \rightarrow uv$.
  - This process may produce new $\epsilon$ rules. For example, if we have the rules $A \rightarrow B$ and $B \rightarrow \epsilon$, then eliminating $B \rightarrow \epsilon$ produces the rule $A \rightarrow \epsilon$. We eliminate these new $\epsilon$ rules in the same way. Because the new rules that we get produces shorter strings of terminals and variables than the ones they were derived from, this process eventually terminates.

  Now we have a grammar where $S \rightarrow \epsilon$ is the only possible $\epsilon$ rule.

- Remove unit rules.

- Fix rules that produce long strings.

- Replace terminals with variables.

- We have produced a CNF grammar that generates $L$. $\square$.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

- ○ Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

- ○ Introduce a new start variable, $S_0$, and the rule $S_0 \to S$

- ○ Eliminate $\epsilon$ transitions.

- ● Remove unit rules.   If $A \to B$ and $B$ is a variable:
  - ● eliminate the rule, $A \to B$,
  - ● for every rule of the form $C \to uAv$, add a new rule $C \to uBv$.

  Now we have a grammar where $S_0 \to \epsilon$ is the only possible $\epsilon$ rule and every rule produces a string of one terminal, or at least two terminals and/or variables.

- ○ Fix rules that produce long strings.

- ○ Replace terminals with variables.

- ○ We have produced a CNF grammar that generates $L$.
  □.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

○ Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

○ Introduce a new start variable, $S_0$, and the rule $S_0 \rightarrow S$

○ Eliminate $\epsilon$ transitions.

○ Remove unit rules.

● Fix rules that produce long strings.
$A \rightarrow u_1 u_2 u_3 \ldots u_k$ becomes $A \rightarrow u_1 A_2$, $A_2 \rightarrow u_2 A_3$, $\ldots A_{k-1} \rightarrow u_{k-1} u_k$. Now, each rule produces a single terminal or a string of length two. or is the rule $S_0 \rightarrow \epsilon$.

○ Replace terminals with variables.

○ We have produced a CNF grammar that generates $L$.
$\Box$.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

○ Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

○ Introduce a new start variable, $S_0$, and the rule $S_0 \to S$

○ Eliminate $\epsilon$ transitions.

○ Remove unit rules.

○ Fix rules that produce long strings.

● Replace terminals with variables. For each rule $A \to Bc$ or $A \to bC$ where $B$ and $C$ are variables and $b$ and $c$ are terminals:

  ● replace $A \to Bc$ with $A \to BU_c$, $A \to bC$ with $A \to U_bC$, and $A \to BC$ with $A \to U_bU_b$.

  ● Introduce new rules: $U_b \to b$, etc. $C \to uAv$,

  Now, each rule produces two variables, one terminal, or is the rule $S_0 \to \epsilon$.

○ We have produced a CNF grammar that generates $L$.
  $\Box$.

# Proof: $\forall$ CFL $\exists$ a CNF Grammar

○ Let $L$ be a language, and let $G$ be a CFG for $L$ with start variable $S$.

○ Introduce a new start variable, $S_0$, and the rule $S_0 \rightarrow S$

○ Eliminate $\epsilon$ transitions.

○ Remove unit rules.

○ Fix rules that produce long strings.

○ Replace terminals with variables.

● We have produced a CNF grammar that generates $L$.
$\square$.

# Example: Converting to CNF

- The original grammar:

$$
\begin{array}{rcll}
S & \rightarrow & ASA & | \quad \mathtt{a}B \\
A & \rightarrow & B & | \quad S \\
B & \rightarrow & \mathtt{b} & | \quad \epsilon
\end{array}
$$

# Example: Converting to CNF

● Introduce a new start variable, $S_0$:

$$
\begin{array}{rcl}
 & & \\
S & \rightarrow & ASA \;\mid\; \mathtt{a}B \\
A & \rightarrow & B \;\mid\; S \\
B & \rightarrow & \mathtt{b} \;\mid\; \epsilon
\end{array}
\qquad
\begin{array}{rcl}
S_0 & \rightarrow & S \\
S & \rightarrow & ASA \;\mid\; \mathtt{a}B \\
A & \rightarrow & B \;\mid\; S \\
B & \rightarrow & \mathtt{b} \;\mid\; \epsilon
\end{array}
$$

# Example: Converting to CNF

- Eliminate $\epsilon$ rules:

$$
\begin{array}{lcl}
S_0 & \to & S \\
S & \to & ASA \mid \mathtt{a}B \\
A & \to & B \mid S \\
B & \to & \mathtt{b} \mid \epsilon
\end{array}
\qquad
\begin{array}{lcl}
S_0 & \to & S \\
S & \to & ASA \mid \mathtt{a}B \mid \mathtt{a} \\
A & \to & B \mid S \mid \epsilon \\
B & \to & \mathtt{b}
\end{array}
$$

# Example: Converting to CNF

- Eliminate $\epsilon$ rules:

$$
\begin{array}{lll}
S_0 & \to & S \\
S & \to & ASA \mid aB \mid a \\
\\
A & \to & B \mid S \mid \epsilon \\
B & \to & b
\end{array}
\qquad
\begin{array}{lll}
S_0 & \to & S \\
S & \to & ASA \mid aB \mid a \mid \\
& & AS \mid SA \mid S \\
A & \to & B \mid S \\
B & \to & b
\end{array}
$$

# Example: Converting to CNF

● Eliminate unit rules:

$$S_0 \rightarrow S$$

$$
\begin{array}{lll}
S & \rightarrow & ASA \;\mid\; \mathtt{a}B \;\mid\; \mathtt{a} \;\mid \\
  &            & AS \;\mid\; SA \;\mid\; S \\
A & \rightarrow & B \;\mid\; S \\
\\
B & \rightarrow & \mathtt{b}
\end{array}
$$

$$
\begin{array}{lll}
S_0 & \rightarrow & ASA \;\mid\; \mathtt{a}B \;\mid\; \mathtt{a} \;\mid \\
    &            & AS \;\mid\; SA \\
S & \rightarrow & ASA \;\mid\; \mathtt{a}B \;\mid\; \mathtt{a} \;\mid \\
  &            & AS \;\mid\; SA \\
A & \rightarrow & \mathtt{b} \;\mid\; ASA \;\mid\; \mathtt{a}B \;\mid \\
  &            & \mathtt{a} \;\mid\; AS \;\mid\; SA \\
B & \rightarrow & \mathtt{b}
\end{array}
$$

# Example: Converting to CNF

● Fix rules that produce long strings:

$$
\begin{array}{llllll}
S_0 & \rightarrow & ASA & |\ \ aB & |\ \ a & | \\
    &            & AS  & |\ \ SA &        &   \\
S   & \rightarrow & ASA & |\ \ aB & |\ \ a & | \\
    &            & AS  & |\ \ SA &        &   \\
A   & \rightarrow & b   & |\ \ AA_1 & |\ \ aB & | \\
    &            & a   & |\ \ AS   & |\ \ SA  &   \\
    &            &     &          &         &   \\
B   & \rightarrow & b   &          &         &
\end{array}
\qquad
\begin{array}{llllll}
S_0 & \rightarrow & AA_1 & |\ \ aB & |\ \ a & | \\
    &            & AS   & |\ \ SA &        &   \\
S   & \rightarrow & AA_1 & |\ \ aB & |\ \ a & | \\
    &            & AS   & |\ \ SA &        &   \\
A   & \rightarrow & b    & |\ \ AA_1 & |\ \ aB \\
    &            & a    & |\ \ AS   & |\ \ SA \\
A_1 & \rightarrow & SA   &          &         \\
B   & \rightarrow & b    &          &
\end{array}
$$

# Example: Converting to CNF

- Replace terminals with variables:

$$
\begin{aligned}
S_0 &\rightarrow AA_1 \mid \mathtt{a}B \mid \mathtt{a} \mid \\
&\quad AS \mid SA \\
S_0 &\rightarrow AA_1 \mid \mathtt{a}B \mid \mathtt{a} \mid \\
&\quad AS \mid SA \\
A_1 &\rightarrow SA \\
A &\rightarrow \mathtt{b} \mid AA_1 \mid \mathtt{a}B \mid \\
&\quad \mathtt{a} \mid AS \mid SA \\
A_1 &\rightarrow SA \\
\\
B &\rightarrow \mathtt{b}
\end{aligned}
\qquad
\begin{aligned}
S_0 &\rightarrow AA_1 \mid U_aB \mid \mathtt{a} \mid \\
&\quad AS \mid SA \\
S_0 &\rightarrow AA_1 \mid U_aB \mid \mathtt{a} \mid \\
&\quad AS \mid SA \\
A_1 &\rightarrow SA \\
A &\rightarrow \mathtt{b} \mid AA_1 \mid U_aB \\
&\quad \mathtt{a} \mid AS \mid SA \\
A_1 &\rightarrow SA \\
U_a &\rightarrow \mathtt{a} \\
B &\rightarrow \mathtt{b}
\end{aligned}
$$