# Lecture 12: Cleaning up CFGs and Chomsky Normal form

3 March 2009

In this lecture, we are interested in transforming a given grammar into a cleaner form. We start by describing how to clean up a grammar. Then, we show how to transform a cleaned up grammar into a grammar in Chomsky Normal Form.

# 1 Cleaning up a context-free grammar

The main problem with a general context-free grammar is that it might be complicated, contained parts that can not be used, and not necessarily effective.

It might be useful to think about CFG as a program, that we would like to manipulate (in a similar way that a compiler handles programs). If want to cleanup a program, we need to understand its structure. To this end, for CFGs, we will show a sequence of algorithms that analyze and cleanup a given CFG. Interestingly, these procedures uses similar techniques to the one used by compilers to manipulate programs being compiled.

Note, that some of the cleanup steps are not necessary if one just wants to transform a grammar into Chomsky Normal Form. In particular, Section 1.2 and Section 1.3 are not necessary for the CNF conversion. Note however, that the algorithm of Section 1.3 gives us am immediate way to decide if a grammar is empty or not, see Theorem 1.2.

## 1.1 Example of a messy grammar

For example, consider the following strange very strange grammar.

$$
\text{(G1)} \quad
\begin{array}{rl}
\Rightarrow & S_0 \to S \mid X \mid Z \\
& S \to A \\
& A \to B \\
& B \to C \\
& C \to Aa \\
& X \to C \\
& Y \to aY \mid a \\
& Z \to \epsilon.
\end{array}
$$

This grammar is bad. How bad? Well, it has several problems:

(i) The variable $Y$ can never be derived by the start symbol $S_0$. It is a ***useless*** variable.

(ii) The rule $S \to A$ is redundant. We can replace any appearance of $S$ by $A$, and reducing the number of variables by one. Rule of the form $S \to A$ is called a ***unit production*** (or ***unit rule***.

(iii) The variable $A$ is also ***useless*** since we can note derive any word in $\Sigma^*$ from $A$ (because once we starting deriving from $A$ we get into an infinite loop).

(iv) We also do not like $Z$, since one can generate $\epsilon$ from it (that is $Z \overset{*}{\Rightarrow} \epsilon$. Such a variable is called ***nullable***. We would like to have the property that only the start variable can be derived to $\epsilon$.

We are going to present a sequence of algorithms that transform the grammar to not have these drawbacks.

## 1.2  Removing useless variables unreachable from the start symbol

Given a grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S)$, we would like to remove all variables that are not derivable from $S$. To this end, consider the following algorithm.

> **compReachableVars** $\Big( \mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S) \Big)$
>
> $\quad V_{old} \leftarrow \emptyset$
> $\quad V_{new} \leftarrow \{S\}$
> $\quad$**while** $V_{old} \neq V_{\text{new}}$ **do**
> $\quad\quad V_{\text{old}} \leftarrow V_{\text{new}}.$
> $\quad\quad$**for** $\quad X \in V_{old}$ **do**
> $\quad\quad\quad$**for** $\quad (X \to w) \in \mathcal{R}$ **do**
> $\quad\quad\quad\quad$Add all variables appearing in $w$ to $V_{\text{new}}.$
> $\quad$**return** $V_{\text{new}}.$

Clearly, this algorithm returns all the variables that are derivable form the start symbol $S$. As such, settings $\mathcal{V}' = \textbf{compReachableVars}(\mathcal{G})$ we can set our new grammar to be $\mathcal{G}' = (\mathcal{V}', \Sigma, \mathcal{R}', S)$, where $\mathcal{R}'$ is the set of rules of $\mathcal{R}$ having only variables in $\mathcal{V}'$.

## 1.3  Removing useless variables that do not generate anything

In the next step we remove variables that do not generate any string.

Given a grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S)$, we would like to remove all variables that are not derivable from $S$. To this end, consider the following algorithm.

$$\boxed{\begin{aligned}
&\textbf{compGeneratingVars}\left(\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, \mathsf{S})\right)\\
&\qquad V_{\text{old}} \leftarrow \emptyset\\
&\qquad V_{\text{new}} \leftarrow V_{\text{old}}\\
&\qquad \textbf{do}\\
&\qquad\qquad V_{\text{old}} \leftarrow V_{\text{new}}.\\
&\qquad\qquad \textbf{for}\quad \mathsf{X} \in \mathcal{V} \textbf{ do}\\
&\qquad\qquad\qquad \textbf{for}\quad (\mathsf{X} \rightarrow w) \in \mathcal{R} \textbf{ do}\\
&\qquad\qquad\qquad\qquad \textbf{if } w \in (\Sigma \cup V_{\text{old}})^* \textbf{ then}\\
&\qquad\qquad\qquad\qquad\qquad V_{\text{new}} \leftarrow V_{\text{new}} \cup \{\mathsf{X}\}\\
&\qquad\qquad \textbf{while } (V_{old} \neq V_{\text{new}})\\
&\qquad\qquad \textbf{return } V_{\text{new}}.
\end{aligned}}$$

As such, settings $\mathcal{V}' = \textbf{compReachableVars}(\mathcal{G})$ we can set our new grammar to be $\mathcal{G}' = (\mathcal{V}', \Sigma, \mathcal{R}', \mathsf{S})$, where $\mathcal{R}'$ is the result of removing all rules that uses variables not in $\mathcal{V}'$. In the new grammar, every variable can derive some string.

**Lemma 1.1** *Given a context-free grammar (CFG) $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, \mathsf{S})$ one can compute an equivalent CFG $\mathcal{G}'$ such that any variable of $\mathcal{G}'$ can derive some string in $\Sigma^*$.*

Note, that if a grammar $\mathcal{G}$ has an empty language, then the equivalent grammar generated by Lemma 1.1 will have no variables in it. Namely, given a grammar we have an algorithm to decide if the language it generates is empty or not.

**Theorem 1.2 (CFG emptiness.)** *Given a CFG $\mathcal{G}$, there is an algorithm that decides if the language of $\mathcal{G}$ is empty or not.*

Applying the algorithm of Section 1.2 together with the algorithm of Lemma 1.1 results in a CFG without any useless variables.

**Lemma 1.3** *Given a CFG one can compute an equivalent CFG without any useless variables.*

# 2  Removing $\epsilon$-productions and unit rules from a grammar

Next, we would like to remove $\epsilon$-**production** (i.e., a rule of the form $\mathsf{X} \rightarrow \epsilon$) and **unit-rules** (i.e., a rule of the form $\mathsf{X} \rightarrow \mathsf{Y}$) from the language. This is somewhat subtle, and one needs to be careful in doing this removal process.

## 2.1  Discovering nullable variables

Given a grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, \mathsf{S})$, we are interested in discovering all the nullable variables. A variable $\mathsf{X} \in \mathcal{V}$ is **nullable**, if there is a way derive the empty string from $\mathsf{X}$ in $\mathcal{G}$. This can be done with the following algorithm.

$$\mathbf{compNullableVars}\left(\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, \mathsf{S})\right)$$

$\quad V_{\text{null}} \leftarrow \emptyset$

$\quad \mathbf{do}$

$\qquad V_{\text{old}} \leftarrow V_{\text{null}}.$

$\qquad \mathbf{for} \quad \mathsf{X} \in \mathcal{V} \ \mathbf{do}$

$\qquad\qquad \mathbf{for} \quad (\mathsf{X} \to w) \in \mathcal{R} \ \mathbf{do}$

$\qquad\qquad\qquad \mathbf{if} \ w = \epsilon \ \ \mathbf{or} \ \ w \in (V_{\text{null}})^* \ \mathbf{then}$

$\qquad\qquad\qquad\qquad V_{\text{null}} \leftarrow V_{\text{null}} \cup \{\mathsf{X}\}$

$\qquad \mathbf{while} \ (V_{\text{null}} \neq V_{\text{old}})$

$\qquad \mathbf{return} \ V_{\text{null}}.$

## 2.2 Removing $\epsilon$-productions

A rule is an $\epsilon$-***production*** if it is of the form $VX \to \epsilon$. We would like to remove all such rules from the grammar (or almost all of them).

To this end, we run **compNullableVars** on the given grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, \mathsf{S})$, and get the set of all nullable variable $V_{\text{null}}$. If the start variable is nullable (i.e., $\mathsf{S} \in V_{\text{null}}$), then we create a new start state $\mathsf{S}'$, and add the rules to the grammar

$$\mathsf{S}' \to \mathsf{S} \mid \epsilon.$$

We also now remove all the other rules of the form $\mathsf{X} \to \epsilon$ from $\mathcal{R}$. Let $\mathcal{G}' = (\mathcal{V}', \Sigma, \mathcal{R}', \mathsf{S}')$ be the resulting grammar. The grammar $\mathcal{G}'$ is not equivalent to the original rules, since we missed some possible productions. For example, if we had the rule

$$\mathsf{X} \to \mathsf{ABC},$$

where $\mathsf{B}$ is nullable, then since $\mathsf{B}$ is no longer nullable (we removed all the $\epsilon$-productions form the language), we missed the possibility that $\mathsf{B} \overset{*}{\Rightarrow} \epsilon$. To compensate for that, we need to add back the rule

$$\mathsf{X} \to \mathsf{AC},$$

to the set of rules.

So, for every rule $\mathsf{A} \to X_1 X_2 \dots X_m$ is in $\mathcal{R}'$, we add the rules of the form $\mathsf{A} \to \alpha_1 \dots \alpha_m$ to the grammar, where

(i) If $X_i$ is not nullable (its a character or a non-nullable variable), then $\alpha_i = X_i$.

(ii) If $X_i$ is nullable, then $\alpha_i$ is either $X_i$ or $\epsilon$.

(iii) Not all $\alpha_i$s are $\epsilon$.

Let $\mathcal{G}'' = (\mathcal{V}, \Sigma, \mathcal{R}', \mathsf{S}')$ be the resulting grammar. Clearly, no variable is nullable, except maybe the start variable, and there are no $\epsilon$-production rules (except, again, for the special rule for the start variable).

Note, that we might need to feed $\mathcal{G}''$ into our procedures to remove useless variables. Since this process does not introduce new rules or variables, we have to do it only once.

# 3 Removing unit rules

A ***unit rule*** is a rule of the form $X \rightarrow Z$. We would like to remove all such rules from a given grammar.

## 3.1 Discovering all unit pairs

We have a grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S)$ that has no useless variables or $\epsilon$-predictions. We would like to figure out all the unit pairs. A pair of variables $Y$ and $X$ is a ***unit pair*** if $X \overset{*}{\Rightarrow} Y$ by $\mathcal{G}$. We will first compute all such pairs, and their we will remove all unit

Since there are no $\epsilon$ transitions in $\mathcal{G}$, the only way for $\mathcal{G}$ to derive $Y$ from $X$, is to have a sequence of rules of the form

$$X \rightarrow Z_1, Z_1 \rightarrow Z_2, \ldots, Z_{k-1} \rightarrow Z_z = Y,$$

where all these rules are in $\mathcal{R}$. We will generate all possible such pairs, by generating explicitly the rules of the form $X \rightarrow Y$ they induce.

> **compUnitPairs** $\left( \mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S) \right)$
>      $R_{\text{new}} \leftarrow \left\{ X \rightarrow Y \mid (X \rightarrow Y) \in \mathcal{R} \right\}$
>      **do**
>          $R_{\text{old}} \leftarrow R_{\text{new}}.$
>          **for** $(X \rightarrow Y) \in R_{\text{new}}$ **do**
>              **for** $(Y \rightarrow Z) \in R_{\text{new}}$ **do**
>                  $R_{\text{new}} \leftarrow R_{\text{new}} \cup \{X \rightarrow Z\}.$
>      **while** $(R_{\text{new}} \neq R_{\text{old}})$
>      **return** $R_{\text{new}}.$

## 3.2 Removing unit rules

If we have a rule $X \rightarrow Y$, and $Y \rightarrow w$, then if we want to remove the unit rule $X \rightarrow Y$, then we need to introduce the new rule $X \rightarrow w$. We want to do that for all possible unit pairs.

> **removeUnitRules** $\left( \mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S) \right)$
>      $U \leftarrow \textbf{compUnitPairs}(\mathcal{G})$
>      $\mathcal{R} \leftarrow \mathcal{R} \setminus U$
>      **for** $(X \rightarrow A) \in U$ **do**
>          **for** $(A \rightarrow w) \in R_{\text{old}}$ **do**
>              $\mathcal{R} \leftarrow \mathcal{R} \cup \{X \rightarrow w\}.$
>      **return** $(\mathcal{V}, \Sigma, \mathcal{R}, S).$

We thus established the following result.

**Theorem 3.1** *Given an arbitrary* **CFG***, one can compute an equivalent grammar* $\mathcal{G}'$*, such that* $\mathcal{G}'$ *has no unit rules, no $\epsilon$-productions (except maybe a single $\epsilon$-production for the start variable), and no useless variables.*

# 4 Chomsky Normal Form

***Chomsky Normal Form*** requires that each rule in the grammar is either

(C1) of the form $A \to BC$, where $A$, $B$, $C$ are all variables and neither $B$ nor $C$ is the start variable.

    (That is, a rule has exactly two variables on its right side.)

(C2) $A \to a$, where $A$ is a variable and $a$ is a terminal.

    (A rule with terminals on its right side, has only a single character.)

(C3) $S \to \epsilon$, where $S$ is the start symbol.

    (The start variable can derive $\epsilon$, but this is the only variable that can do so.)

Note, that rules of the form $A \to B$, $A \to BCD$ or $A \to aC$ are all illegal in a $CNF$.

Also a grammar in $CNF$ never has the start variable on the right side of a rule.

Why should we care for $CNF$? Well, its an effective grammar, in the sense that every variable that being expanded (being a node in a parse tree), is guaranteed to generate a letter in the final string. As such, a word $w$ of length $n$, must be generated by a parse tree that has $O(n)$ nodes. This is of course not necessarily true with general grammars that might have huge trees, with little strings generated by them.

## 4.1 Outline of conversion algorithm

All context-free grammars can be converted to $CNF$. We did most of the steps already. Here is an outline of the procedure:

(i) Create a new start symbol $S_0$, with new rule $S_0 \to S$ mapping it to old start symbol (i.e., $S$).

(ii) Remove nullable variables (i.e., variables that can generate the empty string).

(iii) Remove unit rules (i.e., variables that can generate each other).

(iv) Restructure rules with long righthand sides.

The only step we did not describe yet is the last one.

## 4.2 Final restructuring of a grammar into $CNF$

Assume that we already cleaned up a grammar by applying the algorithm of Theorem 3.1 to it. So, we now want to convert this grammar $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S)$ into $CNF$.

**Removing characters from right side of rules.** As a first step, we introduce a variable $V_c$ for every character $c \in \Sigma$ and it to $\mathcal{V}$. Next, we add the rules $V_c \to c$ to the grammar, for every $c \in \Sigma$.

Now, for any string $w \in (\mathcal{V} \cup \Sigma)^*$, let $\widehat{w}$ denote the string, such that any appearance of a character $c$ in $w$, is replaced by $V_c$.

Now, we replace every rule $X \to w$, such that $|w| > 1$, by the rule $X \to \widehat{w}$.

Clearly, (C2) and (C3) hold for the resulting grammar, and furthermore, any rule having variables on the right side, is made only of variables.

**Making rules with only two variables on the right side.** The only remaining problem, is that in the current grammar, we might have rules that are too long, since they have long string on the right side. For example, we might have a rule in the grammar of the form

$$X \to B_1 B_3 \ldots B_k.$$

To make this into a binary rule (with only two variables on the right side, we remove this rule from the grammar, and replace it by the following set of rules

$$X \to B_1 Z_1$$
$$Z_1 \to B_2 Z_2 \qquad\qquad\qquad\qquad Z_2 \to B_3 Z_3$$
$$\ldots$$
$$Z_{k-3} \to B_{k-2} Z_{k-2}$$
$$Z_{k-2} \to B_{k-1} B_k,$$

where $Z_1, \ldots, Z_{k-2}$ are new variables.

We repeat this process, till all rules in the grammar is binary. This gramamr is now in CNF. We summarize our result.

**Theorem 4.1 (CFG $\to$ CNF.)** *Any context-free grammar can be converted into Chomsky normal form.*

## 4.3 An example of converting a CFG into CNF

Let us look at an example grammar with start symbol $S$.

$$(G0) \qquad \boxed{\begin{aligned} \Rightarrow \quad &S \to ASA \mid aB \\ &A \to B \mid S \\ &B \to b \mid \epsilon \end{aligned}}$$

After adding the new start symbol $S_0$, we get the following grammar.

$$(G1) \qquad \boxed{\begin{aligned} \Rightarrow \quad &S_0 \to S \\ &S \to ASA \mid aB \\ &A \to B \mid S \\ &B \to b \mid \epsilon \end{aligned}}$$

**Removing nullable variables** In the above grammar, both $A$ and $B$ are the nullable variables. We have the rule $S \rightarrow ASA$. Since $A$ is nullable, we need to add $S \rightarrow SA$ and $S \rightarrow AS$ and $S \rightarrow S$ (which is of course a silly rule, so we will not waste our time putting it in). We also have $S \rightarrow aB$. Since $B$ is nullable, we need to add $S \rightarrow a$. The resulting grammar is the following.

$$
\text{(G2)} \qquad
\begin{array}{l}
\Rightarrow \quad S_0 \rightarrow S \\
\quad S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \\
\quad A \rightarrow B \mid S \\
\quad B \rightarrow b
\end{array}
$$

**Removing unit rules.** The unit pairs for this grammar are $\{A \rightarrow B, A \rightarrow S, S_0 \rightarrow S\}$. We need to copy the productions for $S$ up to $S_0$, copying the productions for $S$ down to $A$, and copying the production $B \rightarrow b$ to $A \rightarrow b$.

$$
\text{(G3)} \qquad
\begin{array}{l}
\Rightarrow \quad S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS \\
\quad S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \\
\quad A \rightarrow b \mid ASA \mid aB \mid a \mid SA \mid AS \\
\quad B \rightarrow b
\end{array}
$$

**Final restructuring.** Now, we can directly patch any places where our grammar rules have the wrong form for $CNF$. First, if the rule has at least two symbols on its righthand side but some of them are terminals, we introduce new variables which expand into these terminals. For our example, the offending rules are $S_0 \rightarrow aB$, $S \rightarrow aB$, and $A \rightarrow aB$. We can fix these by replacing the $a$'s with a new variable $U$, and adding a rule $U \rightarrow a$.

$$
\text{(G4)} \qquad
\begin{array}{l}
\Rightarrow \quad S_0 \rightarrow ASA \mid UB \mid a \mid SA \mid AS \\
\quad S \rightarrow ASA \mid UB \mid a \mid SA \mid AS \\
\quad A \rightarrow b \mid ASA \mid UB \mid a \mid SA \mid AS \\
\quad B \rightarrow b \\
\quad U \rightarrow a
\end{array}
$$

Then, if any rules have more than two variables on their righthand side, we fix that with more new variables. For the grammar (G4), the offending rules are $S_0 \rightarrow ASA$, $S \rightarrow ASA$, and $A \rightarrow ASA$. We can rewrite these using a new variable $Z$ and a rule $Z \rightarrow SA$. This gives us the $CNF$ grammar shown on the right.

$$
\text{(G5)} \qquad
\begin{array}{l}
\Rightarrow \quad S_0 \rightarrow AZ \mid UB \mid a \mid SA \mid AS \\
\quad S \rightarrow AZ \mid UB \mid a \mid SA \mid AS \\
\quad A \rightarrow b \mid AZ \mid UB \mid a \mid SA \mid AS \\
\quad B \rightarrow b \\
\quad U \rightarrow a \\
\quad Z \rightarrow SA
\end{array}
$$

We are done!