# Chomsky Normal Form

Prakash Panangaden

18$^{\text{th}}$ October 2011

Any context-free language (CFL) has a context-free grammar (CFG) with all the rules in the following form

$$A \rightarrow BC \text{ or } A \rightarrow a$$

where upper case letters are variables (non-terminals) and lower case letters are terminals. We allow, as a special case, the rule $S \rightarrow \varepsilon$ *for the start symbol only*; otherwise we could never have $\varepsilon \in L$. A CFG in this form is said to be in **Chomsky normal form** (CNF).

How do we prove such a thing? We will assume we have a CFG $G$ for our CFL $L$ and we will systematically transform $G$ so that all the rules have the required forms; taking care that each transformation preserves the language. I will not give formal proofs that the transformations preserve the language, it is usually clear, in any case.

We order the variables and process them one by one. We call rules of the form $A \rightarrow \varepsilon$ null rules. We start by getting rid of null rules. Let us suppose that we have a rule $A \rightarrow \varepsilon$ that we want to eliminate. We throw it out but we introduce some new rules. If we see any rule with $A$ on the right hand side we proceed as follows. Suppose, for example, there is a rule $X \rightarrow aABAc$, we keep this rule and add new rules

$$X \rightarrow aBAc \mid aABc \mid aBc.$$

In general, suppose we have a rule

$$X \rightarrow \alpha_1 A \alpha_2 A \alpha_3 A \ldots \alpha_k A \beta$$

where the $\alpha$s and the $\beta$ are strings of symbols (including terminals and variables) without any occurrence of $A$. We introduce new rules

$$X \rightarrow \alpha_1 \alpha_2 A \alpha_3 A \ldots \alpha_k A \beta$$

$$X \rightarrow \alpha_1 A \alpha_2 \alpha_3 A \ldots \alpha_k A \beta$$

$$\ldots\ldots\ldots\ldots\ldots\ldots$$

$$X \rightarrow \alpha_1 A \alpha_2 A \alpha_3 A \ldots \alpha_k \beta$$

$$X \rightarrow \alpha_1 \alpha_2 \alpha_3 A \alpha_4 \ldots \alpha_k A \beta$$

$$X \rightarrow \alpha_1 \alpha_2 A \alpha_3 \alpha_4 \ldots \alpha_k A \beta$$

$$\ldots\ldots\ldots\ldots\ldots\ldots$$

and so on. We first have to have rules with just one of the $A$s removed, then rules with 2 of the $A$s removed and so on until we have accounted for all possible combinations of $A$s being removed. As you can see writing the most general form is neither pleasant nor particularly illuminating; furthermore it is pretty clear how to do it if someone holds a gun to your head.

We get rid of rules of the form $A \to B$ (these are called unit rules) by adding rules of the form $A \to \alpha$ whenever we have $B \to \alpha$ unless $A \to \alpha$ is a unit rule already removed.

We remove rules of the form $X \to ab$ by introducing new variables $\langle a \rangle$ and $\langle b \rangle$ and then adding rules of the form

$$X \to \langle a \rangle \langle b \rangle \quad \langle a \rangle \to a \quad \langle b \rangle \to b.$$

Similarly we get rid of rules of the form $A \to aB$ by introducing the new variable $\langle a \rangle$ and adding the rules

$$X \to \langle a \rangle B \quad \langle a \rangle \to a.$$

Finally we are left with rules of the form

$$X \to u_1 u_2 u_3 \ldots u_k$$

where $k \geq 3$ and the $u_i$ are either terminal symbols or variables. We introduce new variables called $X_1, X_2, \ldots X_k$ and introducing new rules

$$X \to u_1 X_1, X_1 \to u_2 X_2, X_3 \to u_3 X_3, \ldots$$

Some of these rules are in the right form – the ones where the $u_i$ are variables – but some are not; we get new rules of the form $A \to aB$, but we know how to get rid of them. I claim, without proof, that this process will eventually terminate yielding a CFG in CNF.