# 1 Context-free grammar

A **context-free grammar** (CFG) can be defined as the tuple $G = (\Gamma, \Sigma, P, S)$, where:

- $\Gamma$ : variables (A,B,...)

- $\Sigma$ : terminals (a,b,...)

- $P$ : productions (e.g A $\rightarrow$ aB)

- $S$ : start symbol ($S \in \Gamma$)

  Example (Palindromes):

$P \rightarrow \varepsilon$ (the empty string is a palindrome)
$P \rightarrow 0$
$P \rightarrow 1$
$P \rightarrow 0P0$ (any palindrome surrounded by two 0s is also a palindrome)
$P \rightarrow 1P1$ (any palindrome surrounded by two 1s is also a palindrome)

We could write this in a simpler way as:

$P \rightarrow \varepsilon \mid 0 \mid 1 \mid 0P0 \mid 1P1$
(the vertical line represents OR, allowing you to merge multiple productions with the same head)

Other example:

$\Sigma = \{a, b, 0, 1, (, ), *\}$

$I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$
$E \rightarrow I \mid E * E \mid E + E \mid (E)$

Leftmost derivation $E \Rightarrow_{lm} E * E \Rightarrow_{lm} I * E \Rightarrow_{lm} a * E \Rightarrow_{lm} a * (E) \Rightarrow_{lm} a * (E + E) \Rightarrow_{lm} ...$

Rightmost derivation $E \Rightarrow_{rm} E * E \Rightarrow_{rm} E * (E) \Rightarrow_{rm} E * (E + E) \Rightarrow_{rm} E * (E + I) \Rightarrow_{rm} ...$

Continue the above derivations, and show that $a * (a + b00)$ is in the language of E.

# 2 Chomsky Normal Form

Every context-free language without $\varepsilon$ has a grammar G in which all productions have one of the following forms:

1. $A \rightarrow BC$, where $A$, $B$ and $C$ are all variables, or

2. $A \rightarrow a$, where $A$ is variable and $a$ is a terminal

Further, G also has no useless symbols. Such a grammar is in *Chomsky Normal Form* (CNF).

We say that $X$ is a useful symbol if: $S \Rightarrow^* \alpha X \beta \Rightarrow^* w$
$X$ is generating if: $X \Rightarrow^* w$
$X$ is reachable if: $S \Rightarrow^* \alpha X \beta$
Useful symbols are both generating and reachable.

What do you think of this grammar?

$S \rightarrow AB \mid a$

$A \rightarrow b$

**Answer:**

$B$ is a useless symbol. It's reachable ($S \rightarrow AB$), but not generating. Because of this, we should delete $S \rightarrow AB$. This leaves us with

$S \rightarrow a$

$A \rightarrow b$

Again, A is useless. It's generating ($A \rightarrow b$), but it's not reachable from S. Deleting $A \rightarrow b$ leaves us with

$S \rightarrow a$

which is the final form of our grammar.

Variable $A$ is nullable if $A \Rightarrow^* \varepsilon$

To convert a grammar to CNF, we should:

1. remove useless symbols

2. remove $\varepsilon$ productions (e.g. $A \rightarrow \varepsilon$)

3. remove unit productions (e.g. $A \rightarrow B$)

4. modify bodies of length 2 or more to contain only variables

5. break bodies of length 3 or more

Example 1:

$S \rightarrow AB$

$A \rightarrow aAA \mid \varepsilon$

$B \rightarrow bBB \mid \varepsilon$

**Answer:**

All symbols are useful. $A$ and $B$ are reachable ($S \Rightarrow AB$), and they are both generating ($A \Rightarrow^* a$ and $B \Rightarrow^* b$).

We have to get rid of $\varepsilon$-productions $A \rightarrow \varepsilon$ and $B \rightarrow \varepsilon$

$S \rightarrow AB \mid A \mid B$

$A \rightarrow aAA \mid aA \mid a$

$B \rightarrow bBB \mid bB \mid b$

We should remove unit productions $S \rightarrow A$ and $S \rightarrow B$

$S \rightarrow AB \mid aAA \mid aA \mid a \mid bBB \mid bB \mid b$

$A \rightarrow aAA \mid aA \mid a$

$B \rightarrow bBB \mid bB \mid b$

To modify bodies of length 2 or more to contain only variables, we introduce $A_1 \rightarrow a$ and $B_1 \rightarrow b$

$S \rightarrow AB \mid A_1AA \mid A_1A \mid a \mid B_1BB \mid B_1B \mid b$

$A \rightarrow A_1AA \mid A_1A \mid a$

$B \rightarrow B_1BB \mid B_1B \mid b$

$A_1 \rightarrow a$

$B_1 \rightarrow b$

To break up bodies of length 3 or more, we introduce $A_2 \rightarrow A_1A$ and $B_2 \rightarrow B_1B$

$S \rightarrow AB \mid A_2A \mid A_1A \mid a \mid B_2B \mid B_1B \mid b$

$A \rightarrow A_2A \mid A_1A \mid a$

$B \rightarrow B_2B \mid B_1B \mid b$

$A_1 \rightarrow a$

$B_1 \rightarrow b$
$A_2 \rightarrow A_1 A$
$B_2 \rightarrow B_1 B$
The above grammar is in CNF, because it satisfies the above requirements. If we still had more productions with body length 3 or more, we would have to repeat this last step again in a similar fashion.

Example 2:
$S \rightarrow aXbX$
$X \rightarrow aY \mid bY \mid \varepsilon$
$Y \rightarrow X \mid c$


Example 3:
$S \rightarrow AbA \mid a$
$A \rightarrow Aa \mid \varepsilon$