

Production Grammars

A *production grammar* (also known as a generative grammar, and which I will call simply grammar from now on) is a rewrite system that describes how to *generate* strings using a set of production rules.

Example 1: Here is a simple grammar given by two rules:

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned} \tag{5.1}$$

These rules say that you can rewrite symbol S into aSb , or into the empty string. Here is a sequence of rewrites showing that these rules, starting with symbol S , can generate the string $aaabbb$:

$$\begin{aligned} \underline{S} &\Rightarrow a\underline{S}b \\ &\Rightarrow aa\underline{S}bb \\ &\Rightarrow aaa\underline{S}bbb \\ &\Rightarrow aaabbb \end{aligned}$$

(At every step, I indicated which symbol gets rewritten by underlining it.) It's not too difficult to see that such a grammar can generate all strings of the form $a^n b^n$ for any $n \geq 0$.

Example 2: Here is a slightly more complicated grammar, given by five

rules:

$$\begin{aligned}
 S &\rightarrow TB \\
 T &\rightarrow aTb \\
 T &\rightarrow \epsilon \\
 B &\rightarrow bB \\
 B &\rightarrow \epsilon
 \end{aligned}
 \tag{5.2}$$

Here is a sequence of rewrites showing that these rules, starting with symbol S , can generate the string $aabbb$:

$$\begin{aligned}
 \underline{S} &\Rightarrow \underline{T}B \\
 &\Rightarrow a\underline{T}bB \\
 &\Rightarrow aa\underline{T}bbB \\
 &\Rightarrow aabb\underline{B} \\
 &\Rightarrow aabbb\underline{B} \\
 &\Rightarrow aabbb
 \end{aligned}$$

Symbols S , T , B are intermediate (or nonterminal) symbols used during the rewrites, as opposed to a and b which are symbols in the strings that we care about generating. Again, it is not difficult to see that this grammar generates strings of the form $a^n b^m$ where $m \geq n \geq 0$.

All of the above grammars have the characteristic that the left-hand side of every rule consists of a single nonterminal symbol, that is, every rule is of the form $A \rightarrow \dots$ for some nonterminal A . We call grammars made up of such rules *context-free grammars*, and they are an important class of grammars.

Example 3: Here is a grammar that is *not* context-free (also called unrestricted):

$$\begin{aligned}
 S &\rightarrow AB \\
 B &\rightarrow XbBc \\
 B &\rightarrow \epsilon \\
 bX &\rightarrow Xb \\
 A &\rightarrow AA \\
 A &\rightarrow \epsilon \\
 AX &\rightarrow a \\
 aX &\rightarrow Xa
 \end{aligned}
 \tag{5.3}$$

Intuitively, these rules let us expand the initial A into a sequence of A s, and the initial B into a sequence of b along with the same number of X s on the left of the b s and c s on the right. Rules allow the X s to "migrate" to the nearest A s and interact with them to produce a s.

Here is a sequence of rewrites showing how to generate $aabbcc$:

$$\begin{aligned}
\underline{S} &\Rightarrow \underline{AB} \\
&\Rightarrow \underline{AXbBc} \\
&\Rightarrow \underline{AXbXbBcc} \\
&\Rightarrow \underline{AXbXbcc} \\
&\Rightarrow \underline{AXXbbcc} \\
&\Rightarrow \underline{AAXXbbcc} \\
&\Rightarrow \underline{AaXbbcc} \\
&\Rightarrow \underline{AXabbcc} \\
&\Rightarrow aabbcc
\end{aligned}$$

This grammar generates all strings of the form $a^n b^n c^n$ for $n \geq 0$.

Formal Definitions

Definition: A production grammar is a tuple $G = (N, \Sigma, R, S)$ where

- N is a finite set of nonterminal symbols;
- Σ is a finite set of terminal symbols;
- R is a finite set of rules, each of the form $w_1 \rightarrow w_2$ where $w_1, w_2 \in (N \cup \Sigma)^*$ and w_1 has at least one nonterminal symbol;
- $S \in N$ is a nonterminal symbol called the start symbol.

Rewriting using a grammar G is defined using a relation \Rightarrow_G :

$$uw_1v \Rightarrow_G uw_2v \text{ if } w_1 \rightarrow w_2 \in R$$

and generalizing to the multi-step rewrite relation \Rightarrow_G^* defined by taking $w_1 \Rightarrow_G^* w_2$ if $w_1 = w_2$ or $\exists u_1, \dots, u_k$ such that $w_1 \Rightarrow_G u_1$, $u_1 \Rightarrow_G u_2$, \dots , $u_{k-1} \Rightarrow_G u_k$, and $u_k \Rightarrow_G w_2$. We usually drop the G when it's clear from context.

The language $L(G)$ of grammar G is the set of all strings of terminals that can be generated from the start symbol of the grammar:

$$L(G) = \{w \in \Sigma^* \mid S \Longrightarrow_G^* w\}$$

Context-Free Languages

An important class of languages is the *context-free languages*. A language is *context-free* if there exists a context-free grammar that can generate it.

It is easy to see that regular languages are context-free. To show that, it suffices to show that to every deterministic finite automaton there exists a context-free grammar that generates the language accepted by the automaton.

Let $M = (Q, \Sigma, \delta, s, F)$ be a deterministic finite automaton. Construct the grammar $G_M = (N, \Sigma, R, S)$ by taking $N = Q$ and $S = s$, and by having one rule in R of the form

$$p \rightarrow aq$$

for every transition in M of the form $\delta(p, a) = q$, and one rule in R of the form

$$p \rightarrow \epsilon$$

for every $p \in F$.

Since $\{a^n b^n \mid n \geq 0\}$ is context-free by grammar (1) above but not regular, the class of context-free languages is a strictly larger class of languages than the regular languages.

