

Context-Free Grammars

Informal Comments

- ◆ A *context-free grammar* is a notation for describing languages.
- ◆ It is more powerful than finite automata or RE's, but still cannot define all possible languages.
- ◆ Useful for nested structures, e.g., parentheses in programming languages.

Informal Comments – (2)

- ◆ Basic idea is to use “variables” to stand for sets of strings (i.e., languages).
- ◆ These variables are defined recursively, in terms of one another.
- ◆ Recursive rules (“productions”) involve only concatenation.
- ◆ Alternative rules for a variable allow union.

Example: CFG for $\{ 0^n 1^n \mid n \geq 1 \}$

◆ Productions:

$S \rightarrow 01$

$S \rightarrow 0S1$

◆ **Basis:** 01 is in the language.

◆ **Induction:** if w is in the language, then so is $0w1$.

Example: CFG for palindromes

◆ Productions:

$P \rightarrow \epsilon$

$P \rightarrow 0$

$P \rightarrow 1$

$P \rightarrow 0P0$

$P \rightarrow 1P1$

◆ **Basis:** $\epsilon, 0, 1$ are palindromes.

◆ **Induction:** if w is in the language, then so are $0w0$ and $1w1$.

Shorthand Forms

◆ Several productions can be merged into a shorthand form using the “|” sign

◆ Productions for 0^n1^n :

$$S \rightarrow 01 \mid 0S1$$

◆ Productions for palindromes:

$$P \rightarrow \epsilon \mid 0 \mid 1 \mid 0P0 \mid 1P1$$

CFG Formalism

- ◆ *Terminals* = symbols of the alphabet of the language being defined.
- ◆ *Variables* = *nonterminals* = a finite set of other symbols, each of which represents a language.
- ◆ *Start symbol* = the variable whose language is the one being defined.

Productions

- ◆ A *production* has the form **variable** \rightarrow **string of variables and terminals**.
- ◆ **Convention:**
 - ◆ A, B, C, \dots are variables.
 - ◆ a, b, c, \dots are terminals.
 - ◆ \dots, X, Y, Z are either terminals or variables.
 - ◆ \dots, w, x, y, z are strings of terminals only.
 - ◆ $\alpha, \beta, \gamma, \dots$ are strings of terminals and/or variables.

Example: Formal CFG

- ◆ Here is a formal CFG for $\{ 0^n 1^n \mid n \geq 1 \}$.
- ◆ Terminals = $\{0, 1\}$.
- ◆ Variables = $\{S\}$.
- ◆ Start symbol = S .
- ◆ Productions =
 - $S \rightarrow 01$
 - $S \rightarrow 0S1$

Derivations – Intuition

- ◆ We *derive* strings in the language of a CFG by starting with the start symbol, and repeatedly replacing some variable A by the right side of one of its productions.
 - ◆ That is, the “productions for A ” are those that have A on the left side of the \rightarrow .

Derivations – Formalism

◆ We say $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production.

◆ **Example:** $S \rightarrow 01$; $S \rightarrow 0S1$.

◆ $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111$.



The diagram illustrates the derivation $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111$. Each step is marked with a colored diamond: a green diamond for the first step, a red diamond for the second, and a purple diamond for the third. The symbols S in each intermediate string are circled in the same color as their respective step. Curved arrows connect the circled S in one string to the circled $0S1$ in the next string, showing the replacement process.

Example: Iterated Derivation

\Rightarrow^* means "zero or more derivation steps."


- ◆ $S \rightarrow 01; S \rightarrow 0S1.$
- ◆ $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111.$
- ◆ So $S \Rightarrow^* S; S \Rightarrow^* 0S1; S \Rightarrow^* 00S11;$
 $S \Rightarrow^* 000111.$

Sentential Forms

- ◆ Any string of variables and/or terminals derived from the start symbol is called a *sentential form*.
- ◆ Formally, α is a sentential form iff

$$S \Rightarrow^* \alpha$$

Language of a Grammar

- ◆ If G is a CFG, then $L(G)$, the *language of G* , is $\{w \mid S \Rightarrow^* w\}$.
 - ◆ **Note:** w must be a terminal string, S is the start symbol.
- ◆ **Example:** G has productions $S \rightarrow \epsilon$ and $S \rightarrow 0S1$.
- ◆ $L(G) = \{0^n 1^n \mid n \geq 0\}$. **Note:** ϵ is a legitimate right side.

Context-Free Languages

- ◆ A language that is defined by some CFG is called a *context-free language*.
- ◆ There are CFL's that are not regular languages, such as the example just given.
- ◆ But not all languages are CFL's.
- ◆ **Intuitively**: CFL's can count two things, not three.

Leftmost and Rightmost Derivations

- ◆ Derivations allow us to replace any of the variables in a string.
- ◆ Leads to many different derivations of the same string.
- ◆ By forcing the leftmost variable (or alternatively, the rightmost variable) to be replaced, we avoid these “distinctions without a difference.”

Leftmost Derivations

- ◆ Say $wA\alpha \Rightarrow_{lm} w\beta\alpha$ if w is a string of terminals only and $A \rightarrow \beta$ is a production.
- ◆ Also, $\alpha \Rightarrow_{lm}^* \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{lm} steps.

Example: Leftmost Derivations

- ◆ Balanced-parentheses grammar:

$$S \rightarrow SS \mid (S) \mid ()$$

- ◆ $S \Rightarrow_{lm} SS \Rightarrow_{lm} (S)S \Rightarrow_{lm} (()S \Rightarrow_{lm} (()()$

- ◆ Thus, $S \Rightarrow_{lm}^* (()()$

- ◆ $S \Rightarrow SS \Rightarrow S() \Rightarrow (S)() \Rightarrow (()()$ is a derivation, but not a leftmost derivation.

Rightmost Derivations

- ◆ Say $\alpha Aw \Rightarrow_{rm} \alpha\beta w$ if w is a string of terminals only and $A \rightarrow \beta$ is a production.
- ◆ Also, $\alpha \Rightarrow_{rm}^* \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{rm} steps.

Example: Rightmost Derivations

- ◆ Balanced-parentheses grammar:

$$S \rightarrow SS \mid (S) \mid ()$$

- ◆ $S \Rightarrow_{\text{rm}} SS \Rightarrow_{\text{rm}} S() \Rightarrow_{\text{rm}} (S)() \Rightarrow_{\text{rm}} ((()))()$

- ◆ Thus, $S \Rightarrow_{\text{rm}}^* ((()))()$

- ◆ $S \Rightarrow SS \Rightarrow SSS \Rightarrow S()S \Rightarrow ()()S \Rightarrow ()()()$ is neither a rightmost nor a leftmost derivation.