

Context-Free Languages and Parse Trees

Mridul Aanjaneya



Stanford University

July 12, 2012

Context-Free Grammars

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

Context-Free Grammars

- Basic idea is to use **variables** to stand for sets of strings (i.e., languages).
- These variable 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 $0w1$.

- **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.

- A **production** has the form
variable \rightarrow string of **variables** and **terminals**
- Convention
 - **A, B, C, ...** are **variables**.
 - **a, b, c, ...** are **terminals**.
 - **..., X, Y, Z** are either **terminals** or **variables**.
 - **..., 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 =
 $S \rightarrow 01$
 $S \rightarrow 0S1$

- 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 .

- We say $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production.
- **Example:** $S \rightarrow 01$; $S \rightarrow 0S1$.
- $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111$.

Iterated Derivation

- \Rightarrow^* means **zero or more** derivation steps.
- **Basis:** $\alpha \Rightarrow^* \alpha$ for any string α .
- **Induction:** if $\alpha \Rightarrow^* \beta$ and $\beta \Rightarrow \gamma$, then $\alpha \Rightarrow^* \gamma$.

Example: Iterated Derivation

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

- 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 $S \rightarrow \varepsilon$ and $S \rightarrow 0S1$.
 - **Note:** ε is a valid right hand side.
- $L(G) = \{0^n 1^n \mid n \geq 0\}$.

- 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.
- **Intuition:** CFL's can count **two** things, not **three**.

- Grammars for programming languages are often written in **Backus-Naur Form (BNF)**.
- Variables are words in $\langle \dots \rangle$, e.g., $\langle \text{statement} \rangle$.
- Terminals are often **multicharacter** strings indicated by boldface or underline, e.g., **while** or WHILE.

- Symbol $::=$ is often used for \rightarrow .
- Symbol $|$ is used for **or**.
 - A shorthand for a list of productions with the **same** left side.

Example: $S \rightarrow 0S1 \mid 01$ is a shorthand for
 $S \rightarrow 0S1$ and $S \rightarrow 01$.

BNF Notation: Kleene Closure

- Symbol \dots is used for **one or more**.
- **Example:** $\langle \text{digit} \rangle ::= 0|1|2|3|4|5|6|7|8|9$
 $\langle \text{unsigned integer} \rangle ::= \langle \text{digit} \rangle \dots$
 - **Note:** that's **not** exactly the $*$ for RE's.
- **Translation:** Replace $\alpha \dots$ with a new variable **A** and productions $\mathbf{A} \rightarrow \mathbf{A}\alpha \mid \alpha$.
- Grammar for unsigned integers can be replaced by:
 $\mathbf{U} \rightarrow \mathbf{U}\mathbf{D} \mid \mathbf{D}$
 $\mathbf{D} \rightarrow 0|1|2|3|4|5|6|7|8|9$

BNF Notation: Optional Elements

- Surround one or more symbols by $[\dots]$ to make them **optional**.
- **Example:** $\langle \text{statement} \rangle ::= \mathbf{if} \langle \text{condition} \rangle \mathbf{then} \langle \text{statement} \rangle [\mathbf{;else} \langle \text{statement} \rangle]$
- **Translation:** replace $[\alpha]$ by a new variable **A** with productions
 $A \rightarrow | \epsilon.$
- Grammar for **if-then-else** can be replaced by:
 $S \rightarrow iCtSA$
 $A \rightarrow ;eS \mid \epsilon$

BNF Notation: Grouping

- Use $\{\dots\}$ to surround a sequence of symbols that need to be treated as a **unit**.
 - Typically, they are followed by a \dots for **one or more**.
- **Example:**
 $\langle \text{statement list} \rangle ::= \langle \text{statement} \rangle [\{ ; \langle \text{statement} \rangle \} \dots]$

BNF Notation: Grouping (Translation)

- You may, if you wish, create a new variable **A** for $\{\alpha\}$.
- One production for **A**: $A \rightarrow \alpha$.
- Use **A** in place of $\{\alpha\}$.

Example: Grouping

$$L \rightarrow S[\{;S\}..]$$

- Replace $L \rightarrow S[A..]; A \rightarrow ;S$
 - A stands for $\{;S\}$.
- Then by $L \rightarrow SB; B \rightarrow A.. \mid \varepsilon; A \rightarrow ;S$
 - B stands for $[A..]$ (zero or more A 's).
- Finally by $L \rightarrow SB; B \rightarrow C \mid \varepsilon; C \rightarrow AC \mid A; A \rightarrow ;S$.
 - C stands for $A..$

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

Example: Rightmost Derivations

- Balanced parentheses grammar:

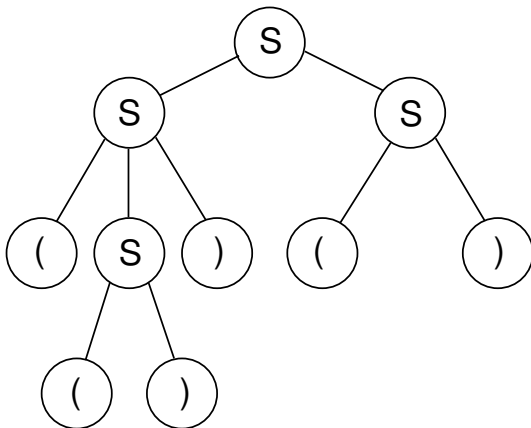
$$S \rightarrow SS \mid (S) \mid ()$$

- $S \Rightarrow_{rm} SS \Rightarrow_{rm} S() \Rightarrow_{rm} (S)() \Rightarrow_{rm} (())()$
- Thus, $S \Rightarrow_{rm}^* (())()$
- $S \Rightarrow SS \Rightarrow SSS \Rightarrow S()S \Rightarrow ()()S \Rightarrow ()()()$ is **neither** a **rightmost** derivation **nor** a **leftmost** derivation.

- **Parse trees** are trees labeled by symbols of a particular CFG.
- **Leaves:** labeled by a **terminal** or ϵ .
- **Interior nodes:** labeled by a **variable**.
 - Children are labeled by the **right** side of a production for the parent.
- **Root:** must be labeled by the **start** symbol.

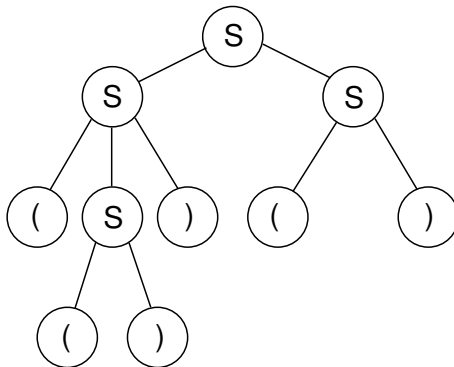
Example: Parse Tree

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



Yield of a Parse Tree

- The concatenation of the labels of the leaves in **left-to-right** order is called the **yield** of the parse tree.
- **Example:** Yield of the given parse tree is $((\))()$.

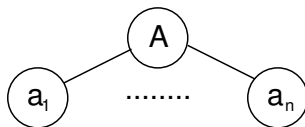


Parse Trees, Leftmost and Rightmost Derivations

- For every parse tree, there is a unique leftmost and a unique rightmost derivation.
- We'll prove:
 - ① If there is a parse tree with root labeled A and yield w , then $A \Rightarrow_{lm}^* w$.
 - ② If $A \Rightarrow_{lm}^* w$, then there is a parse tree with root A and yield w .

Part 1: Basis

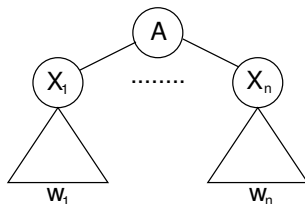
- Induction on the **height** (length of the **longest** path from the root) of the tree.
- **Basis:** height 1. Tree looks like



- $A \rightarrow a_1 \dots a_n$ must be a production.
- Thus, $A \Rightarrow_{\text{lm}}^* a_1 \dots a_n$.

Part 1: Induction

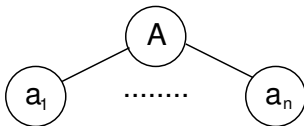
- Assume (1) for trees of height $< h$, and let this tree have height h :
- By **IH**, $X_i \Rightarrow_{lm}^* w_i$.
 - **Note:** If X_i is a terminal, then $X_i = w_i$.



- Thus, $A \Rightarrow_{lm} X_1 \dots X_n \Rightarrow_{lm}^* w_1 X_2 \dots X_n \Rightarrow_{lm}^* w_1 w_2 X_3 \dots X_n \Rightarrow_{lm}^* \dots \Rightarrow_{lm}^* w_1 \dots w_n$

- Given a leftmost derivation of a terminal string, we need to prove the **existence** of a parse tree.
- The proof is an **induction** on the **length** of the derivation.

- If $A \Rightarrow_{lm}^* a_1 a_2 \dots a_n$ by a one-step derivation, then there must be a parse tree

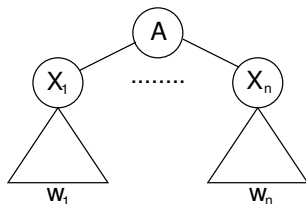


Part 2: Induction

- Assume (2) for derivations of fewer than $k > 1$ steps, and let $A \Rightarrow_{\text{Im}}^* w$ be a k -step derivation.
- First step is $A \Rightarrow_{\text{Im}} X_1 \dots X_n$.
- **Key point:** w can be divided such that the first portion is derived from X_1 , the next is derived from X_2 , and so on.
 - If X_i is a terminal, then $w_i = X_i$.

Part 2: Induction

- That is, $X_i \Rightarrow_{lm}^* w_i$ for all i such that X_i is a variable.
 - And the derivation takes fewer than k steps.
- By the **IH**, if X_i is a variable, then there is a parse tree with root X_i and yield w_i .
- Thus, there is a parse tree



Parse Trees and Rightmost Derivations

- The ideas are essentially the **mirror image** of the proof of the leftmost derivations.
- Left to the **imagination!**

Parse Trees and Any Derivation

- The proof that you can obtain a parse tree from a **leftmost derivation** doesn't really depend on **leftmost**.
- First step **still** has to be $A \Rightarrow_{\text{lm}} X_1 \dots X_n$.
- And **w** can still be divided such that the first portion is **derived** from X_1 , the next is derived from X_2 , and so on.

Ambiguous Grammars

- A CFG is **ambiguous** if there is a string in the language that is the yield of **two or more** parse trees.
- **Example:** $S \rightarrow SS \mid (S) \mid ()$
- **Two** parse trees for $()()()!$

Ambiguity, Leftmost and Rightmost Derivations

- If there are two **different** parse trees, they must produce **two different leftmost derivations** by the construction given in the proof.
- Conversely, two **different** leftmost derivations produce **different parse trees** by the other part of the proof.
- Likewise for **rightmost** derivations.

Ambiguity, Leftmost and Rightmost Derivations

- Thus, **equivalent** definitions for **ambiguous grammar** are:
 - ① There is a string in the language that has **two different leftmost derivations**.
 - ② There is a string in the language that has **two different rightmost derivations**.