

Lec 09. Context-free language and Parse trees

Eunjung Kim

DERIVATION

Consider a CFG $G = (V, \Sigma, R, S)$, $u, v, w \in (\Sigma \cup V)^*$ (a string of terminals and nonterminals) and a variable (nonterminal) $A \in V$.

YIELD, DERIVE, DERIVATION

- We say that uAv **yields** uwv , written $uAv \Rightarrow_G uwv$, if G has the rule $A \rightarrow w$; put another way, uwv is obtained by substituting a variable in the string uAv by the body of a rule whose head is the said variable.
- We say that u **derives** v , written $u \Rightarrow_G^* v$, if $u = v$ or there is a sequence u_1, \dots, u_k for some $k \geq 1$ such that

$$u \Rightarrow_G u_1 \Rightarrow_G \dots \Rightarrow_G u_k \Rightarrow_G v$$

and the sequence is called a **derivation**.

We omit the subscript G in \rightarrow_G and \Rightarrow_G if the CFG under consideration is clear in the context.

DERIVATION BY EXAMPLE

We want to describe, as CFG,

- all arithmetic expressions with $+$ and \times
- over the variables of the form $(a \cup b)(a \cup b \cup 0 \cup 1)^*$.

Consider the following CFG $G_{ari} = (\{E, I\}, \{a, b, 0, 1, +, \times, (,)\}, R, E)$, where R consists of the following rules

- 1 $E \rightarrow I$
- 2 $E \rightarrow E + E$
- 3 $E \rightarrow E \times E$
- 4 $E \rightarrow (E)$
- 5 $I \rightarrow a$
- 6 $I \rightarrow b$
- 7 $I \rightarrow Ia$
- 8 $I \rightarrow Ib$
- 9 $I \rightarrow I0$
- 10 $I \rightarrow I1$

CONTEXT-FREE LANGUAGE

For a CFG $G = (V, \Sigma, R, S)$, the language of G , denoted by $L(G)$ is the set of all strings consisting of terminals (only) that have derivations from the start symbol, i.e.

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

A language is said to be **context-free** if it is the language of a context-free grammar. A context-free languages is often abbreviated as CFL.

SOME REMARKS ON CFG

- In general, the rule of a grammar has the form $u \rightarrow v$ with **both** u and v are strings of terminals and nonterminals.
- A grammar is context-free if the head u is a nonterminal (variable) in all the rules; we do not need to consider the context.
- Different restrictions on the grammar define the hierarchy of formal languages.

Class	Languages	Automaton	Rules	Word Problem	Example
type-0	recursively enumerable	Turing machine	no restriction	undecidable	Post's corresp. problem
type-1	context sensitive	linear-bounded TM	$\alpha \rightarrow \gamma$ $ \alpha \leq \gamma $	PSPACE-complete	$a^n b^n c^n$
type-2	context free	pushdown automaton	$A \rightarrow \gamma$	cubic	$a^n b^n$
type-3	regular	NFA / DFA	$A \rightarrow a$ or $A \rightarrow aB$	linear time	$a^* b^*$

Figure 1: Chomsky Hierarchy

Figure 1, Lecture note on 15-411: Compiler Design, CMU, 2023

GRAMMAR AND MEANING OF THE LANGUAGE

Let us show that $L(G_{pal})$ is precisely the set of palindromes consisting of 0's and 1's.

- (\rightarrow) We want to show that if $S \Rightarrow^* w$, then w is a palindrome.
Induction on the length of derivation.

- 1 Base: if length at most 1, then $w = \epsilon, 0$ or 1 , which is trivially a palindrome.
- 2 Induction: if the derivation has length $n + 1$, then it is of the form

$$S \Rightarrow 0S0 \Rightarrow^* 0x0 = w$$

(or 0 replaced by 1) where $S \Rightarrow^* x$ is a derivation of length n . By I.H. x is a palindrome, and thus $0x0$ is.

GRAMMAR AND MEANING OF THE LANGUAGE

Let us show that $L(G_{pal})$ is precisely the set of palindromes consisting of 0's and 1's.

- (\rightarrow) We want to show that if $S \Rightarrow^* w$, then w is a palindrome.

Induction on the length of derivation.

- 1 Base: if length at most 1, then $w = \epsilon, 0$ or 1 , which is trivially a palindrome.
- 2 Induction: if the derivation has length $n + 1$, then it is of the form

$$S \Rightarrow 0S0 \Rightarrow^* 0x0 = w$$

(or 0 replaced by 1) where $S \Rightarrow^* x$ is a derivation of length n . By I.H. x is a palindrome, and thus $0x0$ is.

- (\leftarrow) We want to show that if w is a palindrome, then $S \Rightarrow^* w$.

Induction on $|w|$.

- 1 Base: $|w| \leq 1$, then $S \Rightarrow w$ by a single application of one of the rules $S \rightarrow \epsilon|0|1$.
- 2 Induction: if w is a palindrome of length ≥ 2 , it is of the form $0x0$ or $1x1$ for some palindrome x . By I.H., $S \Rightarrow^* x$. Therefore,

$$S \Rightarrow 0S0 \Rightarrow^* 0x0 = w$$

GRAMMAR AND DERIVATION

Consider CFG G with the following rules.

- $S \rightarrow XSX \mid R$
- $R \rightarrow aTb \mid bTa$
- $T \rightarrow XTX \mid X \mid \epsilon$
- $X \rightarrow a \mid b$

1 Variables? Terminals? Start Variable?

2 $T \Rightarrow^* T$?

3 $T \Rightarrow^* XXX$?

4 $XXX \Rightarrow^* aba$?

5 $R \Rightarrow^* \epsilon$?

6 $S \Rightarrow^* abaababbaaba$?

GRAMMAR AND MEANING OF THE LANGUAGE

Consider CFG G with the following rules. Describe the language $L(G)$ in plain English.

$$1 \quad S \rightarrow aSb \mid aA \mid bB$$

$$2 \quad A \rightarrow aA \mid \epsilon$$

$$3 \quad B \rightarrow bB \mid \epsilon$$

DESIGNING A CONTEXT-FREE GRAMMAR

- 1 $L = \{0^n 1^n : n \geq 1\}$.
- 2 $L = \{0^n 1^n : n \geq 1\} \cup \{0^n 1^n : n \geq 1\}$.
- 3 $L = \{a^i b^j : i > j\}$.
- 4 $L = \{a^i b^j c^k : i \neq j \text{ or } j \neq k\}$.
- 5 $L = \{\text{the set of all well-formed parentheses}\}$.
- 6 $L = \{\text{all strings with the same number of 0's and 1's}\}$.

HOW TO DESIGN CFG

The design of CFG requires some ingenuity. Some useful tips here.

- Many CFLs are the union of simpler CFLs.
- It is convenient to think of a variable as something that represents a set of strings; those which can be derived from that variable.
- A CFG for a regular language is easy to construct.
- Sometimes you use some nice combinatorial property of the language.

$$L = \{0^n 1^n : n \geq 1\}.$$

$$L = \{0^n 1^n : n \geq 1\} \cup \{1^n 0^n : n \geq 1\}.$$

$$L = \{a^i b^j : i > j\}.$$

$$L = \{a^i b^j c^k : i \neq j \textbf{ OR } j \neq k\}.$$

$L = \{\text{THE SET OF ALL WELL-FORMED PARENTHESES}\}.$

$L = \{\text{THE SAME \# OF 0'S AND 1'S}\}.$

LEFTMOST/RIGHTMOST DERIVATION

DEFINITION

A derivation is a leftmost derivation if a production rule is applied to the leftmost variable in each step. A rightmost derivation is defined similarly.

Example: a leftmost derivation of the string " $a \times (a + b00)$ " in the CFG G_{ari}

$$1 \quad E \rightarrow I \mid E + E \mid E \times E \mid (E)$$

$$2 \quad I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$$

GRAPHIC REPRESENTATION OF THE DERIVATION

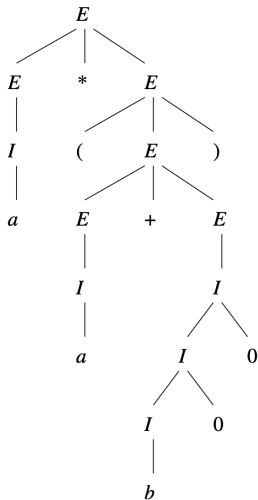


Figure 5.6, Hopcroft et. al. 2006

PARSE TREE

DEFINITION

Let $G = (V, \Sigma, R, S)$ be a context-free grammar. A **parse tree for the grammar G** is a (rooted) tree satisfying the following.

- 1 Each internal node is labelled by a variable in V .
- 2 Each leaf is labelled by a member of $V \cup \Sigma \cup \{\epsilon\}$. If a leaf is labelled by ϵ , it is the only child of its parent.
- 3 If an internal node is labelled by A , and its children are labelled by

$$X_1, \dots, X_k$$

when read from the left to right, then there is a rule $A \rightarrow X_1 \dots X_k$ in R .

YIELD OF A PARSE TREE

DEFINITION

Let $G = (V, \Sigma, R, S)$ be a context-free grammar. The yield of a parse tree is a string in $(V \cup \Sigma \cup \{\epsilon\})^*$ obtained by concatenating the labels on the leaves of the parse tree from left to right.

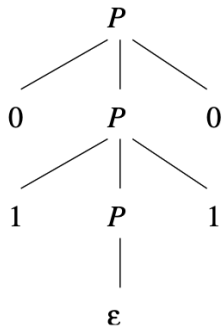


Figure 5.5, Hopcroft et. al. 2006

EQUIVALENCE OF PARSE TREE AND DERIVATION

THEOREM

Let $G = (V, \Sigma, R, S)$ be a context-free grammar. The following are equivalent.

- $S \Rightarrow^* w$ (i.e. $w \in L(G)$).
- *There is a parse tree with root S and yield w .*
- $S \Rightarrow_{lm}^* w$.
- $S \Rightarrow_{rm}^* w$.