

Context-free Grammars

Narges Khakpour

Department of Computer Science,
Linnaeus University

Table of Content

- 1 Introduction
- 2 Context-free Grammars
- 3 Context-free Languages
- 4 Parse Trees
- 5 Language Ambiguity
- 6 Conclusions

Recall...

- Alphabet, e.g. $\Sigma = \{a, b, c\}$
- Strings, e.g. *abbb*
- Σ^* all possible strings defined over Σ
- Languages, i.e. $L \subseteq \Sigma^*$, e.g. $L = \{a, ac, bbc\}$
- Regular expressions to represent languages (the operators of $+$ and $*$)
- Regular languages vs irregular languages
- Language equivalence of NFA, DFA and RE
- Pumping lemma

Introduction

- Regular languages are efficient but very limited in power
- Can we still come up with a recognizer to express non-regular languages?
- The topic of today's lecture.

Table of Content

- 1 Introduction
- 2 Context-free Grammars**
- 3 Context-free Languages
- 4 Parse Trees
- 5 Language Ambiguity
- 6 Conclusions

Context-free Grammars

Definition

A context-free grammar (CFG) G is a quadruple $\langle V, \Sigma, R, S \rangle$ where

- V : a set of non-terminal symbols
- Σ : a set of terminals ($V \cap \Sigma = \emptyset$)
- R : a set of production rules ($R : V \rightarrow (V \cup \Sigma)^*$)
- $S \in V$: a start symbol.

Context-free Grammars

Example

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \emptyset \end{aligned}$$

$$\text{or } S \rightarrow aSb \mid \emptyset$$

- $V = \{S\}$
- $\Sigma = \{a, b\}$
- $R = \{S \rightarrow aSb, S \rightarrow \emptyset\}$
- The left-hand-side of a production rule (e.g. S) is called head, and the right-hand-side (e.g. aSb) is called body
- S is the start symbol.
- Convention: We always use capital letters for non-terminals.

Context-free Grammars: Language

- Derivations, using productions from head to body.
- Let $A \rightarrow B$ be a production rule.
- If we apply this rule on a string xAy , then we obtain xBy . We say that xAy derivatives xBy and denote it by $xAy \rightarrow xBy$.
- In each step, we rewrite one symbol
- If $A_1 \rightarrow A_2 \rightarrow \dots \rightarrow A_n$, then we write $A_1 \rightarrow^* A_n$.

Example 1

$$S \rightarrow aSb \mid \emptyset$$

- Does S derivative a^3b^3 ?

$$\begin{aligned} S &\rightarrow aSb \\ &\rightarrow aaSbb \\ &\rightarrow aaaSbbb \\ &\rightarrow aaabbb \end{aligned}$$

Example 2

Let the grammar $G = (\{E, F\}, \{+, (,), *, a, b, 0, 1\}, P, E)$ that generates some arithmetic expressions where P is the following set of productions.

Derive $(b1) * a + ba1$.

$$1 : E \rightarrow F$$

$$2 : E \rightarrow E + E$$

$$3 : E \rightarrow E * E$$

$$4 : E \rightarrow (E)$$

$$5 : F \rightarrow a$$

$$6 : F \rightarrow b$$

$$7 : F \rightarrow Fa$$

$$8 : F \rightarrow Fb$$

$$9 : F \rightarrow F0$$

$$10 : F \rightarrow F1$$

Example 2

Let the grammar $G = (\{E, F\}, \{+, (,), *, a, b, 0, 1\}, P, E)$ that generates some arithmetic expressions where P is the following set of productions.

Derive $(b1) * a + ba1$.

	$E \rightarrow E * E$	(3)
1 : $E \rightarrow F$	$\rightarrow E * E + E$	(2)
2 : $E \rightarrow E + E$	$\rightarrow (E) * E + E$	(4)
3 : $E \rightarrow E * E$	$\rightarrow (F) * E + E$	(1)
4 : $E \rightarrow (E)$	$\rightarrow (F1) * E + E$	(10)
5 : $F \rightarrow a$	$\rightarrow (b1) * E + E$	(6)
6 : $F \rightarrow b$	$\rightarrow (b1) * F + E$	(1)
7 : $F \rightarrow Fa$	$\rightarrow (b1) * F + F$	(1)
8 : $F \rightarrow Fb$	$\rightarrow (b1) * a + F$	(5)
9 : $F \rightarrow F0$	$\rightarrow (b1) * a + F1$	(10)
10 : $F \rightarrow F1$	$\rightarrow (b1) * a + Fa1$	(7)
	$\rightarrow (b1) * a + ba1$	(6)

Example 2

Let the grammar $G = (\{E, F\}, \{+, (,), *, a, b, 0, 1\}, P, E)$ that generates some arithmetic expressions where P is the following set of productions.

Derive $(b1) * a + ba1$.

	$E \rightarrow E + E$	(2)
1 : $E \rightarrow F$	$\rightarrow E + F$	(5)
2 : $E \rightarrow E + E$	$\rightarrow E + F1$	(10)
3 : $E \rightarrow E * E$	$\rightarrow E + Fa1$	(7)
4 : $E \rightarrow (E)$	$\rightarrow E + ba1$	(6)
5 : $F \rightarrow a$	$\rightarrow E * E + ba1$	(2)
6 : $F \rightarrow b$	$\rightarrow (E) * E + ba1$	(4)
7 : $F \rightarrow Fa$	$\rightarrow (F) * E + ba1$	(1)
8 : $F \rightarrow Fb$	$\rightarrow (F1) * E + ba1$	(10)
9 : $F \rightarrow F0$	$\rightarrow (b1) * E + ba1$	(6)
10 : $F \rightarrow F1$	$\rightarrow (b1) * F + ba1$	(1)
	$\rightarrow (b1) * a + ba1$	(5)

Derivations

- At each step we might have several applicable rules,
- Not all choices lead to successful derivations of a particular string
- Leftmost derivation: always replace the leftmost variable by one of its rule-bodies.
- Rightmost derivation: always replace the rightmost variable by one of its rule

Table of Content

- 1 Introduction
- 2 Context-free Grammars
- 3 Context-free Languages**
- 4 Parse Trees
- 5 Language Ambiguity
- 6 Conclusions

Context-Free Language (CFL)

- A string s is generated by a grammar $G = \langle V, \Sigma, R, S \rangle$, if $S \rightarrow^* s$.
- The languages of G is $L(G) = \{s \in \Sigma^* | S \rightarrow^* s\}$.
- The languages generated by a context-free grammar is called a context-free language, or,
- A language L' is context-free, if and only if there is a grammar G such that $L(G) = L'$

Example

Example

- Consider the following grammar:

$$S \rightarrow abSc \mid c$$

- The derivations are of the following forms:

$$S \rightarrow c$$

or

$$S \rightarrow abSc \rightarrow ababScc \rightarrow abababSccc \rightarrow ababababSccccc \rightarrow \dots \rightarrow (ab)^n c(c)^n$$

- The language is strings with the pattern $(ab)^n c(c)^n$ where $n \geq 0$

$$L(G) = \{s \mid s = (ab)^n c(c)^n, n \geq 0\}$$

Example

Example

- Consider the following grammar:

$$\begin{aligned} S &\rightarrow aSc \mid A \\ A &\rightarrow bA \mid d \end{aligned}$$

- The derivations are of the following forms:

$$S \rightarrow A \rightarrow bA \rightarrow bbA \rightarrow \dots \rightarrow b^n A \rightarrow b^n d$$

$$\text{or} \quad S \rightarrow A \rightarrow d$$

$$\text{or} \quad S \rightarrow aSc \rightarrow aaSc^2 \rightarrow \dots \rightarrow a^m Sc^m \rightarrow a^m Ac^m \rightarrow a^m dc^m$$

$$\begin{aligned} \text{or} \quad S &\rightarrow aSc \rightarrow aaSc^2 \rightarrow \dots \rightarrow a^m Sc^m \rightarrow a^m Ac^m \rightarrow a^m bAc^m \\ &\dots \rightarrow a^m b^k Ac^m \rightarrow a^m b^k dc^m \end{aligned}$$

$$L(G) = \{s \mid s = a^n b^m dc^n, m, n \geq 0\}$$

Example 3

Example

- **BonusPoint** What is the language of the following grammar? S is the start symbol.

$$S \rightarrow A \mid cS \mid a \mid b$$

$$A \rightarrow aSb \mid baS \mid abS \mid Sab \mid bSa \mid Sba$$

where $\Sigma = \{a, b, c\}$

Table of Content

- 1 Introduction
- 2 Context-free Grammars
- 3 Context-free Languages
- 4 Parse Trees**
- 5 Language Ambiguity
- 6 Conclusions

Parse Tree: Example

$$1 : E \rightarrow F$$

$$2 : E \rightarrow E + E$$

$$3 : E \rightarrow E * E$$

$$4 : E \rightarrow (E)$$

$$5 : F \rightarrow a$$

$$6 : F \rightarrow b$$

$$7 : F \rightarrow Fa$$

$$8 : F \rightarrow Fb$$

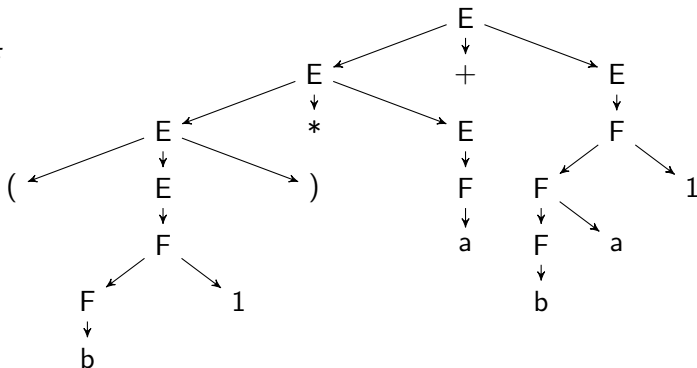
$$9 : F \rightarrow F0$$

$$10 : F \rightarrow F1$$

Parse Tree: Example

The parse tree for $(b1) * a + ba1$

- 1 : $E \rightarrow F$
- 2 : $E \rightarrow E + E$
- 3 : $E \rightarrow E * E$
- 4 : $E \rightarrow (E)$
- 5 : $F \rightarrow a$
- 6 : $F \rightarrow b$
- 7 : $F \rightarrow Fa$
- 8 : $F \rightarrow Fb$
- 9 : $F \rightarrow F0$
- 10 : $F \rightarrow F1$



Parse Tree

- Graphical representation of derivations
- Non-terminal are the interior nodes
- Terminals are the leaves
- A string s belongs to a language, if we can produce a parse tree for it with the root labeled by the start symbol
- The yield of a parse tree is the string of leaves from left to right or the concatenation of its children's yields from left to right.

From Parse Tree to the String

Yields of a parse tree

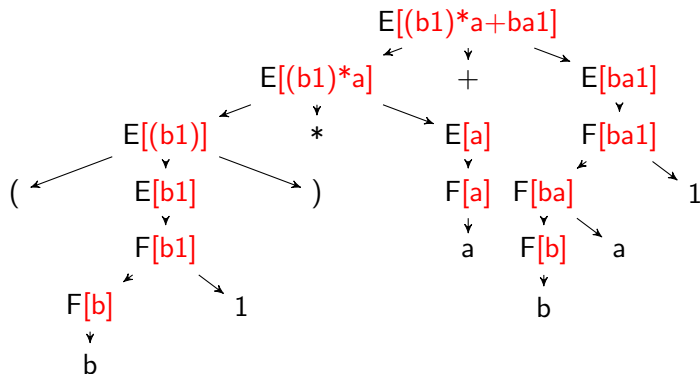


Table of Content

- 1 Introduction
- 2 Context-free Grammars
- 3 Context-free Languages
- 4 Parse Trees
- 5 Language Ambiguity**
- 6 Conclusions

Language Ambiguity

- Consider the following grammar with the start symbol S

$$S \rightarrow S + S \mid S * S \mid a \mid b$$

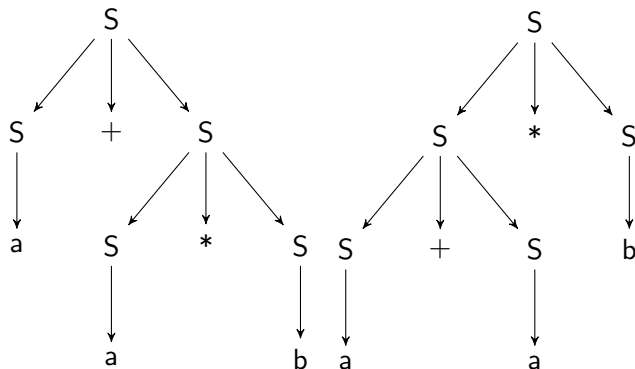
- Derivative the string $a + a * b$

-

$$S \xrightarrow{1} S + S \xrightarrow{3} a + S \xrightarrow{2} a + S * S \xrightarrow{3} a + a * S \xrightarrow{4} a + a * b$$

Language Ambiguity

$$S \rightarrow S + S \mid S * S \mid a \mid b$$



Language Ambiguity

- A language is ambiguous, if it has a string with two different parse trees.
- **BonusPoint** Is the following statement correct?
"if a language contains a string that can be obtained using two different derivations, it's ambiguous"
- No general algorithm to check ambiguity
- No general algorithm to resolve ambiguity
- An unambiguous grammar has unique left-most derivations
- There are some methods to remove ambiguity in **some** languages.

Removing Language Ambiguity: Left-factoring

- $$A \rightarrow \alpha\beta_1|\alpha\beta_2|\dots|\alpha\beta_n|\omega_1|\dots|\omega_m$$
- ω is those rh's that do not start with α
- Rewrite it as the following where A' is non-terminal

$$A \rightarrow \alpha A'|\omega_1|\dots|\omega_m$$

$$A' \rightarrow \beta_1|\beta_2|\dots|\beta_n$$

Removing Language Ambiguity: Left-factoring

Example

$$S \rightarrow \text{if } E \text{ then } S' \text{ else } S'$$

$$S \rightarrow \text{if } E \text{ then } S'$$

can be rewritten as

$$S \rightarrow \text{if } E \text{ then } S' S''$$

$$S'' \rightarrow \text{else } S' \mid \epsilon$$

Removing Language Ambiguity: Left-Recursion

-
- Rewrite it as the following where A' is non-terminal

$$A \rightarrow A\alpha \mid \omega_1 \mid \dots \mid \omega_m$$

$$A \rightarrow \omega_1 A' \mid \dots \mid \omega_m A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

Example

$$E \rightarrow x \mid y \mid z \mid (E) \mid E + E$$

$x + y + z$ has more than one parse tree

After removing left-recursion

$$E \rightarrow xE' \mid yE' \mid zE' \mid (E)E'$$

$$E' \rightarrow +EE' \mid \epsilon$$

Removing Language Ambiguity: Left-Recursion

BonusPoint Is the following grammar ambiguous? Why? If yes, remove ambiguity.

$$A \rightarrow Bw \mid x$$

$$B \rightarrow Ct \mid zC$$

$$C \rightarrow Ax A \mid z$$

Removing Language Ambiguity: Indirect Left-Recursion

- Indirect recursion, e.g.

$$A \rightarrow Bw \mid x$$

$$B \rightarrow Ct \mid zC$$

$$C \rightarrow Ax A \mid z$$

after rewriting and removing C

$$B \rightarrow AxAt \mid zt \mid zAx A \mid zz$$

after rewriting and removing B

$$A \rightarrow AxAtw \mid ztw \mid zAxAw \mid zzw \mid x$$

that has left recursion!

Removing Language Ambiguity: Indirect Left-Recursion

After removing left-recursion

$$A \rightarrow zt w A' \mid z A x A w A' \mid z z w A' \mid x A'$$

$$A' \rightarrow x A t w A' \mid \epsilon$$

Needs left-factoring z in the next step

$$A \rightarrow z Z \mid x A'$$

$$Z \rightarrow t w A' \mid A x A w A' \mid z w A'$$

$$A' \rightarrow x A t w A' \mid \epsilon$$

Needs one further step rewriting to left-factorize z in the rules for Z (via A and the third rule)....

$$Z \rightarrow t w A' \mid z Z x A w A' \mid x A' x A w A' \mid z w A'$$

to

$$Z \rightarrow t w A' \mid z Z' \mid x A' x A w A'$$

$$Z' \rightarrow Z x A w A' \mid w A'$$

Removing Language Ambiguity: Indirect Left-Recursion

Final grammar

$$A \rightarrow zZ \mid xA'$$

$$A' \rightarrow xAtwA' \mid \epsilon$$

$$Z \rightarrow twA' \mid zZ' \mid xA'xAwA'$$

$$Z' \rightarrow ZxAwA' \mid wA'$$

Table of Content

- 1 Introduction
- 2 Context-free Grammars
- 3 Context-free Languages
- 4 Parse Trees
- 5 Language Ambiguity
- 6 Conclusions**

Conclusions

- Context-free grammars and languages
- Parse trees and language ambiguity
- Next lecture is on CFGs vs REs and parsing