

CS112: Theory of Computation (LFA)

Lecture9: Context-free Grammars

Dumitru Bogdan

Faculty of Computer Science
University of Bucharest

April 20, 2022

Table of contents

1. Previously on CS112
2. Context setup
3. Ambiguity
4. Chomsky normal form
5. Examples
6. Pushdown Automaton

Section 1

Previously on CS112

Intuition

- A grammar consists of a collection of **substitution rules**, also called **productions**
- Each rule appears as a line in the grammar, comprising a symbol and a string separated by an arrow
- The symbol is called a **variable**. The string consists of variables and other symbols called **terminals**
- The variable symbols often are represented by capital letters
- The terminals are analogous to the input alphabet and often are represented by lowercase letters, numbers, or special symbols
- One variable is designated as the start variable
- It usually occurs on the left-hand side of the topmost rule

Example

For example, grammar G_1 contains three rules. G_1 's variables are A and B , where A is the start variable. Its terminals are 0, 1, and $\#$

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

Intuition

You use a grammar to describe a language by generating each string of that language in the following manner:

1. Write down the start variable. It is the variable on the left-hand side of the top rule, unless specified otherwise
2. Find a variable that is written down and a rule that starts with that variable. Replace the written down variable with the right-hand side of that rule
3. Repeat step 2 until no variables remain

Example

For example, grammar G_1 generates the string $000\#111$. The sequence of substitutions to obtain a string is called a **derivation**. A derivation of string $000\#111$ in grammar G_1 is:

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

Example

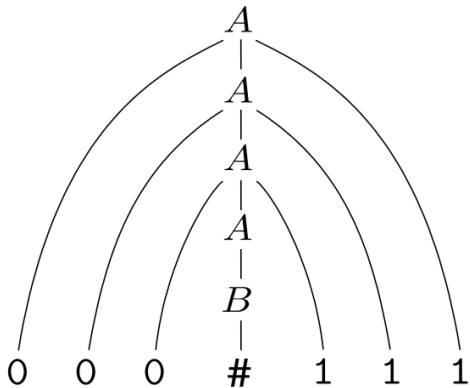


Figure: Parse tree for 000#111 in grammar G_1

Section 2

Context setup

Context setup

Corresponding to Sipser 2.1

Section 3

Ambiguity

Ambiguity

- Sometimes a grammar can generate the same string in several different ways
- Such a string will have several different parse trees and thus several different meanings
- This result may be undesirable for certain applications, such as programming languages, where a program should have a unique interpretation

Ambiguity

- If a grammar generates the same string in several different ways, we say that the string is derived *ambiguously* in that grammar
- If a grammar generates some string ambiguously, we say that the grammar is *ambiguous*
- If we consider grammar G_5

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle \mid (\langle \text{EXPR} \rangle) \mid a$$

we observe that it generates the string $a + a \times a$ ambiguously

Ambiguity

Grammar G_5 doesn't capture the usual precedence relations and so may group the $+$ before the \times or vice versa.

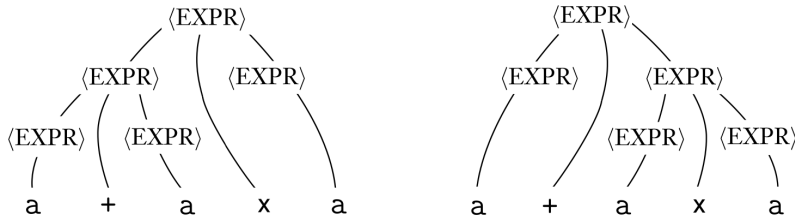


Figure: The two parse trees for the string $a + a \times a$ in grammar G_5

Ambiguity

- Grammar G_2 (see previous slides) is another example of an ambiguous grammar
- The sentence *the girl touches the boy with the flower* has two different derivations
- Give the two parse trees and observe their correspondence with the two different ways to read that sentence (\Leftarrow get a CS112 T-Shirt)

Ambiguity

Now we formalize the notion of ambiguity

- When we say that a grammar generates a string ambiguously, we mean that the string has two different parse trees, not two different derivations
- Two derivations may differ merely in the order in which they replace variables yet not in their overall structure
- A derivation of a string w in a grammar G is a **leftmost derivation** if at every step the leftmost remaining variable is the one replaced

Ambiguity

Definition

A string w is derived **ambiguously** in context-free grammar G if it has two or more different leftmost derivations. Grammar G is **ambiguous** if it generates some string ambiguously

- Sometimes when we have an ambiguous grammar we can find an unambiguous grammar that generates the same language
- Some context-free languages, however, can be generated only by ambiguous grammars. Such languages are called **inherently ambiguous**
- Study why $\{a^i b^j c^k \mid i = j \text{ or } j = k\}$ is inherently ambiguous (\Leftarrow and get a CS112 T-Shirt)

Section 4

Chomsky normal form

Noam Chomsky

Avram Noam Chomsky is an American linguist, **philosopher**, **political activist**, author, and lecturer. He is an Institute Professor and professor emeritus of linguistics at the Massachusetts Institute of Technology

- According to The New York Times, Noam Chomsky is **arguably the most important intellectual alive**
- https://www.goodreads.com/author/list/2476.Noam_Chomsky
- If in doubt, start with this book: **How the World Works**

Chomsky normal form

- Working with context-free grammars, it is often convenient to have them in simplified form
- The simplest and most useful forms is called **the Chomsky normal form**
- Chomsky normal form is useful in giving algorithms for working with context-free grammars

Formal definition

Definition

A context-free grammar is in Chomsky normal form if every rule is of the form

$$A \rightarrow BC$$

$$A \rightarrow a$$

where a is any terminal and A , B , and C are any variables — except that B and C may not be the start variable. In addition, we allow the rule $S \rightarrow \epsilon$, where S is the start variable

Converting to Chomsky normal form

Theorem

Any context-free language is generated by a context-free grammar in Chomsky normal form

- Proof idea is to convert any grammar G into Chomsky normal form
- The conversion has several stages wherein rules that violate the conditions are replaced with equivalent ones that are satisfactory

Proof idea

- First, we add a new start variable
- Then, we eliminate all ϵ -rules of the form $A \rightarrow \epsilon$
- We also eliminate all unit rules of the form $A \rightarrow B$
- In both cases we patch up the grammar to be sure that it still generates the same language
- Finally, we convert the remaining rules into the proper form

Formal proof

Proof.

First, we add a new start variable S_0 and the rule $S_0 \rightarrow S$, where S was the original start variable. This change guarantees that the start variable doesn't occur on the right-hand side of a rule

Second, we take care of all ϵ -rules. We remove an ϵ -rule $A \rightarrow \epsilon$, where A is not the start variable. Then for each occurrence of an A on the right-hand side of a rule, we add a new rule with that occurrence deleted. In other words, if $R \rightarrow uAv$ is a rule in which u and v are strings of variables and terminals, we add rule $R \rightarrow uv$. We do so for each occurrence of an A , so the rule $R \rightarrow uAvAw$ causes us to add $R \rightarrow uvAw$, $R \rightarrow uAvw$ and $R \rightarrow uvw$. If we have the rule $R \rightarrow A$, we add $R \rightarrow \epsilon$ unless we had previously removed the rule $R \rightarrow \epsilon$. We repeat these steps until we eliminate all ϵ -rules not involving the start variable.

Formal proof

Proof.

Third, we handle all unit rules. We remove a unit rule $A \rightarrow B$. Then, whenever a rule $B \rightarrow u$ appears, we add the rule $A \rightarrow u$ unless this was a unit rule previously removed. As before, u is a string of variables and terminals. We repeat these steps until we eliminate all unit rules. Finally, we convert all remaining rules into the proper form. We replace each rule $A \rightarrow u_1 u_2 \dots u_k$, where $k \geq 3$ and each u_i is a variable or terminal symbol, with the rules $A \rightarrow u_1 A_1, A_1 \rightarrow u_2 A_2, A_2 \rightarrow u_3 A_3, \dots, A_{k-2} \rightarrow u_{k-1} u_k$. The A_i 's are new variables. We replace any terminal u_i in the preceding rule(s) with the new variable U_i and add the rule $U_i \rightarrow u_i$ □

Section 5

Examples

Example

- Let G_6 be the following CFG and convert it to Chomsky normal form by using the conversion procedure just given.
- The series of grammars presented illustrates the steps in the conversion.

$$S \rightarrow ASA|aB$$

$$A \rightarrow B|S$$

$$B \rightarrow b|\epsilon$$

Example

1. The result of applying the first step to make a new start variable appears on the right

$$S \rightarrow ASA|aB$$

$$A \rightarrow B|S$$

$$B \rightarrow b|\epsilon$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA|aB$$

$$A \rightarrow B|S$$

$$B \rightarrow b|\epsilon$$

2. Remove ϵ -rules $B \rightarrow \epsilon$, shown on the left and $A \rightarrow \epsilon$, shown on the right

$$S_0 \rightarrow S$$

$$S \rightarrow ASA|aB|a$$

$$A \rightarrow B|S|\epsilon$$

$$B \rightarrow b$$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA|aB|a|SA|AS|S$$

$$A \rightarrow B|S$$

$$B \rightarrow b$$

Example

3. Remove unit rules $S \rightarrow S$, shown on the left and $S_0 \rightarrow S$, shown on the right

$$S_0 \rightarrow S$$

$$S \rightarrow ASA|aB|a|SA|AS$$

$$A \rightarrow B|S$$

$$B \rightarrow b$$

$$S_0 \rightarrow ASA|aB|a|SA|AS$$

$$S \rightarrow ASA|aB|a|SA|AS$$

$$A \rightarrow B|S$$

$$B \rightarrow b$$

4. Remove unit rules $A \rightarrow B$ and $A \rightarrow S$

$$S_0 \rightarrow ASA|aB|a|SA|AS$$

$$S \rightarrow ASA|aB|a|SA|AS$$

$$A \rightarrow S|b$$

$$B \rightarrow b$$

$$S_0 \rightarrow ASA|aB|a|SA|AS$$

$$S \rightarrow ASA|aB|a|SA|AS$$

$$A \rightarrow b|ASA|aB|a|SA|AS$$

$$B \rightarrow b$$

Example

5. Convert the remaining rules into the proper form by adding additional variables and rules. The final grammar in Chomsky normal form is equivalent to G_6 . We simplified the resulting grammar by using a single variable U and rule $U \rightarrow a$

$$S_0 \rightarrow AA_1 | UB | a | SA | AS$$

$$S \rightarrow AA_1 | UB | a | SA | AS$$

$$A \rightarrow b | AA_1 | UB | a | SA | AS$$

$$A_1 \rightarrow SA$$

$$U \rightarrow a$$

$$B \rightarrow b$$

Section 6

Pushdown Automaton

Schematics

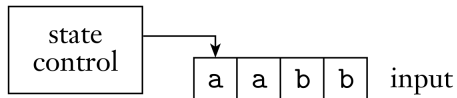


Figure: Schematic of a finite automaton

- The left figure is a schematic representation of a finite automaton
- The control represents the states and transition function
- The tape contains the input string, and the arrow represents the input head, pointing at the next input symbol to be read

Schematics

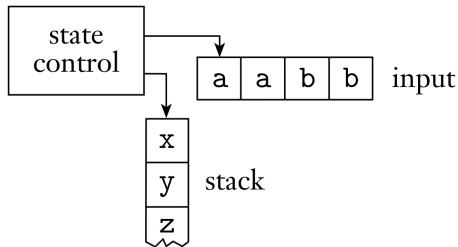


Figure: Schematic of a pushdown automaton

- A **pushdown automaton** (PDA) can write symbols on the stack and read them back later
- Writing a symbol "pushes down" all the other symbols on the stack.
- At any time the symbol on the top of the stack can be read and removed. The remaining symbols then move back up
- Writing a symbol on the stack is often referred to as pushing the symbol, and removing a symbol is referred to as popping it

Schematics

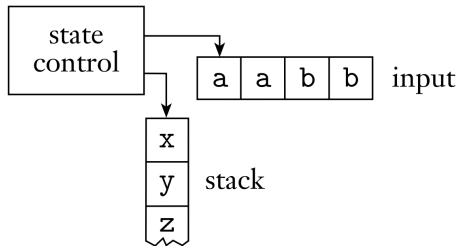


Figure: Schematic of a pushdown automaton

- Note that all access to the stack, for both reading and writing, may be done only at the top. In other words a stack is a "last in, first out" storage device
- If certain information is written on the stack and additional information is written afterward, the earlier information becomes inaccessible until the later information is removed
- A stack is valuable because it can hold an unlimited amount of information

Example

Let's take the language $\{0^n 1^n | n \geq 0\}$. While a DFA is unable to recognize, a PDA is able to recognize this language because it can use its stack to store the number of 0s it has seen. Thus the unlimited nature of a stack allows the PDA to store numbers of unbounded size, This is how:

- Read symbols from the input. As each 0 is read, push it onto the stack
- As soon as 1s are seen, pop a 0 off the stack for each 1 read
- If reading the input is finished exactly when the stack becomes empty of 0s, accept the input
- If the stack becomes empty while 1s remain or if the 1s are finished while the stack still contains 0s or if any 0s appear in the input following 1s, reject the input

Nondeterminism

- PDA may be nondeterministic. Deterministic and nondeterministic PDA are *not* equivalent in power
- Nondeterministic PDA recognize certain languages that no deterministic PDA can recognize
- Recall that deterministic and nondeterministic finite automata do recognize the same class of languages, so the PDA situation is different
- We will focus on nondeterministic PDA because these automata are equivalent in power to context-free grammars

Transition function

- Formal definition of a PDA is similar to that of a finite automaton, except for the stack
- The stack is a device containing symbols drawn from some alphabet
- The machine may use different alphabets for its input and its stack, so now we specify both an input alphabet Σ and a stack alphabet Γ
- Recall that $\Sigma_\epsilon = \Sigma \cup \{\epsilon\}$ and $\Gamma_\epsilon = \Gamma \cup \{\epsilon\}$
- Core aspect of any formal definition of an automaton is the transition function, which describes its behavior
- The domain of the transition function is $Q \times \Sigma_\epsilon \times \Gamma_\epsilon$
- Thus the current state, next input symbol read, and top symbol of the stack determine the next move of a PDA
- Either symbol may be ϵ , causing the machine to move without reading a symbol from the input or without reading a symbol from the stack

Transition function

- We need to consider what to allow the automaton to do when it is in a particular situation
- It may enter some new state and possibly write a symbol on the top of the stack
- The function δ can indicate this action by returning a member of Q together with a member of Γ_ϵ , that is, a member of $Q \times \Gamma_\epsilon$
- Because we allow nondeterminism in this model, a situation may have several legal next moves
- The transition function incorporates nondeterminism in the usual way, by returning a set of members of $Q \times \Gamma_\epsilon$, that is a member of $\mathcal{P}(Q \times \Gamma_\epsilon)$
- Gluing everything together we get: $\delta : Q \times \Sigma_\epsilon \times \Gamma_\epsilon \rightarrow \mathcal{P}(Q \times \Gamma_\epsilon)$

Formal definition

Definition

A pushdown automaton is a 6-tuple $(Q, \Sigma, \Gamma, \delta, q_0, F)$ where Q, Σ, Γ and F are all finite sets, and

1. Q is the set of states
2. Σ is the input alphabet
3. Γ is the stack alphabet
4. $\delta : Q \times \Sigma_{\epsilon} \times \Gamma_{\epsilon} \rightarrow \mathcal{P}(Q \times \Gamma_{\epsilon})$ is the transition function
5. $q_0 \in Q$ is the start state
6. $F \subseteq Q$ is the set of accept states

PDA computation

A PDA $M = (Q, \Sigma, \Gamma, \delta, \delta_0, F)$ computes as follows. It accepts input w if w can be written as $w = w_1 w_2 \dots w_m$ where each $w_i \in \Sigma_\epsilon$ and sequences of states $r_0, r_1, \dots, r_m \in Q$ and strings $s_0, s_1, \dots, s_m \in \Gamma^*$ exist that satisfy the following three conditions:

1. $r_0 = q_0$ and $s_0 = \epsilon$. This condition signifies that M starts out properly, in the start state and with an empty stack
2. For $i = 0, \dots, m - 1$ we have $(r_{i+1}, b) \in \delta(r_i, w_{i+1}, a)$ where $s_i = at$ and $s_{i+1} = bt$ for some $a, b \in \Gamma_\epsilon$ and $t \in \Gamma^*$. This condition states that M moves properly according to the state, stack, and next input symbol
3. $r_m \in F$. This condition states that an accept state occurs at the input end

The strings s_i represent the sequence of stack contents that M has on the accepting branch of the computation