

Context-Free Grammars

Chanathip Namprempre

Computer Science
Reed College

Outline

- 1 Context-Free Grammars
 - Introduction
 - Examples
 - Formal definition
- 2 Designing CFGs
- 3 Ambiguity
 - Introduction
 - Definition
- 4 Chomsky Normal Form
 - Definition
 - Examples
- 5 Closure properties

Outline

1 Context-Free Grammars

- Introduction
- Examples
- Formal definition

2 Designing CFGs

3 Ambiguity

- Introduction
- Definition

4 Chomsky Normal Form

- Definition
- Examples

5 Closure properties

- We learned that some languages are **not regular**. Finite automata cannot handle them.
- A more powerful method to describe some of these languages is to use **context-free grammars**.
- When you learn a new computer language, you must learn the **syntax** for that language (e.g. `printf('%d\n', x);`). The syntax is often specified in the form of grammars.
- A compiler must use grammars to **parse** a given program to figure out what the program means so that it can generate machine code for the program.
- The languages associated with context-free grammars are called **context-free languages**. They include regular languages and many more.

Example of a Context-Free Grammar

This is what a grammar looks like:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

The components are

- substitution rules, aka productions
- variable
- terminals
- start variable

Example of a Derivation

This is what a grammar looks like:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

This grammar generates the string 000#111. The sequence of substitutions to obtain this string is called a [derivation](#).

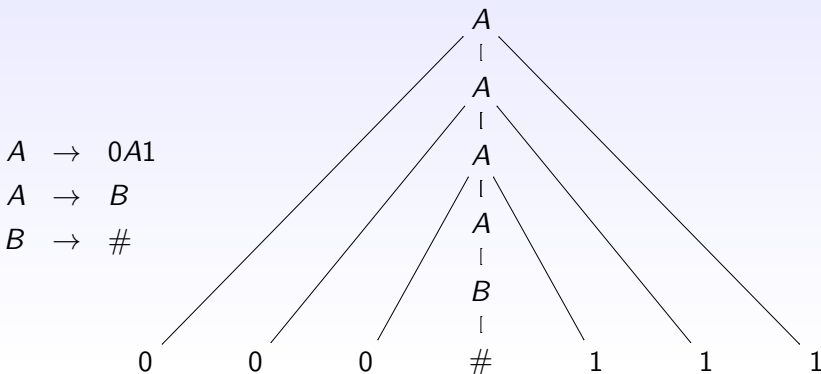
$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000\#111$$

Derivation and parse tree

The derivation

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

for 000#111 can be represented using a [parse tree](#).



Language of a Grammar

All strings generated from a grammar is called the **language** of that grammar.

Let G_1 be the grammar:

$$A \rightarrow 0A1 \mid B$$

$$B \rightarrow \#$$

Then, the language of G_1 , denoted $L(G_1)$, is $\{0^n \# 1^n \mid n \geq 0\}$.

Any language that can be generated by some context-free grammar is called a **context-free language** (CFL).

More example of a Context-Free Grammar

The following grammar G_2 describes a fragment of the English language.

$\langle \text{sentence} \rangle$	\rightarrow	$\langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle$
$\langle \text{noun-phrase} \rangle$	\rightarrow	$\langle \text{cmplx-noun} \rangle \mid \langle \text{cmplx-noun} \rangle \langle \text{prep-phrase} \rangle$
$\langle \text{verb-phrase} \rangle$	\rightarrow	$\langle \text{cmplx-verb} \rangle \mid \langle \text{cmplx-verb} \rangle \langle \text{prep-phrase} \rangle$
$\langle \text{prep-phrase} \rangle$	\rightarrow	$\langle \text{prep} \rangle \langle \text{cmplx-noun} \rangle$
$\langle \text{cmplx-noun} \rangle$	\rightarrow	$\langle \text{article} \rangle \langle \text{noun} \rangle$
$\langle \text{cmplx-verb} \rangle$	\rightarrow	$\langle \text{verb} \rangle \mid \langle \text{verb} \rangle \langle \text{noun-phrase} \rangle$
$\langle \text{article} \rangle$	\rightarrow	$\text{a} \mid \text{the}$
$\langle \text{noun} \rangle$	\rightarrow	$\text{cat} \mid \text{dog} \mid \text{paw}$
$\langle \text{verb} \rangle$	\rightarrow	$\text{touches} \mid \text{likes} \mid \text{sees}$
$\langle \text{prep} \rangle$	\rightarrow	with

More example of a Context-Free Grammar

An example of a derivation for G_2 on the string “a cat sees”.

$\langle \text{sentence} \rangle \Rightarrow \langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle$
 $\Rightarrow \langle \text{cmplx-noun} \rangle \langle \text{verb-phrase} \rangle$
 $\Rightarrow \langle \text{article} \rangle \langle \text{noun} \rangle \langle \text{verb-phrase} \rangle$
 $\Rightarrow a \langle \text{noun} \rangle \langle \text{verb-phrase} \rangle$
 $\Rightarrow a \text{ cat} \langle \text{verb-phrase} \rangle$
 $\Rightarrow a \text{ cat} \langle \text{cmplx-verb} \rangle$
 $\Rightarrow a \text{ cat} \langle \text{verb} \rangle$
 $\Rightarrow a \text{ cat sees}$

More examples of a Context-Free Grammar

The following grammar G_3 describes a language of balanced parentheses.

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

Try deriving $()(())((()))$.

The following grammar G_4 describes a language of mathematical expressions.

$$\begin{aligned} \langle \text{expr} \rangle &\rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle \\ \langle \text{term} \rangle &\rightarrow \langle \text{term} \rangle \times \langle \text{factor} \rangle \mid \langle \text{factor} \rangle \\ \langle \text{factor} \rangle &\rightarrow (\langle \text{expr} \rangle) \mid a \end{aligned}$$

Try deriving $a + a \times a$ and $(a + a) \times a$.

Formal definition of a context-free grammar

Definition

A **context-free grammar** is a 4-tuple (V, Σ, R, S) where

- 1 V is a finite set called the **variables**,
- 2 Σ is a finite set, disjoint from V , called the **terminals**,
- 3 R is a finite set of **rules**, with each rule being a variable or a string of variables and terminals, and
- 4 $S \in V$ is the **start variable**.

Example

How would you write the grammar G_1 formally?

$$\begin{aligned}A &\rightarrow 0A1 \mid B \\ B &\rightarrow \#\end{aligned}$$

$G_1 = (V, \Sigma, R, S)$ where

1. [Variables] $V = \{A, B\}$
2. [Terminals] $\Sigma = \{0, 1, \#\}$
3. [Rules] $R = \{ A \rightarrow 0A1 \mid B, B \rightarrow \# \}$
4. [Start variable] $S = A$

Outline

1 Context-Free Grammars

- Introduction
- Examples
- Formal definition

2 Designing CFGs

3 Ambiguity

- Introduction
- Definition

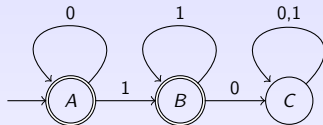
4 Chomsky Normal Form

- Definition
- Examples

5 Closure properties

Let's design CFGs

A grammar for 0^*1^* ?



$$A \rightarrow 0A \mid 1B \mid \varepsilon$$

$$B \rightarrow 1B \mid 0C \mid \varepsilon$$

$$C \rightarrow 0C \mid 1C$$

Idea

- For every transition $\delta(q_i, a) = q_j$, add rule $q_i \rightarrow aq_j$.
- For every accept state q , add rule $q \rightarrow \varepsilon$.

Let's design CFGs

- ❶ the language of balanced parentheses
- ❷ Σ^*
- ❸ $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$
- ❹ $\{w \mid w \text{ has an even number of 1s}\}$
- ❺ $\{w \mid w \text{ contains at least three 1s}\}$
- ❻ $\{w \mid w \text{ starts and ends with the same symbol}\}$
- ❼ $\{w \mid \text{the length of } w \text{ is odd}\}$
- ❽ $\{w \mid \text{the length of } w \text{ is odd and its middle symbol is 0}\}$
- ❾ $\{w \mid w = w^{\mathcal{R}}, \text{ that is, } w \text{ is a palindrome}\}$
- ❿ $\{\}$

Outline

1 Context-Free Grammars

- Introduction
- Examples
- Formal definition

2 Designing CFGs

3 Ambiguity

- Introduction
- Definition

4 Chomsky Normal Form

- Definition
- Examples

5 Closure properties

Ambiguous Grammars

Some grammars are **ambiguous**. This means that there is some string that the grammar generates ambiguously. For example, consider the following grammar G_5 :

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

This grammar generates the string $a + a \times a$ ambiguously. This means that there are at least **two ways to parse** the string using the given grammar.

G_5 generates exactly the same language as that generated by G_4 but G_4 is unambiguous.

Ambiguity Defined

Definition

A string w is derived ambiguously in a CFG G if it has two or more different **leftmost derivations**. Grammar G is **ambiguous** if it generates some strings ambiguously.

G_2 is also ambiguous. Try constructing the parse trees for the sentence “the girl touches the boy with the flower”.

Outline

1 Context-Free Grammars

- Introduction
- Examples
- Formal definition

2 Designing CFGs

3 Ambiguity

- Introduction
- Definition

4 Chomsky Normal Form

- Definition
- Examples

5 Closure properties

Chomsky Normal Form

Definition

A CFG is in **Chomsky Normal Form** if every rule is of the form

$$A \rightarrow BC$$

$$A \rightarrow a$$

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

Theorem

Any CFL is generated by a CFG in Chomsky Normal Form.

Converting CFG into Chomsky Normal Form

- ➊ Add new start variable S_0 and add a rule $S_0 \rightarrow S$ where S is the start variable.
- ➋ Remove all ε -rules $A \rightarrow \varepsilon$ not involving the start variable.
- ➌ Remove all unit rules $A \rightarrow B$.
- ➍ Convert all rules into the two proper forms.

Converting any CFGs into Chomsky Normal Form

1

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

2

$$A \rightarrow BAB \mid B \mid \varepsilon$$

$$B \rightarrow 00 \mid \varepsilon$$

Outline

1 Context-Free Grammars

- Introduction
- Examples
- Formal definition

2 Designing CFGs

3 Ambiguity

- Introduction
- Definition

4 Chomsky Normal Form

- Definition
- Examples

5 Closure properties

Closure properties

The class of context-free languages are closed under **union**, **reversal**, **concatenation**, and **star**.

The class of context-free languages are not closed under **complementation** and **intersection**.

[Consider the languages $\{a^m b^m c^n \mid m, n \geq 0\}$, $\{a^m b^n c^n \mid m, n \geq 0\}$, $\{a^n b^n c^n \mid n \geq 0\}$]