# INT201 Decision, Computation and Language

Lecture 6 – Context-Free Languages (1)

Dr Yushi Li

Xi'an Jiaotong-Liverpool University
西交利物浦大学

---

## Context-Free Languages

- Context-Free Grammar (CFG)

- Chomsky Normal Form (CNF)

---

## Context-Free Languages

- **Finite automata** accept precisely the strings in the language.

  *Perform a computation to determine whether a specific string is in the language.*

- **Regular expressions** describe precisely the strings in the language

  *Describe the general shape of all strings in the language.*

- **Context-free grammar (CFG)** is an entirely different formalism for defining a class of languages.

  *Give a procedure for listing off all strings in the language.*

---

## Context-Free Languages   上下文无关语言

Applications of CFG

- Programming languages: CFGs are used to define the syntax of programming languages, allowing parsers to analyze code structure.

- NLP: CFGs help in parsing sentences, enabling applications like machine translation and speech recognition

- Compilers: CFGs facilitate syntax analysis, ensuring that the source code adheres to the language's grammatical rules.

## Context-Free Grammar

**Example**

- Start variable S with rules:

  $S \to AB$
  $A \to a$
  $A \to aA$
  $B \to b$
  $B \to bB$

$$L = \{a^m | b^m : m \geq 1\}$$

variables: $S, A, B$   terminals: a, b

- Following these rules, we can yield ?

we can infer a language from given rule

$S \Rightarrow AB \Rightarrow aAB \Rightarrow aAbB \Rightarrow aaAbB \Rightarrow$
$aaaAbB \Rightarrow aaaAbbB \Rightarrow aaaa\, bbb.$
aaa bb
aa bbb...

---

## Context-Free Grammar

**Definition**

A context-free grammar is a 4-tuple $G = (V, \Sigma, R, S)$, where

有限集合    变量
1. V is a finite set, whose elements are called **variables**,
有限集合    末端
2. $\Sigma$ is a finite set, whose elements are called **terminals**, (注意和 DFA/NFA 的 $\Sigma$ 区别)
3. $V \cap \Sigma = \emptyset$,    variable $\cap$ terminal   随元素揽
4. S is an element of V ; it is called the **start variable**, 开始
5. R is a finite set, whose elements are called **rules**. Each rule has the form $A \to w$,
where $A \in V$ and $w \in (V \cup \Sigma)^*$.

A is a variable in V     w is the strings constructed from $(V \cup \Sigma)^*$

---

## Context-Free Grammar

**Example**

Language $L = \{0^k 1^k : k \geq 0\}$ has CFG $G = (V, \Sigma, R, S)$,

variable set $V = \{s\}$

Terminal set $\Sigma = \{0, 1\}$

start variable $S$

Rule set R:   $S \to 0S1$            $0^k 1^k$
                   $S \to \epsilon$          11

0S1
$S \to 0\,\textcircled{S}\,1 \to 00S11 \to 000S111 \Rightarrow 0\cdots0S1\cdots1$

---

## Deriving strings and languages using CFG

$\Rightarrow$ : **yeild**    产出

Let $G = (V, \Sigma, R, S)$ be a context free grammar with

- $A \in V$
- $u, v, w \in (V \cup \Sigma)^*$,
- $A \to w$ is a rule of the grammar

The string uwv can be derived in one step from the string uAv, written as

$$uAv \Rightarrow uwv$$

**Example:** $aaAbb \Rightarrow aaaAbb$

## Deriving strings and languages using CFG

⇒* : **derive**    右由左得到

Let G = (V, Σ, R, S) be a context free grammar with

- u, v ∈ (V ∪ Σ)*

得到

The string v can be derived from the string u, written as u ⇒* v, if one of the following conditions holds:

1. u = v

2. there exist an integer $k \geq 2$ and a sequence $u_1, u_2, \ldots, u_k$ of strings in (V ∪ Σ)*, such that

(a) $u = u_1$,

$u_1, u_2, \ldots u_k \in (V \cup \Sigma)^*$

(b) $v = u_k$, and $u_1 \Rightarrow u_2 \Rightarrow \ldots \Rightarrow u_k$.

**Example:** With the rules A → B1 | D0C

$$0AA \stackrel{*}{\Rightarrow} 0D0CB1$$

---

## Language of CFG

**Definition**

The language of CFG G = (V, Σ, R, S) is

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

Such a language is called **context-free**, and satisfies $L(G) \subseteq \Sigma^*$.

**Example**

CFG G = (V, Σ, R, S) with

1. V = {S}

2. Σ = {0, 1}

3. Rules R: S → 0S | ε

L(G) = ?

$$S \to 0S \to 00S \to \cdots$$
$$\to 0 \cdots 0S$$
$$\therefore S \to \varepsilon$$
$$\therefore S \to 0 \cdots 0 \implies L(G) = \{0^n : n \geq 0\}$$

---

**Example (Palindrome)**  回文

CFG G = (V, Σ, R, S) with

1. V = {S}

2. Σ = {a, b}

3. Rules R: S → aSa | bSb | a | b | ε

$$\begin{cases} S \to aSa \\ S \to bSb \\ S \to a \\ S \to b \\ S \to \varepsilon \end{cases}$$

Language of this CFG ?

$$S \Rightarrow aSa \Rightarrow aaSaa \stackrel{*}{\Rightarrow} a\cdots aSa\cdots a$$

$$\Rightarrow \begin{cases} a\cdots aaa\cdots a & S \to a \\ a\cdots aba\cdots a & S \to b \\ a\cdots a\varepsilon a\cdots a & S \to \varepsilon \end{cases}$$

$$S \Rightarrow bSb \Rightarrow bbSbb \cdots\cdots \text{ same measure as above}$$

$$L(G) = \{w \in \Sigma^* \mid w = w^R\} \quad R: reverse$$

---

减字角 算术表达    **Example (Simple Arithmetic Expressions)**

CFG G = (V, Σ, R, S) with

1. V = {S}

2. Σ = {+, −, ×, /, (, ), 0, 1, 2, . . . , 9}

3. Rules R:
   S → S + S | S − S | S × S | S/S | (S) | −S | 0 | 1 | · · · | 9

L(G): valid arithmetic expressions over single-digit integers

S derives string 3 × (5 + 6)?

$$S \Rightarrow S \times S \Rightarrow S \times (S) \Rightarrow S \times (S+S) \Rightarrow 3 \times (S+S) \Rightarrow 3 \times (5+S)$$
$$\Rightarrow 3 \times (5+6)$$

## Regular Languages are context-free

(if) (could say)

**Theorem**  Regular Language => Context free

Let $\Sigma$ be an alphabet and let $L \subseteq \Sigma^*$ be a regular language. Then L is a context-free language (Every regular language is context-free).

**Proof**  (general idea)

因是则语言，有个DFA M接受，L是上下文无美的需要
有个上下文无美语法 G 满足 $L=L(M)=L(G)$

Since L is a regular language, there exists a deterministic finite automaton $M = (Q, \Sigma, \delta, q, F)$ that accepts L. To prove that L is context-free, we have to define a context-free grammar $G = (V, \Sigma, R, S)$, such that $L = L(M) = L(G)$. Thus, G must have the following property:

$$w \in L(M) \Leftrightarrow w \in L(G)$$

For every string $w \in \Sigma^*$,

$w \in L(M)$ if and only if $w \in L(G)$,

which can be reformulated as

$$M \text{ accepts } w \text{ if and only if } S \overset{*}{\Rightarrow} w.$$

G的V就是M的所有 Q

Set $V = \{R_i \mid q_i \in Q\}$ (that is, G has a variable for every state of M). Now, for every transition $\delta(q_i, a) = q_j$ add a rule $R_i \rightarrow aR_j$. For every accepting state $q_i \in F$ add a rule $R_i \rightarrow \varepsilon$. Finally, make the start variable $S = R_0$.

$R_0$ is the initial state of the machine

---

## Regular Languages are context-free

L is regular => L is context free $\nLeftarrow$



**Closure properties of CFLs:** CFLs are closed under operations like union and concatenation but not under intersection or complementation.
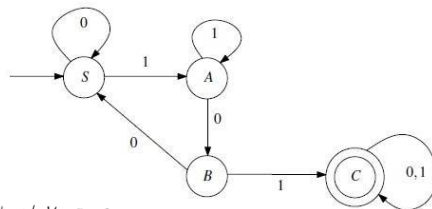
---

## Regular Languages are context-free

**Example**

Let L be the language defined as
$$L = \{w \in \{0, 1\}^* : 101 \text{ is a substring of } w\}.$$

The DFA M that accepts L



将DFA转换为 CFG

How can we convert M to a context-free grammar G whose language is L?

---

## Regular Languages are context-free

**Example**

- $G = \{V, \Sigma, R, S\}$

$V = \{S, A, B, C\}$

$\Sigma = \{0, 1\}$

Start variable : S (initial state of M)

Rules:

$S \rightarrow 0S \mid 1A$

$A \rightarrow 0B \mid 1A$

$B \rightarrow 0S \mid 1C$

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

## Chomsky Normal Form (CNF) 乔姆斯基 公式

**Definition**

A context-free grammar G = (V, Σ, R, S) is said to be in **Chomsky normal form**, if every rule in R has one of the following three forms: 如果Rules满足下面三点条件

- $A \to BC$, where A, B, and C are elements of V, $B \neq S$, and $C \neq S$.
- $A \to a$, where A is an element of V and a is an element of Σ.
- $S \to \varepsilon$, where S is the start variable.

**Why CNF?**

Grammars in Chomsky normal form are far easier to analyze.

**Example**

Rules of CFG in Chomsky normal form with V = {S, A, B}, Σ = {a, b}:

$G_1$ : $S \to AB$, $S \to c$, $A \to a$, $B \to b$ (CNF)

$G_1$ : $S \to aA$, $A \to a$, $B \to c$ (not CNF)

---

## Chomsky Normal Form (CNF)

**Theorem**

Let Σ be an alphabet and let $L \subseteq \Sigma^*$ be a context-free language. There exists a context-free grammar in Chomsky normal form, whose language is L (Every CFL can be described by a CFG in CNF).

**CFL → CNF**

Given CFG G = (V, Σ, R, S). Replace, one-by-one, every rule that is not "Chomsky".

- Start variable (not allowed on RHS of rules)
- ε-rules ($A \to \varepsilon$ not allowed when A isn't start variable)
- all other violating rules ($A \to B$, $A \to aBc$, $A \to BCDE$)

---

context free grammar → chomsky normal form

## Converting CFG into CNF

**Transformation steps**

11首先删除 start variable

Step 1. Eliminate the start variable from the right-hand side of the rules.

- New start variable $S_0$
- New rule $S_0 \to S$

Step 2. Remove **ε-rules** $A \to \varepsilon$, where $A \in V - \{S\}$.

- Before: $B \to xAy$ and $A \to \varepsilon \mid \cdots$
- After: $B \to xAy \mid xy$ and $A \to \cdots$

When removing $A \to \varepsilon$ rules, insert all new replacements:

- Before: $B \to AbA$ and $A \to \varepsilon \mid \cdots$
- After: $B \to AbA \mid bA \mid Ab \mid b$ and $A \to \cdots$

---

In final, All rules must be satisfied with above 3 requirements.

## Converting CFG into CNF

**Transformation steps**

Step 3. Remove **unit rules** $A \to B$, where $A \in V$.

- Before: $A \to B$ and $B \to xCy$
- After: $A \to xCy$ and $B \to xCy$

Step 4. Eliminate all rules having more than two symbols on the right-hand side.

- Before: $A \to B_1B_2B_3$
- After: $A \to B_1A_1$, $A_1 \to B_2B_3$

Step 5. Eliminate all rules of the form $A \to ab$, where a and b are not both variables.

- Before: $A \to ab$
- After: $A \to B_1B_2$, $B_1 \to a$, $B_2 \to b$.

## Converting CFG into CNF

**Example**

Given a CFG G = (V, Σ, R, S), where V = {A, B}, Σ = {0, 1}, A is the start variable, and R consists of the rules:

$$A \to BAB \mid B \mid \varepsilon$$
$$B \to 00 \mid \varepsilon$$

ε-rules:
$$A \to \varepsilon$$
$$B \to \varepsilon$$

Convert this G to CNF:

Step 1. Eliminate the start variable from the right-hand side of the rules.

$$S \to A$$
$$A \to BAB \mid B \mid \varepsilon$$
$$B \to 00 \mid \varepsilon$$

---

## Converting CFG into CNF

$$S \to A$$
$$\checkmark A \to BAB \mid B \mid \varepsilon$$
$$\checkmark B \to 00 \mid \varepsilon$$

**Example**

Step 2. Remove ε-**rules**.

(1) Remove A → ε: S → A, A → BAB

$$\begin{cases} S \to A \mid \varepsilon \\ A \to BAB \mid B \mid BB \\ B \to 00 \mid \varepsilon \end{cases}$$

(2) Remove B → ε: A → BAB, A → B, A → BB

$$S \to A \mid \varepsilon$$
$$A \to BAB \mid B \mid BB \mid AB \mid BA \mid A$$
$$B \to 00$$

---

## Converting CFG into CNF

$$S \to A, \quad \underline{A \to A}$$
$$A \to B$$

$$\begin{cases} S \to A \mid \varepsilon \\ A \to BAB \mid B \mid BB \mid AB \mid BA \mid A \\ B \to 00 \end{cases}$$

**Example**

Step 3. Remove **unit-rules**.

(1) Remove A → A:

$$\begin{cases} S \to A \mid \varepsilon \\ A \to BAB \mid B \mid BB \mid AB \mid BA \\ B \to 00 \end{cases}$$

(2) Remove S → A:

$$S \to B$$
$$A \to B$$
$$\begin{cases} S \to \varepsilon \mid BAB \mid B \mid BB \mid AB \mid BA \\ A \to BAB \mid B \mid BB \mid AB \mid BA \\ B \to 00 \end{cases}$$

---

## Converting CFG into CNF

$$S \to \varepsilon \mid BAB \mid B \mid BB \mid AB \mid BA$$
$$A \to BAB \mid B \mid BB \mid AB \mid BA$$
$$B \to 00$$

**Example**

Step 3. Remove **unit-rules**.

(3) Remove S → B:

$$\begin{cases} S \to \varepsilon \mid BAB \mid BB \mid AB \mid BA \\ \checkmark \underline{A \to BAB \mid B \mid BB \mid AB \mid BA} \\ B \to 00 \end{cases}$$

(4) Remove A → B:

$$S \to \varepsilon \mid BAB \mid BB \mid AB \mid BA \mid 00$$
$$A \to BAB \mid BB \mid AB \mid BA$$
$$B \to 00$$

## Converting CFG into CNF

$S \to \varepsilon \mid BAB \mid BB \mid AB \mid BA \mid 00$
$A \to BAB \mid BB \mid AB \mid BA \mid 00$
$B \to 00$

**Example**

Step 4. Eliminate all rules having more than two symbols on the right-hand side.

(1) Remove $S \to BAB$:   想法迁起 BAB 变为两个symbol

$BAB \xrightarrow{A_1} BA_1$

$S \to \varepsilon \mid BB \mid AB \mid BA \mid 00 \mid BA_1$
$A \to BAB \mid BB \mid AB \mid BA \mid 00$
$B \to 00$
Assume $A_1 \to AB$

(2) Remove $A \to BAB$:

$S \to \varepsilon \mid BB \mid AB \mid BA \mid 00 \mid BA_1$    replace $00 \to A_3A_3$
$A \to BB \mid AB \mid BA \mid 00 \mid BA_2$    replace $00 \to A_4A_4$
$B \to 00$
$A_1 \to AB$
$A_2 \to AB$

---

## Converting CFG into CNF

**Example**

Step 5. Eliminate all rules, whose right-hand side contains exactly two symbols, which are not both variables.

(1) Remove $S \to 00$:   $S \to \varepsilon \mid BB \mid AB \mid BA \mid BA_1 \mid A_3A_3$
$A \to BB \mid AB \mid BA \mid BA_2 \mid A_4A_4$
$B \to 00$ => replace to $B \to A_5A_5$
$A_1 \to AB$
$A_2 \to AB$
$A_3 \to 0$
$A_4 \to 0$
$A_5 \to 0$

(2) Remove $A \to 00$:
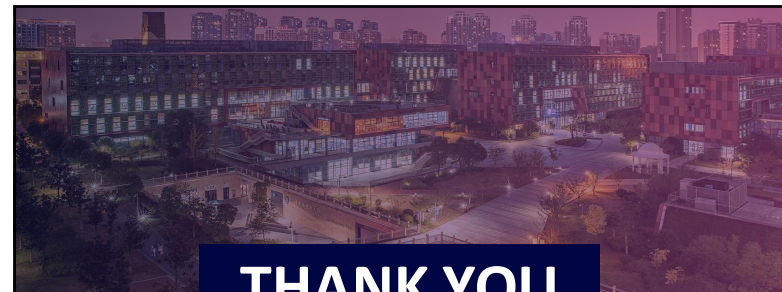
(3) Remove $B \to 00$

---

## Converting CFG into CNF

**Example**

Step 5. Eliminate all rules, whose right-hand side contains exactly two symbols, which are not both variables.

(3) Remove $S \to 00$:

---



THANK YOU