# Finite Automata Theory and Formal Languages
# TMV027/DIT321– LP4 2018

Lecture 13

Ana Bove

May 7th 2018

## Recap: Context-Free Grammars

- Equivalence between recursive inference, (leftmost/rightmost) derivations and parse trees;

- Ambiguous grammars;

- Inherent ambiguity;

- Proofs about grammars and languages.

## Overview of Today's Lecture

- Simplification of CFL;
- Chomsky normal form for CFL.

**Contributes to the following learning outcome:**

- Explain and manipulate the diff. concepts in automata theory and formal lang;
- Simplify automata and context-free grammars;
- Differentiate and manipulate formal descriptions of lang, automata and grammars.

And guest lecture by *Martin Fabian* on *Application of Formal Verification to the Lane Change Module of an Autonomous Vehicle*.

## Generating, Reachable, Useful and Useless Symbols

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.
Let $X \in V \cup T$ and let $\alpha, \beta \in (V \cup T)^*$.

**Definition:** $X$ is *reachable* if $S \Rightarrow^* \alpha X \beta$.
(This is similar to accessible states in FA.)

**Definition:** $X$ is *generating* if $X \Rightarrow^* w$ for some $w \in T^*$.

**Definition:** The symbol $X$ is *useful* if $S \Rightarrow^* \alpha X \beta \Rightarrow^* w$ for some $w \in T^*$.
**Note:** A symbol that is useful should be generating and reachable.

**Definition:** $X$ is *useless* iff it is not useful.

We shall simplify the grammars by eliminating useless symbols.

# Computing the Generating Symbols

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.

The following recursive procedure computes the generating symbols of $G$:

Base Case: All elements of $T$ are generating;

Recursive Step: If a production $A \rightarrow \alpha$ is such that all symbols of $\alpha$ are known to be generating, then $A$ is also generating.
Observe that $\alpha$ could be $\epsilon$.

The recursive step must be applied until no new symbols are found generating.

**Theorem:** *The procedure above finds all and only the generating symbols of a grammar.*

**Proof:** See Theorem 7.4 in the book.

# Example: Generating Symbols

Consider the grammar over $\{a\}$ given by the rules:

$$
\begin{aligned}
S &\rightarrow aS \mid W \mid U \\
W &\rightarrow aW \\
U &\rightarrow a \\
V &\rightarrow aa
\end{aligned}
$$

$a$ is generating.

$U$ and $V$ are generating since $U \rightarrow a$ and $V \rightarrow aa$.

$S$ is generating since $S \rightarrow U$.

No other symbol is found generating so $W$ is not generating.

After eliminating the non-generating symbols and their productions we get

$$S \rightarrow aS \mid U \qquad U \rightarrow a \qquad V \rightarrow aa$$

## Computing the Reachable Symbols

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.

The following recursive procedure computes the reachable symbols of $G$:

Base Case: The start variable $S$ is reachable;

Recursive Step: If $A$ is reachable and we have a production $A \to \alpha$ then all symbols in $\alpha$ are reachable.

The recursive step must be applied until no new symbols are found reachable.

**Theorem:** *The procedure above finds all and only the reachable symbols of a grammar.*

**Proof:** See Theorem 7.6 in the book.

## Example: Reachable Symbols

Consider the grammar given by the rules:

$$S \to aB \mid BC \qquad C \to b$$
$$A \to aA \mid c \mid aDb \quad D \to B$$
$$B \to DB \mid C$$

$S$ is reachable.

Hence $a$, $B$ and $C$ are reachable.

Then $b$ and $D$ are reachable.

No other symbol are found reachable so $A$ and $c$ are not reachable.

After eliminating the non-reachable symbols and their productions we get

$$S \to aB \mid BC \quad C \to b$$
$$B \to DB \mid C \quad D \to B$$

# Eliminating Useless Symbols

It is important in which order we check generating and reachable symbols!

**Example:** Consider the following grammar

$$S \to AB \mid a \qquad A \to b$$

If we first check for generating symbols and then for reachability we get

$$S \to a$$

If we first check for reachability and then for generating we get

$$S \to a \qquad A \to b$$

# Eliminating Useless Symbols

**Theorem:** *Let $G = (V, T, \mathcal{R}, S)$ be a CFG and let $\mathcal{L}(G) \neq \emptyset$.*
*Let $G' = (V', T', \mathcal{R}', S)$ be constructed as follows:*

1. *First, eliminate all non-generating symbols and all productions involving one or more of those symbols;*
2. *Then, eliminate all non-reachable symbols and all productions involving one or more of those symbols.*

*Then $G'$ has no useless symbols and $\mathcal{L}(G) = \mathcal{L}(G')$.*

**Proof:** See Theorem 7.2 in the book.

# Example: Eliminating Useless Symbols

Consider the grammar given by the rules:

$$
\begin{aligned}
S &\rightarrow gAe \mid aYB \mid CY & A &\rightarrow bBY \mid ooC \\
B &\rightarrow dd \mid D & C &\rightarrow jVB \mid gl \\
D &\rightarrow n & U &\rightarrow kW \\
V &\rightarrow baXXX \mid oV & W &\rightarrow c \\
X &\rightarrow fV & Y &\rightarrow Yhm
\end{aligned}
$$

After eliminating non-generating symbols:

$$
\begin{aligned}
S &\rightarrow gAe & A &\rightarrow ooC \\
B &\rightarrow dd \mid D & C &\rightarrow gl \\
D &\rightarrow n & U &\rightarrow kW \\
& & W &\rightarrow c
\end{aligned}
$$

After eliminating non-reachable symbols:

$$
S \rightarrow gAe \qquad A \rightarrow ooC \qquad C \rightarrow gl
$$

What is the language generated by the grammar?

# Nullable Variables

**Definition:** A variable $A$ is *nullable* if $A \Rightarrow^* \epsilon$.

**Note:** Observe that only variables are nullable!

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.

The following recursive procedure computes the nullable variables of $G$:

Base Case: If $A \rightarrow \epsilon$ is a production then $A$ is nullable;

Recursive Step: If $B \rightarrow X_1 X_2 \ldots X_k$ is a production and all the $X_i$ are nullable then $B$ is also nullable.

The recursive step must be applied until no new symbols are found nullable.

**Theorem:** *The procedure above finds all and only the nullable variables of a grammar.*

**Proof:** See Theorem 7.7 in the book.

# Eliminating $\epsilon$-Productions

**Definition:** An *$\epsilon$-production* is a production of the form $A \to \epsilon$.

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.

The following procedure eliminates the $\epsilon$-production of $G$:

1. Determine all nullable variables of $G$;
2. Build $\mathcal{P}$ with all the productions of $\mathcal{R}$ plus a rule $A \to \alpha\beta$ whenever we have $A \to \alpha B \beta$ and $B$ is nullable.
   **Note:** If $A \to X_1 X_2 \ldots X_k$ and all $X_i$ are nullable, we do not include the case where all the $X_i$ are absent;
3. Construct $G' = (V, T, \mathcal{R}', S)$ where $\mathcal{R}'$ contains all the productions in $\mathcal{P}$ except for the $\epsilon$-productions.

**Theorem:** *The grammar $G'$ constructed from the grammar $G$ as above is such that $\mathcal{L}(G') = \mathcal{L}(G) - \{\epsilon\}$.*

**Proof:** See Theorem 7.9 in the book.

# Example: Eliminating $\epsilon$-Productions

**Example:** Consider the grammar given by the rules:

$$S \to aSb \mid SS \mid \epsilon$$

By eliminating $\epsilon$-productions we obtain

$$S \to ab \mid aSb \mid S \mid SS$$

**Example:** Consider the grammar given by the rules:

$$S \to AB \qquad A \to aAA \mid \epsilon \qquad B \to bBB \mid \epsilon$$

By eliminating $\epsilon$-productions we obtain

$$S \to A \mid B \mid AB \qquad A \to a \mid aA \mid aAA \qquad B \to b \mid bB \mid bBB$$

# Eliminating Unit Productions

**Definition:** A *unit production* is a production of the form $A \rightarrow B$.

(This is similar to $\epsilon$-transitions in a $\epsilon$-NFA.)

Let $G = (V, T, \mathcal{R}, S)$ be a CFG.

The following procedure eliminates the unit production of $G$:

1. Build $\mathcal{P}$ with all the productions of $\mathcal{R}$ plus a rule $A \rightarrow \alpha$ whenever we have $A \rightarrow B$ and $B \rightarrow \alpha$;

   Observe that this step might introduce new unit productions that must be expanded!

2. Construct $G' = (V, T, \mathcal{R}', S)$ where $\mathcal{R}'$ contains all the productions in $\mathcal{P}$ except for the unit production.

**Theorem:** *The grammar $G'$ constructed from the grammar $G$ as above is such that $\mathcal{L}(G') = \mathcal{L}(G)$.*

**Proof:** See Theorem 7.13 in the book.

# Example: Eliminating Unit Productions

Consider the grammar given by the rules:

$$
\begin{aligned}
S &\rightarrow CBh \mid D & A &\rightarrow aaC \\
B &\rightarrow Sf \mid ggg & C &\rightarrow cA \mid d \mid C \\
D &\rightarrow E \mid SABC & E &\rightarrow be
\end{aligned}
$$

By eliminating unit productions we obtain:

$$
\begin{aligned}
S &\rightarrow CBh \mid be \mid SABC & A &\rightarrow aaC \\
B &\rightarrow Sf \mid ggg & C &\rightarrow cA \mid d \\
D &\rightarrow be \mid SABC & E &\rightarrow be
\end{aligned}
$$

## Simplification of a Grammar

**Theorem:** *Let $G = (V, T, \mathcal{R}, S)$ be a CFG whose language contains at least one string other than $\epsilon$. If we construct $G'$ by*

1. *First, eliminating $\epsilon$-productions;*
2. *Then, eliminating unit productions;*
3. *Finally, eliminating useless symbols;*

*using the procedures shown before then $\mathcal{L}(G') = \mathcal{L}(G) - \{\epsilon\}$.*

*In addition, $G'$ contains no $\epsilon$-productions, no unit productions and no useless symbols.*

**Proof:** See Theorem 7.14 in the book.

**Note:** It is important to apply the steps in this order!

## Chomsky Normal Form

**Definition:** A CFG is in *Chomsky Normal Form* (CNF) if $G$ has no useless symbols and all the productions are of the form $A \rightarrow BC$ or $A \rightarrow a$.

**Note:** Observe that a CFG that is in CNF has no unit or $\epsilon$-productions!

**Theorem:** *For any CFG $G$ whose language contains at least one string other than $\epsilon$, there is a CFG $G'$ that is in Chomsky Normal Form and such that $\mathcal{L}(G') = \mathcal{L}(G) - \{\epsilon\}$.*

**Proof:** See Theorem 7.16 in the book.

# Constructing a Chomsky Normal Form

Let us assume $G$ has no $\epsilon$- or unit productions and no useless symbols.

Then every production is of the form $A \to a$ or $A \to X_1 X_2 \ldots X_k$ for $k > 1$.

If $X_i$ is a terminal introduce a new variable $A_i$ and a new rule $A_i \to X_i$
(if no such rule exists for $X_i$ with a variable that has no other rules).

Use $A_i$ in place of $X_i$ in any rule whose body has length $> 1$.

Now, all rules are of the form $B \to b$ or $B \to C_1 C_2 \ldots C_k$ with all $C_j$ variables.

Introduce $k - 2$ new variables and break each rule $B \to C_1 C_2 \ldots C_k$ as

$$B \to C_1 D_1 \quad D_1 \to C_2 D_2 \quad \cdots \quad D_{k-2} \to C_{k-1} C_k$$

# Example: Chomsky Normal Form

**Example:** Consider the grammar given by the rules:

$$S \to aSb \mid SS \mid ab$$

We first obtain

$$S \to ASB \mid SS \mid AB \qquad A \to a \qquad B \to b$$

Then we build a grammar in Chomsky Normal Form

$$
\begin{aligned}
S &\to AC \mid SS \mid AB & A &\to a \\
C &\to SB & B &\to b
\end{aligned}
$$

**Example:** Observe however that

$$S \to aa \mid a$$

is NOT equivalent to

$$S \to SS \mid a$$

Instead we need to build

$$S \to AA \mid a \qquad A \to a$$

# Overview of Next Lecture

Sections 7.2–7.4, and notes on *Pumping lemma*:

- Regular grammars;
- Chomsky hierarchy;
- Pumping lemma for CFL;
- Closure properties of CFL;
- Decision properties of CFL.

# Overview of next Week

| Mon 14 | Tue 15 | Wed 16 | Thu 17 | Fri 18 |
|---|---|---|---|---|
| | 10-12 EA Exercise | | 10-12 ES61 Individual help | |
| Lec 13-15 HB3 CFL. | | | Lec 13-15 HB3 PDA. TM. | |
| 15-17 EA Exercise | | 15-17 EL41 Consultation | | |

**Assignment 6:** CFL.
*Deadline:* Sunday May 20th 23:59.