# CSE211 – Formal Languages and Automata Theory

## U2L10_Simplification of CFG

**Dr. P. Saravanan**
**School of Computing**
**SASTRA Deemed University**

# Outline

- Recap of previous class
- Properties of CGL
- Substitution rule
- Simplification of CFG
  - Eliminating useless production
  - Eliminating e-production
  - Eliminating unit production
- Reason for Simplication

# Properties of CFL

- CFG's may be simplified to fit certain special forms, like

  – *Chomsky Normal Form (CNF)* and

  – *Greiback Normal Form (GNF).*

- Some, but not all, properties of RL's are also possessed by the CFL's.

- Unlike the RL, many computational problems about the CFL *cannot* be answered.

- That is, there are many undecidable problems about CFL's.

# A Substitution Rule

$$S \rightarrow aB$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc$$

$$B \rightarrow aA$$

$$B \rightarrow b$$

Substitute
$$B \rightarrow b$$

$$S \rightarrow aB \,|\, ab$$

$$A \rightarrow aaA$$

$$A \rightarrow abBc \,|\, abbc$$

$$B \rightarrow aA$$

# In general

$$A \rightarrow xBz$$

$$B \rightarrow y_1$$

Substitute

$$B \rightarrow y_1$$

$$A \rightarrow xBz \mid xy_1z$$

# Simplification of CFG

- *Every CFG can be transformed into an equivalent grammar in Chomsky Normal Form,* after simplifying the CFG in the following ways:

- 

  – eliminating *useless symbols* (which do not appear in any derivation from the start symbol);

  – eliminating $\varepsilon$-*productions* (of the form $A \rightarrow \varepsilon$);

  – eliminating *unit productions* (of the form $A \rightarrow B$);

# Eliminating Useless Symbols

- We say symbol *X* is *useful* for a grammar *G = (V, T, P, S)* if there is some derivation of the form
  - *S $\overset{*}{\Rightarrow}$ aXb $\overset{*}{\Rightarrow}$ w with w∈T\**.

- A symbol is said to be *useless* if not useful.

- Omitting useless symbols obviously will not change the language generated by the grammar.

- There are two types of *usefulness* ---
  - *X* is *generating* if *X $\overset{*}{\Rightarrow}$ w*;
  - *X* is *reachable* if *S $\overset{*}{\Rightarrow}$ aXb*.

# Eliminating Useless Symbols Example 1

- Eliminate useless symbols in a grammar with the following productions:
  - $S \rightarrow AB \mid a$
  - $A \rightarrow b$.

- *B* is *not generating*, and is so eliminated at first, resulting in $S \rightarrow a$, $A \rightarrow b$, in which *A* is *not reachable*

- and so eliminated too, with $S \rightarrow a$ as the only production left.

- The order of eliminations is *essential*: *eliminate non-generating symbols at first*.

# Eliminating Useless Symbols Thorem

■ Let $G = (V, T, P, S)$ be a CFG, and assume that $L(G) \neq \phi$, i.e., assume that $G$ generates at least one string. Let $G_1 = (V_1, T_1, P_1, S)$ be the grammar obtained by the following steps *in order*:

– eliminate non-generating symbols and all related productions, resulting in grammar $G_2$;

– eliminate all symbols not reachable in $G2$.

■ Then, $G_1$ has no useless symbol and $L(G_1) = L(G)$.

# Computing Generating and Reachable Symbols

- **How to compute generating symbols?**

- *Basis*: Every terminal symbol is generating.

- *Induction*: if every symbol in $a$ in $A \rightarrow a$ is generating, then $A$ is generating.

- 

- **How to compute reachable symbols?**

- *Basis*: the start symbol $S$ is reachable.

- *Induction*: if nonterminal $A$ is reachable, then all the symbols in $A \rightarrow a$ are reachable.

# Eliminating ε -Productions

- **A definition** --- a nonterminal A is said to be *nullable* if
  - A $\overset{*}{\Rightarrow}$ ε.

- **A Theorem** --- We want to prove that
  - if a language *L* has a CFG, then the language $L - \{\varepsilon\}$ can be generated by a CFG without ε -production.

- Two steps for the above proof:
  - find "nullable" symbols;
  - transform productions into ones which generate no empty string using the nullable symbols.

# Eliminating ε -Productions

- Given a grammar with productions as follows:

  $S \rightarrow AB$

  $A \rightarrow aAA \mid ε$

  $B \rightarrow bBB \mid ε$

- then, we can see the following facts:

  – *A* and *B* are *nullable* because they derive empty strings;

  – *S* is also *nullable* because *A* and *B* are nullable.

# Eliminating ε -Productions

- How to find nullable symbols systematically?

- *Algorithm 1 ---*

- *Basis*: if $A \rightarrow \varepsilon$ is a production, then *A* is nullable

- *Induction*: if all $C_i$ in $B \rightarrow C_1 C_2 \ldots C_k$ are nullable, then *B* is nullable, too.

# Eliminating ε -Productions

- How to transform productions into ones which generate no empty string?

- ***Algorithm 2 ---***

- For each production $A \rightarrow X_1X_2...X_k$, in which *m* of the *k* $X_i$'s are nullable, then generate accordingly 2*m* versions of this production where

  - (1) the nullable $X_i$'s in all possible combinations are present or absent; and

  - (2) if $A \rightarrow \varepsilon$ is in the 2*m* ones, eliminate it.

- For $S \rightarrow AB$, $A \rightarrow aAA \mid \varepsilon$, $B \rightarrow bBB \mid \varepsilon$:

  - We know $S$, $A$, $B$ are *nullable.*

  - From $S \rightarrow AB$, we get $S \rightarrow AB \mid A \mid B \mid \varepsilon$ where $S \rightarrow \varepsilon$ should be eliminated.

  - From $A \rightarrow aAA$, we get $A \rightarrow aAA \mid aA \mid aA \mid a$ where the repeated $A \rightarrow aA$ should be removed.

  - And from $B \rightarrow bBB,$ similarly we get $B \rightarrow bBB \mid bB \mid b.$

  - Overall result:

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

# Why do we simplify & Example?

# Summary

- Recap of previous class

- Properties of CGL

- Substitution rule

- Simplification of CFG
  - Eliminating useless production
  - Eliminating e-production
  - Eliminating unit production

- Reason for Simplication

# References

- John E. Hopcroft, Rajeev Motwani and Jeffrey D. Ullman, *Introduction to Automata Theory*, Languages, and Computation, Pearson, 3$^{rd}$ Edition, 2011.

- Peter Linz, An Introduction to Formal Languages and Automata, Jones and Bartle Learning International, United Kingdom, 6$^{th}$ Edition, 2016.

**Next Class:**

# Chomsky Normal Form (CNF)

**THANK YOU**