

Project Three - Report

Converting Context-Free Grammar into Chomsky Normal Form

Nov.17 2019, Harbin Institute of Technology, Shenzhen

Members:

- **Champagne Kwan** *conversion algorithm designation, coding*
- **Yanglin Liu** *test & report*

Programming Language: **Python3**

Part1. Introduction

In this experiment, we try to implementing an algorithm to convert a given context-free grammar to corresponding Chomsky normal form.

In formal language theory, a context-free grammar G is said to be in Chomsky normal form if all of its production rules are of the form:

$$\begin{aligned}A &\rightarrow BC \\A &\rightarrow a \\S &\rightarrow \epsilon\end{aligned}$$

where A , B , and C are nonterminal symbols, a is a terminal symbol (a symbol that represents a constant value), S is the start symbol, and ϵ denotes the empty string.

We can know that all of production rules must be:

- only one nonterminal symbol at left part, and
- all candidate forms (right part of production rules) contain exactly two nonterminal symbol, or only one terminal symbol
- start symbol could product an empty string

Part2. Designation

2.1 brief description of conversion method

Here we referenced an example from textbook^[1]. The original CFG is shown below. And we use colored text to represent updated parts: red for adding, cyan for deleting.

$$\begin{aligned} \text{Original :} & S \rightarrow ASA|aB \\ & A \rightarrow B|S \\ & B \rightarrow b|\epsilon \end{aligned}$$

1. **Start:** introduce a new start symbol to grammar:

Step1 : introduce new start symbol

$$\begin{aligned} S_0 & \rightarrow S \\ S & \rightarrow ASA|aB \\ A & \rightarrow B|S \\ B & \rightarrow b|\epsilon \end{aligned}$$

2. **Epsilon Eliminating:** for each rule, eliminating any empty string in right part, and replace left symbol in other rules with non-empty right parts.

Step2 : eliminate $B \rightarrow \epsilon$, and add to other rules contains B

$$S_0 \rightarrow S$$

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

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

$$B \rightarrow b|\epsilon$$

Step3 : eliminate $A \rightarrow \epsilon$, and add to other rules contains A

$$S_0 \rightarrow S$$

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

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

$$B \rightarrow b$$

3. **Unit Eliminating:** remove unit rules (only one nonterminal in right parts), and replace left symbol in other rules with corresponding right parts.

Step4 : remove $S \rightarrow S$

$$S_0 \rightarrow S$$

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

$$A \rightarrow B|S$$

$$B \rightarrow b$$

Step5 : eliminate $S_0 \rightarrow S$ and replace S use corresponding right part

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

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

$$A \rightarrow B|S$$

$$B \rightarrow b$$

Step6 : eliminate $A \rightarrow B$ and replace B use corresponding right part

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

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

$$A \rightarrow S|b$$

$$B \rightarrow b$$

Step7 : eliminate $A \rightarrow S$ and replace S use corresponding right part

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

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

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

$$B \rightarrow b$$

4. **Folding Replacement of variables:** convert remaining rules which right parts have more than 2 nonterminals into which have exact 2 nonterminals. Here we use *folding strategy* shown below:

Folding Strategy example :

$A \rightarrow BCDE$ will be converted into 3 rules

$$F \rightarrow BC$$

$$G \rightarrow FD$$

$$A \rightarrow GE$$

In step 4, algorithm should ensure that left part should be same once right parts are same.

2.2 Designation of Entity Classes

Here we need three entity classes: *variable, production rule and grammar*. In source code, the data structure is:

- class **Var**: *represents a variable, terminals or nonterminals*
 - **symbol**: string type, used as identifier to represent.
 - **is_nonterm**: boolean type, indicates if variable is a nonterminal symbol.
- class **Product**: *represents a production rule*
 - **left**: a Var object, represents left part of rule, we only need one variable, it's guaranteed by feature of context-free.
 - **right**: list, represents all right parts of rule, elements of this list is list since a rule may contains many right parts and any part may contains many variables.
- class **CFG**: *represents a context-free grammar*
 - **start**: a Var object, represents start symbol of grammar.
 - **nonterm**: set, contains all nonterminal symbols appeared in grammar.
 - **term**: set, contains all terminal symbols appeared in grammar.
 - **products**: list, all Product objects.
 - **nonterminal_cand**: a list of constants, used in step 4 for folding replacement of variables mentioned in Part 2.1

2.3 Core Algorithm

Algorithm of Epsilon Eliminating is shown below, creating this algorithm is challenging since a nonterminal which derives empty string may appear many times in a given candidate form, for example: if two rules below is known

$$A \rightarrow BSB \quad \text{and} \quad B \rightarrow \epsilon$$

then three new candidate forms will be added into first rule because any B can be assigned an empty value.

$$A \rightarrow BSB|BS|SB|S$$

Thus the algorithm is designed to be a recursive procedure to handle all situations.

```
function EpsilonEliminating:
    while not reached end of candidate and current symbol is not
    given symbol:
        StackPush(symbol stack,current symbol)
        current symbol = next symbol
    if reached end of candidate:
        add all symbols in symbol stack into new candidates
        return
    StackPush(symbol stack,given symbol)
    # First recursion, given symbol is not assigned to be empty
    EpsilonEliminating() with next position
    while not StackEmpty(symbol stack) and StackTop(symbol
    stack) != given symbol:
        StackPop(symbol stack)
    StackTop(symbol stack)
    # Second recursion, given symbol is assigned to be empty
    EpsilonEliminating() with next position
```

Part3. Test

Here are two test cases used in testing, and results.

TestCase1 : (from textbook)

$$S \rightarrow ASA|aB$$

$$A \rightarrow B|S$$

$$B \rightarrow b|\epsilon$$

TestCase2 : (from website)

$$S \rightarrow ASB$$

$$A \rightarrow aAS|a|\epsilon$$

$$B \rightarrow SbS|A|bb$$

TestCase 1:

Non-Terminals:A,B,S

Terminals:a,b

Start:S

Products:

$$S \rightarrow DA|CB|a|AS|SA$$

$$S \rightarrow DA|CB|a|AS|SA$$

$$A \rightarrow b|DA|CB|a|AS|SA$$

$$B \rightarrow b$$

$$C \rightarrow a$$

$$D \rightarrow AS$$

TestCase 2:

Non-Terminals:A,B,S

Terminals:a,b

Start:S

Products:

$$S \rightarrow EB|SB|AS$$

$$S \rightarrow EB|SB|AS$$

$$A \rightarrow FS|a|CS$$

$$B \rightarrow GS|DD|FS|a|CS$$

$$C \rightarrow a$$

$$D \rightarrow b$$

$$E \rightarrow AS$$

$$F \rightarrow CA$$

$$G \rightarrow SD$$

