

## Topic: Chomsky Normal Form

### Course: Formal Languages & Finite Automata

### Author: Gîncu Olivia

**Git Repository:** <https://github.com/livia994/LFA/tree/main/Lab5>

### Overview

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

$A \rightarrow BC$ , or

$A \rightarrow a$ , or

$S \rightarrow \epsilon$ ,

where  $A$ ,  $B$ , and  $C$  are nonterminal symbols, the letter  $a$  is a terminal symbol (a symbol that represents a constant value),  $S$  is the start symbol, and  $\epsilon$  denotes the empty string. Also, neither  $B$  nor  $C$  may be the start symbol, and the third production rule can only appear if  $\epsilon$  is in  $L(G)$ , the language produced by the context-free grammar  $G$ .

Every grammar in Chomsky normal form is context-free, and conversely, every context-free grammar can be transformed into an equivalent one which is in Chomsky normal form and has a size no larger than the square of the original grammar's size.

Chomsky Normal Form (CNF) is a standard form used to represent context-free grammars (CFGs) in formal language theory. In CNF, each production rule has one of two forms: either a non-terminal symbol derives exactly two non-terminals, or a non-terminal symbol derives a single terminal. This strict structure simplifies the parsing process and facilitates various computational tasks, such as parsing algorithms and grammar analysis.

CNF offers several advantages in theoretical and practical contexts. Firstly, it enables efficient parsing algorithms, such as the CYK (Cocke-Younger-Kasami) algorithm, which operates efficiently on grammars in CNF. Additionally, CNF provides an unambiguous representation of CFGs, making it easier to reason about and analyze grammar. Moreover, CNF can aid in grammar simplification and optimization, facilitating grammar transformations and ensuring grammatical correctness.

However, it's important to note that not all context-free grammars can be directly converted into CNF without introducing additional non-terminal symbols or restructuring productions. Handling epsilon ( $\lambda$ ) productions, unit productions, and productions with more than two symbols on the right-hand side requires careful consideration and transformation techniques. Despite these challenges, Chomsky Normal Form remains a fundamental concept in formal language theory and computational linguistics, serving as a cornerstone for various language processing tasks and theoretical analyses.

### Implementation:

#### 1. Identifying Nullable Symbols (Step 1):

# Identifying nullable symbols

*nullable\_symbols = set()*

*for non\_terminal, productions in self.rules.items():*

*if " in productions:*

*nullable\_symbols.add(non\_terminal)*

It iterates over the grammar's productions to identify symbols that can derive the empty string ( $\lambda$ ). If a production contains an empty string, we add its non-terminal symbol to the set of nullable symbols.

## 2. Removing Epsilon Productions (Step 2):

# Removing epsilon productions

```
for non_terminal, productions in self.rules.items():
    updated_productions = [p for p in productions if p != ""]
    self.rules[non_terminal] = updated_productions
```

It iterates over the grammar's productions again and removes any epsilon ( $\lambda$ ) productions. This ensures that productions no longer derive the empty string.

## 3. Iteratively Removing Productions Containing Nullable Symbols (Step 3):

# Iteratively remove productions containing nullable symbols

while True:

exists\_new\_null = False

for non\_terminal, productions in self.rules.items():

for production in productions:

null\_test = production

for nullable in nullable\_symbols:

null\_test = null\_test.replace(nullable, "")

if null\_test == "" and non\_terminal not in nullable\_symbols:

nullable\_symbols.add(non\_terminal)

exists\_new\_null = True

if not exists\_new\_null:

break

It iteratively removes productions that contain nullable symbols until no new nullable symbols are found. This process ensures that all nullable symbols are properly accounted for in the grammar.

## 4. Replacing Productions Containing Nullable Symbols (Step 4):

# Replace productions containing nullable symbols

for non\_terminal, productions in self.rules.items():

new\_productions = []

for production in productions:

null\_indices = [i for i, symbol in enumerate(production) if symbol in nullable\_symbols]

if null\_indices:

combinations = itertools.product([0, 1], repeat=len(null\_indices))

for combination in combinations:

new\_production = "".join([symbol for i, symbol in enumerate(production) if  
i not in null\_indices or combination[null\_indices.index(i)] == 1])

if new\_production != "":

new\_productions.append(new\_production)

else:

new\_productions.append(production)

self.rules[non\_terminal] = list(set(new\_productions))

It replaces productions that contain nullable symbols with new productions that account for all possible combinations of nullable symbols. This ensures that the grammar properly reflects the nullable symbols' influence on production derivations.

## 5. Modifying Productions to Avoid Infinite Loops (Step 5):

# Modify productions to avoid infinite loops

if 'C' in self.rules:

self.rules['C'] = ['aC', 'bC'] # Modify the production to avoid infinite loops

It modifies specific productions, such as 'C -> abC', to prevent potential infinite loops in the grammar. This modification ensures that the grammar is well-defined and avoids recursive productions that could lead to non-termination in parsing algorithms.

### Output/Conclusion:

```
S -> bA | b | B | aA | C | a | AC
A -> X | ab | aS | X' | b
B -> bS | a
C -> aC | bC
D -> B | AB
X -> AB
X' -> Ba
```

*Through this laboratory work, I explored how to simplify the structure of grammar used in language processing. By learning about techniques like removing certain types of rules and adjusting others, I gained a better grasp of how computers understand and process language. This experience not only helped me understand the theory behind language processing but also provided practical skills that I can use in real-world applications.*