

# Context-Free Grammar (CFG)

CFG stands for context-free grammar. It is a formal grammar which is used to generate all possible patterns of strings in a given formal language. Context-free grammar  $G$  can be defined by four tuples as:

1.  $G = (V, T, P, S)$

## Where,

**G** is the grammar, which consists of a set of the production rule. It is used to generate the string of a language.

**T** is the final set of a terminal symbol. It is denoted by lower case letters.

**V** is the final set of a non-terminal symbol. It is denoted by capital letters.

**P** is a set of production rules, which is used for replacing non-terminals symbols(on the left side of the production) in a string with other terminal or non-terminal symbols(on the right side of the production).

**S** is the start symbol which is used to derive the string. We can derive the string by repeatedly replacing a non-terminal by the right-hand side of the production until all non-terminal have been replaced by terminal symbols.

## Example 1:

Construct the CFG for the language having any number of a's over the set  $\Sigma = \{a\}$ .

Solution:

As we know the regular expression for the above language is

1. r.e. =  $a^*$

Production rule for the Regular expression is as follows:

1.  $S \rightarrow aS$  rule 1
2.  $S \rightarrow \epsilon$  rule 2

Now if we want to derive a string "aaaaa", we can start with start symbols.

1.  $S$
2.  $aS$
3.  $aaS$  rule 1
4.  $aaaS$  rule 1
5.  $aaaaS$  rule 1

6. aaaaaS      rule 1
7. aaaaaaS     rule 1
8. aaaaaaε     rule 2
9. aaaaaa

The r.e.  $= a^*$  can generate a set of string  $\{\epsilon, a, aa, aaa, \dots\}$ . We can have a null string because  $S$  is a start symbol and rule 2 gives  $S \rightarrow \epsilon$ .

### **Example 2:**

Construct a CFG for the regular expression  $(0+1)^*$

#### **Solution:**

$\{ \epsilon, 0, 1, 01, 10, 11, 00, 010, \dots \}$

The CFG can be given by,

Production rule (P):

1.  $S \rightarrow 0S \mid 1S$
2.  $S \rightarrow \epsilon$

The rules are in the combination of 0's and 1's with the start symbol. Since  $(0+1)^*$  indicates  $\{\epsilon, 0, 1, 01, 10, 00, 11, \dots\}$ . In this set,  $\epsilon$  is a string, so in the rule, we can set the rule  $S \rightarrow \epsilon$ .

### **Example 3:**

Construct a CFG for a language  $L = \{wcwR \mid \text{where } w \in (a, b)^*\}$ .

Solution:

The string that can be generated for a given language is  $\{aacia, bcb, abcb, bacab, abbcbb, \dots\}$

The grammar could be:

1.  $S \rightarrow aSa$       rule 1
2.  $S \rightarrow bSb$       rule 2
3.  $S \rightarrow c$         rule 3

Now if we want to derive a string "abbcbb", we can start with start symbols.

1.  $S \rightarrow aSa$
2.  $S \rightarrow abSba$       from rule 2
3.  $S \rightarrow abbSbba$     from rule 2

4.  $S \rightarrow abbcbbba$  from rule 3

Thus any of this kind of string can be derived from the given production rules.

**Example 4:**

Construct a CFG for the language  $L = a^n b^{2n}$  where  $n \geq 1$ .

Solution:

The string that can be generated for a given language is  $\{abb, aabbbb, aaabbbbbbb, \dots\}$ .

The grammar could be:

5.  $S \rightarrow aSbb \mid abb$

Now if we want to derive a string "aabbbb", we can start with start symbols.

$S \rightarrow aSbb$

$S \rightarrow aabbbb$

**Example 5:**

A CFG for the regular language corresponding to the RE  $0^* 1 1^*$ .

Solution:

The language is the concatenation of two languages: all strings of zeroes with all strings of ones.

The grammar could be:

$S \rightarrow CD$

$C \rightarrow 0C \mid 0$

$D \rightarrow 1D \mid 1$

# BACKUS NAUR FORM(BNF) :

It stands for Backus-Naur Form one notation used to write grammar.

It is a formal, mathematical way to specify context-free grammars

Example:

BNF for number--

$$\langle \text{number} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{number} \rangle \langle \text{digit} \rangle$$
$$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

WHERE

“ $::=$ ” means “is defined as” (some variants use “ $:=$ ” instead)

“ $\mid$ ” means “or”

Angle brackets mean a Non Terminal

Symbols without angle brackets are Terminals

$$\langle \text{expression} \rangle ::= \langle \text{expression} \rangle + \langle \text{term} \rangle$$
$$\mid \langle \text{expression} \rangle - \langle \text{term} \rangle$$
$$\mid \langle \text{term} \rangle$$
$$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$$
$$\mid \langle \text{term} \rangle / \langle \text{factor} \rangle$$
$$\mid \langle \text{factor} \rangle$$
$$\langle \text{factor} \rangle ::= ( \langle \text{expression} \rangle )$$
$$\mid \langle \text{primary} \rangle$$
$$\mid \langle \text{variable} \rangle$$
$$\mid \langle \text{number} \rangle$$

EXAMPLE 2:

There is the production for any grammar as follows:

$$S \rightarrow aSa$$
$$S \rightarrow bSb$$
$$S \rightarrow c$$

In BNF, we can represent above grammar as follows:

$$S \rightarrow aSa \mid bSb \mid c$$

# EXTENDED BACKUS NAUR FORM(EBNF) :

EBNF is a few simple extensions to BNF which make expressing grammars more convenient

- “\*” (The Kleene Star): means 0 or more occurrences
- “+” (The Kleene Cross): means 1 or more occurrences
- “[ ... ]”: means 0 or 1 occurrences
- () Use of parentheses for grouping
- {}\* used to show arbitrary sequence
- If you have a rule such as:

Examples :

- $\langle \text{real no} \rangle ::= \langle \text{integer part} \rangle . \langle \text{fraction part} \rangle \mid . \langle \text{fraction part} \rangle$

You can replace it with:

$$\langle \text{real no} \rangle ::= [\langle \text{integer part} \rangle] . \langle \text{fraction part} \rangle$$

- If you have a rule such as:

$$\langle \text{signed integer} \rangle ::= + \langle \text{integer} \rangle \mid - \langle \text{integer} \rangle$$
$$\langle \text{integer} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{integer} \rangle \langle \text{digit} \rangle$$

You can replace it with:

$$\langle \text{signed integer} \rangle ::= [+|-] \langle \text{digit} \rangle \{ \langle \text{digit} \rangle \}^*$$
$$\langle \text{signed integer} \rangle ::= [+|-] \{ \langle \text{digit} \rangle \}^+$$

- If you have a rule such as:

$$\begin{aligned} \langle \text{exp} \rangle &::= \langle \text{term} \rangle \\ &\quad \mid \langle \text{exp} \rangle + \langle \text{term} \rangle \\ &\quad \mid \langle \text{exp} \rangle - \langle \text{term} \rangle \end{aligned}$$

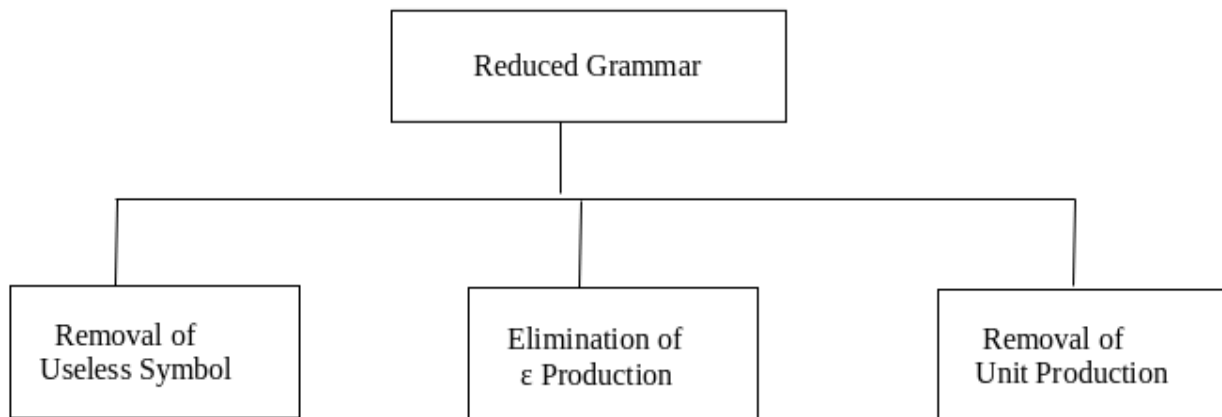
You can replace it with:

$$\langle \text{exp} \rangle ::= \langle \text{term} \rangle [(+|-) \langle \text{term} \rangle]$$

# SIMPLIFICATION OF CFG

As we have seen, various languages can efficiently be represented by a context-free grammar. All the grammar are not always optimized that means the grammar may consist of some extra symbols(non-terminal). Having extra symbols, unnecessary increase the length of grammar. Simplification of grammar means reduction of grammar by removing useless symbols. The properties of reduced grammar are given below:

1. Each variable (i.e. non-terminal) and each terminal of  $G$  appears in the derivation of some word in  $L$ .
2. There should not be any production as  $X \rightarrow Y$  where  $X$  and  $Y$  are non-terminal.
3. If  $\epsilon$  is not in the language  $L$  then there need not to be the production  $X \rightarrow \epsilon$ .



## Removal of Useless Symbols

A symbol can be useless if it does not appear on the right-hand side of the production rule and does not take part in the derivation of any string. That symbol is known as a useless symbol. Similarly, a variable can be useless if it does not take part in the derivation of any string. That variable is known as a useless variable.

For Example:

1.  $T \rightarrow aaB \mid abA \mid aaT$
2.  $A \rightarrow aA$
3.  $B \rightarrow ab \mid b$
4.  $C \rightarrow ad$

In the above example, the variable 'C' will never occur in the derivation of any string, so the production  $C \rightarrow ad$  is useless. So we will eliminate it, and the other productions are written in

such a way that variable C can never reach from the starting variable 'T'.

Production  $A \rightarrow aA$  is also useless because there is no way to terminate it. If it never terminates, then it can never produce a string. Hence this production can never take part in any derivation.

To remove this useless production  $A \rightarrow aA$ , we will first find all the variables which will never lead to a terminal string such as variable 'A'. Then we will remove all the productions in which the variable 'B' occurs.

## Elimination of $\epsilon$ Production

The productions of type  $S \rightarrow \epsilon$  are called  $\epsilon$  productions. These type of productions can only be removed from those grammars that do not generate  $\epsilon$ .

**Step 1:** First find out all nullable non-terminal variable which derives  $\epsilon$ .

**Step 2:** For each production  $A \rightarrow a$ , construct all production  $A \rightarrow x$ , where x is obtained from a by removing one or more non-terminal from step 1.

**Step 3:** Now combine the result of step 2 with the original production and remove  $\epsilon$  productions.

### Example:

Remove the production from the following CFG by preserving the meaning of it.

1.  $S \rightarrow XYX$
2.  $X \rightarrow 0X \mid \epsilon$
3.  $Y \rightarrow 1Y \mid \epsilon$

#### Solution:

Now, while removing  $\epsilon$  production, we are deleting the rule  $X \rightarrow \epsilon$  and  $Y \rightarrow \epsilon$ . To preserve the meaning of CFG we are actually placing  $\epsilon$  at the right-hand side whenever X and Y have appeared.

Let us take

1.  $S \rightarrow XYX$

If the first X at right-hand side is  $\epsilon$ . Then

1.  $S \rightarrow YX$

Similarly if the last X in R.H.S. =  $\epsilon$ . Then

1.  $S \rightarrow XY$

If  $Y = \epsilon$  then

1.  $S \rightarrow XX$

If Y and X are  $\epsilon$  then,

1.  $S \rightarrow X$

If both X are replaced by  $\epsilon$

1.  $S \rightarrow Y$

Now,

1.  $S \rightarrow XY \mid YX \mid XX \mid X \mid Y$

Now let us consider

1.  $X \rightarrow 0X$

If we place  $\epsilon$  at right-hand side for X then,

1.  $X \rightarrow 0$
2.  $X \rightarrow 0X \mid 0$

Similarly  $Y \rightarrow 1Y \mid 1$

Collectively we can rewrite the CFG with removed  $\epsilon$  production as

1.  $S \rightarrow XY \mid YX \mid XX \mid X \mid Y$
2.  $X \rightarrow 0X \mid 0$
3.  $Y \rightarrow 1Y \mid 1$

## Removing Unit Productions

The unit productions are the productions in which one non-terminal gives another non-terminal. Use the following steps to remove unit production:

**Step 1:** To remove  $X \rightarrow Y$ , add production  $X \rightarrow a$  to the grammar rule whenever  $Y \rightarrow a$  occurs in the grammar.

**Step 2:** Now delete  $X \rightarrow Y$  from the grammar.

**Step 3:** Repeat step 1 and step 2 until all unit productions are removed.

For example:

1.  $S \rightarrow 0A \mid 1B \mid C$
2.  $A \rightarrow 0S \mid 00$
3.  $B \rightarrow 1 \mid A$



4.  $C \rightarrow 01$

**Solution:**

$S \rightarrow C$  is a unit production. But while removing  $S \rightarrow C$  we have to consider what  $C$  gives. So, we can add a rule to  $S$ .

1.  $S \rightarrow 0A \mid 1B \mid 01$

Similarly,  $B \rightarrow A$  is also a unit production so we can modify it as

1.  $B \rightarrow 1 \mid 0S \mid 00$

Thus finally we can write CFG without unit production as

1.  $S \rightarrow 0A \mid 1B \mid 01$

2.  $A \rightarrow 0S \mid 00$

3.  $B \rightarrow 1 \mid 0S \mid 00$

4.  $C \rightarrow 01$