

Context Free Grammars

Shakir Ullah Shah

Context Free Grammars

Three fundamental areas covered in the book are

1. **Theory of Automata**
2. **Theory of Formal Languages**
3. **Theory of Turing Machines**

We have completed the first area.

We begin exploring the second area in this chapter.

Syntax as a Method for Defining Languages

Recursively definition of the set of valid arithmetic expressions as follows:

Rule 1: Any number is in the set AE.

Rule 2: If x and y are in AE, then so are (x) , $-(x)$, $(x + y)$, $(x - y)$, $(x * y)$, (x/y) , $(x ** y)$

where $**$ is our notation for exponentiation

Note that we use parentheses around every component factor to avoid ambiguity expressions such as $3 + 4 - 5$ and $8/4/2$.

There is a different way for defining the set AE:

using a set of substitution rules similar to the grammatical rules.

Defining AE by substitution rules

Start \rightarrow AE

AE \rightarrow (AE)

AE \rightarrow -(AE)

AE \rightarrow (AE + AE)

AE \rightarrow (AE - AE)

AE \rightarrow (AE * AE)

AE \rightarrow (AE/AE)

AE \rightarrow (AE ** AE)

AE \rightarrow AnyNumber

Rule 2: If x and y are in AE, then so are

(x), -(x), (x + y), (x - y), (x * y),
(x/y), (x ** y)

Rule 1: Any number is in the set AE.

Example

We will show that $((3 + 4) * (6 + 7))$ is in AE

Start \Rightarrow AE \Rightarrow (AE * AE)

$\Rightarrow ((\text{AE} + \text{AE}) * (\text{AE} + \text{AE}))$

$\Rightarrow ((3 + 4) * (6 + 7))$

1. Start \rightarrow AE
2. AE \rightarrow (AE + AE)
3. AE \rightarrow (AE - AE)
4. AE \rightarrow (AE * AE)
5. AE \rightarrow (AE/AE)
6. AE \rightarrow (AE ** AE)
7. AE \rightarrow (AE)
8. AE \rightarrow -(AE)
9. AE \rightarrow

AnyNumber

Definition of Terms

A word that cannot be replaced by anything is called **terminal**.

- In the above example, the terminals are the phrase AnyNumber, and the symbols $+$ $-$ $*$ $/$ $**$ $($ $)$

A word that must be replaced by other things is called **nonterminal**.

- The nonterminals are Start and AE.

The sequence of applications of the rules that produces the finished string of terminals from the starting symbol is called a **derivation** or a **generation** of the word.

The grammatical rules are referred to as

Symbolism for Generative Grammars

Definition:

A **context-free grammar (CFG)** is a collection of three things:

1. An alphabet Σ of letters called **terminals** from which we are going to make strings that will be the words of a language.

2. A set of symbols called **nonterminals**, one of which is the symbol S , standing for “start here”.

3. A finite set of **productions** of the form:
One nonterminal \rightarrow finite string of terminals and/or nonterminals

where the strings of terminals and nonterminals can consist of only terminals, or of only nonterminals, or of any mixture of terminals and nonterminals, or even the empty string. We require that at least one production has the nonterminal S as its left side.

Definition:

The **language generated by a CFG** is the set of all strings of terminals that can be produced from the start symbol S using the productions as substitutions.

A language generated by a CFG is called a **context-free language (CFL)**.

Notes:

The language generated by a CFG is also called the **language defined by the CFG**, or the **language derived from the CFG**, or the **language produced by the CFG**.

We insist that **nonterminals be designated by capital letters**, whereas **terminals are designated by lowercase letters and special symbols**.

Example

Let the only terminal be a and the productions be

- Prod1 $S \rightarrow aS$
- Prod2 $S \rightarrow \Lambda$

If we apply Prod 1 six times and then apply Prod 2, we generate the following:

- $S \Rightarrow aS \Rightarrow aaS \Rightarrow aaaS \Rightarrow aaaaS$
- $\Rightarrow aaaaaS \Rightarrow aaaaaaS \Rightarrow aaaaaa\Lambda = aaaaaa$

If we apply Prod2 without Prod1, we find that Λ is in the language generated by this CFG.

Hence, this CFL is exactly a^* .

Note: the symbol “ \rightarrow ” means “can be replaced by”, whereas the symbol “ \Rightarrow ” means “can develop to”.

Let the terminals be a and b , the only nonterminal be S , and the productions be

Prod1 $S \rightarrow aS$

Prod2 $S \rightarrow bS$

Prod3 $S \rightarrow \Lambda$

The word ab can be generated by the derivation

$S \Rightarrow aS \Rightarrow abS \Rightarrow ab\Lambda = ab$

The word $baab$ can be generated by

$S \Rightarrow bS \Rightarrow baS \Rightarrow baaS \Rightarrow baabS \Rightarrow$
 $baab\Lambda = baab$

Clearly, the language generated by the above CFG is

$(a + b)^*$.

Let the terminals be a and b , the only nonterminal be S , and the productions be

Prod1 $S \rightarrow aS$

Prod2 $S \rightarrow bS$

Prod3 $S \rightarrow a$

Prod4 $S \rightarrow b$

The language generated by the above CFG is
 $(a + b)^+.$

Example

Let the terminals be a and b , the the nonterminal be S and X , and the productions be

Prod 1 $S \rightarrow XaaX$

Prod 2 $X \rightarrow aX$

Prod 3 $X \rightarrow bX$

Prod 4 $X \rightarrow \Lambda$

We already know from the previous example that the last three productions will generate any possible strings of a 's and b 's from the nonterminal X .

Hence, the words generated from S have the form

anything aa anything

Hence, the language produced by this CFG is

$$(a + b)^*aa(a + b)^*$$

which is the language of all words with a double a in them somewhere.

For example, the word *baabb* can be generated by

$$\begin{aligned} S &\rightarrow XaaX \rightarrow bXaaX \rightarrow baaX \rightarrow baaX \\ &\rightarrow baabX \rightarrow baabbX \rightarrow baabb\Lambda = baabb \end{aligned}$$

Example

Consider the CFG:

$$S \rightarrow aSb$$

$$S \rightarrow \Lambda$$

It is easy to verify that the language generated by this CFG is the **non-regular** language $\{a^n b^n\}$.

For example, the word $a^4 b^4$ is derived by

$$S \Rightarrow aSb \Rightarrow aaSbb \Rightarrow aaaSbbb$$

$$\Rightarrow aaaaSbbbb \Rightarrow aaaa\Lambda bbbb = aaaabbbb$$

Disjunction Symbol |

Let us introduce the symbol | to mean disjunction (or).

We use this symbol to combine all the productions that have the same left side.

For example, the CFG

Prod 1 $S \rightarrow XaaX$

Prod 2 $X \rightarrow aX$

Prod 3 $X \rightarrow bX$

Prod 4 $X \rightarrow \Lambda$

can be written more compactly as

Prod 1 $S \rightarrow XaaX$

Prod 2 $X \rightarrow aX|bX|\Lambda$