

Context Free Grammars and Languages

Grammars are used to generate the words of a language and to determine whether a word is in a language. Formal languages, which are generated by grammars, provide models for both natural languages, such as, English, and for programming languages, such as C, JAVA.

Context free grammars are used to define syntax of almost all programming languages, in context of computer science. Thus an important application of context free grammars occurs in the specification and compilation of programming languages.

Formal Description of Context Free Grammar:

A context free grammar is defined by 4-tuples (V, T, P, S) where,

V = set of variables

T = set of terminal symbols

P = set of rules and productions

S = start symbol and $S \in V$

There are four components in CFG.

1. There is a finite set of symbol that form the strings of language being defined. We call this alphabet the terminal or terminal symbols. In above tuples, it is represented by T
2. There is a finite set of variables, also called sometimes non-terminals or non-terminal symbols or syntactic categories'. Each variable represents a language i.e a set of strings.
3. One of the variables represents the language being defined. It is called the start symbol.
4. There is a finite set of productions or rules that represent the recursive definition of a language. Each production consists of :
 - a. A variable that is being defined by the production. This is called head of production.
 - b. The production symbol \rightarrow
 - c. A string of Zero or more terminals and Variables. This string is called the body of the production, represents one way to form the string in the language of the head.

Note: The variable symbols are after represented by capital letters. The terminals are analogous to the input alphabet and are often represented by lower case letters. One of the variables is designed as a start variable. It usually occurs on the left hand side of the topmost rule.

For example,

$S \rightarrow \epsilon$

$S \rightarrow 0S1$

This is a CFG defining the grammar of all the strings with equal no of 0's followed by equal no of 1's.

Here,

the two rules define the production P,
 $\epsilon, 0, 1$ are the terminals defining T,
 S is a variable symbol defining V
 And S is start symbol from where production starts.

CFG Vs RE

The CFG are more powerful than the regular expressions as they have more expressive power than the regular expression. Generally regular expressions are useful for describing the structure of lexical constructs as identical keywords, constants etc. But they do not have the capability to specify the recursive structure of the programming constructs. However, the CFG are capable to define any of the recursive structure also. Thus, CFG can define the languages that are regular as well as those languages that are not regular.

Compact Notation for the Productions

It is convenient to think of a production as “belonging” to the variable of head. we may use the terms like A-production or the production of A to refer to the production whose head is A. We write the production for a grammar by listing each variable once and then listing all the bodies of the production for that variable, separated by |. This is called compact Notation.

Example

CFG representing the language over $\Sigma = \{a, b\}$ which is palindrome language.

$S \rightarrow \epsilon \mid a \mid b$ Grammar (1)

$S \rightarrow a S a$

$S \rightarrow b S b$

Here the first production is written in compact notation. The detail of this production look like as

$S \rightarrow \epsilon$

$S \rightarrow a$

$S \rightarrow b$

Meaning of context free:

Consider an example:

$S \rightarrow a M b$

$M \rightarrow A \mid B$

$A \rightarrow \epsilon \mid aA$

$B \rightarrow \epsilon \mid bB$

How, consider a String aaAb, which is an intermediate stage in the generating of aab. It is natural to call the strings “aa” and “b” that surround the symbol A, the “context” of A in this particular string. Now, the rule $A \rightarrow aA$ says that we can replace A by the string aA no matter what the surrounding strings are; in other words, independently of the context of A.

When there is a production of form $lw_1r \rightarrow lw_2r$ (but not of the form $w_1 \rightarrow w_2$), the grammar is context sensitive since w_1 can be replaced by w_2 only when it is surrounded by the strings “l” and “r”.

Derivation Using Grammar Rule:

We apply the production rule of a CFG to infer that certain strings are in the language of a certain variable. This is the process of deriving strings by applying productions rules of the CFG.

General Convention used in Derivation

- Capital letters near the beginning of the alphabet A, B and so on are used as variables.
- Capital letters near the end of the alphabet such as X, Y are used as either terminals or variables.
- Lower case letters near the beginning of the alphabet, a, b and so on are used for terminals.
- Lower case letters near the end of alphabet, such as ‘w’ or ‘z’ are used for string of terminals.
- Lower case Greek letters such as α and β are used for string of terminal or variables.

Derivation:

A derivation of a context free grammar is a finite sequence of strings $\beta_0 \beta_1 \beta_2 \dots \beta_n$ such that:

→ For $0 \leq i \leq n$, the string $\beta_i \in (V \cup T)^*$

→ $\beta_0 = S$

→ For $0 \leq i \leq n$, there is a production of P that applied to β_i yields β_{i+1}

→ $\beta_n \in T^*$

There are two possible approaches of derivation:

→ Body to head (Bottom Up) approach.

→ Head to body (Top Down) approach.

Body to head

Here, we take strings known to be in the language of each of the variables of the body, concatenate them, in the proper order, with any terminals appearing in the body, and infer that the resulting string is the language of the variables in the head.

Consider grammar,

$S \rightarrow S + S$

$S \rightarrow S/S$ Grammer(2)

$S \rightarrow (S)$

$S \rightarrow S-S$

$S \rightarrow S*S$

$S \rightarrow a$

Here given $a + (a*a) / a - a$

Now, by this approach.

SN	String Inferred	Variable(For Language of)	Production	String Used
1	a	S	$S \rightarrow a$	
2	a^*a	S	$S \rightarrow S^*S$	String 1
3	(a^*a)	S	$S \rightarrow (S)$	String 2
3	$(a^*a)/a$	S	$S \rightarrow S/S$	String 3 and string 1
4	$a+(a^*a)/a$	S	$S \rightarrow S+S$	String 1 and 3
5	$a+(a^*a)/a-a$	S	$S \rightarrow S-S$	String 4 and 1

Thus, in this process we start with any terminal appearing in the body and use the available rules from body to head.

Head to Body

Here, we use production from head to body. We expand the start symbol using a production, whose head is the start symbol. Here we expand the resulting string until all strings of terminal are obtained. Here we have two approaches:

Leftmost Derivation: Here leftmost symbol (variable) is replaced first.

Rightmost Derivation: Here rightmost symbol is replaced first.

For example: consider the previous example of deriving string

$a+(a^*a)/a-a$ with the grammar(2)

Now leftmost derivation for the given string is:

$S \rightarrow S + S$ rule $\rightarrow S + S$
 $S \rightarrow a + S$ rule $S \rightarrow a$ [Note here we have replaced left symbol of the body]
 $S \rightarrow a + S - S$ rule $S \rightarrow S - S$
 $S \rightarrow a + S / S - S$ rule $S \rightarrow S / S$
 $S \rightarrow a + (S) / S - S$ rule $S \rightarrow (S)$
 $S \rightarrow (S^*S) / S - S;$ rule $S \rightarrow S^*S$
 $S \rightarrow a + (a^*S) / S - S$ rule $S \rightarrow a$
 $S \rightarrow a + (a^*a) / S - S$ “ “
 $S \rightarrow a + (a^*a) / a - S$ “ “
 $S \rightarrow a + (a^*a) / a - a$ “ “

And rightmost derivation is;

$S \rightarrow S - S;$ rule $S \rightarrow S - S$
 $S \rightarrow S - a;$ rule $S \rightarrow a$
 $S \rightarrow S + S - a;$ rule $S \rightarrow S + S$
 $S \rightarrow S + S / S - a;$ rule $S \rightarrow S / S$
 $S \rightarrow S + S / a - a;$ rule $S \rightarrow a$
 $S \rightarrow S + (S) / a - a;$ rule $S \rightarrow (S)$
 $S \rightarrow S + (S^*S) / a - a;$ rule $S \rightarrow S^*S$

$S \rightarrow S + (S^*a) / a - a;$ rule $S \rightarrow a$
 $S \rightarrow S + (a^*a) / a - a;$
 $S \rightarrow a + (a^*a) / a - a;$

Direct Derivation:

$\alpha_1 \rightarrow \alpha_2$: If α_2 can be derived directly from α_1 , then it is direct derivation.

$\alpha_1 \rightarrow^* \alpha_2$: If α_2 can be derived from α_1 with zero or more steps of the derivation, then it is just derivation.

Example

$S \rightarrow aSa \mid ab \mid a \mid b \mid \epsilon$

Direct derivation:

$S \rightarrow ab$

$S \rightarrow aSa$

$S \rightarrow aaSaa$

$S \rightarrow aaabaa$

Thus $S \rightarrow^* aaabaa$ is just a derivation.

Language of Grammar (Context Free Grammar):

Let $G = (V, T, P \text{ and } S)$ is a context free grammar. Then the language of G denoted by $L(G)$ is the set of terminal strings that have derivation from the start symbol in G .

i.e. $L(G) = \{x \in T^* \mid S \rightarrow^* x\}$

The language generated by a CFG is called the Context Free Language (CFL).

Sentential forms

Derivations from the start symbol produce strings that have a special role. We call these "sentential forms". i.e. if $G=(V,T,P,S)$ is a CFG, then any string α is in $(V \cup T)^*$ such that $S \rightarrow^* \alpha$ is a sentential form. If $S \rightarrow_{lm}^* \alpha$, then α is a left sentential forms, and if $S \rightarrow_{rm}^* \alpha$, then α is a right sentential form.

Note: The language $L(G)$ is those sentential forms that are in T^* i.e. they consist solely of terminals.

Derivation Tree / Parse Tree:

Parse tree is a tree representation of strings of terminals using the productions defined by the grammar. A parse tree pictorially shows how the start symbol of a grammar derives a string in the language.

Parse tree may be viewed as a graphical representation for a derivation that filters out the choice regarding the replacement order.

Formally, given a Context Free Grammar $G = (V, T, P \text{ and } S)$, a parse tree is a n-ary tree having following properties:

- The root is labeled by the start symbol.
- Each interior node of parse tree are variables.
- Each leaf node of parse is labeled by a terminal symbol or ϵ .
- If an interior node is labeled with a non terminal A and its childrens are x_1, x_2, \dots, x_n from left to right then there is a production P as:
 $A \rightarrow x_1, x_2, \dots, x_n$ for each $x_i \in T$.

Example

Consider the grammer

$$S \rightarrow aSa \mid a \mid b \mid \epsilon$$

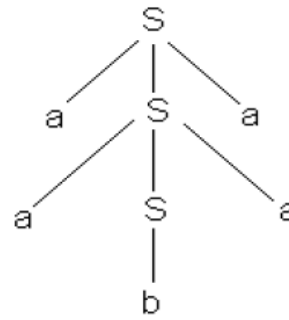
Now for string $S \rightarrow^* aabaa$

We have,

$$S \rightarrow aSa$$

$$S \rightarrow aaSaa$$

$$S \rightarrow aabaa \quad \text{So the parse tree is } \rightarrow:$$



Exercise

a) Consider the Grammar G:

$$S \rightarrow A1B$$

$$A \rightarrow 0A \mid \epsilon$$

$$B \rightarrow 0B \mid 1B \mid \epsilon$$

1. Construct the parse tree for 00101
2. Construct the parse tree for 1001
3. Construct the parse tree for 00011

b) Consider the grammer for the arithmetic expression

$$E \rightarrow EOPE \mid (E) \mid id$$

$$OP \rightarrow + \mid - \mid * \mid / \mid \uparrow$$

1. Construct the parse tree for $id + id * id$
2. Construct the parse tree for $(id + id) * (id + id)$

c) Construct a CFG that generates language of balanced parentheses:

-Show the parse tree computation for

$$1. () ()$$

$$2. ((())) ()$$

Solution:

The CFG is $G=(V,T,P,S)$

$S \rightarrow SS$

$S \rightarrow (S)$

$S \rightarrow \epsilon$

Where,

$V = \{S\}$

$S = \{S\}$

$T = \{(\,)\}$

& $P =$ as listed above

Now the parse tree for $(\)(\)$ is:

And the parse tree for $((\)(\))(\)$

Ambiguity in Grammar:

A Grammar $G = (V, T, P$ and $S)$ is said to be ambiguous if there is a string $w \in L(G)$ for which we can derive two or more distinct derivation tree rooted at S and yielding w . In other words, a grammar is ambiguous if it can produce more than one leftmost or more than one rightmost derivation for the same string in the language of the grammar.

Example

$S \rightarrow AB \mid aaB$

$A \rightarrow a \mid Aa$

$B \rightarrow b$

For any string aab ;

We have two leftmost derivations as;

$S \rightarrow AB$ the parse tree for this derivation is:

$\rightarrow AaB$

$\rightarrow aaB$

$\rightarrow aab$

Also,

$S \rightarrow aaB$ the parse tree for this derivation is:

$\rightarrow aab$

Normal forms and Simplification of CFG

The goal of this section is to show that every context free language (without ϵ) is generated by a CFG in which all production are of the form $A \rightarrow BC$ or $A \rightarrow a$, where A, B and c are variables, and a is terminal. This form is called Chomsky Normal Form. To get there, we need to make a number of preliminary simplifications, which are themselves useful in various ways;

1. We must eliminate “useless symbols”: Those variables or terminals that do not appear in any derivation of a terminal string from the start symbol.
2. We must eliminate “ ϵ -production”: Those of the form $A \rightarrow \epsilon$ for some variable A .
3. We must eliminate “unit production”: Those of the form $A \rightarrow B$ for variables A and B .

Eliminating Useless Symbols:

We say a symbol x is useful for a grammar $G = (V, T, P, S)$ if there is some derivation of the form $S \xrightarrow{*} \alpha x \beta \xrightarrow{*} w$, where w is in T^* . Here, x may be either variable or terminal and the sentential form $\alpha x \beta$ might be the first or last in the derivation. If x is not useful, we say it is useless.

Thus useful symbols are those variables or terminals that appear in any derivation of a terminal string from the start symbol. Eliminating a useless symbol includes identifying whether or not the symbol is “generating” and “reachable”.

Generating Symbol: We say x is generating if $x \xrightarrow{*} w$ for some terminal string w : Note that every terminal is generated since w can be that terminal itself, which is derived by zero steps.

Reachable symbol: We say x is reachable if there is derivation $S \xrightarrow{*} \alpha x \beta$ for some α and β .

Thus if we eliminate the non generating symbols and then non-reachable, we shall have only the useful symbols left.

Example

Consider a grammar defined by following productions:

$S \rightarrow aB \mid bX$
 $A \rightarrow Ba \mid bSX \mid a$
 $B \rightarrow aSB \mid bBX$
 $X \rightarrow SB \mid aBX \mid ad$

Here;

A and X can directly generate terminal symbols. So, A and X are generating symbols. As we have the productions $A \rightarrow a$ and $X \rightarrow ad$.

Also,

$S \rightarrow bX$ and X generates terminal string so S can also generate terminal string. Hence, S is also generating symbol.

B can not produce any terminal symbol, so it is non-generating.

Hence, the new grammar after removing non-generating symbols is:

$S \rightarrow bX$
 $A \rightarrow bSX \mid a$
 $X \rightarrow ad$

Here,

A is non-reachable as there is no any derivation of the form $S \rightarrow^* \alpha A \beta$ in the grammar. Thus eliminating the non-reachable symbols, the resulting grammar is:

$$S \rightarrow bX$$

$$X \rightarrow ad$$

This is the grammar with only useful symbols.

Exercise

1) Remove useless symbol from the following grammar:

$$S \rightarrow xyZ \mid XyzZ$$

$$X \rightarrow Xz \mid xYZ$$

$$Y \rightarrow yYy \mid XZ$$

$$Z \rightarrow Zy \mid z$$

2) Remove useless symbol from the following grammar

$$S \rightarrow aC \mid SB$$

$$A \rightarrow bSCa$$

$$B \rightarrow aSB \mid bBC$$

$$C \rightarrow aBc \mid ad$$

Eliminating ϵ -productions:

A grammar is said to have ϵ -productions if there is a production of the form $A \rightarrow \epsilon$. Here our strategy is to begin by discovering which variables are “nullable”. A variable 'A' is “nullable” if $A \rightarrow^* \epsilon$.

Algorithm (Steps to remove ϵ -production from the grammar):

- If there is a production of the form $A \rightarrow \epsilon$, then A is “nullable”.
- If there is production of the form $B \rightarrow X_1, X_2, \dots, X_n$ And each X_i 's are nullable then B is also nullable.
- Find all the nullable variables.
- If $B \rightarrow X_1, X_2, \dots, X_n$ is a production in P then add all productions P' formed by striking out some subsets of there X_i 's that are nullable.
- Do not include $B \rightarrow \epsilon$ if there is such production.

Example:

Consider the grammar:

$$S \rightarrow ABC$$

$$A \rightarrow BB \mid \epsilon$$

$$B \rightarrow CC \mid a$$

$$C \rightarrow AA \mid b$$

Here,

$$A \rightarrow \epsilon \quad A \text{ is nullable.}$$

$$C \rightarrow AA \rightarrow^* \epsilon, \quad C \text{ is nullable}$$

$$B \rightarrow CC \rightarrow^* \epsilon, \quad B \text{ is nullable}$$

$$S \rightarrow ABC \rightarrow^* \epsilon, \quad S \text{ is nullable}$$

Now for removal of ϵ -production:

In production

$S \rightarrow ABC$, all A, B and C are nullable. So, striking out subset of each the possible combination of production gives new productions as:

$$S \rightarrow ABC \mid AB \mid BC \mid AC \mid A \mid B \mid C$$

Similarly for other can be done and the resulting grammar after removal of ϵ -production is:

$$S \rightarrow ABC \mid AB \mid BC \mid AC \mid A \mid B \mid C$$

$$A \rightarrow BB \mid B$$

$$B \rightarrow CC \mid C \mid a$$

$$C \rightarrow AA \mid A \mid b$$

Exercise:

Remove ϵ -productions for each of grammar;

1)

$$S \rightarrow AB$$

$$A \rightarrow aAA \mid \epsilon$$

$$B \rightarrow bBB \mid \epsilon$$

Eliminating Unit Production:

A unit production is a production of the form $A \rightarrow B$, where A and B are both variables.

Here, if $A \rightarrow B$, we say B is A-derivable. $B \rightarrow C$, we say C is B-derivable.

Thus if both of two $A \rightarrow B$ and $B \rightarrow C$, then $A \rightarrow^* C$, hence C is also A-derivable.

Here pairs (A, B), (B, C) and (A, C) are called the unit pairs.

To eliminate the unit productions, first find all of the unit pairs. The unit pairs are;

(A, A) is a unit pair for any variable A as $A \rightarrow^* A$

If we have $A \rightarrow B$ then (A, B) is unit pair.

If (A, B) is unit pair i.e. $A \rightarrow B$, and if we have $B \rightarrow C$ then (A, C) is also a unit pair.

Now, to eliminate those unit productions for a, gives grammar say $G = (V, T, P, S)$, we have to find another grammar $G' = (V, T, P', S)$ with no unit productions. For this, we may work out as below;

- ✓ Initialize $P' = P$
- ✓ For each $A \in V$, find a set of A-derivable variables.
- ✓ For every pair (A, B) such that B is A-derivable and for every non-unit production $B \rightarrow \alpha$, we add production $A \rightarrow \alpha$ in P' if it is not in P' already.
- ✓ Delete all unit productions from P' .

Example

Remove the unit production for grammar G defined by productions:

$$P = \{ S \rightarrow S + T \mid T \\ T \rightarrow T * F \mid F \\ F \rightarrow (S) \mid a \\ \};$$

Initialize

1) $P' = \{ S \rightarrow S + T \mid T \\ T \rightarrow t * F \mid F \\ F \rightarrow (S) \mid a \}$

2) Now, find unit pairs;

Here, $S \rightarrow T$ So, (S, T) is unit pair.

$T \rightarrow F$ So, (T, F) is unit pair.

Also, $S \rightarrow T$ and $T \rightarrow F$ So, (S, F) is unit pair.

3) Now, add each non-unit productions of the form $B \rightarrow \alpha$ for each pair (A, B);

$$P' = \{ \\ S \rightarrow S + T \mid T * F \mid (S) \mid a \\ T \rightarrow T * F \mid (S) \mid a \mid F \\ F \rightarrow (S) \mid a \\ \}$$

4) Delete the unit productions from the grammar;

$$P' = \{ \\ S \rightarrow S + T \mid T * F \mid (S) \mid a \\ T \rightarrow T * F \mid (S) \mid a \\ F \rightarrow (S) \mid a \\ \}$$

Exercise

1) **Simplify the grammar $G = (V, T, P, S)$ defined by following productions.**

$$S \rightarrow ASB \mid \epsilon$$

$$A \rightarrow aAS \mid a$$

$$B \rightarrow SbS \mid A \mid bb \mid \epsilon$$

Note: Here simplify means you have to remove all the useless symbol, Unit production and ϵ -productions.

2) **Simplify the grammar defined by following production:**

$$S \rightarrow 0A0 \mid 1B1 \mid BB$$

$$A \rightarrow C$$

$$B \rightarrow S \mid A$$

$$C \rightarrow S \mid \epsilon$$

Chomsky Normal Form

A context free grammar $G = (V, T, P, S)$ is said to be in Chomsky's Normal Form (CNF) if every production in G are in one of the two forms;

$A \rightarrow BC$ and

$A \rightarrow a$ where $A, B, C \in V$ and $a \in T$

Thus a grammar in CNF is one which should not have;

ϵ -production

Unit production

Useless symbols.

Theorem: Every context free language (CFL) without ϵ -production can be generated by grammar in CNF.

Proof:

If all the productions are of the form $A \rightarrow a$ and $A \rightarrow BC$ with $A, B, C \in V$ and $a \in T$, we have done.

Otherwise, we have to do two task as:

1. Arrange that all bodies of length 2 or more consist only of variable
2. Break bodies of length 3 or more into a cascade of production, each with a body consisting of two variable.

The construction for task (1) is as follows

if the productions are of the form: $A \rightarrow X_1, X_2, \dots, X_m$, $m > 2$ and if some X_i is terminal a , then we replace the X_i by C_a having $C_a \rightarrow a$ where C_a is a variable itself.

Thus as result we will have all productions of the form:

$A \rightarrow B_1 B_2 \dots B_m$, $m > 2$; where all B_i 's are non-terminal.

The construction for task (2) is as follows

We break those production $A \rightarrow B_1 B_2 \dots B_m$ for $m \geq 3$, into group of production with two variables in each body. We introduce $m-2$ new variables C_1, C_2, \dots, C_{m-2} . The original production is replaced by the $m-1$ productions

$A \rightarrow B_1 C_1$,

$C_1 \rightarrow B_2 C_2$,

.....

.....

...

$C_{k-2} \rightarrow B_{k-1} B_k$

Finally, all of the productions are achieved in the form as:

$A \rightarrow BC$ or $A \rightarrow a$

This is certainly a grammar in CNF and generates a language without ϵ -productions.

Consider an example:

$S \rightarrow AAC$

$A \rightarrow aAb \mid \epsilon$

$C \rightarrow aC \mid a$

1) Now, removing ϵ - productions;

Here, A is nullable symbol as $A \rightarrow \epsilon$

So, eliminating such ϵ -productions, we have;

$$S \rightarrow AAC \mid AC \mid C$$

$$A \rightarrow aAb \mid ab$$

$$C \rightarrow aC \mid a$$

2) Removing unit-productions;

Here, the unit pair we have is (S, C) as $S \rightarrow C$

So, removing unit-production, we have;

$$S \rightarrow AAC \mid AC \mid aC \mid a$$

$$A \rightarrow aAb \mid ab$$

$$C \rightarrow aC \mid a$$

Here we do not have any useless symbol. Now, we can convert the grammar to CNF. For this;

→ First replace the terminal by a variable and introduce new productions for those which are not as the productions in CNF.

i.e.

$$S \rightarrow AAC \mid AC \mid C_1C \mid a$$

$$C_1 \rightarrow a$$

$$A \rightarrow C_1AB_1 \mid C_1B_1$$

$$B_1 \rightarrow b$$

$$C \rightarrow C_1C \mid a$$

Now, replace the sequence of non-terminals by a variable and introduce new productions.

Here, replace $S \rightarrow AAC$ by $S \rightarrow AC_2$, $C_2 \rightarrow AC$

Similarly, replace $A \rightarrow C_1AB_1$ by $A \rightarrow C_1C_3$, $C_3 \rightarrow AB_1$

Thus the final grammar in CNF form will be as;

$$S \rightarrow AC_2 \mid AC \mid C_1C \mid a$$

$$A \rightarrow C_1C_3 \mid C_1b_1$$

$$C_1 \rightarrow a$$

$$B_1 \rightarrow b$$

$$C_2 \rightarrow AC$$

$$C_3 \rightarrow AB_1$$

$$C \rightarrow C_1C \mid a$$

Exercise

Simplify following grammars and convert to CNF;

1)

$$S \rightarrow ASB \mid \epsilon$$

$$A \rightarrow aAS \mid a$$

$$B \rightarrow SbS \mid A \mid bb$$

2) $S \rightarrow AACD$

$$A \rightarrow aAb \mid \epsilon$$

$$C \rightarrow aC \mid a$$

$$D \rightarrow aDa \mid bDa \mid \epsilon$$

3)
 $S \rightarrow aaaaS \mid aaaa$

Left recursive Grammar

A grammar is said to be left recursive if it has a non-terminal A such that there is a derivation $(A \rightarrow^* A \alpha)$ for some string α (which may variable or terminal).

The top down parsing methods can not handle left recursive grammars, so a transformation that eliminates left recursion is needed.

Removal of immediate left recursion:

Let $A \rightarrow A\alpha \mid \beta$, where β does not start with A . Then, the left-recursive pair of productions could be replaced by the non-left-recursive productions as;

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon; \text{ without changing the set of strings derivable from } A.$$

Or equivalently, these production can be rewritten as with out ϵ - production

$$A \rightarrow \beta A' \mid \beta$$

$$A' \rightarrow \alpha A' \mid \alpha;$$

No matter how many A -productions there are, we can eliminate immediate left recursion from them.

So, in general,

If $A \rightarrow A \alpha_1 \mid A \alpha_2 \mid \dots \mid A \alpha_m \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$; With β_i does not start with A .

Then

we can remove left recursion as:

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \epsilon$$

Equivalently, these productions can be rewritten as:

$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_n A' \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_m A' \mid \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_m$$

Consider an example

$$S \rightarrow AA \mid 0$$

$$A \rightarrow AAS \mid 0S \mid 1 \text{ (Where } AS \text{ is of the form } = \alpha_1, 0S \text{ of the form } = \beta_1 \text{ and } 1 \text{ of the form } = \beta_2)$$

Here, the production $A \rightarrow AAS$ is immediate left recursive. So removing left recursion,

we have;

$$S \rightarrow AA \mid 0$$

$$A \rightarrow 0SA' \mid 1A'$$

$$A' \rightarrow ASA' \mid \epsilon$$

Equivalently, we can write it as:

$$S \rightarrow AA \mid 0$$

$$A \rightarrow 0SA' \mid 1A' \mid 0S \mid 1$$

$$A' \rightarrow ASA' \mid AS$$

Exercise

For the following grammar

$$E \rightarrow E+T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid a$$

Remove the Left recursion.

Greibach Normal Form (GNF)

A grammar $G = (V, T, P \text{ and } S)$ is said to be in Greibach Normal Form, if all the productions of the grammar are of the form:

$$A \rightarrow a \alpha, \text{ where } a \text{ is a terminal, i.e. } a \in T^* \text{ and } \alpha \text{ is a string of zero or more variables. i.e. } \alpha \in V^*$$

So we can rewrite as:

$$A \rightarrow aV^* \text{ with } a \in T^*$$

Or,

$$A \rightarrow aV^+$$

$$A \rightarrow a \text{ with } a \in T$$

This form called Greibach Normal Form, after Sheila Greibach, Who first gave a way to construct such grammars. Converting a grammar to this form is complex, even if we simplify the task by, say, starting with a CNF grammar.

To convert a grammar into GNF;

- Convert the grammar into CNF at first.
- Remove any left recursions.
- Let the left recursion tree ordering is $A_1, A_2, A_3, \dots, A_p$
- Let A_p is in GNF
- Substitute A_p in first symbol of A_{p-1} , if A_{p-1} contains A_p . Then A_{p-1} is also in GNF.
- Similarly substitute first symbol of A_{p-2} by A_{p-1} production and A_p production and so on.....

Consider an Example;

$$S \rightarrow AA \mid 0$$

$$A \rightarrow SS \mid 1$$

This Grammar is already in CNF.

Now to remove left recursion, first replace symbol of A-production by S-production (since we do not have immediate left recursion) as:

$$S \rightarrow AA \mid 0$$

$$A \rightarrow AAS \mid 0S \mid 1 \text{ (Where } AS \text{ is of the form } = \alpha 1, 0S \text{ of the form } = \beta 1 \text{ and } 1 \text{ of the form } = \beta 2)$$

Now, removing the immediate left recursion;

$$\begin{aligned} S &\rightarrow AA \mid 0 \\ A &\rightarrow 0SA' \mid 1A' \\ A' &\rightarrow ASA' \mid \epsilon \end{aligned}$$

Equivalently, we can write it as:

$$\begin{aligned} S &\rightarrow AA \mid 0 \\ A &\rightarrow 0SA' \mid 1A' \mid 0S \mid 1 \\ A' &\rightarrow ASA' \mid AS \end{aligned}$$

Now we replace first symbol of S-production by A-production as;

$$\begin{aligned} S &\rightarrow AASA' \mid 1A'A \mid 0SA \mid 1A \mid 0 \\ A &\rightarrow 0SA' \mid 1A' \mid 0S \mid 1 \\ A' &\rightarrow ASA' \mid AS \end{aligned}$$

Similarly replacing first symbol of A'-production by A-production, we get the grammar in GNF as;

$$\begin{aligned} S &\rightarrow 0SA'A \mid 1A'A \mid 0SA \mid 1A \mid 0 \\ A &\rightarrow 0SA' \mid 1A' \mid 0S \mid 1 \\ A' &\rightarrow 0SA'SA' \mid 1A'SA' \mid 0SSA' \mid 1SA' \mid 0SA'S \mid 1A'S \mid 0SS \mid 1S \end{aligned}$$

Regular Grammar:

A regular grammar represents a language that is accepted by some finite automata called regular language. A regular grammar is a CFG which may be either left or right linear. A grammar in which all of the productions are of the form $A \rightarrow wB$ (wB is of the form a) or $A \rightarrow w$ for $A, B \in V$ and $w \in T^*$ is called left linear.

Equivalence of Regular Grammar and Finite Automata

1. Let $G = (V, T, P \text{ and } S)$ be a right linear grammar of a language $L(G)$. we can construct a finite automata M accepting $L(G)$ as;

$$M = (Q, T, \delta, [S], \{[\epsilon]\})$$

Where,

Q – consists of symbols $[\alpha]$ such that α is either S or a suffix from right hand side of a production in P .

T – is the set of input symbols which are terminal symbols of G .

$[S]$ – is a start symbol of G which is start state of finite automata.

$\{[\epsilon]\}$ – is the set of final states in finite automata.

And δ is defined as:

If A is a variable then, $\delta([A], \epsilon) = [\alpha]$ such that $A \rightarrow \alpha$ is a production.

If 'a' is a terminal and α is in $T^* \cup T^*V$ then $\delta([a\alpha], a) = [\alpha]$.

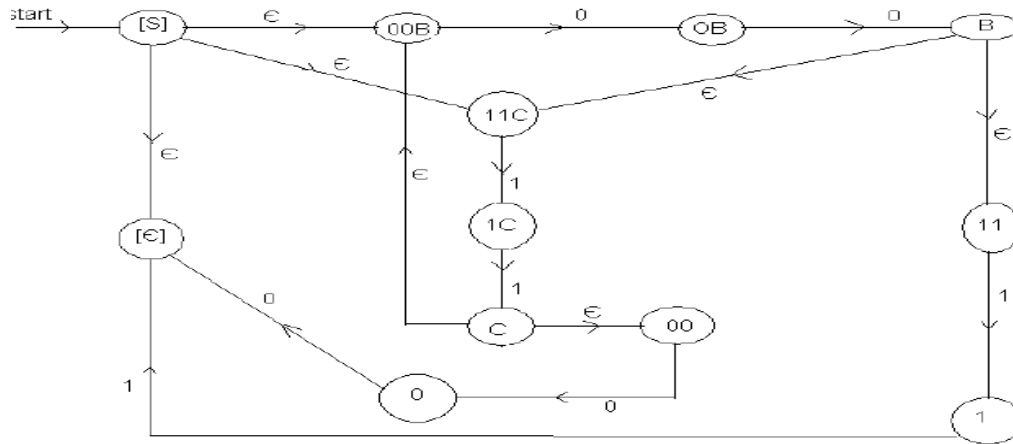
For example;

$$S \rightarrow 00B \mid 11C \mid \epsilon$$

$$B \rightarrow 11C \mid 11$$

$$C \rightarrow 00B \mid 00$$

Now the finite automata can be configured as;



Another way:

If the regular grammar is of the form in which all the productions are in the form as;
 $A \rightarrow aB$ or $A \rightarrow a$, where $A, B \in V$ and $a \in T$ Then, following steps can be followed to obtain an equivalent finite automaton that accepts the language represented by above set of productions.

- a. The number of states in finite automaton will be one more than the number of variables in the grammar. (i.e. if grammar contain 4 variables then the automaton will have 5 states)
 - Such one more additional state in the final state of the automaton.
 - Each state in the automaton is a variable in the grammar.
- b. The start symbol of regular grammar is the start state of the finite automaton.
- c. If the grammar contains ϵ as $S \rightarrow *C$, S being start symbol of grammar, then start state of the automaton will also be final state.
- d. The transition function for the automaton is defined as;
 - For each production $A \rightarrow aB$, we write $\delta(A, a) = B$. So there is an arc from state A to B labeled with a .
 - For each production $A \rightarrow a$, $\delta(A, a) = \text{final state of automaton}$
 - For each production $A \rightarrow \epsilon$, make the state A as final state.

Exercise

Write the equivalent finite automata

1.
 $S \rightarrow abA \mid bB \mid aba$
 $A \rightarrow b \mid aB \mid bA$
 $B \rightarrow aB \mid aA$

2.
 $S \rightarrow 0B \mid 11C \mid \epsilon$
 $B \rightarrow 00C \mid 11$
 $C \rightarrow 00B \mid 0D$
 $D \rightarrow 0B \mid 00$

3.

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

4.

 $A \rightarrow 0B \mid 1D \mid 0$ $B \rightarrow 0D \mid 1C \mid \epsilon$ $C \rightarrow 0B \mid 1D \mid 0$ $D \rightarrow 0D \mid 1D$

5.

 $A \rightarrow 0B \mid E \mid 1C$ $B \rightarrow 0A \mid F \mid \epsilon$ $C \rightarrow 0C \mid 1A \mid 0$ $E \rightarrow 0C \mid 1A \mid 1$ $F \rightarrow 0A \mid 1B \mid \epsilon$ **Pumping lemma for context free Language:**

The “pumping lemma for context free languages” says that in any sufficiently long string in a CFL, it is possible to find at most two short, nearby substrings that we can “pump” in tandem (one behind another). i.e. we may repeat both of the strings I times, for any integer I , and the resulting string will still be in the language. We can use this lemma as a tool for showing that certain languages are not context free.

Size of Parse tree: Our first step in pumping lemma for CFL’s is to examine the size of parse trees, During the lemma, we will use the grammar in CNF form. One of the uses of CNF is to turn parses into binary trees. Any grammar in CNF produces a parse tree for any string that is binary tree. These trees have some convenient properties, one of which can be exploited by following theorem.

Theorem: Let a terminal string w be a yield of parse tree generated by a grammar $G = (V, T, P \text{ and } S)$ in CNF. If length path is n , then $|w| \leq 2^{n-1}$

Proof: Let us prove this theorem by simple induction on n .

Basis Step: Let $n = 1$, then the tree consists of only the root and a leaf labeled with a terminal.

So, string w is a simple terminal.

Thus $|w| = 1 = 2^{1-1} = 2^0$ which is true.

Inductive Step: Suppose n is the length of longest path and $n > 1$. Then the root of the tree uses a production of the form; $A \rightarrow BC$, since $n > 1$

No path in the sub-tree rooted at B and C can have length greater than $n-1$ since B & C are child of A . Thus, by inductive hypothesis, the yield of these sub-trees are of length at most 2^{n-2} .

The yield of the entire tree is the concatenation of these two yields and therefore has length at most,

$$2^{n-2} + 2^{n-2} = 2^{n-1}$$

$$|w| \leq 2^{n-1}. \text{ Hence proved.}$$

Pumping Lemma

It states that every CFL has a special value called pumping length such that all longer strings in the language can be pumped. This time the meaning of pumped is a bit more complex. It means the string can be divided into five parts so that the second and the fourth parts may be repeated together any number of times and the resulting string still remains in the language.

Theorem: Let L be a CFL. Then there exists a constant n such that if z is any string in L such that $|z|$ is at least n then we can write $z=uvwxy$, satisfying following conditions.

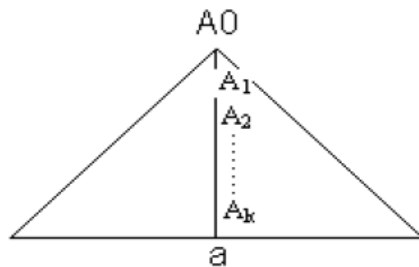
- I) $|vx| > 0$
- II) $|vwx| \leq n$
- III) For all $i \geq 0$, $uv^iwx^iy \in L$. i.e. the two strings v & x can be “pumped” any number of times, including ‘0’ and the resulting string will be still a member of L .

Proof: First we can find a CNF grammar for the grammar G of the CFL L that generates $L - \{\epsilon\}$.

Let m be the number of variables in this grammar choose $n = 2^m$. Next suppose z in L is of length at least n . i.e. $|z| \geq n$

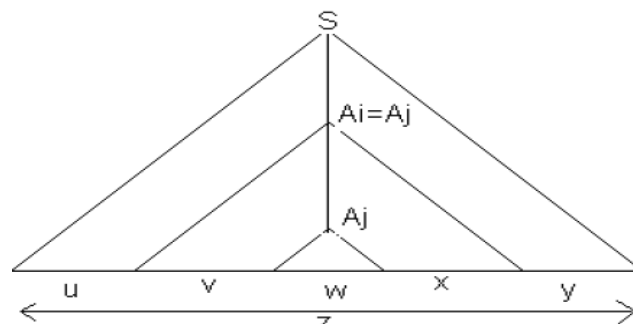
Any parse tree for z must have height $m+1$, other if it would be less than $m+1$, i.e. say m then by the lemma for size of parse tree, $|z| = 2^{m-1} = 2^m/2 = n/2$ is contradicting. So it should be $m+1$.

Let us consider a path of maximum length in tree for z , as shown below, where k is the least m and path is of length k .



Since $k \geq m$. there are at least $m+1$ occurrence of variables $A_0, A_1, A_2, \dots, A_k$ on the path. But there are only m variables in the grammar, so at least two of the last $m+1$ variable on the path must be same, (by pigeonhole principle).

Suppose $A_i = A_j$, then it is possible to divide tree,



String w is the yield of sub-tree noted at A_j , string v and x are to left and right of w in yield of larger sub-tree rooted at A_i . If u and y represent beginning and end z i.e. to right and left of sub-tree rooted at A_i , then

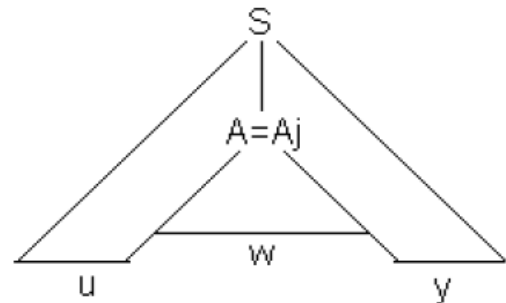
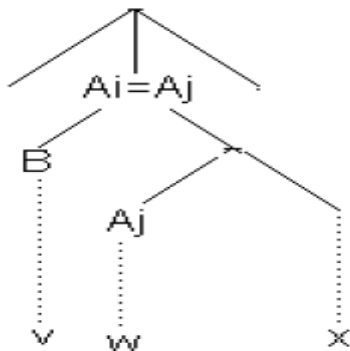
$$Z = uvwxy$$

Since, there are not unit productions so v & x both could not be ϵ , i.e. empty. Means that A_i has always two children corresponding to variables. If we let B denotes the one tree is not ancestor of A_j , then since

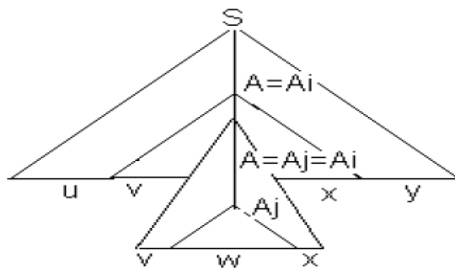
w is derived from A_j then strings of terminal derived from B does not overlap x . it follows that either v or x is not empty. So,

$$|vx| > 0$$

Now, we know $A_i = A_j = A$ say, as we found any two variables in the tree are same. Then, we can construct new parse tree where we can replace sub-tree rooted at A_i , which has yield vwx , by sub-tree rooted at A_j , which has yield w . the reason we can do so is that both of these trees have root labeled A . the resulting tree yields uv^0wx^0y as;



Another option is to replace sub-tree rooted at A_j by entire sub-tree rooted at A_i . Then the yield can be of pattern uv^2wx^2y as



Hence, we can find any yield of the form $uv^iwx^i y \in L$ for any $i \geq 0$ Now, to show $|vwx| \leq n$. the A_i in the tree is the portion of path which is root of vwx .

Since we begin with maximum path of length with height $m+1$. this sub-tree rooted at A_i also have height of less than $m+1$.

So by lemma for height of parse tree,

$$|vwx| \leq 2^{m+1-1} = 2^m = n$$

$$|vwx| \leq n$$

Thus, this completes the proof of pumping lemma.

Example

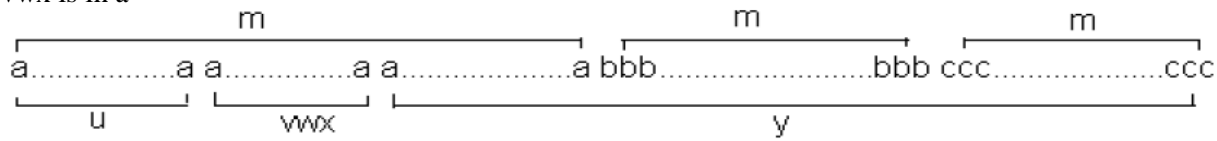
Show that $L = \{a^n b^n c^n : n \geq 0\}$ is not a CFL

To prove L is not CFL, we can make use of pumping lemma for CFL. Let L be CFL. Let any string in the language is $a^m b^m c^m$

$$z = uvwxy, |vwx| \leq m, |vx| > 0$$

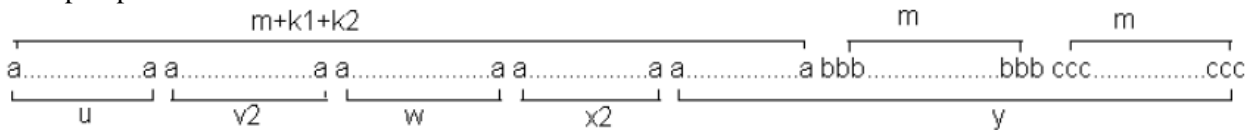
Case I:

vw^kx is in a^m



$v = a^{k_1}, x = a^{k_2}; k_1+k_2 > 0$ i.e. $k_1+k_2 \geq 1$

Now pump v and x then



But by pumping lemma, $uv^2wx^2y = a^{m+k_1+k_2}b^m c^m$ does not belong to L as, k_1+k_2 . Hence, it shows contradiction that L is CFL. So, L is not CFL.

Note: other cases can be done similarly.

Bakus-Naur Form

This is another notation used to specify the CFG. It is named so after John Bakus, who invented it, and Peter Naur, who refined it. The Bakus-Naur form is used to specify the syntactic rule of many computer languages, including Java. Here, concept is similar to CFG, only the difference is instead of using symbol “ \rightarrow ” in production, we use symbol $::=$. We enclose all non-terminals in angle brackets, $\langle \rangle$.

For example:

The BNF for identifiers is as;

- $\langle \text{identifier} \rangle ::= \langle \text{letter or underscore} \rangle | \langle \text{identifier} \rangle | \langle \text{symbol} \rangle$
- $\langle \text{letter or underscore} \rangle ::= \langle \text{letter} \rangle | \langle _ \rangle$
- $\langle \text{symbol} \rangle ::= \langle \text{letter or underscore} \rangle | \langle \text{digit} \rangle$
- $\langle \text{letter} \rangle ::= a | b | \dots | z$
- $\langle \text{digit} \rangle ::= 0 | 1 | 2 | \dots | 9$

Closure Property of Context Free Languages

Given certain languages are context free, and a language L is formed from them by certain operation, like union of the two, then L is also context free. These theorem lemma are often called closure properties of context free languages, since they show whether or not the class of context free language is closed under the operation mentioned.

Closure properties express the idea that one (or several) languages are context free, and then certain related languages are also context free or not.

Here are some of the principal closure properties for context free languages;

A. The context free language are closed under union

i.e. Given any two context free languages L_1 and L_2 , their union $L_1 \cup L_2$ is also context free language.

Proof

The inductive idea of the proof is to build a new grammar from the original two, and from start symbol of the new grammar have productions to the start symbols of the original two grammars.

Let $G_1 = (V_1, T_1, P_1 \text{ and } S_1)$ and $G_2 = (V_2, T_2, P_2 \text{ and } S_2)$ be two context free grammars defining the languages $L(G_1)$ and $L(G_2)$. Without loss of generality, let us assume that they have common terminal set T , and disjoint set of nonterminals. Because, the non-terminals are distinct so the productions P_1 and P_2 .

Let S be a new non-terminal not in V_1 and V_2 . Then, construct a new grammar $G = (V, T, P, S)$ where;

$$V = V_1 \cup V_2 \cup \{S\}$$

$$P = P_1 \cup P_2 \cup \{S \rightarrow S_1, S \rightarrow S_2\}$$

G is clearly a context free grammar because the two new productions so added are also of the correct form, we claim that $L(G) = L(G_1) \cup L(G_2)$.

For this, Suppose that $x \in L(G_1)$. Then there is a derivation of x ;

$$S_1 \rightarrow^* x$$

But in G we have production, $S \rightarrow S_1$

So there is a derivation of x also in G as:

$$S \rightarrow S_1 \rightarrow^* x$$

Thus, $x \in L(G)$. Therefore, $L(G_1) \subseteq L(G)$. A similar argument shows $L(G_2) \subseteq L(G)$.

So, we have,

$$L(G_1) \cup L(G_2) \text{ is subset of } L(G)$$

Conversely, suppose that $x \in L(G)$. Then there is a derivation of x in G as:

$$S \rightarrow \beta \rightarrow^* x$$

Because of the way in which P is constructed, β must be either S_1 or S_2 .

Suppose $\beta = S_1$. Any derivation in G of the form $S_1 \rightarrow^* x$ must involve only productions of G_1 so, $S_1 \rightarrow^* x$ is a derivation of x in G_1 .

Hence, $\beta = S_1 \rightarrow x \in L(G_1)$. A simpler argument proves that $\beta = S_1$ implies $x \in L(G_1)$.

Thus $L(G)$ is subset of $L(G_1) \cup L(G_2)$.

It follows that $L(G) = L(G_1) \cup L(G_2)$

Similarly following two closure properties of CFL can be proved as for the union we have proved.

B. The CFLs are closed under concatenation

C. The CFLs are closed under kleen closure [prove yourself]

However context free languages are not closed under some cases like, intersection and complementation.

The context free languages are not closed under intersection.

To prove this, the idea is to exhibit two CFLs whose intersection results in the language which is not CFL.

For this as we learned in pumping lemma that

$L = \{a^n b^n c^n \mid n \geq 1\}$ is not CFL

However following two languages are context free;

$L1 = \{a^n b^n c^i \mid n \geq 1, i \geq 1\}$

$L2 = \{a^i b^n c^n \mid n \geq 1, i \geq 1\}$

Clearly,

$L = L1 \cap L2$. to see why, observe that L1 requires that there be the same number of a's + b's, while L2 requires the number of b's and c's to be equal. A string in both must have equal number of all three symbols and thus to be in L.

If the CFL's were closed under intersection, then we could prove the false statement that L is context free. Thus, we conclude by contradiction that the CFL's are not closed under intersection.