

# Context Free Grammars

3.1

## Introduction to Context free Grammars

Consider a given sentence she ran fast. The rules to form this sentence can be written as

sentence  $\rightarrow$  <pronoun> <verb> <adverb>

pronoun  $\rightarrow$  she

verb  $\rightarrow$  ran

Adverb  $\rightarrow$  fast

These rules are called production. Here, the words sentence, pronoun, verb and Adverb are called non-terminals or variables. Each production starts with non-terminal followed by an arrow followed by combination of non-terminals or terminals. The left hand side of the 1<sup>st</sup> production is called the start symbol. The set of rules to form a sentence which in turn used to generate a language is called grammar.

A gramm G can be define as a 4 tuple or quadruple,  $G = \{V, T, P, S\}$  where

V is the set of variables or non-terminals

T is set of terminals

P is set of rules called production

S is the start symbol.

In the above example sentence, pronoun, verb, adverb is set of non-terminals

or variables {she, ran, fast} is the set of terminals.

$\langle \text{sentence} \rangle \rightarrow \langle \text{pronoun} \rangle \langle \text{verb} \rangle \langle \text{Adverb} \rangle$  is production and sentence is the start symbol.

Each production is of the form  $\alpha \rightarrow \beta$  where  $\alpha$  is a non-empty string of terminals and/or non-terminals and  $\beta$  is string of terminals and/or non-terminal including the null string i.e.  $\alpha$  is string in  $(VUT)^*$  and  $\beta$  is a string in  $(VUT)^+$ . In production, all the lower case letters, digits or any other special characteristics are considered as terminals and all the upper case letters are considered as non-terminal or variables. It should be noted that non-terminals can be replaced with string of terminals or non-terminals whereas the terminals cannot be replaced or substituted.

## Production

→ Consider the production shown below:

$$A \rightarrow aA/a$$

Now, consider, a string  $w$  of the form,

$$w = xAx$$

Here, the non-terminals  $A$  can be replaced

by ' $\alpha A'$  or ' $\alpha'$ . If  $A$  is replaced by  $\alpha A$ , we have,

$$w = x\alpha A X$$

This can be written as

$$X A X \Rightarrow X \alpha A X$$

This means the  $XAX$  driven  $X\alpha AX$  in one step i.e. by applying one production.

A production can be applied any number of times in an arbitrary order.

In general can be written as,

$$w \Rightarrow w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n$$

A string  $w_n$  is derived from and can be written as

$$w^* \Rightarrow w_n$$

Here, '\*' indicates that unspecified number of production (including applying no production) unspecified number of production have been applied to a get number of string  $w_n$  from  $w$ .

Also,

$$w^+ \Rightarrow w_n$$

Here, '+' indicates at least one or more product are applied to get the string  $w_n$  from  $w$ .

Derivation :-

Consider,

$$A \rightarrow \alpha B \gamma$$

and  $B \rightarrow \beta$  are productions in grammar G. Where,  $\alpha$ ,  $\beta$  and  $\gamma$  are string of terminals/non-terminals and A and B are non-terminals or variables. The non-terminals A directly derives the string  $\alpha \beta \gamma$  by replacing the non-terminal B in  $\alpha B \gamma$  by the string  $\beta$  by applying the production  $B \rightarrow \beta$  and can be written as

$$A \Rightarrow \alpha \beta \gamma$$

Thus the process of obtaining string terminal and/or non-terminals from the start symbol by applying some or all production is called derivations.

If a string is obtained by applying only one production, then it is called one step derivation and is denoted by  $\Rightarrow$ .

If one or more productions are applied to get a string  $\alpha \beta \gamma$  from A, then we write.

$$A \xrightarrow{+} \alpha \beta \gamma$$

If unspecified number of productions are applied (including applying no production) to get string  $\alpha \beta \gamma$  from A, then we write

$$A \Rightarrow^* \alpha \beta \gamma$$

Example:- Consider the grammar shown from which any arithmetic expression can be obtained:-

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

$$E \rightarrow id$$

Obtain the string  $id + id * id - id$  from above grammar.

$$\text{Ans} \rightarrow E \rightarrow E + E$$

$$\rightarrow id + E$$

Applying  $E \rightarrow id$

$$\rightarrow id + E * E$$

$$E \rightarrow E * E$$

$$\rightarrow id + id * E - E$$

$E \rightarrow id$  and  $E \rightarrow E - E$

$$\rightarrow id + id * id - E$$

"  $E \rightarrow id$

$$\rightarrow id + id * id - id$$

"  $E \rightarrow id$

## \* Language :-

The given grammar can be generated the set of string consisting of only terminals by applying productions in different order. The set of such string is called language.

Let,  $G = (V, T, P, S)$  be a grammar.

The language  $L(G)$  generated by a grammar  $G$  is

$$L(G) = \{ w \mid s \xrightarrow{*} w \in T^* \}$$

Date \_\_\_\_\_  
Page \_\_\_\_\_

i.e.,  $w$  is string of only terminals obtained from the start symbol  $s$  by applying an arbitrary number of productions. The intermediate of terminals or/and non-terminals obtained during the derivation process is called sentential form of  $G$ .

Q.L Let,  $G = (V, T, P, S)$

where,  $V = \{S, C\}$

$T = \{a, b\}$

$S$  is the start symbol.

$P$  is given as

$S \rightarrow aCa$

$C \rightarrow aCa \mid b$

Obtained the language generated by above grammar

Soln

$S \rightarrow aCa$

$S \rightarrow aba$  (By applying  $C \rightarrow b$ )

Again,

$S \rightarrow aCa$

$\rightarrow aaCaa$  (By applying  $C \rightarrow aCa$ )

$\rightarrow aaaCaaa$

$\rightarrow aaa...b...aaa$  (by applying  $C \rightarrow b$ )

Therefore the language is generated by a given grammar is

$$L(G) = \{a^n b a^n \mid n \geq 1\}$$

Q.2. Write a regular grammar that generates all set of palindromes over the  $\Sigma = \{a, b\}$ .

Soln → The recursive definition of palindrome is

(i)  $E$  or  $\lambda$

(ii)  $a$  and  $b$  are palindromes

(iii) If  $w$  is palindromes, then string  $awa$ ,  $bwb$  are palindromes.

Production corresponding to above definition are given as

$$S \rightarrow E$$

$$S \rightarrow a/b$$

$$S \rightarrow aSa/bSb$$

Hence, the grammar,  $G = (V, T, P, S)$  where,

$V$  = non-terminal/variable,  $V = \{S\}$

$T$  = Terminal,  $T = \{a, b\}$

$P$  = Production

$$S \rightarrow E$$

$$S \rightarrow a/b$$

$$S \rightarrow aSa/bSb$$

$S$  = start symbol

Q.3: Obtain a Grammar to generate the language  
 $L = \{0^n 1^{n+1} \mid n \geq 0\}$

Solu<sup>n</sup> → The production to generate,  $L = \{0^n 1^{n+1} \mid n \geq 0\}$   
 is  $A \rightarrow 0A1 \mid E$   
 To get the string  $0^n 1^{n+1}$  and extract  $0^n$  should be placed. Production of this is given as

$$S \rightarrow A1$$

Hence, the grammar G to generate the language L, will be  $G = (V, T, P, S)$  where,

$$V = \{S, A\}$$

$$T = \{0, 1\}$$

$$P = \{S \rightarrow A1, A \rightarrow 0A1 \mid E\}$$

where, S is the start symbol.

Q.4: Obtain the grammar to generate the language

$$L = \{w \mid n_a(w) = n_b(w)\}$$

Solu<sup>n</sup> → The production to generate the language  $L = \{w \mid n_a(w) = n_b(w)\}$  is

$$S \rightarrow aSb \mid bSa$$

To get the string  $n_a(w) = n_b(w)$  i.e. n no. of a equal to n no. of b. The production of this

Q.3. Obtain a Grammar to generate the language

$$L = \{0^n 1^{n+1} \mid n \geq 0\}$$

Solu<sup>n</sup> → The production to generate,  $L = \{0^n 1^{n+1} \mid n \geq 0\}$  is

$$A \rightarrow 0A1 \mid E$$

To get the string  $0^n 1^{n+1}$  and extract  $0_1$  is  
should be placed. Production of this is  
given as

$$S \rightarrow SA1$$

Hence, the grammar  $G$  to generate the  
language  $L$ , will be  $G = (V, T, P, S)$  where,

$$V = \{S, A\}$$

$$T = \{0, 1\}$$

$$P = \{$$

$$S \rightarrow A1$$

$$A \rightarrow DA1 \mid E$$

where,  $S$  is the start symbol.

Q.4. Obtain the grammar to generate the language

$$L = \{w \mid n_a(w) = n_b(w)\}$$

Solu<sup>n</sup> → The production to generate the language

$$L = \{w \mid n_a(w) = n_b(w)\}$$
 is

$$S \rightarrow aSb \mid bSa$$

To get the string  $n_a(w) = n_b(w)$  i.e. no. of  $a$   
equal to no. of  $b$ . The production of this

given as

$$S \rightarrow SS$$

Hence, the grammar  $G$  to generate the language  $L$ , will be.

$$G = (V, T, P, S) \text{ where}$$

$$V = \{S\}$$

$$T = \{a, b\}$$

$$P = \{$$

$$S \rightarrow SS$$

$$S \rightarrow aSb \mid bSa$$

3

where  $S$  is the start symbol.

Q.5. Obtain grammar to obtain integers.

→ The sign of number can  $+ | - | \epsilon$ . The production for this can be written as

$$S \rightarrow + \mid - \mid \epsilon$$

The production to obtained the digits of a number can be written as

$$D \rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 8 \mid 9$$

The production to obtained a number  $N$  can be written as,

$$N \rightarrow D \text{ i.e. a no. } N \text{ is a digit}$$

$$N \rightarrow DN \text{ i.e. the no. } N \text{ followed by } D$$

also a number.

We know that an integer number can be a number N or sign of a number followed by number N. The production for this can be written as

$$I \rightarrow N|SN$$

Hence, the grammar G to obtain integer number can be written as,

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

where,

$$V = \{S, D, N, I\}$$

$$T = \{+, -, 0, 1, \dots, 9\}$$

$$P = \{$$

$$I \rightarrow N|SN$$

$$N \rightarrow D$$

$$N \rightarrow DN$$

$$S \rightarrow +|-|E$$

$$D \rightarrow 0|1|\dots|8|9$$

Where, I is the start symbol.



## Context free Grammar (CFG)

The production in the regular grammar are restricted in two ways : ~~left~~ the left side must be a single variable i.e. non terminal while the right side have special form. The right hand side of the production

may contain 0 or more terminals. If non-terminals is present, that non-terminals should be left most symbol of the right most symbol ~~by~~ but, in a Context Free Grammar there is only one restriction that is the symbol on the left hand side of the production should be a single non-terminal but there is no restriction on the right hand side of the production.

A Grammar  $G = \{V, T, P, S\}$  is said to be Context Free Grammar (CFG) or type 2 grammar if all the productions are of the form:

$$A \rightarrow \alpha$$

where,  $A \in V$  and

$$\alpha \in (V \cup T)^*$$

Here, the null string i.e.  $\epsilon$  can appear on the R.H.S of the any production. There is only one symbol A on the left hand side of the production and that symbol must be non-terminal.

$\alpha \in (V \cup T)^*$  implies that R.H.S of the string may contain any no. of terminals and non-terminal including the null character.

Every regular grammar is a Context Free Grammar and hence a regular language is Context Free language (CFL) but the reverse is not true always. Regular grammar is subset of Context

Free Grammar and the Regular language is subset of Context free language.

The language generated by Context free Grammar (CFG) is called Context free language (CFL). CFL is the most important aspect of formal language because it is applied to programming language. Actual programming languages have many features that can be described ~~most~~ elegantly by means of Context free languages. What formal language theory tells us about context free languages, has important application in the design of programming languages, as well as in the constructions of efficient compilers.

### Q.1. Obtain CFG for the language

$$L = \{a^m b^n \mid m \neq n\}$$

Soln It is possible to construct a CFG for the given language. Let us assume that  $m=n$ . The production for this given that

$$S \rightarrow aSb|e$$

Now, consider the following two cases for the given language.

case-I ( $m > n$ )

When no. of a's is more than that of no. of b's. The extra a's can be generated using the production.

$$\textcircled{a} \quad A \rightarrow aA|a$$

This production append towards the left of the string generated from the production,

$$S \rightarrow aSb|E$$

Then, we get production as,

$$S' \rightarrow AS$$

case-II ( $m < n$ )

When ( $m < n$ ) i.e. no of a's is less than that of no. of b's. The extra b's can be generated using the production

$$B \rightarrow bB|b$$

This production should be appended towards right of the string generated from the production

$$S \rightarrow aSb|E$$

Then, we get the production as,

$$S' \rightarrow S B$$

The context free Grammar  $G = (V, T, P, S)$

where,

$$V = \{S, A, S', B\}$$

$$T = \{a, b\}$$

$$P = \{$$

$$S' \rightarrow AS|SB$$

$$S \rightarrow aSb|E$$

$$A \rightarrow aA|a$$

$$B \rightarrow bB|b$$

y

$$S \rightarrow SB$$

$$\rightarrow aSbB$$

$$\rightarrow abB$$

$$\rightarrow abbB$$

$$\rightarrow abbb$$

where,  $S'$  is the start symbol.

Q.2. Obtain CFG for the regular expression  
 $(011+1)^*(01)^*$

Solu<sup>n</sup> → The production to generate  
 $011+1$  can be written as

$$A \rightarrow 011|1$$

$$B \rightarrow E$$

$$B \rightarrow AB$$

$$A \rightarrow 011|1$$

The regular expression  $(011+1)^*$

may be an empty language or

string consisting of any combination

of 011's and 1's. The production

for this can be written as

$$S \rightarrow BD$$

$$B \rightarrow AB|E$$

$$A \rightarrow 011|1$$

$$D \rightarrow CD|E$$

$$B \rightarrow E$$

$$B \rightarrow AB$$

$$A \rightarrow 011|1|E$$

Now, the production to generate the string  
 $01$  can be written as

$C \rightarrow 01$

The production to generate the string  $(01)^*$  can be written as

$D \rightarrow E$

$D \rightarrow CD$

$C \rightarrow 01 | E$

Now, the production to get the regular expression can be written as

$S \rightarrow BD$

Hence, the context free grammar for the regular expression is given by

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

where,

$$V = \{S, A, B, C, D\}$$

$$T = \{0, 1\}$$

$$P = \{S \rightarrow BD, B \rightarrow AB | E, A \rightarrow 011 | 1, C \rightarrow 01, D \rightarrow CD | E\}$$

$$S \rightarrow BD$$

$$B \rightarrow AB | E$$

$$A \rightarrow 011 | 1$$

$$D \rightarrow CD | E$$

$$C \rightarrow 01$$

?

where,  $S$  is the start symbol.



## left most and Right most derivation:

A derivation is said to be left most if left most variable, i.e. left most non-terminal in the sentential form which replaced every steps. Derivation is said to be right most derivation if right most non-terminal in sentential form is replaced at every steps.

$$S \rightarrow BD$$

$$B \rightarrow AB | E$$

$$A \rightarrow OI | I$$

$$D \rightarrow CD | E$$

$$C \rightarrow OI$$



## Derivation Tree (Parse Tree)

The left derivation or right derivation can be represented in the form of a tree such a tree is known as derivation tree or parse tree. The formal definition of Derivation Tree (or Parse Tree) is given as:

$$\text{Let } G = (V, T, P, S)$$

be a Context-free Grammar and an ordered tree for  $G$  if and only if it has a following property

- (i) The root has the label  $S$ .
- (ii) Every vertex has label which is in  $(V \cup T)^*$
- (iii) Every leaf node has label from terminal ( $T$ ) and an interior vertex has a label from  $V$ .

(iv) If a vertex is labeled A and if

$x_1, x_2, x_3, \dots, x_n$  are all children of left, then P must contain production of the form

$$A \rightarrow x_1 x_2 x_3 x_4 \dots x_n$$

Example:

The production of the grammar

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow id$$

To obtain the expression  $(id * id) + id$  the derivation is given as,

$$E \Rightarrow E + E$$

$$\Rightarrow (E) + E$$

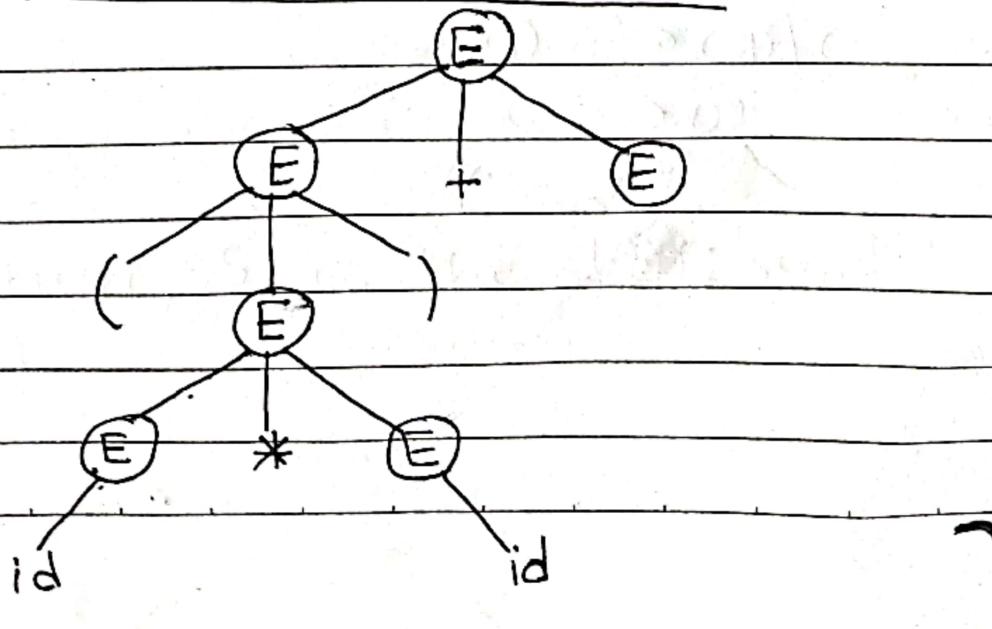
$$\Rightarrow (E * E) + E$$

$$\Rightarrow (id * E) + E$$

$$\Rightarrow (id * id) + E$$

$$\Rightarrow (id * id) + id$$

Left most derivation tree



## ★ Partial Parse Tree:-

Let  $G = (V, T, P, S)$  be CFG. The tree is partial derivation tree if it has the following properties:

- (i) The root has the label  $S$ .
- (ii) Every vertex has a label which is in ~~(VUTUE)~~ (VUTUE)
- (iii) Every leaf node has a label from (VUTUE). This property differs from that a derivation tree.
- (iv) If a vertex is labelled  $A$  and if  $x_1, x_2, x_3$  are all children of left, then  $P$  must contain production of the form.

$$A \rightarrow x_1 x_2 x_3 x_4$$

Example:- Consider the production of the given grammar as

$$E \rightarrow (\text{E})$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E / E$$

$$E \rightarrow id$$

If the left most partial derivation is given

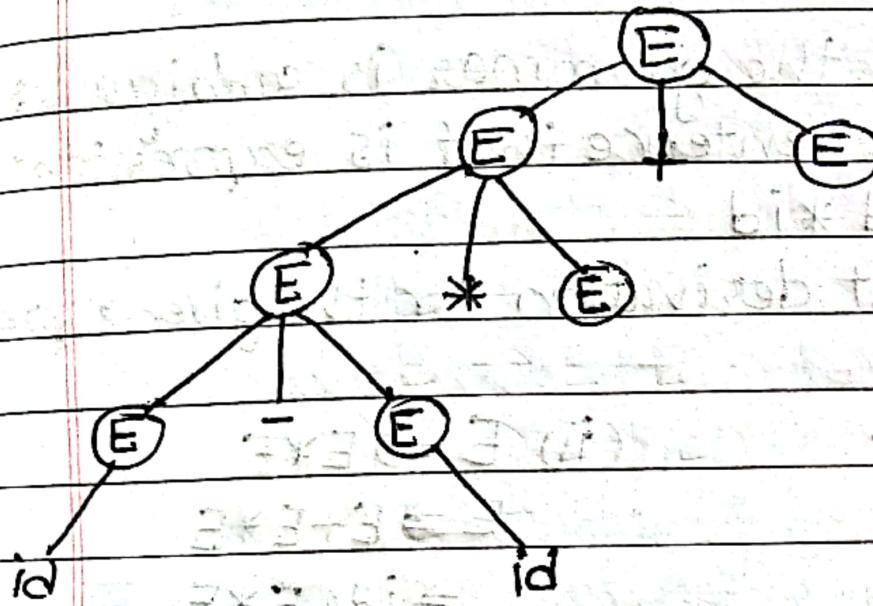
$$E \rightarrow E + E$$

$$\Rightarrow E * E + E$$

$$\Rightarrow E - E * E + E$$

$$\Rightarrow id - E * E + E$$

$$\Rightarrow id - id * E + E$$



## Ambiguous Grammar

A Grammar that has more than one leftmost derivation or more than one rightmost derivation for the same sentence is called Ambiguous Grammar. This means that an Ambiguous grammar has more than one parse tree for the same sentence either by applying left most derivation twice or by applying right most derivation twice. Consider the production for the grammar for arithmetic expression.

$$P = \{$$

$$E \rightarrow (E)$$

$$E \rightarrow E + E$$

$$E \rightarrow E - E$$

$$E \rightarrow E * E$$

$$E \rightarrow E/E$$

$$E \rightarrow id$$

Show that the grammar is ambiguous.

→ Consider the sentence that is expression.

$id + id * id$

The left most derivation for the given expression is given as

(i)  $E \Rightarrow E+E$

$$\Rightarrow id + E$$

$$\Rightarrow id + E * E$$

$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

(ii)  $E \Rightarrow E * E$

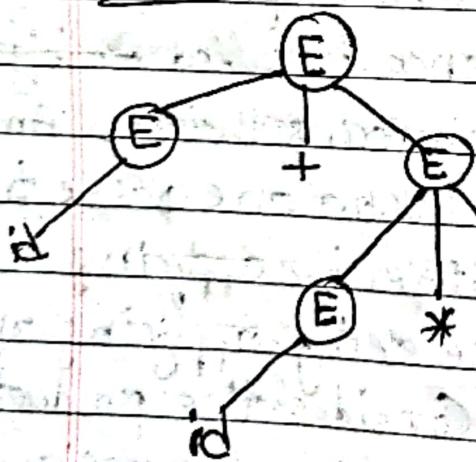
$$E \Rightarrow E+E * E$$

$$\Rightarrow id + E * E$$

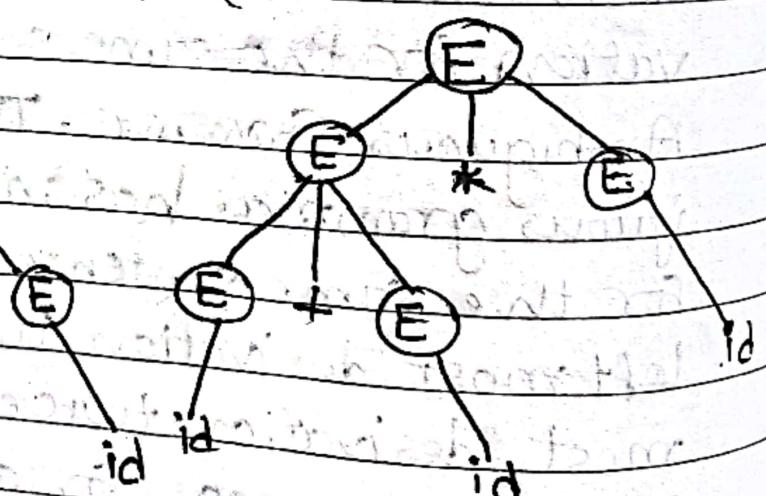
$$\Rightarrow id + id * E$$

$$\Rightarrow id + id * id$$

Tree (i)



Tree (ii)



Since two left most derivation trees are different for the same expression, the given grammar is ambiguous.

\* Parsing and Simple grammars i.e. S-Grammar.  
 Parsing is the method to check whether w can be derived from the grammar G. If we can be derived by applying leftmost derivation from the start symbol, then parsing is successful otherwise not. Parsing can be easily done by special type of Context Free Grammars called simple grammar or S-Grammar which can be formally defined as:

The simple grammar or S-Grammas is a special type of grammar where all production are of the form

$$A \rightarrow aX \quad \text{where, } A \in V, a \in T \text{ & } X \in V^*$$

where, a pair  $(A, a)$  can occur at most once in p i.e., production i.e., if A is there on the left hand side of the production then there can be maximum of production where a is the 1st symbol on the right hand side of the production. Consider the grammar shows below which is S-Grammar. Where as the grammar

$$S \rightarrow aABb$$

$$S \rightarrow aA/a$$

$$B \rightarrow bB/b$$

$$A \rightarrow aA/a$$

is not an S-Grammar because there are two production where the left hand side

of the production is  $S$  and the first symbol on the right hand side of the production is  $a$ .

**Q** Determine the given grammar is ambiguous or not.

$$(i) S \rightarrow aS \mid X$$

$$X \rightarrow aXia$$

Consider the two left most derivations for the string  $aaaa$ .

$$S \Rightarrow aS$$

$$\Rightarrow aaS$$

$$\Rightarrow aaaS$$

$$\Rightarrow aaaX$$

$$\Rightarrow aaaa$$

$$S \Rightarrow X$$

$$\Rightarrow aX$$

$$\Rightarrow aaX$$

$$\Rightarrow aaX$$

$$\Rightarrow aaaa$$

Since, there are two left most derivations for the same sentence  $aaaa$ , the given grammar is ~~not~~ ambiguous.

$$(ii) S \rightarrow AB \mid aAB$$

$$A \rightarrow a \mid Aa$$

$$B \rightarrow b$$

$$(a) \quad S \rightarrow AB$$

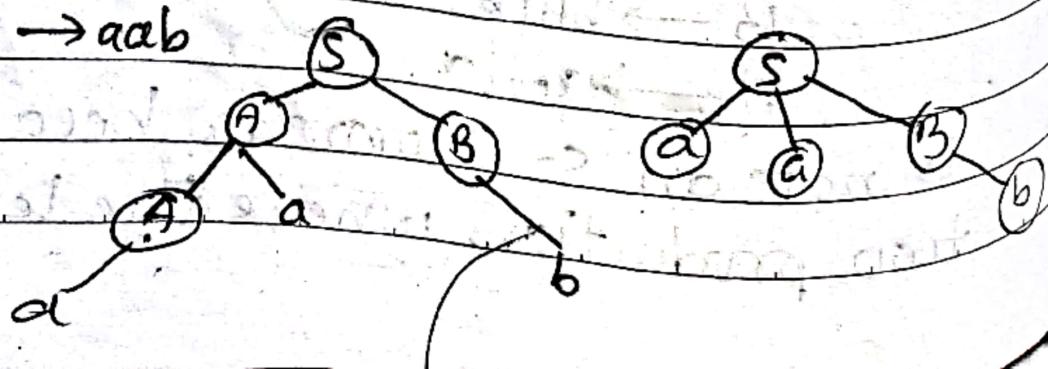
$$\rightarrow Aab$$

$$\rightarrow aAB$$

$$\rightarrow aab$$

$$(b) \quad S \rightarrow aAB$$

$$\rightarrow aab$$



$$\begin{aligned}
 (iii) \quad S &\rightarrow aSbS \\
 S &\rightarrow bSaS \\
 S &\rightarrow \epsilon
 \end{aligned}$$

## A Simplification of Context Free Grammar (CFG)

The CFG imposes no restrictions whatsoever on the right side of the production. However, complete freedom is not necessary. There are several ways in which one can restrict the format of production without reducing the generative power of Context Free Grammar (CFG).

In CFG, it may be necessary to eliminate some of the useless symbol and productions and we have the following method of simplifying the CFG.

- (i) Elimination of useless variable or production
- (ii) Elimination of  $\epsilon$  or  $\lambda$  production.
- (iii) Elimination of Unit production
- (iv) Chomsky Normal Form (CNF)
- (v) Greiback Normal Form (GNF)

### I. Elimination of useless variable

Let  $G = (V, T, P, S)$  be a context free Grammar. A variable  $A \in V$  is said to be useful if and only if there is at least  $w \in L(G)$ . Such that

$$S \xrightarrow{*} x A y \xrightarrow{*} w$$

~~With  $x$  and  $y$  in  $(VUT)^*$~~

Any variable or symbol which is not reachable from the start symbol and which are not used while deriving a string are useless and all the productions which contain these symbol are useless production.

Eg. 1. Eliminate useless symbols and productions from  $G = (V, T, P, S)$  where,

$$V = \{S, B, C, A\}$$

$$T = \{a, b\}$$

$$P = \{S \rightarrow Bb, B \rightarrow bB\}$$

$$B \rightarrow bB \cup \epsilon$$

$$C \rightarrow aq$$

$$A \rightarrow ab$$

Here, C and A are useless variables and the productions

$$C \rightarrow aq$$

$$A \rightarrow ab$$

are useless productions.

Eg. 2. Eliminate useless symbols and productions from  $G = (V, T, P, S)$  where

$$V = \{S, A, B, C\}$$

$$\text{and } T = \{a, b\}$$

with P consisting of

$$A \rightarrow a$$

$$B \rightarrow aa$$

$$C \rightarrow aCb$$

Solns  $\rightarrow$  Step1:- Identify the non-terminal which derives the terminal string.

Here, A and B are variable that are non-terminals which derive the terminal strings. Therefore removing C from production we get production-  $P_1$  as:

$$P_1 = \{$$

$$S \rightarrow aS | A | \epsilon C$$

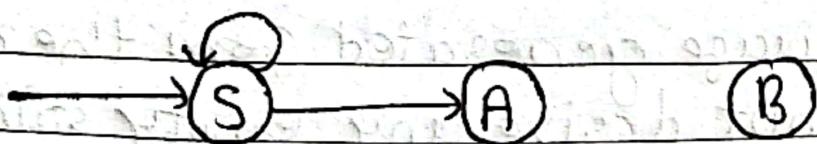
$$A \rightarrow a$$

$$B \rightarrow aa$$

?

Step2: Identify the variables which are reachable from the starting states.

To identify this, draw the dependency graph from the set of productions i.e.  $P_1$ . For the above productions the dependency graph is as follows.



Here, B is independent node and is not reached by the starting state. Eliminate that production. Then,  $P_1'$  becomes:

$$P_1' = \{$$

$$S \rightarrow aS | A | C$$

$$A \rightarrow a$$

Hence, the simplified Grammar is

$$G = V^*, T^*, P^*$$

$$G = (V^*, T^*, P^*, S)$$

where,

$$V^* = \{S, A\}^*$$

$$T^* = \{a\}^*$$

$$P^* = S$$

$$S \rightarrow aS \mid A$$

$$A \rightarrow a$$

?

Here, S is the start symbol.

## (2) Elimination of $\lambda$ or $\epsilon$ Production

If an empty string is not derived from the start symbol, then a production of the form  $A \rightarrow \epsilon$  is not desirable. If the language generated from the grammar G does not derive any empty string and the grammar consists of  $\epsilon$  or  $\lambda$  production then such  $\epsilon$  production can be removed.

An production of a Context Free Grammar of the form  $A \rightarrow \epsilon$  is called  $\epsilon$  or  $\lambda$  production or null production.

Any variable A for which the derivation  $A \rightarrow \epsilon$  is possible is called nullable.

Consider the production for the grammar.

$$S \rightarrow aCb$$

$$C \rightarrow aCb | \epsilon$$

which generates  $\epsilon$ -free language  $a^n b^n | n \geq 1$ ,

Here,  $\epsilon$  production  $C \rightarrow \epsilon$  can be eliminated after adding new production obtained by substituting  $\epsilon$  for  $C$ . Then we get grammar as

$$S \rightarrow aCb | ab$$

$$C \rightarrow aCb | ab$$

Example:- Find a CFG without  $\epsilon$  or  $\lambda$  production equivalent to the grammar defined by

$$S \rightarrow ABCa | bD$$

$$A \rightarrow BC | b$$

$$B \rightarrow b | \epsilon$$

$$C \rightarrow c | \epsilon$$

$$D \rightarrow d$$

Solution:- Step 1:- Find all the nullable variables from the given production, we have

$$B \rightarrow \epsilon$$

$$A \rightarrow BC$$

$$C \rightarrow \epsilon$$

Here, the nullable variables are  $(B, C, A)$

Step 2:- Construct production without  $\epsilon$

Production P	Production $P'$ without $\epsilon$
$S \rightarrow ABCa$	$S \rightarrow ABCa   Bca   Aca   ABa   Ca   Ba   Aa$
$S \rightarrow bD$	$S \rightarrow bD$
$A \rightarrow BC$	$A \rightarrow BC   C   B$
$A \rightarrow b$	$A \rightarrow b$
$B \rightarrow b   \epsilon$	$B \rightarrow b$
$C \rightarrow c   \epsilon$	$C \rightarrow c$
$D \rightarrow d$	$D \rightarrow d$

Hence, the required grammar  $G' = (V', T', P', S)$   
where,

$$V' = \{A, B, C, D\}$$

$$T' = \{a, b, c, d\}$$

$$P' = \{$$

$$S \rightarrow ABCa | Bca | Aca | ABa | Ca | Ba | Aa | a$$

$$S \rightarrow bD$$

$$A \rightarrow BC | C | B$$

$$A \rightarrow b$$

$$B \rightarrow b$$

$$C \rightarrow c$$

$$D \rightarrow d$$

?

where, S is the start symbol.

### 3. Elimination of unit production

Consider the production  $A \rightarrow B$ . The left hand side and right hand side of the production contained only one variable. Such a production is called unit production. The formal definition of unit production is given as: let  $G = (V, T, P, S)$  be a context free grammar. Any production in  $G$  of the form  $A \rightarrow B$  where  $A$  and  $B \in V$  is a unit production. Consider this production

$$A \rightarrow B$$

$$B \rightarrow bB | b$$

Where,  $A \rightarrow B$  is unit production. Therefore, the production  $A \rightarrow B$  can be eliminated and the resulting production can be written as

$$A \rightarrow bB$$

$$B \rightarrow \cancel{bB} | b$$

- Q.1. Eliminate all unit production of the Grammar given below:

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow C | b$$

$$C \rightarrow D$$

$$D \rightarrow E | bC$$

$$E \rightarrow d | Ab$$

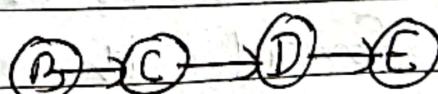


Fig: Dependency graph of unit production

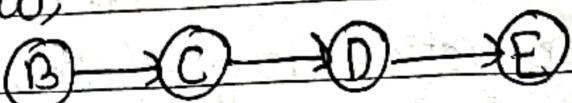
Solution :-

$B \rightarrow C$

$C \rightarrow D$

$D \rightarrow E$

The dependency graph of unit production is shown below,



It is clear from the dependency graph that all null unit production E can be generated from D. The null unit production from E are

$E \rightarrow d | Ab$

Since,  $D \not\Rightarrow^* E$ , we have

$D \rightarrow d | Ab | bC$

From dependency graph. It is clear that  
 $C \not\Rightarrow^* D$

So, null unit productions from D can be generated from C.

$C \rightarrow d | Ab | bC$

Consider the production

$B \not\Rightarrow^* C$

So, null unit production

$B \rightarrow d | Ab | bC | b$

Hence, the final grammar after eliminating unit productions is given as

$$G' = (V', T', P', S)$$

where,  $V' = \{S, A, B, C, D, E\}$

$$T' = \{a, b\}$$

$$P' = \{$$

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow d | Ab | bC | b$$

$$C \rightarrow d | Ab | bC$$

$$D \rightarrow d | Ab | bC$$

$$E \rightarrow d | Ab$$

?  
y

Where,  $S$  is the start symbol.

## Normal Forms

In context free Grammar (CFG) the right hand side of a production can be any string of variables and terminals i.e.  $(VUT)^*$  when production in  $G$  satisfy certain restrictions then  $G$  is said to be in normal form. The two normal forms are:

- (i) Chomsky Normal Form (CNF)
- (ii) Greibach Normal Form (GNF)

## 1. Chomsky Normal Form (CNF)

In CNF, R.H.S. of the production of a Grammar should ~~not~~ contain two or one symbols. If there are two symbols on the R.H.S. of the production, these two symbols must be non-terminal and if there is only one symbol that symbol must be a terminal.

CNF is defined as follows:

Let,  $G = (V, T, P, S)$  be CFG. A grammar  $G$  is said to be CNF if all production are of the form:

$$A \rightarrow BC$$

$$\text{or, } A \rightarrow a \quad \text{or, } S \rightarrow E$$

where,  $A, B$  and  $C \in V$

and  $a \in T$

To convert the grammar into CNF grammar should not have any  $E$  production and unit production. If such production exist, we need to eliminate those production before converting a given grammar into CNF.

Q.1. Convert the grammar given below into CNF.

$$S \rightarrow OA | LB$$

$$A \rightarrow OAA | LS | I$$

$$B \rightarrow LBB | OS | O$$

Sdu<sup>n</sup> Consider the productions which are not in CNF. Replace the terminal on the R.H.S. of the production by a non-terminal and introduce a new production. This step has to be carried out for each production which are not in CNF.

$$\begin{array}{l} A \rightarrow 1 \\ B \rightarrow 0 \end{array} \quad \text{+} \quad \text{(i)}$$

Given Production	Action	Final Production (P)
$S \rightarrow OA   LB$	Replace O with $B_0$ i.e. $B_0 \rightarrow O$ and L with $B_1$ $B_1 \rightarrow L$	$S \rightarrow B_0A   B_1B$ $B_0 \rightarrow O$ $B_1 \rightarrow L$
$A \rightarrow OAA   LS$	Replace O by $B_0$ and I with $B_L$ $B_0 \rightarrow O$ $B_L \rightarrow I$	$A \rightarrow B_0AA   B_L S$ $B_0 \rightarrow O$ $B_L \rightarrow I$
$B \rightarrow LBB   OS$	$B_0 \rightarrow L$ $B_1 \rightarrow O$	$B \rightarrow B_0BB   B_1S$ $B_0 \rightarrow L$ $B_1 \rightarrow O$

Now, we need to restrict the no. of variables on R.H.S. of two production.

The production obtained above which are in CNF are

$$\begin{array}{l} S \rightarrow B_0 A \mid B_1 B \\ B_0 \rightarrow D \\ B_1 \rightarrow 1 \\ A \rightarrow B_0 S \\ B \rightarrow B_0 S \end{array} \quad \text{--- (ii)}$$

The production which are not in CNF are

$$A \rightarrow B_0 AA$$

$$B \rightarrow B_1 BB$$

Considering  $A \rightarrow B_0 AA$

Replace AA on R.H.S. of production with variables  $D_0$  i.e.

$$D_0 \rightarrow AA$$

So, resulting production

$$\begin{array}{l} A \rightarrow B_0 D_0 \\ D_0 \rightarrow AA \end{array} \quad \text{--- (iii)}$$

Similarly for  $B \rightarrow B_1 BB$

$$D_1 \rightarrow BB$$

So, resulting production is

$$\begin{array}{l} B \rightarrow B_1 D_1 \\ D_1 \rightarrow BB \end{array} \quad \text{--- (iv)}$$

Hence, the final grammar which is in CNF is given as  $G' = (V', T', P', S)$   
where,

$$V' = \{ S, B_0, A, B, B_1, D_0, D_1 \}$$

$$T' = \{ 0, 1 \}$$

$P'$  = Combing (ii), (iii) and (iv)

i.e.  $P' = \{$

$$S \rightarrow B_0 A \mid B, B$$

$$A \rightarrow B, S \mid B_0 D_0 \mid 1$$

$$B \rightarrow B, S \mid B, D_1 \mid 0$$

$$B_0 \rightarrow 0$$

$$B_1 \rightarrow 1$$

$$D_0 \rightarrow AA$$

$$D_1 \rightarrow BB$$

?

and  $S$  is the start symbol.

Example 2 Convert the given CFG to CNF:

~~$S \rightarrow S \rightarrow G = \{ S \rightarrow ASA \mid aB \}$~~

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid E$$

Sol'n:

Note:- Steps for converting CFG to CNF  
Step 1: If the start symbol  $S$  occurs on the right hand side of a grammar rule, create a new start symbol  $S'$  and a new production or

grammar rule  $S' \rightarrow S$

Step 2:- Remove all production rules as well as unit production rules from the grammar.

Step 3:- Replace each production rule  $A \rightarrow B_1 B_2 \dots B_n$  where  $n > 2$ , with  $A \rightarrow B_1 C$  where  $C \rightarrow B_2 \dots B_n$ .

Repeat this step for all production rules of the CFG having two or more symbols on the right hand side.

Step 4:- If the R.H.S. of any grammar rule is in the form  $A \rightarrow \alpha B$  where ' $\alpha$ ' is a terminal symbol and  $A, B$  are non-terminals, then the production rule is replaced by  $A \rightarrow \alpha B$  and  $\alpha \rightarrow A$ .

Repeat this step for every production rule of the grammar which is of the form  $A \rightarrow \alpha B$

$$G = \{ S \rightarrow ASA \mid AB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid E$$

?

Step 1:- If we observe the production rules, we can find that in the first two rules starting state  $S$  appear R.H.S. of the productions. We need to add a new starting state  $S'$  and a new production rule  $S' \rightarrow S$ .

Step 2: Removing null production rules  $B \rightarrow \epsilon$  &  $A \rightarrow \epsilon$

After removing  $B \rightarrow \epsilon$ , grammar G.

$$S' \rightarrow S$$

$$S \rightarrow ASA | \alpha B | \alpha$$

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

$$B \rightarrow b$$

After removing  $A \rightarrow \epsilon$ , grammar G.

$$S' \rightarrow S$$

$$S \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$A \rightarrow B | S$$

$$B \rightarrow b$$

Removing unit production rules,  $S' \rightarrow S$ ,  $A \rightarrow B$  &  $A \rightarrow S$

After removing  $S' \rightarrow S$ , grammar G.

$$S \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$S \rightarrow ASA | \alpha B | \alpha | \alpha S | SA$$

$$A \rightarrow B | S$$

$$B \rightarrow b$$

After removing  $A \rightarrow B$ , grammar G

$$S' \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$S \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$A \rightarrow b | S$$

$$B \rightarrow b$$

After removing  $A \rightarrow S$ , grammar G

$$S' \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$S \rightarrow ASA | \alpha B | \alpha | AS | SA$$

$$A \rightarrow b | ASA | \alpha B | \alpha | AS | SA$$

$$B \rightarrow b$$

Step 3: We find some production rules that have more than two symbols on the R.H.S.

$$S' \rightarrow ASA$$

$$S \rightarrow ASA$$

$$A \rightarrow SA$$

Let's replace  $SA$  by  $X$  i.e.  $X \rightarrow SA$

~~$$S' \rightarrow AX | aB$$~~

$$S' \rightarrow AX | aB | a | AS | SA$$

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

$$A \rightarrow b | AX | aB | a | AS | SA$$

$$B \rightarrow b$$

$$X \rightarrow SA$$

Now we will replace terminal symbol  $a$  by  $Y$  from these rules,  $S' \rightarrow aB$ ,  $S \rightarrow aB$  and  $A \rightarrow aB$ .

Now, the grammar  $G$  in CNF

$$S' \rightarrow AX | YB | a | AS | SA$$

$$S \rightarrow AX | YB | a | AS | SA$$

$$A \rightarrow b | AX | YB | a | AS | SA$$

$$B \rightarrow b$$

$$X \rightarrow SA$$

$$Y \rightarrow a$$

Advantages of CNF

- Simple and Standardized grammar
- Efficient Parsing Algorithm

- less Grammar Size
- Elimination of Empty string Rules

### Disadvantages of CNF

- Increased Size
- Loss of Readability
- Complexity of Transformation

H/W :

(Q1). Convert the following CFG into CNF

$$(i) \quad S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \epsilon$$

$$(ii) \quad S \rightarrow aA$$

$$A \rightarrow a$$

$$B \rightarrow c$$

$$(iii) \quad S \rightarrow a \mid aA \mid B$$

$$A \rightarrow aB \mid B \mid \epsilon$$

$$B \rightarrow Aa \mid b$$

⊗

$$(iv) \quad S \rightarrow aXX$$

$$X \rightarrow aS \mid bS \mid a$$

## 2. Greibach Normal Forms (GNF)

GNF stands for Greibach Normal form - A CFG is in GNF if all the production rules satisfy one of the following conditions

$$S \rightarrow E$$

$$A \rightarrow a$$

$$S \rightarrow \alpha d$$

where,  $\alpha \in V^*$

Sample

$$S \rightarrow aAB | ab$$

$$A \rightarrow aA | a$$

$$B \rightarrow bB | b$$

## 2. Greibach Normal Forms (GNF)

GNF stands for Greibach Normal form. A CFG is in GNF if all the production rules satisfy one of the following conditions

$$S \rightarrow E$$

$$A \rightarrow a$$

$$S \rightarrow ad$$

where,  $d \in V^*$

Sample

$$S \rightarrow aAB|ab$$

$$A \rightarrow aA|a$$

$$B \rightarrow bB|b$$

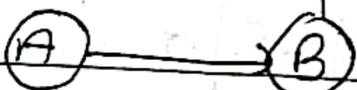
Q1. Convert the given grammar into GNF:

$$S \rightarrow AB|O$$

$$A \rightarrow OO|B$$

$$B \rightarrow 1A1$$

Ans  $\rightarrow$  Since the given grammar does not have any  $E$  production, so we now eliminate unit production  $A \rightarrow B$ . The dependency graph is as follows:



It is clear from dependency graph that  $A \rightarrow B$ . So, all the symbols derivable from  $B$  are also derivable from  $A$ . So in the production

$A \rightarrow 00A|B$ , the variable B can be replaced by the string 1A1 using the production  $B \rightarrow 1A1$ . Then the given grammar becomes

$$S \rightarrow ABA|B$$

$$A \rightarrow 00A|1A1$$

$$B \rightarrow 1A1$$

Next convert the above grammar into CNF, first replace the terminal on the RHS of the production which is not in CNF with variable. Here, replace 1 with  $A_1$  and 0 with  $A_0$  in productions which are not in CNF.

Now the grammar becomes

$$S \rightarrow ABA_1|B$$

$$A \rightarrow A_0 A_0 A_1 | A_0 A A_1$$

$$B \rightarrow A_1 A A_1$$

$$\begin{array}{ll} A_0 \rightarrow 0 & A_1 \rightarrow 1 \\ A_1 \rightarrow 1 & A_0 \rightarrow 0 \end{array}$$

Now, restrict the number of variables.

In RHS of production to two. Then the resulting grammar becomes

$$S \rightarrow AD_1|B$$

$$A \rightarrow A_0 D_2 | A_0 D_3$$

$$B \rightarrow A_1 D_3$$

$$A_1 \rightarrow 1$$

$$A_0 \rightarrow 0$$

$$D_1 \rightarrow BA_1$$

$$D_2 \rightarrow A_0 A$$

$$D_3 \rightarrow AA_1$$

Now, Converting this grammar which is in CNF into GNF. Rename all the variables as shown below

$$\text{let, } S \Rightarrow A_1$$

$$A \Rightarrow A_2$$

$$B \Rightarrow A_3$$

$$A_1 \Rightarrow A_4$$

$$A_0 \Rightarrow A_5$$

$$D_1 \Rightarrow A_6$$

$$D_2 \Rightarrow A_7$$

$$D_3 \Rightarrow A_8$$

The given production becomes

$$A_1 \rightarrow A_2 A_6 | 0$$

$$A_2 \rightarrow A_5 A_7 | A_4 A_8$$

$$A_3 \rightarrow A_4 A_8$$

$$A_4 \rightarrow 1$$

$$A_5 \rightarrow 0$$

$$A_6 \rightarrow A_3 A_4$$

$$A_7 \rightarrow A_5 A_9$$

$$A_8 \rightarrow A_2 A_4$$

Consider production  $A_3 \rightarrow A_4 A_8$

By using the production  $A_4 \rightarrow 1$ , we get

$$A_3 \rightarrow 1 A_8, \text{ which is in GNF.}$$

Consider production:  $A_6 \rightarrow A_3 A_4$

Using  $A_3 \rightarrow 1 A_8$ , we get

*discuss*  
Date \_\_\_\_\_  
Page \_\_\_\_\_

$A_6 \rightarrow 1A_8A_4$ , which is in GNF.

Consider:  $A_2 \rightarrow A_5A_7 | A_4A_5$

Using  $A_5 \rightarrow 0$ ,  $A_4 \rightarrow 1$

$A_2 \rightarrow 0A_7 | 1A_8$ , which is in GNF.

Consider:  $A_1 \rightarrow A_2A_6$

Using  $A_2 \rightarrow 0A_7 | 1A_8$ , we get

$A_1 \rightarrow (0A_7 | 1A_8)A_6$

$A_1 \rightarrow 0A_7A_6 | 1A_8A_6$  which is in GNF

Consider:  $A_7 \rightarrow A_5A_2$

Using  $A_5 \rightarrow 0$ , we get

$A_7 \rightarrow 0A_2$  which is in GNF

Consider:  $A_8 \rightarrow A_2A_4$

Using  $A_2 \rightarrow 0A_7 | 1A_8$ , we get

$A_8 \rightarrow (0A_7 | 1A_8)A_4$

$A_8 \rightarrow 0A_7A_4 | 1A_8A_4$

Finally, the converted grammar in GNF be:

$$G' = (V', T', P', S)$$

$$T' = \{0, 1\}$$

$$P' = S$$

$$A_1 \rightarrow 0A_7A_6 | 1A_8A_6$$

$$A_2 \rightarrow 0A_7 | 1A_8$$

$$A_3 \rightarrow 1A_8$$

$$A_4 \rightarrow 1$$

$$A_5 \rightarrow 0$$

$$A_6 \rightarrow 1 A_8 A_4$$

$$A_7 \rightarrow 0 A_2$$

$$A_8 \rightarrow 0 A_7 A_4 \quad | \quad 1 A_8 A_4$$

γ

Where  $A_1$ , i.e.  $S$  is the start symbol.

Q.2. Convert the following grammar

$$A \rightarrow BC$$

$$B \rightarrow CA|b$$

$$C \rightarrow AB|a$$

into GNF