

3. Context Free Grammar

Formal Languages

Grammars

Classification of Grammar

Chomsky Hierarchy Theorem

Context Free Grammar

Leftmost & Right Most Derivation

Parse Tree

Ambiguous Grammar

Simplification of Context Free Grammar - elimination of Useless Symbol

E-production and Unit Productions

Normal forms of Context Free Grammar - Chomsky Normal Form &

Greibach Normal Form.

Pumping Lemma

Closure properties

Application of Context Free Grammar

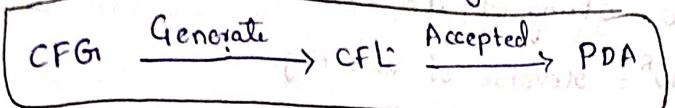
Context free grammar is a formal system of describing language. It consists of four components: Non-terminal symbols, Terminal symbols, Production rules, and Start symbol. The production rules define how non-terminal symbols can be expanded into sequences of terminal symbols. The start symbol indicates the beginning of a sentence. A context-free grammar is called "context-free" because the rules do not depend on the position of symbols in the string being generated. This makes it easier to parse and generate strings. Context-free grammars are used to describe languages like English, Spanish, and French. They are also used in computer science to describe programming languages like C, Java, and Python. Context-free grammars are useful for generating large amounts of data quickly and efficiently. They are also used in natural language processing to understand and generate human language. Context-free grammars are a powerful tool for describing languages, but they have some limitations. For example, they cannot describe all languages, such as Chinese or Japanese, which have complex grammatical structures. They also cannot handle certain types of errors, such as misspellings or punctuation mistakes. Despite these limitations, context-free grammars are a valuable tool for understanding and generating languages.

Rules to follow

Definition of Context Free Grammar:

In Formal Language, CFG generates Context Free Language

Set of CFL are accepted by Push Down Automata



Context Free Grammars is defined by 4-tuples

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

where

N = Set of Non-Terminals

T = Set of Terminals

P = Production rule

S = Start Symbol

The production rule in CFG is in form of

$$\alpha \rightarrow \beta$$

where $\alpha \in NT$ (non terminal)

$$|\alpha| = 1$$

$$\beta \in (NT \cup T)^*$$

Who discovered all those production rules used here?

Backus - Naur Form:

In 1960, John Backus and Peter Naur introduced
Formal notation method for describing Syntax of programming

language, which is known as Backus Naur Form or simply BNF.

Basically, this BNF is designed for ALGOL-60

BNF and context free grammar were nearly identical.

BNF is metalanguage for programming language.

Meta language is a language that is used to describe another language.

BNF symbols to describe a syntax

::= "is defined as"

<> means "Can be defined as"
described

/ means "or"

Rules in BNF:

LHS

<NT> $\longrightarrow \{T\} \{NT\}$

Start symbol

(or)

<NT> ::= {T} {NT}

Problem

Construct a CFG for the language $L = \{www^R | w \in \{a,b\}^*\}$

$$\textcircled{1} \quad L = \{www^R | w \in \{a,b\}^*\}$$

here $w = \text{string of } a's \& b's$

$w^R = \text{reverse of } w \text{ string.}$

$$L = \{c, aba, bcb, abcba, bacab, \dots\}$$

CFG defines 4 tuples $(e, q, T, A) \in \Theta$

$$G_1 = (N, T, P, S)$$

$$N = \{S\}$$

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

Production rule for language: $babbab \rightarrow abab \rightarrow aabb \rightarrow aa$

$$S \rightarrow C$$

$$S \rightarrow aSa$$

$$S \rightarrow bSb$$

$$q \leftarrow s$$

(longest non-Terminal string)

$$L = \{a\}$$

$$^*(TUTU) \in L$$

Derivations and Parse Tree

Derivation

28/12/21 (class) ~~multiple direct left recursion~~ ~~multiple indirect left recursion~~ ~~multiple direct left recursion~~

Left recursion: ~~multiple direct left recursion~~ ~~multiple indirect left recursion~~ ~~multiple direct left recursion~~

CFG

$$A \xrightarrow{\text{CFG}} A\alpha \rightarrow \text{ex of direct left recursion}$$

If LHS of RHS leftmost is same then we get left recursion

left recursion is not supported by grammar. So we can remove this by adding two new production rules.

Eg: $\text{direct } A \xrightarrow{\text{CFG}} A\alpha / B\beta$

$$A \xrightarrow{\text{CFG}} \beta A' / \epsilon$$

equivalent to:

remove redundancy of grammar at end of string

Indirect left recursion:

$$A \xrightarrow{\text{CFG}} B\alpha / \beta$$

$$B \xrightarrow{\text{CFG}} A\beta / \emptyset$$

"to benefit us" \Rightarrow $B \xrightarrow{\text{CFG}} A\beta / \emptyset$

"to reuse" \Rightarrow $A \xrightarrow{\text{CFG}} B\alpha / \beta$

$$A \xrightarrow{\text{CFG}} AB\alpha / \beta / \emptyset$$

RHS

$$A \xrightarrow{\text{CFG}} \gamma A\alpha / \beta$$

$$A' \xrightarrow{\text{CFG}} \beta A' / \emptyset$$

$$A \xrightarrow{\text{CFG}} \delta A'$$

$\{\text{CFG}\} [T] = :: \langle T \rangle$

8.11 factoring • Lookback • Backtracking

$$n \rightarrow \alpha A_1 | \alpha A_2 | \alpha A_3 | \dots | \alpha A_n$$

$$\downarrow$$

$$n \rightarrow \alpha n'$$

$$n' \rightarrow A_1 | A_2 | A_3 | \dots | A_n$$

$$n \rightarrow abd | abc | abt | cfgx | cfgy | cfgz$$

$$\downarrow$$

$$n \rightarrow abn' , n \rightarrow cfgn''$$

$$n' \rightarrow d | left$$

$$n'' \rightarrow x | y | z$$

Simplification of (Context-Free) Grammar

1) Reduce the complexity of program by reducing production rules.

2) Easily generate grammar with less no. of production rules.

how we can simplify:

1) by eliminating useless symbols \rightarrow Non-generating Symbols

2) by eliminating unit production

3) by eliminating ϵ -production.

A non-terminal which doesn't produce a terminal = Non-generating symbol.

If it is available in start symbol, we can't remove it if we have only one start symbol.

$$S \rightarrow AB$$

$$B \rightarrow b$$

$$S \rightarrow AB$$

$$A \rightarrow aB$$

$$B \rightarrow b$$

$$C \rightarrow AB$$

Unit production:

$$\begin{array}{l} A \rightarrow D \\ D \rightarrow C \\ C \rightarrow a \\ B \rightarrow C \\ C \rightarrow b \end{array}$$

$$\begin{array}{l} S \rightarrow AB \\ A \rightarrow aB \\ B \rightarrow b \\ C \rightarrow aC \\ C \rightarrow b \end{array}$$

without unit production

$$\begin{array}{l} S \rightarrow aB \\ A \rightarrow B \\ B \rightarrow T \\ T \rightarrow ab \\ S \rightarrow ab \end{array}$$

+ with unit production

ϵ -production:

$$S \rightarrow Ab$$

$$A \rightarrow BA$$

$$A \rightarrow a$$

$$A \rightarrow \epsilon$$

$$S \rightarrow ab$$

$$S \rightarrow ab | aA$$

$$A \rightarrow b | \epsilon$$

$$S \rightarrow ab | a$$

$$S \rightarrow AB$$

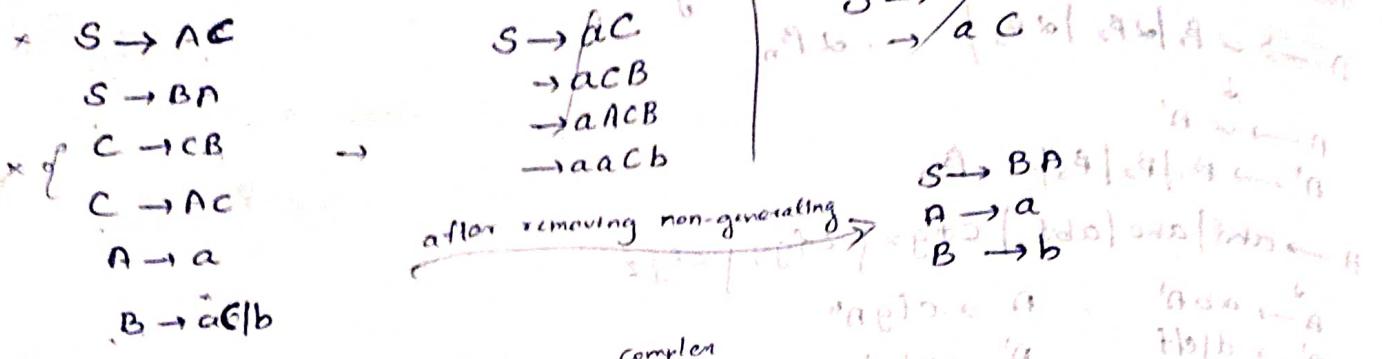
$$A \rightarrow E$$

$$B \rightarrow C$$

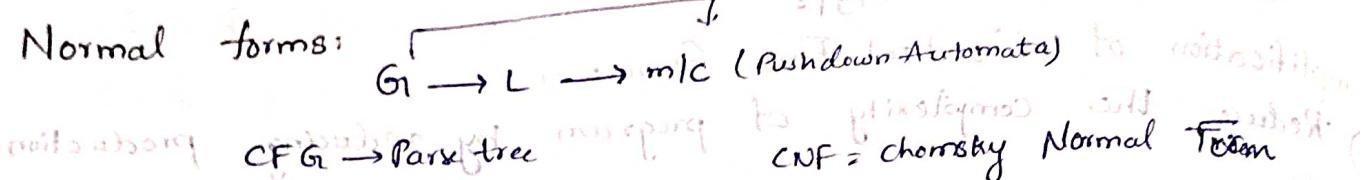
$$C \rightarrow D$$

$$D \rightarrow B$$

$$E \rightarrow a$$



Normal forms:



$CFG \rightarrow$ non-binary Parse tree \rightarrow GNF = Greibach Normal Form

$CFG \rightarrow$ binary Parse tree \rightarrow CNF = Chomsky Normal Form

We get binary Parse tree

$CNF = NT \rightarrow$ Exactly two NT
 $NT \rightarrow$ Single terminals

$S \rightarrow ABa/b$ as follows shows NFA
 $A \rightarrow CC$
 $B \rightarrow BA$
 $A \rightarrow a$
 $B \rightarrow b$
 $C \rightarrow c$

Step 1: Eliminate useless units

Step 2: If

$NT \rightarrow$ two NT

$NT \rightarrow$ single T

stop

else go to next step

Step 3: $NT \rightarrow T_1 T_2 T_3 \dots T_n$ ($\because A \rightarrow abc$)

Step 4: $NT \rightarrow NT_1 NT_2 NT_3 \dots$

Step 5: $NT \rightarrow T_1 T_2 \dots T_n \cdot NT_1 NT_2 \dots NT_n$

$S \rightarrow ABa$
 $A \rightarrow aab$
 $B \rightarrow Ac$

Consider

$C_a \rightarrow a$
 $C_b \rightarrow b$
 $C_c \rightarrow c$

$do \leftarrow 2$

$S \rightarrow ABCa$
 $A \xrightarrow{do \leftarrow 2} C_a \xrightarrow{do \leftarrow 2} C_b$
 $B \xrightarrow{do \leftarrow 2} AC_c$

$BA \leftarrow 2$
 $C_a \rightarrow a$
 $C_b \rightarrow b$
 $C_c \rightarrow c$
 $do \leftarrow 3$
 $S \rightarrow AP$
 $A \rightarrow CaE$
 $B \rightarrow ACC$
 $D \rightarrow BCa$
 $E \rightarrow CaCb$

$A \rightarrow CaCaCb$

$Ao \xrightarrow{do \leftarrow 2} CaEc$

$E \xrightarrow{do \leftarrow 2} CaCb$

$do \leftarrow 2$

$d \leftarrow a$

NT \rightarrow Single terminal & string of TUNT

NT \rightarrow Single terminal

$$A_i \rightarrow A_j x$$

$i < j$ No problem

$i \geq j$ Problem

$$S \rightarrow AA | a$$

$$A \rightarrow SS | b$$

} it is in conflict if $i \geq j$ because a can't be derived from b

$$S \rightarrow A_1$$

$$A \rightarrow A_2$$

$$A_1 \rightarrow A_2 A_2 | a$$

$$A_2 \rightarrow A_1 A_1 | b$$

$$A_2 \rightarrow A_2 A_2 A_1 | a A_1 | b_3 | \text{now } i = j \text{ and } d = gA$$

Now introduce new terminals which is remaining terms other than equal NT.

$$z \rightarrow A_1 A_2 | A_2 A_1 | A_2 A_1 z | eA_1 A_2 gA_1 d \leftarrow gA$$

$$A_2 \rightarrow a A_1 z | b z | a A_1 | b$$

$$A_1 \rightarrow a A_1 z A_2 | b z A_2 | a A_1 A_2 | b A_2 | a$$

$$z \rightarrow a A_1 z A_1 | b z A_1 | a A_1 A_1 | b A_1 | a_1 A_1 z A_1 z | b z A_1 z | a A_1 A_1 z | b A_1 z$$

CFG is not closed under Complement.

Applications:-

used in

CFG \Rightarrow It is a Syntax Analyzer

$$\text{Ex:- } S \rightarrow a | aA | B$$

$$A \rightarrow aBB | \epsilon$$

$$B \rightarrow Aa | b | a$$

After Remove ϵ -production.

$$S \rightarrow a$$

$$S \rightarrow aA x$$

$$S \rightarrow B x$$

$$A \rightarrow aBB x$$

$$B \rightarrow Aa x$$

$$B \rightarrow b$$

$$B \rightarrow a$$

let $\epsilon a = a$

$$S \rightarrow a$$

$$S \rightarrow CaA$$

$$S \rightarrow B x$$

$$A \rightarrow aD$$

$$D \rightarrow BB$$

$$B \rightarrow ACa$$

$$B \rightarrow b$$

$$B \rightarrow a$$

$$S \rightarrow a$$

$$S \rightarrow CaA$$

$$S \rightarrow b | a$$

$$A \rightarrow aD$$

$$D \rightarrow BB$$

$$B \rightarrow ACa$$

$$B \rightarrow a | b$$

$$C \rightarrow a$$

$$S \rightarrow a | aA | b$$

$$A \rightarrow aD$$

$$D \rightarrow BB$$

$$B \rightarrow ACa | aL$$

$$C \rightarrow a$$

GNF: $A_1 \rightarrow A_2 A_3$

$A_2 \rightarrow A_3 A_1 | b$

$A_3 \rightarrow A_1 A_2 | a$

Substitute A_1 value

$A_3 \rightarrow A_2 A_3 A_2 | a$

subst A_2

$A_3 \rightarrow A_3 A_1 A_3 A_2 | b A_3 A_2 | a$

Repeating step 2

$Z \rightarrow A_1 A_3 A_2 | A_1 A_3 A_2 Z$

$A_3 \rightarrow$ Substitute $b A_3 A_2 | a$ in $A_3 A_1 A_3 A_2$

$A_3 \rightarrow b A_3 A_2 A_1 A_3 A_2 | a A_1 A_3 A_2$

$\rightarrow b A_3 A_2 Z | a Z$

$A_3 \rightarrow b A_3 A_2 Z | a Z | b A_3 A_2 | a A_1 A_3 A_2$

It's time to write into

$A_2 \rightarrow b A_3 A_2 Z A_1 | a Z A_1 | b A_3 A_2 A_1 | a A_1 | b$. substitute value

$A_1 \rightarrow b A_3 A_2 Z A_1 A_3 | a Z A_1 A_3 | b A_3 A_2 A_1 A_3 | a A_1 A_3 | b A_3$.

$Z \rightarrow d | A_0 | s d | s | A_0 \leftarrow s A$

$d | c A^2 | c A A^2 | c A s d | c A s | A_0 \leftarrow s A$

$s | A A^2 | s A^2 | s A s | A A^2 | A s | A_0 \leftarrow s$

After all the steps

Step 3, 4, 5, 6

Final result obtain from step 2 & 3

Step 7, 8, 9, 10, 11

Step 12, 13, 14, 15

Step 16, 17, 18, 19

Step 20, 21, 22, 23

Step 24, 25, 26, 27

Step 28, 29, 30, 31

Step 32, 33, 34, 35

Step 37, 38, 39, 40

Step 43, 44, 45, 46

Step 49, 50, 51, 52

Step 55, 56, 57, 58

Step 61, 62, 63, 64

Step 67, 68, 69, 70

Step 73, 74, 75, 76

Step 81, 82, 83, 84

Step 87, 88, 89, 90

$d | A_0 | s \leftarrow e$

$e | B B | s B \leftarrow A$

$a | d | s A \leftarrow B$

so it's E-grammar now

$d | A_0 | s \leftarrow e$

$e | B B | s B \leftarrow A$

$a | d | s A \leftarrow B$

$\cdot d \leftarrow e$

~~Definition of Context Free Grammar:~~

According to the Chomsky's hierarchy, context free grammar or in short CFG is type 2 Grammar.

In Mathematical description, we can describe it as all production are in the form

where $\alpha \in V_N$ (Set of non-terminals)

$|Q|=1$ i.e there will be only one Nonterminal at the left hand side (LHS)

$\beta \in V_N \cup \Sigma$ (Set of non-terminals & terminals)

{A string consists of atleast one non-terminal} \rightarrow {A string of terminals and non-terminals}

Backus Naur Form (BNF)

This is a formal notation to describe the syntax of a given language

This notation was first introduced by John Backus and Peter Naur.

The symbols used in BNF are as follows:

$::=$ denotes "is defined as"

$/$ denotes "Or"

$<>$ used to hold category names

Eg: $<\text{numbers}> ::= <\text{digit}> | <\text{number}> <\text{digit}>$

$<\text{digit}> ::= 0|1|2|3|4|5|6|7|8|9$

This can be described as no is a digit or a num followed by ato digit

A digit is any of the characters from 0to9.

From BNF comes the extended BNF which uses some notations of regular expression such as + or *

As an example, the identifier of a programming language

is defined as

$<\text{identifier}> ::= \text{letter} (\text{letter} | \text{digit})^*$

$<\text{letter}> ::= a|...|z|A|.|_z$

$<\text{digit}> ::= 0|1|2|3|4|5|6|7|8|9$

BNF is widely used in CFG.

- Q) Construct a grammar for the language $L = \{www\mid w \in \{a, b\}^*\}$
- 2) also as string of any combination of a & b .

SOLN: So w^n is also as string of any combination of a and b , but it is a string which is reverse of w . c is a terminal symbol like a, b .

$$L = \{ c, aca, bcb, abcba, bacab, abbccbba, \dots \}$$

we $(a, b)^*$

The null symbols are also accepted in the place of a, b . That means only c is accepted by the grammar to point to $\{\dots\}$ \rightarrow $\{\text{language set}\}$.

\therefore The production rules are (not null except)

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

The grammar becomes $G_1 = \{V_N, T, P, S\}$, where $V_N = \{S\}$

$$T | \Sigma = \{a, b, c\}$$

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

$$S = \{S\}$$

- 2) Construct a CFG for the RE $(0+1)^* 0^*$

Any combination of 0 and 1 followed by a single 0 and ending with any number of 1 .

In RE, a single 0 is between $(0+1)^m$ and 1^n . Consider two NT A & B , which can be replaced multiple times on an NFA.

The production rules of the grammar for constructing this regular expression are:

$$S \rightarrow A S B | 0$$

$$A \rightarrow 0 A \mid 1 A \mid \epsilon$$

The grammar is $G_1 = (V_N, T, P, S)$

$$\text{where, } V_N = \{S, A, B\}$$

$$S = \{S\}$$

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

$$P = S \rightarrow A S B \mid 0, A \rightarrow 0 A \mid 1 A \mid \epsilon$$

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

$$S = \{S\}$$

Construct a CFG for the regular expression $(0111)^* (01)^*$.
 The regular expression consists of two parts $(0111)^*$ followed by $(01)^*$.
 Here from the start symbol, two NT each of them producing one part, are taken (0111) can be written as $0111/1$. As * represents any combination, null string is also included in the language set.

The production rules (P) of the grammar are:

$$S \rightarrow BC$$

$$B \rightarrow AB/\epsilon$$

$$A \rightarrow 0111/1$$

$$C \rightarrow DC/\epsilon 10$$

$$D \rightarrow 01$$

The grammar is

$$\text{where } V_N = \{S, A, B, C, D\}$$

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

$$S = \{S\}$$

4) Construct a CFG for a string in which there are equal number of binary numbers.

$$S \rightarrow 0S1/1S0/\epsilon$$

The grammar is $G_1 = \{V_N, \Sigma, P, S\}$

$$V_N = \{S\}$$

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

$$S = \{S\}$$

Derivation and Parse Tree:

derivation:

In the process of generating a language from the given production rules of grammar, the NT are replaced by corresponding strings of the right hand side (RHS) of the production.

But if there are more than one NT, then which of the ones will be replaced must be determined.

Depending on the selection, the derivation is divided into 2 parts

1) LMD (left most derivation)

2) RMD (Right Most derivation)

i) Left Most Derivation:
A derivation is called Left Most Derivation if we replace only the left most non-terminal by some production rules at each step of the generating process of the language from the grammar.

ii) Right Most Derivation:
A derivation is called Right Most Derivation if we replace only the right most NT by some production rules at each step of the generating process of the language from the grammar.

Eg:-

i) Construct the string 0100110 from the following grammar by using LMD and RMD

$$S \rightarrow OS | 1A | A$$

$$A \rightarrow 0 | 1A | 0B$$

$$B \rightarrow 1 | 0BB$$

ii) Left Most Derivation

$$S \rightarrow OS$$

$$\rightarrow O\cancel{1}AA$$

$$\rightarrow O\cancel{1}BA$$

$$\rightarrow O100\cancel{B}BA$$

$$\rightarrow O1001\cancel{B}A$$

$$\rightarrow O10011\cancel{A}$$

$$\rightarrow O100110$$

$$3/08, /120 \leftarrow 2$$

$$\{2\} = 4V$$

$$\{1, 0\} = 3$$

$$\{0\} = 2$$

NT that are replaced are underlined.

ii) Right Most Derivation

$$S \rightarrow OS$$

$$\rightarrow O1AB$$

$$\rightarrow O1A\cancel{B}$$

$$\rightarrow O1A\cancel{B}O$$

$$\rightarrow O100B\cancel{B}O$$

$$\rightarrow O100B\cancel{B}1O$$

$$\rightarrow O1001\cancel{B}1O$$

$$\rightarrow O10011\cancel{B}O$$

$$\rightarrow O100110$$

Construct the string abbbb from the grammar given below by using i) LMD ii) RMD

i) LMD: $S \rightarrow aAB$ (Non-rightmost in RHS)
 $A \rightarrow bBb$
 $B \rightarrow A|\epsilon$

ii) RMD: $S \rightarrow aAB$ (Non-leftmost in RHS)
 $\rightarrow abBbB$
 $\rightarrow abAbb$
 $\rightarrow abbBbbB$
 $\rightarrow abbbb$

iii) RMD:
 $S \rightarrow aAB$
 $\rightarrow aAA$
 $\rightarrow aAbBb$
 $\rightarrow aAbEb$
 $\rightarrow aAbb$
 $\rightarrow abBbbb$
 $\rightarrow abbbb$

Parse Tree:-
"Parsing a string is finding a derivation for that string from a given grammar."

A parse tree is a tree representation of deriving a CFL from a given context grammar. These type of trees are called as derivation trees.

* A parse tree is an ordered tree in which the LHS of the production represents a parent node and the RHS of a production represents a children node.

* There are certain conditions for constructing a parse tree from a given CFG. Those are as follows:

- Each vertex of the tree must have a label. The label is a non-terminal or terminal or null
- The root of the tree is the start symbol 'S'.
- The label of internal vertices is a NT $\in V_n$
- If there is a production $A \rightarrow x_1 x_2 x_3 \dots x_k$, then for a vertex label A, the children of the node will be $x_1 x_2 \dots x_k$
- A vertex n is called a leaf of parse tree if its label is a terminal symbol $\in \Sigma$ or null (λ)

- * We are describing the similar properties of a tree and a CFG.
 - * For a tree, there must be some root. For every CFG, there is a single start symbol.
 - * Each node of a tree has a single label. For every CFG at LHS there is a single NT.
 - * A child node is derived from a single parent. For constructing a CFL, a NT is replaced by a suitable string at RHS. Each of the character of the string is generating a node. That is, for each single node there is a single parent.
 - * For these similarities, a parse tree can be generated only for a CFG.
- Note:- A parse tree is not possible with context sensitive Grammar because there are productions such as $bB \rightarrow aa$ or $AA \rightarrow B$.
- i) Find the parse tree for generating the string 0100110 from the following grammar.
 - For generating the string 0100110 from the given CFG, the leftmost derivations will be as follows:

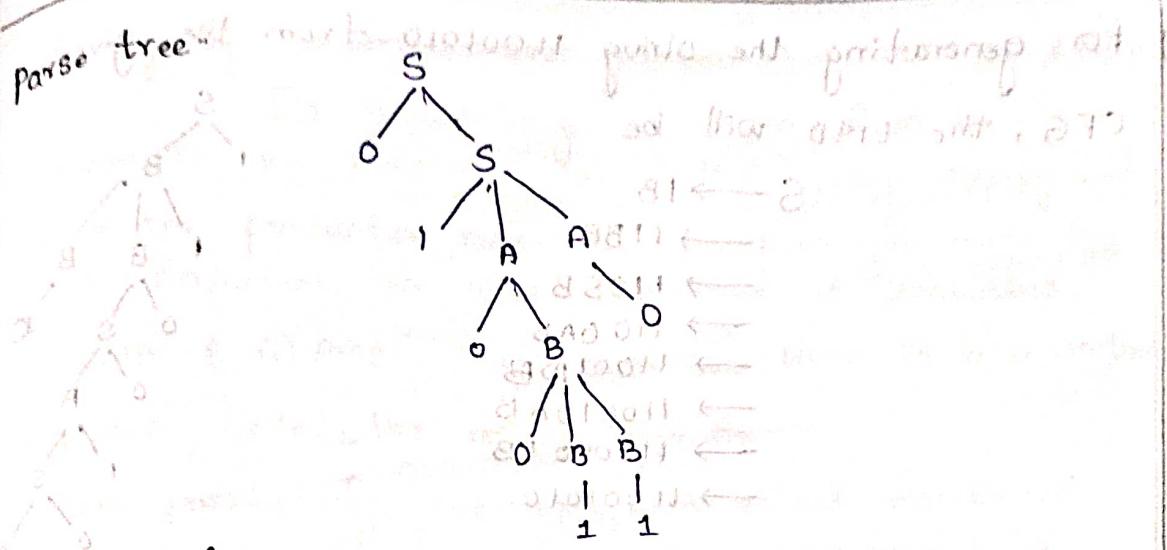
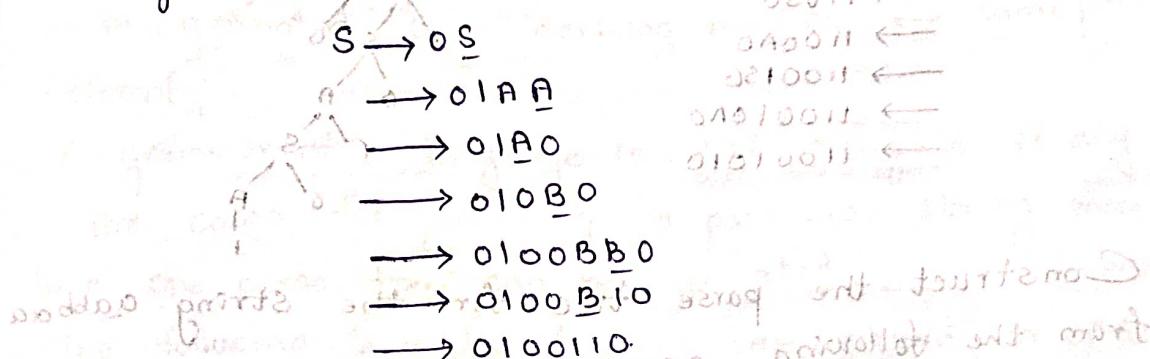
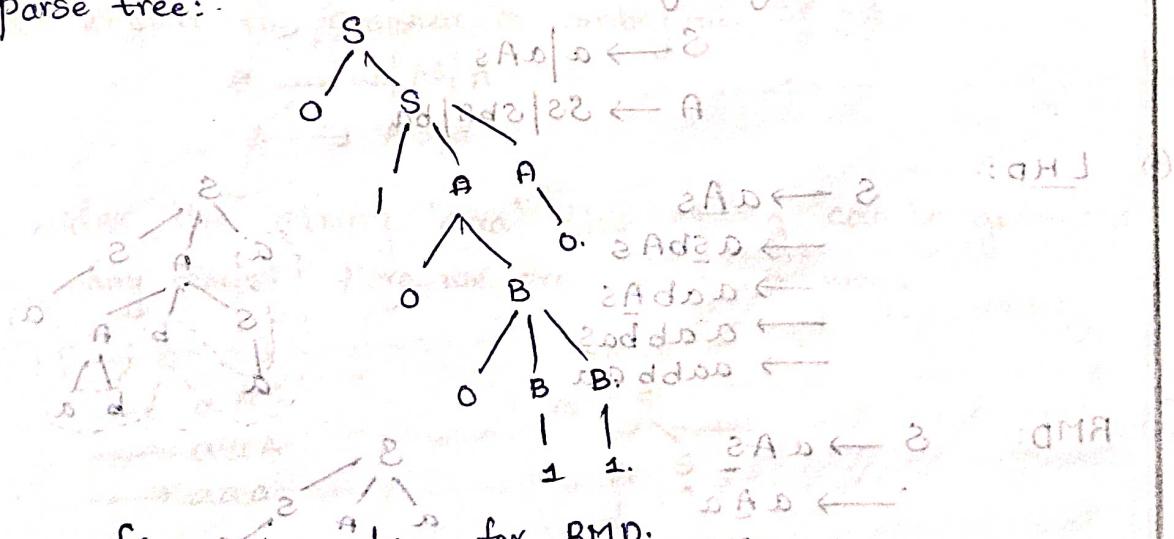


fig: parse tree for LMD

For generating the string 0100110 from the given CFG
the rightmost derivation will be.



parse tree: -



~~fig: parse tree for RMD.~~

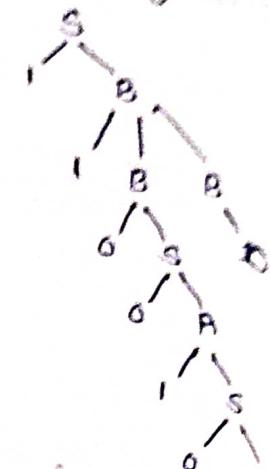
Q) Find the parse tree for generating the string 11001010 from the given grammar.

S → IB | OA

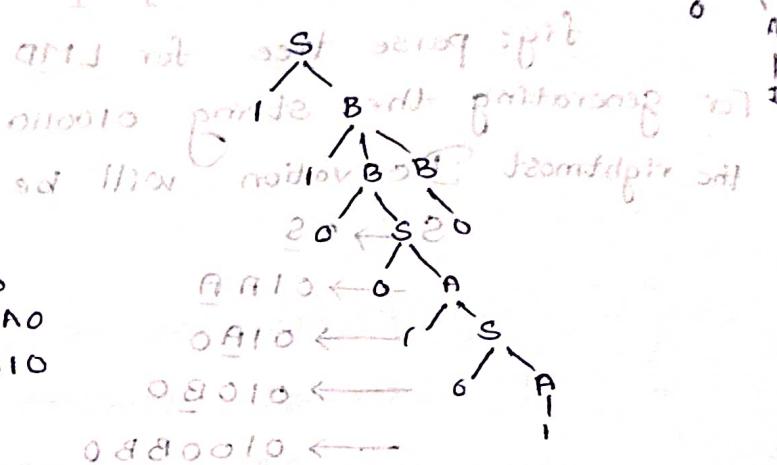
$A \rightarrow 1 | is | oAA$

B → 0|0S|1BB

A) For generating the string 11001010 from the given CFG, the LMD will be

$$\begin{aligned}
 S &\rightarrow 1B \\
 &\rightarrow 11BB \\
 &\rightarrow 110SB \\
 &\rightarrow 1100AB \\
 &\rightarrow 11001SBB \\
 &\rightarrow 110010AB \\
 &\rightarrow 1100101B \\
 &\rightarrow 11001010.
 \end{aligned}$$


RMD:

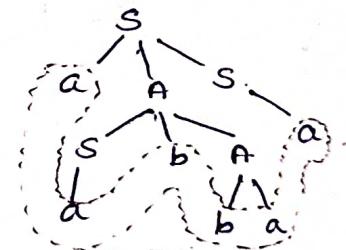
$$\begin{aligned}
 S &\rightarrow 1B \\
 &\rightarrow 11BB \\
 &\rightarrow 11B0 \\
 &\rightarrow 110S0 \\
 &\rightarrow 1100AO \\
 &\rightarrow 11001S0 \\
 &\rightarrow 110010AO \\
 &\rightarrow 11001010
 \end{aligned}$$


B) Construct the parse tree for the string aabbba from the following grammar.

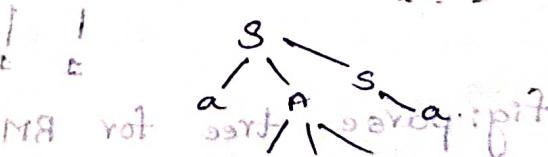
$S \rightarrow a | aAs$

$A \rightarrow SS | sba | ba$

LHD:-

$$\begin{aligned}
 S &\rightarrow aAs \\
 &\rightarrow asbAs \\
 &\rightarrow aabAs \\
 &\rightarrow aabbas \\
 &\rightarrow aabbba
 \end{aligned}$$


RMD:-

$$\begin{aligned}
 S &\rightarrow aAs \\
 &\rightarrow aAa \\
 &\rightarrow asba \\
 &\rightarrow asbbaa \\
 &\rightarrow aabbba
 \end{aligned}$$


$$\begin{aligned}
 &A | a \leftarrow 2 \\
 &AA | aa \leftarrow 1 \\
 &AA | aa | 1 \leftarrow 1 \\
 &AA | aa | 1 | 0 \leftarrow 0
 \end{aligned}$$

Ambiguity in Context free grammar

For generating a string from a given grammar, we have to derive the string step by step from the production rule of the given grammar. For this derivation, we know two types of derivations:

i) LMD & ii) RMD. Except these two there is also another approach called the mixed approach.

More particularly the LMD/RMD is not maintained in each step whereas any of the NT present in deriving string B is replaced by suitable production rule.

By this process, different types of derivation can be generated from deriving may be the same or different.

→ A grammar of a language is called ambiguous if any of the cases for generating a particular string more than one parse tree can be generated.

→ The following examples describe the ambiguity in CFG:

i) Check whether the grammar is ambiguous or not.

$$S \rightarrow aS \mid AS \mid A$$

$$A \rightarrow AS \mid a$$

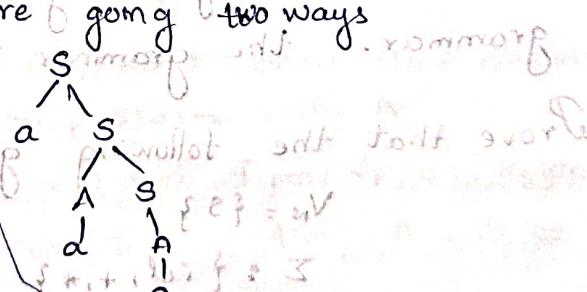
Consider the string "aaa" the string can be generated in many ways. Here we are going two ways.

i) $S \rightarrow aS \mid AS \mid A$

$$\rightarrow aAS$$

$$\rightarrow aaA$$

$$\rightarrow aaa$$



ii) ~~RMD~~: $S \rightarrow AS$

$$\rightarrow A$$

$$\rightarrow AS$$

$$\rightarrow ASS$$

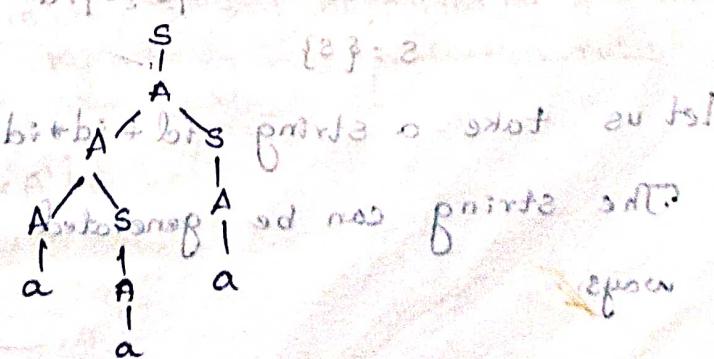
$$\rightarrow aSS$$

$$\rightarrow aAS$$

$$\rightarrow aas$$

$$\rightarrow aaA$$

$$\rightarrow aaa$$



Here for the same string, derived from the same grammar we are getting more than one parse tree. This according to definition, the grammar is an ambiguous grammar.

- 2) Prove that the following grammar is ambiguous.

$$S \rightarrow a \mid absb \mid aAb$$

and also

$$A \rightarrow bs \mid aAb$$

- Take a string abababb. The string can be generated in the following ways.

LMD: $S \rightarrow ab \underline{s} b$ form 1 gives more a boundary b of no

$$\rightarrow abaabb$$

$$\rightarrow ababsbb$$

form 2 gives abababb

and 3 gives abababb

RMD: $S \rightarrow aAb$ form 1 gives more s boundary a of no

$$\rightarrow ab \underline{s} b$$

$$\rightarrow ababsbb$$

$$\rightarrow abababb$$

As we are getting two parse trees of generating a string from the given grammar, the grammar is ambiguous.

- 3) Prove that the following grammar is ambiguous.

$$V_n = \{s\}$$

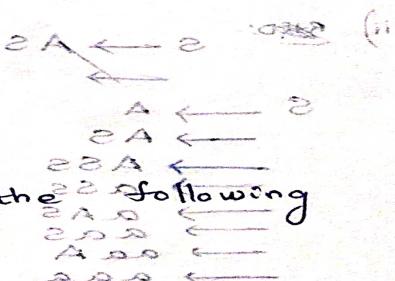
$$\Sigma = \{id, +, *\}$$

$$P: S \rightarrow S + S / S * S / id$$

$$S = \{s\}$$

- 4) Let us take a string id + id * id

The string can be generated in the following ways.



1^{MD}

$$\begin{aligned}
 S &\rightarrow S + S \\
 &\rightarrow id + S \\
 &\rightarrow id + S + S \\
 &\rightarrow id + id + S \\
 &\rightarrow id + id + id
 \end{aligned}$$

A most string matching with tokens will be the one which starts with id.

2^{MD}

$$\begin{aligned}
 S &\rightarrow S + S \\
 &\rightarrow S + id \\
 &\rightarrow id + S + S \\
 &\rightarrow id + S + id \\
 &\rightarrow id + id + S \\
 &\rightarrow id + id + id
 \end{aligned}$$

The first string matching with tokens will be the one which starts with id.

As we are getting two parse trees for generating a string from the given grammar, the grammar is ambiguous.

~~Left Recursion & Left factoring:~~

~~Left Recursion:~~

A context free grammar is called left recursive if a non-terminal 'A' as a leftmost symbol appears alternatively at the times of derivations either immediately or after some other non-terminal.

→ In other words, a grammar is left recursive if it has a NT 'A' such that there is a derivation.

~~A → Aα for something~~

~~Direct left Recursion:~~

Let the grammar rule $A \rightarrow A\alpha | B$, where α consists of terminals and but B doesn't start with A .

At time of derivation, if we start from A and replace A by production $A \rightarrow A\alpha$, then for all time A will be the leftmost NT symbol.

∴ left recursion is not supported by grammar. So we can remove this by adding new production rules

$$\begin{aligned}
 A &\rightarrow BA' \\
 A' &\rightarrow \alpha A' | \epsilon
 \end{aligned}$$

$$A \rightarrow A' \alpha$$

$$A \rightarrow A' \alpha$$

Indirect left Recursion:

Take a grammar in the form $A \rightarrow \alpha A \beta$ or $A \rightarrow \alpha B \beta$ or $B \rightarrow \alpha B \beta$ or $B \rightarrow \alpha A \beta$

where α, β consists of terminals / NT.

At the time of derivation if we start from A, replace A by production $A \rightarrow B\alpha$ & after that replace the B by the production $B \rightarrow A\beta$ then A will appear against as a leftmost NT symbols. This is called the indirect left recursion.

In general, a grammar in the form

$A_0 \rightarrow A_0\alpha_1 | \beta$
 $A_1 \rightarrow A_2\alpha_2 | \beta$
 $A_n \rightarrow A_n\alpha_n | \beta$

is the left recursive grammar

Note: The left recursive grammar is to be converted into a non-left recursive grammar because the top-down parsing technique can't handle left recursion.

Immediate left recursions can be removed by the following process. This is called the Moore's proposal.

For a grammar in the form $A \rightarrow A\alpha | \beta$ where β doesn't start with A

the equivalent grammar after removing left recursion is

$A \rightarrow B_1 A' | \beta$ $A \rightarrow B_n A' | \beta$

for indirect left recursion removal, there is algorithm

Arrange non-terminals in some order: A_1, \dots, A_n

- for i from 1 to n do {

- for j from i to i-1 do {

replace each production

$A_i \rightarrow A_j \beta$ by

$A_i \rightarrow \alpha_1 \alpha_2 \dots \alpha_n \beta$

where

$A_i \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$

- eliminate immediate left recursion among A_1 productions

Consider the following example to make the Moore's proposal. Clear with all the steps with the help of Remove the left recursion from the following grammar

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

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

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

In the grammar there are two immediate left recursion

$$E \rightarrow E + T \quad \text{and} \quad T \rightarrow T * F$$

By using Moore's proposal the left recursion

$$E \rightarrow E + T \quad \text{is removed as}$$

$$E \rightarrow T E' \quad \text{and} \quad E' \rightarrow + T E' \mid \epsilon$$

the left recursion $T \rightarrow T * F$ is removed as

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \epsilon$$

The CFG after removing the left recursion becomes

$$E \rightarrow T E' \quad \text{left recursion is removed}$$

$$E' \rightarrow T E' \mid \epsilon \quad \text{left factor is removed}$$

$$T \rightarrow F T' \quad \text{left recursion is removed}$$

$$T' \rightarrow * F T' \mid \epsilon \quad \text{left recursion is removed}$$

$$F \rightarrow \text{id}(E) \quad \text{left recursion is removed}$$

Left Factoring:

Let us assume that is in a grammar, there is a production rule in the form $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n$. The parses generated from this kind of grammar is not efficient as it requires backtracking. To avoid this problem we need left factor for grammar.

After left factoring, the previous grammar is transformed into

$$A \rightarrow \alpha A_1$$

$$A_1 \rightarrow \beta_1 \mid \beta_2 \mid \beta_3 \dots \mid \beta_n$$

1) Left factor the following grammar statements.

$$A \rightarrow abB | aB | cdg | cdcB | cdfB$$

4) In the previous grammar, the RHS productions abB , aB and ab both start with 'a'. So they can be left factorized. In the same way, cdg , $cdcB$, $cdfB$ all start with 'cd'. So they can be left factorized.

$$A \rightarrow aA'$$

$$A' \rightarrow bB | B$$

$$A \rightarrow cdA''$$

$$A'' \rightarrow g | cb | fb$$

Removing Ambiguity:

→ An ambiguous grammar may be converted into an unambiguous grammar by implementing the precedence and associativity constraints.

These constraints are implemented using the following rules:-

Rule 1:-

Rule 1:-

The precedence of the operators is implemented using the following rules:

→ The level at which the production is present defines the priority of the operator contained in it

→ The higher the level of the production, the lower the priority of operator

→ The lower the level of production, the higher the priority of operator

Rule 2:- Rule 2:- The associativity of operators is implemented by using the following rules:-

→ The operator is left associative, induce left recursion in its production

→ If the operator is right associative, induce right recursion in its production

$$A \rightarrow A \circ A | A \circ$$

Simplification of Context-Free Grammar:-

CFG may contain different types of useless symbols, unit productions and null production. These types of symbols and productions increase the number of steps in generating a language from a CFG. Reduced Grammar contains less number of non-terminals and productions. So the time complexity for the language generated process between less from the reduced grammar. CFG can be simplified in the following process:

1. Removal of useless symbols

- Removal of non-generating Symbols.
- Removal of non-reachable Symbols.

2. Removing unit production Symbols

3. Removing Null production

Removal of Useless Symbols:

Useless Symbols are of two types:-

- Non generating Symbol: They are those symbols in which do not produce any terminal string.
- Non-reachable Symbol: Those symbols which cannot be reached at any time starting from the Start Symbol.

These two types of useless symbols can be removed according to the following process:-

→ Find the non-generating symbols i.e. the symbols which do not generate any terminal string. If the start symbol is found non-generating leave that symbol. Because the start symbol cannot be removed, the language generating process starts from symbol

→ For removing non-generating symbol, remove those production whose right side and left side contain those symbols.

→ Now find the non-reachable symbols i.e. the symbols which cannot be reached, starting from the start symbol

→ Remove non-reachable symbols such as rule(2) etc.

$$S \leftarrow D, D \leftarrow C, C \leftarrow B, B \leftarrow A, A \leftarrow \emptyset$$

i) Remove the useless symbols from the given CFG.

$$S \rightarrow A$$

leads to $S \rightarrow BA$

$$C \rightarrow CB$$

$$C \rightarrow NC$$

$$A \rightarrow a$$

$$B \rightarrow a/b$$

→ There are two types of useless symbols: non-generating & non-reachable symbols. First, we are finding non-generating symbols.

→ Those symbols which do not produce any terminal string are non-generating symbol as it does not produce any terminal string.

→ So, we have to remove the symbol C. To remove C, all the productions containing C as a symbol must be removed.

By removing the productions, the minimized grammar will be

$$S \rightarrow BA$$

$$A \rightarrow a/b$$

$$B \rightarrow b$$

→ Now we have to find non-reachable symbols, the symbols which cannot be reached at any time starting from start symbol.

→ There is no non-reachable symbol in the grammar. So the minimized form of the grammar by removing useless symbols is

$$S \rightarrow BA/A$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Removal of Unit Productions, Production in the form Non-terminal \rightarrow Single non-terminal is called unit production.

Unit production increases the number of steps as well as the complexity at the time of generating M.

language from the grammar.

This will be clear if we take an example.

Let there be a grammar

$$S \rightarrow AB, A \rightarrow E, B \rightarrow E, C \rightarrow D, D \rightarrow b, E \rightarrow a$$

From here, if we are going to generate a language, then it will be generated by the following way:

$$S \rightarrow AB \rightarrow EB \rightarrow aB \rightarrow aC \rightarrow aD \rightarrow ab.$$

The grammar by removing unit production and as well as minimizing will be:

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

From this the language will be generated by the following way:

$$S \rightarrow AB \rightarrow aB \rightarrow ab$$

i) Remove the unit productions from the following grammar.

$$\begin{aligned} S &\rightarrow AB \\ A &\rightarrow E \\ B &\rightarrow C \\ C &\rightarrow D \\ D &\rightarrow b \\ E &\rightarrow a \end{aligned}$$

In the given grammar, the unit productions are:

$$A \rightarrow E$$

$$B \rightarrow C$$

For removing the unit production $A \rightarrow E$, we have to find a non-unit production in the form $A \rightarrow E$.

Ex: Some strings of terminal or NT or both.

There is a production $E \rightarrow a$.

So, $A \rightarrow a$ will be added to the production rule.

Similarly $B \rightarrow C$ will be added to the production rule.

Now, the modified production rules are:

$$S \rightarrow AB$$

$$S \rightarrow A$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$C \rightarrow b$$

$$D \rightarrow b$$

This is the grammar without unit production but it is not minimized.

There are useless symbols (non-reachable symbol) C, D, E .

By removing the useless symbols, the grammar will be

$$S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

Removal of Null Productions:

A production in the form $NT \rightarrow \epsilon$ is called

Null production.

If ϵ (null string) is the language set generated from the grammar, then the null production cannot be

removed. That is, if we get $S \rightarrow \epsilon$, then that null production

cannot be removed from the production rules.

Procedure to remove Null production:-

A procedure in the form of $NT \rightarrow \epsilon$ is called

Null production.

If ϵ is the string language set generated from the grammar, then the null production can be removed.

That is if we get $S \rightarrow \epsilon$, then the null production can't be removed from the production rule.

Procedure to Remove Null production:-

Step1: If $A \rightarrow \epsilon$ is a production to be eliminated, then we look for all productions whose right side contains

Step2: Replace each occurrence of A in each of those productions to obtain the non- ϵ production.

Step3: These non-null productions must be added to the grammar to keep the grammar language generating power the same.

1) Remove the ϵ production from the following grammar

$$S \rightarrow aA$$

$$A \rightarrow b|\epsilon$$

Ed: In the previous production rules, there is a null production $A \rightarrow \epsilon$. the grammar doesn't produce null string as a language & so this null production can be removed.

According to the first step we have to look for all productions, whose right side contains there is

one step $A \rightarrow \epsilon$.

Remove the ϵ -production

so we get the same set of productions.

$$S \rightarrow ABaC$$

as well, $A \rightarrow BC$ will produce which is same as the

productions of $B \rightarrow b|c$ and $C \rightarrow D|c$ in A to consider now for

$$C \rightarrow D|c$$

$$D \rightarrow d$$

The productions will be: $S \rightarrow ABaC$

$$A \rightarrow BC$$

$$B \rightarrow b|c$$

$$C \mid A \rightarrow E$$

$$C \rightarrow D|c$$

$$D \leftarrow E$$

$$D \rightarrow d$$

replacing ϵ in place of C & B .

$$S \rightarrow ABa$$

$$S \rightarrow Aa$$

$$S \rightarrow AAC$$

$$A \rightarrow B$$

$$A \rightarrow C$$

the grammar will be

$$S \rightarrow ABaC \mid ABa \mid AAC \mid Aa$$

$$A \rightarrow BC \mid B \mid c$$

$$C \rightarrow D$$

$$D \rightarrow d$$

$$B \rightarrow b$$

Remove the ϵ -production

$$S \rightarrow ABAC$$

$$A \rightarrow aA \mid \epsilon$$

$$B \rightarrow bB \mid \epsilon$$

$$C \rightarrow c \mid \epsilon$$

a) the given productions are:-

$$S \rightarrow ABAC$$

$$A \rightarrow aA, A \rightarrow \epsilon$$

$$B \rightarrow bB, B \rightarrow \epsilon$$

$$C \rightarrow c \mid d \mid xd \leftarrow \epsilon$$

replacing ' ϵ '

$$S \rightarrow aABAC$$

$$S \rightarrow BAC$$

$$S \rightarrow AbBAC$$

$$S \rightarrow AAC$$

$$S \rightarrow ABAC$$

$$S \rightarrow BAC$$

$$S \rightarrow ABC$$

$$S \rightarrow AAC$$

$$S \rightarrow BC$$

$$S \rightarrow C$$

$$S \rightarrow AC$$

Linear Grammar:-

The grammar will be

$$S \rightarrow BAC \mid A \bar{B} C \mid A \bar{A} C \mid B \bar{C} \mid C \mid A \bar{C} \mid A B A C$$

$$A \rightarrow aA$$

$$B \rightarrow bB$$

$$C \rightarrow c$$