

Unit-4 : Context Free Grammar

- Ankit Pangani

The term grammar applied to a language refers to the mechanism for constructing phrases & sentences that belong to a language. It consists a set of rules, by which strings in a language can be generated.

Context free grammars (CFGs)

CFGs are used to describe context-free languages. It is a set of recursive rules used to generate patterns of strings.

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

{
 V: finite set of variables (Non Terminal). (In Capital letters)
 T: finite set of Terminals ($N_T \neq \emptyset$) (small letters)
 P: Production Rules (substitution Rules) {
 S: Start variable.
 where.
 $\alpha \rightarrow \beta$
 $\alpha \in V$ which is
 only one variation

Herr,

Ques.
P is the set of production rules which is used in replacing non-terminal symbol 'X' on the left side of production in a string with other terminal or non-terminal symbols ie. 'P' on the right side of production

S is the start symbol used to derive the string.

Example:

Rule 1: $S \rightarrow \epsilon$

Rule 2: $S \rightarrow OS\perp$

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, 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 where production starts

Formally, grammar is defined as $\{ V = \{S\}, T = \{\epsilon, 0, 1\}, P = \{ S \rightarrow \epsilon, S \rightarrow OS\perp \}, S = \{S\} \}$

$$P = \{ S \rightarrow \epsilon | OS\perp \}$$

Q1: CFG for language $0^n 1^n$ where $n > 0$

$$\begin{array}{l} S \rightarrow \epsilon \\ S \rightarrow OS\perp \end{array} \quad \left(\begin{array}{l} S \xrightarrow{0^n} OS\perp \\ \vdots \\ S \xrightarrow{n1} 1^n \end{array} \right)$$

Q2: Construct CFG for the language having any no. of 0's over the set $\Sigma = \{a\}$

$$\Sigma = \{a\}$$

$$L = \{ \epsilon, a, aa, aaa, \dots \}$$

$$R.E = a^*$$

So, production rules are

$$S \xrightarrow{*} \epsilon$$

$$S \xrightarrow{*} S a$$

Rule 1: $S \rightarrow \epsilon$ Rule 2: $S \rightarrow aS$

$$\therefore V = \{S\}$$

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

$$S = \{S\}$$

Let's derive a string "aaaaaa"

$$S \rightarrow aS$$

$$\Rightarrow aas \quad (S \rightarrow aS)$$

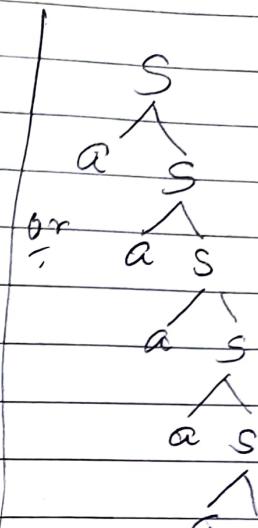
$$\Rightarrow aaas \quad (S \rightarrow aS)$$

$$\Rightarrow aaas \quad (S \rightarrow aS)$$

$$\Rightarrow aaaaas \quad (S \rightarrow aS)$$

$$\Rightarrow aaaaaas \quad (S \rightarrow aS)$$

$$\Rightarrow aaaaaaa \quad (S \rightarrow \epsilon)$$



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

$\Rightarrow L = \{abb, aabbabb, aaabbbaabb, \dots\}$

The production rules are:

Rule 1: $S \rightarrow aSbb$

rule 2: $S \rightarrow abb$

$$\therefore G = (\{S\}, \{a, b\}, P, S)$$

Ex: Any no. of a's and b's

$\Rightarrow L = \{\epsilon, a, b, aa, ab, \dots\}$

$$RE = (a+b)^*$$

The production rules are:

$$S \rightarrow aS$$

$$S \rightarrow bS$$

or

$$S \rightarrow aS | bS | \epsilon$$

$$S \rightarrow \epsilon$$

$$\therefore G = (\{S\}, \{a, b\}, P, S)$$

Ex: At least 2 a's

$\Rightarrow L = \{aa, abaaab, aaaabb, \dots\}$

$$RE = (a+b)^* a (a+b)^* a (a+b)^*$$

The production rules are:

$$S \rightarrow AaAaA$$

$$A \rightarrow aA | bA | \epsilon$$

Ex: One occurrence of 000 over $S = \{0, 1\}$

$$\Rightarrow L = \{000, 0001, 00001, 111000, \dots\}$$

$$R.F = (0+1)^* 000 (0+1)^*$$

The production rules are.

$$S \rightarrow A \quad 000 \quad A$$

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

$$Ex: L = \{a^n b^n c^m d^m \mid n, m \geq 1\}$$

$$\Rightarrow S \rightarrow aSb \mid a$$

$$S \rightarrow AB$$

$$A \rightarrow aAb \mid ab$$

$$B \rightarrow cBd \mid cd$$

$$Ex: L = \{a^n b^n \mid n \geq 1\}$$

$$\Rightarrow L = \{abb, aabb, aaabbb, \dots\}$$

$$\therefore S \rightarrow aSb$$

$$S \rightarrow ab$$

$$\therefore L = \{a^n b^{n+2}, n \geq 0\}$$

$$\Rightarrow S \rightarrow aSb \mid bb$$

$$Ex: L = \{a^n b^m c^n \mid n, m \geq 1\}$$

$$\therefore L = \{a^{2n+3} b^n, n \geq 0\}$$

\Rightarrow

$$S \rightarrow aSc \mid aAc$$

$$A \rightarrow bA \mid b$$

$$\Rightarrow S \rightarrow aaSb \mid aaa$$

$$Ex: L = \{a^n b^m c^n d^n \mid n \geq 1, m \geq 1\}$$

$$\Rightarrow S \rightarrow AS,$$

$$S \rightarrow aSb \mid \epsilon$$

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

$$\therefore L = \{a^{m+n} b^m c^n \mid n, m \geq 1\}$$

$$\Rightarrow L = \{a^m a^n b^m c^n \mid n, m \geq 1\}$$

$$S \rightarrow aSc \mid aAc$$

$$A \rightarrow aAb \mid ab$$

$$Ex: L = \{a^m b^m c^n \mid m, n \geq 0\}$$

Ex: CFG for strings having length ^{at least} two

$$\Rightarrow RE = (0+1)(0+1)(0+1)^*$$

$$S \rightarrow AAB$$

$$A \rightarrow 0 \mid \perp$$

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

Uses of CFG.

- Context-free grammars are used to describe context-free languages.
- CFGs are studied in the fields of theoretical computer science, compiler design, and linguistics.
- CFGs are used to describe programming languages and parsers programs in compilers can be generated automatically from context-free grammars.
- Document Type Definition in XML is a context-free grammar which describes the HTML tags & the rules to use the tags in a nested fashion.

Context Free Language (CFL) (Language of Grammar)

A language L is a context free language if there exists a context free grammar G such that $L = L(G)$. where $L(G) = \{x \in T^* \mid S \Rightarrow^* x\}$

Note: $G = (V, T, P \text{ and } S)$

The set of all context free language is identical to the set of languages accepted by pushdown automata, & the set of regular languages is a subset of context-free language.

CFL & CFG have applications in Computer science and linguistics such as natural language processing and computer language design.

* Derivation Using Grammar rule.

The process of producing strings using the production rules of grammar is called derivation.

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

- o For $0 \leq i \leq n$, the string $\beta_i \in (VUT)^*$
- o $\beta_0 = S$
- o For $0 \leq i \leq n$, there is a production of P that applied to β_i yields β_{i+1} .
- o $\beta_n \in T^*$

There are two possible approaches of derivation.

- i) Body to head (Bottom up) approach.
- ii) Head to body (Top down) approach. [leftmost rightmost]

* Body to head (Bottom up) approach.

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.

Ex: Consider grammar. $S \rightarrow S + S \mid S/S \mid (S) \mid S-S \mid S^* \mid \alpha$

Here to generate $\alpha + (\alpha * \alpha) \mid \alpha - \alpha$

S.N.	String Inferred	Variable	Production	String used
1.	α	S	$S \rightarrow \alpha$	α ; string ₁
2.	$\alpha * \alpha$	S	$S \rightarrow S * S$	$\alpha * \alpha$; string ₂
3.	$(\alpha * \alpha)$	S	$S \rightarrow (S)$	string ₁ & string ₂ ; string ₃
4.	$(\alpha * \alpha) / \alpha$	S	$S \rightarrow S / S$	string ₁ & string ₃ ; string ₄
5.	$\alpha + (\alpha * \alpha) / \alpha$	S	$S \rightarrow S + S$	string ₁ & string ₄ ; string ₅
6.	$\alpha + (\alpha * \alpha) / \alpha - \alpha$	S	$S \rightarrow S - S$	string ₁ & string ₅ ; string ₆

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 (Top down) approach

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:

1. Leftmost derivation.
2. Rightmost derivation.

1. Leftmost derivation.

A derivation in CFG is a leftmost derivation (LMD) if, at each step, a production is applied to the leftmost variable occurred in the current string. It is denoted by \xrightarrow{lm}

Ex: Derive the string "aababbba" for leftmost derivation using a CFG.

$$\begin{aligned} S &\rightarrow aB \mid bA \\ A &\rightarrow a \mid as \mid bAA \\ B &\rightarrow b \mid bs \mid aBB \end{aligned}$$

→ Leftmost derivation of the string "aababbba" is shown:-
Here, 'S' is start variable, so:-

$$\begin{aligned} S \xrightarrow{lm} aB & \quad (\text{Rule: } S \rightarrow aB) \\ S \xrightarrow{lm} aABB & \quad (\text{Rule: } B \rightarrow aBB) \\ S \xrightarrow{lm} aabB & \quad (\text{Rule: } B \rightarrow b) \\ S \xrightarrow{lm} aab\ bs & \quad (\text{Rule: } B \rightarrow bs) \\ S \xrightarrow{lm} aab\ bAB & \quad (\text{Rule: } S \rightarrow aB) \\ S \xrightarrow{lm} aabb\ bs & \quad (\text{Rule: } B \rightarrow bs) \\ S \xrightarrow{lm} aabb\ bBA & \quad (\text{Rule: } S \rightarrow bA) \\ S \xrightarrow{lm} aabbabbba & \quad (\text{Rule: } A \rightarrow a) \end{aligned}$$

Ex: Consider a grammar defined as

$$S \rightarrow S^*S \quad | \quad S^*S \quad | \quad S-S \quad | \quad S/S \quad | \quad (S) \quad | \quad a$$

Derive the string $a + (a * a)$ using leftmost derivation

→ Leftmost derivation of string $a + (a * a)$ is shown as

Here, S' is start variable so,

$S \xrightarrow{r_m} S + S$	(Rule: $S \rightarrow S + s$)
$S \xrightarrow{r_m} a + S$	(Rule: $S \rightarrow a$)
$S \xrightarrow{r_m} a + (S)$	(Rule: $S \rightarrow (S)$)
$S \xrightarrow{r_m} a + (S * S)$	(Rule: $S \rightarrow S * S$)
$S \xrightarrow{r_m} a + (a * S)$	(Rule: $S \rightarrow a$)
$S \xrightarrow{r_m} a + (a * a)$	(Rule: $S \rightarrow a$)

2. Rightmost derivation

A derivation in context-free grammar is a RMD if, at each step, a production is applied to the rightmost variable occurred in the current string. It is denoted by \xrightarrow{rm}

Ex: Derive the string "aabbabba" for ^{right} leftmost derivation using a CFG

$$\begin{aligned} S &\rightarrow aB \mid bA \\ A &\rightarrow a \mid as \mid bAA \\ B &\rightarrow b \mid bs \mid aBB \end{aligned}$$

$S \xrightarrow{rm} aB$	(Rule: $S \rightarrow aB$)
$S \xrightarrow{rm} aaBB$	(Rule: $B \rightarrow aBB$)
$S \xrightarrow{rm} aaBs$	(Rule: $B \rightarrow bs$)
$S \xrightarrow{rm} aaBbba$	(Rule: $B \rightarrow ba$)
$S \xrightarrow{rm} aaBBba$	(Rule: $A \rightarrow a$)
$S \xrightarrow{rm} aabsbba$	(Rule: $B \rightarrow bs$)
$S \xrightarrow{rm} aabbAbba$	(Rule: $S \rightarrow ba$)
$S \xrightarrow{rm} aabbabba$	(Rule: $A \rightarrow a$)

Ex: Consider a grammar defined as

$$S \rightarrow S * S \mid S + S \mid S - S \mid S / S \mid (S) \mid a$$

Now the rightmost derivation of the string $aabaaa$ is shown as:

$$\begin{aligned} S &\xrightarrow{r_m} S + S & (\text{Rule: } S \rightarrow S + S) \\ S &\xrightarrow{r_m} S + (S) & (\text{Rule: } S \rightarrow (S)) \\ S &\xrightarrow{r_m} S + (S * S) & (\text{Rule: } S \rightarrow S * S) \\ S &\xrightarrow{r_m} S + (S * a) & (\text{Rule: } S \rightarrow a) \\ S &\xrightarrow{r_m} S + (a * a) & (\text{Rule: } S \rightarrow a) \\ S &\xrightarrow{r_m} a + (a * a) & (\text{Rule: } S \rightarrow a) \end{aligned}$$

Ex: Consider a grammar : $S \rightarrow aAs \mid a$
 $A \rightarrow Sba \mid ss \mid ba$

Given a string $aabaaa$, give left most & right most derivation

Leftmost derivation

$$\begin{aligned} S &\xrightarrow{l_m} aAs & ("") \\ S &\xrightarrow{l_m} ass & ("") \\ S &\xrightarrow{l_m} aass & ("") \\ S &\xrightarrow{l_m} aaAss & ("") \\ S &\xrightarrow{l_m} aabass & ("") \\ S &\xrightarrow{l_m} aaabaas & ("") \\ S &\xrightarrow{l_m} aaabaaa & ("") \end{aligned}$$

Rightmost derivation

$$\begin{aligned} S &\xrightarrow{r_m} aAs & ("") \\ S &\xrightarrow{r_m} aAa & ("") \\ S &\xrightarrow{r_m} assa & ("") \\ S &\xrightarrow{r_m} asaAsa & ("") \\ S &\xrightarrow{r_m} asataa & ("") \\ S &\xrightarrow{r_m} asabaaa & ("") \\ S &\xrightarrow{r_m} aabaaa & ("") \end{aligned}$$

X Sentential forms

- Derivations from the start symbol are called sentential forms.
- That is, given a CFG $G = (N, T, P, S)$, if $S \xrightarrow{*} \alpha$ where $\alpha \in (VUT)^*$, then α is a sentential form.
- If $S \xrightarrow{*} \alpha$ where $\alpha \in (VUT)^*$, then α is a left-sentential form.
- If $S \xrightarrow{*} \alpha$ where $\alpha \in (VUT)^*$, then α is a right-sentential form.

Note:- we already used sentential form in previous topic of leftmost derivation & rightmost derivation.

There are basically two methods:
Sentential method, another is tree method.

An example of tree method is:

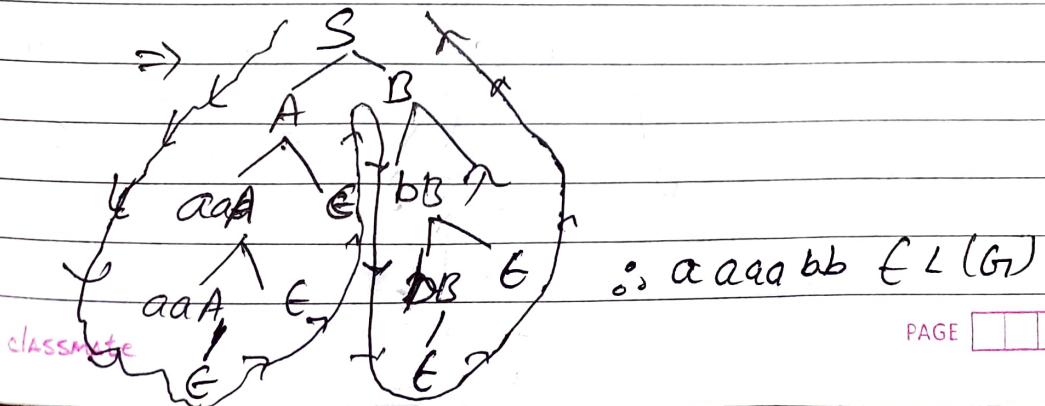
$$S \rightarrow AB$$

Check whether

$$A \rightarrow aaA / \epsilon$$

$$w = aaaabb \in L(G).$$

$$B \rightarrow bB / \epsilon$$



* Parse Tree / Derivation Tree:

A parse tree or parsing tree or derivation tree is an ordered, rooted tree that represents the syntactic structure of an input string according to some context-free grammar.

It is a hierarchical representation of terminals or non terminals. The leaves of parse tree represent terminals and each interior node represents production of grammar. The in-order traversal gives original input string.

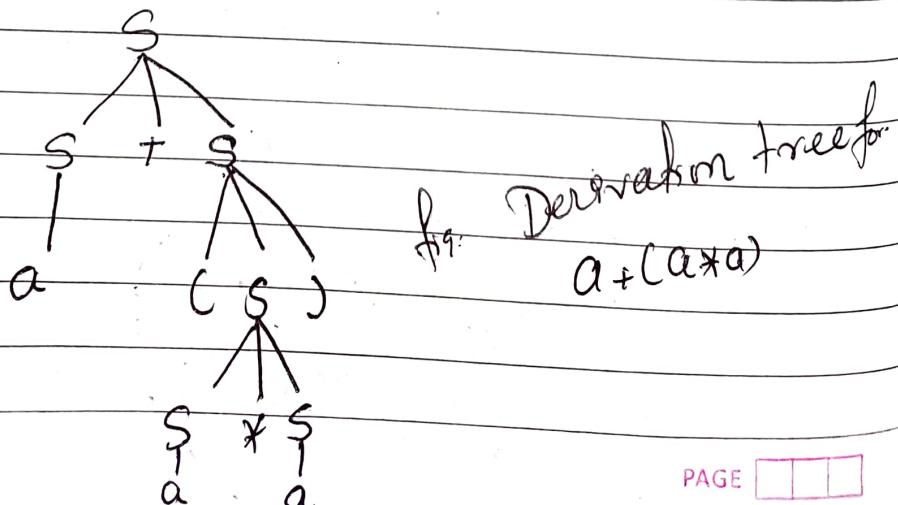
Note: The root is labeled by the start symbol.

Formally, given a context free grammar (CFG) $G = \{V, T, P, S\}$, a parse tree is a n-ary tree having the following properties:

Ex: Consider the grammar:

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

The parse tree for the string $a + (a * a)$ is shown as follows:



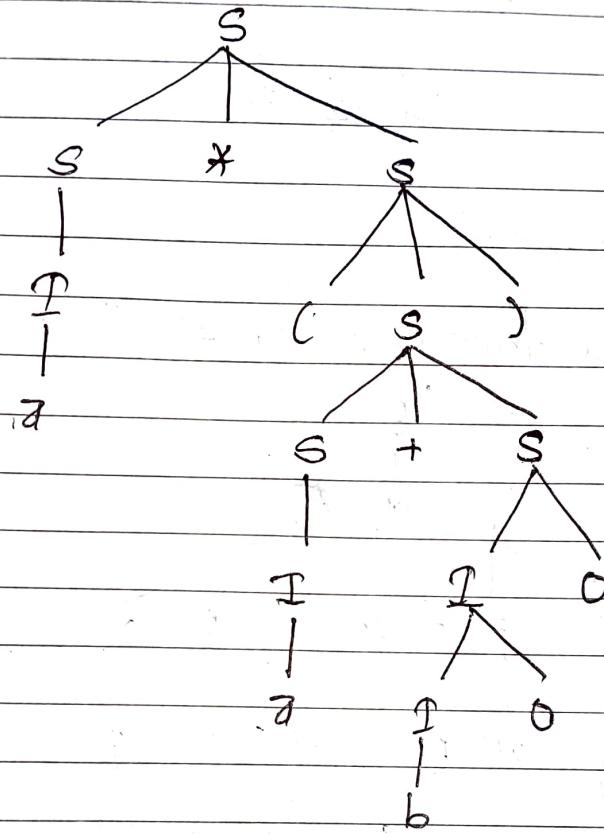
Ex. Consider the Grammar

$$S \rightarrow I \mid S+S \mid S*S \mid (S)$$

$$I \rightarrow a \mid b \mid Ia \mid Ib \mid Io \mid II$$

Construct the parse tree for $a * (a+b00)$

\Rightarrow

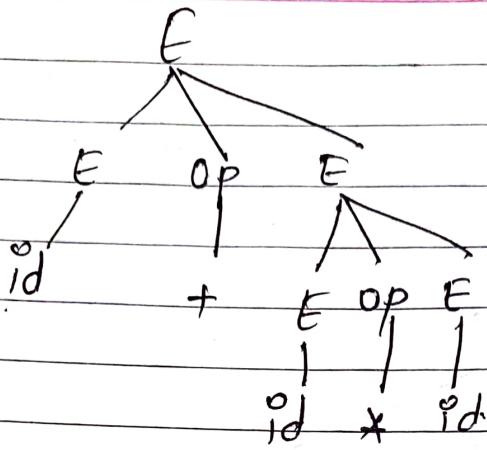


Ex. Construct a grammar defining arithmetic expression and generate parse tree for $id + id * id$ and $(id + id)^*$ $(id + id)$

$$\Rightarrow E \rightarrow E \text{ OP } E \mid (E) \mid id$$

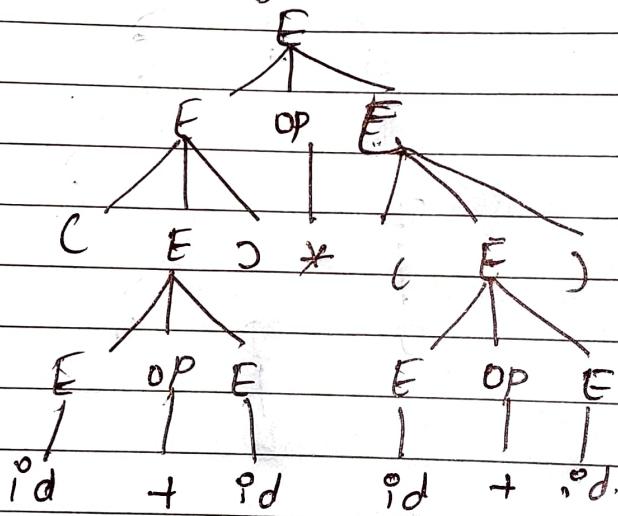
$$\text{OP} \rightarrow + \mid - \mid * \mid / \mid \uparrow$$

Now, the parse tree for $id + id * id$ is



$\text{id} + \text{id} * \text{id}$

Now, parse tree for $(\text{id} + \text{id})^*$



$(\text{id} + \text{id})^*$

* Ambiguous grammar

A grammar $G = (V, T, P, S)$ is said to be ambiguous if there is a string $w \in L(G)$ for which we can derive two or more distinct derivation trees rooted at S and yielding the string w .

In other words, a grammar is ambiguous if it can produce more than one leftmost and more than one rightmost derivation for the same string in the language of the grammar.

$$\text{Ex: } S \rightarrow AB \mid aaB \\ A \rightarrow a \mid Aa \\ B \rightarrow b$$

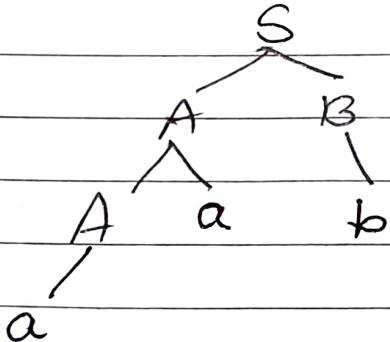
For any string aab , we have two leftmost derivations:

$$\begin{aligned} S &\rightarrow AB \\ S &\rightarrow AaB \quad (\text{Rule: } S \rightarrow aAB) \\ S &\rightarrow aaB \quad (\text{Rule: } A \rightarrow a) \\ S &\rightarrow aab \quad (\text{Rule: } B \rightarrow b) \end{aligned}$$

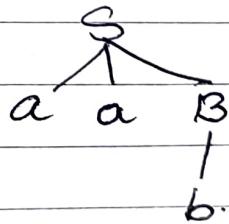
Also,

$$\begin{aligned} S &\rightarrow aAB \\ S &\rightarrow aab \quad (\text{Rule: } B \rightarrow b) \end{aligned}$$

The parse tree for the derivation is:



The parse tree for this derivation is:



Hence, the above grammar is ambiguous grammar.

* Inherently ambiguity

A context free language L is said to be inherently ambiguous if all its grammars are ambiguous. e.g. $L = \{0^i 1^j 2^k \mid i=j \text{ or } j=k\}$

* Regular Grammar

A regular grammar represents a language which is represented by regular expressions. The regular grammar is accepted by NFA and DFA and the language of the grammar is called regular language. A regular grammar is a subset of CFG. The regular grammar may be either left or right linear.

1) Right Linear Regular Grammar

A regular grammar in which all of the productions are of the form $A \rightarrow wB$ or $A \rightarrow w$ where $A, B \in V$ and $w \in T$ is called right linear.

$$\text{Ex: } \begin{array}{l} S \rightarrow 00B \mid 11C \\ B \rightarrow 11C \mid S \mid 1 \\ C \rightarrow 00B \mid 00 \end{array}$$

2) Left Linear Regular Grammar

A regular in which all of the productions are of the form $A \rightarrow Bw$ or $A \rightarrow w$

Where $A, B \in V$ and $w \in T$ is called left linear.

$$\text{Ex: } S \rightarrow B00 \mid C11 \\ B \rightarrow C11 \mid \$ \mid 1 \\ C \rightarrow B00 \mid 00$$

* Equivalence Of Regular Grammar and Finite Automata

Let $G = (V, T, P, S)$ be a right linear grammar of a language $L(G)$, we can construct a finite automaton M accepting $L(G)$ as:

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

where:

α : Consists of symbols $\{\alpha\}$ such that α is either S or a suffix from right hand side of a production in P .

T : Set of input symbols Σ which are terminal symbols
 $[S]$: Start symbol of grammar G , which is start state of FA
 $\{[e]\}$: Set of final states in finite automata

δ : transition function defined as

If A is a variable then,

$\delta([A], e) = [\alpha]$ such that $A \rightarrow \alpha$ is a production.

If ' a ' is a terminal & α is in $(T \cup V)^*$ then

$$\delta([a\alpha], a) = [\alpha]$$

Note: No. of states in automata will be equal to no. of non terminals plus one.

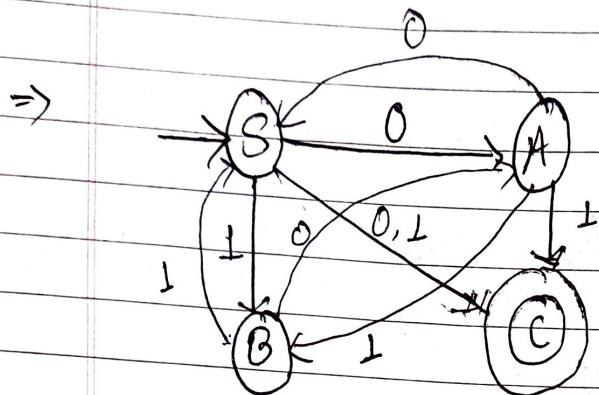
Each state in automata represents each non-terminal in regular grammar. Additional classmate state will be final state

Ex: Convert the following EFGP into FA

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

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

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



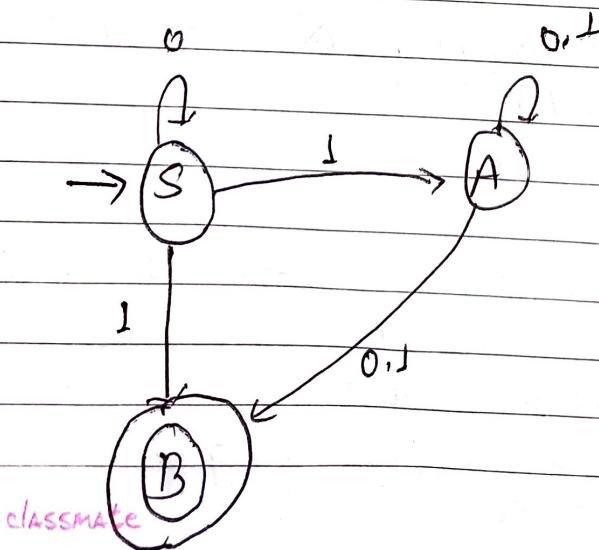
Hence, we can see
no. of non terminals
in LFGP = 3 so, no. of
states in FA = 3 + 1 = 4.

Note: (Imp)

1. For every production $A \rightarrow aB$, make $\delta(A, a) = B$
(Make an edge labelled 'a' from A to B)
2. For every production $A \rightarrow a$, make $\delta(A, a) = \text{Final state}$
3. For every production $A \rightarrow \epsilon$, make $\delta(A, \epsilon) = A$
and A will be Final state

$$Ex: S \rightarrow 0S \mid 1A \mid 1$$

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

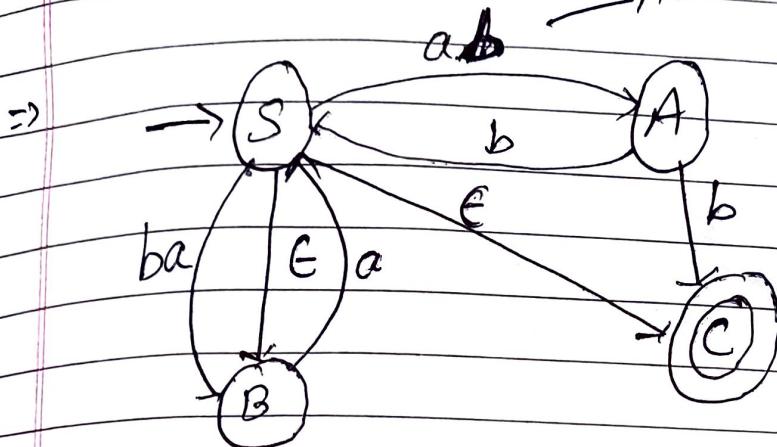


Ex: Convert the following CFG into FA. Also, trace the transition to accept the string abba

$$S \rightarrow aba \mid b \mid bab \mid \epsilon$$

$$A \rightarrow bs \mid b$$

$$B \rightarrow as$$



Now, to trace the transition to accept string abba.

$$S \xrightarrow{ab} A \xrightarrow{b} S \xrightarrow{\epsilon} B \xrightarrow{a} S \xrightarrow{\epsilon} C$$

∴ The transition required to traverse string abba is
 $S \rightarrow A \rightarrow S \rightarrow B \rightarrow S \rightarrow C$.

Ex: Convert the following CFG into FA.

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

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

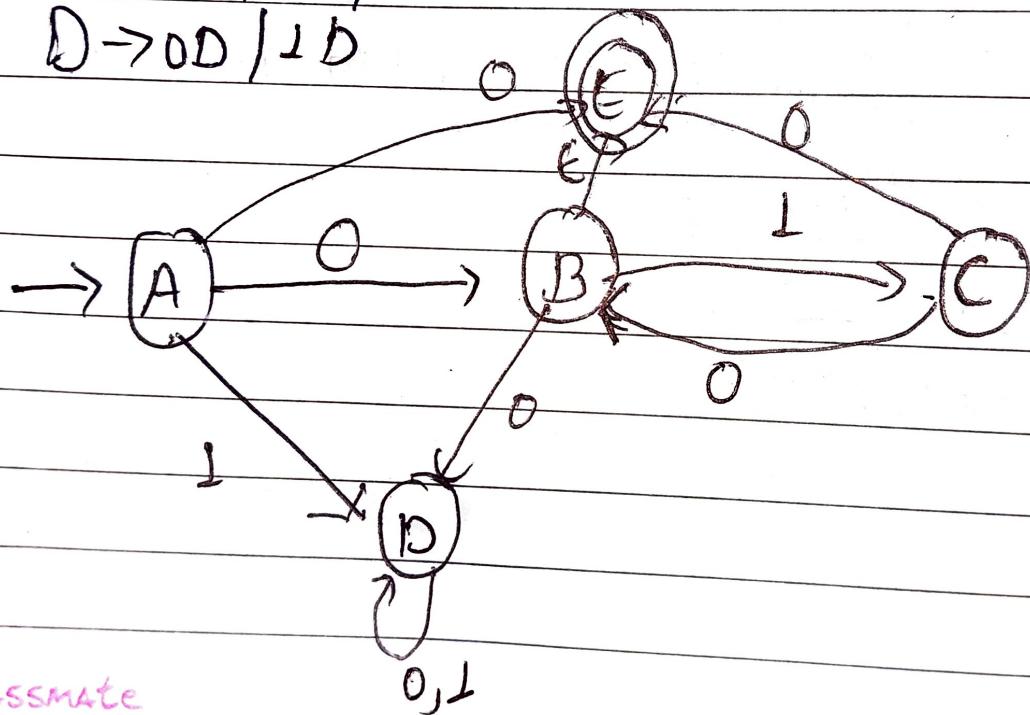
$$C \rightarrow 00B \mid 00$$

Ex: $A \rightarrow_0 B \mid \perp D \mid 0$

$B \rightarrow_0 D \mid \perp C \mid \epsilon$

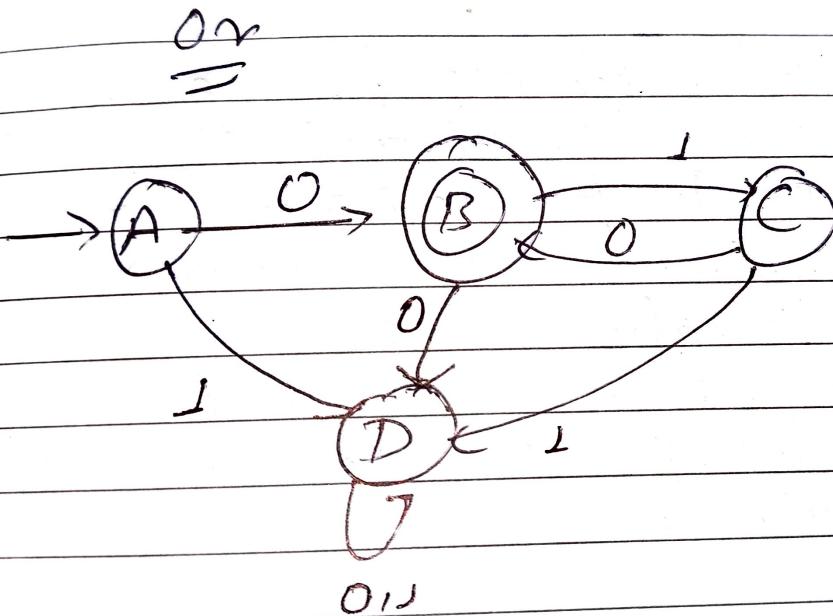
$C \rightarrow_0 B \mid \perp D \mid 0$

$D \rightarrow_0 D \mid \perp b$



classmate

DATE



Simplification of CFG

DATE

In a CFG, it may not be necessary to use all the symbols in V and T , or all the productions in P for deriving strings from the grammar. Thus simplification of CFG refers to elimination of those symbols and productions in G which are not useful or are unnecessary for the derivation of strings. It involves following:

- i) Removal of useless symbols.
- ii) Removal of nullable symbols (ϵ)
- iii) Removal of unit production.

Eliminating ϵ -productions (removal of nullable symbols)

To eliminate ϵ -productions we have to find nullable variables. Nullable variables are the variables that produce ϵ . i.e. Every variable A for which there is a production $A \rightarrow \epsilon$ is nullable. If $A_1, A_2, A_3, \dots, A_k$ are nullable variables and $B \rightarrow A_1, A_2, A_3, \dots, A_k$, then B is also nullable variable.

Algorithm / procedure:

Step-1: To remove $A \rightarrow \epsilon$, look for all productions whose right side contains A .

Step-2: Replace each occurrence of ' A ' in each of these productions with B .

Step-3: Add the resultant productions to the grammar.

Note: ① If $A \rightarrow E$, then A is nullable variable.

② Given $X \rightarrow Y_1 Y_2 \dots Y_n$ if $Y_1, Y_2, Y_3, \dots, Y_n$

are nullable variables then X is also nullable variable.

DATE: _____

Ex: Consider the grammar:- Remove all null prod.

$$S \rightarrow AB$$

$$A \rightarrow \alpha A A / \epsilon$$

$$B \rightarrow b B B / \epsilon$$

∴ Here

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

$$B \rightarrow \epsilon \quad B \text{ is nullable}$$

$S \rightarrow AB$ Since A & B are both nullable S is too.

Now, removing ϵ -production, we get.

$$\begin{array}{|c|c|c|} \hline S & \rightarrow & AB | A | B \\ \hline A & \rightarrow & \alpha AA | \alpha A | \alpha \\ \hline B & \rightarrow & b BB | b B | b \\ \hline \end{array}$$

To remove $A \rightarrow \epsilon$

$$S \rightarrow AB \Rightarrow S \rightarrow B$$

$$A \rightarrow \alpha AA / \epsilon \Rightarrow A \rightarrow \alpha A / \alpha$$

To remove $B \rightarrow \epsilon$

$$S \rightarrow AB | B, A \rightarrow \alpha A / \alpha, B \rightarrow b B / b$$

AAA |

Ex: Remove Null productions from the following

$$S \rightarrow ABAC$$

$$A \rightarrow \alpha A / \epsilon$$

$$B \rightarrow b B / \epsilon$$

$$C \rightarrow c$$

∴ Here

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

1) To eliminate $A \rightarrow \epsilon$

$$S \rightarrow ABAC$$

$$S \rightarrow ABC | BAC | BC$$

$$A \rightarrow \alpha A$$

$$A \rightarrow \alpha$$

classmate new product: $S \rightarrow ABAC | ABC | BAC | BC$

$A \rightarrow \alpha A / \alpha, B \rightarrow b B / \epsilon, C \rightarrow c$

PAGE: _____

2) To eliminate $B \rightarrow e$

$$S \rightarrow AAC \mid AC \mid (\textcircled{Ac}) \mid c$$

$B \rightarrow b$

no need to write

New production:

$$S \rightarrow ABAC \mid ABC \mid BAc \mid BC \mid AAC \mid Ac \mid c$$

$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

$$C \rightarrow C$$

Ex: Remove non productions from following grammar

$$S \rightarrow AB$$

$$A \rightarrow \alpha AA \mid \epsilon$$

$$B \rightarrow \beta BB \mid \epsilon$$

→ Here nullable variables are

$$A \rightarrow \epsilon$$

$$B \rightarrow \epsilon$$

1) To eliminate $A \rightarrow \epsilon$

$$S \rightarrow AB$$

$$S \rightarrow B$$

$$A \rightarrow \alpha AA$$

$$A \rightarrow \alpha A \mid \alpha$$

New production:

$$S \rightarrow AB \mid B$$

$$A \rightarrow \alpha AA \mid \alpha A \mid \alpha$$

$$B \rightarrow \beta BB \mid \epsilon$$

2) To eliminate $B \rightarrow E$

$$S \rightarrow AB|B, \quad S \rightarrow A$$

$$B \rightarrow bBB$$

$$B \rightarrow bB|b$$

New production:

$$S \rightarrow AB|A|B$$

$$A \rightarrow \alpha AA|\alpha A|\alpha$$

$$B \rightarrow bBB|bB|b$$

Elimination (removal) of Unit Productions.

Any production rule of the form $A \rightarrow B$ where $A, B \in N$ Non Terminals is called Unit Production.

The procedure for Removal

Step-1: To remove $A \rightarrow B$, add production $A \rightarrow x$ to the grammar rule whenever $B \rightarrow x$ occurs in the grammar [$x \in T$ Terminal, x can be Null]

Step-2: Delete $A \rightarrow B$ from the grammar

Step-3: Repeat from Step 1 until all unit productions are removed.

Example:

Remove unit productions from the Grammar
whose production rule is given by

$$P: S \rightarrow XY, X \rightarrow a, Y \rightarrow Z|b, Z \rightarrow M, \\ M \rightarrow N, N \rightarrow a$$

⇒ The unit productions are:

$$Y \rightarrow Z, Z \rightarrow M, M \rightarrow N$$

1) To remove $M \rightarrow N$

Here, $M \rightarrow N$ and $N \rightarrow a$
So, we add $M \rightarrow a$

New production rule:

$$P: S \rightarrow XY, X \rightarrow a, Y \rightarrow Z|b, Z \rightarrow a, M \rightarrow a, N \rightarrow a$$

2) To remove $Z \rightarrow M$

Since, $M \rightarrow a$ So, we add $Z \rightarrow a$

New production rule:

$$P: S \rightarrow XY, X \rightarrow a, Y \rightarrow Z|b, Z \rightarrow a, M \rightarrow a, N \rightarrow a$$

3) To remove $Y \rightarrow Z$

Since, $Z \rightarrow a$ So, we add $Y \rightarrow a$

New production rule is:

$$P: S \rightarrow XY, X \rightarrow a, Y \rightarrow a|b, Z \rightarrow a, M \rightarrow a, N \rightarrow a$$

Since, we can't reach the symbols Z, M, N .

So, we remove the unreachable symbols

Final production rule :-

$$P: S \rightarrow XY, X \rightarrow a, Y \rightarrow a/b.$$

Qn. Remove the unit production from following grammar

$$\begin{array}{l} S \rightarrow S + T \mid \pi \\ T \rightarrow T * F \mid F \\ F \rightarrow (S) \mid a \end{array}$$

∴ The unit productions are: $S \rightarrow T$, $T \rightarrow F$
Also, $S \rightarrow T$ & $T \rightarrow F$ so, (S, F) is a unit pair.

Now, after removing unit production, new production rule is:

$$\begin{array}{l} S \rightarrow S + T \mid T * F \mid (S) \mid a \\ T \rightarrow T * F \mid (S) \mid a \\ F \rightarrow (S) \mid a \end{array}$$

Qn. Simplify the grammar $G = (V, T, P, S)$ defined by following production

$$S \rightarrow AaB \mid \epsilon$$

$$A \rightarrow a \mid SA \mid a$$

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

∴ First removing ϵ -production

Here, $S \rightarrow \epsilon$ and $B \rightarrow \epsilon$ are null variables

→ To eliminate $S \rightarrow \epsilon$

$S \rightarrow ASB$

$S \rightarrow AB$

$A \rightarrow \alpha S A$ ~~(1)~~

$A \rightarrow \alpha A$

$B \rightarrow S b S$

$B \rightarrow Sb$

$B \rightarrow bs$

$B \rightarrow b$

New production rule is

$S \rightarrow ASB | AB$

$A \rightarrow \alpha S A | \alpha A | \alpha$

$B \rightarrow Sbs | Sb | bs | A | bb | b$

→ To eliminate $B \rightarrow \epsilon$

$S \rightarrow ASB | AB$

$S \rightarrow AS | A$

~~(1)~~

New Production rule is

$S \rightarrow ASB | AB | AS | A$

$A \rightarrow \alpha S A | \alpha A | \alpha$

$B \rightarrow Sbs | Sb | bs | A | bb | b$

Second, Removing E-products

The unit productions are

$$S \rightarrow A, \quad A \rightarrow B \quad B \rightarrow A$$

So, $S \rightarrow B$ is also unit production by transit.

Now, new production is

$$S \rightarrow ASB \mid AB \mid As \mid ASA \mid aA$$

$$A \rightarrow AsA \mid zA \mid z$$

$$B \rightarrow .Sbs \mid sb \mid bs \mid ASA \mid zA \mid bb \mid z$$

Removal of useless symbols: (Reduction of grammar)

Note: Symbols that are not derivable or unreachable are useless.

(CFG are reduced in two phases:-

Phase 1: Derivation of an equivalent grammar G_1 from the CFG G , such that each variable derives some terminal string.

(Removal of undervivable symbols & their production)

Procedure :-

Step 1: Include all the symbols w_i , that derives some terminal and initialize $i=1$

Step 2: Include symbols w_{i+1} , that derives w_i

Step 3: Increment i and repeat Step 2, until $w_i = w$

Step 4: Include all production rules that have w_i in it.

(Removal of unreachable symbols & their products)

Phase 2: Derivation of an equivalent grammar G''_1 , from the CFG G_1 , such that each symbol appears in a sentential form.

Procedure:

- Step 1: Include the start symbol in Y_1 and initialize i_1 .
- Step 2: Include all symbols Y_{i+1} , that can be derived from Y_i and include all production rules.
- Step 3: Increment i and repeat Step 2, until $Y_{i+1} = Y_i$.

Note: all symbols: both terminals & non-terminals
not present in Y_{i+1} can be removed. (variables)

Example: Find a reduced grammar equivalent to the grammar G_1 , having production rules P

$$P: S \rightarrow AaB, A \rightarrow a, C \rightarrow cBC, E \rightarrow aAe$$

⇒ Phase 1: (Removal of Undeivable symbols)
Terminal Symbols $T = \{a, c, e\}$

$$W_1 = \{A, c, E\}$$

is the set which includes all the symbols that derive the terminal symbols.

$$W_2 = \{A, c, E, S\}$$

is the set which includes all the symbols that can derive the symbols present in the set W_1 .

$$W_3 = \{A, c, E, S\}$$

, , " , " , " , " W_2

$\therefore G' = \{ (A, C, E, S), \{a, c, e\}, P, \{S\} \}$ where,

$P: S \rightarrow AC, A \rightarrow a, C \rightarrow c, E \rightarrow aA/e$

Phase 1: (Removal of unreachable symbols)

$$Y_1 = \{S\}$$

$$Y_2 = \{S, A, C\}$$

includes all the symbols that can be derived from the symbol in Y_1 (see G')

$$Y_3 = \{S, A, C, a, c\}$$

includes all the " . . ." symbols in Y_1

$$Y_4 = \{S, A, C, a, c\}$$

Since $Y_3 \times Y_4$ have same symbols so, we stop here

$G'' = \{ (A, C, S), \{a, c\}, P, \{S\} \}$ which is reduced form of grammar G' , where:

$P: S \rightarrow AC, A \rightarrow a, C \rightarrow c$

Q. Let's take another example for phase 2 only:

$$S \rightarrow aAB/bA$$

$$A \rightarrow aAa/bCa$$

$$B \rightarrow Ca/AC$$

$$C \rightarrow Ab/cB$$

$$D \rightarrow aAbB/dD$$

$$Y_1 = \{S\}$$

$$Y_2 = \{S, A, B, a, b\}$$

$$Y_3 = \{S, A, B, C, a, b, d\}$$

$$Y_4 = \{S, A, B, C, a, b, d\}$$

$\therefore Y_2 \neq \{S, A, B\}$ PAGE YS {same}

P: $S \rightarrow \alpha A\beta / bA$
 $A \rightarrow \alpha A\gamma / b(\alpha)$
 $\beta \rightarrow c\alpha / \alpha C$

$C \rightarrow Ab / CB$

DATE

Note:

Formal definition of reachability.

$A \xrightarrow{*} w$, A is non-terminal ($A \in V$)
 w is a terminal string ($w \in T^*$)
Such a symbol A is derivable!"

$S^* \rightarrow \alpha A\beta$, $\alpha, \beta \in (V \cup T)^*$
 S is a start symbol. $A \in V$
Such a symbol A is reachable.

$S \xrightarrow{*} \alpha A\beta \xrightarrow{*} w$
reachability derivability

Useless productions:

The productions of grammar involving containing useless symbols are called useless productions. Such productions must be removed from grammar to obtain a simplified or reduced grammar.

Chomsky Normal Form (CNF)

In Chomsky Normal Form (CNF) we have a restriction on the length of RHS; which is; elements in RHS should either be two variables or a Terminal

A CFG is in Chomsky Normal Form if the productions are in the following form:

$$A \rightarrow a$$

$$A \rightarrow BC$$

Where, A, B, C are non-terminal & a is a terminal

Thus, a grammar in CNF is one which shouldn't have

$\rightarrow \epsilon$ -production

\rightarrow Unit production

\rightarrow Useless symbols.

\Rightarrow Steps to convert a given CFG to CNF.

Step 1: If the start symbol S occurs on some right side, create a new start symbol S' and a new production $S' \rightarrow S$

Step 2: Remove Null Productions.

Step 3: Remove Unit Productions.

Step 4: Replace each production $A \rightarrow B_1 \dots B_n$ (where $n > 2$) with $A \rightarrow B_1 C$ where $C \rightarrow B_2 \dots B_n$.

Repeat this step for all productions having two or more symbols on the right side.

Step 5: If the right side of any production is in the form $A \rightarrow aB$ where 'a' is a terminal & A and B are non-terminals, then the production is replaced by $A \rightarrow XB$ & $X \rightarrow a$. Repeat this for every production $A \rightarrow a$

Ex: Convert the following CFG to CNF:
 $P: S \rightarrow ASA | aB | \epsilon, A \rightarrow B | S, B \rightarrow b | \epsilon$

\Rightarrow Step-1: Since S appears in RHS, we add a new state S' and $S' \rightarrow S$ is added to the production rule.

$P: S' \rightarrow S, S \rightarrow ASA | aB | \epsilon, A \rightarrow B | S, B \rightarrow b | \epsilon$

Step-2: Remove the null production.

Removal of $A \rightarrow \epsilon$

$S \rightarrow ASA | aB$

$S \rightarrow As | SA | S | aB$

New production rule:

$S' \rightarrow S, S \rightarrow ASA | aB | As | SA | S |, A \rightarrow B | S, B \rightarrow b | \epsilon$

Removal of $B \rightarrow \epsilon$

$S \rightarrow aB$

$S \rightarrow a$

$A \rightarrow B | S$

$A \rightarrow S$

\Rightarrow New production rule:

$S' \rightarrow S, S \rightarrow ASA | aB | a | As | SA | S, A \rightarrow B | S, B \rightarrow b$

Step-3: Removing the unit production.

$S \rightarrow S, S' \rightarrow S, A \rightarrow B, A \rightarrow S$

After removing $S \rightarrow S$: $P: S' \rightarrow S, S \rightarrow ASA | aB | a | As | SA | S, A \rightarrow B | S, B \rightarrow b$

After Removing $S' \rightarrow S$: P: $S' \rightarrow ASA | aB | a | AS | SA$,
 $S \rightarrow ASA | aB | a | AS | SA$
 $A \rightarrow B | S, B \rightarrow b$

After Removing $A \rightarrow B$: P: $S' \rightarrow ASA | aB | a | AS | SA$,
 $S \rightarrow ASA | aB | a | AS | SA$
 $A \rightarrow b | S, B \rightarrow b$

After Removing $A \rightarrow S$: P: $S' \rightarrow ASA | aB | a | AS | SA$,
 $S \rightarrow ASA | aB | a | AS | SA$
 $A \rightarrow b | ASA | aB | a | AS | SA$
 $B \rightarrow b$

Step: Now find out the productions that has more than two variables in RHS.

$S' \rightarrow ASA$, $S \rightarrow ASA$ and $A \rightarrow ASA$

After removing these, we get

P: $S' \rightarrow Ax | aB | a | AS | SA$,
 $S \rightarrow Ax | aB | a | AS | SA$,
 $A \rightarrow b | Ax | aB | a | AS | SA$,
 $B \rightarrow b$,
 $X \rightarrow SA$

Now Change the productions $S' \rightarrow aB$, $S \rightarrow aB$ & $A \rightarrow aB$

Finally, P: $S' \rightarrow Ax | yB | a | AS | SA$,
 $S \rightarrow Ax | yB | a | AS | SA$

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

$$B \rightarrow b$$

$$X \rightarrow SA$$

$$Y \rightarrow a$$

Which is the required CNF for the given CFL

Qn: Remove the nullable symbols

$$S \rightarrow ACbB$$

$$A \rightarrow BC$$

$$B \rightarrow b | \epsilon$$

$$C \rightarrow D | \epsilon$$

$$D \rightarrow d$$

\Rightarrow Nullable variables are

$$B \rightarrow \epsilon, C \rightarrow \epsilon, A \rightarrow BC$$

Removing $B \rightarrow \epsilon$

$$S \rightarrow ACbB | Acb$$

$$A \rightarrow BC | \epsilon$$

$$B \rightarrow b$$

$$C \rightarrow D$$

$$D \rightarrow d$$

Removing $C \rightarrow \epsilon$

$$S \rightarrow A \cancel{C} bB | AbB | A \cancel{C} b | Ab$$

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

$$B \rightarrow b$$

$$C \rightarrow D$$

$$D \rightarrow d$$

Removing $A \rightarrow BC$ (ie. $A \rightarrow \epsilon$)

$$S \rightarrow ACbB / AB\beta / ACb / Ab / CbB / bB / cb / b$$

$$A \rightarrow BC / C / B$$

$$\beta \rightarrow b$$

$$C \rightarrow D$$

$$D \rightarrow d$$

Ex: Remove the nullable symbol.

$$S \rightarrow aSb / aAb$$

$$A \rightarrow \epsilon$$

→ Nullable variable is $A \rightarrow \epsilon$, removing it

$$S \rightarrow asb / aAb / ab$$

Ex: Remove the nullable:

$$S \rightarrow AB$$

$$A \rightarrow aAA / G$$

$$B \rightarrow bBB / \epsilon$$

Note: Here $S \rightarrow AB$, so $S \leftrightarrow \epsilon$ so, the language itself contains ϵ , so no need to remove $S \rightarrow \epsilon$

Removing $A \rightarrow \epsilon$

$$S \rightarrow AB / B$$

$$A \rightarrow aAA / aA / a$$

$$B \rightarrow bBB$$

Removing $B \rightarrow \epsilon$

$$S \rightarrow AB / B / A / \epsilon$$

$$S \rightarrow aAA / aA / a$$

$$B \rightarrow bBB / bB / b$$

Ex: Remove the unit production

$$S \rightarrow aA | B$$

$$A \rightarrow ba | b.b.$$

$$B \rightarrow A | bba$$

$$\Rightarrow S \rightarrow aA | ba | bb | bba$$

$$A \rightarrow ba | bb$$

$$B \rightarrow ba | bb | bba$$

Removing unreachable symbols.

$$\Rightarrow S \rightarrow aA | ba | bb | bba$$

$$A \rightarrow ba | bb$$

Ex: Convert following grammar into CNF

$$S \rightarrow aA | bB$$

$$A \rightarrow aAAA | a$$

$$B \rightarrow bBB | b$$

Find out the productions that has more than two variables in RHS.

$$A \rightarrow aAAA, B \rightarrow bBB$$

After removing these, we get modified rule as:

$$S \rightarrow aA | bB$$

$$A \rightarrow XY | a$$

$$B \rightarrow ZB | b$$

$$X \rightarrow aA$$

$$Y \rightarrow AA$$

$$Z \rightarrow bB$$

Now, change the productions:-

$$S \rightarrow aA, S \rightarrow bB, X \rightarrow aA, Z \rightarrow bB$$

Modified rule:-

$$S \rightarrow MA | NB$$

$$A \rightarrow XY | a$$

$$B \rightarrow ZB | b$$

$$X \rightarrow MA$$

$$Y \rightarrow AA$$

$$Z \rightarrow NB$$

$$M \rightarrow a$$

$$N \rightarrow b$$

Since, it is in CNF. This is the final grammar.

Greibach Normal Form (GNF)

A CFG is in Greibach Normal Form if all its production rules are in the form

$$A \rightarrow aV^*$$

$$A \rightarrow a.$$

where, A is a non terminal,
a is a terminal,
V is a set of non terminals

Ex. $X \rightarrow bABCDE$ or $X \rightarrow bX_1X_2X_3X_4$
 $X \rightarrow c$ $X \rightarrow b$

Steps to convert a given CFG to GNF:

Step-1: Check if the given CFG has any unit production or null production and remove if there are any.

Step-2: Check whether the CFG is already in CNF and convert it to CNF if it is not.

Step-3: Change the name of the Non-Terminal Symbols into some A_i in descending order of;

Ex: $S \rightarrow EA \quad | \quad BB$

 $B \rightarrow b \quad | \quad SB$
 $C \rightarrow b$
 $A \rightarrow a$

CNF

Note: $S \rightarrow S_1$

 $S \rightarrow CA_1BB$
 $B \rightarrow b/SB$
 $C \rightarrow b$
 $A \rightarrow a$

→ Step 1 & 2 are satisfied already

Now, replace S with A₁, C with A₂, A with A₃, B with A₄

We get.

$$\begin{aligned}A_1 &\rightarrow A_2 A_3 / A_4 A_4 \\A_4 &\rightarrow b / A_1 A_4 \\A_2 &\rightarrow b \\A_3 &\rightarrow a\end{aligned}$$

Step 4: Alter the rules so that the Non-Terminals are in ascending order, such that, If the Production is of the form $A_i \rightarrow A_j x$, then $i \leq j$ and should never be $i > j$

$$\begin{aligned}A_4 &\rightarrow b / A_1 A_4 \\A_4 &\rightarrow b / A_2 A_3 \quad \cancel{A_1 A_4} \quad \cancel{A_3 A_4} \\A_4 &\rightarrow b / b A_3 \\A_1 &\rightarrow A_2 A_3 / A_4 A_4 \\A_4 &\rightarrow b / A_1 A_4 \\A_4 &\rightarrow b / A_2 A_3 A_4 / A_4 A_4 A_4 \\A_4 &\rightarrow b / b A_3 A_4 / A_4 A_4 A_4\end{aligned}$$

Step 5: Remove left recursion. (of the form $A_i \rightarrow A_i x$)

Introduce a new variable to remove the left recursion.

Here, $A_4 \rightarrow b / b A_3 A_4 / A_4 A_4 A_4$

introduce. $Z \rightarrow A_4 A_4 Z / A_4 A_4$

New production rule PS:

$$A_4 \rightarrow b / b A_3 A_4 / b Z / b A_3 A_4 Z$$

1000 with the grammar

$$A_1 \rightarrow A_2 A_3 \mid A_4 A_4$$

$$A_4 \rightarrow b \mid b A_3 A_4 \mid b z \mid b A_3 A_4 z$$

$$z \rightarrow A_4 A_4 \mid A_4 A_4 z$$

$$A_2 \rightarrow b$$

$$A_3 \rightarrow a$$

Since, it's not PNF in GNF,

$$\text{Now, } A_1 \rightarrow b A_3 \mid b A_4 \mid b A_3 A_4 A_4 \mid b z A_4 \mid b A_3 A_4 z A_4$$

$$A_4 \rightarrow b \mid b A_3 A_4 \mid b z \mid b A_3 A_4 z$$

$$z \rightarrow b A_4 \mid b A_3 A_4 A_4 \mid b z A_4 \mid b A_3 A_4 z A_4 \\ b A_4 z \mid b A_3 A_4 A_4 z \mid b z A_4 z \mid b A_3 A_4 z A_4 z$$

$$A_2 \rightarrow b$$

$$A_3 \rightarrow a$$

Which is in GNP.

Ex. Convert the following CFG to GNF (Another method)

$$S \rightarrow S S$$

$$S \rightarrow a S b$$

$$S \rightarrow a b$$

If not in CNF, do this method.

⇒ Step-1: Introduce non-terminals A, B for a, b

$$A \rightarrow a$$

$$B \rightarrow b$$

$$S \rightarrow S S$$

$$S \rightarrow A S B$$

$$S \rightarrow A B$$

Step-2: Introduce an order among non-terminals by renaming them as follows

Note: A grammar $G(N, T, P, S)$ is left recursive if it has a production in the form $A \rightarrow Ax \mid B$
where $A \in N$ non terminals
 $x = \text{string composed from } (V \cup T)^*$

$$\begin{aligned}A_1 &\rightarrow A_1 A_1 \\A_1 &\rightarrow A_2 A_1 A_3 \\A_1 &\rightarrow A_2 A_3 \\A_2 &\rightarrow a \\A_3 &\rightarrow b\end{aligned}$$

Step-4: Check if the non-terminals are in ascending order.

$$\begin{aligned}A_1 &\rightarrow \cancel{A_2 A_1 A_3} A_1 A_1 \\A_1 &\rightarrow A_2 A_1 A_3 \\A_1 &\rightarrow A_2 A_3 \\A_2 &\rightarrow a \\A_3 &\rightarrow b\end{aligned}$$

It can be written as

$$\begin{aligned}A_1 &\rightarrow \cancel{A_2 A_1 A_3} A_1 A_1 \mid A_2 A_1 A_3 \mid A_2 A_3 \\A_2 &\rightarrow a \\A_3 &\rightarrow b\end{aligned}$$

Step-5: Remove the left recursion, introduce new non-terminal Z such that,

$$Z \rightarrow A_1 Z \mid A_1$$

Then, new rule is

$$\begin{aligned}A_1 &\rightarrow A_2 A_1 A_3 Z \mid A_2 A_3 Z \mid A_2 A_1 A_3 \mid A_2 A_3 \\Z &\rightarrow A_1 Z \mid A_1 \\A_2 &\rightarrow a \\A_3 &\rightarrow b\end{aligned}$$

Which is still not in GNF.

$$\begin{aligned}A_1 &\rightarrow a A_1 A_3 Z \mid a A_3 Z \mid a A_1 A_3 \mid a A_3 \\Z &\rightarrow a A_1 A_3 Z \mid a A_3 Z \mid a A_1 A_3 \mid a A_3 \\&\quad \mid a A_1 A_3 Z \mid a A_3 Z \mid a A_1 A_3 \mid a A_3 \\A_2 &\rightarrow a, A_3 \rightarrow b\end{aligned}$$

Ex: Remove left recursive production from following.

$$\begin{aligned} S &\rightarrow AA \mid 0 \\ A &\rightarrow AA \mid S \mid 1 \end{aligned}$$

→ Introduce a new variable Z such that

$$Z \rightarrow ASZ \mid AS$$

Now, the new rules

$$\begin{aligned} S &\rightarrow AA \mid 0 \\ A &\rightarrow \cancel{A} OSZ \mid 1Z \mid OS \mid 1 \\ Z &\rightarrow ASZ \mid AS \end{aligned}$$

Ex: Remove left recursive production.

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \end{aligned}$$

→ Since, there are two left recursions $E \rightarrow E + T$ and $T \rightarrow T * F \mid F$.

Introduce two new non-terminals Y and Z such that

$$\begin{aligned} Y &\rightarrow +TY \mid +T \\ Z &\rightarrow *FZ \mid *F \end{aligned}$$

So, new rule is

$$\begin{aligned} E &\rightarrow TY \mid T \\ Y &\rightarrow +TY \mid +T \\ T &\rightarrow FZ \mid F \\ Z &\rightarrow *FZ \mid *F \\ F &\rightarrow (E) \mid a \end{aligned}$$

Ex: Remove the left recursion and convert it into GNF

$$S \rightarrow AA10$$

$$A \rightarrow SS11$$

∴ The grammar is already in CNF, so do by LR method.

Replace S with A, A with A₂

$$A_1 \rightarrow A_2 A_2 | 0$$

$$A_2 \rightarrow A_1 A_1 | 1$$

Alter the rules so that the non-terminals are in ascending order.

$$A_1 \rightarrow A_2 A_2 | 0$$

$$A_2 \rightarrow A_2 A_2 A_1 | OA_1 | 1$$

Now, removing the left recursion

Introduce new variable z such that

$$z \rightarrow A_2 A_1 z | A_2 A_1$$

Now, new rule is

$$A_1 \rightarrow A_2 A_2 | 0$$

$$A_2 \rightarrow OA_1 z | 1 z | OA_1 | 1$$

$$z \rightarrow A_2 A_1 z | A_2 A_1$$

It's still not in GNF

$$A_1 \rightarrow OA_1 z A_2 | 1 z A_2 | OA_1 A_2 | 1 A_2 | 0$$

$$A_2 \rightarrow OA_1 z | 1 z | OA_1 | 1$$

$$z \rightarrow OA_1 z A_1 z | 1 z A_1 z | OA_1 A_1 z | 1 A_1 z | OA_1 z A_1 | 1 A_1$$

$$1 z A_1 | OA_1 A_1 | 1 A_1$$

(no need to do this)

DATE

Now, you can replace A_1 with S and A_2 with A to match with original grammar.

$$S \rightarrow OSZA | IZ | OSAI | JA | \emptyset$$

$$A \rightarrow OSZ | IZ | OS | I$$

$$Z \rightarrow OSZSZ | IZSZ | OSSZ | JSZ | OSZS | IZS | OSZS | IS$$

IS

Ex: Remove left recursion and convert into CNF

$$S \rightarrow AB$$

$$A \rightarrow BS | b$$

$$B \rightarrow SA | a$$

⇒ It is already in CNF So do by 1st method

$$A_1 \rightarrow A_2 A_3$$

$$A_2 \rightarrow A_3 A_1 | b$$

$$A_3 \rightarrow A_1 A_2 | a$$

Check if non-terminals are in ascending order

$$A_1 \rightarrow A_2 A_3$$

$$A_2 \rightarrow A_3 A_1 | b$$

$$A_3 \rightarrow A_2 A_3 A_2 | a$$

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

→ New rule

$$A_1 \rightarrow A_2 A_3$$

$$A_2 \rightarrow A_3 A_1 | b$$

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

Since, there is a left recursion $A_3 \rightarrow A_3A_1, A_3A_2$

So, introduce a new variable Z such that

$$Z \rightarrow A_1A_3A_2 \cup A_1A_3A_2$$

So, new rules

$$\begin{aligned} A_1 &\rightarrow A_2A_3 \\ A_2 &\rightarrow A_3A_1 \cup b \end{aligned}$$

$$\begin{aligned} A_3 &\rightarrow bA_3A_2Z \cup aZ \cup bA_3A_2 \cup a \\ Z &\rightarrow A_1A_3A_2 \cup A_1A_3A_2 \end{aligned}$$

It is still not in GNF.

$$\begin{aligned} A_1 &\rightarrow A_3A_1A_3 \cup bA_3 \\ A_2 &\rightarrow bA_3A_2ZA_1 \cup ZaA_1 \cup bA_3A_2A_1 \cup ZaA_1 \cup b \\ A_3 &\rightarrow bA_3A_2Z \cup aZ \cup bA_3A_2 \cup a \\ Z &\rightarrow A_1A_3A_2Z \cup A_1A_3A_2 \end{aligned}$$

It is still not in GNF.

$$\begin{aligned} A_1 &\rightarrow bA_3A_2ZA_1A_3 \cup aZA_1A_3 \cup bA_3A_2A_1A_3 \cup \\ &\quad ZaA_1A_3A_2 \not\cup bA_1A_3 \not\cup bA_3 \end{aligned}$$

$$\begin{aligned} A_2 &\rightarrow bA_3A_2ZA_1 \cup ZaA_1 \cup bA_3A_2A_1 \cup ZaA_1 \cup b \\ A_3 &\rightarrow bA_3A_2Z \cup ZaZ \cup bA_3A_2 \cup a \\ Z &\rightarrow bA_3A_2ZA_1A_3A_2Z \cup ZaZA_1A_3A_2Z \cup \\ &\quad bA_3A_2A_1A_3A_2Z \cup ZaA_1A_3A_2A_1A_2 \\ &\quad bA_3A_2ZA_1A_3A_2 \cup ZaA_1A_3A_2A_1A_2 \\ &\quad bA_3A_2A_1A_3A_2Z \cup ZaA_1A_3A_2A_1A_2 \end{aligned}$$

Baekus-Naur Form (BNF)

It is a notation technique used for context free grammar. This is used for specifying the syntax of the language. It can be expressed as

$\langle \text{Symbol} \rangle ::= \text{Exp}^1 \mid \text{Exp}^2 \mid \text{Exp}^3 \dots \dots$
 Where,

Symbol is non-terminal

and Exp is a sequence of symbols
 (terminal & non-terminal)

Here the concept is similar to CFG, only the difference, instead of using symbol (\rightarrow) in a product, we use symbol ($::=$). We enclose all non terminals in curly brackets $\langle \rangle$.

Ex:

$\langle \text{integer} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{integer} \rangle \langle \text{digit} \rangle$
 $\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid \dots \mid 9$

Ex: The BNF for identifiers as :-

$\langle \text{Identifier} \rangle ::= \langle \text{letter or underscore} \rangle \text{Identifier} \langle \text{symbol} \rangle$
 $\langle \text{letter or underscore} \rangle ::= \langle \text{letter} \rangle \mid$

$\langle \text{symbol} \rangle ::= \langle \text{letter or underscore} \rangle \mid \langle \text{digit} \rangle$

$\langle \text{letter} \rangle ::= a \mid b \mid \dots \mid z$

$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid \dots \mid 9$

Context Sensitive Grammar (CSG)

DATE

A CSG is a formal grammar on which left hand sides and right hand sides of any production may be surrounded by a context of terminal and non-terminal symbols.

A formal grammar, $G = (V, T, P, S)$ is context sensitive if all the production ' p ' are of the form

$$\alpha A \beta \rightarrow \alpha Y \beta \quad \text{where } A \in V$$

$$\alpha \in V$$

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

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

The name Context sensitive is explained by α and β that form the context of A & determine whether A can be replaced with Y or not. Thus the production $A \rightarrow Y$ can only be applied in contexts where α occurs to left of A & β occurs to right of A .

$\{a^n b^n c^n \mid n \geq 1\}$ is a context sensitive language defined as:

$$S \rightarrow \alpha SBC \mid \alpha BC$$

$$CB \rightarrow BC$$

$$\alpha B \rightarrow \alpha b$$

$$bB \rightarrow bb$$

$$bc \rightarrow bc$$

$$cC \rightarrow cc$$

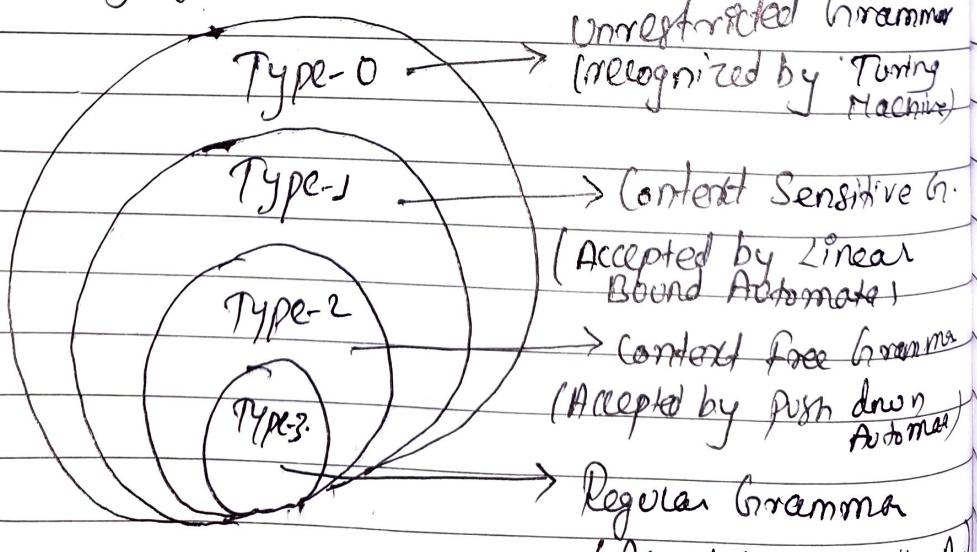
Chomsky hierarchy

Chomsky hierarchy represents the class of languages that are accepted by the different machines. It was described by Noam Chomsky in 1956.

According to Chomsky hierarchy grammar is divided onto 4 types as follows:

1. Type 0 is known as unrestricted grammar
2. Type 1 ~~is~~ is known as context-sensitive grammar
3. Type 2 is known as context-free grammar
4. Type 3 Regular grammar

Here, the restriction of grammar increases from going to type 0 from type 3. The more restriction increases, the more ~~complex~~ easier the grammar becomes and requires less powerful machine. If grammar is simple (type-0), it requires powerful machine as Turing machine to accept the language generated by grammar.



classmate Chomsky hierarchy

~~X Type -0- grammar~~ (Unrestricted grammar / phase. Structured grammar.)

They include all formal grammars. They are recognized by Turing machine. It is used to generate recursive enumerable language (REL) accepted by TM. Formally, $G = \{V, T, P, S\}$

where $P: \alpha \rightarrow \beta$

where, $\alpha : (V + T)^*$ $V (V + T)^*$
where, V: Variables
T: Terminals.

and β is $(V + T)^*$.

Here, there must be at least one variable on the left side of the production.

Ex: $Sab \rightarrow ba$
 $A \rightarrow S$

Variables are S, A, & Terminals.

~~+ Type-1-grammar~~ (Context Sensitive grammar / length increasing grammar / non-contracting grammar)

It generates CS which is accepted by Linear Bounded Automata (LBA).

Here, all type 1 grammar should be Type 0. Grammar production is in the form of: $\alpha \rightarrow \beta$ where,

β is $(V + T)^*$ & $|\alpha| \leq |\beta|$

Here, α cannot be ϵ and $|\alpha| \leq |\beta|$

Ex: $S \rightarrow AB$ $B \rightarrow b$
 $AB \rightarrow abc$

~~+ Type-2-grammar~~ (Content free grammar)

It generates CFG accepted by PDA. It should be type 1. The left-hand side of production can have only one variable & there is no restriction PAGE NO. P. _____

Production is in the form of $X \rightarrow \beta$ where
 $|X| = 1$

$$\text{Ex: } S \rightarrow AB$$

$$A \rightarrow a$$

$$B \rightarrow b$$

$$X \in V$$

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

* Type-3 Grammar

→ It generates regular languages accepted by FA. These languages are exactly all languages that can be accepted by FA. It is the most restricted form of grammar.
 $G = \{V, T, P, S\}$

It should be in the given form only.

$$V \rightarrow VT \mid T \quad (\text{left-regular grammar})$$

or

$$V \rightarrow T V \mid T \quad (\text{right-regular grammar})$$

Ex: $S \rightarrow a$. This is strictly regular grammar.

Note: There is another form of regular grammar. Called extended regular grammar.

$$V \rightarrow VT^* \mid T^*$$

(extended left-regular gr.)

$$V \rightarrow T^* V \mid T^*$$

(extended right-regular gr.)

Pumping Lemma for Context Free Language (CFL)

If A is a context free language, then, A has a pumping length ' p ' such that any string ' s ', where $|s| > p$, may be divided into 5 pieces $S = UVXYZ$ such that the following conditions might be true.

- i) $Uv^ixy^iz \in A$ for every $i \geq 0$
- ii) $|vy| > 0$
- iii) $|Vxy| \leq p$

Pumping Lemma (for CFL) is used to prove that a language is NOT context free.

To prove that a language is Not Context free
Using pumping lemma we follow the below steps:-
(We prove using contradiction)

- Assume that A is context free
- It has a pumping length (say p)
- All strings longer than p can be pumped $|s| > p$
- Now find a string ' s ' in A such that $|s| > p$
- Divide S into $UVXYZ$
- Show that $Uv^ixy^iz \notin A$ for some i
- Then consider the ways that S can be divided into $UVXYZ$.
- Show that none of these can satisfy all the 3 pumping conditions at the same time.
- S cannot be pumped == CONTRADICTION

Ex: Show that $L = \{a^N b^N c^N \mid N \geq 0\}$ is not Context free.

→ i) Suppose that L is context free
 ii) Let pumping length p exists such that any string $s \in L$ is
 $s = a^p b^p c^p$, $|s| = 3p > p$

iii) Let's divide s into $uvxyz$.

Assume that $p=4$, Then

$$s = a^4 b^4 c^4$$

$$\therefore s = aaaa bbbb cccc$$

Now, let's see all possible ways in which we can divide s into $uvxyz$

Case 1: v and y each contain only 1 type of symbol

$$s = \underbrace{aaaa}_{U} \underbrace{b}_{V} \underbrace{bbb}_{X} \underbrace{cc}_{Y} \underbrace{ccc}_{Z}$$

Now:
 $Wx = \underbrace{aaaa}_{U} \underbrace{b}_{V} \underbrace{bbb}_{X} \underbrace{ccc}_{Z} \Rightarrow p=4$
 but $|Wx| = 9 > p$
 so, make sure $|Wx| \leq p$

Case 2: Either v or y has more than one kind of symbol

$$s = \underbrace{aaaa}_{U} \underbrace{b}_{V} \underbrace{bbb}_{X} \underbrace{ccc}_{Y} \underbrace{ccc}_{Z}$$

Now, we show that $uv^ixy^iz \notin L$ for some $i \geq 2$
 let's take $i=2$

Case 1: $s = aaaa abbb bc ccc \rightarrow a^6 b^4 c^5 \notin L$

Case 2: $s = aa aabb aabb bbb cccc \notin L$

Hence, ourie isn't satisfied. So, it contradicts our assumption. Hence, L is not context free

Ex: Show that $L = \{ww \mid w \in \{0,1\}^*\}$ is not Context free.

- ⇒ p) Suppose that L is context free.
- p) Let pumping length = P exists such that any string $s \in L$ is: $s = P_1 P_0 P_1 P$
- iii) Let's divide s into $uvxyz$

Assume that $P = 5$, then

$$s = 0^5 1^5 0^5 1^5$$

$$s = 000001111000001111$$

Now, let's see the possible ways in which we can divide s into $uvxyz$.

Case 1: vxy does not straddle a boundary.

$$\underbrace{00000}_{U}, \underbrace{1111}_{vxy} \underbrace{000001111}_{Z}$$

$$\text{Here, } |vxy| = 3 \leq P = 5$$

so, condition (3) is satisfied.

$$|vy| = 2 > 0 \quad (\text{condition (2) is satisfied})$$

Now, we show that $uv^pxyz \notin L$ for some p .

uv^pxyz lets take $p = 2$

$$s = uv^2xy^2z$$

$$= 00001111000001111$$

$$= 0^5 1^7 0^5 1^5 \notin L$$

so, it doesn't satisfy condition (1)

Here, case 1 already doesn't satisfy, so no need to check other cases.

So, it contradicts our assumption. And hence

L cannot be pumped so it is not Context free.

Application of pumping lemma for CFL

- Pumping lemma is used to check whether a grammar is context free or not.
- Used in some rules of thumb.

Closure properties of CFL

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 / lemmas are often called closure properties of context-free language.

Some of the principle closure properties for context-free languages are:

1) Union (closed)

The CFL is closed under union.

Let's suppose two CFL L_1 and L_2 and corresponding grammars are G_1 and G_2 . Let start symbol for G_1 is S_1 and G_2 is S_2 .

Let's take a string for S_1 where a followed by equal no. of b. & for S_2 is b followed by equal no. of c (i.e. $L_1: a^n b^n$, $L_2: b^n c^n$)

$$\begin{array}{l} L_1 \rightarrow G_1 \\ L_2 \rightarrow G_2 \end{array}$$

$S_1 \rightarrow a S_1 b / \epsilon$
$S_2 \rightarrow b S_2 c / \epsilon$

Their union is: $S \rightarrow S_1 / S_2$

iii) Concatenation (closed)

- ⇒ The CFL is closed under concatenation.
Let's take same example as above.

$$S_1 \rightarrow aS_1b | \epsilon$$

$$S_2 \rightarrow bS_2c | \epsilon$$

Then, Concatenation: $S \rightarrow S_1.S_2$

So, It is closed under concatenation.

iv) Kleene closure (closed)

- ⇒ The CFL is closed under Kleene closure.
In above example.

$$L_1 \rightarrow G_{L_1}$$

$$L_2 \rightarrow G_{L_2}$$

$$\boxed{L_1^* \rightarrow G_{L_1} \\ S \rightarrow S_1S | \epsilon}$$

which is closed.

v) Intersection (not closed)

- ⇒ The CFL is not closed under intersection.

Let's suppose L_1 is $a^n b^n c^m$

$$L_2 \text{ is } a^m b^n c^n$$

Then, $L_1 \cap L_2$ is $a^n b^n c^n \rightarrow \text{CSL}$

context sensitive.

So, It can't be context free. So, It is not closed under intersection.

vi) Complementation (not closed)

- ⇒ Take above same example: Let's suppose CFL is closed

Under complement and $L_1 \cap L_2 = \overline{L_1} \cup \overline{L_2}$ under complement

Suppose, $\overline{L_1} \cap \overline{L_2}$ is closed, then $\overline{L_1}$ is CFL, $\overline{L_2}$ is CFL, then $\overline{L_1} \cup \overline{L_2}$ is CFL, whose

DATE



Complement becomes CFL.

So, $L_1 \cap L_2$ becomes CFL. Which
contradicts the property of intersection.
So, it is not closed.

Questions asked from this Chapter

Q. Explain about Chomsky's Hierarchy about the language and grammars. (2028-5 marks) (2021-5 marks) (2067-5 marks)

Q. Convert the following grammar into Chomsky Normal Form.

$$\begin{array}{l} S \rightarrow abS \quad | \quad a \quad | \quad aAb \\ A \rightarrow bS \quad | \quad aAAb \quad | \quad \epsilon \end{array} \quad (2028-5 \text{ marks})$$

Q. Define Chomsky Normal form and Greibach Normal form in reference to CFG. Give a suitable example each. (2026-5 marks)

Q. Convert following to CNF: $S \rightarrow abSb|aa$. Explain CNF

Q. Convert the following CFG into CNF

$$\begin{array}{ll} S \rightarrow aAa \quad | \quad bBb \quad | \quad \epsilon \\ A \rightarrow c \quad | \quad a & D \rightarrow A \quad | \quad B \quad | \quad ab \\ B \rightarrow c \quad | \quad a & \\ C \rightarrow CDEFG & \end{array} \quad (2023-5 \text{ marks})$$

Q. Explain closure properties of CFL with examples. (2069-10 marks)

Q. Explain about Unrestricted Grammar. How they differ with CFG? Explain. (2021-5 marks, 2067-5 marks)

Q. Give a detailed description of ambiguity in context free grammar. (2069-5 marks)

Q. Prove that CFG is ambiguous. (2020-5 marks)

Q. What is CFG? Define CFG for palindromes with alphabets {0,1,3}. (2023-5 marks)

Q. Define regular grammar? What is the relation of regular grammar with other grammar? (2024-5 marks)