

## UNIT-II

### Regular expressions

Regular Expressions :- A regular expression is a string that describes the whole set of strings according to a certain syntax rule.

(or)

Regular expressions are useful for representing certain set of strings in an algebraic fashion. RE describes the language accepted by finite state automata.

The following are the notations used in regular expressions

- $\phi, \epsilon$  and  $a$  belonging to  $\Sigma$ , are all regular expressions and are called the primitive regular expressions.
- If  $R_1$  and  $R_2$  are two regular expressions then union of these represented as  $R_1 \cup R_2$  (or)  $R_1 + R_2$  also regular expression.
- If  $R_1$  and  $R_2$  are two regular expressions then concatenation of these represented as  $R_1 R_2$  is also regular expressions.
- The Kleene closure of a regular expression  $R$  is denoted by  $R^*$  is also a regular expression.
- If  $R$  is a regular expression then  $(R)$  is also regular exp.

Example :- if  $\Sigma = \{a|b\}$ , if  $R_1 = c$  and  $R_2 = \phi$

then  $c + \phi$  i.e  $R_1 + R_2$  is also a regular expression.

Regular set :- Any set represented by a regular expression is called regular set.

If  $a, b$  are the elements of  $\Sigma$  then regular expressions

- $a$  denotes the set  $\{a\}$
- $ab$  denotes the set  $\{a, b\}$
- $ab$  denotes the set  $\{ab\}$
- $a^*$  denotes the set  $\{\epsilon, a, aa, aaa, \dots\}$
- $(a+b)^*$  denotes the set  $\{\epsilon, a, b, aa, bb, ba, aaa, \dots\}$

Regular language :- A language represented by a regular expression defines a regular language. If ' $R$ ' is a regular expression, then  $L(R)$  denotes the language associated with  $R$ .

The language  $L(R)$  denoted by any regular expression  $R$  is defined by the following rules.

- $\emptyset$  is a regular expression denoting the empty set.
- $\epsilon$  is a regular expression denoting the set  $\{\epsilon\}$
- For every  $a \in \Sigma$ ,  $a$  is regular expression denoting  $\{a\}$ .

Example :-  $R_E = \{abta\}$

$$\text{i.e } L(R) = \{aa, aba, abb, a, abbba, \dots\}$$

The set of all strings of  $a$ 's and  $b$ 's that begin and end with ' $a$ '.

$$\rightarrow R_E = a^* b^*$$

$$L(R) = \{\epsilon, a, ab, ab, bb, aa, abb, aabb, \dots\}.$$

## Identity rules (or) Algebra of R.E's :-

Let  $R, S$  and  $T$  be any arbitrary regular expressions.

Then, the following properties are true:

$$1) (R+S)+T = R+(S+T)$$

$$2) R+R = R$$

$$3) R+\emptyset = \emptyset + R = R$$

$$4) R+S = S+R$$

$$5) R\emptyset = \emptyset R = \emptyset$$

$$6) R \cdot E = \epsilon R = R$$

$$7) (RS)T = R(ST)$$

$$8) R(S+T) = RS+RT$$

$$9) (S+T)R = SR+TR$$

$$10) \emptyset^* = \epsilon^* = \epsilon$$

$$11) R^*, R^* = R^* = (R^*)^*$$

$$12) R \cdot R^* = R^* \cdot R = R^* = \epsilon + RR^*$$

$$13) (R+S)^* = (R^*S^*)^* = (R^*+S^*)^*$$

$$14) (RS)^* = (R^*S^*)^* = (R^*+S^*)^*$$

## Manipulations of Regular Expressions 8- (Ardem's theorem)

Let  $p$  and  $Q$  be the two regular expressions over the input set  $\Sigma$ . The regular expression  $R$  is given as

$$R = Q + RP$$

which has unique solution as  $R = QP^*$

proof :- if 'P' does not contain  $\epsilon$  then there exists  $R$  such that  $R = Q + RP \quad \dots \textcircled{1}$

we will replace  $R$  by  $QP^*$  in eq  $\textcircled{1}$

Consider R.H.S of equation  $\textcircled{1}$

$$\begin{aligned}
 &= Q + RP \\
 &= Q + QP^* P \\
 &= Q(\epsilon + PP^*) \\
 &= QP^* \quad \boxed{\epsilon + RR^* = R^*}
 \end{aligned}$$

thus  $R = QP^*$  is proved.

To prove that  $R = QP^*$  is a unique solution, we will now replace L.H.S of equation  $\textcircled{1}$  by  $Q + RP$ , then it becomes,

$$R = Q + RP$$

$$Q + RP = Q + RP$$

But again  $R$  can be replaced with  $Q + RP$

$$= Q + (Q + RP)P$$

$$= Q + QP + QP^2$$

again replace  $R$  by  $Q + RP$

$$= Q + QP + (Q + RP)P^2$$

$$= Q + QP + QP^2 + RP^3$$

thus if we go on replacing  $R$  by  $Q + RP$  then we get,

$$Q + RP = Q + QP + QP^2 - QP^i + RP^{i+1}$$

$$= Q(\epsilon + P + P^2 + \dots P^i) + RP^{i+1}$$

From equation  $\textcircled{1}$

$$R = Q + RP \quad \text{and also}$$

$$R = Q + RP$$

$$R = Q(\epsilon + P + P^2 + \dots + P^i) + RP^{i+1} \quad \text{--- (2)}$$

where  $i \geq 0$ , consider eq (2)

$$R = Q(\underbrace{\epsilon + P + P^2 + \dots + P^i}_{P^*}) + RP^{i+1}$$

$$R = QP^* + RP^{i+1}$$

let  $w$  be a string of length  $i$ , In  $RP^{i+1}$  has no string of less than  $i+1$  length, hence  $w$  is not in  $RP^{i+1}$

hence  $R$  and  $QP^*$  represent the same set, hence it is proved that

$$R = Q + RP \text{ has unique solution } R = QP^*.$$

Equivalence of two regular expressions - one important theorem named arden's theorem which helps in checking the equivalence of two regular expressions.

Ex :-

$$\text{prove } (1+00^*1) + (1+00^*1)(0+10^*1)^* (0+10^*1) = 0^*1(0+10^*1)^*$$

Sol :- let us solve L.H.S first

$$(1+00^*1) + (1+00^*1)(0+10^*1)^* (0+10^*1)$$

we will take  $(1+00^*1)$  as a common factor

$$= (1+00^*1) [\epsilon + (0+10^*1)^* (0+10^*1)]$$

$$(\epsilon + RR^*) \text{ where, } R = (0+10^*1)^*$$

$$\text{As we know } \epsilon + RR^* = \epsilon + R^*R = R^*$$

$(1+00^*)((0+10^*)^*)$  out of this consider,

$$\Rightarrow 1(\varepsilon + 00^*)(0+10^*)^*$$

Apply  $\varepsilon + 00^* = 0^*$

$$= 10^*(0+10^*)^*$$

$$= \text{R.H.S}$$

Hence two regular expressions are equivalent.

2) Simplify the following regular expression

$$r = \varepsilon + a^*(abb)^* (a^*(abb)^*)^*$$

Sol:- The given regular expression is in the following format

$$r = \varepsilon + rr^*$$

considering  $r = a^*(abb)^*$

This format is simplified according to identity rule,

$$\varepsilon + rr^* = \varepsilon + r^*r = r^*$$

$$r = (a^*(abb)^*)^*$$

again this can be simplified using identity rule,

$$r = (r_1^* r_2^*)^* \quad (\because r_1 = a, r_2 = abb)$$

using  $(p+q)^* = (p^*q^*)^*$

$$= (P+Q)^*$$

$$\Rightarrow (P+Q)^*$$

thus  $r = (a+abb)^*$

Equivalence between finite automata and regular expressions:-

Equivalence of DFA and regular expressions :-

If  $L$  is accepted by a DFA, then  $L$  can be expressed by a regular language. In other words, for any given DFA, we can obtain a RE and vice-versa.

Following are the two methods of constructing R.E from a given DFA.

- By solving equations
- By using transition diagrams

Construction of RE from DFA by solving equations :- following are the steps required to construct a R.E.

Step 1 :- for each of the states  $q_1, q_2, \dots, q_n$  in DFA, write down the equations by considering all edges, that enter into that state.

Step 2 :- for the initial state of DFA, the equation is added with  $\epsilon$ .

Step 3 :- Compute the equations for each state.

Step 4 :- substitute the results of each state equation into the final state equation of DFA, to get a R.E for DFA.

Example :- construct R.E from the DFA



So:- writing the equations for each state, by considering all edges that enter into that state:

$$A = \epsilon + A_0$$

$$B = A_1 + B_1$$

$$C = B_0 + C(C_0 + I)$$

Since,

$$R = Q + RP \Rightarrow R = QP^*$$

Consider,

$$A = \epsilon + A_0$$

$$R = Q + RP$$

$$R = A \quad Q = \epsilon \quad P = O$$

$$R = QP^* \Rightarrow \epsilon \cdot O^*$$

$$R = O^*$$

Now substitute A in B

$$B = A_1 + B_1$$

$$B = O^* I + B_1$$

$$R = Q + RP \Rightarrow QP^* \Rightarrow O^*$$

$$B = O^* I + O^*$$

since A and B are final states in the above DFA, the addition of A and B gives the required result.

$$A + B = O^* + O^* I + O^*$$

$$= O^*(\epsilon + I + I)$$

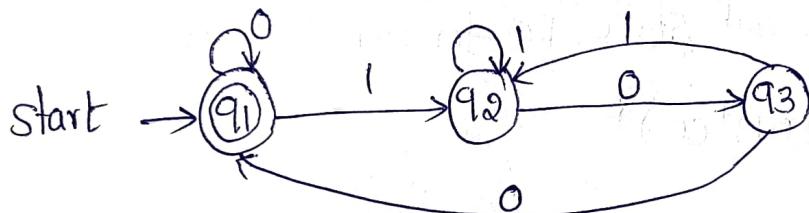
$$= O^* I^* \quad [\epsilon + RR^* = R^*]$$

$$R \cdot E = O^* I^*$$

→ Construct a RE corresponding to the DFA represented in below in below table.  $q_1$  is both the initial and final state.

$\Sigma$	Present IP.	
$q$	0	1
$q_1$	$q_1$	$q_2$
$q_2$	$q_3$	$q_2$
$q_3$	$q_1$	$q_2$

Sol:- the transition diagram is



Now the equations for each state can be written, by considering all edges that enter into that state,

Thus  $q_1 = q_1 0 + q_3 0 + \epsilon$  ( $\epsilon$ , only for initial state)

$$q_2 = q_2 1 + q_3 1 + q_1 1$$

$$q_3 = q_2 0$$

Substituting  $q_3$  in  $q_2$

$$\Rightarrow q_2 = q_1 1 + q_2 1 + q_3 1$$

$$= q_1 1 + q_2 1 + q_2 0 1$$

$$= q_1 1 + q_2 (1 + 0 1)$$

$$= q_1 1 (1 + 0 1)^*$$

$(R = Q + RP \text{ is } R = QP^*)$

Now,

$$\begin{aligned}
 q_1 &= q_1 0 + q_3 0 + \epsilon \\
 &= \epsilon + q_1 0 + q_2 0 . 0 \\
 &= \epsilon + q_1 (0 + 1(1+01)^* . 00) \\
 &= \epsilon + q_1 (0 + 1(1+01)^* . 00)
 \end{aligned}$$

$$[R = Q + RP \Rightarrow R = QP^*] \quad Q = \epsilon$$

$$q_1 = \epsilon \cdot (0 + 1(1+01)^* . 00)^* \quad P = (0 + 1(1+01)^* . 00)$$

since  $q_1$  is the final state, RE for DFA will be

$$RE = (0 + 1(1+01)^* . 00)^*$$

## 2. Construction of RE from DFA by using Transition diagram :-

The following are some transition diagrams and their regular exp's.

→ if a transition diagram has no final state, then  $RE = \emptyset$   
i.e FA does not accept any string.



$$RE = \emptyset$$

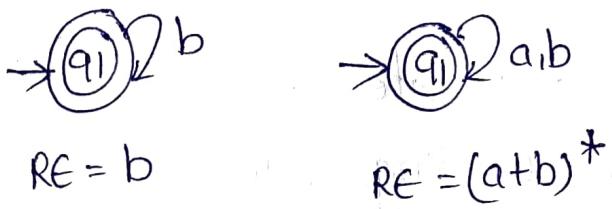
→ if the TD has the same initial and final states, then  $RE = \epsilon$  i.e FA accepts one of the string as  $\epsilon$



$$RE = \epsilon$$

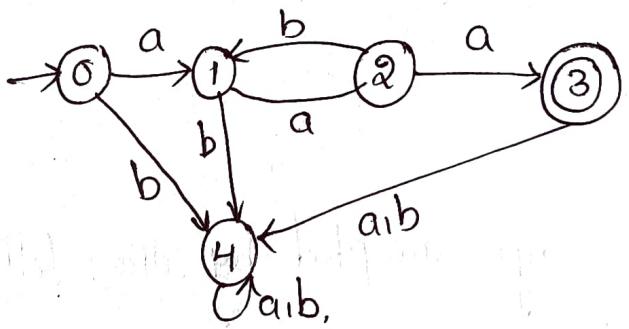
→ for a transition diagram, with transition to itself

$$RE = b^* \text{ (or) } RE = (a+b)^*$$



Example :-

Construct the R.E from the given DFA is.



Sol:-

$$RE = a(ab)^* aa$$

first from the start state it goes to state 1 by reading one 'a'. then from state 1, it goes through the cycle 1-2-1 any no. of times by reading the substring ab and come back to state 1. this is represented by  $(ab)^*$ . then from state 1, it goes to state 2 and then to state 3 by reading aa. thus

$$RE = a(ab)^* aa$$

language acceptance of FA from transition diagram :-

if there is a path in the transition diagram such that, it

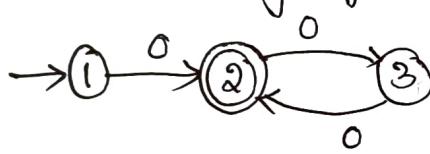
- begins at start state
- ends at an accepting state.
- has a sequence of labels  $w_1, w_2 \dots w_n$

In general, the language accepted by FA, first find the RE from the given FA say 'M'.

then  $L(M) = RE$ .

Example :-

what is the language accepted by the following DFA's ?



Sol:-  $RE = 0 \cdot (00)^*$

$$\Rightarrow L(M) = 0(00)^*$$

i.e machine which accepts all zero's.

Construction of DFA for a given RE :-

Example:- Construct a DFA for the RE given below.

$$RE = \emptyset$$

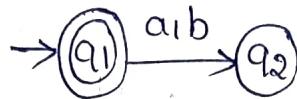
Sol:-



This DFA accepts nothing, since it has no final state.

$$\rightarrow RE = \epsilon \text{ over } \Sigma = \{a, b\}$$

Sol :-



DFA for accepting  $\epsilon$  (since initial and final states are same)

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

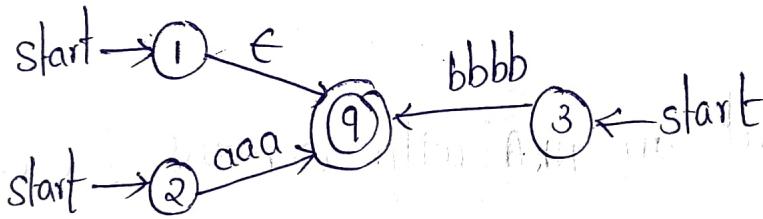
Sol :-



DFA that accepts  $\epsilon$  or a's or b's.

$$\rightarrow RE = \epsilon + aaa + bbbb \text{ over } \Sigma = \{a, b\} \text{ accept only words } \epsilon, aaa, bbbb$$

Sol :-



equivalence of NFA and regular Expressions :-

The rules for conversion of regular expressions into finite automata.

Reg. Exp

meaning

finite automata.

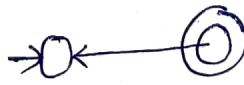
1)  $\epsilon$

null



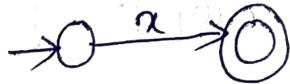
2)  $\emptyset$

-



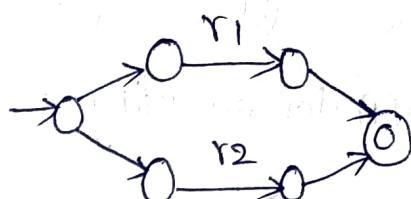
3)  $x$

$x \in \Sigma$



4)  $r_1 + r_2$

parallel  
alternate path



5)  $r_1 \cdot r_2$

series (or)  
continuous path



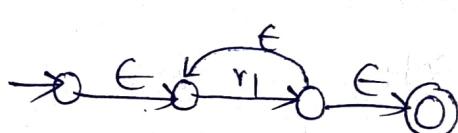
6)  $r_1^*$

Kleene closure



7)  $r_1^+$

positive closure



Example:-

Construct an NFA with  $\epsilon$ -moves for  $00^* + 1$

Sol:-

$$\text{Let } r = 00^* + 1$$

$$r = r_1 + r_2$$

$$\Rightarrow r_1 = 00^* \quad r_2 = 1$$

$\uparrow \quad \uparrow$   
 $r_3 \quad r_4$

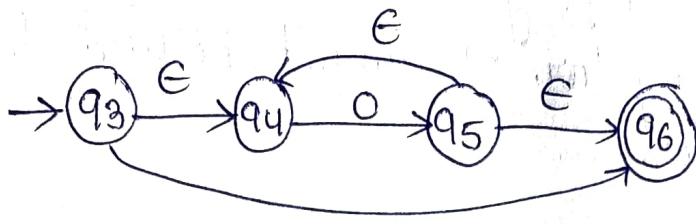
$$r_3 = 0 \quad r_4 = 0^*$$

Construct a machine for  $r_3$ :  $r_3 = 0$

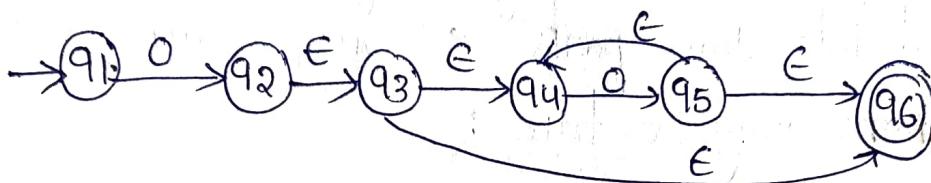
$\rightarrow$  To accept 0 :



$$ru = 0^*$$



NOW construct machine for  $r_1 = r_3 \cdot r_4$

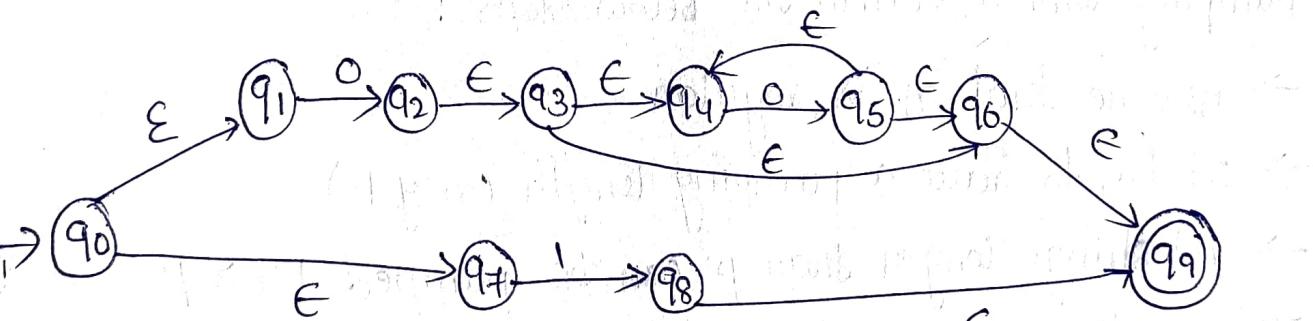


$$r_2 = 1$$



$$\text{NOW } r = r_1 + r_2$$

$$r = 00^* + 1$$



This is the final NFA for the given regular expression.

Pumping Lemma :- pumping lemma is used to prove that a language is not regular. It can not be used to prove that a language is regular.

Statement :- If  $A$  is a regular language, then  $A$  has a pumping length ' $p$ ' such that any string ' $s$ ', where  $|s| \geq p$  may be divided into three parts  $s = xyz$  such that the following conditions may be true :

- i)  $xy^i z \in A$  for every  $i \geq 0$
- ii)  $|y| > 0$
- iii).  $|xy| \leq p$

Proof:- To prove that a language is not regular using Pumping lemma, follow the below steps :

- Assume that  $A$  is regular
- It has to have a pumping length (say  $p$ )
- All strings longer than  $p$  can be pumped  $|s| \geq p$
- Now find a string ' $s$ ' in  $A$  such that  $|s| \geq p$
- Devide  $s$  into  $xyz$
- Show that  $xy^i z \notin A$  for some  $i$
- Then consider all ways that can be devide into  $xyz$
- Show that none of these can satisfy all 3 pumping conditions at the same time,
- $s$  cannot be pumped = contradiction.

Example :- Using pumping lemma prove that the following language  $A = \{a^n b^n \mid n \geq 0\}$  is not regular.

Proof :- Assume that A is regular

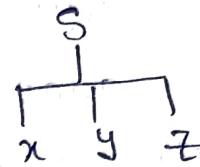
pumping length = P

$$s = a^P b^P$$

divide s into  $xyz$

$$\text{then } p=7, s = a^7 b^7$$

$$s = aaaaaaaabbbaaaaa$$



Case 1 :- the y is in the a part

$$\underline{aaaaaaaa} \underline{bbbbb} \underline{b}$$
  
 $x \quad y \quad z$

$$\text{Now } xy^i z \Rightarrow xy^2 z \text{ (where } i=2) \\ = xyyz$$

$$= aa\underline{aaaaa}a\underline{ag}abbbbaaaa$$

Here no. of a's are not equal to no. of b's

$$= 11 \neq 7$$

Case 2 :- the y is in the b part

$$\underline{aaaaaaaa} \underline{bb} \underline{bbbb} \underline{b}$$
  
 $x \quad y \quad z$

$$\text{Now } xy^i z = xy^2 z \text{ (i=2)}$$

$$\underline{aaaaaaaa} \underline{bb} \underline{bbbb} \underline{bb} \underline{bb}$$

no. of a's  $\neq$  no. of b's

$$7 \neq 11$$

case 3:- the y is in the a part and b part

aaaaaaa a bbbbbbb  
x y z

NOW  $xy^iz = xy^2z$

aaaaaaaabbaaabb bbbbbb

Here also no. of a's not equal to no. of b's. hence the rule  $a^n b^n$ .

Condition 2 :-  $|y| > 0$

$|y| > 0$ , is true

Condition 3 :-  $|xy| \leq p$  where  $p \leq 7$

$|3| \leq 7$  false

$|6| \leq 7$  true

$|9| \leq 7$  false hence false

So assume A is regular, it is contradiction

hence the language A is not regular from pumping lemma.

$xy^iz \notin A$

Show that the language  $L = \{a^i a b^{2i} \mid i > 0\}$  is not regular.

Sol :- The set of strings accepted by language L is

$$L = \{abb, aabbbb, aaabbbb bbbb \dots\}$$

Applying pumping lemma for any of the strings above

take the string abb

it is of the form  $xyz$ , where  $|xy| \leq i, |y| \geq 1$

$$\begin{array}{c} \text{ab} \\ \text{xy} \\ \text{z} \\ \text{x=a y=b z=b} \end{array}$$

To find  $i$  such that  $u^i x y^i z \notin L$  (where  $i=2$ )

$$\begin{array}{lcl} u^i = & a(bb)b \\ xy^i z = & abbb \end{array}$$

hence  $u^i x y^i z = abbb \notin L$

since abbb is not present in the strings of language.

$\therefore L$  is not regular.

### Applications of Regular expressions

- 1) Validation :- checking the correctness of inputs i.e. whether the given string compiles with a set of formatted constraints/not.
- 2) Search and selection :- Identification of a subset of items from a larger set, on the basis of pattern match i.e regular expression.
- 3) Tokenisation :- Conversion of a string of characters into a sequence of words, for later interpretation i.e generation of program called lexical analyzer.
- 4) Regular expression offers a declarative way, to express set of strings.
- 5) Regular expressions are used to define languages. They define the languages accepted by finite automata exactly.

## Closure properties of regular sets / languages :-

Theorem :- Regular languages are closed under union, concatenation, and kleene closure.

Proof :- method 1 :-

Given two regular languages -  $L_1$  and  $L_2$

If  $R_1$  and  $R_2$  are two regular expressions, then

$\rightarrow R_1 + R_2$ , denotes the union of two regular languages

$L_1 \cup L_2$ ,  $L_1 \cup L_2$  is also regular.

$\rightarrow R_1 \cdot R_2$  denotes the union concatenation of two regular languages  $L_1$  and  $L_2$ ,  $L_1 \cdot L_2$  is also regular.

$\rightarrow R_1^*$  denotes the kleene closure of  $L_1 = L_1^*$  is regular.

If  $R_1$  and  $R_2$  are two regular expressions:

with  $R_1 = \{ab, c\}$ ,  $R_2 = \{d, ef\}$ .

$$\begin{aligned} \text{Then } R_1 \cup R_2 &= \{ab, c\} \cup \{d, ef\} \\ &= \{ab, c, d, ef\} \end{aligned}$$

$$R_1 \cdot R_2 = \{ab, c\} \cdot \{d, ef\}$$

$$= \{abd, abef, cd, cef\}$$

$$\text{If } R_1 = \{ab, c\}$$

$$\text{then } R_1^* = \{\epsilon, c, ab, abc, abab, cab, \dots\}$$

→ if  $r = a + b$ , then the language of  $r$  is

method 1:-

$$L(r) = L(a + b)$$

$$= L(a) \cup L(b)$$

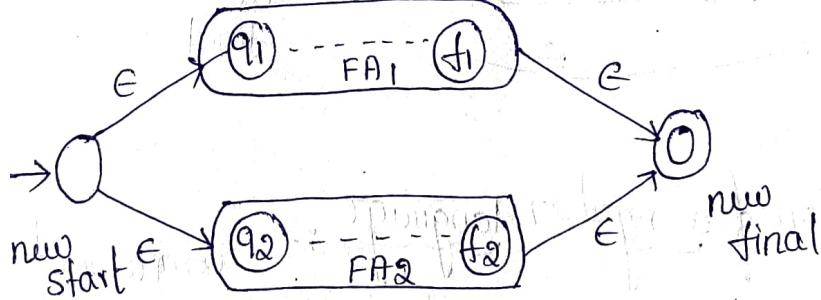
$$= \{a\} \cup \{b\}$$

$L(r) = \{a, b\}$  is also a regular.

method 2:- Assume that  $FA_1$  and  $FA_2$  are two finite automata's accepting languages  $L_1$  and  $L_2$  defined by regular expressions  $r_1$  and  $r_2$  respectively, as.

$$FA_1 = (Q_1, \Sigma_1, \delta_1, q_1, f_1), FA_2 = (Q_2, \Sigma_2, \delta_2, q_2, f_2)$$

Construction of  $L_1 + L_2$ :



→ if  $r = ab^*a$ , then the language of  $r$  is

method 1:-  $L(r) = L(ab^*a)$

$$= L(a) \cdot L(b^*) \cdot L(a)$$

$$= \{a\} \cdot \{b(b)\}^* \cdot \{a\}$$

$$= \{a\} \cdot \{\epsilon, b, bb, bbb, \dots\} \cdot \{a\}$$

$L(r) = \{aa, aba, abba, \dots\}$  is also regular.

method 2 :- construction of  $L_1 L_2$



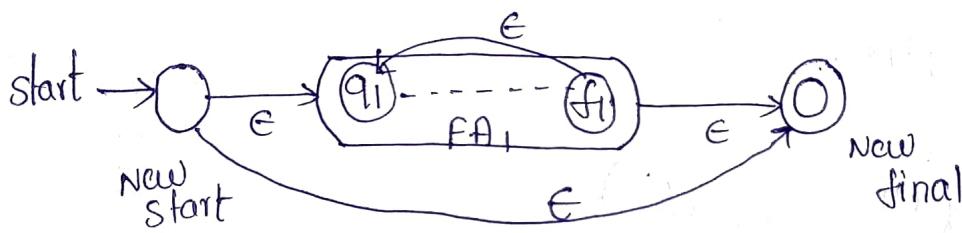
→ if  $r = a^*$ , then the language of  $r$  is

method 1:-  $L(r) = L(a^*)$

$$= (L(a))^*$$

$L(r) = \{ \epsilon, a, aa, \dots \}$  is also regular.

method 2:- construction of  $L_1^*$



complementation of a regular language :-

language  $L$  is a language, such that,

$$\bar{L} = \Sigma^* - L$$

e.g.:  $\Sigma = \{a, b\}$   $L = \{a, b, aa\}$  then  $\bar{L}$  as follows

$$\begin{aligned}\bar{L} &= \{\epsilon, a, b, aa, bb, aaa, bbb, ba, ba b, abbab, \dots\} - \{a, b, aa\} \\ &= \{\epsilon, bb, aaa, bbb, \dots\}\end{aligned}$$

Intersection :- Intersection of two regular languages  $L_1$  and  $L_2$ , is defined as  $L_1 \cap L_2$ , is also regular.

Theorem :- Regular languages are closed under complementation.

Proof :- If  $L$  is a regular language over the alphabet  $\Sigma$ , then there is a DFA,  $D = \{Q, \Sigma, \delta, q_0, F\}$  accepting  $L = L(D)$ .

To accept  $\bar{L} = \Sigma - L$ , Complement the final states of  $D$ . for this consider another DFA,

$$D' = (Q, \Sigma, \delta, q_0, F')$$

where  $F' = Q - F$  i.e  $D$  and  $D'$  differ only in final states. Now, if  $x$  is a input symbol, then  $x \in L(D)'$

$$\Rightarrow \delta(q_0, x) \in F'$$

$$\Rightarrow \delta(q_0, x) \in Q - F$$

$$\Rightarrow \delta(q_0, x) \notin F$$

$$\Rightarrow x \notin L$$

$$\Rightarrow x \in L'$$

$$\Rightarrow x \in \Sigma^* - L.$$

$$\therefore L(D') = \bar{L} = \Sigma^* - L.$$

Hence  $\bar{L}$  is a regular language since it is accepted by DFA.

Intersection of two regular languages :- Intersection of two regular languages  $L_1$  and  $L_2$  is defined as  $L_1 \cap L_2$ .

Theorem :- Regular languages are closed under intersection.

Proof :- Given two regular languages  $L_1$  and  $L_2$

applying demorgan's law,

$$L_1 \cap L_2 = \overline{\bar{L}_1 \cup \bar{L}_2}$$

since  $L_1, L_2$  is regular.

$\Rightarrow \bar{L}_1, \bar{L}_2$  is a regular (regular sets closed under complementation)

$\Rightarrow \bar{L}_1 \cup \bar{L}_2$  is a regular (" " " " union)

$\Rightarrow \overline{\bar{L}_1 \cup \bar{L}_2}$  is a regular (" " " " complementation)

$\Rightarrow L_1 \cap L_2$  is regular (by deMorgan's law)

Formal Grammar :- A grammar is a notation for defining a language through finite number of rules.

Two main categories of formal grammar are :

- Generative grammars
- Analytic grammars.

Generative grammars :- These grammars are set of rules for generation of strings in a language. It describes how to write only those strings in the set.

Analytic grammars :- These grammars are set of rules to determine whether a string is member of the language. It describes how to write only those strings to recognize the strings which are members of a set.

Generative grammars :-

Components of generative grammar :- The generative grammar  $G_1$  consists of the following components:

- A finite set  $V$  of non-terminal symbols

Formally, the generative grammar  $G_1$  is quad-tuple,

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

where

- $V$  is finite set of non terminals

- $T$  is a finite set of terminals.

→  $P$  is the finite set of production rules, each of the form,  $(TUV)^* \cdot V(TUV)^* \rightarrow (TUV)^*$

each production rule maps from one string of symbols to another.

→  $S$  is the start symbol,  $S \in V$ .

The language of a generative grammar  $G_1$ , denoted by  $L(G_1)$ .

Classification of context free grammars: In 1956, Noam Chomsky classified the generative grammars into types known as the Chomsky hierarchy. Each type of generative grammar distinguishes the other types in the form of production rules. The following are the types,

- Unrestricted (or) phrase structure grammar
- Context-sensitive grammar.
- Context-free grammar
- Regular grammar.

Unrestricted grammar :- An unrestricted grammar is a formal grammar, on which no restrictions are made on the left and right sides of the grammar's productions.

Formally, unrestricted grammar  $G_1$  is a quad-tuple,

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

where  $V, T, S$  are same as in the case of generative grammar and  $P$  is the finite set of production rules, each of the form,

$$(VUT)^+ \rightarrow (VUT)^*$$

i.e no  $\epsilon$  on the left hand side of any productions.  
However  $\epsilon$  can appear on the right hand side of production.

Ex:- Let  $G_1 = (V, T, P, S)$  be unrestricted grammar whose production rules are,

$$\begin{aligned} P = \{ & S \rightarrow AS | \epsilon \\ & A \rightarrow aA | a \\ & aaaA \rightarrow ba \} \end{aligned}$$

The derivation for  $w = bbbaa$  is as follows,

$$\begin{aligned} S &\Rightarrow AS \quad [S \rightarrow AS] \\ &\Rightarrow AAS \quad [A-S \rightarrow AS] \\ &\Rightarrow AAAS \quad [S \rightarrow \epsilon] \\ &\Rightarrow AAA \quad [A \rightarrow Aa\bar{A}] \\ &\Rightarrow aAA \quad [A \rightarrow aa\bar{A}] \\ &\Rightarrow aa\bar{A}A \quad [A \rightarrow a\bar{A}] \\ &\Rightarrow \underline{aaaA}AA \quad [aaaA \rightarrow ba] \\ &\Rightarrow ba\bar{A}A \quad [A \rightarrow Aa\bar{A}] \\ &\Rightarrow baa\bar{A}A \quad [A \rightarrow a\bar{A}] \\ &\Rightarrow baaa\bar{A}A \quad [aaaA \rightarrow ba] \\ &= bba\bar{A} \quad [A \rightarrow a] \\ &\Rightarrow bbbaa \end{aligned}$$

In chomsky hierarchy, an unrestricted grammar is called type 0 grammar. The language accepted/generated by this grammar is recursively enumerable language and the recognizers are Turing machines.

$\{w_1 w_2 \{a\}^b\}, \{a^n b^n | n \geq 0\}$  and  $\{a^n | n \geq 0\}$  are some of the languages that are generated by unrestricted grammar.

### Context-sensitive grammar :-

A context sensitive grammar is a formal grammar, in which the left hand sides and right side of the productions may be surrounded by a context of terminals and non-terminal symbols.

Formally, a context-sensitive grammar  $G_1$  is a quad-tuple  $G_1 = (V, T, P, S)$

where  $V, T, S$  are same as generative grammar, and  $P$  is the finite set of productions, each of them,

$$(VUT)^+ \rightarrow (VUT)^+$$

i.e no  $\epsilon$  on the left and right-hand sides of any productions hence this grammar is  $\epsilon$ -free.

Ex :- Let  $G_1 = (V, T, P, S)$  be a context sensitive grammar, whose production rules are,

$$P = \{ S \rightarrow aBb, aB \rightarrow bBB, bB \rightarrow aa, B \rightarrow b \}$$

the derivation for  $w = aaabb$  as follows,

$$\begin{aligned} S &\Rightarrow \underline{aBb} \\ &\Rightarrow \underline{bBBb} \\ &\Rightarrow a\underline{aBb} \\ &\Rightarrow ab\underline{BBb} \\ &\Rightarrow aaa\underline{Bb} \\ &\Rightarrow aaabb. \end{aligned}$$

In chomsky hierarchy, context sensitive grammar is called type 1 grammar. The language generated by this grammar is context sensitive language and the recogniser is LBA (linear bounded automata).

The language  $\{a^n b^n c^n | n \geq 1\}$  is context-sensitive lang.

3) Context-free grammar :-

A context free grammar is a grammar in which the left hand side of each production rule consists of only a single non terminal symbol. This restriction does not matter.

formally, a context free grammar  $G_1$  is a quad-tuple

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

where  $V, T, S$  are same as generative grammar and  $P$  is the finite set of production, as each of the form,

$$A \rightarrow (VUT)^*$$

Any no. of terminals & non terminals including  $\epsilon$  on right side of production

Ex:- let  $G_1 = (V_1, T_1, P, S)$  be a context free grammar whose production rules are,

$$P = \{ S \rightarrow aSa, S \rightarrow bSb | a | b | \epsilon \}$$

The derivation for  $w = ababa$  is as follows,

$$\begin{aligned} S &\Rightarrow aSa \\ &\Rightarrow abSba \\ &\Rightarrow ababa \end{aligned}$$

In Chomsky hierarchy, Context free grammar is called type 2 grammar. The language generated by this grammar is context free language and the corresponding recogniser is PDA (push down automata).

Ex:- The language  $\{a^n b^n | n \geq 1\}$  is context free language.

Regular grammar :- A regular grammar is a formal grammar with the restrictions on both left and right hand sides of the productions.

- left hand side of each production rule consists of only a single non-terminal symbol.
- Right side of production rule may be nothing, (or) a single terminal (or) a single terminal symbol followed by non-terminal symbols, but nothing else.

Formally, a regular grammar  $G_1$  is a quad-tuple,

$$G_1 = (V_1, T_1, P, S)$$

where  $V, T, S$  are same as generative grammar and  $P$  is the finite set of productions, each of the form,

$$A \rightarrow t \text{ (or) } A \rightarrow tB \text{ (or) } A \rightarrow \epsilon$$

where  $A, B \in V$  and  $t \in T$ .

Ex :- let  $G_1 = (V, T, P, S)$  be a regular grammar, whose production rules are  $P = \{ S \rightarrow aS | a \}$

The derivation for  $w = aaa$  is as follows

$$S \Rightarrow aS$$

$$\Rightarrow aaS$$

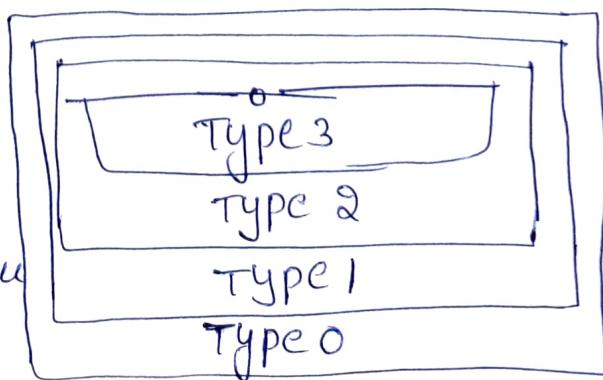
$$\Rightarrow aaa$$

In chomsky hierarchy, regular grammar is called a type 3 grammar. The language generated by this grammar is a regular language and the recogniser is FA (finite automata).

The language  $\{a^n b^m \mid m, n \geq 1\}$  is a regular language, and

Regular grammars are commonly expressed using regular expressions.

Regular grammar  
Context free  
Context sensitive  
unrestricted



Left linear and Right linear grammars :- we discuss

some of the commonly used grammars,

Linear and Nonlinear grammars :-

A grammar, with at most one variable (Non terminal) at the right side of a production, is a linear grammar, otherwise it is non-linear grammar.

Example: (Linear grammars) (Non-linear grammar)

i)  $S \rightarrow asb$   
 $S \rightarrow \epsilon$

ii)  $S \rightarrow Ab$   
 $A \rightarrow aAb$   
 $A \rightarrow \epsilon$

iii)  $S \rightarrow A$   
 $A \rightarrow aB/\epsilon$   
 $B \rightarrow Ab$

i)  $S \rightarrow SS$   
 $S \rightarrow \epsilon$

$S \rightarrow asb$   
 $S \rightarrow bsa$

ii)  $S \rightarrow aA$   
 $A \rightarrow \epsilon$   
 $S \rightarrow asAs$

Linear grammars are further classified into 2 types

1. Right linear grammar

2. Left linear grammar

1. Right linear grammar :-

A grammar  $G_1 = (V, T, P, S)$  is a right linear, if all productions are of the form

$A \rightarrow xB$

$A \rightarrow x$

where  $A, B \in V$  and  $x \in T^*$

Ex :-

$S \rightarrow abS$

$S \rightarrow a$

left linear grammar :- A Grammar  $G_1 = \{V, T, P, S\}$  is left linear, if all productions are of the form,

$A \rightarrow Bx$

$A \rightarrow x$

where  $A, B \in V$  and  $x \in T^*$

Ex :-  $S \rightarrow Aab$

$S \rightarrow Aab | B$

$B \rightarrow a$

3) Regular Grammar :- A Grammar  $G_2 = \{V, T, P, S\}$  is said to be regular, if it is either right linear or left linear.

Ex :-  $S \rightarrow abS$      $S \rightarrow Aab$

$S \rightarrow a$      $A \rightarrow Aab | B$

$B \rightarrow a$

4) non regular grammar :- A grammar  $G_3 = \{V, T, P, S\}$  is non-regular, if it is neither right linear nor left linear.

Ex :-  $S \rightarrow asb$      $S \rightarrow \epsilon$

$S \rightarrow \epsilon$      $S \rightarrow asb$

$S \rightarrow bsa$

5) simple (or) S-grammar :- A grammar  $G_1 = \{V, T, P, S\}$  is said to be a simple grammar (or) S-grammar, if all productions are of the form

$$A \rightarrow a\alpha$$

where  $A \in V$ ,  $a \in T$  and  $\alpha \in V^*$  and any pair  $(A, \alpha)$  occurs at most once in  $P$ .

Ex :-  $S \rightarrow aS | bSS | c$  is simple grammar, since pairs  $(S, a)$ ,  $(S, b)$  occur only once in production.

$S \rightarrow aS | bSS | ass | c$  is not simple grammar, since pair  $(S, a)$  occurs twice in production as  $aS$  &  $ass$ .

6) Recursive grammar :- A Grammar  $G_1 = \{V, T, P, S\}$  is recursive, if it contains either a recursive production (or) an indirectly recursive production.

Ex :-  $S \rightarrow aS \quad S \rightarrow SS$   
 ~~$S \rightarrow SS$~~   $S \rightarrow \epsilon$   
 $S \rightarrow asb$ .

Derivation Tree :- When deriving the string  $w$  from  $s$ , if every derivation is considered to be a step in the tree construction, then the graphical display of the derivation of string  $w$  results in a tree structure.

This is called a derivation tree (or) parse tree (or) generation tree (or) production tree.

## → Relation between regular grammar and finite automata:

17

the following relation exists between regular grammar and finite automata.

- 1) Finite automata from regular grammar :- for a given right linear grammar  $G_1$ , there exists a language  $L(G_1)$ , which is accepted by a finite automata.
- 2) Regular grammar from finite automata :- For a given finite automata  $M$ , there exists a right linear grammar  $G_1$ , such that  $L(M) = L(G_1)$ .
- 3) Regular expression from regular grammar :- for a given regular grammar  $G_1$ , there exists a regular expression that specifies  $L(G_1)$ .

### 1) Finite automata from regular grammar:

Let  $G_1 = (V, T, P, S)$  be a right linear grammar  
let  $V = \{v_0, v_1, \dots\}$  be the variables and productions  $P(G_1)$  be

$$P = \{ v_i \rightarrow a_1 a_2 \dots a_m v_j \mid$$

(or)

$$\} \\ v_i \rightarrow a_1 a_2 \dots a_m$$

we construct the finite automata  $M = (Q, \Sigma, \delta, q_0, F)$   
from these productions, by using the following steps.

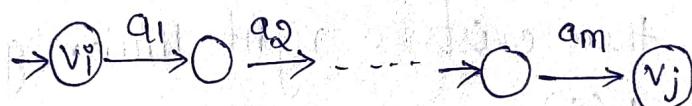
K. Lavanya

step1 :- start symbol of grammar is the start symbol of  $M (q_0 = v_0)$ .

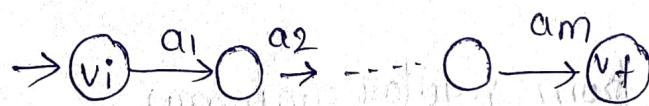
step2 :- each variable  $v_i$  of  $G$ , corresponds to a state in  $M$ .



step3 :- for each production of the form  $v_i \rightarrow a_1 a_2 \dots a_m v_j$ , if the transitions of  $M$  are of the form  $s(v_i, a_1 a_2 \dots a_m) = v_j$ , add the transitions and intermediate states.



step4 :- for each production of the form  $v_i \rightarrow a_1 a_2 \dots a_m$  if the transitions of  $M$  are of the form  $(v_i, a_1 a_2 \dots a_m) = v_f$  (where  $v_f$  is the final states of  $M$ ), add the transitions and intermediate states.



Example :- Construct the finite automata to accept the language generated by the following grammar,

$$S \rightarrow aA|B$$

$$A \rightarrow aAB$$

$$B \rightarrow bB|a$$

so :- let  $G$  be the right linear grammar and  $M$  is the finite automata

step1 :- the start symbol of  $G$  is the start symbol for  $M$ , i.e  $q_0 = \{S\}$

Step 2

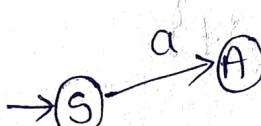
$\Rightarrow$  Every variable  $v = \{S, A, B\}$  of  $G_1$ , corresponds to state  $v$ .

Step 3 :-

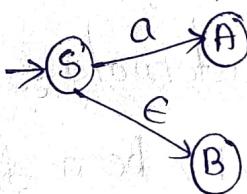
Compute the transitions of  $M$  for all production of the form

$$V_i \rightarrow a_1 a_2 \dots a_m V_j.$$

a)  $S \rightarrow aA$

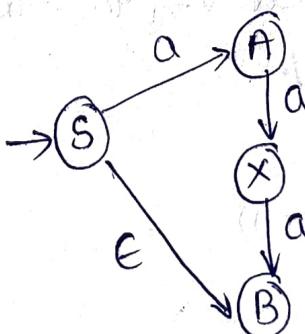


b)  $S \rightarrow B$

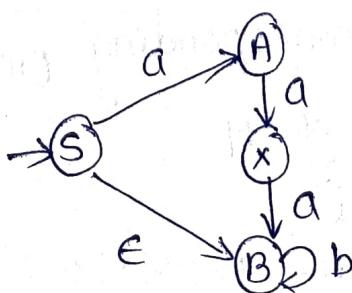


c)  $A \rightarrow a a B$

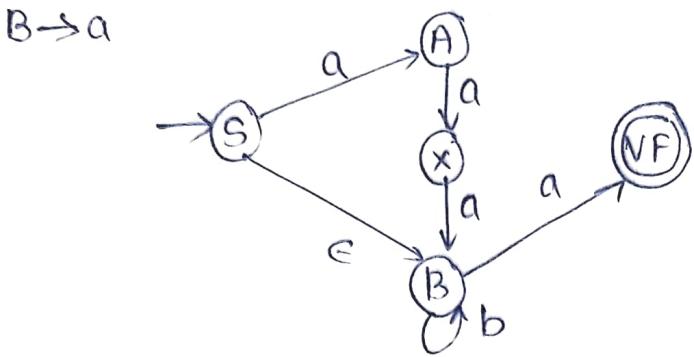
add transitions and intermediate states ( $x$ )



d)  $B \rightarrow b B$



Step 4:- Compute for all productions of the form  $v_i \rightarrow a_1 a_2 \dots a_m$



Final N with  $L(N) = L(G_1) = aaab^*a + b^*a$

thus  $N = (Q, \Sigma, \delta, q_0, F)$

where  $Q = \{S, A, X, B, NF\}$

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

$q_0 = S$ ,  $F = NF$ .

## 2) Regular grammar from finite automata :-

Let  $N = (Q, \Sigma, \delta, q_0, F)$  be a finite automata

where  $Q = \{q_0, q_1, \dots, q_n\}$  is the finite number of states and  $\Sigma = \{q_1, q_2, \dots, q_m\}$  is the set of input symbols. Construct a regular grammar  $G = (V, T, P, S)$ , where  $V = \{q_0, \dots, q_n\}$ ,  $T = \Sigma$  and  $S = q_0$ . by using the following steps.

Step 1:- for any transition of the form

$$(q_i) \xrightarrow{a} (q_j)$$

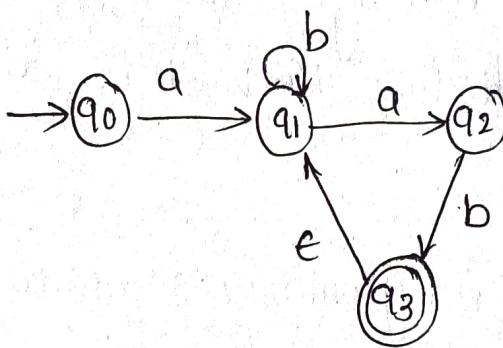
i.e  $\delta(q_i, a) = q_j$ , the corresponding production is  $q_i \rightarrow a q_j$

Step 2:- for any final state of N

$$(q_f)$$

the production added is  $q_f \rightarrow \epsilon$

Example :- Construct a regular grammar from the following finite automata.



∴ for each transition of FA, the corresponding productions are

$q_i \xrightarrow{a} q_j$  productions

$$q_0 \xrightarrow{a} q_1 \quad q_0 \rightarrow aq_1$$

$$q_1 \xrightarrow{b} q_1 \quad q_1 \rightarrow bq_1$$

$$q_1 \xrightarrow{a} q_2 \quad q_1 \rightarrow aq_2$$

$$q_2 \xrightarrow{b} q_3 \quad q_2 \rightarrow bq_3$$

$$q_3 \xrightarrow{e} q_1 \quad q_3 \rightarrow q_1$$

$$q_3 \Rightarrow q_f \quad q_3 \rightarrow \epsilon$$

thus  $G_1 = (V, T, P, S)$  is the required grammar, where  
 $V = \{q_0, q_1, q_2, q_3\}$ ,  $T = \{a, b\}$ ,  $S = \{q_0\}$ ,  $F = \{q_3\}$  and

$$P = \{$$

$$q_0 \rightarrow aq_1$$

$$q_1 \rightarrow bq_1$$

$$q_1 \rightarrow aq_2$$

$$q_2 \rightarrow bq_3$$

$$q_3 \rightarrow q_1$$

$$q_3 \rightarrow \epsilon$$

}

3) Regular Expressions from regular grammar :-

Let  $G_1 = (V, T, P, S)$  be a right linear grammar.

A regular expression  $R$ , that satisfies  $L(G_1)$ , can be directly obtained as follows,

- Replace the ' $\rightarrow$ ' symbols in the grammar's productions with '=' symbols to get a set of equations.
- Convert the equations of the form  $A = aAlb$ , as  $A = a^k b$ .
- solve the set of equations obtained, to obtain the value of variable  $s$ . The result is the regular expression,  $L(G_1)$ .

Example:- obtain the regular expression for the grammar given below.

$$S \rightarrow 0IB|0$$

$$B \rightarrow IB|II$$

Sol:-

Replace ' $\rightarrow$ ' by '=' in the above productions.

so that,

$$S = 0IB|0 \quad -①$$

$$B = IB|II \quad -②$$

Now from  $B = IB|II$

$$\Rightarrow B = I^*II \quad [A \neq aAlb \text{ is } A = a^k b]$$

substituting for  $B$  in equation ① we get,

$$S = 0IB|0$$

$$\Rightarrow S = 0I^*II|0$$

$$(\text{or}) \quad S = (0I^*II + 0) \quad \text{this is the regular expression}$$

## Conversion of LLG to RLG :-

Procedure :-

- 1) Convert the given left linear grammar to FA
- 2) Interchange the initial and final states of FA
- 3) Reverse the directions of all the transitions.
- 4) Construct the regular grammar from this FA.

Example :- Convert the given left linear grammar to right linear grammar.

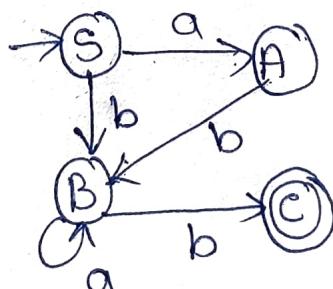
$$S \rightarrow AaB.b$$

$$A \rightarrow B.b$$

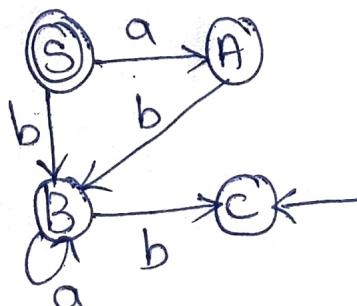
$$B \rightarrow BaB.b$$

Sol:- Let  $G_1$  be the left linear grammar.

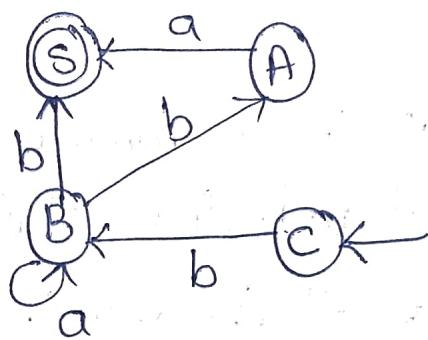
Step 1 :- Convert the given left linear grammar to FA.



Step 2 :- Interchange the initial and final states of FA



Step 3 :- Reverse the directions of all the transitions



Step 4 :- Construct regular grammar from this FA

$$\begin{array}{l} C \rightarrow bB \\ B \rightarrow aB \mid ba \mid b \\ A \rightarrow a \\ S \rightarrow \epsilon \end{array} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{required RLG}$$

Homework problem :-

$$\begin{array}{l} S \rightarrow ca \mid Aa \mid Bb \\ A \rightarrow Ab \mid ca \mid Bb \mid a \\ B \rightarrow Bb \mid b \\ C \rightarrow Aa \end{array}$$

K. Jeevunay o.