

# Regular Expression

A regular expression is a concise way to describe a regular language. These are symbolic notations used to define search patterns in strings.

For a given character set  $\Sigma$  of a language  $L$ , the regular expression is defined by the following rules:

Rule 1: Every character belonging to  $\Sigma$  is a regular expression.  $r = a$ .

Rule 2: Null string  $\epsilon$  is a regular expression.  $r = \epsilon$ .

Rule 3: Empty language is a regular expression.  $r = \emptyset$ .

Rule 4: If  $R_1$  &  $R_2$  are two regular expressions, then  $R_1 R_2$  is also a regular expression (concatenation of two regular expressions is a RE).

Rule 5: If  $R_1$  &  $R_2$  are two RE, then  $R_1 + R_2$  is a RE (Union of two RE is a RE).

Rule 6: If  $R$  is a RE, then  $(R)$  is a RE.

Rule 7: If  $R$  is a RE, then  $R^*$ ,  $R^+$ ,  $R^2$  is also a RE. (Power of a RE is a RE).

Rule 8: Any combination of the preceding rules is also a RE.

i.e. there are three types of basic (atomic) RE.  
 $'a'$ ,  $'\epsilon'$ ,  $'\emptyset'$

& there are 4 types of operations on RE.  
 $'+'$ ,  $'.'$ , Power ( $^*$ ,  $^+$ , exact value), Parenthesis  
Union, Concatenation

High to low Precedence of Operators:

( ), Power, . , +

Highest

Lowest

Regular Set: Set of strings generated by regular expression is known as Regular Set. It is the language corresponding to a regular expression & is also known as regular language.

Some Regular Language corresponding to R.E.

R.E.

a

a + b + c

ab + ba

ab(d + e)

(aa)\*

(a+b)(c+d+e)u

abc\*de

(a+b)c\*

(abc)\*

(0+1)2(33+1)

a\*a

(abc)\*k

∅

(a+b+E)

E

(abc)k\*

ab(c+d)e\*

Regular Set

{a}

{a, b, c}

{ab, ba}

{abd, abe}

{ε, aa, aaa, aaaa, ...}

{acu + adu + aeu + bu + bua, ...}

{abde, abcde, abccde, ...}

{a, b, ac, bc, acc, ...}

{ε, abc, abcabc, ...}

{021, 0233, 11233, 1121}

{a, aa, aaa, ...}

{k, abck, abcabck, ...}

{}

{ε, a, b}

{E}

{abc, abck, abckk, ...}

{abc, abd, abce, abde, ...}

Regular expression corresponding to regular lang

Language

R.E.

$L = \{aaa, aaaaa, \dots\}$  set of odd length of a, min length 3

$(aa)^+a$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{0, 1\}, w \text{ begins with } 00 \text{ & ends with } 11\}$

$00(0+1)^*11$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{0, 1\}, \text{ such that } w \text{ contains alternates } 0s \text{ & } 1s\}$

$0(10)^* + 0(10)^*1 + 1(01)^* + 1(01)^*0$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{0, 1\}, w \text{ ends with } 11\}$

$(0+1)^*11$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{a, b, c\}, 'ccc' \text{ is a substring of } w\}$

$(a+b+c)^*ccc(a+b+c)^*$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{a, b\}, w \text{ contains exactly}$

2 a's

$b^*a^*b^*a^*b^*$

atleast 2 a's

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

atmost 2 a's

$b^*(a+t)^*b^*(a+t)^*b^*$

$L = \{0^i 1^j 2^k, \text{ where } i, j, k \in \mathbb{N}\}$

$0^+1^+2^+$

$L = \{0^i 1^j 2^k, \text{ where } i, j, k \geq 0\}$

$0^*1^*2^*$

$L = \{0^i 1^j 2^k, \text{ where } i, j \geq 1, k \geq 0\}$

$0^+1^+2^*$

$L = \{w \mid w \in \Sigma^*, \Sigma = \{0, 1\}, w \text{ contains no. of } 1's$   
divisible by 3}

$0^*(10^*10^*10^*)^*$

## Identities related to Regular Expression :

Two regular expressions  $R_1$  &  $R_2$  are equivalent if their corresponding languages are same, i.e.  $L(R_1) = L(R_2)$ . following are some such Regular expression which are equivalent:

- 1.)  $\phi + R = R + \phi = R$
- 2.)  $\epsilon + R = R + \epsilon$
- 3.)  $\epsilon R = R\epsilon = R$
- 4.)  $R + R = R$
- 5.)  $\epsilon^* = \epsilon$
- 6.)  $RR^* = R^*R = R^+$
- 7.)  $\epsilon + RR^* = \epsilon + R^*R = R^*$
- 8.)  $(P + Q)R = PR + QR$
- 9.)  $(PQ)^*P = P(QP)^*$
- 10.)  $(R^*)^* = R^*$
- 11.)  $(P + Q)^* = (P^* + Q^*)^* = (P^*Q^*)^*$
- 12.)  $\phi R = R\phi = \phi$

- Q. Which of the following statements is correct for  $R_1$  &  $R_2$ .  $R_1 = ((aa)^* + a^*)^*$ ,  $R_2 = (a^*)^+$
- a)  $L(R_1) \subset L(R_2)$
  - b)  $L(R_1) \supset L(R_2)$
  - c)  $L(R_1) \neq L(R_2)$
  - d)  $L(R_1) = L(R_2)$

Solution : d) since  $R_1$  &  $R_2$  both are  $a^*$ .

- Q. Identify which of the following R.E. can generate all the strings which do not contain the substring 100.

a)  $0^*(0+1)^*$   
c)  $0^*1^*01$

b)  $0^*1010^*$   
d)  $0^*(10+1)^*$

Solution: d) although c) also does not create strings containing 100, but it also cannot create  $\epsilon, 0, 1$  strings which also does not contain 100.

Q. Identify which of the R.E. cannot generate a string 1101.

a)  $110^*(0+1)^*$   
c)  $(10)^*(01)^*(00+11)^*$

b)  $1(0+1)^*101$

d)  $(00+(11)^*01^*)^*$

Solution: c)

Q. Identify which of the following R.E. correctly represent  $L = \{w \in \{0+1\}^* \mid w \text{ has even no. of 1's}\}$ .

a)  $(0^*10^*1)^*$   
c)  $0^*(10^*1)^*0^*$

b)  $0^*(10^*10^*)^*$

d)  $0^*1(10^*1)^*10^*$

Solution: b), a) cannot generate only 0's.

c) can't generate 11011.

d) can't generate  $\epsilon, 0, 00 \dots$ .

Q. Identify which of the following R.E. represent a language which contains atleast one a and one b.

a)  $(a+b)^+$   
c)  $(a+b)^*b(a+b)^*a(a+b)^*$

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

d)  $(a+b)^*a(a+b)^*b(a+b)^* + (a+b)^*b(a+b)^*a(a+b)^*$

Sol. d)

## Regular Expression From Finite Automata

A finite automata contains following types of paths:

Serial path, Parallel path, loop, cycle.

Corresponding to these path the regular expression will be : concatenation for serial path, Union for parallel path.

Kleene closure for loop. Similarly for cycle also it is a Kleene closure.

For converting a regular expression, we need to find path from initial state to final state.

There are two methods to convert FA to Regular Expression :

- 1.) Arden's Theorem.
- 2.) State Elimination.

Arden's Theorem : Let  $P$ ,  $Q$ , and  $R$  be the regular expressions over  $\Sigma$ .

If  $P$  is not  $\epsilon$ , then the equation

$$R = Q + RP$$

has a unique solution  $QP^*$ .

Proof : let's first check  $QP^*$  is actually a sol.

Let's substitute the value of  $QP^*$  in  $R = Q + RP$ .

$$QP^* = \underbrace{Q + (QP^*)P}_{Q(\epsilon + P^*P)}$$

$$Q(\epsilon + P^*P) = Q(P^*) = QP^*$$

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

So yes  $QP^*$  is a sol of  $R = Q + RP$ .

Now, let's check if it is the only sol.

$$R = Q + RP$$

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

$$= Q + QP + (Q + RP)P^2 = Q + QP + QP^2 + RP^3$$

⋮  
⋮

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

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

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

$$\Rightarrow R = Q + RP$$

∴ has two solutions

$$QP^* \text{ and }$$

$$RP^{i+1}$$

Since  $P \neq \epsilon$ , and  $i \geq 0$  hence

$|R| < |RP^{i+1}|$  {Length of  $R$  is less than  $RP^{i+1}$ }  
& that is not possible, so  $RP^{i+1}$  cannot be a solution to  $R$ .

Hence the only feasible solution for  $R$  is  $QP^*$ .

$\therefore R = QP^*$  is a unique sol<sup>n</sup> for  $R = Q + RP$ .

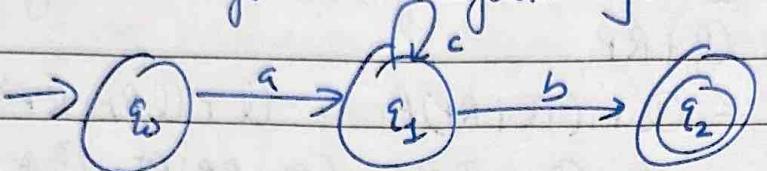
Conditions to apply Arden's theorem:

1) The finite automata should not contain null ( $\epsilon$ ) moves

2) There should be exactly one initial state.

For finding regular expression using arden's theorem,  
first we create equations for each state, i.e.  
find out all the ways to reach at that state.  
for initial state,  $\epsilon$  is always a way to reach there.

Q. Find RE for the following FA automata.



Sol. The transitions required to reach various states in the automata.

$$q_0 = \epsilon \quad -\textcircled{1}$$

$$q_1 = q_0 a + q_1 c \quad -\textcircled{2}$$

$$q_2 = q_1 b \quad -\textcircled{3}$$

Regular expression of the FA will be same as that of final state, so we need to find the value of  $q_2$ .

Putting value of  $\textcircled{1}$  in  $\textcircled{2}$ .

$$q_1 = (\epsilon)a + q_1 c = a + q_1 c$$

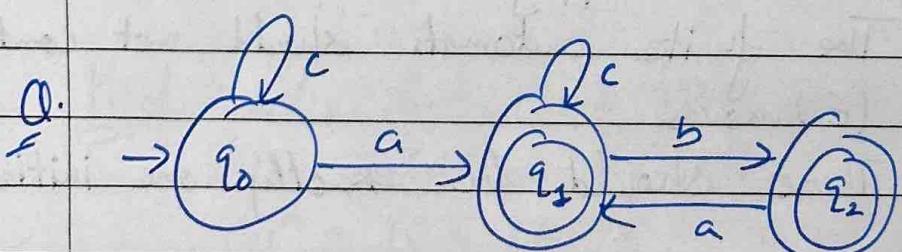
( $R = Q + RP^*$ ) Applying Arden's Theorem

$$R = QP^* \quad q_1 = ac^* \quad -\textcircled{4}$$

Put the value of  $\textcircled{4}$  in  $\textcircled{3}$ .

$$q_2 = (ac^*)b = \underline{ac^*b}$$

↳ R.E. of FA.



Sol.

$$q_0 = \epsilon + q_0 c \quad -\textcircled{1}$$

$$q_1 = q_0 a + q_1 c + q_2 a \quad -\textcircled{2}$$

$$q_2 = q_1 b \quad -\textcircled{3}$$

Applying Arden's on ①

$$q_0 = \epsilon(c^*)^* = c^* \quad -④$$

Put value of  $q_0$  from ④ in ②

$$q_1 = c^*a + q_1c + q_2a \quad -⑤$$

Put value of  $q_2$  from ③ in ⑤

$$q_1 = c^*a + q_1c + (q_1b)a$$

$$q_1 = c^*a + q_1(c+ba) \quad -⑥$$

Applying Arden's theorem on ⑥

$$q_1 = c^*a(c+ba)^* \quad -⑦$$

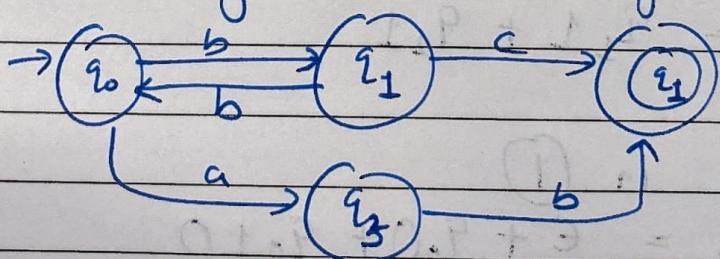
Putting value of ⑦ in ③

$$q_2 = c^*a(c+ba)^*b$$

Since there are two final states, so R.E. for the automata is union of R.E. of both the final states. i.e.

$$\underline{c^*a(c+ba)^* + c^*a(c+ba)^*b}$$

Q. Find the regular expression for



Sol.

$$q_0 = \epsilon + q_1b \quad -①$$

$$q_1 = q_0b \quad -②$$

$$q_2 = q_1c + q_3b \quad -③$$

$$q_3 = q_0a \quad -④$$

Putting the value of ② in ①

$$q_0 = \epsilon + (q_0 b) b$$

$$q_0 = \epsilon + q_0 b b$$

Applying Arden's

$$q_0 = (bb)^* - ③$$

Putting the value of ③ in ② & ④

$$q_1 = (bb)^* b$$

$$q_3 = (bb)^* a$$

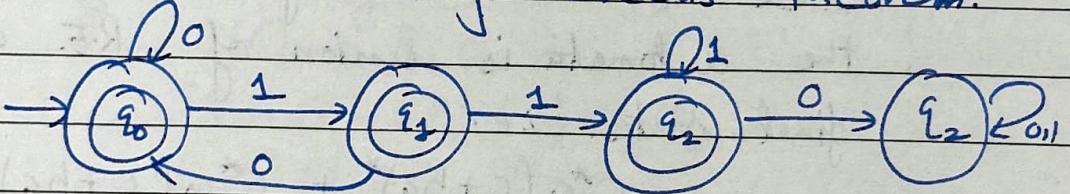
Putting these values in ③.

$$q_2 = (bb)^* bc + (bb)^* ab$$

Since  $q_2$  is the only final state, the R.E. for FA is  $(bb)^* bc + (bb)^* ab$

Q. Create a regular expression for '110' is not a subset using Arden's theorem.

Sol



$$q_0 = \epsilon + q_0 0 + q_1 0 \quad - ①$$

$$q_1 = q_0 1 \quad - ②$$

$$q_2 = q_1 1 + q_2 1 \quad - ③$$

① in ①

$$q_0 = \epsilon + q_0 0 + q_0 1 0$$

$$q_0 = \epsilon + q_0 (0 + 1 0)$$

Applying Arden's

$$q_0 = \epsilon (0 + 1 0)^* = (0 + 1 0)^* \quad - ④$$

Put it in ②

$$q_1 = (0 + 1 0)^* 1 \quad - ⑤$$

Put value of  $q_5$  in  $q_3$   
 $q_2 = (0+10)^* 11 + q_2 L$

Applying Arden's .

$$q_2 = (0+10)^* 11 1^*$$

Since  $q_0$ ,  $q_1$  &  $q_2$  all are final states  
Regular expression of automata is

$$\begin{aligned} & q_0 + q_1 + q_2 \\ & (0+10)^* + (0+10)^* 1 + (0+10)^* 11 1^* \\ & = \underline{(0+10)^* (\epsilon + 1 + 11 1^*)} \end{aligned}$$

## State Elimination Method