

a=1 b=2 c=3 d=4 e=5 f=6 g=7
h=8 i=9 j=10

PAGE NO.: / /

Rabin-Karp Algorithm

m=3

n=6

Text: a a a a a b pattern: a a b
 $\begin{array}{c} \text{3} \\ \text{3} \neq \text{4} \end{array}$ $\begin{array}{c} \text{1} \\ \text{1} + \text{1} + \text{2} = \text{4} \end{array}$

a a a a a b
 $\begin{array}{c} \text{3} \\ \text{3} \neq \text{4} \end{array}$

a a a a a b
 $\begin{array}{c} \text{3} \\ \text{3} \neq \text{4} \end{array}$

a a a a a b
 $\begin{array}{c} \text{4} \\ \text{4} = \text{4} \end{array}$

- Keep on checking for the hash function of each m letters
- This sliding is called as rolling hash func

E.g. 2

~~n=5~~ n=8

a b c d a b c e pattern: b c e
 $\begin{array}{c} \text{1} \\ \text{6} \times \end{array}$ $\begin{array}{c} \text{2} \\ \text{2} + \text{3} + \text{5} = \text{10} \end{array}$

$\begin{array}{c} \text{6} \\ \text{6} \times \end{array}$

$\begin{array}{c} \text{8} \\ \text{8} \times \end{array}$

$\begin{array}{c} \text{7} \\ \text{7} \times \end{array}$

$\begin{array}{c} \text{6} \\ \text{6} \times \end{array}$

$\begin{array}{c} \text{10} \\ \text{10} \checkmark \end{array}$

a=1 b=2 c=3 d=4 e=5 f=6 g=7
h=8 i=9 j=10

Rabin-Karp Algorithm

$m = 3$

$n = 6$

Text: a a a a a b pattern: a ab
 | # #
 | 1 + 1 + 2 = 4
 | 2
 | hash hash
 3 ≠ 4 funcn code

a a a a a b
 |
 3 ≠ 4

a a a a a b
 |
 3 ≠ 4

a a a a a b
 |
 4 = 4

- Keep on checking for the hash or hashfunction of each m letters
- This sliding is called as rolling hash func

E.g. 2

~~$n = 5$~~ $n = 8$

a b c d a b c e
 |
 6 x

pattern: b c e
 |
 2 + 3 + 5 = 10

 |
 9 x

 |
 8 x

 |
 7 x

 |
 6 x

 |
 10 ✓

Drawback:

Spurious hits: Hash codes are matching but pattern is not

$$n=11$$

c c a c c a a e d b a
 $\overbrace{\quad\quad\quad}^{7x}$

$$\overbrace{\quad\quad\quad}^{7x}$$

$$\overbrace{\quad\quad\quad}^{7x}$$

$$\overbrace{\quad\quad\quad}^{7x}$$

$$m=3$$

d b a
 4 2 1 = 7

Hence the hash codes match but the pattern isn't the same, \Rightarrow there are spurious hits.

In case of worst case, the time complexity is $\Theta(mn)$ & in the best case $\Theta(n-m+1)$

So, we can't use such a weak hash function we need to use a strong hash function.

Let's consider the same e.g.

$$n=11$$

Text: c c a c c a a e d b a
 $\overbrace{\quad\quad\quad}^{300}$ $\overbrace{\quad\quad\quad}^{30}$ $\overbrace{\quad\quad\quad}^1$

$$300 + 30 + 1 = 331$$

$$\overbrace{\quad\quad\quad}^1$$

$$313$$

$$m=3$$

pattern: d b a
 4 2 1

$$421$$

$$4 \times 10^2 + 2 \times 10 + 1 \times 10^0 \\ p[1] \times 10^{m-1} + p[2] \times 10^{m-2} + \\ p[3] \times 10^{m-3}$$

$$400 + 20 + 1 = 421$$

$$331$$

By using this algo. we have reduced
the chances of worst-case.

Best case TC : $O(n-m+1)$

Worst case TC : $O(mn)$

Naive String matching

algo:

- It is the simplest method using brute force approach.
- It compares 1st character of the pattern with searchable text. If match is found, pointers in both strings are advanced.
- If match not found, pointer of text is shifted one place to right from previous starting position & pointer of pattern is reset.
- The process is repeated until the end of text.

$$T[s+1 \dots s+m] = P[1 \dots m]$$

$$T[s+j] = P[j] \quad 1 \leq j \leq m \\ 0 \leq s \leq n-m$$

E.g. $T = \text{acaabc}$ $P = \text{aab}$

Iter-1 of

a	c	a	a	b	c
---	---	---	---	---	---

 \checkmark ~~x~~

$s=0$

a	a	b
---	---	---

 $\leftarrow P$

Iter-2 ~~x~~

a	c	a	a	b	c
---	---	---	---	---	---

 $s=1$ ~~x~~

a	a	b
---	---	---

gter-3

$s=2$

a	x	c	a	a	b	c
			a	a	b	

\therefore Pattern q_1 found at $s=2$

Now we need to check whether there are four patterns ahead
gter-4 there are four patterns ahead

gter-4

a	x	c	a	x	a	b	c	l	r	m	n	o	p	q	r	s	t	u	v	w	x	y	z

\downarrow
 $\not x$

$s=3$	a	a	b
-------	---	---	---

Stop here

\therefore Pattern found at $s=2$

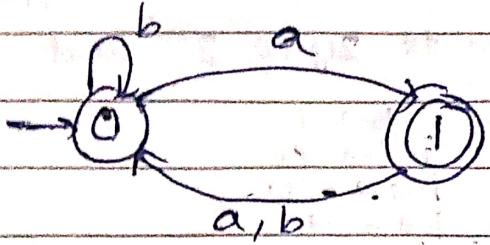
Finite Automata

- Idea of this approach is to build finite automata to scan text T for finding all occurrences of pattern P
- This approach examines each character of text exactly once to find the pattern. Thus it takes linear time for matching but pre-processing time may be large.

A finite automata M is a 5-tuple $(Q, q_0, \Sigma, \delta, F)$

- ① Q is a finite set of states
- ② $q_0 \in Q$ is the starting state
- ③ $F \subseteq Q$ is the set of accepting states
- ④ Σ is the input alphabet (finite)
- ⑤ $\delta : Q \times \Sigma \rightarrow Q$ is the transition function

e.g. of finite automata.



Starting state = 0

Final state = 0

\therefore Here $Q = \{0, 1\}$

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

$$\begin{aligned} s(a, \delta(0, a)) &\rightarrow 1 & \delta(1, a) &\rightarrow 0 \\ s(0, b) &\rightarrow 0 & \delta(1, b) &\rightarrow 0 \end{aligned}$$

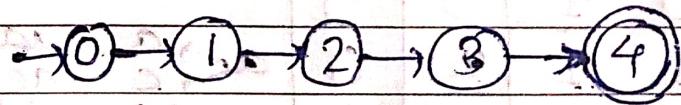
$$q_0 = \{0\} \quad A = \{1\}$$

Table : Input

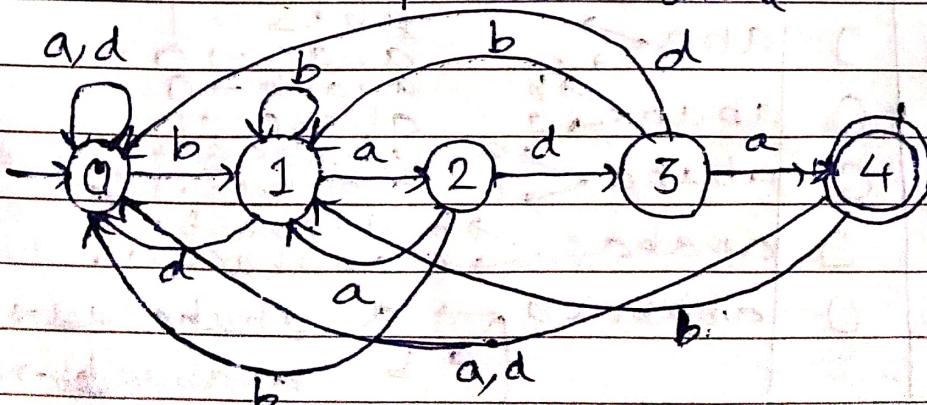
State	a	b
0	1	0
1	0	0

Construction of transition table with given string.

String: bada



If the 1st step make the



Check how many prefixes are matching
how many suffixes

PAGE NO.:

	a	b	d
→ 0	0	<u>1</u>	0
1	<u>2</u>	1	0
2	1	0	<u>3</u>
3	<u>4</u>	1	0
* 4	0	1	0

Note: Prefix starts from 1st letter &
goes up to 2nd last letter

b → b
d → d
ba → ba
bab → b

baa → a
bab → b
badb → b

bada → a
badab → b
badad → d

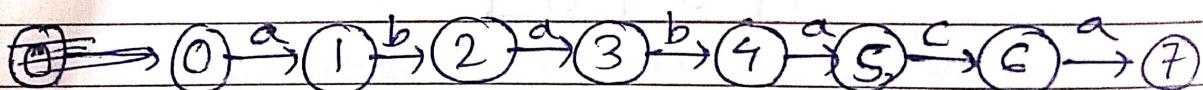
Suffix starts from last letter & goes upto
2nd letter

Now with the help of transition table
complete the automata

E.g. P = ababaca

T = abababacaba

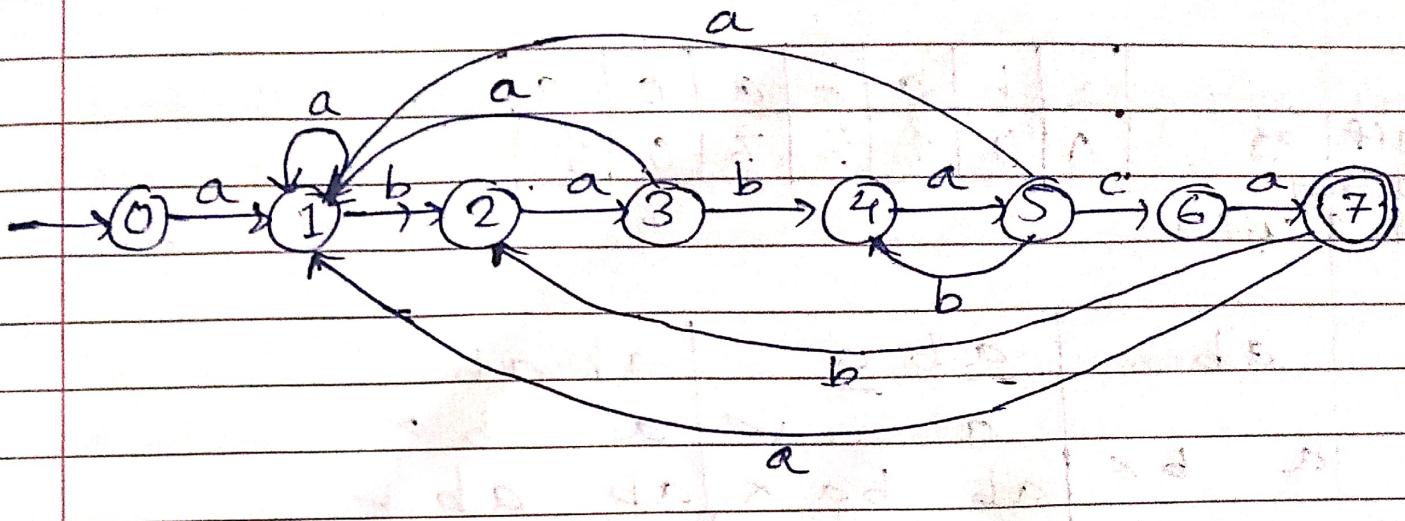
Make finite automata for pattern



	a	b	c
0	<u>1</u>	0	0
1	<u>2</u>	0	0
2	<u>3</u>	0	0
3	<u>4</u>	0	0
4	<u>5</u>	0	0
5	<u>6</u>	0	0
6	<u>7</u>	0	0
7	<u>8</u>	0	0

Transitions:

- 0 → 1 (a), 1 → 2 (b), 2 → 3 (a), 3 → 4 (b), 4 → 5 (a), 5 → 6 (c), 6 → 7 (a)
- 1 → a (2), 2 → b (3), 3 → a (4), 4 → b (5), 5 → c (6)
- 2 → b (3), 3 → b (4), 4 → a (5), 5 → b (6)
- 3 → a (4), 4 → a (5), 5 → a (6), 6 → a (7)
- 4 → c (5), 5 → c (6), 6 → c (7)
- 5 → a (6), 6 → a (7)



Write the test string given

1 2 3 4 5 6 7 8 9 10 11

State	a	b	a	b	a	b	a	c	a	b	a
0	1	2	3	4	5	4	5	6	7	2	3
	1	2	3	4	5	4	5	6	7	2	3

Take these values from transition table
of previous value in the state

We get $q = 7$ = length of substring
check the index above which is 9.

$$\therefore \text{shift} = 9 - 7 = 2$$

\therefore Shift 2 digits from left & you get the
substring from index

$$T.C. = \Theta(m^2 |\Sigma|)$$

Prefix table construction:

Length of
matching
prefix & suffix

a	b	a	b	a	c	a
0	0	1	2	3	0	1

ab	<u>ab a</u>	ab ab
a b x	a <u>a</u> ✓	a b x
ab ba x	ab ab ✓	
	aba bab x	

ababa	<u>ababac</u>
a a ✓	a c x
ab ba x	ab ac x
aba aba ✓	aba bac x
abab baba x	abab abac x
	ababa babac x

<u>ababac a</u>
a a ✓
ab ca x
aba ac a x
babab bac a x
ababa abac a x
ababac babac a x

Knuth-Morrison-Pratt

KMP is a linear time string matching algorithm run in $\Theta(n)$. KMP uses concept of prefix & suffix for generation of Π table.

Π table longest prefix that is same as suffix (LPS)

e.g.

Π ,	a	b	c	d	a	b	e	a	b	f					
	0	0	0	0	1	2	0	1	2	0					

If is the same as prefix table

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

T a b a b c a b c a b a b a b d

P a b a b d
 Π

	0	1	2	3	4	5
P	a	b	a	b	d	
Π	0	0	1	2	0	

Steps : S-1) Take two pointers i & j
 $i=1$ (1st character of text) & $j=0$ (not the 1st character of pattern)

S-2) Compare i & $j+1$

→ If matching move i & j ahead

→ Else back move j up $\Pi[j]$ & i remains where it was

→ Repeat this until ~~$i=1+j$~~ i is not equal to length of string

① $\overset{i}{\overbrace{1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5}} \overset{j}{\overbrace{6 \rightarrow 7 \rightarrow 8 \rightarrow 9 \rightarrow 10 \rightarrow 11 \rightarrow 12 \rightarrow 13 \rightarrow 14 \rightarrow 15}}$
 $a b a b c a b c a b a b a b d$

$j \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$

a	b	a	b	d
0	0	1	2	0

$b \neq c$

put $j \geq 7$ to $T[4] = 2$

$$\therefore j = 2$$

② $\overset{i}{\overbrace{1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5}} \overset{j}{\overbrace{6 \rightarrow 7 \rightarrow 8 \rightarrow 9 \rightarrow 10 \rightarrow 11 \rightarrow 12 \rightarrow 13 \rightarrow 14 \rightarrow 15}}$
 $a b a b c a b c a b a b a b d$

$j \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$

a	b	a	b	d
0	0	1	2	0

a	b	a	b	d
0	0	1	2	0

$b \neq c$

put $j \neq 0$

③ $\overset{i}{\overbrace{1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5}} \overset{j}{\overbrace{6 \rightarrow 7 \rightarrow 8 \rightarrow 9 \rightarrow 10 \rightarrow 11 \rightarrow 12 \rightarrow 13 \rightarrow 14 \rightarrow 15}}$
 $a b a b c a b c a b a b a b d$

$j \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$

a	b	a	b	d
0	0	1	2	0

Now again $a \neq c$

& $j = 0$, so move i ahead

Now $a \neq c$

\therefore Move j to 0

④ 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
ab ab c ab c ab a b a b a b d

$j_0 \rightarrow j_1 \rightarrow j_2 \rightarrow j_3 \rightarrow j_4 \rightarrow j_5$

a	b	a	b	d
0	0	1	2	0

$$c \neq a$$

\therefore Move i ahead

$$a \neq d$$

\therefore Move j to 2

⑤

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15
ab ab c ab c ab a b a b d

j

0	i	$j \rightarrow$	3	4	5
a	b	a	b	d	
0	0	1	2	0	

As $j = \text{length of pattern}$,
we have found the pattern in given
text

$$T.C = O(N+M)$$

$$S.C = O(M)$$