

$T[1 \dots m]$; $m \leq n$ like finding P in T .
 $P[1 \dots m]$

* Alphabet $\Sigma = \{0,1\}$ or $\Sigma = \{a,b,c \dots z\}$.

T & P are strings of Σ

* Pattern occurs with shift s in text T .

eg: \rightarrow abacaabacaca. $\rightarrow T$
 $P \rightarrow$ aaba

$s = 4 \rightarrow P$ starts from $s+1$. {indexing '1' based}.

$$\therefore 0 \leq s \leq m - n$$

* Valid Shift \rightarrow If P occurs with shift s . \rightarrow valid shift.

* String Matching \rightarrow finding all valid shifts.

Prefix (\sqsubset) and suffix (\supset)

* If w is a prefix of string x , $x = wy$. $\nearrow w \sqsubset x$
 $w \rightarrow$ suffix, $x = yw \rightarrow w \supset x$.

eg: \rightarrow ab \sqsubset abacca
 ab \supset acab

empty string;

\downarrow

ϵ

(in Σ^*)

Concatenation of 2 string

\downarrow

$x + y$.

$\Rightarrow |T| = |x| + |y|$

Σ^*

\downarrow

all finite length
string using Σ .

Main string matching algo

for ($S=0$; $S \leq n-m$)

if $P[1..m] = T[S+1..S+m]$
 print "matched".

T.C $\rightarrow O((n-m+1)m)$

String Matching using Finite automaton.

Simple Machine.

process Texts in 'n' time, $O(n)$
 each character exactly once

\rightarrow Processing time $= O(n)$

* Finite automaton M has.

i.) Q is a finite set of states.

ii.) $q_0 \in Q$ is the start state

iii.) $A \subseteq Q$ is accepting state

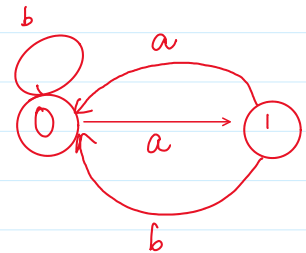
iv.) Σ is finite input alphabet

v.) $\delta \rightarrow$ transition function.

eg $\rightarrow Q = \{0, 1\}$
 $q_0 = 0$
 $\Sigma = \{a, b\}$

State

	input	
	a	b
0	1	0
1	0	0



$\varphi(w) \rightarrow$ state M ends up after scanning the string w

$\delta(wa) = \delta(\varphi(w), a)$ for $w \in \Sigma^*$, $a \in \Sigma$

$P = abaca$

States = $0 \rightarrow 4$

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

Tabular representation of transition funcⁿ.

	a	b	c
0	1	0	0
1	1	2	0
2	3	0	0
3	1	0	4
4	5	0	0
5	1	2	0

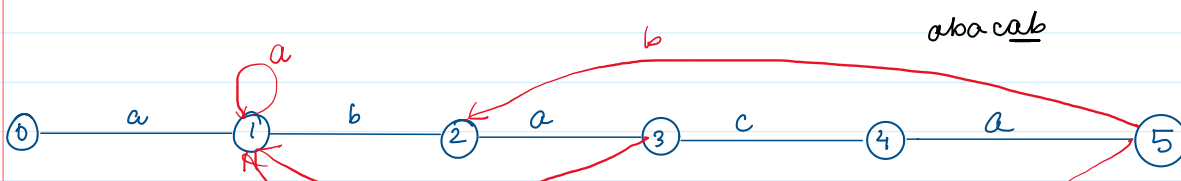
Longest prefix of P that is a suffix of a

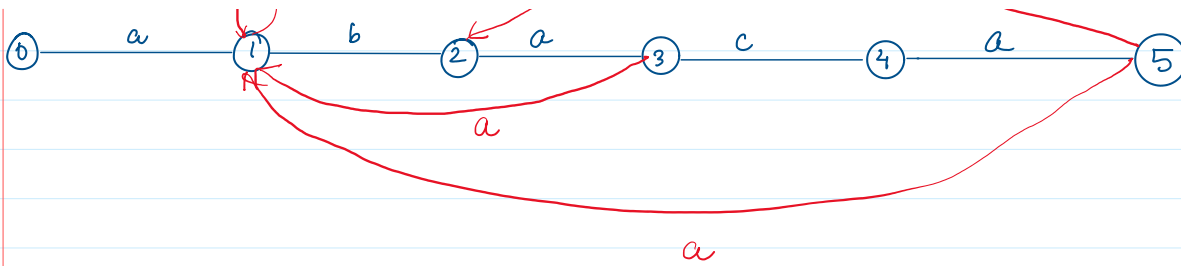
3: $aba|a$
 $abab|b$
 $abac|c$

4: $abac|a$
 $b|b$
 $c|c$

$abaca$

0: a 1: $a|a$ 2: $ab|a$
 b $a|b$ $ab|b$
 c $a|c$ $abc|c$





$$Q = \{0, 1, 2, \dots, m\}$$

start q_0 , accepting state = m .

$\delta(q, a) = r(P_q a)$ because we need to keep track of longest prefix that is matched so far.

$P = ababaca$.

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

$$Q = \{0, \dots, 7\}$$

0: $\frac{a}{b}$
c

1: $\frac{a|b}{a|a}$
 $\frac{a|c}{ab|c}$

3: $\frac{aba|a}{aba|b}$
 $\frac{aba|c}{aba|c}$

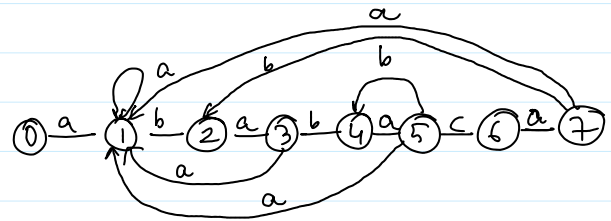
4: $\frac{abab|a}{abab|b}$
 $\frac{abab|c}{abab|c}$

5: $\frac{abab|a}{ab|a|b}$
a|c

6: $\frac{abab|a}{ac|b}$
 $\frac{ac|c}{ac|c}$

$\frac{aca|a}{aca|b}$
c|c

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



$$i = 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12$$

$$f[i] = a \ b \ a \ b \ a \ b \ a \ c \ a \ b \ a$$

$$\text{state } \varphi(i) = 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 4 \ 5 \ 6 \ 7 \ 2 \ 3$$

Accepting state.

KNUTH - MORRIS - PRATT ALGO

* No computation of δ .

* Has a preprocessed array $\pi[1 \dots m] \rightarrow$ order of $\Theta(m)$

$\pi[q] \rightarrow$ info to compute $\delta(q, a)$ but not dependent on a .

In automata $\delta \rightarrow \Theta(m|\Sigma|)$

$\pi \rightarrow \Theta(m)$

Prefix function (π) for a pattern: \rightarrow stores info about how pattern matches with shifts

eg: $T = abcababcaac.$

$P = abab$

$q = 4$, Success

but see $s+1$ is not a valid shift $a \neq b$.

$P = ababa$
 $s' = s+2$

but $s+2$

given $P[1..q]$ match $T[s+1..s+q]$

least shift s' for some $K < q$

$$P[1..K] = T[s'+1..s'+K] \text{ where } s'+K = s+q. \quad \text{--- (i)}$$

In above case $s' = s+2$

$K = 2$

In other words given $P_q \supset T_{s+q}$ find longest proper prefix P_K of P_q that is also suffix of T_{s+q}

Here $K = 2$.

Finding smallest shift $s' =$ finding longest prefix.

$$s' = s + q - K$$

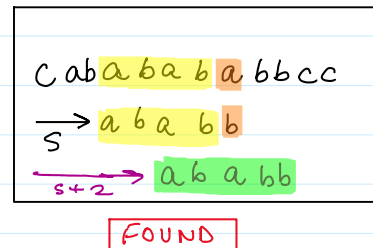
$q \rightarrow$ no. of success matching.

$K \rightarrow$ another longest prefix, also a suffix of T_{s+q}

Best case, $K = 0$, mtlb ki q tak aur koi match nhi mila.

$(s+1, s+2, s+3 \dots s+q) \rightarrow (X)$ Rule out.

egⁿ⁽ⁱ⁾ generate $K < q$, $P_K \supset P_q$



$q = 4$

$K = 2$

$$s' \rightarrow s + (q - K) = s + (4 - 2) = s + 2$$

given a pattern $P[1..m]$ the prefix function $\pi: \{1, 2, \dots, m\} \rightarrow \{0, \dots, m-1\}$

such that

$$\pi[q] = \max \{K: K < q \text{ and } P_K \supset P_q\} \quad \text{Prefix function.}$$

$$\pi[q] = \{K: K < q \text{ and } P_K \supset P_q\}.$$

eg: for a pattern $\rightarrow ababbaca$

$\pi[q] \rightarrow$ longest prefix of P that is a suffix of P_q

$\pi[q] \rightarrow$ longest prefix of P that is a suffix of P_q

eg: $\rightarrow P = ababaea$

1 $\rightarrow a \underline{e}$

2 $\rightarrow a \underline{b} (x)$

3 $\rightarrow a \underline{ba}$

4 $\rightarrow a \underline{bab}$

5 $\rightarrow abab \underline{a}$

6 $\rightarrow ababae$

7 $\rightarrow ababae \underline{a}$

i	1	2	3	4	5	6	7
$P[i]$	a	b	a	b	a	e	a
$\pi[i]$	0	0	1	2	3	0	1

KMP matcher (T, P)

1.) $n \leftarrow T \cdot \text{length}$

2.) $m \leftarrow P \cdot \text{length}$

3.) $q \leftarrow 0$

4.) $\pi = \text{COMPUTE_PREFIX_FUNCTION}(P)$

5.) for $i = 1$ to n .

6.) while ($q > 0$ and $P[q+1] \neq T[i]$)

7.) $q = \pi[q]$

8.) if $P[q+1] == T[i]$

9.) $q = q + 1$

10.) if $q == m$

11.) print "Pattern occurs with shift" $s - m$.

12.) $q = \pi[q]$

i
 abacabab **aa** b e b a b
 ababac **a**

$q = \pi[5] = 3$, checks again
 if $P[4] = T[i]$
 (X)

so $q = \pi[3] = 1$.
 if $P[2] = T[i]$ ✓

COMPUTE PREFIX FUNCTION (P)

1.) $m \leftarrow P \cdot \text{length}$

2.) $\pi[1 \dots m] = \text{new array}$

3.) $\pi[1] = 0$

4.) $k = 0$

5.) for $q = 2$ to m

6.) while ($k > 0$ and $P[k+1] \neq P[q]$)

7.) $k = \pi[k]$

Running time = $\Theta(m)$

$$8.) \quad \forall p(K+1) = p(q)$$

$$9.) \quad K = K + 1$$

$$10.) \quad \pi(q) = K$$

11.) return π

Running time analysis.

* K increments only in line 9.

↳ at most $(m-1)$ increments

* On entering for loop $K=0$, $q=2$, $\Rightarrow K < q$

only assignment $\pi[i] = 0$

$$\pi[q] = K \quad \Rightarrow \quad \pi[q] < q.$$

for all $q \in \{1, 2, \dots, m\}$

\therefore each iteration of while loop decreases value of K

* K is always +ve

\therefore The total decrease in K is bounded by the total increase in it.

which is $(m-1)$

\therefore while loop runs in $\Theta(m)$.

By similar analysis, KMP matcher is $\Theta(n)$

Preprocessing $O(m|\Sigma|) \rightarrow \Theta(m)$

Actual $\Theta(n) \rightarrow \Theta(n)$