

Unit-8 :- String matching

DOMS

Page No. I

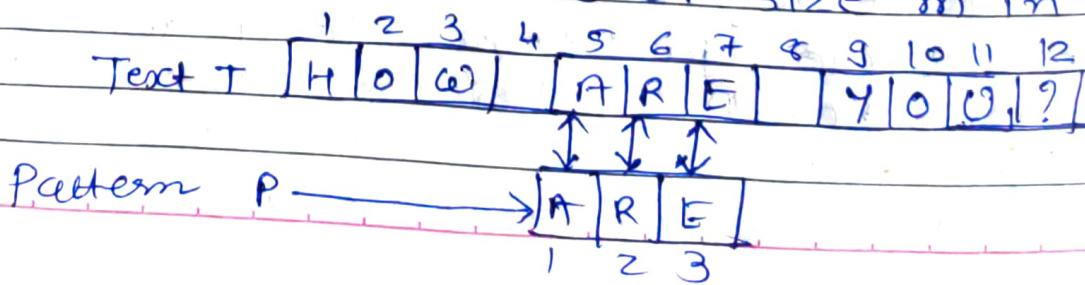
Date / /

* Topics :-

- Introduction
- The naive string matching algorithm
- The Rabin - Karp algorithm
- String matching with finite automata
- The Knuth - Morris - Pratt algorithm.

* Introduction :-

- String matching operation is core part in many text processing applications.
- Objective of string matching algorithm is to find pattern P from given text T.
- Typically $|P| \ll |T|$.
 - In design of compilers and text editors, string matching operation is crucial. So locating P in T efficiently is very important.
- Problem :- Given some text string $T[1...n]$ of size n, find all occurrences of pattern $P[1...m]$ of size m in T.



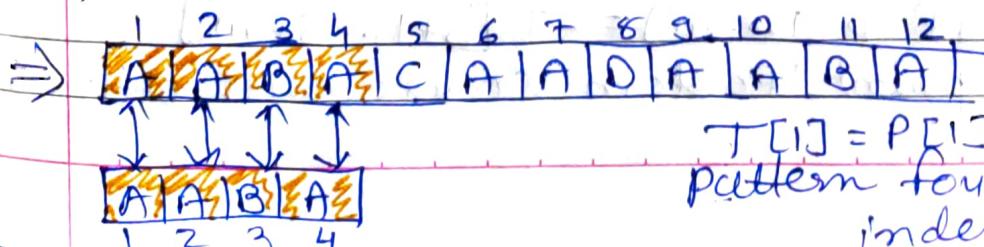
* Naive String Matching Algorithm :-

- This is simple and inefficient brute force approach.
- It compares first character of pattern with searchable text. If match is found, pointers in both strings are advanced.
- If match is not found, pointer of text is incremented and pointer of pattern is reset. This process is repeated till the end of text.
- Naive approach does not require any pre-processing.
- Given text T and pattern P, it directly starts comparing both strings character by character.
- After each comparison, it shifts pattern string one position to the right.

\Rightarrow Example:-

Text : A A B A A C A A D A A B A

Pattern : A A B A



$T[1] = P[1], T[2] = P[2]$
 $T[3] \neq P[3]$
 pattern found at $T[4] = P[4]$
 index 1

Text

- [A] A B | A | A | C | A | A | D | A | A | B | A

↓ ↓
[A] A B | A
Pattern

then shift one position right.

Text

- [A] A B | A | A | C | A | A | D | A | A | B | A

↓ x
[A] A B | A
Pattern

- we shift one position right till end.

- At least,

1 2 3 4 5 6 7 8 9 10 11 12 13
[A] A B | A | A | C | A | A | D | A | A | B | A
↑ ↓ ↓ ↑
[A] A B | A

- So, pattern found at position 7 & 10.

⇒ Algorithm :-

Naïve - string matching (T, P)

1. $n = T.length$

2. $m = P.length$

3. for $s = 0$ to $n-m$

4. if $P[1..m] == T[s+1..s+m]$

5. print "Pattern occurs with shift" s

⇒ Time Complexity :- $O((n-m+1)m)$

* Rabin-Karp Algorithm :-

- The Rabin-Karp string matching algo. calculates a hash value for the pattern, as well as for each character subsequences of text to be compared.
- If the hash values are unequal, the algorithm will determine the hash value for next character sequence.
- If the hash values are equal, the algorithm will analyze the pattern and the character sequence.
- In this way, there is only one comparison per text subsequence, and character matching is only required when the hash values match.
- If pattern is found then it is called hit. otherwise it is called spurious hit.

⇒ Example :-

$$T = 31415926535$$

$$P = 26$$

$$q = 11$$

- First find $P \bmod q = 26 \bmod 11 = 4$

- Now finding hash values of text

$31 \bmod 11 = 9$ 4 4 4 9

3	1	4	1	5	9	1	2	6	5	3	5
int											
3	8			4		10			2		

- $15 \bmod 11 = 4$ equal to 4
 is spurious hit
 but string is not match
 $15 \neq 26$
 - $59 \bmod 11 = 4$ equal to 4
 is spurious hit
 but string is not match
 $59 \neq 26$
 - $92 \bmod 11 = 4$ equal to 4
 but string is not match
 $92 \neq 26$
 - $26 \bmod 11 = 4$ equal to 4
 exact match $26 = 26$
- (Text) (pattern)
 is called hit.

→ Pattern found at position 7.

\Rightarrow Algorithm :-

Rabin-Karp matching (T, P, d, q)

$n \leftarrow T.length$;

$m \leftarrow P.length$;

$h \leftarrow d^{m-1} \text{ mod } q$;

$p \leftarrow 0$;

$to \leftarrow 0$;

for $i \leftarrow 1$ to n do

$p \leftarrow (dp + P[i]) \text{ mod } q$

$to \leftarrow (d * to + T[i]) \text{ mod } q$

for $s \leftarrow 0$ to $n-m$ do

if $p = ts$ then

if $P[1...m] == T[s+1...s+m]$ then

print "pattern occurs with
shift's"

if $s < n-m$ then

$ts+1 \leftarrow (d(ts - T[s+1]h) + T[s+m+1]) \text{ mod } q$

* String matching with Finite Automata :-

- In this approach, 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 preprocessing time may be large.
 - finite automaton is defined by triple,
- $$M = \{Q, \Sigma, q_0, F, \delta\}$$

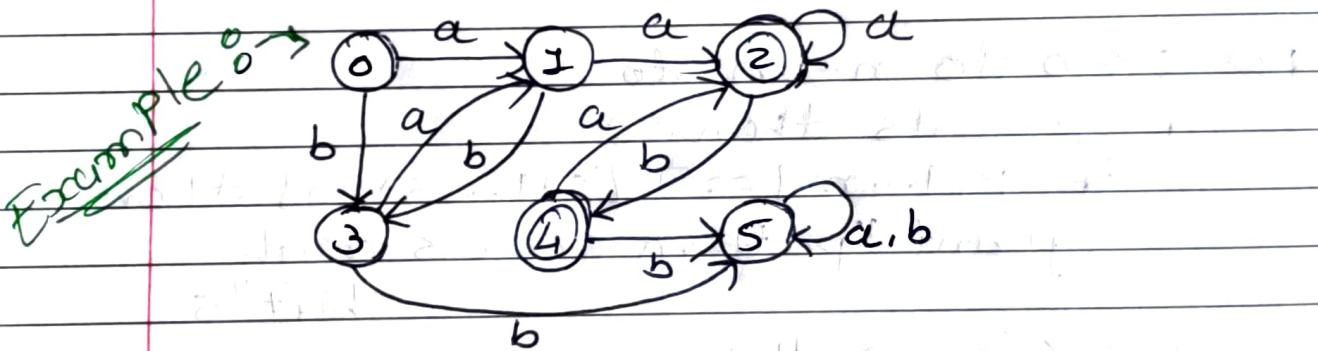
where, Q = Set of states in finite automaton

Σ = set of input symbols

q_0 = Initial state

F = set of final states

δ = Transition function



$$Q = \{0, 1, 2, 3, 4, 5\}$$

$$q_0 = 0$$

$$F = \{2, 3\}$$

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

Transition function:

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

\Rightarrow Algorithm :-

$n \leftarrow \text{length}[T]$

$q \leftarrow 0$

for $i \leftarrow 1$ to n do

$q \leftarrow s(q, T[i])$

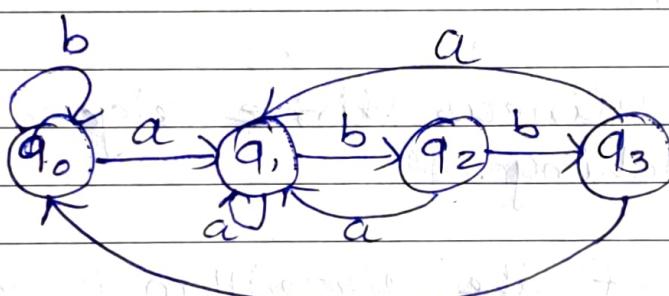
if $q == m$ then

print "Pattern occurs with
shift" $i-m$

\Rightarrow Example :-

text $T = ababbacababba$

pattern $P = aabb$



$$\mathcal{Q} = \{q_0, q_1, q_2, q_3\}$$

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

$q_0 \rightarrow$ starting state

$q_3 \rightarrow$ End state

Transition table (8) :-

s	a	b
q_0	q_1	q_0
q_1	q_1	q_2
q_2	q_1	q_3
q_3	q_1	q_0

a	b	a	b	b	a	a	b	a	b	b	a
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									
a	b	b									

* The Knuth-Morris-Pratt Algorithm :- (KMP Algo.)

- KMP is the linear time algorithm for string matching.
- Main idea of the algorithm is to avoid computation of transition function δ and reducing useless shifts performed in naive approach.
- This improvement is achieved by using auxiliary function Π .
- Function Π is very useful, it computes the numbers of shifts of pattern by finding pattern matches within itself.

⇒ Example 8-

text T = bacbabababcceeb

text T = bacbabababcceeb.

pattern P = abababcceeb

→ First find out prefix table,

↳ Consider string "a"

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

prefix = ϵ
suffix = ϵ

↳ Consider string "ab"

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

prefix = ϵ, ab
suffix = ϵ, b

↳ Consider string "ceba"

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

prefix = $\epsilon, ce, ceb, ceba$
suffix = ϵ, e, eb, bae

↳ Consider string "abceb"

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

prefix = $\epsilon, a, ab, abc, abce, abceb$
suffix = $\epsilon, b, eb, ceb, abceb$

↳ Consider string "ababa"

prefix = ^, a, ab, aba, abab

suffix = ^, a, ba, cba, abab

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

↳ Consider string "ababac"

prefix = ^, a, ab, aba, abab, ababa

suffix = ^, c, ac, bc, bac, abac, babac

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

↳ Consider string "ababacce"

prefix = ^, a, ab, abc, abac, abaca, ababac

suffix = ^, ce, ca, acc, bacc, cebacce, babacce

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

- Now compute $T[0]$ with $P[0]$

$T : \boxed{b|a|c|b|a|b|a|b|a|c|a|b}$
~~#~~
 $P : \boxed{a|b|a|b|a|c|q}$

here $a \neq b \Rightarrow$ check prefix table
at $\pi[0] = 0$; take $P[0]$
compute $T[1]$ with $P[0]$

$T : \boxed{b|a|c|b|a|b|a|b|a|c|a|b}$
~~#~~
 $P : \boxed{a|b|a|b|a|c|q}$

Compare $T[2]$ with $P[1]$

No match $C \neq b$

- Now, $P[0]$ with $T[2]$ compare,

No match,

- Now, compare $T[3]$ with $P[0] \Rightarrow$ No match

- Compare $P[0]$ with $T[4]$

$T : \boxed{b|a|c|b|a|b|a|b|a|c|a|b}$
 $P : \boxed{a|b|a|b|a|c|q}$

- Pattern found at $T[4]$.

Algorithm 8-

Algorithm Knuth-morris-Pratt (T, P)

$n = T.length$

$m = P.length$

$\pi = \text{compute-prefix-function}(P)$

$q = 0$

for $i = 1$ to $n+m$

 while $q > 0 \& P[q+1] \neq T[i]$

$q = \pi[q]$

 if $P[q+1] == T[i]$

$q = q + 1$

 if $q == m$

 print "occurs shift" $i - m$

$q = \pi[q]$

Compute-Prefix (P)

$m = P.length$

let $\pi[1..m]$ be a new array

$\pi[1] = 0$

$K = 0$

for $q = 2$ to m

 while $K > 0 \& P[K+1] \neq P[q]$

$K = \pi[K]$

 if $P[K+1] == P[q]$

$K = K + 1$

$\pi[q] = K$

return π