

# **CHAPTER 06 – PART II**

## **STRING SEARCH**

## **ALGORITHMS**

**Data Structures and Algorithms (17ECSC204)**

**Prakash Hegade**  
**School of CSE, KLE TECH**

**2017-18**

## Brute Force String Search Algorithm:

A brute force string matching algorithm is quite obvious. Align the pattern against the first  $m$  characters of the text and start matching the corresponding pairs of characters from left to right until either all  $m$  pairs of characters match or a mismatching pair is encountered. In the latter case, shift the pattern one position to the right and resume character comparisons, starting again with the first character of the pattern and its counterpart in the text.

### Example:

Text: N O B O D Y \_ S A W \_ M E

Pattern: S A W

N O B O D Y \_ S A W \_ M E

S A W

  S A W

    S A W

      S A W

        S A W

          S A W

            S A W

              S A W

**ALGORITHM** BruteForceStringMatch( $T[0 \dots n-1]$ ,  $P[0 \dots m-1]$ )

// Implements brute force string match

// Input: An array  $T[0 \dots n-1]$  of  $n$  characters representing a text and an array  $P[0 \dots m-1]$  of  $m$  characters representing a pattern

// Output: The index of the first character in the text that starts a matching substring or -1 if the search is unsuccessful

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

$j \leftarrow 0$

  while  $j < m$  and  $P[j] = T[i+j]$  do

$j \leftarrow j+1$

  if  $j = m$

    return  $i$

return -1

The worst case would be that the algorithm would have to make all the  $m$  comparisons before shifting the pattern and this can happen for  $n-m+1$  tries.

Thus, in the worst case, the algorithm efficiency is in  $O(nm)$ .

## Boyer-Moore Algorithm:

Algorithm involves constructing two tables:

### Bad Symbol Shift table:

Given a character  $c$ ,  $T(c)$  is computed as:

- the pattern length  $m$ , if  $c$  is not among the first  $m-1$  characters of the pattern
- the distance from the rightmost  $c$  among the first  $m-1$  characters of the pattern to its last character, otherwise.

The size of shift is computed by:  $T(c) - k$

where  $k$  is number of matched characters.

If  $T(c) - k$  is  $\leq 0$ , then shift by one position to right.

### Good Suffix Shift Table:

→ check if there is another occurrence of matched pattern not preceded by same character as in its last occurrence.  $d_2$  is distance between such second rightmost occurrence of pattern & its rightmost occurrence.

→ If not, find the longest prefix of size  $\leq k$  ( $k$  = matched pattern length), that matches the suffix of same size  $\leq k$ . If such a prefix exists, shift  $d_2$  is computed as the distance between

this prefix & the corresponding suffix.

→ otherwise  $d_2$  is set to pattern length  $m$ .

The final shift distance  $d$  is computed by:

$$d = \begin{cases} d_1 & \text{if } k=0 \\ \max\{d_1, d_2\} & \text{if } k>0 \end{cases}$$

$$\text{where } d_1 = \max\{TCC) - k, 1\}$$

Examples:

1. Construct bad symbol shift table for the pattern BARBER

C	A	B	C	D	E	F	...	R	...	Z	-
TCC)	4	2	6	6	1	6	6	3	6	6	6

2. Construct a good suffix shift table for: ABCBAB

k	pattern	d2
1	ABC <u>B</u> AB	2
2	A <u>BCB</u> AB	4
3	A <u>BCBA</u> B	4
4	A <u>BCBAB</u>	4
5	<u>ABCBA</u> B	4

3. Apply Boyer-Moore algorithm on the given Text & pattern.

Text: BESS-KNEW-ABOUT-BAOBABBS

Pattern: BAOBABBS

Bad Symbol table:

C	A	B	C	D	...	O	...	Z	-
T(c)	1	2	6	6	6	3	6	6	6

Good suffix table:

k	pattern	$d_2$
1	B A O <u>B A B</u>	2
2	<u>B A O B A B</u>	5
3	<u>B A O B A B</u>	5
4	<u>B A O B A B</u>	5
5	<u>B A O B A B</u>	5

Iterations:

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
B	E	S	S	-	K	N	E	W	-	A	B	O	U	T	-	B	A	O	B	A	B	S
B	A	O	B	A	B																	
$d_1 = 6 - 0 = 6$									X	✓	✓											
						B	A	O	B	A	B											
						$d_1 = 6 - 2 = 4$										X	✓					
						$d_2 = 5$																
						$d = \max\{4, 5\} = 5$																
										B	A	O	B	A	B							
										$d_1 = 6 - 1 = 5$												
										$d_2 = 2$												
										$d = \max\{5, 2\} = 5$												
																✓	✓	✓	✓	✓	✓	✓
																B	A	O	B	A	B	-
match.																						

4. Construct Bad Symbol Table for the pattern:  
CONSISTING

C	A	B	C	...	G	H	I	...	N	O	...	S	T	...	Z	-
T(c)	10	10	9	10	10	10	2	10	1	8	10	4	3	10	10	10

A Better Undustable table:

C	O	N	S	I	S	T	I	N	G
9	8				4	3	2	1	10

5. Construct bad symbol table for the pattern:  
DISGUSTING

D	I	S	G	U	S	T	I	N	G
9			6	5	4	3	2	1	

6. Construct bad symbol table and good-suffix  
table for: 00001

C	0	1
T(c)	1	5

k	pattern	d <sub>2</sub>
1	0000 <u>1</u>	5
2	0000 <u>1</u>	5
3	0000 <u>1</u>	5
4	0000 <u>1</u>	5

78. Compute the two Shift tables for the pattern  
String: 10000

Bad-Symbol table

c	o	l
T(c)	1	4

Good Suffix table

k	pattern	d2
1	1000 <u>0</u>	3
2	100 <u>00</u>	2
3	10 <u>000</u>	1
4	1 <u>0000</u>	5

Prefixes & Suffixes:

Text: School

Prefixes: S  
Sc  
Sch  
Scho  
School  
School

Suffixes: School  
chool  
hool  
ool  
ol  
L

For the string "School": S, Sc, sch, scho, school  
are proper prefixes. Similar explanation holds  
good for proper Suffixes too.

i.e; A proper prefix or proper suffix of a string  
is not equal to the string itself.

## Knuth-Morris-Pratt Algorithm:

Principle to Generate the prefix table  $\pi$ :

"Find the length of the longest proper prefix in the subpattern that matches a proper suffix in the same subpattern."

Example:

Pattern: ABAB

char	A	B	A	B
index	0	1	2	3
value	0	0	1	2

a)

Substring: A

proper prefix: NULL

proper suffix: NULL

b) Substring: AB

proper prefix = A

proper suffix = B

c) Substring: ABA

proper prefix: A, AB

proper suffix: BA, A

d) Substring: ABAB

proper prefix: A, AB, ABA

proper suffix: BAB, AB, B

Shift is computed by the formula:

$k - \pi[k-1]$  where  $k$  is number of matched chars.

Eg:- If  $k=2$  then  $shift = 2 - \pi[2-1]$

$$= 2 - \pi[1]$$

$$= 2 - 0$$

$$= \underline{2}$$



Example:

1. Search for pattern ababaca in the text  
bacbababababacaab using KMP algorithm.

Filling the prefix table  $\pi$ , for the pattern:

char	a	b	a	b	a	c	a
index	0	1	2	3	4	5	6
value	0	0	1	2	3	0	1

<u>Substring</u>	<u>Proper Prefix</u>	<u>Proper Suffix</u>	<u><math>\pi</math>-value</u>
a	Null	Null	0
ab	a	b	0
aba	a, ab	a, ba	1
abab	a, ab, aba	b, <del>ab</del> , bab	2
ababa	a, ab, aba, abab	a, ba, aba, baba	3
ababac	a, ab, aba, abab, ababa	c, ac, bac, abac, babac	0
ababaca	a, ab, aba, ababa, ababa, ababac	a, ca, aca, baca, abaca, babaca	1

### Tracing:

i)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & c & a & a & b \\ & \# & & & & & & & & & & & \\ a & b & a & b & a & c & a & & & & & & \end{array}$   
shift by 1

ii)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & b & a & c & a & a & b \\ & \# & \# & & & & & & & & & & & \\ a & b & a & b & a & c & a & & & & & & & \end{array}$   
 $k=1$   
 $shift = 1 - \pi[0]$   
 $= 1 - 0 = 1$

iii)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & b & a & c & a & a & b \\ & \# & & & & & & & & & & & & \\ a & b & a & b & a & c & a & & & & & & & \end{array}$   
shift by 1

iv)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & b & a & c & a & a & b \\ & \# & & & & & & & & & & & & \\ a & b & a & b & a & c & a & & & & & & & \end{array}$   
shift by 1

v)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & b & a & c & a & a & b \\ & \# & \# & \# & \# & \# & \# & \# & & & & & & \\ a & b & a & b & a & c & a & & & & & & & \end{array}$

$k=5$   
 $shift = 5 - \pi[4] = 5 - 3 = 2$

vi)  $\begin{array}{cccccccc} b & a & c & b & a & b & a & b & a & b & a & c & a & a & b \\ & \# & \# & \# & \# & \# & \# & \# & \# & \# & & & & \\ a & b & a & b & a & c & a & & & & & & & \end{array}$

- match found.

2. Apply KMP for pattern AABAB in the text  
AABBAABBAABBAABA.

### Efficiency Analysis:

1. When searching for the first occurrence of the pattern in Boyer-Moore algorithm, the worst case efficiency is known to be linear.

If implemented as presented in original paper has worst case running time of  $O(n+m)$  only if the pattern does not appear in text. When appears, the worst case is  $O(nm)$ .

### 2. KMP Analysis

$O(m)$  to compute prefix function values

$O(n)$  to compare pattern to text

Total =  $O(n+m)$

## Space & Time tradeoff Technique:

A way of solving a problem in less time by using more storage or by solving a problem in very little space by spending a long time.

### Examples:

- Boyer Moore algorithm
- Hashing.

The most common situation is an algorithm involving a lookup table.

\* — \* — \* — \* — \*