

Design and analysis of Algorithms

Sl. (Prof. Parminder Kaur Wadhwa)

String Matching Algorithms

Part A :- includes:-

1. Simple / Naiive Algorithm] same
2. Brute Force Algorithm.
3. Rabin - Karp Algorithm.

Part B :- includes:-

4. Boyer - Moore Algorithm
5. Knuth - Morris - Pratt Algorithm
6. Application of string matching algorithms.

Design and Analysis of Algorithms

(Prof. Parminder Kaur Wadhwa)

String Matching Algorithms

(Part - B)

Topics covered :-

4. Boyer-Moore Algorithm
5. Knuth-Morris-Pratt Algorithm
6. Applications of string Matching algorithms.

4. Boyer-Moore Algorithm

4.1 The Basic Boyer-Moore Algorithm

Involving :-

1. Looking-glass heuristic
2. character-jump heuristic
(or Bad-character Rule)

4.2. Boyer-Moore Algorithm using Good-Suffix Rule (along with looking-glass heuristic)

4.3. Boyer-Moore Algorithm combining all the three components together:-

1. looking-glass heuristic
2. Bad Character Rule
3. Good-Suffix Rule

4.4. Boyer-Moore Hopcroft Algorithm

4.1.

The Basic Boyer-Moore Pattern Matching algorithm

Boyer-Moore algorithm learn from the past character comparisons and skip the unnecessary comparisons in future.

It works on the following two heuristics:-

1. Looking glass heuristic
2. Character-jump heuristic (Bad-character rule)

1. Looking glass heuristic :-

It says that, when testing a possible placement of P against T, begin the comparison from the end of P and move backward to the front of P.

In other words, we can say that

"Try alignments in 'left-to-right' order,
and

try character-comparisons in 'right-to-left' order."

For example:-

(LEFT)

T: f o \times m e f o r
 P: f o r

(RIGHT)

Compare \times and r first
 (comparisons will
 be started
 from right-to
 -left)

Pattern
 will move
 from left-to-right
 for alignment

This is opposite to naive or simple or brute
 force pattern matching algorithm.

Boyer-Moore
Algorithm

Left

T: f o \times m e f o r

Right

(left)

P: f o r \leftarrow 1st comparison

(from Right to
 left)

Naive / Simple
Brute Force
Algorithm

(Right)

T: f o \times m e f o r

P: f o r

1st comparison
 (from left to right)

2. Character-jump heuristic (or Bad-character Rule)

Let T is the text given

P is the pattern to be matched

$T[i] = c$, a character at i th index due to which mismatch has occurred.

The Bad-character Rule says that:-

Upon mismatch, skip alignments/placements until:-

(a) mismatch becomes a match

Or

(b) pattern P moves past mismatched character

Thus, the Character-jump heuristic says that:-

during the testing of a possible placement of P against T , a mismatch of test character $T[i] = c$ with the corresponding pattern character $P[j]$ is handled as follows:-

If c is not contained anywhere in P , then shift P completely past

$T[i]$ (because it cannot match any character in P).

Otherwise, shift P until an occurrence of character c in P gets aligned with $T[i]$, (making the mismatch a match).

The actual search for the pattern, in basic Boyer-Moore algorithm takes $\mathcal{O}(nm)$ time in the worst case,
(where $n \rightarrow$ is the size of the given text T)

$m \rightarrow$ is the size of the pattern P to be matched / found)

Example :- of :-

The Basic Boyer-Moore Algorithm for Pattern Matching

T: fox me for

P: for

1. T: f o ~~x~~ m e f o r

P: f o ~~r~~ _{1st comparison}

X mismatch, 'x' is not required
so skipping 'x' all together
(otherwise, it will go on creating
mismatches).

2. T: f o x i m e f o r

P: f o ~~r~~ _{2nd comparison}

X mismatch, but 'f' is
required

so, aligning 'f' of text T with
'f' of pattern P,

(thus, making the mismatch, a match)

3. T: f o x m e f o r

P:

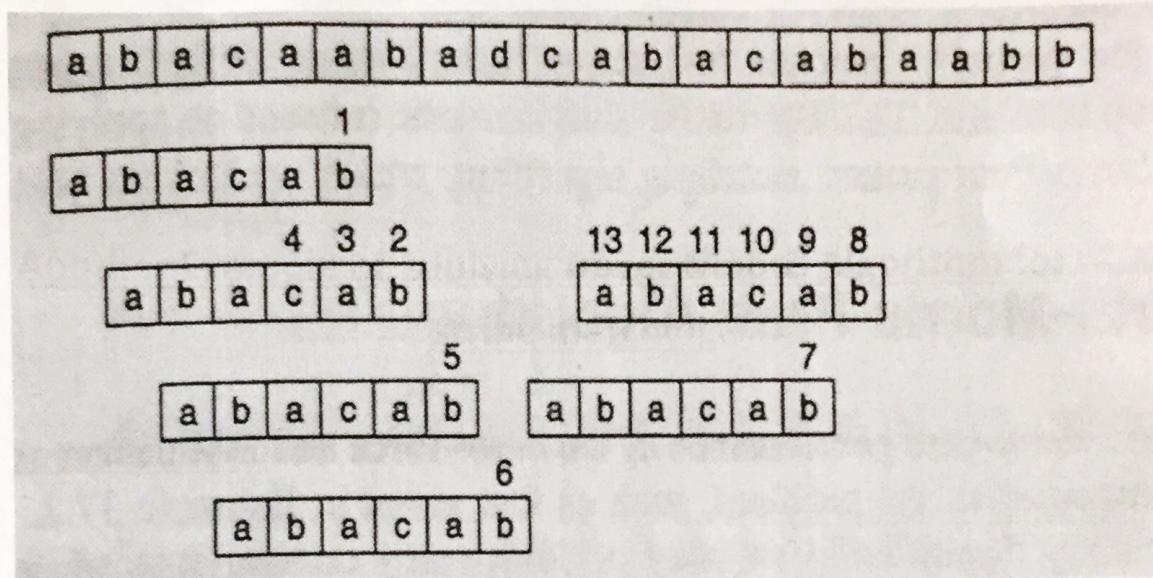
f o r _{5th 4th 3rd comparisons}
_{✓ ✓ ✓} are successful.

(MATCH FOUND)

(BASIC)

Another Example of applying Boyer Moore Algorithm for string matching

To find pattern: a b a c a b from the given text



Due to mismatch in **comparison 1**, due to 'a' , both 'a' of text and pattern are aligned.

Comparisons 2nd and 3rd are successful.

4th comparison is a mismatch, due to 'a'. The immediate unchecked 'a' of pattern is aligned with 'a' of pattern p.

5th comparison is a mismatch, due to a. The immediate unchecked 'a' of pattern is aligned with 'a' of pattern p.

6th comparison is a mismatch, due to 'd'. But 'd' is not required. So, cross 'd' otherwise, it will go on creating mismatches.

7th comparison is a mismatch due to a . The immediate unchecked 'a' of pattern is aligned with 'a' of pattern p.

The comparisons **8th, 9th, 10th, 11th, 12th and 13th** are successful and the pattern is matched in given text .

Applying Bad-character rule (Boyer-Moore Algorithm) for DNA pattern Matching :-

DNA is Deoxyribonucleic acid

It is the molecule that carries genetic information for the development and functioning of an organism. It is a molecule that supplies the genetic instructions that tell living creatures how to develop, live and reproduce. DNA can be found inside every cell and is passed down from parents to their offsprings.

DNA is made of two linked strands that wind around each other to resemble a twisted ladder - a shape known as double helix.

Each strand has a backbone made of alternating sugar (deoxyribose) and phosphate groups.

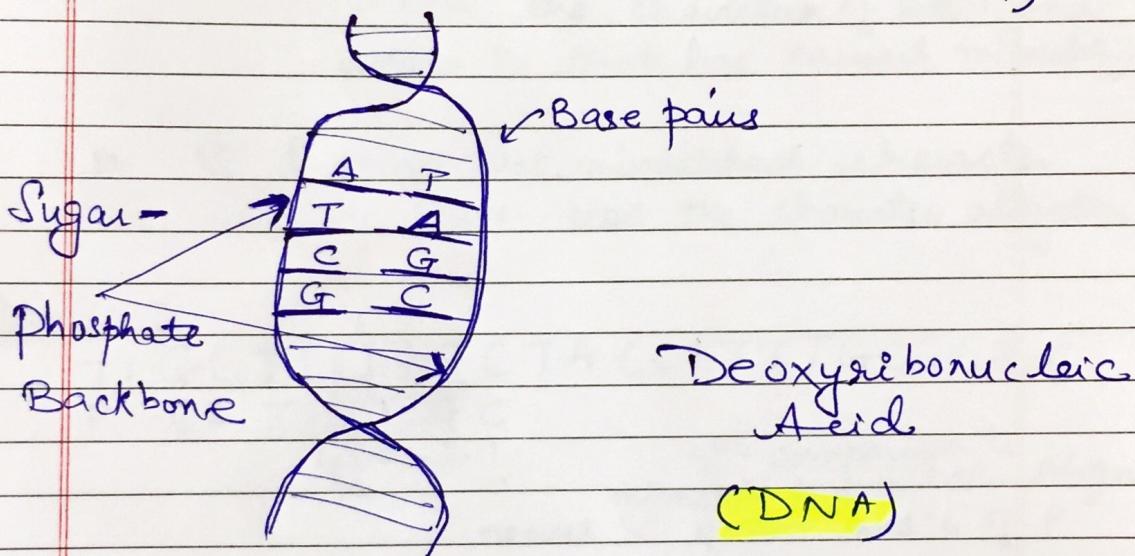
Attached to each sugar is one of four bases:-

1. Adenine (A)
2. Cytosine (C)
3. Guanine (G)
4. Thymine (T)

The two strands are connected by chemical bonds between the bases:-

adenine bonds with thymine (A bonds with T)

cytosine bonds with guanine (C with G)



The sequence of the bases along DNA's backbone encodes biological information, such as the instructions for making a protein or RNA molecule.

The DNA pattern matching, helps the scientists to find the locations of particular DNA subsequences in a DNA sequence.

Applying Boyer-Moore: Bad character Rule for DNA pattern matching

The bad character rule is:- upon mismatch :-
skip alignments until :-

(a) mismatch becomes a match

(align the characters of text T and pattern P that has caused mismatch)

or (b) P moves past mismatched character

(i.e. cross the character altogether)

Step 1:

T: G C T T C T G C T A C C T T T T G C G C G C
 P: C C T T T T G C
 4 3 2 1 X

4th comparison is mismatch due to 'C'. Align nearest 'C' of Text and 'C' of P.

Step 2: T: G C T T C T G C T A C C T T T T G C G G G E

P: C C T T T T G C
 6 X 5 ✓

6th comparison is mismatch due to 'A'. But it is not required. So, cross it altogether.

Step 3: T: G C T T C T G C T A C C T T T T G C G G G E

P:
 CC TTTT G C
 14 13 12 11 10 9 8 7
 v v v v v v v v

from 7th to 14th, all comparisons are successful. [Match is Found]

4.2.

Boyer-Moore Algorithm using good-suffix rule (along with looking glass heuristic)

The bad-character rule says that upon mismatch we want to shift P such that the mismatch becomes a match
but,

the good-suffix rule (or good-character rule) says that for the characters we did match, we want to shift P such that we don't turn any of those matches into a mismatch at the end of that shift.

Bad-character Rule

Upon mismatch,
shift P such
that mismatch
becomes a match
i.e. (try to convert a
mismatch into a
match)

Good-character rule (or good-suffix rule)

Upon mismatch,
shift P such that
old matches done should
not become new mismatches
i.e. (try to maintain
matches already
done)

Boyer-Moore: Good Suffix Rule

Let t = substring matched by inner loop;

Skip until:-

(a) there are no mismatches between P and t .

(i.e. old matches are still maintained to maximum extent).

Or

(b) P moves past t .

(cross altogether)

Example:-

T: C G T G C @ \overleftarrow{t} T A C T T A C T T A C
P: C T T A C T T A C
 4 3 2 1
 X ✓ ✓ ✓ mismatch

try to maintain matches already done.
Don't make old matches as new mismatches.

Applying Boyer-Moore - Good Suffix Rule for DNA pattern matching

Step 1 :-

T: CG T G C **O** **TAC** TT ACTT ACTT ACTT AC
 P: CTT A **CTT AC**
 4 3 2 1
 X v v v

Try to look occurrence of **TAC** i.e. 't' already matched in pattern P (on left side) and align them.

Step 2 :-

T: CG T G C **O** **TAC TT ACTT ACTT ACTT AC**,
 P: CTT A **CTT AC** TT AC
 12 11 10 9 8 7 6 5
 X v v v v v v v

12th comparison has caused mismatch.

But 5th to 11th comparisons are successful.

Try to find occurrence of already matched character sequence on left of P.

X **TACTTAC**
 CT T A C TTAC

The sequence "CTTAC" should align with each other to maintain old matches to their maximum level.

Step 3 :-

T: CG T G C C T A **CTT AC TTAC** TT ACTT AC
 P: CTT A C TTAC
 21 20 19 18 17 16 15 14 13
 v v v v v v v v v

Comparisons from 13th to 21st are successful.

(MATCH FOUND)

4.3.

Boyer-Moore Algorithm combining all the three components together:-

1. Looking glass heuristics (Compare from right side)
2. Bad-character Rule
3. Good-Suffix Rule

In practice, both rules:-
bad-character rule and good-suffix rule can be used together.

Whenever, we have a mismatch, we can try both rules. Each rule will tell us some amount that we can shift the pattern P , in other words:-

Some number of alignments/placements that we can skip (i.e. to avoid comparisons which are unnecessary), and

we are simply going to take the maximum of the two.

i.e. Upon mismatch:-

Calculate:- how many alignments are skipped by both rules.

The rule which gives the maximum number of alignments/placements is chosen to apply for that particular mismatch, since it is avoiding more comparisons than other.

Applying Boyer-Moore Algorithm supporting both rules

Bad-character rule
Good-Suffix rule

for DNA pattern matching

Step 1 :-

T: GTT ATA G C T G A T C G C G G C G T A G C G G C G A A
P: G T A G C G G C G

! mismatch

Since, no character is matched, the good suffix rule is not applicable here as there is no matching characters. So, we cannot use the good-suffix rule.

We have to use the bad-character rule.
So, 'T' has caused the mismatch.

But 'T' is required. So, aligning 'T' of text with 'T' of pattern P,
(making a mismatch as a match)

Step 2 :-

T: GTT ATA G e T G A T C G C G G C G T A G C G G C G A A
P: G T A G C G G C G
i 2 3 4 5 6

If we align 'T', and make the mismatch as a match, we are skipping 6 alignments (or comparisons).

On the other hand, good-suffix rule, has no matching characters and thus, nothing is there to maintain (and

is not applicable).

i.e. $bc = 6, gs = 0$

($bc \Rightarrow$ bad-character
 $gs \Rightarrow$ good suffix)

Here, $bc > gs$.

So, Bad-character rule is applied.

Step 2:-

Alignments are skipped (by gs)

T: G T T A T A G C T G A T C G C G G C G T A G C G G C G A A
P: G T A G C G G C G
 5 4 3 2
 X ✓ ✓ ✓

If bad-character rule is applied, i.e. 'c' which has caused mismatch are aligned (of text T and pattern P)

⇒ then, no alignment (or comparison) is avoided or skipped.

But, if we apply good-suffix rule:-

GCG is already matched. We need to find the occurrence of GCG on the left of P (to their maximum level) and GCG another ~~another~~ occurrence exists.

If we align GCG of T and P, then, it skips 2 alignments.

$$\therefore bc = 0, gs = 2$$

$$gs > bc$$

\therefore good-suffix rule is chosen to apply for this mismatch.

Step 3:-

T: G T T A T A G C T G A T E G C G G C G T A G C G G C G A A

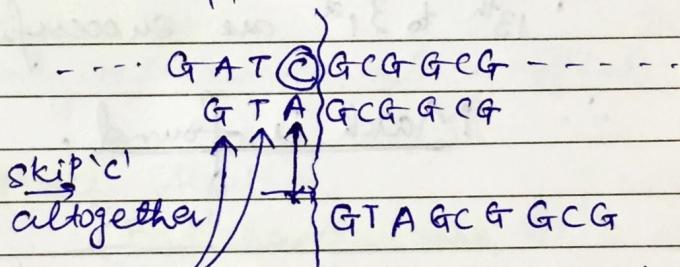
p: G T A G C G G C G

12 11 10 9 8 7 6
X ✓ ✓ ✓ ✓ ✓ ✓

12th comparison is a mismatch, due to 'C'.

If bad-character-rule is applied, then 'C' is not required and we need to skip 'C' altogether.

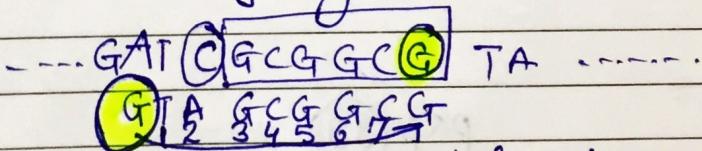
If bc rule is applied then :-



Then, these two alignments (or comparisons) are skipped.

$$\therefore bc = 2$$

But if , we apply good-suffix rule, then:-



G can be maintained as match, leading to

skipping of 7 alignments.

$$\therefore g_s = 7$$

Now, $bc = 2$, $g_s = 7$

$g_s > bc \quad \therefore$ good-suffix-rule
is chosen to be applied
for this mismatch.

∴

Step 4 :-

T: G T T A T A G C T G A T C G C G G C G T A G C G G C G A A
G T A G G G G G
21 20 19 18 17 16 15 14 13
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

All comparisons from :-
13th to 21st are successful.

∴ Match is found.

The worst-case running time of basic Boyer-Moore algorithm [(using bad-character-rule) only] is $O(nm)$

But, the Boyer-Moore algorithm achieves running time = $O(n+m)$

If the good-suffix-rule is also used (in combination with bad-character-rule).

Brute-force Pattern Matching Algo = $O(nm)$

Rabin-Karp " " " " = $O(n+m)$

Basic Boyer-Moore " " " " = $O(nm)$
 (Only bad-character-rule is applied)

Boyer Moore Algorithm = $O(n+m)$
 (Bad-character-rule &
 Good-Suffix-rule
 are applied)

- Brute force algorithm is good for smaller texts and small patterns.
- Rabin-Karp Algo is good when it is used for multiple-pattern matching

(i.e. finding all occurrences of a set of patterns in a fixed text).

- Boyer-Moore Algorithm is good for single-pattern matching.
 → (using bad-character rule & good-suffix rule)

Boyer-Moore algorithm is extremely fast on large text (relative to the length of the pattern).

An attempt has been made by the researchers to have a variant of Boyer-Moore algorithm (involving both rules): -

in order to achieve an average time complexity of :-

$O(n)$

which is linear, ($n \rightarrow$ size of text T) and is independent of:-

size of pattern m. This algorithm is Boyer-Moore Hopcroft Algorithm.

4.4

Boyer-Moore Horspool Algorithm

In 1980, Nigel Horspool (a retired professor of computer science, of the University of Victoria, Canada) invented a variant of Boyer-Moore Algorithm, which is called as:- (or known as):-

Boyer-Moore Horspool Algorithm.

[Note:- Boyer-Moore algorithm was developed

- a) by Robert S. Boyer (professor of computer science at University of Texas, US)
- b) and J. Strother Moore (Computer Scientist, at University of Texas, US)

They developed Boyer-Moore Algorithm in the year 1977.

The Boyer-Moore Horspool algorithm, was able to give the average time complexity
 $= O(n)$

which is linear & is independent of the size of the pattern m.

But, its worst case is $O(nm)$ which is same as that of brute force algorithm.

Boyer-Moore Horspool

Algorithm

Boyer-Moore Algorithm

(using bad-character
1 good-suffix rule)

Average-case time
complexity = $O(n)$

Average case time
complexity = $O(n+m)$

Worst case time
complexity
 $= O(nm)$

Worst case time
complexity
 $= O(n+m)$

The explanation of Boyer-Moore Horspool Algorithm

- It includes: the step of:- making a
 - Bad Match Table
- Compare pattern to text, starting from rightmost character in the pattern
- If mismatch, move pattern forward corresponding to value in the Bad Match Table.

For example:- (Applying Boyer-Moore Horspool Algorithm).

T: trust hard toothbrushes

P: tooth

Constructing Bad Match Table:- (use the formula)

$$\text{Value} = \text{length} - \text{index} - 1$$

and \rightarrow (keep last letter's value = length)

and \rightarrow (and every other letter = length)

TOOTH Length = 5 (Total characters in P)
 Index 0 1 2 3 4

Step 1

letter	T	O	H	*
Value				(Any other)

We need to fill this table. There are three different letters in pattern p: "Tooth". These are T, O and H.

Step 2:-

$$\text{for letter T} \rightarrow \text{length} - \text{index} - 1$$

$$= 5 - 0 - 1$$

$$= 4$$

TOOTH
0 1 2 3 4

(0 is the index of T)

Letter	T	O	H	*
Value	4			

Step 3:-

T	O	T	H
index = 0	1	2	3

For letter O at index 1

$$\begin{aligned} &\Rightarrow \text{length} - \text{Index} - 1 \\ &= 5 - 1 - 1 \\ &= 3 \end{aligned}$$

letter	T	O	H	*
value	4	3		

Step 4:-

T	O	T	H
0	1	2	3

For letter O at index 2

$$\begin{aligned} &\Rightarrow \text{length} - \text{Index} - 1 \\ &= 5 - 2 - 1 \\ &= 2 \end{aligned}$$

Updating the value for 'O' in the Bad-Match-Table as follows: —

letter	T	O	H	*
value	4	2	.	

Step 5:-

T	O	T	H
0	1	2	3

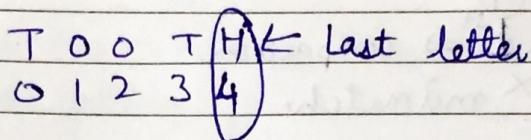
↓
index

for letter T at index 3 $\Rightarrow \text{length} - \text{Index} - 1$
 $= 5 - 3 - 1$
 $= 1$

Updating the value for letter 'T' in the Bad Match Table as follows:-

Letter	T	O	H	*
Value	1	2		

Step 6:-



Keep the value for the LAST LETTER as length i.e. = 5

Letter	T	O	H	*
Value	1	2	5	

↓
last letter's value = length

Step 7:-

for "any other letter" i.e. *, keep the value = length = 5.

Letter	T	O	H	*
Value	1	2	5	5

So, this is the Bad Match Table that will be used when, Boyer-Moore Hopcroft algorithm will be applied.

Proceeding with the algorithm;

Bad Match Table :-

Letter	T	O	H	*
Value	1	2	5	5

Step 1

T: TRUST HARD TOOTH BRUSHES
 P: TOOTH
 $\begin{matrix} \text{I}^{\text{st}} \\ \times \end{matrix} \Rightarrow \text{comparison}$
 X mismatch

The mismatch is due to 'T'. Look the value for letter 'T' in the Bad Match Table :-
 The value = 1 (for letter T)

So, shift the pattern by '1' character.

Step 2 :-
 T: TRUST HARD TOOTH BRUSHES
 P: $\begin{matrix} \text{I} \\ \text{TOOT} \\ \text{H} \end{matrix}$
 $\begin{matrix} \text{X}^{\text{2nd}} \\ \text{V} \end{matrix} \text{ comparison}$
 4th comparison is a mismatch.

for this alignment we, compared "H" in the text first.

TRUST HARD . . .

TOOTH $\begin{matrix} \text{X} \\ \text{V} \end{matrix}$ This was the first comparison done in this particular alignment

(But, further it led to a mismatch)

So, we will look for the value of letter H in the table, i.e. 5

[Or, you may look * in the table, as S does not exist]

So, the pattern will be shifted by 5 characters.

Step 3:-

T: TRUST HARD TOOTH BRUSHES
P: ~~T~~~~R~~~~U~~~~S~~~~T~~~~H~~~~A~~~~R~~~~D~~~~T~~~~O~~~~O~~~~T~~~~H~~

X

5th comparison is a mismatch due to O

X

look 'O' in the table. Value for 'O' is 2
So, shift the pattern by 2.

Step 4:-

T: TRUST HARD TOOTH BRUSHES
P: ~~T~~~~R~~~~U~~~~S~~~~T~~~~H~~~~A~~~~R~~~~D~~~~T~~~~O~~~~O~~~~T~~~~H~~

X

6th comparison is a mismatch due to T

look for value of letter 'T' in the table.
It is 1.

So, shift the pattern by 1.

Step 5:-

T: TRUST HARD TOOTH BRUSHES
P: ~~T~~~~R~~~~U~~~~S~~~~T~~~~H~~~~A~~~~R~~~~D~~~~T~~~~O~~~~O~~~~T~~~~H~~

~~T~~~~O~~~~O~~~~T~~~~H~~

11 10 9 8 7
VVVVV

all comparisons
are successful.

(Match is found)

5. Knuth - Morris - Pratt Algorithm

(KMP) algorithm

The KMP algorithm, avoids the waste of information in case a pattern character does not match in the text. Thus, it achieves a running time of $O(n+m)$.

This is the worst case time complexity which is acceptable as,

even in the worst case, any pattern matching algorithm with this time complexity is examining all the characters of the text and all the characters of the pattern at least once.

The main idea of the KMP algorithm is to preprocess the pattern string P so as to compute a failure function ρ that indicates the proper shift of P so that, to the largest extent possible, we can reuse previously performed comparisons.

The algorithm for calculating prefix function is as follows:-

Compute - Prefix-function (π)

{ m is length of P

$$\pi(1) = 0$$

(Assign π of first character as 0)

$$k=0$$

for $q \leftarrow 2$ to m do

{

while ($R > 0$. and $p[k+1] \neq p[q]$) do

$$R \leftarrow \pi(R);$$

// R when not matching

if ($p[k+1] = p[q]$) then

$$\{ k = k + 1 \}$$

// increment R when

$$\pi[q] = R$$

matching

return π

}

g

Applying KMP algorithm for given text T and pattern P

$T = \text{a b a c a a b a c c a b a c a b a a b b}$
 $P = \text{a b a c a b}$, $m = \text{size of pattern} = 6$

First we need to preprocess the pattern P:—

$P \Rightarrow P[1] P[2] P[3] P[4] P[5] P[6]$
a b a c a b

Let $\pi[1] = 0$

(for $q=1$)

$\pi[2] = \pi[1]$
 \Rightarrow is set as
Zero)

and $R=0$

<u>R</u>	<u>Q</u>	<u>P[R+1]</u>	<u>P[Q]</u>	<u>P[1] = a.</u>
O	1			
O	2	$P[0+1] = P[1] = a$	$P[2] = b$	$k > 0? \text{ No}$
O	3	$P[0+1] = P[1] = a$	$P[3] = a$ $k > 0? \times$ $a = a, k = k+1$	$\pi[3] = 1$
I	4	$P[1+1] = P[2] = b$	$P[2] = P[4] = c$ $k > 0? \text{ Yes}$ $b \neq c$ $k = \pi[k] = \pi[1] = 0$	$\pi[4] = 0$
O	5	$P[0+1] = P[1] = a$	$P[2] = P[5] = a$ $k > 0? \times$ $a = a$ $k = k+1 = 0+1 = 1$	$\pi[5] = 1$
I	6	$P[1+1] = P[2] = b$	$P[6] = b$ $k > 0? \text{ Yes}$ but $b = b$ $k = k+1 = 1+1 = 2$	$\pi[6] = 2$

<u>new R</u>	<u>$\pi[Q]$</u>	<u>$\pi[1] = 0$</u> ← initialized
O	$\pi[2] = 0$	
I	$\pi[3] = 1$	
O	$\pi[4] = 0$	

<u>$P[Q]$</u>
1 2 3 4 5 6 a b a c a b
0 0 1 0 1 2

The prefix function for pattern abacab

is computed as follows:-

	1	2	3	4	5	6
P[]	a	b	a	c	a	b
π[]	0	0	1	0	1	2

In KMP algorithm, the comparison is from Left to Right:-

T: aba ca a b a cc a ba c a b a a b b
 P: $\begin{matrix} a & b & a & c & a \\ 1 & 2 & 3 & 4 & 5 \end{matrix}$ (b) $\begin{matrix} & & & & 6 \\ \checkmark & \checkmark & \checkmark & \checkmark & \checkmark & X \end{matrix}$

6th comparison is a failure for 'b' in pattern

The value in table, for this 'b' is 2.

So, start comparing from 2nd character in the pattern from this position.

T: aba ca a b a cc a ba c a b a a b b
 P: $\begin{matrix} a & b & a & c & a & b \\ & & & & & \end{matrix}$ (b) $\begin{matrix} & & & & & \\ & & & & & \end{matrix}$

Already \checkmark
 5th matched
 (no comparison
 will be done
 here)

Value is '0', so, More the string altogether

T: a b a ca a b a cc a b a c a b a a b b
 P: $\begin{matrix} a & b & a & c & a & b \\ & & & & & \end{matrix}$ (b) $\begin{matrix} & & & & & \\ & & & & & \end{matrix}$

But 7th comparison is a mismatch.

9th comparison is a mismatch. Value for this is zero, so more the pattern altogether.

T: aba caab a c } aba cabaabb

P:

ab a cab
10th

comparison is a mismatch
due to a.

value for this a is zero. So move
the string altogether.

T: aba caab a c c } aba cabaabb

P:

ab a cab
11 12 13 14 15 16
✓ ✓ ✓ ✓ ✓ ✓

Match is found.

C. Applications of String Matching Algorithm

The string matching algorithms provide / play key role in various real world problems or applications.

A few of its important applications are:- as:-

- a) Spell checkers
- b) Spam filters
- c) Intrusion Detection System
- d) Search Engines
- e) Plagiarism Detection
- f) Bioinformatics
- g) Digital forensics
- h) Information Retrieval systems,

etc. . . .