

String Matching Algorithm

In Computer Science, a string matching algorithm (sometimes called string-searching algorithm or pattern matching algorithm) are an important class of string algorithms that try to find a place where one or several strings (also called patterns) are found within a larger string or text.

Text Array $T[1 \dots n]$

Pattern Array $P[1 \dots m]$

$$m \leq n$$

Elements of 'P' and 'T' can be $\{0, 1, \dots, 9\}$ or $\{a, b, c, \dots, z\}$

Example 3

T

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

P

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

Applications of String Matching Algorithms

- (1) Plagiarism Detection
- (2) Bioinformatics and DNA Sequencing
- (3) Digital Forensics
- (4) Spelling checker
- (5) Spam filters
- (6) Search engines
- (7) Intrusion Detection System
- (8) Content search in large databases.

String Matching Algorithms

- (1) Naive Algorithm
- (2) Rabin Karp Algorithm
- (3) Knuth Morris Pratt (KMP) Algorithm
- (4) String matching with Finite Automata.
- (5) Boyer Moore Algorithm

Naive String Matching Algorithm

It is the simplest method for pattern matching.

The naive approach tests all the possible placements of pattern $P[1 \dots m]$ relative to text $T[1 \dots n]$.

We try shifts $S = 0, 1, \dots, n-m$ successively and for each shift S , compare $T[S+1 \dots S+m]$ to $P[1 \dots m]$.

The naive algorithm finds all valid shifts using a loop that checks the condition $P[1 \dots m] = T[S+1 \dots S+m]$ for each of the $n-m+1$ possible value of S .

Naive - String - Matcher(T, P)

1. $n \leftarrow \text{length}[T]$
2. $m \leftarrow \text{length}[P]$
3. for $S \leftarrow 0$ to $n-m$
4. do if $P[1 \dots m] = T[S+1 \dots S+m]$
5. then print "Pattern occurs with shift" S

Analysis

Here, for-loop executes for $n-m+1$ (we need at least m characters at the end) times and in iteration we are doing m comparisons.

So the total complexity is $O(n-m+1)$.

Example \Rightarrow

T

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

P

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

Soln

T

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

P

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

 $S=0$

T

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

P

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

 $S=1$

T

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

P

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

 $S=2$

T

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

P

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

 $S=3$

\therefore valid shifts are $S=2$

\therefore P is found in T at $S=2$.

Example \Rightarrow

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

Soln

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=0$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=1$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=2$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=3$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=4$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=5$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=6$

T

1	0	1	1	1	0	1	1	1	0
---	---	---	---	---	---	---	---	---	---

P

1	1	1
---	---	---

 $S=7$

∴ valid shifts are $S=2$ and $S=6$

∴ P is found in T at $S=2$ and $S=6$.

Rabin-Karp String Matching Algorithm

- Developed by Michael O Rabin and Richard M Karp.
- This algorithm match the pattern in given text using hash function.

Examples \Rightarrow

Text (T) = 31234862

Pattern (P) = 234

Hash function used is $P \bmod q$, where P is pattern and q is any random prime number.

let $q = 13$

$\therefore P \bmod q = 234 \bmod 13$

$= 0$
 \uparrow
hash code

Note \Rightarrow If hash code match, then

- if all characters of 'P' and 'T' matches
then it is valid Hit

- if all the characters of 'P' and 'T' does not
matches, then it is spurious hit

$$312 \bmod 13 = 0 \leftarrow \text{Spurious Hit}$$

$$123 \bmod 13 = 6$$

$$234 \bmod 13 = 0 \leftarrow \text{valid Hit}$$

$$348 \bmod 13 = 10$$

$$486 \bmod 13 = 5$$

$$862 \bmod 13 = 4$$

Q $T = 2359023141526739921$

$$P = 31415$$

$$q = 13$$

Soln $P \bmod q = 31415 \bmod 13$
 $= 7$

$$23590 \bmod 13 = 8$$

$$35902 \bmod 13 = 9$$

$$59023 \bmod 13 = 3$$

$$90231 \bmod 13 = 11$$

$$02314 \bmod 13 = 0$$

$$23141 \bmod 13 = 1$$

$$31415 \bmod 13 = 7 \leftarrow \text{valid Hit}$$

$$14152 \bmod 13 = 8$$

$$41526 \bmod 13 = 4$$

$$15267 \bmod 13 = 5$$

$$52673 \bmod 13 = 10$$

$$26739 \bmod 13 = 11$$

$$67399 \bmod 13 = 7 \leftarrow \text{Spurious Hit}$$

$$73992 \bmod 13 = 9$$

$$39921 \bmod 13 = 11$$

Q $T = 3141592653589793$

$$P = 26$$

$$q = 11$$

Soln $P \bmod q = 26 \bmod 11$
 $= 4$

$$\therefore 31 \bmod 11 = 9$$

$$14 \bmod 11 = 3$$

$$41 \bmod 11 = 8$$

$$15 \bmod 11 = 4 \leftarrow \text{Spurious Hit}$$

$$59 \bmod 11 = 4 \leftarrow \text{Spurious Hit}$$

$$92 \bmod 11 = 4 \leftarrow \text{Spurious Hit}$$

$$26 \bmod 11 = 4 \leftarrow \text{Valid Hit}$$

$$65 \bmod 11 = 10$$

$$53 \bmod 11 = 9$$

$$35 \bmod 11 = 2$$

$$58 \bmod 11 = 3$$

$$89 \bmod 11 = 1$$

$$97 \bmod 11 = 9$$

$$79 \bmod 11 = 2$$

$$93 \bmod 11 = 5$$

Rabin - Karp (T, P)

$$n = T. \text{ length}$$

$$m = P. \text{ length}$$

$$h^P = \text{Hash}(P[...]) \quad \text{i.e. } P \bmod q$$

$$h^T = \text{Hash}(T[...]) \quad \text{i.e. } T \bmod q$$

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

$$\text{if } (h^P == h^T)$$

matching all characters of
P and T

$$\text{if } (P[0...m-1] == T[s+0, ... s+m-1])$$

Print "Pattern found with shift" s

$$\text{if } (s < n-m)$$

$$h^T = \text{Hash}(T[s+1...s+m])$$

Time Complexity

If no spurious hit, then T.C = $O(n-m+1)$

If there is spurious hit, then T.C = $O(mn)$

Knuth - Morris - Pratt Algorithm

- KMP is a linear time string matching algorithm.
- KMP uses concept of prefix and suffix for generation of π table.
- Worst case running time of algorithm is $O(m+n)$

Let Pattern : a b c d a b c

∴ Prefix can be a, ab, abc, abcd, ----

Suffix can be c, bc, abc, dabc, ----

Here, abc is the prefix same as suffix.

How to find π table (or) Prefix function (or) Prefix Table

Pattern : a b c d a b c a

Solution ⇒
1 2 3 4 5 6 7 8
a b c d a b c a
0 0 0 0 1 2 3 1

Pattern : a a b a a b a a a

Solution ⇒
1 2 3 4 5 6 7 8 9
a a b a a b a a a
0 1 0 1 2 3 4 5 2

Pattern : a b c d a b e a b f

Solution ⇒
1 2 3 4 5 6 7 8 9 10
a b c d a b e a b f
0 0 0 0 1 2 0 1 2 0

Pattern: $a a b a a c a a b a a$
 0 1 2 3 4 5 6 7 8 9 10

0	1	0	1	2	0	1	2	3	4	5
---	---	---	---	---	---	---	---	---	---	---

Pattern: $a a a c a a a d a c$
 0 1 2 3 4 5 6 7 8 9

0	1	2	0	1	2	3	3	3	4
---	---	---	---	---	---	---	---	---	---

Initially j is at 0 and i is at 1.
 $[j=0]$ if $P[j] \neq P[i]$, then put 0 and increment i by 1.
 if $P[j] = P[i]$, then put $j+1$ and increment both i and j by 1.

$[j \neq 0]$ if $P[j] \neq P[i]$, then check value at $(j-1)$ character.
 and place j there and compare with $P[i]$
 again, now (j) matched, then put $j+1$ and
 increment both i and j by 1.
 (else) repeat again by placing j at
 $(j-1)$ character value and compare.

Que

Pattern: $a c a c a b a c a c a b a c a c a c$
~~0 1 2 3 4 5 6 7 8 9 10~~
 $\pi = 0 0 1 2 3 0 1 2 3 4 5 6 7 8 9 10 11 4$

Example :

String: a b a b c a b c a b a b a b d

Pattern: a b a b d

Soln:

String: ⁰a ¹b ²a ³b ⁴c ⁵a ⁶b ⁷c ⁸ab ⁹ab ¹⁰ab ¹¹d

(P)
pattern:

	0	1	2	3	4	5
	a	b	a	b	d	
π	0	0	1	2	0	

compare $T(i)$ with $P(j+1)$

↳ if matches then increment i and j by 1.

↳ if does not matches then move j to the index given under current- j position (π value). and again compare $P(j+1)$ with $T(i)$ else ~~and~~ increment i .

string: a b a b c a b c a b a b a b d

pattern:

π	0	0	1	2	0
-------	---	---	---	---	---

Pattern found in string.

Time complexity = $O(m+n)$
 $\downarrow \quad \downarrow$
 a table search

Note \Rightarrow

The implementation of KMP algorithm is efficient because it minimizes the total number of comparisons of the pattern against the input string.

Compute - Prefix - Function (P)

1. $m \leftarrow \text{length}[P]$
2. $\pi[1] \leftarrow 0$
3. $k \leftarrow 0$
4. for $q \leftarrow 2$ to m
5. while $k > 0$ and $P[k+1] \neq P[q]$
6. $k \leftarrow \pi[k]$
7. if $P[k+1] = P[q]$
8. $k \leftarrow k+1$
9. $\pi[q] \leftarrow k$
10. return π

Example : Pattern : a b a b a c a

$$m = \text{length}[P] = 7$$

$$\pi[1] = 0$$

$$k = 0$$

P: a b a b a c a

q: 1 2 3 4 5 6 7

$\pi: 0$

Step 1: $q=2, k=0$

$$\pi[2] = 0$$

P: a b a b a c a

q: 1 2 3 4 5 6 7

$\pi: 0 0$

Step 2: $q=3, k=0$ (here $P[1] == P[3]$)
 $\pi[3] = 1$ So, $k=1$

p : a b a b a c a
 q : 1 2 3 4 5 6 7
 π : 0 0 1

Step 3: $q=4, k=1$ (here $P[2] \neq P[4]$ is False)
 also $P[2] == P[4]$
 $\Rightarrow k=2$

$\therefore \pi[4] = 2$

p : a b a b a c a
 q : 1 2 3 4 5 6 7
 π : 0 0 1 2

Step 4: $q=5, k=2$ ($2 > 0$ and $P[3] \neq P[5]$)
 (T) (F)
 but $P[3] == P[5]$
 $\Rightarrow k=3$

$\therefore \pi[5] = 3$

p : a b a b a c a
 q : 1 2 3 4 5 6 7
 π : 0 0 1 2 3

Step 5: $q=6, k=3$

here, $3 > 0$ and $p[4] \neq p[6]$
(T) (T)

$$\Rightarrow k = \pi[3]$$

$$k = 1$$

$$\therefore \pi[6] = 1$$

p :	a	b	a	b	a	c	a
q :	1	2	3	4	5	6	7
π :	0	0	1	2	3	1	

Step 6: $q=7, k=1$

here, $1 > 0$ and $p[2] \neq p[7]$
(T) (T)

$$\Rightarrow k = \pi[1]$$

$$k = 0$$

Now, $p[1] == p[7]$ is True

$$\therefore k = k + 1$$
$$= 1$$

$$\pi[7] = 1$$

p :	a	b	a	b	a	c	a
q :	1	2	3	4	5	6	7
π :	0	0	1	2	3	1	1

The running time of prefix function is $O(m)$.

KMP-Matcher (T, P)

1. $n \leftarrow \text{length}[T]$
2. $m \leftarrow \text{length}[P]$
3. $\pi \leftarrow \text{Compute-Prefix-Function}(P)$
4. $q \leftarrow 0$
5. for $i \leftarrow 1$ to n
6. while $q > 0$ and $P[q+1] \neq T[i]$
7. $q \leftarrow \pi[q]$
8. if $P[q+1] = T[i]$
9. $q \leftarrow q+1$
10. if $q = m$
11. print "Pattern occurs with shift" $i-m$
12. $q \leftarrow \pi[q]$

Q Compute the prefix function π for the pattern
ababbabbabbabbabb.

Soln

$$m = \text{length}[P] = 19$$
$$\pi[4] = 0$$
$$k = 0$$

Step 1: $q=2, k=0$

Here, $P[1] == P[2]$ (False)

$$\lambda [2] = 0$$

$p: a \ b \ a \ b \ b \ a \ b \ b \ a \ b \ b \ a \ b \ a \ b \ b \ a \ b \ b$
 $q: 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 11 \ 12 \ 13 \ 14 \ 15 \ 16 \ 17 \ 18 \ 19$
 $\pi: 0 \ 0$

After Compute-Prefix-Function (P) for q values 2 to 19, we have

$P:$	a	b	a	b	b	a	b	b	a	b	b	a	b	a	b	b	a	b	b
$q:$	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
$\pi:$	0	0	1	2	0	1	2	0	1	2	0	1	2	3	4	5	6	7	8

Q Apply KMP algorithm.

String: bacbabababacaaab

Pattern: ababaca

Soln String: bacbabababacaaaab
Pattern: ababaca

The total number of shifts that took place for the match are $= i - m$
 $= 13 - 7$
 $= 6$ shifts.

Note \Rightarrow The KMP algorithm never needs to move backwards in the input text T . It makes the algorithm good for processing very large files.