

Design and analysis of Algorithms

(Prof. Parminder Kaur Wadhwa)

String Matching Algorithms

Part A :- includes:-

1. [Simple/Naive Algorithm]
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 A)

Topics covered :-

1. Simple / Naive Algorithm
 2. Brute Force Algorithm
 3. Rabin-Karp Algorithm.
- } are same

String Matching Algorithms

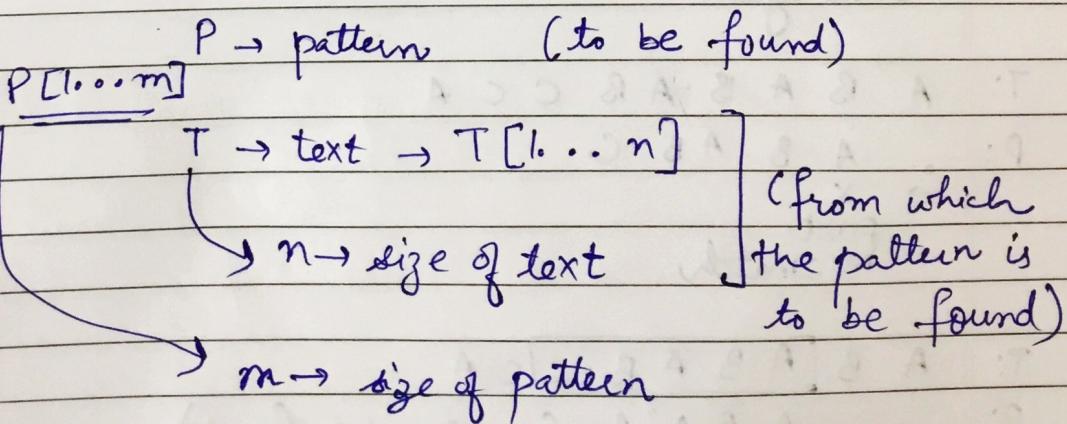
1. Simple or Naive String Matching
 2. Brute Force Pattern matching
 3. Rabin-Karp Algorithm
 4. Boyer-Moore Algorithm
 5. Knuth-Morris-Pratt (KMP) Algorithm.
 6. Applications of string matching algorithms.
- These are same algorithms

The string matching problem:-

It is the problem of detecting the occurrence of a particular substring called pattern, in another string, called the text.

The array P is considered to contain the pattern.

The array $T[1 \dots n]$, is considered to contain the text.



1. Simple or Naive String Matching

In this, comparisons are done (in left to right order) on the pairs of characters in array P and T.

When first mismatch occurs, the pattern is moved one position forward with respect to text, and comparisons start again at the left end of pattern.

For example :-

T: A B A B A B C C C A

P: A B A B C

⇒ T: A B A B A B C C C A

P: A B A B C
✓ ✓ ✓ ✓ X

Mismatch occurred

Now, Shift P towards left by one position and start again

T: A B A B A B C C C A

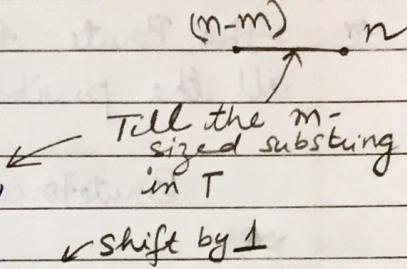
P: . A B A B C
→ X

first
Mismatch

T: A B [A B A B C] C A

P: . → A B A B C (match found)

Simple String Matching (T, P)

$n \leftarrow \text{length}[T]$
 $m \leftarrow \text{length}[P]$
 for $i \leftarrow 0$ to $n-m$
 do if $P[1..m] = T[i+1 .. i+m]$


then print "Match found with shift" i

Running time $O[(n-m+1)m] \approx O(nm)$

If $m = \frac{n}{2}$ (i.e. if substring is half the size of text T)

then $\Rightarrow O(n^2)$

2. Brute-Force Pattern Matching

The Brute-Force pattern matching, tests all the possible placements of P relative to T .

BruteForceMatch (T, P)

```
for i ← 0 to n-m (for each candidate
{ index in T)
    do
        j ← 0
        compare characters
        while (j < m and T[i+j] = P[j]) do
            {
                j ← j+1 // shift by one
            }
        if j = m then
            return i // substring
            found at
        else
            return "There is no
            substring of T matching
            P".
    }
```

This algorithm consists of two nested loops, with the outer loop indexing through all possible starting indices of the pattern in the text, and the inner loop indexing through each character of the pattern, comparing it to its potentially corresponding character in the text.

In this algorithm, for each candidate index in T , we can perform m character comparisons to decide that P does not match T at the current index.

The outer-for-loop is executed $n-m+1$ times, and the inner loop is executed at most m times.

Thus, the running time of the brute-force method is $O((n-m+1)m)$ which is $O(nm)$.

When, $m = \frac{n}{2}$, this algorithm has quadratic running time $O(n^2)$.

The Naive string matching algorithm is same as Brute-force Algorithm.

Some time we say:— the algorithm is:—

Naive or Simple or Brute Force pattern Matching Algorithm.

This algorithm, always shifts the window by exactly one position to the right.

For example :- Naive / Simple / Brute Force
Pattern Matching

T = ANPANMAN

P = MAN

T:	A	N	P	A	N	M	I	M	I	A	N
P	M	A	N								
	x	M	A	N							
	x	M	A	N							
	x	M	A	N							
	x	M	A	N							
	x	M	A	N							

An example of Naive or simple
 or Brute Force Pattern
 Matching Algorithm.

3. Rabin-Karp String Matching Algorithm

In this algorithm, first a calculation is done before doing exact comparisons of characters.

Like, if :-

T is text $T[1..n]$

P is pattern $[1..m]$

assume any number q , (a prime number)

then compute $p \bmod q$

$$r = p \bmod q$$

Now, the substring of text T that is required to check in order to detect p, is first ~~so~~ involves calculation of mod with q .

like; if substring S_1 is to be checked in T,

then first compute $S_1 \bmod q$

if this value matches with r , that,

there is some chance to get the pattern only then the exact comparisons of characters

is done,

If $(s_i \text{ mod } q)$ does not match with r , then there is no chance of finding the occurrence of pattern p for s_i , so, the actual comparison is not done.

The next sequence to be checked is obtained by dropping the higher order bit or number in the text.

For example:-

Applying Rabin Karp string matching algorithm

(Assuming numerical values of characters)

T: 2 3 5 9 0 2 3 1 4 1 5 2 6 7 3 9 9 2 1

P = 3 1 4 1 5 , m = size of p = 5

q = 13 (given)

window of size $m=5$ to be checked

T: 1 2 3 5 9 0 2 3 1 4 1 5 2 6 7 3 9 9 2 1

first compute

$$r = p \bmod q$$

$$= 31415 \bmod 13$$

$$= 7 \quad (\text{is remainder})$$

$$\begin{array}{r} 2416 \\ 13) 31415 \\ \underline{26} \\ 54 \\ \underline{52} \\ 21 \\ \underline{13} \\ 8 \end{array}$$

7 8 17

thus,

check $S_1 = 23590$

$(S_1 \text{ mod } q) == r ?$

$23590 \text{ mod } 13 = r = 7 ?$

remainder
comes
as

$8 \neq 7 \therefore \text{No exact matching}$
 $\text{of characters will}$
 be done.

(Since, there is no chance of getting matching of pattern).

Compute next sequence to be matched:-

$$\begin{aligned} S_2 &= 10 \times (23590 - 10000 \times 2) + 2 \\ &= 10 \times (23590 - 20,000) + 2 \\ &= 10 \times 3590 + 2 \\ &= 35900 + 2 \\ &= 35902 \end{aligned}$$

Already checked S_1

T: $[23590] 23141526739921$

next window S_2
to be checked

$$S_2 = 35902 \text{ mod } 13 = 9 \neq 7 \quad \text{No match}$$

$$\begin{aligned} S_3 &= 10(35902 - 10,000 \times 3) + 3 \\ &= 10(35902 - 30,000) + 3 \\ &= 10(5902) + 3 \\ &= 59020 + 3 = 59023 \end{aligned}$$

$s_3 \bmod 9 ?$

$$59023 \bmod 13 = 3 \neq 7 \quad \text{no match}$$

removing higher order bit

$$\begin{aligned} s_4 &= 10(59023 - 10,000 \times 5) + 1 \\ &= 10(59023 - 50,000) + 1 \\ &= 10(9023) + 1 \end{aligned}$$

$$\begin{aligned} &= 90230 + 1 \\ &= 90231 \end{aligned}$$

window is滑动 at this new sequence

$$s_4 \bmod 9 ? \quad 90231 \bmod 13 = 11 \neq 7$$

$$\begin{aligned} s_5 &= 10(90231 - 10,000 \times 9) + 4 \\ &= 10[0231] + 4 \\ &= 02314 \end{aligned}$$

for removing higher order bit

This is obtained by looking next digit number in text

$$s_5 \bmod 13 = ? \quad 02314 \bmod 13 = 0 \neq 7$$

no match

0

0

0

$$s_7 = 31415$$

$$s_7 \bmod 9 ? = 31415 \bmod 13 = 7$$

Match found

$T = \boxed{2 \ 3 \ 5 \ 9 \ 0} \ 2 \ \boxed{3 \ 1 \ 4 \ 1 \ 5 \ 2 \ 6 \dots}$

s_1 s_2

$T = 2 \ 3 \ \boxed{5 \ 9 \ 0 \ 2 \ 3} \ \boxed{1 \ 4 \ 1 \ 5 \ 2 \ 6 \dots}$

s_3 s_4

$T = 2 \ 3 \ 5 \ 9 \ \boxed{0 \ 2 \ 3 \ 1} \ 4 \ 1 \ 5 \ 2 \ 6$

s_5

$T = 2 \ 3 \ 5 \ 9 \ 0 \ \boxed{2 \ 3 \ 1 \ 4 \ 1 \ 5} \ 2 \ 6$

s_6

$T = 2 \ 3 \ 5 \ 9 \ 0 \ 2 \ \boxed{3 \ 1 \ 4 \ 1 \ 5} \ 2 \ 6$

s_7

This method can be viewed as **sliding window of size m** which moves along the text array, T, and numbers in this window are matched with window containing the pattern.

The window slides along text till it get a window ~~with~~ with same modulo as that of pattern window.

If modulo is same but pattern is not matched, then it is called as **Spurious Hit** (or a false hit) (or a fake match of modulo as pattern is not found).

The calculations involved in finding next window
 e.g.

$$S_2 = 10 \xrightarrow{\text{To remove higher order bit}} (23590 - 10,000 \times 2) + 2 \xrightarrow{\text{↑ next digit}}$$

This is called as Horner's Rule.

Running time for calculating n different windows in array $T[1..n]$ is $O(n)$ and

for matching the pattern $O(m)$, so

the running time of Rabin-Karp method is

$$O(n+m)$$

HASHING in Rabin-Karp Algo:-

The Rabin-Karp algorithm involves the calculation of hash values (tickle) to convert text into a numerical value.

Following is the example, how Rabin-Karp algorithm is applied on the following text T and pattern P :-

T :- A A B A A C B A A

P :- B A A

T: AAB AAC B AA

P: BAA

Alphabets are:-

A, B, C, D, ..., Z
 Value = 1 2 3 4 ... 26

Total Alphabet size = $d = 26$

Assume a prime number (large value)
 (As large value ensures less collision) :-

Let, prime number = $q = 5381$

$n = \text{size of text} \Rightarrow 9$

$m = \text{size of pattern} = m_p$
 $\Rightarrow \text{BA A}$

AAB AAC B AA
 1 2 3 4 5 6 7 8 9 $\circlearrowright n$

∴ Window size to be checked = $m = 3$

~~Slow Hash Function~~ Computing Hash value of pattern

$B \rightarrow 2 \times 26^2$ $A \rightarrow 1 \times 26^1$ $A \rightarrow 1 \times 26^0$

$A-1$
 $B-2$
 $C-3$

$$[(2 \times 26^2) \bmod q + (1 \times 26^1)] \bmod q + [1 \times 26^0] \bmod q$$

$$= \left[\left[(2 \times 26^2) \bmod 5381 + (1 \times 26^1) \right] \bmod 5381 + (1 \times 26^0) \right] \bmod 5381$$

$$= \left[\left[(1352 \bmod 5381) + 26 \right] \bmod 5381 + (1 \times 1) \right] \bmod 5381$$

$$= \left[[1352 + 26] \bmod 5381 + 1 \right] \bmod 5381$$

$$= ((1378 \bmod 5381) + 1) \bmod 5381$$

$$= (1378 + 1) \bmod 5381$$

1379

Hash value = $q \Rightarrow p \bmod q$
for pattern

has been made stronger

\therefore here = hash value for pattern

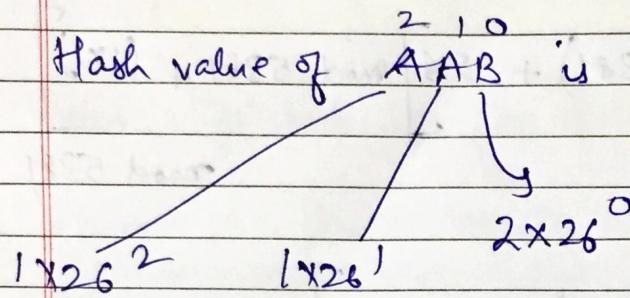
(pattern \Rightarrow BAA) $\rightarrow = q = 1379$

Now! - take window size = $m=3$,
 & start comparing:-

T: **AA B** A A C B A A
 $P = B A A = \boxed{1379}$ is 2

i.e. Hash value
 of pattern

Hash value of **AAB** is



A - 1
 B - 2
 C - 3
 :
 :

$$\therefore [((1 \times 26^2) \text{ mod } 5381 + (1 \times 26^1)) \text{ mod } 5381 \\ + (2 \times 26^0)] \text{ mod } 5381$$

$$= [(676 \text{ mod } 5381) + 26] \text{ mod } 5381 + (2 \times 1) \text{ mod } 5381$$

$$= [(676 + 26) \text{ mod } 5381 + 2] \text{ mod } 5381 \\ = (702 + 2) \text{ mod } 5381$$

$$= 704 \text{ mod } 5381 = 704$$

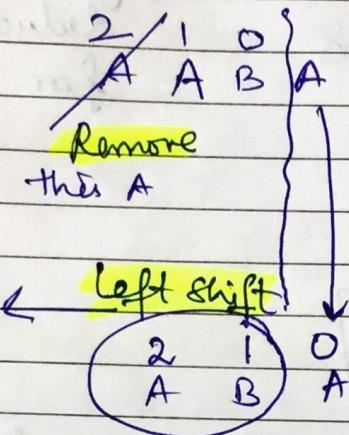
704 is not matching with 1379
 (i.e. 2)

match not found.

so, More the window,

T: A A B A C B A A
 P: B A A

In order to move the window:-



left shift AB and then add new A
 So, steps are:-

(1) Remove A A B

(2) left shift AB

left shift
 $\begin{array}{c} 2 \\ | \\ A B \end{array}$

(3) Add new character

Step 1
 remove A A B

$\begin{array}{c} 2 \\ | \\ A B \end{array}$ Add character

So,

(1×2^2) is deducted from hash value:-

$$\text{Hash value (AAB)} = 704$$

$$\therefore \text{Hash value (AAB)} - (1 \times 2^2) = 704 - 626$$

$$= 28$$

Step 2:-

left shift

$\begin{array}{r} 2 \ 1 \ 0 \\ \leftarrow A \ B \\ \text{involves } \downarrow \end{array}$ left shift

hash value of AB
 $= 28$

[got after
 deducting A
 from AAB]

to be left shifted

This is obtained by :-

$\begin{array}{r} 2 \ 1 \ 0 \\ \leftarrow A \ B \end{array}$

by multiplying it with base value $= d = 26$

$\frac{28 \times 26}{\uparrow}$
 base value $= d$

Note:-

If normal number is shifted

$\leftarrow 28$

it is multiplied by 10
 $\Rightarrow 28 \times 10$

ie 280

left shifted

Step 3 Now: $\begin{array}{r} 2 \ 1 \ 0 \\ \leftarrow A \ B \end{array}$ is 728

add new character

$\begin{array}{r} 2 \ 1 \ 0 \\ \leftarrow A \ B \ A \end{array}$ ← for this

add (1×2^0) to 728

$$\therefore 728 + (1 \times 2^0) = 728 + 1 = 729$$

But $729 \neq 1379$

So, match not found

\therefore move the window.

T: A A B A C B A A
P: B A A

and the process will go on, till we find the match or text is completely scanned and pattern is not there.