
Algorithms

14. String-Pattern Searching

Handwritten [by pankaj kumar](#)

String - pattern Searching

Date _____

Page No. _____

(91)

- ① Naive pattern Searching (192)
- ② Rabin Karp Algorithm (193)
- ③ KMP Algorithm (f) Longest prefix suffix (195)
- ④ Z-algorithm (198)
- ⑤ Longest palindromic substring (202)

① Naive Pattern Searching

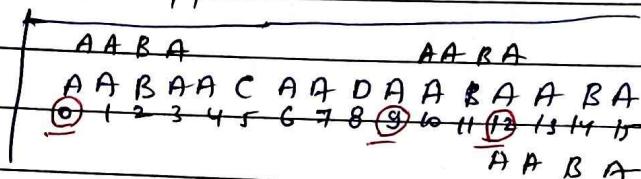
Given a text $txt[0..n-1]$ and a pattern $pat[0..m-1]$. write a function $search(\text{char} pat[], \text{char} txt[])$ that prints all occurrences of $pat[]$ in $txt[]$. you may assume that $n \geq m$.

Ex:- ① $txt[] = "THIS IS A TEST TEXT"$, $pat[] = "TEST"$

O/P:- 10 (pattern found at index 10)

② $txt[] = "AABAA CAADAA BAA BA"$, $pat[] = "AABA"$,

O/P:- 0, 9, 12



Note:- pattern searching is an important problem in computer science. when we do search for a string in notepad/word file or browser or db.

- **Naive Pattern Searching:** - The idea is to slide the pattern over text one by one and check for a match. If a match is found, then slide by 1 again to check for subsequent matches.

- i.e. check for the match of the first character of the pattern in the string, if it matches then check for the subsequent characters of the pattern with the respective characters of the string, if a match found then move forward in the string.

* code:-

```
void search(char* pat, char* txt) {
    int m = strlen(pat);
    int n = strlen(txt);
    for (int i=0; i<n-m; i++) {
        int j;
        for (j=0; j<m; j++)
            if (txt[i+j] != pat[j]) break;
        if (j == m) cout << i << endl;
    }
}
```

Tc: $O(n+m)$

[O(n)]

② **Best Case TC:** The best case occurs, when first character not present at $i \in [0, m \times (n-m+1))$

③ **Worst Case:** when all character of text and pattern are same. or only last character is different

④ $AAA \dots AAA = txt[]$

⑤ $AAA \dots AAA AAAB = pat[]$

(2) Rabin-Karp Algorithm for pattern Searching:-

- Like the Naive Algorithm, Rabin-Karp algorithm also slides the pattern one by one. But unlike the naive algorithm, Rabin-Karp algo. matches the hash value of the pattern with the hash value of current substring of text, and if the hash value match then only it starts matching individual characters. So Rabin-Karp algo. needs to calculate hash values for following string.

① Pattern itself. ② All the substring of the text of length m , i.e. equal to length of pattern string.

- Since we need to efficiently calculate hash values for all the substrings of size m of the text, we must have a hash function which has the following property

⇒ Hash at the next shift must be efficiently computable from the current hash value & next character in text or we can say that, hash($txt[s+1..s+m]$) must be efficiently computable from hash($txt[s..s+m-1]$) and $txt[s+m]$ i.e., $\text{hash}(txt[s+1..s+m]) = \text{rehash}(txt[s+m], \text{hash}(txt[s..s+m-1]))$ and rehash must be $O(1)$ operation.

⇒ The hash function suggested by Rabin-Karp calculates an integer value. the integer value for a string is numeric value of a string. for ex, if all possible characters are from 1 to 10, the numeric value of "22" will be 122. the no. of possible characters is higher than 10 (256) and pattern length can be large. so numeric value cannot be practically stored as an integer, therefore, the numeric value is calculated using modular arithmetic to make sure that the hash values can be stored in an integer variable (can fit in memory words). To do rehashing, we need to take off the most significant digit and add the new least significant digit for in hash value.

$$\text{hash}(txt[s+1..s+m]) = (d[\text{hash}(txt[s..s+m-1] - txt[s]*q)] + txt[s+m]) \bmod q.$$

here, d = no. of character in alphabet

q = A prime number

$h = d^{(m-1)}$ or $d^m - 1$

ex:- "1223"

$\text{hash}("122") = 122$

$d=10$ $10^3 = 1000$

$\text{hash}("1223") = 1 * (10^3 * 1) + 2 = 1223$

* Code

```

#include <bits/stdc++.h>
using namespace std;
#define d 256

void search (char pat[], char txt[], int q) {
    int m = strlen (pat);
    int n = strlen (txt);
    int i, j, p = 0, t = 0, h = 1;

    // hash value for pattern.
    for (i = 0; i < m - 1; i++)
        h = (h * d) + pat[i];
    h = h * q; // hash value for txt

    for (i = 0; i < m; i++) {
        p = (d * p + pat[i]) * q; // calculate the hash value of
        t = (d * t + txt[i]) * q; // pattern and first window
        if (p == t) { // if hash value of current window & pattern
            for (j = 0; j < m; j++) { // match them only check for
                if (txt[i + j] != pat[j]) // next character one by one
                    break;
            }
            if (j == m) { // if the pattern found
                cout << "Pattern found at index " << i + 1;
            }
        }
        if (i <= m - 1) { // calculate hash value for next window
            t = (d * (t - txt[i] * h) + txt[i + m]) * q; // digit
            if (t < 0) t = t + q; // we might get negative value of t,
                                    // so convert it to positive.
        }
    }
}

```

int main() {

char txt[] = "GEEKS FOR GEEKS";

char pat[] = "GEEK";

int q = 203;

search (pat, txt, q);

return 0;

* average time complexity.

is $O(n+m)$

* but worst case is $O(nm)$

when all character of pattern & text are same & the hash value of all substring of txt with hash value of pat[].

txt[] = "AAAAAABBBBBB"
pat[] = "AAAAA"

(3) KMP algorithm for pattern searching.

- the TC of Naive pattern searching and Rabin-Karp for searching patterns. The worst case complexity of both of the algorithms is $O(n+m)$.
 - The TC of KMP in worst case is $O(n)$.
- ④ KMP (Knuth Morris Pratt) pattern searching.
- The KMP matching algo. uses degenerating property (pattern having the same sub-pattern appearing more than once in the pattern) of the pattern and improves the worst case complexity to $O(n)$. The basic idea behind KMP's algo is: whenever we detect a mismatch (after some matches), we already know some of the character in the text of the next window. We take advantage of this information to avoid matching characters that we know will anyway match.

⑤ Matching Overview

- $\text{txt} = \text{"AAAAAABAAAABA"}$, $\text{pat} = \text{"AAAAA"}$
we compare first window of txt with pat
 $\text{txt} = \text{"AAAAAABAAAABA"}$, $\text{pat} = \text{"AAAAA"}$ (we find a match)
- In the next step we compare next window of text with pat
 $\text{txt} = \text{"AAAAA(A)BAAAABA"}$
 $\text{pat} = \text{"AAAA(A)"}$ [pattern shifted one position]

Note:- here we don't need to compare whole string. Here KMP does optimization over Naive. Here, we only compare **Firstly A** of pattern with four character of current window of text to decide whether current window matches or not. Since we know **First three characters will anyway match**, we skipped matching first three characters.

⑥ Need of preprocessing

here, a question arises how to know how many characters to be skipped. To know this, we pre-process pattern and prepare an integer array LPS[] that tells us the count of characters to be skipped.

- LPS[] indicates longest proper prefix which is also suffix.

For ABC:- prefix is "", "A", "AB" & "ABC". Proper prefix is "A" & "AB", "", suffix are: "", "C", "BC", & "ABC".

Preprocessing overview :-

- construct an auxiliary $lps[i]$ of size m (same as size of pattern) which is used to skip characters while matching
- we search for lps in sub-patterns. more clearly we focus on sub-strings of patterns that are either prefix and suffix
- For each sub-pattern $pat[0..i]$ where $i=0$ to $m-1$, $lps[i]$ stores length of the maximum matching proper prefix which is also a suffix of the sub-pattern $pat[0..i]$.

$lps[i] = \text{the longest proper prefix of } pat[0..i]$
which is also a suffix of $pat[0..i]$.

Note:- $lps[i]$ could also be defined as longest prefix which is also proper suffix.

Ex:- $lps[]$ construction:

i) $pat = "AAAAA"$, $lps[] = [0, 1, 2, 3]$

ii) $pat = "ABCDE"$, $lps[] = [0, 0, 0, 0, 0]$

iii) $pat = "AABAACAAABAA"$, $lps = [0, 1, 0, 1, 2, 0, 1, 2, 3, 4, 5]$

A $\xrightarrow{\quad}$ AA $\xrightarrow{\quad}$ AAB $\xrightarrow{\quad}$ AABA $\xrightarrow{\quad}$ AA $\xrightarrow{\quad}$ BAA $\xrightarrow{\quad}$ AABA $\xrightarrow{\quad}$ AABAACAAABAA

iv) $pat = "PAAACAAAPAC"$, $lps = [0, 1, 2, 0, 1, 2, 3, 2, 3, 4]$

* Searching algorithm / How to use $lps[]$

- we start comparison of $pat[i]$ with $j=0$ with character of current window of text

- we keep matching characters $txt[i]$ & $pat[i]$ and keep incrementing $i+j$ while $pat[i]$ & $txt[i]$ keep matching.

- when we see a mismatch.

• we know that characters $pat[0..j-1]$ match with $txt[i..j-1]$

• since, we know $[lps[j-1]]$ is count of character of $pat[0..j-2]$ that are both proper prefix and suffix.

• from above two points, we can conclude that we do not need to match $lps[j-1]$ character with $txt[i-1..i-1]$ because known that these character will anyway match.

ex:- pat = "AAAAA", lps = [0, 1, 2, 3]

txt = "AAAAA B A A A B A"

In this situation we

don't need to match pat[0]

with txt[i], we will match

pat[lps[i]] with txt[i] i.e. match pat[0] with txt[i]

∴ txt = AAAA A B A A A B A

* Code :-

```
#include <bits/stdc++.h>
```

```
using namespace std;
```

```
void computeLPSArray (char* pat, int m, int* lps);
```

```
void kmpSearch (char* pat, char* txt) {
```

```
int m = strlen (pat);
```

```
int n = strlen (txt);
```

```
int lps[m];
```

```
computeLPSArray (pat, m, lps);
```

```
int i=0, j=0;
```

```
while (i < n) {
```

```
if (pat[j] == txt[i]) {
```

```
    j++; i++;
```

```
if (j == m) {
```

```
cout << " found pattern at index " << i - j + 1 << endl;
```

```
j = lps[j-1];
```

```
else if (i < n && pat[j] != txt[i]) {
```

```
if (j != 0)
```

```
j = lps[j-1];
```

```
else i = i+1;
```

```
i = i+1;
```

```
void computeLPSArray (char* pat, int m, int* lps) {
```

```
int clen=0; lps[0]=0; i=0; j=1; lps[0]=0;
```

```
lps[0]=0;
```

```

int i = 1;
while (i < m) {
    if (pat[i] == pat[i + len]) {
        len++;
        if (i + len == n)
            i += len;
    } else {
        if (len != 0) {
            len = lps[len - 1];
        } else {
            lps[i] = 0;
            i++;
        }
    }
}

```

④ Creating LPS Array

longest prefix suffix
O(m)

```

int main() {
    char txt[] = "ABABDABACDABABCABAB";
    char pat[] = "ABABCABAB";
    kmpSearch(pat, txt);
    return 0;
}

```

④ Z-algorithm :-

- This algorithm finds all occurrences of a pattern in text in linear time. Let length of text be n and of pattern be m , then total time taken is $O(n+m)$ with linear space.
- Since KMP & Z have both same time complexity as linear but Z-algorithm is simpler to understand.
- In this algo, we construct a Z array.

⑤ Z-array :-

- Z-array is of same length of string. An element $Z[i]$ of Z array stores length of the longest prefix of string from $str[0..n-1]$ which is also a prefix of $str[0..n-1]$. The first entry of Z array is meaningless as complete string is always prefix of itself.

Ex:- Index

0	1	2	3	4	5	6	7	8	9	10	11
① Text	a	a	b	c	a	a	b	x	a	a	z
2-values	x	1	0	0	3	1	0	0	2	2	1

② str = "aaaaaa", z[] = {x, 5, 4, 3, 2, 1}

③ str = "aabaaacd", z[] = {x, 1, 0, 2, 1, 0, 0}

④ str = "abababab", z[] = {x, 0, 6, 0, 4, 0, 2, 0}

⑤ How z-array used to search pattern in Linear time?

- The idea is to concatenate pattern and text, and create a string "PFT" where P is pattern & F is special character should not be present in pattern and text, and T is text.

Build the z-array for concatenated string. In z-array if z-value at any point is equal to pattern length, then pattern is present at that point.

Ex:-

pattern P = "aab", Text T = "baabaa"

The concatenated string is = "aabbaabaa"

z-array for concatenated string is {x, 1, 0, 0, 0, 1, 3, 1, 0, 2, 1}.

here, the value of ~~over~~ 3 in z-array indicates present of pattern

⑥ To construct z-array

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
a	a	b	x	a	a	b	a	c	a	b	x	a	a	b	a	b	a	z

z-array

For index 5 we don't need to compare since the pattern is same which is at index ① ~~bx~~

and ~~abx~~ b₀ 1 2 3 = 4 5 6 7

∴ it is obvious ~~23 = 56~~

similarly 6 & 7 will be equal to z[2] & z[3] i.e. 0, 0

before using ~~1~~ of index 13 by comparing with index 2[4] which is 4, we need

to calculate if ~~13 + 4~~ should be lies in the range ~~10-12~~ of ~~9~~ since ~~13 + 4~~ is 17 and

it is not in the range so again we have to compare from S[17] and S[18] and change ~~17~~ to 16.

Q) Code:-

```
#include <bits/stdc++.h>
using namespace std;

void getzarr(string str, int zarr[]);
void search(string text, string pattern) {
    string concat = pattern + "$" + text;
    int l = concat.length();
    int zarr[l];
    getzarr(concat, zarr); // construct zarray
    for (int i = 0; i < l; ++i) {
        if (zarr[i] == pattern.length()) // IF z[i] equal to pattern
            cout << "pattern found at index " << i - pattern.length() + 1 << endl;
    }
}
```

Time complexity :- O(n+m)

```
void getzarr(string str, int zarr[]); // 2-array construction
int n = str.length(); // code
int l = 0, r = 0, k;
```

// [l, r] make a window which matches with prefix of string (concatenate)

```
for (int i = 1; i < n; ++i) {
    if (i > r) // if nothing matches so we will
    L = R = i; // calculate z[i] using naive way.
```

while (R < n && str[R - L] == str[R]) // if R-L=0

R++;

z[i] = R - L;

// starting, so it will start
// checking from 0th index,

else of

k = i - L; // k corresponds to no. which matches in [L, R]

if (z[i] < R - i + 1) // if z[i] is less than remaining interval

z[i] = z[k]; // then z[i] will be equal to z[k].

else

L = i; // if out of interval then start from R and

while (R < n && str[R - L] == str[R]) // check manually,

R++;

z[i] = R - L;

5. Longest Palindromic Substring

Given a string s , return the longest palindromic substring in s .

ex:- I/p: $s = "babab"$
O/P: $"bab"$ or $"aba"$

I/p: $s = "bb"$
O/P: \underline{bb}

③ $s = "abc"$
ans = \underline{a} or \underline{b} or \underline{c}

* Two pointer Approach :- (code)

① $\rightarrow abcdefedcg$

(while($i < n$)) $\leftarrow \dots i \leftarrow k \rightarrow$ while($k < n$) \rightarrow do this
for every $i \in (0, n)$

② $\rightarrow abcdefeffedcg$

(while($i < n$)) $\leftarrow \dots i \leftarrow k \rightarrow$ while($k < n$) \rightarrow do this
for every $i \in (0, n)$
skip similar character

string longestPalindrome(string s){

if ($s.empty()$) return "";

if ($s.size() == 1$) return s;

int start = 0, len = 1;

for (int i = 0; i < s.size();) {

if ($s.size() - i <= len / 2$) break;

int j = i, k = i;

while ($k < s.size() - 1$ && $s[k] == s[k + 1]$) ++k;

i = k + 1;

while ($k < s.size() - 1$ && $s[j + 1] == s[k + 1]$)

{
 ++k;
 --j;
}

int temp = k - j + 1;

if ($temp > len$) {

start = j;

len = temp;

{

return s.substr(start, len);

} // Expand palindrome

max length palindrome
is already calculated

skip duplicate
character

T.C: ~~0(n^2)~~ $1 + 2 + 3 + 4 + 5 + \dots + \frac{n}{2} + \frac{n}{2} - 1 + \dots + 3 + 2 + 1 = \frac{2n(n+1)}{2} =$

$O(n^2)$