

String Algorithms

(1) Naive Algo.

- Pattern = p } To find
- Text = t } To search in

$$m = \text{length}(\text{Pattern})$$

$$n = \text{length}(\text{Text})$$

Assume $n > m$

```

Algo:
for (int i=0; i < n-m; i++)
    int count = 0;
    for (int j=0; j < m; j++)
        if (p[j] == text[i+j])
            count++;
    if (count == m)
        return true; // found at index i
    }
return false;

```

(2) Time Complexity

Best case

$$O(n)$$

Worst case

$$O(m \times (n-m+1))$$

→ When not found & other first character is not present

→ When only last character is diff rest all are same
 - P = AAAAAAAAAA, P = AAB

② Optimised Naive

→ In this if the characters are diff in the pattern to search then instead of iterating 1 by 1 we can directly jump " M " [length of pattern] times.

Algo:

$P = \text{Pattern}$

$T = \text{Text}$

$M = \text{length}(P)$

$N = \text{length}(T)$

[Assume: $n > m$]

Ex:

$T = \text{ABCEABCEABCEABCE}$

$P = \text{ABCD}$

$M = 4$

$N = 12$

```
while (i <= N - M) {  
    for (int j = 0; j < M; j++)  
        if (T[i + j] != P[j])  
            break;  
    if (j == M)  
        cout << "pattern at index: " << i << endl;  
    else if (j == 0)  
        i = i + 1;  
    else { i = i + j; } // jump by j  
}
```

* Complexity → Worst case: Some.

* Rolling hash : Gives you the ability to calculate the hash value without rehashing whole string.

Ex: ~~hash(1234) = hash(123)~~

Page No.			
Date			

③ Robin Karp Algo. : It is better than naive as it first calculates the hash value of the substring of length 'm' (length of pattern) & only if it matches then ~~we~~ compute compare the string.

Time - complexity

(a) Avg. case : $O(n+m)$

(b) Worst case : $O(nm)$: When hash value

↳ P = AAA matches.

↳ T = AAAAAAAAA

Based on rolling-hash technique

↳ string = ~~abc~~ ~~1234~~ abc

~~5678901234~~ =

→ We need rolling hash so that we can fastly & efficiently calculate hash.

Ex: String = abcd. (7ent)
 Pattern = abcd (pattern)

so we need to create all possible strings
 of length 3 (length of pattern) &
 calculate their hash value.

#define ll long long int

ll createHashValue (string str, int n)
 {
 ll result = 0;
 for (int i = 0; i < n; i++)
 result += (ll) str[i] * (ll) pow(prime, i);
 return result;
 }
 // Prime no : say : 3,

O(n)
 complexity

// Rolling hash code down.

ll recalculateHash (string str, int oldIndex, int newIndex,
 ll oldhash, int patLength)

{
 ll newhash = oldhash - str[oldIndex];
 newhash = newhash / prime;
 newhash += (ll) (str[newIndex] * (ll) (pow(prime, patLength - 1)));
 return newhash;
 }

O(1)
 complexity

In rolling hash we simply ~~have~~

Exmpl: patt length = 3, string = 23456

hash value for first window '234' = 234.

value of next window '345' = $(234 - 2 \times 100) \times 10 + 5$

→ The algo we have used works as:

string = abcd, prime = 5, length = 3

$$\text{hash}(abc) = 97 \times 5^0 + 98 \times 5^1 + 99 \times 5^2$$

↓

$$\text{hash}(bcd) = \frac{\text{hash}(abc) - 97 \times 5^{\text{prime}}}{5} + 100 \times 5^{\text{pattern length} - 1}$$

$\xrightarrow{5 \rightarrow \text{prime}}$
 $\xrightarrow{\text{first}} \xrightarrow{\text{AICI of d}}$
 $\xrightarrow{2}$

Actual
String
Comparison

```
bool checkEqual (string text, string pattern,
                 int s1, int e1, int s2, int e2)
```

```
{
    if (e1 - s1 != e2 - s2)
        return false; // length mismatch
```

```
    while (s1 <= e1 && s2 <= e2) {
```

```
        if (text[s1] != pattern[s2])
```

```
            return false;
```

```
        s1++;
```

```
        s2++;
```

```
    }
```

```
    return true;
```

```
}
```

// DRIVER CODE

```
int main ()
```

```
{
```

```
    string str = "obobc obd ob";
```

```
    string patt = "obd";
```

```
    ll patHash = createHashValue (patt, patt.length());
```

```
    ll strHash = createHashValue (str, str.length());
```

```
    for (int i=0; i<= str.length() - patt.length(); i++)
```

```
    {
        if (patHash == strHash &&
```

```
            checkEqual (str, patt, i, i + patt.length() - 1,
                        0, patt.length() - 1))
```

```
            cout << "Match at : " << i << endl;
```

```
            return 0;
```

// Only for run-time check on boundaries of length

```
    if (i < str.length() - patt.length())
```

```
    {
```

```
        strHash = calculateHash (str, i, i + patt.length()
```

```
                                strHash, patt.length());
```

```
    }
```

```
    return 0;
```

```
}
```


→ KMP also acts smart whenever we detect a mismatch i.e. when a mismatch occurs it knows already some of the characters in text window & it tries to skip them.

* In the reset table we will store the value of length of longest suffin present till now which is also present as a prefix in text.

Reset table

A	B	C	X	A	B	X

→ In this what we basically do is that till any point that we have matched in the string (pattern) does there exist any suffix which is also a prefix, if yes then we start to compare from character after the prefix match.

* It helps us to ~~do~~ reduce the ~~no~~ total number of comparisons by considering ^{or-occurring} ~~an~~ sub-patterns.

NOTE:

Creating the Reset table is the imp part of KMP as it ~~is~~ only helps to reduce the complexity.

ALGO - Reset Table Creation.

```
void createResetTable(int *arr, int n, char *pat)
// length of "arr" is 'n' i.e. length of pattern.
```

```
{
    int i=1, j=0; // Initiative
    arr[0] = 0; // Always, b'coz nothing before it.
```

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

```
    {
```

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

```
        {
            arr[i] = ++j;
```

```
            i++;
```

```
        }
```

```
    else // When the characters don't match
```

```
    {
```

```
        // When j is still equal to 0.
```

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

```
        {
```

```
            arr[i] = 0;
```

```
            i++;
```

```
        }
```