# String Patterns: Rabin Karp

# Rabin Karp

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

- Pattern itself
- All the substrings of the text of length m which is the size of the pattern.

## *How is Hash Value calculated in Rabin-Karp?*

Hash value is used to efficiently check for potential matches between a pattern and substrings of a larger text. The hash value is calculated using a rolling hash function, which allows you to update the hash value for a new substring by efficiently removing the contribution of the old character and adding the contribution of the new character. This makes it possible to slide the pattern over the text and calculate the hash value for each substring without recalculating the entire hash from scratch.

Here's how the hash value is typically calculated in Rabin-Karp:

Step 1: Choose a suitable base and a modulus:

- Select a prime number 'p' as the modulus. This choice helps avoid overflow issues and ensures a good distribution of hash values.
- Choose a base 'b' (usually a prime number as well), which is often the size of the character set (e.g., 256 for ASCII characters).

Step 2: Initialize the hash value:

- Set an initial hash value 'hash' to 0.

Step 3: Calculate the initial hash value for the pattern:

- Iterate over each character in the pattern from left to right.
- For each character 'c' at position 'i', calculate its contribution to the hash value as 'c * (bpattern_length – i – 1) % p' and add it to 'hash'.
- This gives you the hash value for the entire pattern.

Step 4: Slide the pattern over the text:

- Start by calculating the hash value for the first substring of the text that is the same length as the pattern.

Step 5: Update the hash value for each subsequent substring:

- To slide the pattern one position to the right, you remove the contribution of the leftmost character and add the contribution of the new character on the right.
- The formula for updating the hash value when moving from position 'i' to 'i+1' is:

hash = (hash - (text[i - pattern_length] * (bpattern_length - 1)) % p) * b + text[i]

Step 6: Compare hash values:

- When the hash value of a substring in the text matches the hash value of the pattern, it's a potential match.
- If the hash values match, we should perform a character-by-character comparison to confirm the match, as hash collisions can occur.

# Find count of A inside B

Problem 1's

Q> Given a string A (of length N) & a string B (of length M) find the count of substrings of A which are equal to B

eg. A = "abcbdcacba"    ans ≡ 1
    B = "cba"

A₂ = "abcaba bac"    ans ≡ 2
B = "aba"

① Bruteforce's
   $\leftarrow$ (N-m+1)
   + substring & A of length m, check if P+ is B.
   (sliding window)    O(m)
       T.C ≡ O(N-M+1) × O(M)
          ≡ O(N*M)

② Using Hashing + s.w

   a → 1    ⓪ h('abc') ≡ sum & value & characters  ✗ wont work!
   b → 2                                               two strings can
   c → 3                                               have same val
   :
   2 → 26   ② Convert to number with base next
            to 26 which is prime (P ≡ 29)
            f: h(aba) ≡ (1*p² + 2*p' + 1*p⁰)    ( P ≡ 29 )

taking base ≡ 10

A = "abca babac"
B = "aba" → h(aba) → 1*10² + 2*10' + 1*10⁰ = 121
h: 1*10² + 2*10' + 3*10⁰ ≡ 123    ✗

using sliding
calculate h:
   h₂ ≡ [[(1*10² + 2*10' + 3*10²) - 1*10²] + 1*10²]

   T.C ≡ O(M+N) ≡ O(N)              keep checking hash
              form  for N elements only s.w    with key hash

   S.C ≡ O(1)

edge case   if length of string is large ⇒ h() ≡ over flow
                                            ≡ % M
                                              (10⁹ + 7)
Rolling hash function
   h(s) = (∑ᴺ⁻¹ (s[i] * pᴺ⁻¹⁻ⁱ)) % M
          i=0

But as we reduced the size by % M, there can be a chance
of collision.

## Probability of Collision

| Strings | $h(s)$ | Probability of Collision |
|---|---|---|
| $S_1$ | | $0$ |
| $S_2$ | | $1/M$ |
| $S_3$ | | $2/M$ |
| $\mid$ | | |
| $\mid$ | | |
| $\mid$ | | |
| $S_N$ | $h(s_N)$ | $N - 1/M$ |
| $S_{N+1}$ | $h(s_{N+1})$ | $N/M \equiv 10^5/10^9 \equiv 0.0001$ |
| | | $\quad\quad \hookrightarrow (10^9 + 7)$ small probability |

③ Using Probabilistic Approach

$A \equiv$ "abcababac"   $B =$ "aba"

$(1*b^2 + 2*p + 1*p^0) \% M$

$(1*b^2 + 2*p^1 + 3*p^0) \% M$

Case1 : $(h(s_1) \mathrel{!=} h(s_2)) \Rightarrow s_1 \mathrel{!=} s_2$

Case2 : $h(s_1) == h(s_2) \Rightarrow s_1 == s_2$

$\quad\quad \hookrightarrow (s_1 \equiv s_2)$ with high probability. but not 100%

$\quad\quad\quad \hookrightarrow$ Do char by char comparision to
verify if $[s_1 == s_2]$ cost $\equiv O(M)$

Total T.C $\equiv O(N)$ best Case.

$\quad\quad\quad O(N*M)$ worst Case.

# Boring Substring

**Problem**
Given a string check whether it is possible to rearrange the characters s.t there is nt boring substring in S.

Boring Substring → ( len ≡ 2) f Consecutive chars.
eg: ab, bc, xy, yz, de, ----
[a , z] → 26 chars.

eg:-
"abc" false
"abcd" → cadb    True
           bdac

① Brute force
Check each permutations O(N!)

②
"aabccdb" ⟶ "cc aadbb"        " aab cc b"
      / ‾‾ divide with logic split
  |aacc| |bdb|                |aacc| |bb|
  Reverse and concatenate to get new string
      |ccaa bdb|                        false
                                X will not work

Correct Approach

Step1 :- Split the characters into two sets based on odd/even
         f keep track of smallest/largest.

Step2 :- Try each at (s₁ , s₂), (s₁, l₂), (l₁, s₂)
              If valid ⟶ True
              else ⟶ false.

       T.C : O(N)
       S.C : O(1)