

STRING PATTERN MATCHING

"Believe you can and you're halfway there."

~ Theodore Roosevelt



BRIGHT
DROPS
.com



Today's content

01. Boring substring
02. Pattern searching (Rabin Karp Algo)
 - Probability of collisions

Q. Given a string s ($\text{char}[]$), check whether it is possible to rearrange the characters such that there is no boring substring in s .

Boring substring \rightarrow length = 2 & consecutive alphabetically

eg: ab, cd, ed, zy, dc, bc

eg: $s = "abc" \rightarrow \text{false}$

$s = "abcd" \rightarrow \text{true} \rightarrow bdac \text{ or } cadb$

$s = aabccb \rightarrow \text{false}$

$s = "aabccdb" \rightarrow \text{True} \Rightarrow ccaadbbb$

Total no. of permutations for length $n = n!$

Brute force \rightarrow Generate each and every permutation of string length n & then check if it is forming any boring substring or not.

$$TC = O(n! * n)$$

All possible alphabets



odd chars

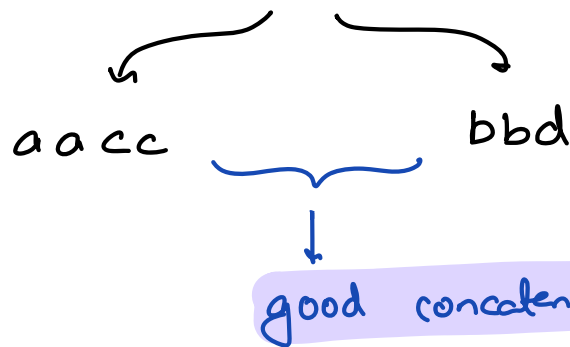
even chars

{ a c e g i ... w y }

{ b d f h ... x z }

Obs → Splitting the string in odd-even can give two sets which will never form boring substring

s = "a a b c c d b"

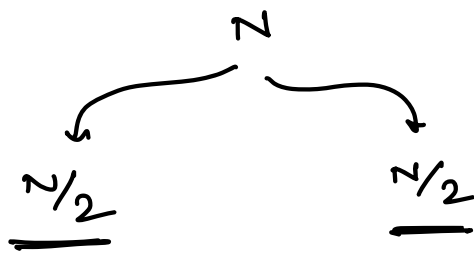


→ if you can find an alphabet from LHS & an alphabet from RHS which when concatenated don't form a boring substring

a c c a d b b → True

* For concatenation

01. Check every character from LHS & every character from RHS



$$TC = O(N^2)$$

02 Use hashMap

→ Only check for unique character

LHS = { 13elb }

RHS = { 13el 3 }

13 unique

13 unique

$$TC = O(13 * 13) = O(1)$$

03

Odd set →

Min

max

Even set →

min

max

s = "ccebdf"

cce

bdf

mino = c

maxo = e

mine = b

maxe = f

Min-Min

ce cb df X

Min - max

ce cf db \rightarrow True

Max - min

cc eb df \rightarrow True

max - max

c c ef d b X

Q String s = "acefdd"

Min odd = a

max odd = e

min even = d

max even = f

TC = $O(n)$

SC = $O(1)$

min odd X min even = ad \rightarrow True

min odd X max even

max odd X min even

max odd X max even

String s = " - - - - - " ;

$\left\{ \begin{array}{l} \text{minO} = \infty \\ \text{maxO} = -\infty \\ \text{mine} = \infty \\ \text{maxe} = -\infty \end{array} \right.$

for (i = 0 ; i < n ; i++) {

 char ch = s.charAt(i);

 if (ch % 2 == 0) {

 mine = Math.min (mine , ch);

 maxe = Math.max (maxe , ch);

 }

 else {

 minO = Math.min (minO , ch);

 maxO = Math.max (maxO , ch);

 }

}

if (minodd - mineven) > 1 return true

if (minodd - maxeven) > 1 return true

if (maxodd - mineven) > 1 "

if (maxodd - maxeven) > 1 "

→ abs value

Pattern matching

02 Given a large text (string A of length N) & small pattern (string B of length M). Find the **count of times B is present as a substring in A.**

Note :- $M < N$ & all lowercase characters

A = "abc**bc**abca" Ans = 2

B = "bca"

A = "abab**ca**babac" Ans = 3

B = "aba"

A = "abc**ab**abac" Ans = 2

B = "aba"

Brute force → For all substring of length M in A,
check if it is equal to B or not.

Sliding window X

No. of substrings of length M = $n - m + 1$

arr[n] =

--	--	--	--	--	--

Subarr = m size

A = "abcababac"

B = "aba"

→ sliding window when comparing strings is of
no use.

$$TC = O(n-m+1) * m \rightarrow O(n \cdot m)$$

Problem → Comparing two strings takes linear time

Solve → Two integers takes $O(1)$ time

Robin Karp Algorithm

B = "acd"

$$h("acd") = 1 + 3 + 4 = 8$$

a → 1

b → 2

c → 3

⋮

z → 26

A = "bacdef"

$$B = "acd" \rightarrow 1 + 3 + 4 = 8$$

$$\text{First window} = "bac" = 2 + 1 + 3 = 6$$

if $(h(s_1) \neq h(s_2)) \rightarrow \text{implies } s_1 \neq s_2$

Second window = "acd" = $1 + 3 + 4 = 8$

if $(h(s_1) == h(s_2)) \rightarrow \text{imply } s_1 = s_2 \quad \times$

$$\left. \begin{aligned} h(acd) &= h(adc) \\ &= h(cad) \\ &= h(cda) \\ &\vdots \end{aligned} \right\} = 8 \quad \begin{array}{l} \text{Multiple inputs} \\ \text{are producing the} \\ \text{same hash value} \end{array}$$

Collision

Avoid
Collision \rightarrow Weightage to the order

$$\begin{aligned} 326 &= 3 * 10^2 + 2 * 10^1 + 6 * 10^0 \\ 263 &= 2 * 10^2 + 6 * 10^1 + 3 * 10^0 \end{aligned} \quad \left. \vphantom{\begin{aligned} 326 \\ 263 \end{aligned}} \right\} \begin{array}{l} \text{Order has some} \\ \text{weightage} \end{array}$$

$$\begin{aligned} h(\text{"acd"}) &= a * 10^2 + c * 10^1 + d * 10^0 \\ &= 1 * 10^2 + 3 * 10^1 + 4 * 10^0 \\ &\Rightarrow 134 \end{aligned}$$

$$h("acd") = a * p^2 + c * p^1 + d * p^0$$

$P=1 \rightarrow$ collision

$P=2 \rightarrow$ x

Eg:- $S_1 = \begin{matrix} 2^1 & 2^0 \\ \downarrow & \downarrow \\ "aa" \end{matrix} = 1 * 2^1 + 1 * 2^0 \Rightarrow 3$

$S_2 = \begin{matrix} "c" \\ \downarrow \\ 2^0 \end{matrix} = 3 * 2^0 = 3$

Usually $P=29$ \rightarrow prime no. are always preferred in the hash fn
 \downarrow
 Min collision

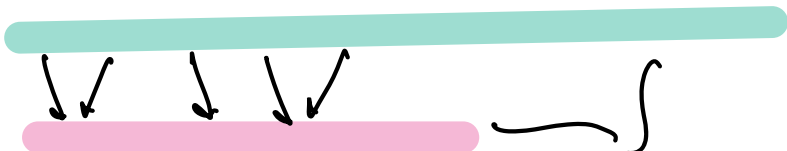
Limitation

$h("yogeshyadau") \rightarrow$ impossible to store it in an integer or long (Integer overflow)

To make sure, we are always in range \rightarrow Mod

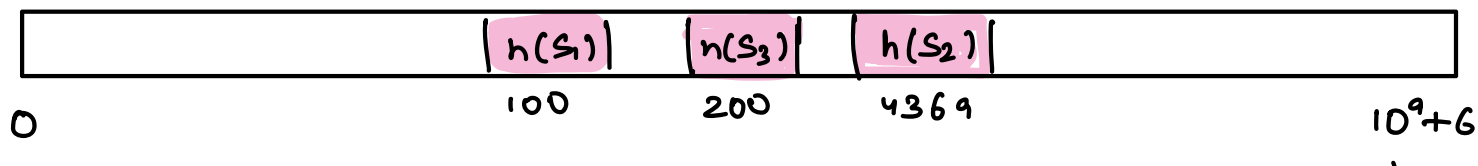
Usually $m = 10^9 + 7$

large prime no.



pigeon hole \rightarrow Collisions

Probability of collision



<u>String</u>	<u>$h(s)$</u>	<u>Probability of collision</u>
s_1	$h(s_1) = 100$	0
s_2	$h(s_2) = 4369$	$\frac{1}{M}$ or $\frac{1}{10^9+6}$
s_3	$h(s_3) = 200$	$\frac{2}{M}$
s_4	$h(s_4)$	$\frac{3}{M}$
\vdots	\vdots	\vdots
s_n	$h(s_n)$	$\frac{(n-1)}{m}$
s_{n+1}	$h(s_{n+1})$	$\frac{n}{m}$

$$m = 10^9 + 7 \rightarrow \text{large}$$

$$\text{Probability of collision} = \frac{n}{m} = \frac{10^5}{10^9} \approx 0.0001$$

↓
neglect the

Hash function

$$h(s) = \sum_{i=0}^n s[i] * p^{n-1-i} \% M$$

$$M = 10^9 + 7$$

$$p = 29$$

Rolling hash fn

B = "acd"

A = "bacd..."

$$\begin{aligned} h(B) = h("acd") &= 1 * 29^2 + 3 * 29 + 4 * 1 \\ &= 932 \end{aligned} \rightarrow TC = O(M)$$

$$\begin{aligned} h("bac") &= 2 * 29^2 + 1 * 29 + 3 * 1 \\ &= 1714 \end{aligned}$$

$$\text{if } (h(s_1) \neq h(s_2)) \rightarrow s_1 \neq s_2$$

$$h(\text{"bac"}) = 2 * 29^2 + 1 * 29 + 3 * 1$$

slide window

$$- \left[2 * 29^2 + 1 * 29^1 + 3 * 29^0 - 2 * 29^2 \right] * 29 + 4$$

↓
Remove first
↓
update place value of other chars
→ add next char

$$\left[(2 * p^2) + (1 * p^1) + (3 * p^0) - (2 * p^2) \right] * p + 4 * p^0$$

if $(h(s_1) == h(s_2)) \rightarrow 99.99\%$ of time
 $s_1 = s_2$

To get 100% accurate result

if $(h(s_1) == h(s_2)) \rightarrow$ compare both strings
 char by char

Steps

01. Calculate the hash value \rightarrow String B $\Rightarrow TC = O(m)$
02. Calculate the hash value $\rightarrow TC = O(m)$
for the first window
03. Sliding window $\rightarrow TC = O(n - m + 1)$

Overall $TC = O(n)$
 $SC = O(1)$

Code \rightarrow {TODO}

— α — α — α —