

★ Rabin karp. String matching algorithm

It is an algorithm that uses hashing to find patterns in strings. Let n be the length of the Text T and m be the length of the Pattern P . We need to find at which all positions in T , Pattern P exists.

⇒ Rabin karp method calculates hash value of pattern P .

⇒ It also calculates hash value for a substring of length m from T

⇒ If hash values are unequal algorithm will ~~compare~~ calculate the hash value for next m character sequence.

⇒ if the hash values are equal, algorithm will compare the pattern and the

m Character sequence.

⇒ In this way there is only one comparison per text subsequence & character matching is only needed when hash values match.

text = abdac
pattern = abc

$$h(abc) = x$$

a b d a b c
1 2 3 4 5 6

$$\text{hash}(abd) = y_1$$

$$\text{Is } y_1 = x ?$$

No

$$\text{hash}(bda) = x$$

$$\text{Is } x = x ?$$

↓
yes

now compare each character since there is mismatch take next substring

$$\text{hash}(dab) = y_2$$

$$\text{Is } y_2 = x ?$$

No

$$\text{hash}(abc) = x$$

$$x = x ?$$

yes

compare each character now. all characters are matching, so pattern is found at index 3.

Now, How to efficiently generate hash values is a question.

We are using Rolling hash function.
We are considering

Prime = 3 to keep calculations simple at this moment - but in reality we should use higher value for prime such as 111 and so on.

- 1) $x = \text{old hash} - \text{val}(\text{old char})$
- 2) $x = x / \text{Prime}$
- 3) $\text{new hash} = x + \text{Prime}^{m-1} \times \text{val}(\text{newChar})$

a \rightarrow 1

b \rightarrow 2

c \rightarrow 3

d \rightarrow 4

e \rightarrow 5

g \rightarrow 26

to keep calculations simple
we have assumed these values ^ but in reality you should choose your values of a, b, c, ...

Consider $T = \overbrace{a b e d}^{\text{window of size 4}} a$

a b e

$$1 + 2 \times 3^1 + 5 \times 3^2 = 52$$

While calculating hash for

b e d

$$2 + 5 \times 3^1 + 4 \times 3^2 = 2 + 15 + 36 =$$

Since now a is leaving we subtract 'a' value from 52

$$\text{i.e. } 52 - 1 = 52 / 3$$

$$= 17 + 4 (\text{value of new Char}) \times 3^{\text{length of P} - 1}$$

(2)

$$= 17 + 36$$

$$= \underline{\underline{53}}$$

So, In 3 steps we calculate hash for bed from hash of a be

hash(eda)

$$e \quad d \quad a$$

$$5 \times 4 \times 3^1 + 1 \times 3^2 = \cancel{53} 26$$

hash(eda) b

$$53 - 2 = 51$$

Subtract b value

$$53 - 2$$

$$= 51$$

$$51 / 3 = 17 + 1 \times 3^2$$

$$= 17 + 9$$

$$= 26$$

divide by 3

$$= 51 / 3$$

$$= \cancel{17} 8 + 1 \times 3^2$$

new character

$$= 8 \times 3 + 9$$

$$= 26$$

example: Pattern = abc
text = abedabc

$$h(abc) = 1 + 2 \times 3^1 + 3 \times 3^2$$

$$= 1 + 6 + 27$$

$$= 34$$

← pattern hash value

Now take substring a be

$$\text{hash}(a be) = 1 + 2 \times 3^1 + 5 \times 3^2$$

$$= 1 + 6 + 45$$

$$= 52$$

$$52 \neq 34$$

Now roll forward

$$\begin{aligned}\text{hash}(\text{bed}) &= (s_2 - 1) / 3 = 17 + 4 \times 3^2 \\ &= 17 + 36 \\ &= 53\end{aligned}$$

$$53 \neq 34$$

Now roll forward

$$\begin{aligned}\text{hash}(\text{eda}) &= s_3 - 2 = 51 / 3 = 17 + 1 \times 3^2 \\ &= 26\end{aligned}$$

$$26 \neq 34$$

roll forward

$$\begin{aligned}\text{hash}(\text{dab}) &= 26 - 5 \\ &= 21 / 3 \\ &= 7 + 2 \times 3^2 \\ &= 7 + 18 \\ &= 25\end{aligned}$$

$$25 \neq 34$$

roll forward

$$\begin{aligned}\text{hash}(\text{abc}) &= 25 - 4 = 21 / 3 = 7 + 3 \times 3^2 \\ &= 7 + 27 \\ &= 34\end{aligned}$$

$$34 == 34 \text{ ?}$$

Now Compare each Character of the pattern & substring. if all characters are equal return Index.

Time Complexity in worst case $O(mn)$

Applications - ① plagiarism check.

② if multiple patterns are to be searched in a string Rabin Karp is useful.

text = abcgabcflmxyz

Pattern 1 = gab = hash(gab) = x_1

Pattern 2 = xyz = hash(xyz) = x_2

Pattern 3 = abc = hash(abc) = x_3

Calculate hash for each substring in a text.
If it matches with any of the 3 hash values
characters from both are compared & if
matching is successful we can say that pattern is
present in a Text +.