

Day 15 Pattern Searching Algos.

Q.

- (i) Basic problem
- (ii) Pattern matching
- (iii) problems on (ii)



Q. Given strings s_1 and s_2 , we need to tell whether the two strings are isomorphic.



i.e. if there is a one to one mapping possible for each character in the string with the other one, and

all occurrences of every character in s_1 map to some char in s_2 .

characters of s_2 can be replaced by characters of s_1 to get a similar string.



example: $egg \rightarrow add$ ✓
 $foo \rightarrow ban$ ✗

$geg \rightarrow add$ ✗

relative position matters.



2 hashmaps to store mapping from $s_1 \rightarrow s_2$ and $s_2 \rightarrow s_1$. If while traversing the strings, any anomaly occurs, we will return false,

* Pattern matching

text: a b c b c g l x

pattern: b c g l

↓

Brute force: for each character in the text, we start matching the pattern, one by one.

Time complexity for which will be $O(m * l)$

length of text length of pattern

↓

this is inefficient & needs to be optimised.

text: a b c n a b c d a b m a b c d a b a b c d e y

pattern: a b c d a b c y

↓ ↑

in brute force the characters in the pattern string are repeatedly matched with chars in the text string. No account of already matched chars in the previous check is

maintained, which information could be of use in the next search & can save us a lot of operations.

↓

This is utilised in KMP approach.

$O(m+n)$

↑

so now the matching can begin from c.

↑

for example a b c d a b
matched with text, since a b

is a prefix and a suffix in the substring

(Knott - Morris - Pratt)

↓

basic premise or hypothesis is that if in the pattern string we have a suffix substring which is a prefix then for the next search, the matching can start from element next to the prefix

i) This involves a pre-processing step to build a lookup array which will give us the length of the ^{longest} substring which is a suffix and a prefix at all index of pattern string.

lps
(longest
prefix &
suffix)

1	2	3	4	5	6	7
0	1	2	3	4	5	6
a	b	c	d	a	b	c

(pattern preprocessing step)

0	0	0	0	1	2	3	0
---	---	---	---	---	---	---	---

* we take 2 variables i & j and check if the characters at the two indices match, if they do then $lps[j] = i + 1$, if they don't then i is transferred to $lps[j-1]$ if i is already not 0, if it is then $lps[j] = 0$, and j is incremented. This is repeated till j reaches the end of the array.



(kind of dp!) the movement of i to $lps[j-1]$ is done to check the current string with the previous string, which is

a prefix & suffix, and check for that.

	↓		↓				↓
a	y	a	b	a	y	a	y
0	0	1	0	1	2	3	2

a	y	a	y	a	b	a	b	a	y	a	y
0	0	1	2	3	0	1	0	1	2	3	4

↓	↓	↓	↓	↓	↓		
a	a	a	a	a	a	b	
0	1	2	3	4	5	0	
↓	↓	↓	↓	↓	↓		
a	b	a	b	a	b	a	b
0	0	1	2	3	4	5	6
0	1	2	3	4	5	6	7

suffix which is also a prefix for the substring ending at the i th index.

			↓				↓	↓
0	1	2	3	4	5	6	7	8
a	y	c	d	a	y	c	a	y
0	0	0	0	1	2	3	1	2

ii) using the lookup table to find the pattern in the text.

	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓	
text :	a	b	x	a	b	c	a	b	c	a	b	y
pattern :		↓	↓	↓	↓	↓	↓					
		a	b	c	a	b	y					
		0	0	0	1	2	0					
					↑	↑	↑					

hence pattern matches

* on a mismatch the lookup array takes its pointer to the index next to the lps of the previous subarray to check if the match happens for that substring or not & the search continues.

↓

constructing lps $\rightarrow O(m)$

search for pattern in text using lps $\rightarrow O(n)$

↓
 $O(m+n)$

space: $O(m)$ auxiliary space.

Q. Given a ~~string~~ string, how many characters need to be added in front to make it a palindrome.

edef \$fede
001000123

str.length() - lps.back()

for a palindrome, $rev(str) = org(str)$,

↓
 problem can be to find how
 far from rev, our org. string
 is.

↓
 longest prefix in org. which is
 suffix in reverse, and the
 length is subtracted from the
 total length.

f edef

↓
 length of the longest palindrome in
 the string starting at ~~index~~ 0th index
 the next letters can be added to
 the string accordingly.

cc a a e b a a b e a a c c

cc a a b a a c c

for a palindrome, the
 string, all prefix and
 suffixes are same.

← a a b a a a a b a a
 0 1 0 1 2 0 1 2 3 4 5

a b c d e f e d c b a
 0 0 0 0 0 0 0 0 0 0 0

d a b a e f e a b a d
 0 0 0 0 0 0 0 0 0 0 0

d a b a e d f d e a b a d
 0 0 0 0 0 1 0 1 0 0 0 0 1

d a b a e o d f d a e a b o d
 0 0 0 0 0 0 1 0 2 2 0 0 0 0 1

longest prefix which is
 a suffix

Another pattern matching algorithm

↓
Rabin Karp algorithm

↓
uses the concept of hashing.

C++, JAVA,

prime for hash
function; 31

↓
i) compute the hash of the
pattern.

ii) calculate hash of all substrings
of length of the pattern.

{ If the hashes match, then do character
by character matching, else
move on

→ if we calculate
hash code for
all substrings of
length m, then
 $O(m \times n)$ time
will be elapsed
efficient method req.

for example: a b d a b c : text

a b c : pattern

let hash: $\sum_{i=0}^m p^i \cdot \text{str}[i]$ where p is a prime no

$$\begin{aligned}\text{hash}(\text{abc}) &= 1 \times 3^0 + 2 \times 3^1 + 3 \times 3^2 \\ &= 1 + 6 + 27 \\ &= 34\end{aligned}$$

↓
let p = 3 here
and let ascii
values for
alphabets be 0-26
index, i.e.
-26 for lowercase
alphabets.

$$\begin{aligned}\text{hash}(\text{abd}) &= 1 \times 3^0 + 2 \times 3^1 + 4 \times 3^2 \\ &= 1 + 6 + 36 \\ \text{hash}(\text{abd}) &\neq \text{hash}(\text{abc}) = 43\end{aligned}$$

now hash(bda) is not calculated from scratch
rather hash(abd) is recycled for bda

$$\begin{aligned} \text{hash}(\text{bda}) &= 2 + 3^0 + 4 \times 3 + 1 \times 3^2 \\ &= 2 + 12 + 9 \\ &= 23 \end{aligned}$$

$$\frac{n(\text{abd}) - a}{p} + a \times p^2$$

$$\frac{n(\text{abd}) - a}{3} + a \times 9 \rightarrow O(1) \text{ operation, rather than looping on pattern or hash.}$$

$$\frac{43 - 1}{3} + 1 \times 9 = 14 + 9 = 23$$



in worst case, the time complexity
 $O(m \times n)$



when all characters of pattern are same & all characters of text are same.



hashes will be same, and character by character checks would be required repeatedly, which was avoided in KMP.

if the hash function is not good, and collisions are frequent, then time elapsed will be more



implementation of hash code fn. and the i/p.

* How do we know that prime number hash is a good hash function?

i) distribution ^{should} be good (well spread) ^{uniform}

ii) if non prime, then the numbers fall in the same range, which is undesirable

Q. Given a stream of characters, we need to return whether the updated string is a palindrome or not.

↓

i) create partition, calculate hash for the rev of RHS, if they hash for rev(RHS) \Rightarrow hash for

a b c b a k

$h(abc)$ $h(kab)$

LHS, then it becomes a candidate for character by

character search, if it doesn't then we return no.

ii) calculating the hash each time would be expensive and to answer the query in $O(1)$, we need to use the previous hash values to calculate the current value.

a b c | b a k

99
5

$$h(abc) = 1 \times 3^0 + 2 \times 3^1 + 3 \times 3^2$$

$$= 34 \times 3 = 2$$

$$h(kab) = 11 + 3 + 18 = 32 \times 3 = 2$$

↓

a b c | a k d

even

↓

odd

$$h(abc) = 34 \times 3 = 2$$

$$h(dka) = 4 + 11 \times 3^1 + 1 \times 3^2 = 4 + 33 + 9 = 46$$

$$h(dka) = (h(kab) - 2 \times 3^2) \times 3 + d(4)$$

$$= (32 - 18) \times 3 + 4$$

$$= 14 \times 3 + 4 = 42 \times 3 = 0.$$

$$\begin{array}{cccccc} a & b & c & b & a & k & d & e \\ & & & & & & & \text{odd} \\ & & & & & & & \downarrow \\ & & & & & & & \text{even} \end{array}$$

$$\begin{aligned} h(abc b) &= h(abc) + 2 + 3^4 \\ &= 34 + 2 \times 81 \\ &= 34 + 162 = 196 \div 3 = 0 \end{aligned}$$

$$\begin{aligned} h(edka) &= h(dka) \times 3 + 5 \\ &= 42 \times 3 + 5 \\ &= 126 + 5 = 131 \div 3 = 2 \end{aligned}$$

* prime number for hash: consider key $\in [0, 100]$
 and a hash table has $m = 12$. Since 3 is
 a factor of 12, the keys which are different
 order multiples of 3 will fall in
 buckets that are multiples of 3.

$$\begin{aligned} \text{keys } [0, 12, 24, 36] &\rightarrow 0 \\ [3, 15, 27, 39] &\rightarrow 3 \\ [6, 18, 30, 42] &\rightarrow 6 \\ [9, 21, 33, 45] &\rightarrow 9 \end{aligned}$$

\downarrow
 if the elements are biased to being 3
 centre then the chaining in these
 buckets will be very high, which
 is unfavourable. \therefore the factors of m ,
 should be minimal, and that is best
 represented or taken care of by a
 prime number.