

pattern matching algorithms

(1)

the problem of string matching

Given a string 's', the problem of string matching deals with finding whether a pattern 'p' occurs in 's' and if 'p' does occur then returning position in 's' where 'p' occurs.

Brute force

Time complexity: $O(mn)$

Here, compare the first element of the pattern to be searched 'p' with the first element of the string 's'. If the first element matches, compare the second element of 'p' with the second element of 's'. If match found proceed likewise until entire 'p' is found. If a mismatch is found at any position, shift 'p' to the position to the right and repeat comparisons beginning from first element of 'p'.

[3] Text: abcabaabcabac

[4] Pattern: abaa

Text:

0	1	2	3	4	5	6	7	8	9	10	11	12
a	b	c	a	b	a	a	b	c	a	b	a	c

Pattern:

a	b	a	a
---	---	---	---

 mismatch, shift 'p' to the right (one position)

a	b	a	a
---	---	---	---

 mismatch

a	b	a	a
---	---	---	---

a	b	a	a
---	---	---	---

 match found, return the position

Drawback: comparisons are more

4 program

#include <stdio.h>

#include <string.h>

int search(char *t, char *p)

{ int i, j; ~~flag = 0~~;

int M = strlen(t);

int N = strlen(p);

for(i=0; i<=M-N; i++)

{ for(j=0; j<N; j++)

{ if(p[j] != t[i+j])
break;

}

if(j==N)

{ printf("Pattern found at index %d", i);

flag = 1; }

}

return i; if(flag == 0) printf("Pattern not found");

void main()

{ ~~int f;~~

char text[100], pat[20];

printf("Enter the main string:");

gets(text);

printf("Enter the pattern:");

gets(pat);

~~search~~ search(text, pat);

if(f != 1)

printf("Pattern not found\n");

else

printf("Pattern found at index %d\n", f);

}

01 2 3 4 5 6 7 8 9 10
text: abcabacabac d
pat: abaa

M=11, N=4

for(i=0; i<=M-N; i++)

{ for(j=0; j<N; j++)

{ p[j] != t[i+j]

break;

if(j==N)

{ printf("Pattern found at index %d", i);

flag = 1; }

return i; if(flag == 0) printf("Pattern not found");

② KMP Algorithm (Knuth-Morris-Pratt) $O(m+n)$ ①

- It is one of the most popular patterns matching algorithms.
- It was the first linear time complexity algorithm for string matching and used to find a pattern in a Text.
- This algorithm compares character by character from left to right. But whenever a mismatch occurs, it uses a preprocessed table called 'prefix table' to skip characters comparison while matching. It is also known as LPS Table. (Longest proper prefix which is also Suffix)

Steps for creating LPS table (prefix table)

- Step
- ① Define a one dimensional array with the size equal to the length of the pattern i.e. $LPS[SIZE]$
 - ② Define variables i & j . set $i=0, j=1$ and $LPS[0]=0$
 - ③ Compare the characters at $pattern[i]$ and $pattern[j]$.
 - ④ if both are matched then set $LPS[i]=i+1$ and increment both i and j by one. goto step ③.
 - ⑤ if both are not matched then check the value of variable ' i '. if it is '0' then set $LPS[j]=0$ and increment ' j ' value by one, if it is not '0' then set $i = LPS[i-1]$. Goto step ③.
 - ⑥ Repeat above steps until all the values of $LPS[]$ are filled.

Q) create LPS table for the following patterns

pattern:

0	1	2	3	4	5	6
A	B	C	D	A	B	D

Sol: Define LPS[] array with size '7'

0	1	2	3	4	5	6

Now, $i=0$, $j=1$ and $LPS[0]=0$

0	1	2	3	4	5	6
0						

 $i=0, j=1$

Compare pattern[0] with pattern[1] i.e 'A' with 'B'

mismatch and $i=0$, $LPS[1]=0$ and increment 'j' by one i.e $j=2$

0	1	2	3	4	5	6
0	0					

 $i=0, j=2$

Compare pattern[0] with pattern[2] i.e 'A' with 'C'.

mismatch and $i=0$, $LPS[2]=0$ and increment 'j' i.e $j=3$

0	1	2	3	4	5	6
0	0	0				

 $i=0, j=3$

Compare pattern[0] with pattern[3] i.e 'A' with 'D'

mismatch and $i=0$, $LPS[3]=0$ and increment 'j' i.e $j=4$

0	1	2	3	4	5	6
0	0	0	0			

 $i=0, j=4$

Compare pattern[0] with pattern[4] i.e 'A' with 'A', match

So, $LPS[4]=0+1=1$ and increment both i & j i.e $i=1, j=5$

0	1	2	3	4	5	6
0	0	0	0	1		

 $i=1, j=5$

Compare pattern[1] with pattern[5], 'B' == 'B', match

So, $LPS[5]=1+1=2$. increment 'i' & 'j'. i.e $i=2, j=6$

0	1	2	3	4	5	6
0	0	0	0	1	2	

 $i=2, j=6$

Compare pattern[2] with pattern[6] i.e 'C' with 'D'.

mismatch and $i \neq 0$. set $i = LPS[i-1] = LPS[1] = 0$

0	1	2	3	4	5	6
0	0	0	0	1	2	

 $i=0$ and $j=6$

Compare pattern[0] with pattern[6] i.e 'A' with 'D', mismatch and $i=0$. So, $LPS[6]=0$ and increment 'j'

LPS[]

0	1	2	3	4	5	6
0	0	0	0	1	2	0

Use of LPS table a) prefat when a mismatch has occurred with the first character in the pattern move the pattern one position right 3

It is used to decide how many characters are to be skipped for comparison when a mismatch has occurred.
 ⇒ when a mismatch occurs, check the LPS value of the previous character of the mismatched character in the pattern. if it is '0' then start comparing the first character of the pattern with the next character to the mismatched character in the text. if it is not '0' then start comparing the character which is at an index value equal to the LPS value of the previous character to the mismatched character in the pattern with the mismatched character in the text.

① Consider the following Text and pattern

Text: ABC ABCDAB ABCDABCDABDE

pattern: ABCDABD

	0	1	2	3	4	5	6
LPS	A	B	C	D	A	B	D
	0	0	0	0	1	2	0

Text:

A	B	C	X	A	B	C	D	A	B	X	A	B	C	D	A	B	C	D	A	B	D	E
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 Pattern:

A	B	C	D	A	B	D
---	---	---	---	---	---	---

Here mismatch occurred at pattern[3]. So take LPS[2] value i.e '0'.
 we must compare first character in the pattern with the mismatched character of text.
 i.e X with A mismatch and LPS[0]=0, move right.

Text:

A	B	C	X	A	B	C	D	A	B	X	A	B	C	D	A	B	C	D	A	B	D	E
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 Pattern:

A	B	C	D	A	B	D
---	---	---	---	---	---	---

Here mismatch occurred at pattern[6]. So take LPS[5] i.e '2' (≠ 0).
 we compare pattern[2] with mismatched character in the text

Text:

A	B	C	A	B	C	D	A	B	X	A	B	C	D	A	B	C	D	A	B	D	E
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 Pattern:

A	B	C	D	A	B	D
---	---	---	---	---	---	---

Here, mismatch occurred at pattern[2]. So take LPS[2] i.e '0'.
 compare pattern[0] with the mismatched character in the text
 i.e 'A' with X, mismatch LPS[0]=0. So, compare pattern[0] with X, move right

Text:

A	B	C	A	B	C	D	A	B	X	A	B	C	D	A	B	C	D	A	B	D	E
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 Pattern:

A	B	C	D	A	B	D
---	---	---	---	---	---	---

mismatch occurred at pattern[6]. So take LPS[5] i.e '2' (≠ 0).
 compare the pattern[2] with the mismatched character in the text

Text:

A	B	C	A	B	C	D	A	B	X	A	B	C	D	A	B	C	D	A	B	D	E
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 Pattern:

A	B	C	D	A	B	D
---	---	---	---	---	---	---

if matched, return 'is' as the index

compute failure function (or) prefix table for the following strings

① abcaby ② aabaabaaa ③ abcdabcy

④ abcdabca

①

	0	1	2	3	4	5
	a	b	c	a	b	y
LPS[]	0	0	0	1	2	0

initial value
 $i=0, j=0$
 $LPS[0]=0$

②

	0	1	2	3	4	5	6	7	8
	a	a	b	a	a	b	a	a	a
	0	1	0	1	2	3	4	5	2

③

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

④

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

Text: abcxabcdabx abcdabcyab
 pattern: abcdabcy

Text: abxabcabcabyabcd

pattern: abcaby

Advantages: optimal & fast $O(m+n)$
 disadvantages: does not work so well as the size of the alphabet increases. By which more chances of mismatch occurs.

③ Boyer-Moore pattern matching Algorithm

It preprocesses the pattern based on two approaches

- ① Bad character rule
- ② Good Suffix rule

At every step, it moves the pattern by the maximum of the moves suggested by the two rules. So, it uses best of the two rules at every step. Unlike the previous pattern searching algorithms, Boyer-Moore algo starts matching from the last character of the pattern.

Bad character rule:

The character of the text which doesn't match with the current character of the pattern is called the Bad character. Upon mismatch, we shift the pattern until—

- Case ① The mismatch becomes a match
- or ② pattern 'p' move past the mismatched character.

eg: case 1: (The mismatch becomes a match)

Text[]: G C A A T G C C T A T G T G A C C
Pattern[]: T A T G T G

↑ ↑ ↑ ↑ ↑ ↑
Mismatch

Start comparison from 5, we got a mismatch at index '3'. Here bad character is 'A'. So we will search for last occurrence of 'A' in the pattern. we got 'A' at index 1. Now we will shift pattern so that 'A' in the pattern get aligned with 'A' in the text. i.e.

Text[]: G C A A T G C C T A T G T G A C C
Pattern[]: T A T G T G

↑ ↑ ↑ ↑ ↑ ↑
Mismatch becomes match

eg: case 2: (pattern 'p' move past the mismatched character)

Text[]: G C A A T G C C T A T G T G A C C
Pattern[]: T A T G T G

↑ ↑ ↑ ↑ ↑ ↑
Mismatch

Text[]: G C A A T G C C T A T G T G A C C
Pattern[]: T A T G T G

↑ ↑ ↑ ↑ ↑ ↑
Mismatch

Here, we have a mismatch at index 7 i.e. C. 'C' does not exist in the pattern, so shift pattern past to the position 7.

Case 3: 'p' moves past 't'

T: A A C A B A B A C B A
 p: C B A A B

C B A A B

Here, there exist no occurrence of t("AB") in p and also there is no prefix in 'p' which matches with the suffix of t. So, we will shift the 'p' past the 't'.

Ex: T: C G T G C C T A C T T A C T T A C T T A C T T A C
 p: C T T A C T T A C

another occurrence of 't'

C T T A C T T A C

can't

T: C G T G C C T A C T T A C T T A C T T A C T T A C
 p: C T T A C T T A C

C T T A C T T A C

prefix 'p' and suffix of 't'

Ex: T: G T T A T A G C T G A T C G C G G C G T A G C G G C G A A
 p: G T A G C G G C

best of two

Step 1: bc: 6 steps forward, gs: 0

Step 2: bc: 0, gs: 2

Steps: bc: 2, gs: 7

t = G C G G C G

prefix G, suffix G

match prefix of pattern & suffix of text.

pattern found