

Xiaoting Lian

CISC3220

String matching

The first solution for string match is called **Naïve**

Suppose we have a string A[A, S, D, F, G, H, J, K, L], the second string is B[D, F, G]

The idea of this algorithm is from beginning shift to the right until you match the string.

1. $n \leftarrow \text{length}[T]$
2. $m \leftarrow \text{length}[P]$
3. for $s \leftarrow 0$ to $n-m$ do
4. $j \leftarrow 1$
5. while $j \leq m$ and $T[s+j] = P[j]$ do
6. $j \leftarrow j+1$
7. If $j > m$ then
8. return valid shift s
9. return no valid shift exist

the for-loop in line 3 is executed at most $n - m + 1$ times, and the while-loop in line 5 is executed at most m times, The run time for naïve string matching is $O((n - m + 1)m)$

Example for Naïve string-matching algorithm

SHIFT = 0

A = A. S. D. F. G. H, J, K, L

B = D, F, G

SHIFT = 1

A = A. S. D. F. G. H, J, K, L

B = D, F, G

SHIFT = 2

A = A. S. D. F. G. H, J, K, L

B = D, F, G

At this point the string is matched.

The second algorithm for string matching is **Rabin-Karp String Matching Algorithm**

Suppose we have a string $S[1, 3, 4, 1, 4, 1, 2, 6]$

Pattern $P[4, 1, 4]$

The basic idea is we consider $P[4, 1, 4]$ as $4+1+4=9$, then we compare this value to every 3 elements in S string, until we find the matching string

To keep going on,

We take 1,2,4 from S , $1+2+4=7$ which is not equal to 9,

we shift to right 3,4,1, $3+4+1=8$ which is not equal to 9

next we compare 4, 1, 4 $4+1+4=9$ valid match

as we keep going on, we can also find that 1+2+6 this is also equal to 9, but the string does not match, we call this **spurious hit**.

To solve this problem we need more variables for this algorithm to work

Example,

we set $(4*10^2)+(1*10^1)+(4*10^0)$ this will give us $400+10+4=414$

what we do next is take 1, 3, 4, $(1*10^2)+(3*10^1)+(4*10^0)=100+30+4=134$,

next is $(134 - 1*10^2)*10 + 1*10^0$ we minus $1*10^2$ is to remove 1 from previous high order, multiply by 10 is to move all power by 1 with the previous then add the new digit.

As we keep shifting, we will notice that there will be only one substring in S that is match to P .

The book is giving different example that instead of adding all numbers, each character is converted in to decimal digit and strings of k -consecutive characters are represented as length- k decimal number, and find module for all numbers

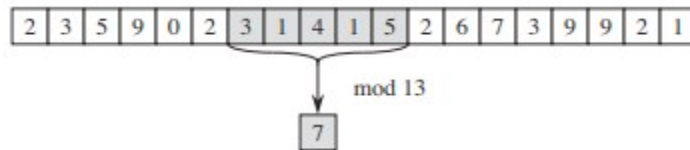
Example: string 31415 is represented as 31,415

RABIN-KARP-MATCHER(T, P, d, q)

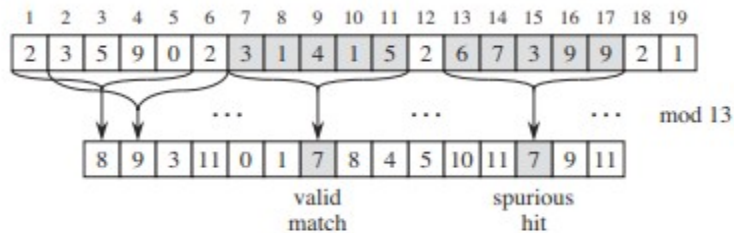
```

1   $n = T.length$ 
2   $m = P.length$ 
3   $h = d^{m-1} \bmod q$ 
4   $p = 0$ 
5   $t_0 = 0$ 
6  for  $i = 1$  to  $m$            // preprocessing
7       $p = (dp + P[i]) \bmod q$ 
8       $t_0 = (dt_0 + T[i]) \bmod q$ 
9  for  $s = 0$  to  $n - m$        // matching
10     if  $p == t_s$ 
11         if  $P[1..m] == T[s+1..s+m]$ 
12             print "Pattern occurs with shift"  $s$ 
13     if  $s < n - m$ 
14          $t_{s+1} = (d(t_s - T[s+1]h) + T[s+m+1]) \bmod q$ 

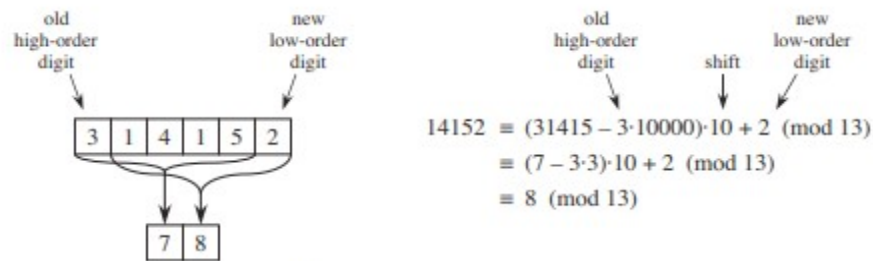
```



(a)



(b)



(c)

Average run time $O(n + m)$.

String matching with finite automata

Example

String S [a, b, a, b, a, b, a, c, a, b, a]

Pattern P [a, b, a, b, a, c, a]

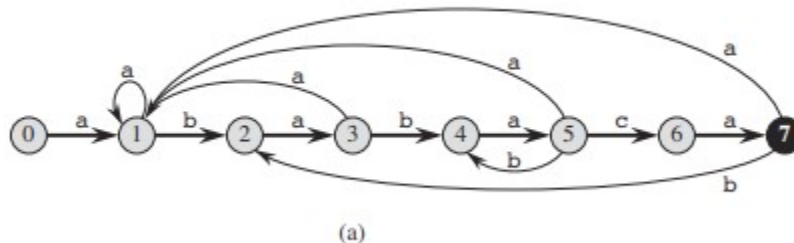
Numbers of type of element is a, b, c, d, 4 types

Total numbers of element of P is 6, the size of the table we will create is $6+1=7$.

State	a	b	c	P
0	1	0	0	a
1	1	2	0	b
2	3	0	0	a
3	1	4	0	b
4	5	0	0	a
5	1	4	6	c
6	7	0	0	a
7	1	2	0	

a match with a therefore the first box we put 1, b and c does not match we put zero.

Next row, we have ab and aa, only we only a and a match, we put 1 in the first box in sec row, then ab and ab, we get 2 match, so we put 2, do the same to the rest



Every time that arrow turn back that means string does not matching

FINITE-AUTOMATON-MATCHER(T, δ, m)

```

1   $n = T.length$ 
2   $q = 0$ 
3  for  $i = 1$  to  $n$ 
4       $q = \delta(q, T[i])$ 
5      if  $q == m$ 
6          print "Pattern occurs with shift"  $i - m$ 

```

The Knuth-Morris-Pratt algorithm

Example

We have string $S[a, b, a, b, a, b, a, b, a, b, c]$

The pattern is $P[a, b, c, d]$

Let S has a pointer i and P has pointer j

S

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

P

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

The second row for P table is when a and a are matching in P we note that second a with 1 and b appears twice we note that 2 . c only appears once, we put 0 , we consider ab as prefix.

Now we start comparing, when i is 1 , we have a , $j + 1$ is a too, then we move to the next, b and b are match as well, then we move again.... We see that when $i=5$ is a , $j+1= 5$ is c , doesn't match what we do here is not go all the way back, we move j to the noted number 2 in the p box which is the first b . then we keep comparing, j will go back to 2 three time. Then we find the matched string.

KMP-MATCHER(T, P)

```
1   $n = T.length$ 
2   $m = P.length$ 
3   $\pi = \text{COMPUTE-PREFIX-FUNCTION}(P)$ 
4   $q = 0$  // number of characters matched
5  for  $i = 1$  to  $n$  // scan the text from left to right
6      while  $q > 0$  and  $P[q + 1] \neq T[i]$ 
7           $q = \pi[q]$  // next character does not match
8      if  $P[q + 1] == T[i]$ 
9           $q = q + 1$  // next character matches
10     if  $q == m$  // is all of  $P$  matched?
11         print "Pattern occurs with shift"  $i - m$ 
12          $q = \pi[q]$  // look for the next match
```

COMPUTE-PREFIX-FUNCTION(P)

```
1   $m = P.length$ 
2  let  $\pi[1..m]$  be a new array
3   $\pi[1] = 0$ 
4   $k = 0$ 
5  for  $q = 2$  to  $m$ 
6      while  $k > 0$  and  $P[k + 1] \neq P[q]$ 
7           $k = \pi[k]$ 
8      if  $P[k + 1] == P[q]$ 
9           $k = k + 1$ 
10      $\pi[q] = k$ 
11 return  $\pi$ 
```

Questions:

1. $T[a,a,a,a,a,a,b]$, $P[a,b]$ how many shift do we needed in naïve string matching algorithm.
Ans: $8-2=6$ shift
2. Rabin Karp Algorithm makes use of elementary number theoretic notions. T or F
Ans: T
3. What is worst case run time for Rabin-Karp string matching algorithm
Ans: $O(n+m)$
4. What is average case run time for Rabin-Karp string matching algorithm
Ans: $O(nm)$
5. What is it called when if the rolling hash produces a candidate match due to this hash collision, which turns out not to be a string match?
Ans: spurious hit