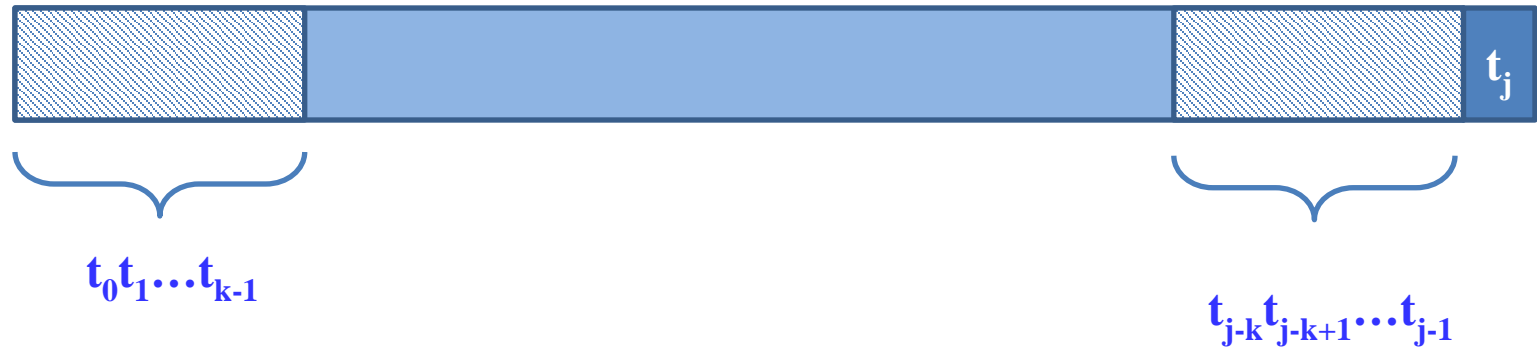


Knuth–Morris–Pratt Algorithm

KMP 算法

最大子串的概念

- 所谓真子串是指模式串 t 存在某个 k ($0 < k < j$)，使得“ $t_0t_1\dots t_{k-1}$ ” = “ $t_{j-k}t_{j-k+1}\dots t_{j-1}$ ”成立。



所有真子串中最长的称为最大子串，其长度 k 记为 $next[j]$

j	0	1	2	3	4	5	6	7	8
$t[j]$	A	B	A	C	A	B	A	B	C
$next[j]$	-1	0	0	1	0	1	2	3	2

目标串



模式串



目标串



模式串



j	0	1	2	3	4	5	6
t[j]	A	B	C	D	A	B	D
next[j]	-1	0	0	0	0	1	2

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

i: 01234567890123456789012
 S: ABC_ABCDAB_ABCDABCDABDE
 t: ABCDABD
 j: 0123456

归纳起来，定义next[j]函数如下：

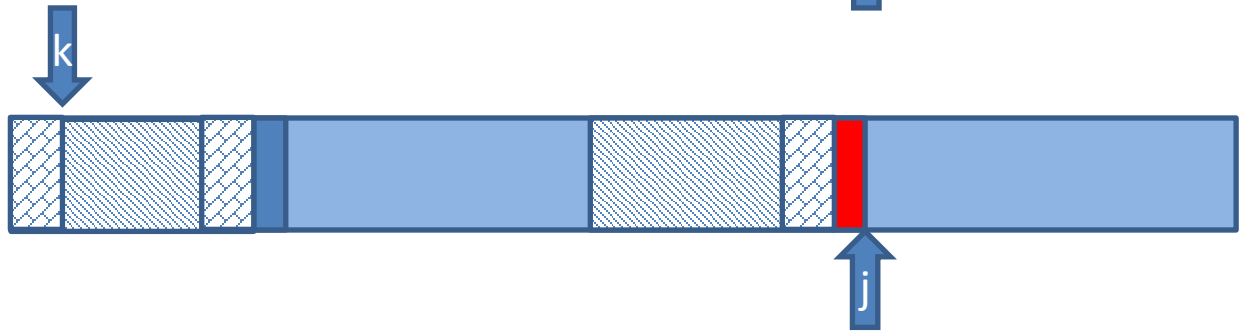
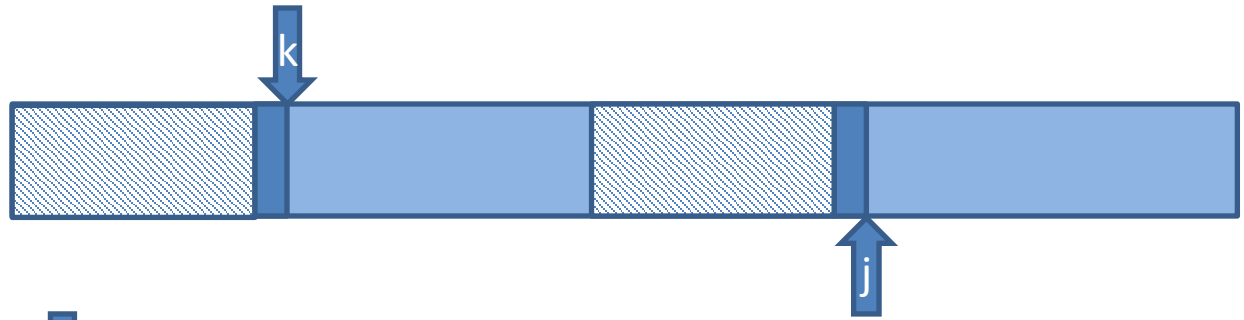
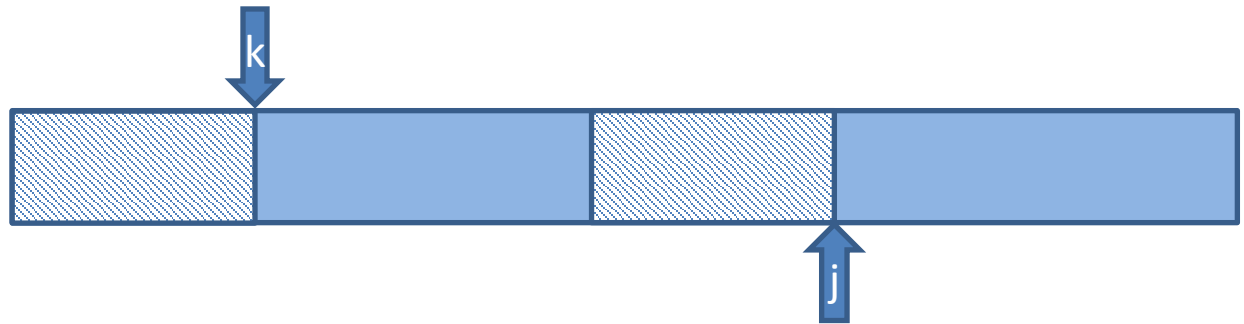
$$\text{next}[j]= \begin{cases} \max\{k|0 < k < j, \text{且 } "t_0t_1\dots t_{k-1}" = "t_{j-k}t_{j-k+1}\dots t_{j-1}" \} & \text{当此集合非空时} \\ -1 & \text{当} j=0 \text{时} \\ 0 & \text{其他情况} \end{cases}$$

t="abab"对应的next数组如下：

<i>j</i>	0	1	2	3
t[j]	a	b	a	b
next[j]	-1	0	0	1

由模式串t求出next值:

```
void GetNext(SqString t,int next[])
{  int j,k;
   j=0;k=-1;next[0]=-1;
   while (j<t.length-1)
   {  if (k==-1 || t.data[j]==t.data[k])
      {  j++;k++;
         next[j]=k;
      }
      else k=next[k];
   }
}
```



KMP算法:

```
int KMPIndex(SqString s,SqString t)
{  int next[MaxSize],i=0,j=0;
   GetNext(t,next);
   while (i<s.length && j<t.length)
   {   if (j==-1 || s.data[i]==t.data[j])
       {   i++;
           j++;           //i,j各增1
       }
       else j=next[j];   //i不变,j后退
   }
   if (j>=t.length)
       return(i-t.length); //返回匹配模式串的首字符下标
   else
       return(-1);        //返回不匹配标志
}
```


设主串s的长度为 n ，子串t长度为 m 。

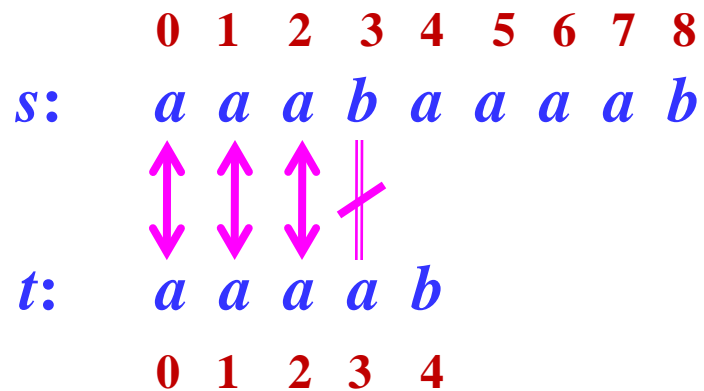
在KMP算法中求next数组的时间复杂度为 $O(m)$ ，在后面的匹配中因主串s的下标不减即不回溯，比较次数可记为 n ，所以KMP算法总的时间复杂度可以认为是 $O(n+m)$ 。

设目标串s=“*aaabaaaab*”，模式串t=“*aaaab*”。KMP模式匹配过程。

求t的next:

j	0	1	2	3	4
t[j]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>	<i>b</i>
next[j]	-1	0	1	2	3

j	0	1	2	3	4
t[j]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>	<i>b</i>
next[j]	-1	0	1	2	3




失败:
i=3
j=3, j=next[3]=2

j	0	1	2	3	4
t[j]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>	<i>b</i>
next[j]	-1	0	1	2	3

i=3
j=2

0 1 2 3 4 5 6 7 8
s: *a a a b a a a a b*

t: *a a a a b*
0 1 2 3 4




失败:
i=3
j=2, j=next[2]=1

j	0	1	2	3	4
t[j]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>	<i>b</i>
next[j]	-1	0	1	2	3

i=3
j=1

0 1 2 3 4 5 6 7 8
s: *a a a b a a a a b*

t: *a a a a b*
0 1 2 3 4



失败:
i=3
j=1, j=next[1]=0

j	0	1	2	3	4
t[j]	a	a	a	a	b
next[j]	-1	0	1	2	3

i=3
j=0

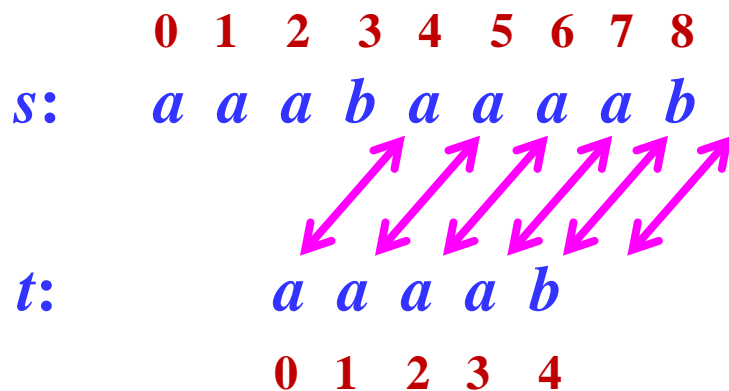
0 1 2 3 4 5 6 7 8
s: *a a a b a a a a b*

t: *a a a a b*
 0 1 2 3 4

失败:
i=3
j=0, j=next[0]=-1

j	0	1	2	3	4
t[j]	a	a	a	a	b
next[j]	-1	0	1	2	3

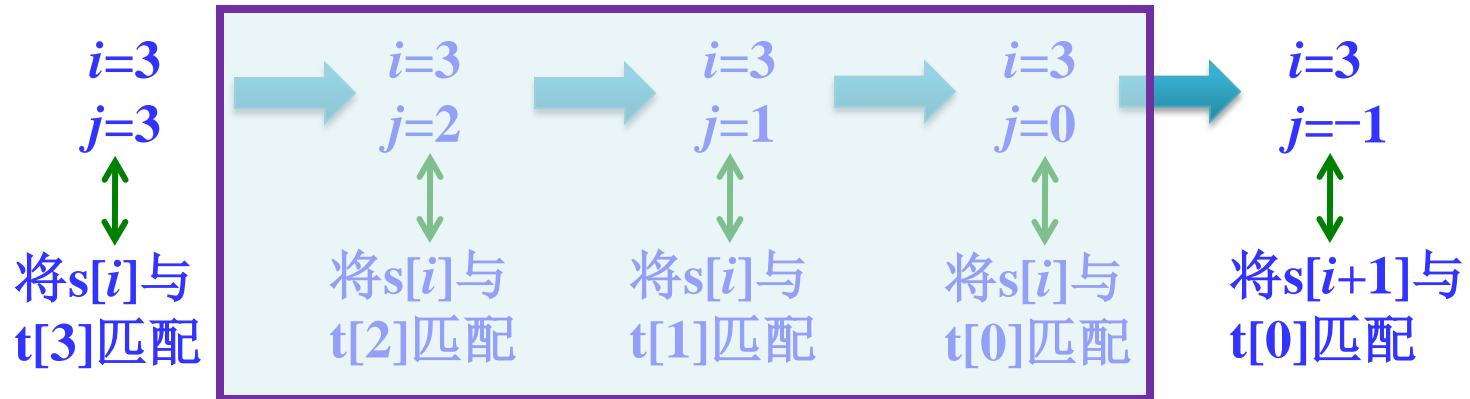
因为j=-1:
i++;
j++;



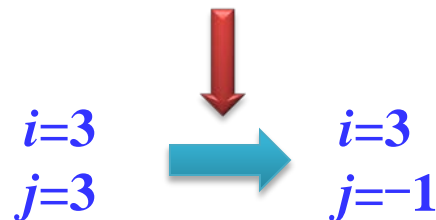
成功:
返回4

j	0	1	2	3	4
t[j]	a	a	a	a	b
next[j]	-1	0	1	2	3

前面的匹配过程：



因为t[3]=t[2]=t[1]=t[0]='a' 是不必要的



将next改为nextval:

j	0	1	2	3	4
t[j]	a	a	a	a	b
next[j]	-1	0	1	2	3
nextval[j]	-1	-1	-1	-1	3

$$\text{next}[1]=0$$

$$t[1]=t[\text{next}[1]]=t[0]='a'$$

$$\therefore \text{nextval}[1]=\text{nextval}[0]=-1$$

$$t[4]='b' \neq t[\text{next}[4]]=t[3]='a'$$

$$\therefore \text{nextval}[4]=\text{next}[4]$$

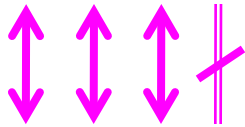


- $\text{nextval}[0]=-1$
- 当 $t[j]=t[\text{next}[j]]$ 时: $\text{nextval}[j]=\text{nextval}[\text{next}[j]]$
- 否则: $\text{nextval}[j]=\text{next}[j]$

用nextval取代next, 得到改进的KMP算法。

使用改进后的KMP算法示例:

j	0	1	2	3	4
t[j]	a	a	a	a	b
nextval[j]	-1	-1	-1	-1	3

0 1 2 3 4 5 6 7 8
s: **a a a b a a a a b**

t: **a a a a b**
0 1 2 3 4

失败:

$i=3$

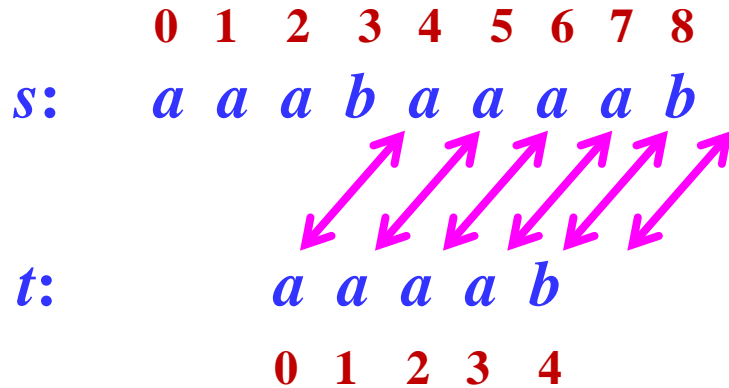
$j=3, j=\text{nextval}[3]=-1$

j	0	1	2	3	4
t[j]	a	a	a	a	b
nextval[j]	-1	-1	-1	-1	3

因为j=-1:

i++;

j++;



成功:
返回4



改进后的KMP算法进一步提高模式匹配的效率。

由模式串t求出nextval值:

```
void GetNextval(SqString t,int nextval[])
{  int j=0,k=-1;
   nextval[0]=-1;
   while (j<t.length)
   {  if (k==-1 || t.data[j]==t.data[k])
       {  j++;k++;
          if (t.data[j]!=t.data[k])
              nextval[j]=k;
          else
              nextval[j]=nextval[k];
        }
       else
           k=nextval[k];
   }
}
```

修改后的KMP算法:

```
int KMPIndex1(SqString s,SqString t)
{  int nextval[MaxSize],i=0,j=0;
   GetNextval(t,nextval);
   while (i<s.length && j<t.length)
   {  if (j==-1 || s.data[i]==t.data[j])
      {  i++;
         j++;
      }
      else
         j=nextval[j];
   }
   if (j>=t.length)
      return(i-t.length);
   else
      return(-1);
}
```

数据结构经典算法的启示

BF算法



KMP算法

利用模式串中部分匹配信息

课堂练习：目标串s=“abcaabbcaaabababaabca”
模式串t=“babab”

求KMP算法匹配过程

求改进的KMP算法的匹配过程

作业：第4章，5，6

——本讲完——