

串和线性结构的区别

串的各种概念

串的匹配 (KMP算法) (next数组)

2.4 String

- Definition and operations
- Applications



2.4.1. Definition



- **串**: 零个或多个字符组成的有限序列。
- **串长度**: 串中所包含的字符个数。
- **空串**: 长度为0的串, 记为: “”。
- **非空串**通常记为: $S = "s_1 s_2 \dots s_n"$
 - 其中: S 是串名, 双引号是**定界符**, 双引号引起来的部分是串值, $s_i (1 \leq i \leq n)$ 是一个任意字符。
 - 字符集: ASCII码、扩展ASCII码、Unicode字符集
- **子串**: 串中任意个连续的字符组成的子序列。
- **主串**: 包含子串的串。
- **子串的位置**: 子串的第一个字符在主串中的序号。

1

2

a= "Welcome to Beijing"
b= "Welcome"
c= "Bei"
d= "welcometo"
子串的位置: 子串在主串中第一次出现的第一个字符的位置。
两个串相等: 两个串的长度相等, 并且各个对应的字符也都相同。
例如, 有下列四个串a, b, c, d:
a= "program"
b= "Program"
c= "pro"
d= "program "



串的基本操作:

- (1) 创建串 StringAssign (s,string_constant)
- (2) 判断串是否为空 StringEmpty(s)
- (3) 计算串长度 Length(s)
- (4) 串连接 Concat(s1,s2)
- (5) 求子串 SubStr(s1,s2,start,len)
- (6) 子串的定位 Index(s1,s2)
- (7) 子串的插入和删除



3

4

2.4.2 串的模式匹配算法 (Pattern Matching)



子串定位运算又称为模式匹配(Pattern Matching)或串匹配(String Matching), 此运算的应用在非常广泛。例如, 在文本编辑程序中, 我们经常要查找某一特定单词在文本中出现的位置。显然, 解此问题的有效算法能极大地提高文本编辑程序的响应性能。

5

在串匹配中, 一般将主串称为目标串, 子串称之为模式串。设 S 为目标串, T 为模式串, 且不妨设:

$$S = "s_1 s_2 \dots s_{n-1}" \quad T = "t_1 \dots t_{m-1}"$$

朴素模式匹配算法(Brute-Force算法): 枚举法

从主串 S 的第一个字符开始和模式 T 的第一个字符进行比较, 若相等, 则继续比较两者的后续字符; 否则, 从主串 S 的第二个字符开始和模式 T 的第一个字符进行比较, 重复上述过程, 直到 T 中的字符全部比较完毕, 则说明本趟匹配成功; 或 S 中字符全部比较完, 则说明匹配失败。



6

No.1 ababcabcabab i=3,j=3 时失败

No.2 ababcbabcabab i=2,j=1 时失败

No.3 ababca**b**cacbab i=7,j=5 时失败

No.4 ababcbabcabab i=4,j=1 时失败

No.5 ababcbabcabab i=5,j=1 时失败

No.6 ababcbabcabab abcac 成功



```
int index-1(sstring s,sstring t,int pos)
```

```
{ int i,j,k;
  int n=s.length;
  int m=t.length;
  for(i=1;i<=n-m+1;i++)
  { j=1;k=i;
    while(j<=m && s.ch[k]==t.ch[j])
    { k++;j++; }
    if (j>m) return i;
  }
  return -1;
}
```



8

```
int index-2(sstring s,sstring t,int pos)
```

```
{ int i=1,j=1;
  int n=s.length;
  int m=t.length;
  while (i<=n)&&(j<=m)
  { if (s.ch[i]==t.ch[j])
    { i++;j++; }
    else
    { i=i-j+1;j=1; }
  }
  if (j>m) return i-j+2;
  else return -1;
}
```



9

KMP 算法：改进的模式匹配算法



为什么BF算法时间性能低？

- 在每趟匹配不成功时存在大量回溯，没有利用已经部分匹配的结果。

如何在匹配不成功时主串不回溯？

- 主串不回溯，模式就需要向右滑动一段距离。

如何确定模式的滑动距离？

- 利用已经得到的“部分匹配”的结果
- 将模式向右“滑动”尽可能远的一段距离(next[j])后，继续进行比较

— 出发点：利用前面匹配的结果，进行无回溯匹配

No.1 ababcabcabab i=3,j=3 时失败

No.2 ababcbabcabab i=2,j=1 时失败

No.3 ababca**b**cacbab i=7,j=5 时失败

No.4 ababcbabcabab i=4,j=1 时失败

No.5 ababcbabcabab i=5,j=1 时失败

No.6 ababcbabcabab abcac 成功



11

KMP 算法



思考：

- 假定：主串为 $S_1S_2...S_n$
- 模式串为 $P_1P_2...P_m$
- 当主串中的第 i 个字符和模式串中的第 j 个字符出现不匹配，主串中的第 i 个字符应该和模式串中的哪个字符匹配（无回溯，主串指针不回溯，模式串右移动）？

12

• 进一步思考

- 假定主串中第*i*个字符与模式串第*j*个字符相比较失败，则应有 $S_i \neq P_j$

$S_{i-j+1} S_{i-j+2} \dots S_{i-k+1} S_{i-k+2} \dots S_{i-1} S_i$

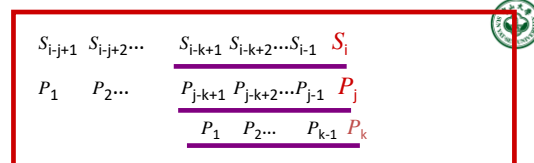
$P_1 P_2 \dots P_{j-k+1} P_{j-k+2} \dots P_{j-1} P_j$

$$k < j$$

$S_{i-k+1} S_{i-k+2} \dots S_{i-1} S_i$
 $P_1 P_2 \dots P_{k-1} P_k$ 成立

S_i 与 P_k 进行比较

13



- 而根据已有的匹配，有

$$P_{j-k+1} P_{j-k+2} \dots P_{j-1} = S_{i-k+1} S_{i-k+2} \dots S_{i-1}$$

- 因此

$$P_{j-k+1} P_{j-k+2} \dots P_{j-1} = P_1 P_2 \dots P_{k-1}$$

- 因此*k*值只和*P*以及*j*有关，定义为Next[j]

意义：当*P*₁比较失败时，右移动模式串，让*P*_{*k*}与当前*S*_{*i*}元素进行比较

14

• Next[j]的定义

$$\text{Next}[j] = \begin{cases} 0, & j=1 \text{ 时} \\ \text{Max}\{k \mid 1 < k < j \text{ and } P_1 P_2 \dots P_{k-1} = P_{j-k+1} \dots P_{j-1}\}, & \\ 1, \text{其它情况} \end{cases}$$

Next数组的实质是模式串中的最长相同的前缀和后缀。
 $P_1 P_2 \dots P_{k-1} = P_{j-k+1} \dots P_{j-1}$

<i>j</i>	1	2	3	4	5	6	7	8
Next[j]	0	1	1	2	2	3	1	2

15

<i>j</i>	1	2	3	4	5
Next[j]	0	1	1	1	2

No.1 $ababcbacbacbab$
 $abacac$ $i=3, j=3$ 时失败

No.2 $ababcbacbacbab$
 $abacac$ $i=7, j=5$ 时失败

No.3 $ababcbacbacbab$
 $abacac$ 成功

16

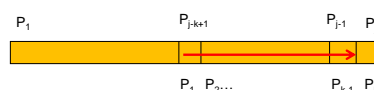
```
int index(ssstring s,ssstring t,int pos)
```

```
{ int i,j;
  int n=s.length;
  int m=t.length;
  while (i<=n)&&(j<=m)
  { if ( s.ch[i]==t.ch[j]) (i=0) || ( s.ch[i]==t.ch[j])
    { i++; j++; }
    else
    { i=i-j+1; j=1; } j=Next[j];
  }
  if (j>m) return i-j+2;
  else return -1;
}
```

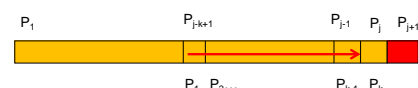
17

计算Next数组的方法

- 1) Next[1]=0;
- 2) 设 Next[j]=k; 则意味着
 $P_{j-k+1} P_{j-k+2} \dots P_{j-1} = P_1 P_2 \dots P_{k-1}$

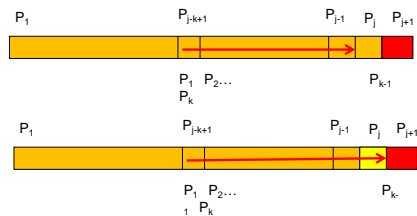


- 3) 求Next[j+1]



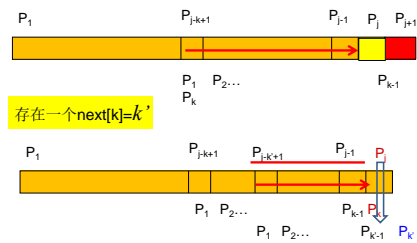
18

3) 求Next[j+1]

当 $P_k = P_j$ 时 $P_{j-k+1}P_{j-k+2}\dots P_{j-1}$ $P_j = P_1P_2\dots$ $P_{k-1}P_k$ 

Next[j+1]=Next[j]+1;

19

当 $P_k \neq P_j$ 时 $P_{j-k+1}P_{j-k+2}\dots P_{j-1}$ $P_j \neq P_1P_2\dots$ $P_{k-1}P_k$ 当 $P_j = P_{k'}$, Next[j+1]=Next[k']+1;

同理，若 $P_j = P_{k'}$ ，则将模式串继续向右滑动至模式串的第next[k']个字符与 P_j 对齐，
，依次类推，直到 P_j 和模式串中某个字符匹配成功或者不存在 k' ($1 < k' < j$)，则
 Next[j+1]=1.

20

Next数组的无回溯匹配计算

```
void Makenext(String p,int *pNext)
{
    int i,k; i=1;k=0;Next[1] = 0;
    while(i<Length(p))
    {
        if (k==0) || (p[i]==p[k])
        { j=j+1; k=k+1;    Next[j]=k;  }
        else k=Next[k];
    }
}
```

21