

字符串 String

1 基本信息

1.1 介绍

字符串，特殊的 线性表，即元素为 **char** 的线性表

$n (\geq 0)$ 个字符的有限序列， $n \geq 1$ 时，一般记作 $S: "c_0c_1c_2...c_{n-1}"$

- S 是串名
- $"c_0c_1c_2...c_{n-1}"$ 是串值
- c_i 是串中的字符
- N 是串长（串的长度）：一个字符串所包含的字符个数
- **空串**：长度为 **0** 的串，它不包含任何字符内容（注意与空格串“ ”的区别）

1.2 特殊线性结构

| | 基本线性结构 | 字符串 |
|------|--------------------|----------------------|
| 数据对象 | 无特殊限制 | 串的数据对象为字符集 |
| 基本操作 | 线性表的大多以“单个元素”为操作对象 | 串通常以“串的整体（子串）”作为操作对象 |
| 存储方式 | 数组、链表等 | 数组、链表等 |

1.3 子串

子串的定义

假设 s_1, s_2 是两个串：

$$s_1 = a_0a_1a_2 \dots a_{n-1}$$

$$s_2 = b_0b_1b_2 \dots b_{m-1}$$

其中 $0 \leq m \leq n$, 若存在整数 $i (0 \leq i \leq n - m)$

使得 $b_j = a_{i+j}, j = 0, 1, \dots, m - 1$ 同时成立

则称串 s_2 是串 s_1 的子串， s_1 为串 s_2 的主串，或称 s_1 包含串 s_2

特殊子串

- **空串** 是任意串的子串
- 任意串 S 都是 **S 本身** 的子串
- **真子串**：非空且不为自身的子串

1.4 字符串的基本操作

假设

s1= "SUSTech"

s2 = "10"

| 操作 | 解释 | 示例 |
|-----------|----------------|--|
| 拼接append | 向字符串后面添加 | s1.append(s2) s1 = "SUSTech10" |
| 赋值assign | 把内容赋值给字符串 | s1.assign(s2) s1 = "10" |
| 插入insert | 插入到字符串的某个位置 | s1.insert(s2,0) s1 = "10SUSTech" |
| 删除erase | 清空字符串的内容 | s1.erase() s1 = "" |
| 替换replace | 替换子字符串某一特定位置的值 | s1.replace('s',0) s1 = "sUSTech" |
| 交换swap | 交换字符串的值 | s1.swap(s2) s1 = "10", s2 = "SUSTech" |
| 查找find | 找到字符串里面某个特定的字符 | s1.find(0) char_find = 'S' |

🔗问题：求子串数量

Give string s="SUSTechCS203", how many sub string it has?

$s.len = 12$
因为包括注意空字符串
 $substring = \frac{(1+s.len)(s.len)}{2} + 1(空字符串) = 79$

🔗问题：S1 + S2 == S2 + S1 成立的所有可能条件

Set S1 and S2 as strings, please give the possible conditions for S1+S2 == S2+S1 to be true (where + is the join operation)

- 1.S1和S2均为空串
- 2.S1为空串，S2为任意字符串
- 3.S2为空串，S1为任意字符串
- 4.S1 == S2两个字符串完全相等
- 5.S1 ≠ S2，则长串由整数个短串组成

1.5 字符串的应用

- 文字处理器Word processors

- 病毒扫描Virus scanning

定义产生病毒的软件都是凭借在软件里查找是否有与病毒软件定义的字符串相似的进行匹配

如果出现可以的字符就会报警

- 文本检索Text retrieval

- 自然语言处理Natural language processing

机器人识别、处理语言

- 网络搜索引擎Web search engine

2 字符串的匹配

2.1 介绍

模式匹配 (pattern matching)

用给定的模式P, 在目标字符串T 中搜索与模式P 全同的一个子串, 并求出T 中第一个和P 全同匹配的子串 (简称为“配串”), 返回其首字符位置

输入

- 一个目标对象 T (字符串)
- search pattern P (字符串)

输出

匹配的首字符位置

2.2 应用

文本编辑时的特定词、句的查找

- 字(词)处理技术
- Web搜索
- 桌面搜索

DNA 信息的提取 (计算生物学的应用)

确认是否具有某种结构

在Unix中使用的Grep操作

2.3 穷举法 Brute Force

思路

逐个字符匹配直到全部匹配完成

T =

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| a | a | a | a | a | a | a | a | a | b |
|---|---|---|---|---|---|---|---|---|---|

P =

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| a | a | a | a | a | a | b | | | | |
| | a | a | a | a | a | a | b | | | |
| | | a | a | a | a | a | a | b | | |
| | | | a | a | a | a | a | a | b | |
| | | | | a | a | a | a | a | a | b |

T =

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | b | a | b | a | b | a | b | a | b | a | b | b |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

P =

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | b | a | b | a | b | b | | | | | | |
| | a | b | a | b | a | b | b | | | | | |
| | | a | b | a | b | a | b | b | | | | |
| | | | a | b | a | b | a | b | b | | | |
| | | | | a | b | a | b | a | b | b | | |
| | | | | | a | b | a | b | a | b | b | |
| | | | | | | a | b | a | b | a | b | b |

伪代码

```

1  // 一个目标对象 T (字符串)
2  //search pattern P (字符串)
3  Algorithm: Brute Force(T, P):
4  n ← len(T)
5  m ← len(P)
6  for i ← 0 to n-m-1 //注意是n-m-1
7      for j ← 0 to m-1
8          if P[j] != T[i + j] then
9              break
10         if j = m-1 then
11             print('pattern occurs with in' + i)
12         end
13 end

```

时间复杂度分析 $O(mn)$

在最坏的情况下，每一次循环都不成功，则一共要进行比较 $(n-m+1)$ 次

每一次“相同匹配”比较所耗费的时间，是 P 和 T 逐个字符比较的时间，最坏情况下，共 m 次

因此，整个算法的最坏时间开销估计为 $O(mn)$

2.4 Rabin-Karp算法

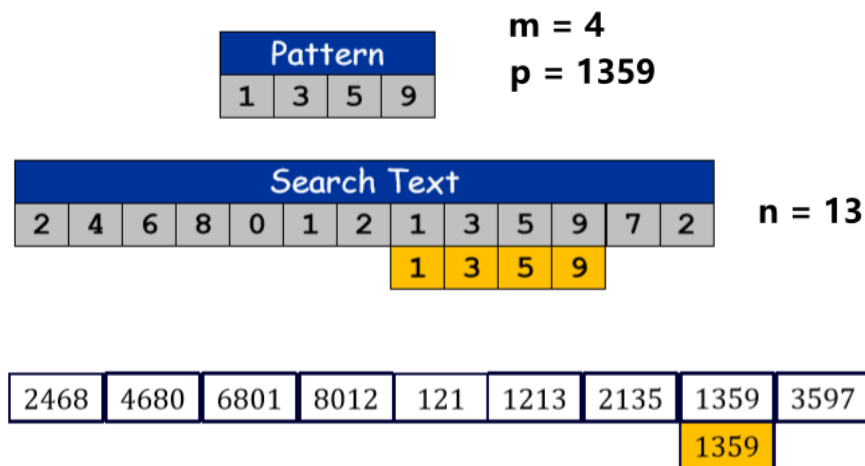
思路: Rabin-Karp

将 P 转换成某一个数字 p

将 T 分成 $n-m+1$ 个数组, $t[0], t[1], \dots, t[n-m]$

将 p 与 $t[i]$ 依次比较

如果 $p = t[i]$ 则 P 出现



思路: 将字符转换为数字 $O(nm) \rightarrow O(n)$

假设条件

- 目标对象: T
- String Pattern: P
- 目标对象的转换数字每个小隔间: t
- String Pattern 转换成数字: p
- 目标对象的总长度: n
- String Pattern 的总长度: m
- 字母表: Σ
- 字母表大小 (全体字符的个数为) d
- 取模的模数 q

转换数字 (进制转换) $O(1)$

$$p = P[m-1] + d(P[m-2] + d(P[m-3] + \dots + d(P[1] + dP[0])))$$

$$p = d^{m-1}P[0] + d^{m-2}P[1] + \dots + d^1P[m-2] + d^0P[m-1]$$

当 m 很大的时候, p 可能太大了, 把 p 取模

$$p = P[m-1] + dP[m-2] + \dots + d^{m-1}P[0] \mod q$$

每一个 $t[i]$ 的转换 $O(1)$ (可以用十进制推算)

$$t[0] = d^{n-1}T[0] + d^{n-2}T[1] + \dots + d^1T[m-2] + d^0T[m-1]$$

$$t[i+1] = d(t[i] - d^{m-1}T[i]) + T[i+m]$$

如果 d^{m-1} 太大, t 太大, 还可以进行取模

$$t[i+1] = d(t[i] - hT[i]) + T[i+m] \mod q$$

$$\text{where } h = d^{m-1} \mod q$$

时间复杂度分析

$t[0] \rightarrow t[n - m]$

故时间复杂度为 $O(n-m)$

伪代码：Rabin-Karp

注意

因为我们取模了！！

即使 $p \equiv t[i] \pmod{q}$ ，这也不能说明 $p = t[i]$

还需要进一步的验证

```
1  Algorithm:Rabin-Karp(T,P,d,q):
2
3  n ← len(T)
4  m ← len(P)
5  h ← d^{m-1} mod q
6  p ← 0
7  t[0] ← 0
8
9  //计算初始p和t0
10 for j ← 0 to m-1
11     p ← (d * p + P[j]) mod q
12     t[0] ← (d * t[0] + T[j]) mod q
13 end
14 //顺序比较
15 for i ← 0 to n-m
16     //不相同，必然有错
17     if p != t[i] then
18         t[i+1] = (d * (t[i] - h * T[i]) + T[i+m]) mod q
19         //相同，还需要进一步对比
20     else
21         //内部比较出现不同
22         for k ← 0 to m-1
23             if P[k] != T[i + k] then
24                 break
25         end
26         //内部比较完全相同
27         if k = m-1 then
28             print("pattern occurs with in" + i)
29 end
```

时间复杂度分析 $O(mn)$

在最坏的情况下，每一次匹配都成功，则一共要进行比较 $(n-m+1)$ 次

每一次“相同匹配”比较所耗费的时间，是 P 和 $t[i]$ 逐个字符比较的时间，最坏情况下，共 m 次

因此，整个算法的最坏时间开销估计为 $O(mn)$

2.5 有限状态自动机FSA

有限状态自动机的定义

Q : 一系列状态集

$q_0 \in Q$: 初始状态

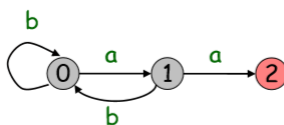
$A \subseteq Q$: 结束状态 (注意是集合)

Σ : 输入字母表

δ : 状态转移函数 (数量为 $Q \times \Sigma$)

| | 0 | 1 |
|---|---|---|
| a | 1 | 2 |
| b | 0 | 0 |

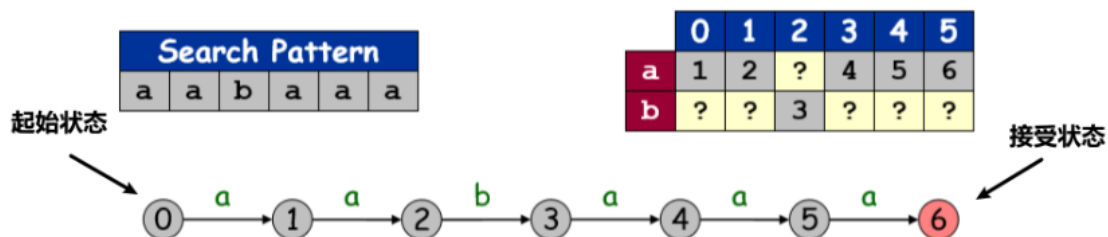
状态转移函数表



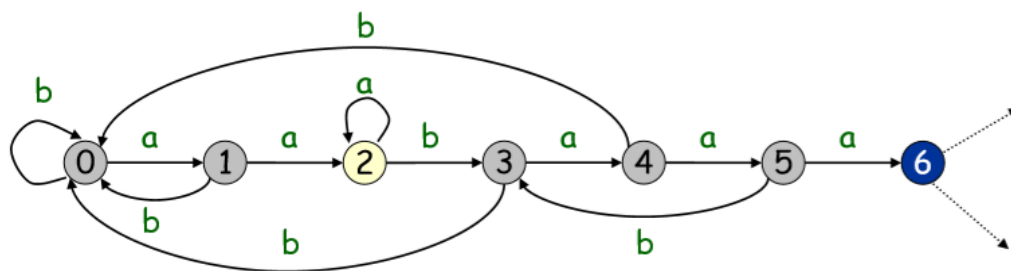
有限状态自动机的图示

思路

核心就是填表



如何填写状态转移表/如何画状态转移函数



问题：绘制状态机

Build FSA for aabaaabb

$Q = \{0, 1, 2, 3, 4, 5, 6, 7, 8\}$

$q_0 = 0$

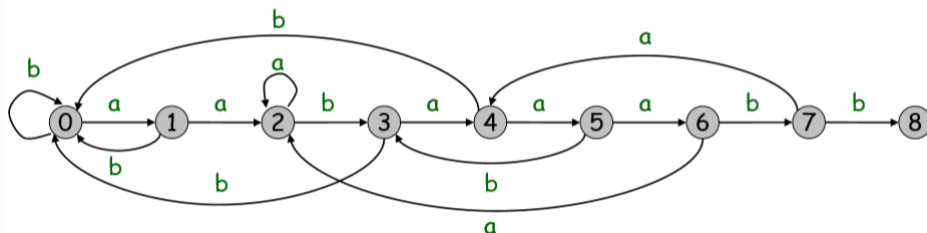
$A = \{8\}$

$\Sigma = \{a, b\}$

状态转移函数如下图所示

| Search Pattern | | | | | | | |
|----------------|---|---|---|---|---|---|---|
| a | a | b | a | a | a | b | b |

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



伪代码

```

1  Algorithm: FSA(T, P):
2    n ← len(T)
3    m ← len(P)
4    δ ← Transition(P, Σ) // 建立δ的时间复杂度是O(mΣ)，但是通常比较小，所以可以估计为O(m)
5    q ← 0
6    for i ← 0 to n-1
7      //新的状态 = 当前状态+新遇到的数
8      q ← δ(q, T[i])
9      if q = m
10         print('pattern occurs with in' + (i - m))
11    end

```

时间复杂度分析 $O(n + m\Sigma) \rightarrow O(n + m)$

建立 δ 的时间复杂度是 $O(m\Sigma)$ ，但是通常比较小，所以可以估计为 $O(m)$

遍历T的时间是 $O(n)$

因为两者是顺序建立的，所以总时间复杂度为 $O(n + m\Sigma) \rightarrow O(n + m)$

FSA用于KMP算法

如果匹配，进入下一个状态

只需要跟踪字符不匹配时去哪里

- 如果状态j中的字符不匹配，进入状态next[j]

2.6 KMP算法

思路

前缀串($P[0, \dots]$)

- 字符串S的aprefix是S的子字符串发生在S的开头
- $P[0, \dots, 0] = \text{"H"}$ (note that $P[0] = \text{"H"}$, $P[0, \dots, 1] = \text{"He"}$, $P[0, \dots, 4] = \text{"Hello"}$, 记为: **$P[0, \dots]$**)

后缀串($P[\dots, m]$)

- 字符串S的后缀是发生在S末尾的S的子字符串

- $P[9,...,9] = "3"$, $P[7,...,9] = "203"$, $P[5,...,9] = "CS203"$, 记为: $P[... ,m]$

通常, FSA的构造是为了让状态数告诉我们有多少前缀P被匹配

FSA的转移函数

查找P的最长前缀也是T的后缀 $[... , i]$, 表示为k, 即, $P[1, ..., k] = T[i-k, ..., i]$

(注意不包括 $P[0]$, 直接从1开始)

读下一个字母 $"k+1"$, 有两种过渡

- $P[k+1] = T[i+1]$, 它匹配, 继续
- 否则,它是不匹配的,去 $\delta(k, T[m+1])$

KMP就是相当于算出来 $\delta(k, T[m+1])$, 其它都不管

❓问题: 填表 $P[i]$ 和 $\pi[i]$

Example: $P = "ababaca"$

本质上求 $a[0,...,k]$ 和 $a[n-k+1...,n-1]$ 有没有匹配

前缀最后一个不能到最后一个 $a[n-1]$

后缀最初一个不能到最前面一个 $a[0]$

| i | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---|---|---|---|---|---|---|
| $P[i]$ | a | b | a | b | a | c | a |
| $\pi[i]$ | 0 | 0 | 1 | 2 | 3 | 0 | 1 |

ababa

a-a

aba-aba 2个

abab

ab

| Seq | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---|----|-----|------|-------|--------|---------|
| $P[i]$ | a | b | a | b | a | c | a |
| x前 | x | a | ab | aba | abab | ababa | ababac |
| x | a | ab | aba | abab | ababa | ababac | ababaca |
| $\pi[i]$ | 0 | 0 | 1 | 2 | 3 | 0 | 1 |

我们看**x**行的字符串, 看它的前缀(不包括最后一个)和后缀(不包括最初一个)可以回到**x**前行的哪一个状态

假如

| Seq | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------|---|---|---|---|---|---|---|---|---|---|----|
| T | a | b | a | b | b | a | b | a | b | a | b |
| P[i] | a | b | a | b | a | c | a | | | | |
| $\Pi[i]$ | 0 | 0 | 1 | 2 | 3 | 0 | 1 | | | | |

在Seq=0的时候, $T[0]=P[0]$, TP指针均右移

在Seq=1的时候, $T[1]=P[1]$, TP指针均右移

在Seq=2的时候, $T[2]=P[2]$, TP指针均右移

在Seq=3的时候, $T[3]=P[3]$, TP指针均右移

在Seq=4的时候, $T[4] \neq P[4]$

P往回看一个 $\Pi[3]=2$

在Seq=2的时候, $T[4] \neq P[2]$

P往回看一个 $\Pi[1]=0$

在Seq=0的时候, $T[4] \neq P[0]$

没有办法了, T指针右移动

在Seq=0的时候, $T[5]=P[0]$, TP指针均右移

在Seq=1的时候, $T[6]=P[1]$, TP指针均右移

在Seq=2的时候, $T[7]=P[2]$, TP指针均右移

在Seq=3的时候, $T[8]=P[3]$, TP指针均右移

在Seq=4的时候, $T[9]=P[4]$, TP指针均右移

在Seq=5的时候, $T[10] \neq P[5]$

P往回看一个 $\Pi[4]=3$

在Seq=3的时候, $T[10]=P[3]$, TP指针均右移

...

在Seq=n的时候, $T[n]=P[6]$, 输出n-6的值为连续值, T右移, P去 $\Pi[6]=1$ 处

在Seq=n+1的时候, 比较 $T[n+1]$ 与 $P[1]$

...

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|----------|---|----|-----|------|-------|--------|---------|----------|
| T | a | b | c | d | a | b | c | a |
| x | a | ab | abc | abcd | abcda | abcdab | abcdabc | abcdabca |
| $\pi[i]$ | 0 | 0 | 0 | 0 | 1 | 2 | 3 | 1 |

$\pi[i]$ 的意义也在于在x行中到此为止的字符串的最长公共前后缀的大小是多少

❗伪代码—— $\pi[i]$ 的填写

```

1  Algorithm:NextArray(P):
2
3  m ← len(p)
4  Let  $\pi[1,...,m]$  be a new array
5   $\pi[1] = 0$ 
6  k ← 0//初始状态
7
8  for j ← 2 to m
9
10     while(k > 0 && P[k+1] != P[j])
11         k ←  $\pi[k]$ 
12     end
13
14     if P[k+1] == P[j] then
15         k ← k + 1
16
17      $\pi[j] = k$ 
18
19 end
20
21 return  $\pi[]$ 

```

❗伪代码——KMP实现

```

1  Algorithm:KMP(T,P):
2  n ← len(T)
3  m ← len(P)
4  q ← 0//初始状态
5
6  for i ← 1 to n
7
8     while q > 0 and P[q+1] != T[i]
9         q ←  $\pi[q]$ 
10    end
11
12    if P[q+1] = T[i] then
13        q ← q + 1
14
15    if q = m
16        print("pattern occurs with shift" + (i - m))
17        q =  $\pi[q]$ 

```

18

19

end