

算法分析与设计

字符串匹配

主要内容

- Brute Force方法
- Rabin-Karp算法
- Knuth-Morris-Pratt算法

字符串匹配问题

- $T[1..n]$ 是一个长度为 n 的数组
- $P[1..m]$ 是一个长度为 m 的数组, $m \leq n$
- P 和 T 中的元素都是属于有限字母表 Σ 的字符
- 字符数组 T 和 P 通常称为字符串
- 字符串匹配问题可以理解为在 T 中寻找所有与 P 匹配的位置

Brute Force方法

●思想

- 用一个循环找出所有匹配之处（有效位移），该循环对 $n-m+1$ 个可能的每一个 s 值进行检查，看 $P[1..m]=T[s+1..s+m]$ 是否满足

Brute Force方法

NAIVE-STRING-MATCHER(T, P)

1 $n \leftarrow \text{length}[T]$

2 $m \leftarrow \text{length}[P]$

3 for $s \leftarrow 0$ to $n - m$

4 do if $P[1 .. m] = T[s + 1 .. s + m]$

5 then print "Pattern occurs with shift" s

Brute Force方法

●时间复杂度分析

- 在最坏情况下，运行时间是 $\Theta((n-m+1)m)$

●实例

- 文本 $T=000010001010001$
- 模式 $P=0001$

Rabin-Karp算法

●基本思想

- 将字符串对应于数值，字符串匹配问题即可归结为数值是否相等的问题
- 字符串如何对应于数值？

Rabin-Karp算法

●基本思想

- 为便于说明，假设 $\Sigma = \{0, 1, 2, \dots, 9\}$ ，这样每个字符都是一个十进制数字（可以推广到一般情况，可以假定每个字符都基数为 $d = |\Sigma|$ 表示法中的一个数字）
- 可以用一个长度为 k 的十进制数表示由 k 个字符组成的字符串
- 模式 $P[1..m]$ ， p 为其对应的十进制数的值
- 对于给定的文本 $T[1..n]$ ，用 t_s 表示其长度为 m 的子串 $T[s+1..s+m]$ 对应的十进制数的值
- $t_s = p$ 当且仅当 $T[s+1..s+m] = P[1..m]$

Rabin-Karp算法

- 如果可以在 $\Theta(m)$ 时间内计算 p 的值，并在 $\Theta(n-m+1)$ 时间内计算出所有 t_s 的值，那就可以通过把 p 值和每个 t_s 值进行比较，就可以在 $\Theta(m) + \Theta(n-m+1) = \Theta(n)$ 的时间内，求出所有的匹配位置（即有效位移，此处不考虑 p 和 t_s 是很大的数的情况）

Rabin-Karp算法

- 可以运用霍纳法则(Horner's rule)在 $\Theta(m)$ 的时间内计算 p 的值, 类似地, 也可以在 $\Theta(m)$ 时间内根据 $T[1..m]$ 计算 t_0 的值
- 对于 $d=10$, 为在 $\Theta(n-m)$ 时间内计算出 t_1, t_2, \dots, t_{n-m} , 可在常数时间内根据下式计算
- $t_{s+1} = 10(t_s - 10^{m-1}T[s+1]) + T[s+m+1]$

Rabin-Karp算法

●存在的问题

- p 和 t_s 的值很大时怎么办？
- 取模运算，选择 q ，对 p 和 t_s 进行取模运算，可在在 $\Theta(m)$ 时间内计算出 p 取模后的值，可在 $\Theta(n-m+1)$ 时间内计算出 t_s 取模后的值
- 通常 q 为一个素数
- 取模，带来的一个问题就是存在伪命中点，也就是取模后的值相等，但实际上并不相等；可以进行进一步测试，显式检查 $P[1..m]$ 和 $T[s+1..s+m]$ 是否相等

Rabin-Karp算法

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

1 $n \leftarrow \text{length}[T]$

2 $m \leftarrow \text{length}[P]$

3 $h \leftarrow d^{m-1} \bmod q$

4 $p \leftarrow 0$

5 $t_0 \leftarrow 0$

6 **for** $i \leftarrow 1$ **to** m \triangleright Preprocessing.

7 **do** $p \leftarrow (dp + P[i]) \bmod q$

8 $t_0 \leftarrow (dt_0 + T[i]) \bmod q$

9 **for** $s \leftarrow 0$ **to** $n - m$ \triangleright Matching.

10 **do if** $p = t_s$

11 **then if** $P[1 .. m] = T[s + 1 .. s + m]$

12 **then** print "Pattern occurs with shift" s

13 **if** $s < n - m$

14 **then** $t_{s+1} \leftarrow (d(t_s - T[s + 1])h + T[s + m + 1]) \bmod q$

Rabin-Karp算法

●复杂度分析

- 预处理时间为 $\Theta(m)$ ，匹配时间最坏情况下为 $\Theta((n-m+1)m)$
- 在实际应用中，有效位移很少（可能只有常数 c 个），算法的期望匹配时间为 $O((n-m+1)+cm)=O(n+m)$ 加上处理伪命中点所需的时间

Knuth-Morris-Pratt算法（KMP算法）

●基本思想

- 根据模式P的信息，在某个位置匹配失败时，利用已经得到的“部分匹配”结果将模式向右移动尽可能远的距离

Knuth-Morris-Pratt算法（KMP算法）

- 模式移动的距离由next[]决定
- next[]如何确定？
- next[]只与模式字符串相关，与目标字符串无关

$$next[j] = \begin{cases} 0, & j = 1 \text{ 时} \\ \text{Max}\{k | 1 < k < j \text{ 且 } 'p_1 p_2 \dots p_{k-1}' = 'p_{j-k+1} \dots p_{j-1}'\}, & j \neq 1 \text{ 时} \\ 1, & j \neq 1 \text{ 时 其它情况} \end{cases}$$

Knuth-Morris-Pratt算法（KMP算法）

●对于如下模式串，求next[]值

➤ a b c a b a a

➤ a b a a b c a c

➤ a b a b a a b a b