

# 后缀数组

MengChunlei

February 27, 2021

- 定义
- 一般做法
- 改进做法
- 计数排序
- 基数排序
- 最后的改进
- 常数优化
- height 数组
- 应用

## 1 定义

对于一个长度为  $n$  的字符串  $s : s[0], s[1], \dots, s[n-2], s[n-1]$ , 它对应的后缀数组也是一个长度为  $n$  的数组  $sa$ . 其中  $sa[i]$  表示在所有  $n$  个后缀中, 排第  $i$  名的后缀是哪个 (第一个字母在  $s$  中的位置). 比如  $s = \text{aabaaaab}$ , 那么对应的  $sa$  为:

<i>index</i>	0	1	2	3	4	5	6	7
<i>s</i>	a	a	b	a	a	a	a	b
<i>sa</i>	3	4	5	0	6	1	7	2

第 0 名: aaaab

第 1 名: aaab

第 2 名: aab

第 3 名: aabaaaab

第 4 名: ab

第 5 名: abaaaab

第 6 名: b

第 7 名: baaaab

另外还有一个数组称之为 rank 数组  $rk$ , 它的长度也是  $n$ ,  $rk[i]$  表示以  $s[i]$  开始的后缀的排名. 那么有以下关系:  $sa[rk[i]] = rk[sa[i]] = i$

## 2 一般做法

很容易想到  $n^2 \log(n)$  的做法. 代码如下:

Listing 1: Normal

```
1 std::vector<int> ComputeSA1(const std::string &s) {
2     int n = static_cast<int>(s.size());
3     std::vector<int> sa(n);
4     for (int i = 0; i < n; ++i) {
5         sa[i] = i;
6     }
7     std::sort(sa.begin(), sa.end(), [&](int p1, int p2) {
8         for (; p1 < n && p2 < n && s[p1] == s[p2]; ++p1, ++p2);
9         return p1 < n && p2 < n ? s[p1] < s[p2] : p1 == n;
10    });
11    return sa;
12 }
```

### 3 改进做法

优化的思路是这样的: 设  $rk_w[i]$  表示以  $i$  开始的长度为  $w$  的子串的排名, 即  $s[i] \rightarrow s[i+w-1]$ . 那么对于  $rk_{2w}$  来说, 它的计算可以借助于  $rk_w$  来进行. 简答来说, 就是以  $rk_w[i]$  为第一关键字, 以  $rk_w[i+w]$  为第二关键字进行排序, 就可以求出  $rk_{2w}$  数组.

Listing 2: Normal

---

```
1 std::vector<int> ComputeSA2(const std::string &s) {
2     int n = static_cast<int>(s.size());
3     std::vector<int> sa(n);
4     std::vector<int> rk(n);
5     for (int i = 0; i < n; ++i) {
6         sa[i] = i;
7         rk[i] = s[i];
8     }
9     for (int w = 1; w < n; w <= 1) {
10        std::sort(sa.begin(), sa.end(), [&](int x, int y) {
11            return rk[x] == rk[y]
12                ? (x + w < n && y + w < n ? rk[x + w] < rk[y + w] : x + w >= n)
13                : rk[x] < rk[y];
14        });
15        std::vector<int> rk_tmp = rk;
16        rk[sa[0]] = 0;
17        for (int p = 0, i = 1; i < n; ++i) {
18            if (rk_tmp[sa[i]] == rk_tmp[sa[i - 1]] &&
19                ((sa[i] + w >= n && sa[i - 1] + w >= n) ||
20                 (sa[i] + w < n && sa[i - 1] + w < n &&
21                  rk_tmp[sa[i] + w] == rk_tmp[sa[i - 1] + w]))) {
22                rk[sa[i]] = p;
23            } else {
24                rk[sa[i]] = ++p;
25            }
26        }
27    }
28    return sa;
29 }
```

---

这个算法的复杂度为  $O(n \log^2(n))$