# String Processing Algorithms 2015 - Week 3 Exercises

## Rodion Efremov

### November 9, 2015

## Exercise 1

### Solution

The time complexity mentioned in the Theorem 1.27 is $\mathcal{O}(||\mathcal{R}|| + m\sigma)$. All we need to do is to modify the COUNTING-SORT procedure:

Now, the desired LSD radix sort for variable-length strings is

## Exercise 2

Use the lcp comparison technique to modify the standard insertion sort algorithm so that it sorts strings in $\mathcal{O}(\Sigma LCP(\mathcal{R}) + n^2)$ time.

### Solution

**Why do we set LCP to zero?**

Consider the following:

| $i$ | $S_i$ | $\text{LCP}_{\mathcal{R}}[i]$ |
|-----|-------|------------------------------|
| 1 | $aaaba$ | 0 |
| 2 | $abaaa$ | 1 |
| 3 | $aaa$ | 1 |
| 4 | $baaa$ | 0 |
| 5 | $aabba$ | 0 |
| 6 | $abaaa$ | 1 |

## Exercise 3

Give an example showing that the worst case time complexity of string binary search without precomputed lcp information is $\Omega(m \log n)$.

**Algorithm 1:** COUNTING-SORT($\mathcal{R} = \{S_1, S_2, \ldots, S_n\}, \ell$)

---

**1** **for** $i = 0$ **to** $\sigma - 1$ **do**
**2** $\quad$ $C[i] = 0$

**3** $s = 0$
**4** **for** $i = 1$ **to** $n$ **do**
**5** $\quad$ **if** $|S_i| < \ell$ **then**
**6** $\quad\quad$ $s = s + 1$
**7** $\quad$ **else**
**8** $\quad\quad$ $C[S_i[\ell]] = C[S_i[\ell]] + 1$

**9** $sum = s$
**10** **for** $i = 0$ **to** $\sigma - 1$ **do**
**11** $\quad$ $tmp = C[i]$
**12** $\quad$ $C[i] = sum$
**13** $\quad$ $sum = sum + tmp$

**14** $p = 0$
**15** **for** $i = 1$ **to** $n$ **do**
**16** $\quad$ **if** $|S_i| < \ell$ **then**
**17** $\quad\quad$ $J[p] = S_i$
**18** $\quad\quad$ $p = p + 1$
**19** $\quad$ **else**
**20** $\quad\quad$ $J[C[S_i[\ell]]] = S_i$
**21** $\quad\quad$ $C[S_i[\ell]] = C[S_i[\ell]] + 1$

**22** $\mathcal{R} = J$

---

**Algorithm 2:** LSDRadixSort($\mathcal{R} = \{S_1, S_2, \ldots, S_n\}$)

---

**1** $m = \max_{i=1,2,\ldots,n}\{|S_i|\}$
**2** **for** $\ell = m$ **to** $1$ **do**
**3** $\quad$ COUNTING-SORT($\mathcal{R}, \ell$)

---

**Algorithm 3:** INSERTIONSORT($\mathcal{R}, \text{LCP}_{\mathcal{R}}$)

---

**1** **for** $i = 2$ **to** $n$ **do**
**2** $\quad$ $s = S_i$
**3** $\quad$ $j = i - 1$
**4** $\quad$ **if** $LCP_{\mathcal{R}}[i - 1] = 0$ **then**
**5** $\quad\quad$ $\text{LCP}_{\mathcal{R}}[i] = 0$
**6** $\quad$ **while** $j > 0$ **and** LCPCOMPARE($S_j, s, LCP_{\mathcal{R}}[i]$) $> 0$ **do**
**7** $\quad\quad$ $S_{j+1} = S_j$
**8** $\quad\quad$ $j = j - 1$
**9** $\quad$ $S_{j+1} = s$

---

## Solution

Suppose the sorted list of strings is

$$\langle \overbrace{aaa\ldots aaa}^{m} a_1,$$

$$\overbrace{aaa\ldots aaa}^{m} a_2,$$

$$\ldots$$

$$\overbrace{aaa\ldots aaa}^{m} a_n \rangle$$

and the string to search for is $\overbrace{aaa\ldots aaa}^{m} a'$, where $a' \neq a_i$ for any $i \in \{1, 2, \ldots, m\}$ and $a_1 \leq a_2 \leq \cdots \leq a_n$. Now as there is no match, the binary search will do $\Omega(\log n)$ string comparisons, and as we assume the naïve implementation of the string comparison, each comparison will have to iterate through $m$ first characters $a$ before getting to the last character that fails the search, which leads to the worst case time complexity of $\Omega(m \log n)$.

# Exercise 4

Let $S[0..n)$ be a string over an integer alphabet. Show how to build a data structure in $\mathcal{O}(n)$ time and space so that afterwards the Karp-Rabin hash function $H(S[i..j))$ for the factor $S[i..j)$ can be computed in constant time for any $0 \leq i \leq j \leq n$.

## Solution

We will need the following identity:

$$H(B) = (H(AB) - H(A) \cdot r^{|B|}) \mod q,$$

where

$$H(S[0..n)) = \left( \sum_{i=0}^{n-1} S[i] r^{n-1-i} \right) \mod q.$$

Since we store only the hash values of all the prefixes of the input integer string (and a matrix indexed by the starting and ending indices), we need an efficient way for finding hash values of the factors of the input integer string. If we want to ask for a hash value of a factor $B$, which is preceded by a factor $A$, we can compute efficiently $H(B)$ simply by looking up the hashes for $H(A)$ and $H(AB)$. Also we need to cache the values $1, r, r^2, \ldots, r^n$.

Assuming $r = 37, q = 11$ and $S[0..5) = \langle 2, 1, 1, 4, 3 \rangle$, we consider all possible prefixes of $S$:

| $i$ | $S[0..i)$ | $H(S[0..i))$ |
|---|---|---|
| 0 | $\varepsilon$ | 0 |
| 1 | $\langle 2 \rangle$ | 2 |
| 2 | $\langle 2, 1 \rangle$ | 9 |
| 3 | $\langle 2, 1, 1 \rangle$ | 4 |
| 4 | $\langle 2, 1, 1, 4 \rangle$ | 9 |
| 5 | $\langle 2, 1, 1, 4, 3 \rangle$ | 6 |

The actual query routine follows  The preprocessing step is

---

**Algorithm 4:** QUERY$(E, F, i, j, q)$

---

**1** $l = j - i$
**2** $f = F[l]$
**3** $a = E[j]$
**4** $b = E[i]$
**5** $r = (a - bf) \mod q$
**6** **if** $r < 0$ **then**
**7** $\quad$ **return** $r + q$
**8** **else**
**9** $\quad$ **return** $r$

---

**Algorithm 5:** PREPROCESS$(S[0..n), r)$

---

**1** $E[0..n] = [0, 0, \ldots, 0]$
**2** $F[0..n] = [0, 0, \ldots, 0]$
**3** $h = 0$
**4** **for** $i = 1$ **to** $n$ **do**
**5** $\quad h = h \times r$
**6** $\quad h = h + S[i - 1]$
**7** $\quad E[i] = h$
**8** $f = 1$
**9** **for** $i = 0$ **to** $n$ **do**
**10** $\quad F[i] = f$
**11** $\quad f = f * r$
**12** **return** $(E, F)$

---

# Exercise 5

$\Omega(\Sigma LCP(\mathcal{R}))$ is a lower bound for string sorting for any algorithm if characters can be accessed only one at a time. However, for a small alphabet, it is possible to pack several characters into one machine word.  Then multiple characters can

be accessed simultaneously and treated as if they were a single *super-character*. For example, the string `abbaba` over the alphabet $\Sigma = \{a, b\}$ can be thought of as the string $\{ab, ba, ba\}$ over the alphabet $\Sigma^2$. Algorithms taking advantage of this are called *super-alphabet* algorithms.

Develop a super-alphabet version of MSD radix sort. What is the time complexity.

**Solution**

---

**Algorithm 6:** MSDRADIXSORT$(\mathcal{R}, \ell)$

---

**1** **if** $|\mathcal{R}| < \sigma^2$ **then**
**2** $\quad$ **return** STRINGQUICKSORT$(\mathcal{R}, \ell)$

**3** $\mathcal{R}_\perp = \{S \in \mathcal{R} \colon |S| = \ell\}$
**4** $\mathcal{R} = \mathcal{R} \setminus \mathcal{R}_\perp$
**5** $(\mathcal{R}_{00}, \mathcal{R}_{01}, \ldots, \mathcal{R}_{0,\sigma-1}, \mathcal{R}_{10}, ..., \mathcal{R}_{\sigma-1,\sigma-1}) = $ COUNTINGSORT$(\mathcal{R}, \ell)$
**6** **for** $i = (00)$ **to** $(\sigma - 1, \sigma - 1)$ **do**
**7** $\quad$ MSDRADIXSORT$(\mathcal{R}_i, \ell + 1)$

**8** **return** $\mathcal{R}_\perp \cdot \mathcal{R}_{00} \cdot \cdots \cdot \mathcal{R}_{\sigma-1,\sigma-1}$

---