

String Processing Algorithms 2015 - Week 2

Exercises

Rodion Efremov

October 30, 2015

Exercise 1

Outline algorithms that find the most frequent symbol in a give string.

- (a) for ordered alphabet, and
- (b) for integer alphabet.

The algorithms should be as fast as possible. What are their (worst case) time complexities? Consider also the case where $\sigma \gg n$.

Solution

Algorithm 1: MOSTFREQUENTSMBOL(S)

```
1 let  $f$  be an empty map  $f: \Sigma \rightarrow \mathbb{N}$ 
2  $\mu = \text{nil}$ 
3  $L_\mu = 0$ 
4 for  $i = 1$  to  $|S|$  do
5   if  $S[i]$  is not mapped in  $f$  then
6      $f(S[i]) = 1$ 
7     if  $L_\mu = 0$  then
8        $L_\mu = 1$ 
9        $\mu = S[i]$ 
10  else
11     $f(S[i]) = f(S[i]) + 1$ 
12    if  $L_\mu < f(S[i])$  then
13       $L_\mu = f(S[i])$ 
14       $\mu = S[i]$ 
15 return  $\mu$ 
```

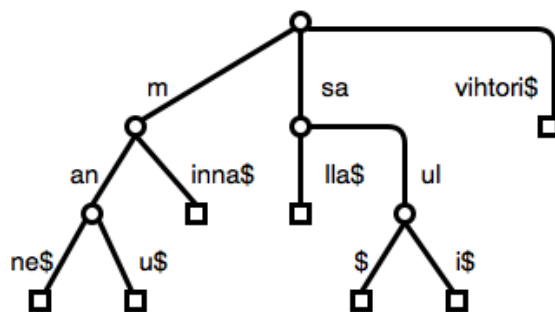
Exercise 2

Let $\mathcal{R} = \{\text{manne}, \text{manu}, \text{minna}, \text{salla}, \text{saul}, \text{sauli}, \text{vihtori}\}$.

- (a) Give the compact trie of \mathcal{R} .
- (b) Give the balanced compact ternary trie of \mathcal{R} .

Solution

(a)



(b)

See the drawing.

Exercise 3

Exercise 4

Prove

- (a) Lemma 1.14: For $i \in [2..n]$, $LCP_{\mathcal{R}}[i] = lcp(S_i, \{S_1, \dots, S_{i-1}\})$.
- (b) Lemma 1.15: $\Sigma LCP(\mathcal{R}) \leq \Sigma lcp(\mathcal{R}) \leq 2 \cdot \Sigma LCP(\mathcal{R})$.

Solution

Both by induction.

(a)

We need an auxiliary lemma first:

Lemma 1 (Distance lemma). *If $S_1 < S_2 < S_3$, $lcp(S_1, S_3) \leq lcp(S_2, S_3)$.*

Proof. The proof is by contradiction: Let $\text{prefix}(A, B)$ be the longest common prefix of A and B . Assume that $l_1 > l_2$, where $l_1 = |\text{prefix}(S_1, S_3)|$ and $l_2 = |\text{prefix}(S_2, S_3)|$. Now $S_1 = a_1 a_2 \dots a_{l_2} \dots a_{l_1} b_1 b_2 \dots$ and $S_2 = a_1 a_2 \dots a_{l_2} b'_1 b'_2 \dots$. Because $S_3 = a_1 a_2 \dots a_{l_2} \dots a_{l_3} c_1 c_2 \dots$, we must have that $S_2 < S_1 < S_3$, which is a contradiction. \square

Example

$S_1 : aaaab$

$S_2 : aaaba$

$S_3 : aaabc$

Base step: Assume $i = 2$. Now $LCP_{\mathcal{R}}[2] = lcp(S_2, S_1) = lcp(S_2, \{S_1\})$.

Induction step: Assume that $LCP_{\mathcal{R}}[i] = lcp(S_i, \{S_1, \dots, S_{i-1}\})$. Now

$$\begin{aligned} lcp(S_{i+1}, \{S_1, \dots, S_i\}) &= \max\{lcp(S_{i+1}, S_i), lcp(S_{i+1}, \{S_1, \dots, S_{i-1}\})\} \\ &\stackrel{DI}{=} lcp(S_{i+1}, S_i), \end{aligned}$$

which concludes the proof.