

58093 String Processing Algorithms (Autumn 2010)

Exercises 1 (11 November)

1. Outline algorithms that find the most frequent symbol in a given string
 - (a) for ordered alphabet, and
 - (b) for integer alphabet.

The algorithms should be as fast as possible. What are their time complexities?

2. The Knuth–Morris–Pratt algorithm differs from the Morris–Pratt algorithm only in the failure function, which can be defined as

$fail_{\text{KMP}}[i] = k$, where k is the length of the longest proper border of $P[0..i]$ such that $P[k] \neq P[i]$, or -1 if there is no such border.

- (a) Compute both failure functions for the pattern *ananassana*.
 - (b) Give an example of a text, where some text character is compared three times by the MP algorithm but only once by the KMP algorithm when searching for *ananassana*.
3. Modify Algorithm 1.6 on the lecture notes to compute $fail_{\text{KMP}}$ instead of $fail_{\text{MP}}$.
 4. Let us analyze the average case time complexity of the Horspool algorithm, where the average is taken over all possible patterns of length m and all possible texts of length n for the integer alphabet $\Sigma = \{0, 1, \dots, \sigma - 1\}$ where $\sigma > 1$. This is the same as the expected time complexity when each pattern and text character is chosen independently and randomly from the uniform distribution over Σ .
 - (a) Show that the average time spent in the loop on line 7 is $\mathcal{O}(1)$.
 - (b) Show that the probability that the shift is shorter than $\min(m, \sigma/2)$ is at most $1/2$.
 - (c) Combine the above results to show that the average time complexity is $\mathcal{O}(n/\min(m, \sigma))$.
 5. Give a *deterministic* finite automaton corresponding to the nondeterministic automaton in Example 1.14 on the lecture notes.
 6. Two string x and y are *rotations* of each other if there exists strings u and v such that $x = uv$ and $y = vu$. For example *abcde* and *deabc* are rotations of each other. Describe a *linear time* algorithm for determining whether given two strings are rotations of each other. (Hint: use a linear time exact string matching algorithm as a subroutine.)