



May 17, 2021

ASSIGNMENT 5

ALGORITHMS FOR SEQUENCE ANALYSIS, SUMMER 2021

Algorithmic Bioinformatics · Prof. Dr. Sven Rahmann

Hand in date: Monday, May 24, before 20:00

Exercise 1: Number of distinct substrings (4 Theory)

Give a linear-time algorithm to compute the total number of distinct substrings of a string s . Think in terms of $s\$$ and the enhanced suffix array. For example, `baaa` has 7 distinct substrings: `a`, `b`, `aa`, `ba`, `aaa`, `baa`, `baaa`.

Exercise 2: Supermaximal repeats (4 Theory)

Give a linear-time algorithm to find all supermaximal repeats of a string s , using its enhanced suffix array and the definitions below. To develop the algorithm, characterize supermaximal repeats in terms of intervals in the enhanced suffix array.

A triple (p, p', m) of two positions $p < p'$ and a length $m \geq 1$ is called

- a *repeated pair* if and only if $s[p : p + m] = s[p' : p' + m] =: w$
(using Python half-open indexing, i.e., the right boundary is not included),
- *right-maximal* if $s[p + m] \neq s[p' + m]$,
- *left-maximal* if $p = 0$ or $s[p - 1] \neq s[p' - 1]$,
- *maximal* if it is both left-maximal and right-maximal.

Then the repeated substring w is called a [left-/right-]maximal *repeat* of length m at p, p' . A *supermaximal repeat* is a maximal repeat that is not a substring of any other repeat.

Example: $s = \text{mambamba}\$$.

Maximal repeats with corresponding maximal repeated pairs:

a: (1, 7, 1), **amba:** (1, 4, 4), **m:** (0, 2, 1), (0, 5, 1).

Supermaximal repeats: **amba** only; the other maximal repeats are substrings of **amba**.

Note that there can be several maximal repeated pairs for the same string; there can be also other (non-maximal) repeated pairs for the string; it suffices that there is at least one maximal repeated pair for a string to call the string a maximal repeat – the repeated pair (2, 5, 1) is not maximal for $w = \text{m}$, but $w = \text{m}$ is still a maximal repeat because other maximal repeated pairs exist for it, such as (0, 2, 1), see above.

It may help to think about this problem in terms of the suffix tree first (which inner nodes are supermaximal repeats?) and then translate your insights to suffix array intervals.