

Text Indexing and Information Retrieval

Übungsblatt 3

Besprechung: 29.10.2018

Aufgabe 1

In der Vorlesung haben wir 2 Details im Linearzeit-Algorithmus für Suffix Arrays übersprungen. Ziel dieser Aufgabe ist es, diese nachzuholen.

- Die Sortierung der S^* -Substrings kann auch ohne einen Standard-String-Sortierer bewerkstelligt werden; dies steht im letzten Absatz vor Thm. 1 im Skript. Führen Sie dieses Verfahren auf beiden in der Vorlesung behandelten Texten aus.
- Betrachten Sie den Text `cacbabc$` und seine S^* -substrings. Der Algorithmus aus der Vorlesung arbeitet hier nicht korrekt. Was ist das Problem? Wie muss der Algorithmus modifiziert werden, so dass er korrekt arbeitet?

Aufgabe 2

- Führen Sie den in der Vorlesung besprochenen Linearzeit-Algorithmus für Suffix Arrays anhand des Beispieltexes $T = \text{antanarivo\$}$ aus.
- Die folgende Grammatik erzeugt einen Text der Länge $n = 2^p$ über dem Alphabet $\{\$, \sigma_1, \dots, \sigma_p\}$: $S \rightarrow T_1\$, T_i \rightarrow T_{i+1}\sigma_i T_{i+1}$ für $i = 1, \dots, p-1$ und $T_p \rightarrow \sigma_p$. Führen Sie den in der Vorlesung besprochenen Linearzeit-Algorithmus für Suffix Arrays auf dem aus $p = 4$ resultierenden Text S aus.
- Was ist das Besondere an dem in Aufgabenteil (b) erzeugten Text in Hinblick auf den in der Vorlesung besprochenen Linearzeit-Algorithmus für Suffix Arrays?

Aufgabe 3

Beweisen die Korrektheit des Linearzeit-Algorithmus für Suffix Arrays formal.

Aufgabe 4 (Praxis)

- a) Auf der Seite <https://sites.google.com/site/yuta256/sais> befinden sich gute Implementierungen des in der Vorlesung besprochenen Linearzeit-Algorithmus für Suffix Arrays. Stellen Sie diesen Code mit seinen Implementierungstricks vor (Hauptdateien `sais.c` bzw. `sais.java`).
- b) Lassen Sie den Algorithmus aus Teil (a) auf möglichst großen Texten von der Seite <http://pizzachili.dcc.uchile.cl/texts.html> laufen und messen Sie Laufzeiten und Speicherplatzbedarf. Wenn möglich, vergleichen Sie die Laufzeit und den Speicherbedarf mit Ihrem in Übungsblatt 2 programmierten $O(n \lg n)$ -Zeit Algorithmus für Suffix-Arrays, und/oder dem dort programmierten $O(n^2 \lg n)$ -Algo.