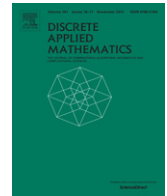




Contents lists available at ScienceDirect

Discrete Applied Mathematics

journal homepage: www.elsevier.com/locate/damComputing covers using prefix tables[☆]Ali Alatabbi^{a,*}, M. Sohel Rahman^b, W.F. Smyth^{a,c,d}^a Department of Informatics, King's College London, United Kingdom^b ALEDA Group, Department of Computer Science & Engineering (CSE), Bangladesh University of Engineering & Technology (BUET), Dhaka-1205, Bangladesh^c Algorithms Research Group, Department of Computing & Software, McMaster University, Canada^d School of Engineering & Information Technology, Murdoch University, Western Australia, Australia

ARTICLE INFO

Article history:

Received 17 October 2014

Received in revised form 17 February 2015

Accepted 11 May 2015

Available online xxxx

Keywords:

Strings

Prefix tables

Covers

Cover array

Indeterminate strings

Linear time

Algorithm

ABSTRACT

An **indeterminate string** $x = x[1..n]$ on an alphabet Σ is a sequence of nonempty subsets of Σ ; x is said to be **regular** if every subset is of size one. A proper substring u of regular x is said to be a **cover** of x iff for every $i \in 1..n$, an occurrence of u in x includes $x[i]$. The **cover array** $\gamma = \gamma[1..n]$ of x is an integer array such that $\gamma[i]$ is the longest cover of $x[1..i]$. Fifteen years ago a complex, though nevertheless linear-time, algorithm was proposed to compute the cover array of regular x based on prior computation of the border array of x . In this paper we first describe a linear-time algorithm to compute the cover array of regular x based on the prefix table of x . We then extend this result to indeterminate strings.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

The idea of a **quasiperiod** or **cover** of a string x was introduced almost a quarter-century ago by Apostolico & Ehrenfeucht [4]: a proper substring u of x such that every position in x lies within an occurrence of u . Thus, for example, $u = aba$ is a cover of $x = ababaababa$. In [5] a linear-time algorithm was described to compute the shortest cover of x ; this contribution was followed by linear-time algorithms to compute

- the shortest cover of every prefix of x [9];
- all the covers of x [17,18];
- all the covers of every prefix of x [16].

A **border** of a string x is a possibly empty proper prefix of x that is also a suffix of x . (Thus a cover of x is necessarily also a border of x .) In the **border array** $\beta = \beta[1..n]$ of the string $x = x[1..n]$, $\beta[i]$ is the length of the longest border of $x[1..i]$. Since for $\beta[i] \neq 0$, $\beta[\beta[i]]$ is the length of a border of x as well as the length of the longest border of $x[1..\beta[i]]$ [2,20], it follows that β provides all the borders of every prefix of x . For example:

$$\begin{array}{cccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \mathbf{x} = & a & b & a & b & a & b & a & a & b & a \\ \mathbf{\beta} = & 0 & 0 & 1 & 2 & 3 & 4 & 5 & 1 & 2 & 3 \end{array} \quad (1)$$

[☆] This work was supported in part by the Natural Sciences & Engineering Research Council of Canada.

* Corresponding author.

E-mail addresses: ali.alatabbi@kcl.ac.uk (A. Alatabbi), msrahman@cse.buet.ac.bd (M. Sohel Rahman), smyth@mcmaster.ca (W.F. Smyth).

As shown in [16], the **cover array** γ has a similar cascading property, giving the lengths of all the covers of every prefix of \mathbf{x} in a compact form:

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \gamma = & 0 & 0 & 0 & 2 & 3 & 4 & 5 & 0 & 0 & 3 \end{array}$$

Here $\mathbf{x}[1..7]$ has covers $\mathbf{u}_1 = \mathbf{x}[1..5] = ababa$ and $\mathbf{u}_2 = \mathbf{x}[1..3] = aba$, while the entire string \mathbf{x} has cover \mathbf{u}_2 . The main result of [16] is an algorithm that computes $\gamma = \gamma[1..n]$ from $\beta = \beta[1..n]$ in $\Theta(n)$ time, while making no reference to the underlying string \mathbf{x} .

The results outlined above all apply to a **regular string** – that is, a string \mathbf{x} such that each entry $\mathbf{x}[i]$ is constrained to be a one-element subset of a given set Σ called the **alphabet**. In this paper we show how to extend these ideas and algorithms to an **indeterminate string** \mathbf{x} – that is, such that each $\mathbf{x}[i]$ can be any nonempty subset of Σ . Observe that every regular string is indeterminate.

The idea of an indeterminate string was first introduced in [12], then studied further in the 1980s as a “generalized string” [1]. Over the last 15 years Blanchet-Sadri has written numerous papers on the properties of “strings with holes” (each $\mathbf{x}[i]$ is either a one-element subset of Σ or Σ itself), together with a monograph on the subject [8]; while other authors have studied indeterminate strings in their full generality, together with related algorithms [6,19,14,15,21–23,10]. In the specific context of this paper, Voráček & Melichar [24] have done pioneering work on the computation of covers and related structures in generalized strings using finite automata.

For indeterminate strings, equality of letters is replaced by the idea of a “match” [14]: $\mathbf{x}[i]$ **matches** $\mathbf{x}[j]$ (written $\mathbf{x}[i] \approx \mathbf{x}[j]$) if and only if $\mathbf{x}[i] \cap \mathbf{x}[j] \neq \emptyset$, while $\mathbf{x} \approx \mathbf{y}$ if and only if $|\mathbf{x}| = |\mathbf{y}|$ and corresponding positions in \mathbf{x} and \mathbf{y} all match. It is important to note that matching is nontransitive: $b \approx \{b, c\} \approx c$, but $b \not\approx c$.

It is [10] that provides the point of departure for our contribution, as we now explain. The **prefix table** $\pi = \pi[1..n]$ of $\mathbf{x}[1..n]$ is an integer array such that $\pi[1] = n$ and, for every $i \in 2..n$, $\pi[i]$ is the length of the longest substring occurring at position i of \mathbf{x} that matches a prefix of \mathbf{x} . Thus, for our example (1):

$$\begin{array}{cccccccccc} & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \mathbf{x} = & a & b & a & b & a & b & a & a & b & a \\ \pi = & 10 & 0 & 3 & 0 & 3 & 0 & 1 & 3 & 0 & 1 \end{array}$$

It turns out [7] that the prefix table and the border array are “equivalent” for regular strings; that is, each can be computed from \mathbf{x} in linear time, and each can be computed from the other, without reference to \mathbf{x} , also in linear time. However, for indeterminate strings, this is not true: the prefix table continues to determine all the borders of every prefix of \mathbf{x} , while the border array, due to the intransitivity of matching, is no longer reliable in identifying borders shorter than the longest one. Consider, for example:

$$\begin{array}{ccc} & 1 & 2 & 3 \\ \mathbf{x} = & a & \{a, b\} & b \\ \beta = & 0 & 1 & 2 \end{array}$$

Here \mathbf{x} does not have a border of length $\beta[\beta[3]] = 1$; on the other hand, $\pi = 320$ correctly identifies all the borders of every prefix of \mathbf{x} .

Moreover, it was shown in [10] that every **feasible** array – that is, every array $\mathbf{y} = \mathbf{y}[1..n]$ such that $\mathbf{y}[1] = n$ and for every $i \in 2..n$, $\mathbf{y}[i] \in 0..n - i + 1$ – is a prefix table of some (indeterminate) string. Thus there exists a many-many correspondence between all possible prefix tables and all possible indeterminate strings. Furthermore, [21] describes an algorithm to compute the prefix table of any indeterminate string, while [3] gives an algorithm to compute a lexicographically least indeterminate string corresponding to a given prefix table.

At this point let us discuss our motivation more precisely. First, realize that to exploit the fullest functionality of a border array of an indeterminate string we need to resort to the extended definition of the border array which in fact requires quadratic space [14,19,6]: unlike the border array of a regular string, which is a simple array of integers, the border array of an indeterminate string is an array of lists of integers. Here at each position, the list gives all possible borders for that prefix. On the other hand, the prefix array, even for the indeterminate string, remains a simple one-dimensional array, just as for a regular string. It thus becomes of interest to make use of the prefix table rather than the border array whenever possible, in order to extend the scope of computations to indeterminate strings.

In Section 2 of this paper, we describe a linear-time algorithm to compute the cover array γ of a regular string \mathbf{x} directly from its prefix table π . Then, Section 3 describes a limited extension of this algorithm to indeterminate strings. Finally, Section 4 outlines future research directions, especially making use of prefix tables to extend the utility and applicability of other data structures to indeterminate strings.

2. Prefix-to-cover for a regular string

In this section we describe our basic $\Theta(n)$ -time Algorithm PCR to compute the cover array $\gamma = \gamma[1..n]$ of a regular string $\mathbf{x} = \mathbf{x}[1..n]$ directly from its prefix table $\pi = \pi[1..n]$. In fact, as noted in the Introduction, γ actually provides all the covers of every prefix of \mathbf{x} . Central to our algorithm are the following definitions:

Download English Version:

<https://daneshyari.com/en/article/4949997>

Download Persian Version:

<https://daneshyari.com/article/4949997>

[Daneshyari.com](https://daneshyari.com)