Introduction

Regular & Indeterminate Strings

Maximal Palindrome Array

Open Problems

References

A New Approach to Regular & Indeterminate Strings

Felipe A. Louza^a Neerja Mhaskar^b W. F. Smyth^{b,c,d}

^a Dept. of Computing and Mathematics, University of Sao Paulo, Brazil

^b Dept. of Computing and Software, McMaster University, Canada

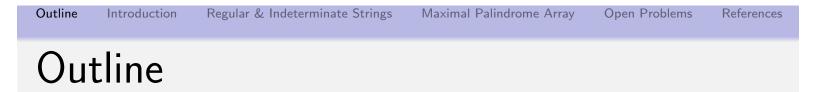
^c Dept. of Informatics, King's College London, UK ^d School of Engineering & Information Technology, Murdoch University, Perth, Australia

LSD & LAW 2019, London, UK

Louza, Mhaskar and Smyth

LSD & LAW 2019. London

Outline - 1 / 26



- Abstract
- Regular and Indeterminate Strings
- Palindromes and Maximal Palindrome Array
- Open Problems

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Outline - 1 / 26



Abstract

We propose a new, more appropriate definition of a regular string; that is, one that is isomorphic to a string whose entries all consist of a single letter. A string that is not regular is said to be indeterminate. We describe an algorithm to determine whether or not a string x is regular and, if so, to replace it by a lexicographically least string string y whose entries are all single letters. We then introduce the idea of a feasible palindrome array MP of a string, and show that every feasible MP corresponds to some (regular or indeterminate) string – perhaps, surprisingly, both! We describe an algorithm that constructs a string x corresponding to given feasible MP, lexicographically least whenever x is regular.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Outline - 2 / 26



Introduction

The idea of a string as something other than a sequence of single letters has been discussed for almost half a century.

In 1974 Fischer & Paterson [FP74] studied pattern-matching on strings x whose entries could be **don't-care** letters; that is, letters matching any single letter in the alphabet Σ on which the string is defined, hence matching every position in x.

In 1987 Abrahamson [Abr87] extended this model by considering pattern-matching on **generalized** strings whose entries could be arbitrary subsets of Σ .

Both of these models have been intensively studied in this century, notably by Blanchet-Sadri ("strings with holes") and Iliopoulos ("degenerate strings").

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Introduction - 3 / 26

Introduction Regular & Inde

Regular & Indeterminate Strings

Maximal Palindrome Array

Array Open Problems

References

Regular & Indeterminate Strings

In this paper we redefine an indeterminate string in a context that we believe captures the idea in a more appropriate way — at once more general and more precise.

A letter ℓ is a finite list of s distinct characters c_1, c_2, \ldots, c_s , each drawn from a set Σ of size $\sigma = |\Sigma|$ called the alphabet.

In the case that Σ is ordered, ℓ is said to be in **normal form** if its characters occur in the ascending order determined by Σ .

The integer $s = s(\ell)$ is called the **scope** of ℓ . For s = 1, ℓ is said to be **regular**, otherwise **indeterminate**.

Two letters ℓ_1, ℓ_2 are said to **match**, written $\ell_1 \approx \ell_2$, if and only if $\ell_1 \cap \ell_2 \neq \emptyset$. In the case that matching ℓ_1 and ℓ_2 are both regular, we may write $\ell_1 = \ell_2$.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 4 / 26

For $n \ge 1$, a string x = x[1..n] is a sequence $x[1], x[2], \ldots, x[n]$ of letters, where n = |x| is the length of x, and every $i \in 1..n$ is a position in x.

If every letter in x is in normal form, then x itself is said to be in **normal form**.

A tuple $T = (i, j_1, j_2)$ of distinct positions i, j_1, j_2 in x such that $x[j_1] \approx x[i] \approx x[j_2]$

is said to be a **triple**.

A triple T is **transitive** if $x[j_1] \approx x[j_2]$, otherwise **intransitive**.

If every triple T in x is transitive, then we say that x is **regular**; otherwise, x is **indeterminate**.

The **scope** of x is given by

$$S(\boldsymbol{x}) = \max_{i \in 1..n} s(\boldsymbol{x}[i]).$$

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 5 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References
---------	--------------	---------------------------------	--------------------------	---------------	------------

Two strings x and y of equal length n are said to be **isomorphic** if and only if for every $i, j \in 1..n$,

$$\boldsymbol{x}[i] \approx \boldsymbol{x}[j] \Longleftrightarrow \boldsymbol{y}[i] \approx \boldsymbol{y}[j].$$
 (1)

Lemma (1)

Every regular string is isomorphic to a string of scope 1.

Lemma (2)

Given a regular string x[1..n], then, corresponding to every triple (i, j_1, j_2) , we can assign a regular letter to $y[i], y[j_1], y[j_2]$ in such a way that the resulting string y[1..n] is isomorphic to x[1..n].

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 6 / 26



We propose the algorithm (function regular) outlined below to determine whether a given string x[1..n] on alphabet Σ is regular.

If x is regular, on exit the string y is the lex-least regular string of scope 1 on the integer alphabet $\Sigma' = \{1, 2, \dots, \sigma'\}$ that is isomorphic to x.

The runtime complexity of regular is $\mathcal{O}(n^2\sigma^2)$

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 7 / 26

Regular & Indeterminate Strings

Maximal Palindrome Array

Open Problems

References

Function regular

Input: String $\boldsymbol{x}[1..n]$

Output: If x is regular, returns true; otherwise, false.

(And if \boldsymbol{x} is regular, also constructs a lex-least string $\boldsymbol{y}[1..n]$.)

Outline of function regular

- Initialize each letter in y[1..n] to 0.
- Scan x from left to right, using y to record previous matches.
- During this scan the following condition holds as long as x is regular:

$$C: \boldsymbol{x}[i] \approx \boldsymbol{x}[j] \Leftrightarrow \boldsymbol{y}[i] = \boldsymbol{y}[j] \wedge \boldsymbol{y}[i] \neq 0.$$

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 8 / 26

Outline Introduction Regular & Indeterminate Strings Maximal Palindro	ome Array Open Problems References
-----------------------------------------------------------------------	------------------------------------

- If at a position $i \in 1..n$, we have y[i] = 0 that is, it was not part of a previous match we fill it with a new character σ' .
- We then scan the rest of the strings x[i+1..n] and y[i+1..n]to see if condition C continues to hold. If it does not, we mark x as indeterminate and exit; otherwise, whenever $x[j] \approx x[i]$ and y[j] = 0, we assign $y[j] \leftarrow \sigma'$.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Regular & Indeterminate Strings - 9 / 26

Outline Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References
----------------------	---------------------------------	--------------------------	---------------	------------

Palindromes

A substring u = x[i..j], $1 \le i \le j \le n$, of length $\ell = j-i+1$ is said to be a palindrome if $x[i+h] \approx x[j-h]$ for every $h \in 0..\lfloor \ell/2 \rfloor$.

A palindrome u = x[i..j] is said to be a maximal palindrome if one of the following holds: i = 1, j = n, or $x[i-1] \not\approx x[j+1]$.

The **centre** of a palindrome \boldsymbol{u} is at position $i + \frac{\ell-1}{2}$.

Since this is not an integer for odd ℓ , we form the string x^* , where $\# \notin \Sigma$ and m = 2n+1.

 $\boldsymbol{x}^*[1..m] = \#x_1 \# x_2 \# \cdots \# x_n \#,$

Now every palindrome in x^* has an integer centre c. We call $d = 2\ell + 1$ the **diameter** and $r = \lfloor d/2 \rfloor$ the **radius** of a palindrome in x^* .

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 10 / 26

Regular & Indeterminate Strings

Maximal Palindrome Array

Open Problems

References

Maximal Palindrome Array

We can now define the maximal palindrome array $MP = MP_{x^*}$ of x^* :

For every $i \in 1..m$, if $x^*[i] = \#$ and $x^*[i-1] \not\approx x^*[i+1]$, then MP[i] = 0 (radius zero); otherwise, $MP[i] \ge 1$ is the radius of the maximal palindrome centred at position i.

For example, MP_{x^*} derived from x = aabac is as follows:

$$\mathbf{x}^{*} = \# \ a \ \# \ a \ \# \ b \ \# \ a \ \# \ c \ \#$$
(2)
$$\mathsf{MP}_{\mathbf{x}^{*}} = 0 \ 1 \ 2 \ 1 \ 0 \ 3 \ 0 \ 1 \ 0 \ 1 \ 0$$

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 11 / 26

Outline Introduction Regular & Indeterminate Strings Maximal Palindrome Array Open Problems Re	eferences
------------------------------------------------------------------------------------------------	-----------

The most general form of the palindrome array is given by

$$\mathsf{MP} = 0i_2i_3\cdots i_{m-1}0,\tag{3}$$

where for every $j \in 2..m - 1$:

(a)
$$i_j \in (1-j \mod 2) \dots \min(j-1, m-j);$$

(b) i_j is odd if and only if j is even.

Any array satisfying (3) is said to be **feasible**.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 12 / 26



Lemma (3)

There exists a string corresponding to every feasible palindrome array.

To prove the above lemma, we construct a string x^* corresponding to the input feasible array $MP = MP_{x^*}[1..m]$ as follows:

- First, for every odd $c \in 1..m$, we assign $x^*[c] \leftarrow \#$, while every even position remains empty.
- For every c ∈ 3..m−2 such that MP[i] = r ≥ 2, we add a unique character to each pair of positions c−k, c+k in x*, where
 - (c even, r odd) $k = 2, 4, \ldots, r-1;$
 - (c odd, r even) k = 1, 3, ..., r-1.
- Finally, we assign a unique character to each position that remains empty.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 13 / 26

Outline Introduction Regular & Indeterminate Strings Maximal Palindrome Array Open Problems Refe	erences
--------------------------------------------------------------------------------------------------	---------

Example of construction for the previous Lemma:

However, as we have seen in (2), the regular string #a#a#b#a#c# has the same palindrome array: a palindrome array can correspond to both a regular and an indeterminate string!

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 14 / 26



Forbidden pair: To each position $c \in 1..m$ of a feasible MP array we associate a pair of integers (i, j) such that

$$i = c - MP[c] - 1$$
, and $j = c + MP[c] + 1$.

Then, provided 0 < i, j < m + 1, we must have $x^*[i] \not\approx x^*[j]$. We call (i, j) the **forbidden pair** with respect to c.

Assuming that $x^*[0] = x^*[m+1] = \emptyset$, we denote by FP the set of all forbidden pairs with respect to each centre $c \in 1..m$.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 15 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References

Example of forbidden pair for each centre c:

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 16 / 26



To characterize MP arrays we use Manacher's condition [Man75], restated in [IIBT10], and restated again below:

In MP = MP_{x^*}, we consider each centre c of a palindrome of radius r = MP[c], where for c even, $r \ge 3$, and for c odd, $r \ge 2$.

Then we must have $x^*[c-k] \approx x^*[c+k]$, where $1 \le k \le r$. For each k, let $r_{\ell} = MP[c-k]$, $r_r = MP[c+k]$. We then have

Definition (Manacher's condition)

Every position c in $MP = MP_{x^*}$ satisfies the following:

(a) if
$$r_{\ell} \neq rk$$
 then $r_r = \min(r_{\ell}, rk)$ else $r_r \geq r_{\ell}$;

(b) if $r_r \neq rk$ then $r_\ell = \min(r_r, rk)$ else $r_\ell \geq r_r$.

LSD & LAW 2019, London

Maximal Palindrome Array - 17 / 26



In [Man75] it is shown that, for every palindrome in string x^* such that $S(x^*) = 1$, Manacher's condition must hold.

Thus, by Lemma (1), Manacher's condition holds for every regular string; that is, for every string whose triples are all transitive.

We propose procedure *construct* which on an input MP produces a lex-least regular string if MP is regular; otherwise produces an indeterminate string.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 18 / 26

Outline Introduction Regular & Indeterminate Strings Maximal Palindrome Array Open Problem	s References
--------------------------------------------------------------------------------------------	--------------

Procedure construct

Input: Feasible maximal palindrome array MP[1..m].

Main Data Structure: FS[1..m] — an array giving forbidden pairs (i, j) that have been identified w.r.t centres c: for every (even) i and $\forall j \in FS[i]$, $\boldsymbol{x}^*[i] \not\approx \boldsymbol{x}^*[j]$.

Outline of procedure *construct*

- For every odd $i \in 1..m$, assign $\boldsymbol{x}^*[i] \leftarrow \#$; for every even $j \in 1..m$, assign $\boldsymbol{x}^*[j] \leftarrow 0$. Set $\boldsymbol{x}^*[2] \leftarrow 1$.
- For each centre $c \in 3..m 1$, compute its forbidden pair (i, j) and update the FS[i] and FS[j] sets accordingly.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 19 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References
---------	--------------	---------------------------------	--------------------------	---------------	------------

- If the string x*[1..c-1] constructed so far is regular and the centre c and range k satisfy Manacher's condition, we continue to construct the regular string by copying the previously filled letter x*[c-k] to its corresponding matching position c+k.
- Whenever there is a choice of filling an empty position that is, when x*[c] = 0 — the lex-least character which is not in the forbidden set of characters FS[c] is chosen.
- If a given centre c and range k do not satisfy Manacher's condition, we mark the string as indeterminate, and every subsequent letter match including the current one is achieved by adding a new character to x*[c-k] and x*[c+k], where c-k and c+k are even.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 20 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References

Example 1 (MP[12] = 3):

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
MP = 0	1	0	3	0	1	0	$\overline{7}$	0	1	0	3	0	1	0	(6)
$\boldsymbol{x^*}=\#$	1	#	2	#	1	#	3	#	1	#	2	#	1	#	

 x^* is the string produced by construct.

The input MP array is regular, therefore the resulting string x^* is regular.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 21 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	References

Example 2 (MP[12] = 1):

1	2	3	4	5	6	$\overline{7}$	8	9	10	11	12	13	14	15	
MP = 0	1	0	3	0	1	0	7	0	1	0	1	0	1	0	(7)
$\boldsymbol{x^*}=\#$	1	#	$\{2, 3\}$	#	$\{1, 4\}$	#	5	#	4	#	3	#	1	#	

 x^* is the string produced by construct.

The input MP array is not regular, therefore the resulting string x^* is indeterminate.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 22 / 26



Results

Lemma (4)

Let x^* be the string produced by procedure construct. Then $S(x^*) = 1 \Leftrightarrow x^*$ is regular.

Theorem (1)

Let x^* be the string produced by the procedure construct on an input MP. Then x^* is regular \Leftrightarrow MP is regular.

Theorem (2)

Given an MP array of length m = 2n+1, procedure construct executes in $\mathcal{O}(n^2\sigma)$ time, where σ is the size of the generated alphabet.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Maximal Palindrome Array - 23 / 26



Open Problems:

- **1** Procedures regular and construct both have worst-case time complexity $\mathcal{O}(n^2)$. It should be possible to reduce this.
- Q Given the new definition of regularity in strings, what is its effect on the computation of data structures such as border array, cover array, prefix array, suffix array, Lyndon array?
- 3 Is it possible to compute the above mentioned data structures without first using the $\mathcal{O}(n^2)$ time function regular?
- What effect does this new definition have on the "reverse engineering" problem related to these arrays?
- There exist algorithms to reverse engineer the other arrays noted above — can they be modified/extended to yield equivalent results?

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

Open Problems - 24 / 26

Outline	Introduction	Regular & Indeterminate Strings	Maximal Palindrome Array	Open Problems	Reference
---------	--------------	---------------------------------	--------------------------	---------------	-----------

References I

[Abr87] K. Abrahamson. Generalized string matching. *SIAM Journal of Computing*, 16(6):1039–1051, 1987.
[FP74] M.J. Fischer and M.S. Paterson. String matching and other products. In R.M. Karp, editor, *Complexity of Computation*,, pages 113–125. American Mathematical Society, 1974.
[IIBT10] Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Counting and verifying maximal palindromes. *Proc. 17th Symposium on String Proocessing & Information Retrieval* (SPIRE), LNCS 6393, pages 135–146, 2010.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

References - 25 / 26



References II

[Man75] G. Manacher.

Anew Linear– Time "On– Line" Algorithm for Finding the Smallest Initial Palindrome of a String.

J. Assoc. Comput. Mach., 22:346-351, 1975.

Louza, Mhaskar and Smyth

LSD & LAW 2019, London

References - 26 / 26