



ELSEVIER

Contents lists available at ScienceDirect

Theoretical Computer Science

www.elsevier.com/locate/tcs



Constructing an indeterminate string from its associated graph [☆]

Joel Helling ^a, P.J. Ryan ^b, W.F. Smyth ^{b,c,1}, Michael Soltys ^{a,*}

^a Department of Computer Science, California State University Channel Islands, United States

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

^c School of Engineering and Information Technology, Murdoch University, Australia

ARTICLE INFO

Article history:

Received 21 May 2016

Received in revised form 2 February 2017

Accepted 14 February 2017

Available online xxxx

Keywords:

String algorithms

Indeterminate strings

Cliques

Graph labeling

ABSTRACT

As discussed at length in Christodoulakis et al. (2015) [3], there is a natural one-many correspondence between simple undirected graphs \mathcal{G} with vertex set $V = \{1, 2, \dots, n\}$ and **indeterminate strings** $\mathbf{x} = \mathbf{x}[1..n]$ – that is, sequences of subsets of some alphabet Σ . In this paper, given \mathcal{G} , we consider the “reverse engineering” problem of computing a corresponding \mathbf{x} on an alphabet Σ_{\min} of minimum cardinality. This turns out to be equivalent to the NP-hard problem of computing the **intersection number** of \mathcal{G} , thus in turn equivalent to the **clique cover** problem. We describe a heuristic algorithm that computes an approximation to Σ_{\min} and a corresponding \mathbf{x} . We give various properties of our algorithm, including some experimental evidence that on average it requires $\mathcal{O}(n^2 \log n)$ time. We compare it with other heuristics, and state some conjectures and open problems.

© 2017 Elsevier B.V. All rights reserved.

1. Introduction

In this paper we seek to extend the connections between graph theory and stringology explored in [3]. We consider a **string** $\mathbf{x} = \mathbf{x}[1..n]$ to be a sequence of **letters** $\mathbf{x}[i]$, $1 \leq i \leq n$, that are nonempty subsets of a given finite set Σ , called the **alphabet**. If $\mathbf{x}[i]$ is a subset of cardinality 1, it is said to be a **regular** letter; otherwise, **indeterminate**. Similarly, if \mathbf{x} contains only regular letters, it is said to be **regular**; otherwise, **indeterminate**. For example, on $\Sigma = \{a, b, c\}$, $\mathbf{x} = ababc$ is regular,² while $\mathbf{y} = \{a, b\}ba\{b, c\}b$ is indeterminate. Indeterminate strings are useful in various application areas, notably bioinformatics, where under certain circumstances DNA sequences can be regarded as indeterminate strings on nucleotides $\{a, c, g, t\}$. In recent years indeterminate strings have been the subject of much study [12,11,17,1].

Given string $\mathbf{x} = \mathbf{x}[1..n]$, we say that for $1 \leq i, j \leq n$, $\mathbf{x}[i]$ **matches** $\mathbf{x}[j]$ (written $\mathbf{x}[i] \approx \mathbf{x}[j]$) if and only if $\mathbf{x}[i] \cap \mathbf{x}[j] \neq \emptyset$. Thus in particular $\mathbf{x}[i] = \mathbf{x}[j] \implies \mathbf{x}[i] \approx \mathbf{x}[j]$. As defined in [3], the **associated graph** $\mathcal{G}_{\mathbf{x}} = (V_{\mathbf{x}}, E_{\mathbf{x}})$ of \mathbf{x} is the simple graph whose vertices are positions $1, 2, \dots, n$ in \mathbf{x} and whose edges are the pairs (i, j) such that $\mathbf{x}[i] \approx \mathbf{x}[j]$. Suppose that for some position $i_0 \in 1..n$, $\mathbf{x}[i_0]$ matches $\mathbf{x}[i_1], \mathbf{x}[i_2], \dots, \mathbf{x}[i_k]$ for some $k \geq 0$, and matches no other elements of \mathbf{x} . We say that position i_0 is **essentially regular** if and only if the entries in positions i_1, i_2, \dots, i_k match each other pairwise. If every

[☆] We are grateful to Manolis Christodoulakis and Shu Wang for helpful discussions.

* Corresponding author.

E-mail addresses: joel.helling904@csuci.edu (J. Helling), ryanpj@mcmaster.ca (P.J. Ryan), smyth@mcmaster.ca (W.F. Smyth), michael.soltys@csuci.edu (M. Soltys).

¹ Supported in part by a Grant No. 14284 from the Natural Sciences and Engineering Research Council of Canada.

² For singleton sets such as $\{a\}$, $\{b\}$, $\{c\}$, we write a, b, c for simplicity.

position in \mathbf{x} is essentially regular, we say that \mathbf{x} itself is **essentially regular**. Hence every essentially regular string can be replaced by an equivalent regular one, and the associated graph $\mathcal{G}_{\mathbf{x}}$ is a collection of disjoint cliques if and only if \mathbf{x} is essentially regular.

For general indeterminate strings, however, $\mathcal{G}_{\mathbf{x}}$ is more interesting. In Section 2 we discuss a conjecture stated in [3], that given a finite simple graph \mathcal{G} whose maximal cliques have basis \mathcal{B} , $|\mathcal{B}|$ is the minimum alphabet size of any string \mathbf{x} whose associated graph $\mathcal{G}_{\mathbf{x}} = \mathcal{G}$. We discover that this conjecture is just a reformulation, in a slightly different context, of the problem of computing the intersection number of \mathcal{G} , which is NP-hard, and hence computing the minimum alphabet size of any string is also NP-hard. Section 3 describes an algorithm that approximates a basis of \mathcal{G} by assigning symbols to the vertices of cliques until all vertices are labeled, thus effectively computing a string \mathbf{x} whose associated graph $\mathcal{G}_{\mathbf{x}} = \mathcal{G}$. This is an example of the “reverse engineering” of a data structure, a class of problems initiated in [8,7] for the border array, and extended to other structures in, for example, [2,9,4]. In Section 4 we discuss our algorithm’s results and execution, especially in the context of other algorithms that perform closely-related computations. Section 5 discusses a few conjectures and open problems.

2. Maximal cliques in the associated graph $\mathcal{G}_{\mathbf{x}}$

Suppose a collection $\mathcal{F} = F_1, F_2, \dots, F_n$ of sets is given. Then the **intersection graph** $\mathcal{G}_{\mathcal{F}}$ of \mathcal{F} is a simple undirected graph on $|\mathcal{F}| = n$ vertices $1, 2, \dots, n$, with an edge (i, j) , $1 \leq i, j \leq n$, if and only if $F_i \cap F_j \neq \emptyset$. Conversely, it was shown in [18] that, given a simple undirected graph \mathcal{G} on vertices $1, 2, \dots, n$, a collection \mathcal{F} of n sets can be found such that \mathcal{G} is the intersection graph of \mathcal{F} . (For example, for each (i, j) in \mathcal{G} , $i < j$, place a unique symbol $\lambda_{i,j}$ in F_i and F_j .) The **intersection number** $\theta(\mathcal{G})$ of \mathcal{G} is the smallest number of distinct symbols that can be placed in the sets of \mathcal{F} such that $\mathcal{G} = \mathcal{G}_{\mathcal{F}}$. In our application, the collection \mathcal{F} becomes a string $\mathbf{x} = \mathbf{x}[1..n]$ with $\mathbf{x}[i] = F_i$ (necessarily nonempty). The associated graph and the intersection graph are thus the same, and we seek a smallest alphabet Σ_{\min} that produces it. Let $\sigma_{\min} = |\Sigma_{\min}|$, the cardinality of such an alphabet.

The standard way of efficiently representing a finite simple graph \mathcal{G} as an intersection graph is by covering the graph by cliques. Take any set of cliques covering all edges of \mathcal{G} . For each vertex v , let F_v be the set of those cliques containing the vertex v . Then the intersection graph of $\{F_v\}$ coincides with \mathcal{G} . As a result, the intersection number $\theta(\mathcal{G})$ is equal to the **edge clique cover number** $ec(\mathcal{G})$, the cardinality of a minimum size set of the cliques that covers all the edges of \mathcal{G} . Erdős et al. [6,16] proved that $\theta(\mathcal{G}) \leq \lfloor n^2/4 \rfloor$, an upper bound that is achieved when \mathcal{G} is a triangle-free graph on $\lfloor n^2/4 \rfloor$ edges [15]. An instructive example is given by the complete bipartite graphs $K_{m,m}$ (for even $n = 2m$) and $K_{m,m+1}$ (for odd $n = 2m + 1$) for which the minimal covering by cliques consists of all edges, the number of which is precisely $\lfloor n^2/4 \rfloor$.

Erdős et al. use coverings that cover all vertices as well as all edges. If \mathcal{G} has no isolated points, this is equivalent to the “edge covering” approach discussed above. For this case, they prove that $ec(\mathcal{G}) \leq \lfloor n^2/4 \rfloor$ and that one need only use 2-cliques and 3-cliques (edges and triangles) in a minimal covering.³

More recently, Conjecture 21 in [3] formulated the problem in a slightly different way, in terms of the **maximal** cliques in \mathcal{G} ; that is, those that are not proper subgraphs of any other clique. We provide here a proof of this conjecture, thus validating the several following remarks made in that paper. To be consistent with [3], we use the notion of **basis**. A basis is a minimum size set of maximal cliques that covers all edges and all vertices of \mathcal{G} .

Lemma 1. *Suppose that a finite simple graph \mathcal{G} with vertex set $V = \{1, 2, \dots, n\}$ has a basis \mathcal{B} of maximal cliques of cardinality σ_{\min} . Then there is a string \mathbf{x} on a base alphabet of size σ_{\min} whose associated graph $\mathcal{G}_{\mathbf{x}} = \mathcal{G}$. No string on a smaller alphabet has this property.*

Proof. Let $\mathcal{B} = \{C_1, C_2, \dots, C_{\sigma}\}$. Let $\{\lambda_s\}_{s=1}^{\sigma}$ be distinct letters. We construct a string \mathbf{x} as follows. For each ordered pair (s, i) with $1 \leq s \leq \sigma$ and $1 \leq i \leq n$, assign λ_s to $\mathbf{x}[i]$ if vertex i occurs in the maximal clique C_s . It is clear from the definitions that the string \mathbf{x} so constructed satisfies $\mathcal{G}_{\mathbf{x}} = \mathcal{G}$.

Now consider any string \mathbf{x} of length n for which $\mathcal{G}_{\mathbf{x}} = \mathcal{G}$ and let τ be the number of distinct (ordinary) letters occurring in \mathbf{x} . For each such letter λ , there is a clique C_{λ} of \mathcal{G} whose vertices are those i for which $\lambda \in \mathbf{x}[i]$. Of course, these cliques may not be maximal, but each C_{λ} can be extended to a maximal clique C'_{λ} . Note that every vertex and edge of \mathcal{G} occurs in one of the cliques C_{λ} and *a fortiori* in one of the maximal cliques C'_{λ} . However, the C'_{λ} might not all be distinct. Let τ' be the number of distinct C'_{λ} . Then $\tau \geq \tau' \geq \sigma$, the latter inequality following from the fact that there is a basis of cardinality σ_{\min} . This shows that τ cannot be less than σ_{\min} and completes the proof. \square

It turns out that maximality is irrelevant in the specification of basis. Let $\phi'(\mathcal{G})$ be the cardinality of a basis (of maximal cliques) in \mathcal{G} and let $\phi(\mathcal{G})$ be the cardinality of a smallest set of cliques that cover all edges and vertices of \mathcal{G} . Then:

Observation 2. $\phi'(\mathcal{G}) = \phi(\mathcal{G})$.

³ The authors thank an anonymous referee who drew their attention to the available material on intersection graphs and clique edge covers.

متن کامل مقاله

دریافت فوری ←

ISIArticles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات