Applied Mathematics Letters

Applied Mathematics Letters 25 (2012) 1226-1229



Contents lists available at SciVerse ScienceDirect

Applied Mathematics Letters

journal homepage: www.elsevier.com/locate/aml

A simple characterization of the minimal obstruction sets for three-state perfect phylogenies

Brad Shutters*, David Fernández-Baca

Department of Computer Science, Iowa State University, Ames, IA 50011, USA

ARTICLE INFO

Article history: Received 18 July 2011 Received in revised form 23 February 2012 Accepted 24 February 2012

Keywords: Computational biology Phylogenetics Perfect phylogeny

ABSTRACT

We give a characterization of the minimal obstruction sets for the existence of a perfect phylogeny for a set of three-state characters that can be inferred by testing each pair of characters. This leads to a $O(m^2n + p)$ time algorithm for outputting all p minimal obstruction sets for a set of m three-state characters over a set of n taxa.

© 2012 Elsevier Ltd. All rights reserved.

1. Introduction

The *k*-state perfect phylogeny problem is a classic problem in computational biology [1]. The input is an *n* by *m* matrix \mathcal{M} of integers from the set $K = \{1, \ldots, k\}$. We refer to each row of \mathcal{M} as a *taxon* (plural *taxa*), each column of \mathcal{M} as a *character*, and each value that occurs in a column *c* of \mathcal{M} as a *state* of character *c*. A perfect phylogeny for \mathcal{M} is a tree *T* with *n* leaves such that each leaf is labeled by a distinct taxon of \mathcal{M} , each internal node is labeled by a vector in K^m , and, for every character *c* and every state *i* of *c*, the nodes labeled with state *i* for character *c* form a connected subtree of *T*. The problem is to decide whether there exists a perfect phylogeny for \mathcal{M} . If so, then the characters of \mathcal{M} are *compatible*. Otherwise, they are *incompatible*. The general *k*-state perfect phylogeny problem is NP-complete [2,3]. However, for fixed *k*, the problem is solvable in $O(m^2n)$ time [4,5], and in O(mn) time when k = 2 [6].

In this note, we focus on the three-state perfect phylogeny problem, and thus, fix \mathcal{M} to be an n by m matrix of integers from the set {1, 2, 3}. We remark that several specialized algorithms have been developed for the specific case where k = 3that can construct a perfect phylogeny for \mathcal{M} (when one exists) in $O(m^2n)$ time [7–11]. However, our main concern here is the case where the characters of \mathcal{M} are incompatible. Since every subset of a compatible set of characters is compatible, it follows that if the characters of \mathcal{M} are incompatible, there must be some minimal subset of the characters of \mathcal{M} that are incompatible. A minimal obstruction set for \mathcal{M} is a minimal subset of the characters of \mathcal{M} that are incompatible. A recent breakthrough [12] showed that every minimal obstruction set for \mathcal{M} has cardinality at most 3; implying an $O(m^3n)$ time algorithm for outputting all minimal obstruction sets for \mathcal{M} . The main result of this work is a characterization of the minimal obstruction sets for \mathcal{M} that can be inferred by testing each pair of the characters of \mathcal{M} . We show that this leads to a $O(m^2n+p)$ time algorithm for outputting all p minimal obstruction sets for \mathcal{M} . Although there can be $O(m^3)$ minimal obstruction sets for \mathcal{M} , in practice we expect the number of minimal obstruction sets to be small.

We conclude with a theorem relating our characterization of the minimal obstruction sets for three-state perfect phylogenies to monochromatic pairs of vertices of the partition intersection graph of \mathcal{M} with no legal minimal separator. By the results of [10–12] the existence of such a pair of vertices certifies that no perfect phylogeny for \mathcal{M} exists.

* Corresponding author. Tel.: +1 515 257 6855.

E-mail addresses: shutters@iastate.edu (B. Shutters), fernande@iastate.edu (D. Fernández-Baca).

^{0893-9659/\$ –} see front matter 0 2012 Elsevier Ltd. All rights reserved. doi:10.1016/j.aml.2012.02.060



Fig. 1. The "forbidden" sets of edges in the intersection graph of three three-state characters that admit a perfect phylogeny. In [12], four forbidden subgraphs are given. However, one of the subgraphs given is a subgraph of another. For our purposes, the supergraph is not needed and so it is not shown here. The colored edges are used in the proof of Theorem 6.

2. Preliminaries

We fix \mathcal{M} to be an n by m matrix of integers from the set {1, 2, 3}. For a subset C of the characters of \mathcal{M} , the matrix $\mathcal{M}|C$ is obtained by restricting \mathcal{M} to the characters in C. $\mathcal{G}(\mathcal{M})$ is the *intersection graph* of \mathcal{M} which has a vertex c_i for each character c of \mathcal{M} and each state i of c, and an edge $c_i d_j$ precisely when there is a taxon of \mathcal{M} having state i for character c and state j for character d. Note that $\mathcal{G}(\mathcal{M})$ cannot have an edge between vertices associated with different states of the same character of \mathcal{M} .

In [7], a matrix $\overline{\mathcal{M}}$ of two-state characters is derived from \mathcal{M} by adding, for each character *c* of \mathcal{M} , two-state characters c(1), c(2), c(3) to $\overline{\mathcal{M}}$. All taxa having state *i* for *c* are given state 1 for c(i), and all other taxa are given state 2 for c(i). Since every character of $\overline{\mathcal{M}}$ has two states, it follows from the splits equivalence theorem (also known as the four-gamete condition) that two characters c(i) and d(j) of $\overline{\mathcal{M}}$ are *incompatible* if and only if the two columns corresponding to c(i) and d(j) contain all four of the pairs (1, 1), (1, 2), (2, 1), and (2, 2); otherwise c(i) and d(j) are *compatible* [13]. Note that this implies that c(i) and d(j) are compatible if and only if $\mathcal{G}(\overline{\mathcal{M}}|\{c(i), d(j)\})$ is not a cycle. A set of two-state characters is compatible if and only if each pair of characters in the set is compatible [13]. Theorem 1 shows that we can test for the existence of a perfect phylogeny for \mathcal{M} by finding a compatible subset of the characters in $\overline{\mathcal{M}}$.

Theorem 1 (See [7]). There is a perfect phylogeny for \mathcal{M} if and only if there is a subset C of the characters in $\overline{\mathcal{M}}$ such that both of the following hold: (i) every pair of characters in C is compatible; and (ii) for every character c of \mathcal{M} , C contains at least two characters from the set {c(1), c(2), c(3)}.

Theorem 2 generalizes the splits equivalence theorem to three-state characters.

Theorem 2 (See [12]). There is a perfect phylogeny for \mathcal{M} if and only if both of the following hold: (i) for every pair $\{a, b\}$ of characters of \mathcal{M} , $\mathcal{G}(\mathcal{M}|\{a, b\})$ is acyclic; and (ii) for every triple $\{a, b, c\}$ of characters of \mathcal{M} , $\mathcal{G}(\mathcal{M}|\{a, b, c\})$ does not contain, up to relabeling of characters and states, any of the subgraphs shown in Fig 1.

3. The main results

For a character c of \mathcal{M} and a state i of c, if there is a character d of \mathcal{M} and two states j and k of d such that c(i) is incompatible with both d(j) and d(k), then we say that i is a *dependent* state of c with d as a *witness*. We give a complete characterization of the obstruction sets for \mathcal{M} in terms of the dependent states of its characters.

Lemma 3. Let *c* be a character of \mathcal{M} and let *i* be a dependent state of *c*. No subset of the characters in $\overline{\mathcal{M}}$ satisfying both conditions (i) and (ii) of Theorem 1 contains *c*(*i*).

Proof. Let *C* be a subset of the characters of $\overline{\mathcal{M}}$ containing c(i). Since *i* is a dependent state, there is a character *d* of \mathcal{M} and two states *j*, *k* of *d* such that c(i) is incompatible with both d(j) and d(k). It follows that either *C* contains at most one of $\{d(1), d(2), d(3)\}$, or *C* contains a pair of incompatible characters. Thus, *C* cannot satisfy both conditions (i) and (ii) of Theorem 1. \Box

Theorem 4. Let a, b, and c be three characters of \mathcal{M} where c has two dependent states i and j such that a is a witness that state i is dependent, and b is a witness that state j is dependent. The set $\{a, b, c\}$ is an obstruction set for \mathcal{M} . Furthermore, if the characters in $\{a, b, c\}$ are pairwise compatible, then $\{a, b, c\}$ is a minimal obstruction set for \mathcal{M} .

Proof. Let $M = \mathcal{M}|\{a, b, c\}$ and suppose that *C* is a subset of the characters in \overline{M} satisfying both conditions (i) and (ii) of Theorem 1. By Lemma 3, $c(i) \notin C$ and $c(j) \notin C$. Then *C* contains at most one of $\{c(1), c(2), c(3)\}$. This contradicts that *C* satisfies condition (ii) of Theorem 1. Hence, no such *C* exists. So by Theorem 1, there is no perfect phylogeny for *M* and $\{a, b, c\}$ is an obstruction set for \mathcal{M} . \Box

We now give a characterization of when a state is dependent using the intersection graph of two characters of \mathcal{M} , and then show that every minimal obstruction set for \mathcal{M} contains a character with two dependent states. For a path $p : p_1p_2p_3p_4p_5$ of length 4, we write mid(p) to denote p_3 , the "middle" vertex of p. We consider a 4-cycle to be a path of length 4 and mid(p) is allowed to be any vertex on the cycle. Note that at least four vertices lie on any cycle in $\mathcal{G}(\mathcal{M})$, and for a pair {c, d} of characters of \mathcal{M} , every state associated with a vertex on a cycle in $\mathcal{G}(\mathcal{M}|\{c, d\})$ is dependent. **Lemma 5.** A state *i* of a character *c* of \mathcal{M} is a dependent state *i* f and only *i* f there is a character *d* of \mathcal{M} and a path *p* of length 4 in $\mathcal{G}(\mathcal{M}|\{c, d\})$ with $c_i = \text{mid}(p)$. Furthermore, if such a *d* exists, then *d* is a witness that *i* is a dependent state of *c*.

Proof. (\Rightarrow) Suppose that *i* is a dependent state of a character *c* of \mathcal{M} with *d* as witness. W.l.o.g. relabel the states of *c* and *d* so that *i* = 1 and that *c*(1) is incompatible with both *d*(1) and *d*(2). Let $G = \mathcal{G}(\mathcal{M}|\{c, d\})$. Then, c_1d_1 and c_1d_2 are edges of *G*. Since, *c*(1) is incompatible with *d*(1), either c_2d_1 or c_3d_1 is an edge of *G*. Since *c*(1) is incompatible with *d*(2), either c_2d_1 or c_3d_1 is an edge of *G*. Since *c*(1) is incompatible with *d*(2), either c_2d_2 or c_3d_2 is an edge of *G*. We show that in every case *G* contains a path *p* such that $mid(p) = c_1$. If c_2d_1 and c_2d_2 are edges of *G*, then $c_1d_1c_2d_2c_1$ is a cycle containing c_1 . If c_2d_1 and c_3d_2 are edges of *G*, then $c_3d_2c_1d_1c_2$ is the required path of length 4. If c_3d_1 and c_2d_2 are edges of *G*, then $c_1d_1c_3d_2c_1$ is a cycle containing c_1 . (\Leftarrow) Let *d* be a character of \mathcal{M} such that there is a path *p* of length 4 in $G = \mathcal{G}(\mathcal{M}|\{c, d\})$ with $c_i = mid(p)$. W.l.o.g. relabel the states of *c* and *d* so that *i* = 1. *G* cannot contain edges between two states of the same character. So either *p* is a cycle containing c_1 , or the path $c_2d_1c_1d_2c_3$ (up to possibly renaming states d_1 and d_2). In both cases, the four-gamete condition shows that c(1) is incompatible with two of d(1), d(2), d(3). Hence, state 1 is a dependent state of *c* with *d* as witness. \Box

Theorem 6. Let C be a minimal obstruction set for *M*. Then C contains a character with two dependent states.

Proof. By Theorem 2, the cardinality of *C* is either 2 or 3. *Case* 1. If the cardinality of *C* is 2, then it follows from Theorem 2 that $\mathcal{G}(\mathcal{M}|C)$ contains a cycle. Since there cannot be an edge in $\mathcal{G}(\mathcal{M}|C)$ between vertices associated with states of the same character, it follows that any cycle $\mathcal{G}(\mathcal{M}|C)$ has at least four vertices, and every state associated with a vertex on this cycle is a dependent state. The theorem follows. *Case* 2. If the cardinality of *C* is 3, then it follows from Theorem 2 that $\mathcal{G}(\mathcal{M}|C)$ contains one of the graphs of Fig 1 as a subgraph (after possibly renaming the characters and states of $\mathcal{M}|C$). If Fig 1(a) is a subgraph of $\mathcal{G}(\mathcal{M}|C)$, then $c_3b_1c_1b_2c_2$ (colored red) is a path showing that state 1 of *c* is dependent, and $c_3a_1c_2a_3c_1$ (colored blue) is a path showing that state 2 of *c* is dependent. If Fig 1(b) is a subgraph of $\mathcal{G}(\mathcal{M}|C)$, then $c_3b_1c_1b_2c_2$ (colored red) is a path showing that state 2 of *c* is dependent. If Fig 1(b) is a subgraph of $\mathcal{G}(\mathcal{M}|C)$, then $c_3b_1c_1b_2c_2$ (colored red) is a path showing that state 2 of *c* is dependent. If Fig 1(c) is a subgraph of $\mathcal{G}(\mathcal{M}|C)$, then $c_3a_2c_1a_1c_2$ (colored red) is a path showing that state 2 of *c* is dependent. If Fig 1(c) is a path showing that state 2 of *c* is dependent. If every case, we have shown that states 1 and 2 are dependent states of *c*. Thus, \mathcal{M} contains a character with two dependent states. \Box

Theorems 4 and 6 together with Theorem 2 give us the following test for the existence of a perfect phylogeny for \mathcal{M} .

Theorem 7. There is a perfect phylogeny for \mathcal{M} if and only if there is at most one dependent state of each character c of \mathcal{M} .

Proof. If there is a perfect phylogeny for \mathcal{M} , then there is no obstruction set for \mathcal{M} . Thus, by Theorem 4, there can be no character of \mathcal{M} with more than one dependent state. If there is no perfect phylogeny for \mathcal{M} , then there must exist some minimal obstruction set for \mathcal{M} . By Theorem 6, there is a character of \mathcal{M} with two or more dependent states. \Box

An immediate consequence of Theorem 7 is that every set *C* of three-state characters has a canonical subset that *does* have a perfect phylogeny, namely the subset { $c \in C : c$ has at most one dependent state}.

We now describe an algorithm, denoted by A, which outputs all of the minimal obstruction sets for M. Step 1 of A computes for each character *c* of M the following.

- A set $\mathcal{B}(c)$ of all characters *d* of \mathcal{M} such that $\mathcal{G}(\mathcal{M}|\{c, d\})$ contains a cycle.
- For each state *i* of *c*, a set $\mathcal{D}(c, i)$ of all characters *d* of \mathcal{M} such that $d \notin \mathcal{B}(c)$ and there is a path *p* of length 4 in $\mathcal{G}(\mathcal{M}|\{c, d\})$ with $c_i = \text{mid}(p)$.

Step 2 of A visits each character c of M and outputs the following.

- For each character d in $\mathcal{B}(c)$, the set $\{c, d\}$.
- For each pair of states $\{i, j\}$ of c with both $\mathcal{D}(c, i)$ and $\mathcal{D}(c, j)$ non-empty, each element of the set $\{\{c, x, y\} : x \in \mathcal{D}(c, i), y \in \mathcal{D}(c, j)\}$.

Theorem 8. A outputs all p minimal obstruction sets for \mathcal{M} in $O(m^2n + p)$ time.

Proof. We first establish the following claim.

Claim 1. For each character c of \mathcal{M} , the sets $\mathcal{B}(c)$, $\mathcal{D}(c, 1)$, $\mathcal{D}(c, 2)$, and $\mathcal{D}(c, 3)$ are pairwise disjoint.

Proof of Claim 1. Clearly $\mathcal{B}(c) \cap (\mathcal{D}(c, 1) \cup \mathcal{D}(c, 2) \cup \mathcal{D}(c, 3)) = \emptyset$, so it suffices to show that for each character c of \mathcal{M} , the sets $\mathcal{D}(c, 1)$, $\mathcal{D}(c, 2)$, and $\mathcal{D}(c, 3)$ are pairwise disjoint. W.l.o.g. let c be a character of \mathcal{M} and let $d \in \mathcal{D}(c, 1) \cap \mathcal{D}(c, 2)$. Let $G = \mathcal{G}(\mathcal{M} | \{c, d\})$. Since $d \in \mathcal{D}(c, 1) \cap \mathcal{D}(c, 2)$, $d \notin \mathcal{B}(c)$, and there are paths p_1 and p_2 of length 4 in G with $c_1 = \operatorname{mid}(p_1)$ and $c_2 = \operatorname{mid}(p_2)$. Since $d \notin \mathcal{B}(c)$, G is acyclic. W.l.o.g. suppose that p_1 is the path $c_2d_1c_1d_2c_3$. Since $\operatorname{mid}(p_2) = c_2$, there must be two edges from c_2 to vertices associated with states of d. We have that c_2d_1 is an edge of p_1 . If c_2d_2 is an edge of p_2 , then we have a cycle in G. So c_2d_3 and either d_3c_3 or d_3c_1 are edges of p_2 . In either case, there is a cycle in G.

By Lemma 5, \mathcal{A} finds all dependent states, and hence, by Theorems 2, 4 and 6, outputs all of the minimal obstruction sets for \mathcal{M} . By Claim 1, every obstruction set output by \mathcal{A} is minimal. We now establish the runtime. Step 1 of \mathcal{A} takes $O(m^2n)$ time to construct the intersection graphs of each pair of characters of \mathcal{M} . Since each intersection graph has exactly six vertices and at most nine edges, it follows that it all cycles and paths of length 4 can be found in O(1) time. Hence step 1 takes $O(m^2n)$ time. Step 2 of \mathcal{A} visits each of the *m* characters of \mathcal{M} and takes O(1) time per set output. Any minimal obstruction set of cardinality 2 will be output twice. If follows from Claim 1 that each minimal obstruction set of cardinality 3 will be output at most three times. Thus, step 2 takes O(m + p) time where *p* is the number of minimal obstruction sets. Hence, \mathcal{A} takes $O(m^2n + p)$ time to complete both steps 1 and 2. \Box

Several approaches to determining the existence of a perfect phylogeny for \mathcal{M} studied in the literature make use of separating sets in $\mathfrak{G}(\mathcal{M})$ [10,11]. For two vertices a and b of $\mathfrak{G}(\mathcal{M})$, an a-b separator is a set of vertices whose removal separates a from b. An a-b separator is minimal if no subset of it is an a-b separator. A minimal separator is a separator that is a minimal a-b separator for some pair a, b of vertices of $\mathfrak{G}(\mathcal{M})$. A minimal separator S of $\mathfrak{G}(\mathcal{M})$ is legal if, for each character c of \mathcal{M} , S contains at most one vertex corresponding to a state of c. A pair of vertices of $\mathfrak{G}(\mathcal{M})$ representing different states of the same character is monochromatic.

Theorem 9 (See [10–12]). There is a perfect phylogeny for \mathcal{M} if and only if both of the following hold: (i) the characters of \mathcal{M} are pairwise compatible; and (ii) every monochromatic pair of vertices in $\mathcal{G}(\mathcal{M})$ is separated by a legal minimal separator.

We conclude with Theorem 10 that relates dependent states to legal minimal separators. A consequence of Theorem 10 is that algorithm \mathcal{A} can be easily modified to output monochromatic pairs of vertices of $\mathcal{G}(\mathcal{M})$ with no legal minimal separator.

Theorem 10. Suppose that the characters of \mathcal{M} are pairwise compatible. Two states *i* and *j* of a character *c* of \mathcal{M} are dependent *i*f and only if there is no legal minimal separator for c_i and c_j in $\mathcal{G}(\mathcal{M})$.

Proof. By Theorems 2 and 9, it suffices to show that the theorem holds for every minimal obstruction set, i.e., for each graph in Fig 1, a monochromatic pair of vertices has no legal minimal separator if and only if they correspond to a pair of dependent states. This is verified by inspection. In the graph of Fig 1(a): $\{c_1, c_2\}$ is the only *monochromatic* pair of vertices with no legal minimal separator; 3 is the only dependent state of *a*; 2 is the only dependent state of *b*; and 1 and 2 are the only dependent states of *c*. In the graph of Fig 1(b): (b_1, b_2) and (c_1, c_2) are the only dependent states of *b*; and 1 and 2 are the only dependent states of *c*. In the graph of Fig 1(b): (b_1, b_2) and (c_1, c_2) are the only dependent states of *b*; and 1 and 2 are the only dependent states of *c*. In the graph of Fig 1(c): $(a_1, a_2), (b_1, b_3), (c_1, c_2)$, and are the only monochromatic pairs of vertices with no legal minimal separator; 1 and 2 are the only dependent states of *a*; 1 and 3 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *b*; and 1 and 2 are the only dependent states of *c*. \Box

Acknowledgments

We thank an anonymous reviewer for several improvements to the exposition of our results. This work was supported in part by the National Science Foundation under grants CCF-1017189 and DEB-0829674.

References

- [1] D. Fernández-Baca, The perfect phylogeny problem, in: Steiner Trees in Industry, Kluwer, 2001, pp. 203–234.
- [2] H. Bodlaender, M. Fellows, T. Warnow, Two strikes against perfect phylogeny, in: ICALP 1992, in: Lect. Notes Comput. Sci., vol. 623, Springer, 1992, pp. 273–283.
- [3] M. Steel, The complexity of reconstructing trees from qualitative characters and subtrees, J. Classif. 9 (1992) 91–116.
- [4] R. Agarwala, D. Fernández-Baca, A polynomial-time algorithm for the perfect phylogeny problem when the number of character states is fixed, SIAM J. Comput. 23 (6) (1994) 1216–1224.
- [5] S. Kannan, T. Warnow, A fast algorithm for the computation and enumeration of perfect phylogenies, SIAM J. Comput. 26 (6) (1997) 1749–1763.
- [6] D. Gusfield, Efficient algorithms for inferring evolutionary trees, Networks 21 (1) (1991) 19–28.
- [7] A. Dress, M. Steel, Convex tree realizations of partitions, Appl. Math. Lett. 5 (3) (1992) 3-6.
- [8] S. Kannan, T. Warnow, Inferring evolutionary history from DNA sequences, SIAM J. Comput. 23 (4) (1994) 713-737.
- [9] D. Gusfield, Y. Wu, The three-state perfect phylogeny problem reduces to 2-SAT, Commun. Inf. Syst. 9 (4) (2009) 195–301.
- [10] D. Gusfield, The multi-state perfect phylogeny problem with missing and removable data: solutions via integer-programming and chordal graph theory, J. Comput. Biol. 17 (3) (2010) 383–399.
- [11] R. Gysel, F. Lam, D. Gusfield, Constructing perfect phylogenies and proper triangulations for three-state characters, in: WABI 2011, in: Lect. Notes Comput. Sc., vol. 6833, Springer, 2011, pp. 104–115.
- [12] F. Lam, D. Gusfield, S. Sridhar, Generalizing the splits equivalence theorem and four gamete condition: perfect phylogeny on three-state characters, SIAM J. Discrete Math. 25 (3) (2011) 1144–1175.
- [13] C. Semple, M. Steel, Phylogenetics, Oxf. Lect. S. Math. Appl., Oxford University Press, 2003.