

# On The Hardness of Inferring Phylogenies from Triplet-Dissimilarities

Ilan Gronau and Shlomo Moran

July 26, 2007

## Abstract

This work considers the problem of reconstructing a phylogenetic tree from *triplet dissimilarities*, which are dissimilarities defined over taxon-triplets. Triplet dissimilarities are possibly the simplest generalization of pairwise dissimilarities, and were used for phylogenetic reconstructions in the past few years. We study the hardness of finding a tree best fitting a given triplet-dissimilarity table under the  $\ell_\infty$  norm. We show that the corresponding decision problem is NP-hard and that the corresponding optimization problem cannot be approximated in polynomial time within a constant multiplicative factor smaller than 1.4. On the positive side, we present a polynomial time constant-rate approximation algorithm for this problem. We also address the issue of best-fit under *maximal distortion*, which corresponds to the largest *ratio* between matching entries in two triplet-dissimilarity tables. We show that it is NP-hard to approximate the corresponding optimization problem within any constant multiplicative factor.

## 1 Introduction

Phylogenetic reconstruction methods attempt to find the evolutionary history of a given set of extant species (taxa). This history is usually described by an edge-weighted tree whose internal vertices represent past speciation events (extinct species) and whose leaves correspond to the given set of taxa. The amount of evolutionary change between two subsequent speciation events is indicated by the weight of the edge connecting them. It is usually assumed (for uniqueness of representation) that internal edges<sup>1</sup> have strictly positive weights. Distance-based phylogenetic reconstruction methods typically try to reconstruct this evolutionary tree from estimates of distances (sum of weights) along edges in this tree.

Most common distance-based reconstruction algorithms receive as input a *dissimilarity matrix*  $D$ , where  $D(i, j)$  is an estimate of the distance between

---

<sup>1</sup>An edge is *external* if it is adjacent to a leaf, and is *internal* otherwise.

taxa  $i$  and  $j$ . A dissimilarity matrix is said to be *additive* if it can be realized by distances along the edges of a tree whose leaves are the elements of  $S$  [3]. There are numerous algorithms which reconstruct a tree given its additive metric, the earliest of which appeared in [3, 13, 14]. However, in reality we are unable to obtain accurate distance estimates, and the input dissimilarity matrix is rarely additive. In such a case, a natural goal is to reconstruct a tree fitting the input matrix in some way. One approach is to return a tree whose implied metric is ‘close’ to the input under a certain distance norm. Unfortunately, finding a tree closest to a given dissimilarity matrix was shown to be NP-hard under the  $\ell_1$  and  $\ell_2$  norms in [4], and under the  $\ell_\infty$  norm in [1]. [1] also presents a 3-approximation algorithm for the problem of finding the tree closest, under  $\ell_\infty$ , to an arbitrary metric; another 3-approximation algorithm for this problem was presented later in [8]<sup>2</sup>.

In this paper we study the problem of reconstructing a phylogenetic tree based on estimates of *triplet distances*. Given an edge-weighted tree  $T$  and three taxa  $i, j, k$ , we denote by  $C(i, j, k)$  the inner vertex of degree 3 in the claw spanned by  $i, j, k$  (see Fig.1), and by  $D_T(i; jk)$  the weight of the path connecting  $i$  and  $C(i, j, k)$ . Note that for all  $k \in S$ ,  $D_T(i; jk) + D_T(j; ik) = D_T(i, j)$ , and in particular  $D_T(i; jj) = D_T(i, j)$ . Hence, triplet-distances generalize the classical notion of pairwise-distances.

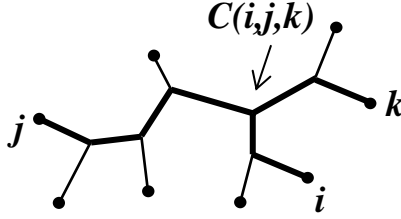


Figure 1:  $C(i, j, k)$  is the inner vertex of degree 3 in the claw spanned by  $i, j, k$ .  $D_T(i; jk)$  is the weight of the path connecting  $i$  and  $C(i, j, k)$

A *triplet-dissimilarity table* contains estimates of all triplet distances over a given taxon-set. A function  $\tau : S \times S \times S \rightarrow \mathbb{R}^+$  is a valid triplet-dissimilarity table iff it satisfies the following properties:

1.  $\tau(i, i, j) = 0$
2.  $\tau(i, j, k) = \tau(i, k, j)$
3.  $\tau(i, j, j) = \tau(j, i, i)$

For such a function we denote:  $\tau(i; jk) \triangleq \tau(i, j, k)$  and  $\tau(i, j) \triangleq \tau(i, j, j)$ .

There are several previous works which propose algorithms for reconstructing trees from triplet-dissimilarity tables. In [11], triplet-dissimilarities are used to

<sup>2</sup>The 3-approximation ratio of the algorithms in [1, 8] is proved under the assumption that the input dissimilarity matrix is a *metric*, meaning that it satisfies the triangle inequality  $[\mathcal{D}(x, y) + \mathcal{D}(y, z) \geq \mathcal{D}(x, z)]$

obtain more accurate estimates of pairwise-distances for Saitou&Nei's Neighbor Joining algorithm (commonly referred to as NJ) [12]. [10] generalizes NJ to receive as input *m-dissimilarity maps*, which contain the *total weights* of all subtrees spanned by subsets of  $m$  taxa. In [8] we present a family of algorithms (DLCA) which construct trees from estimates of triplet-distances from a single root-taxon  $r$ , meaning that the input is a symmetric matrix  $L_r$ , where  $L_r(i, j)$  is an estimate of  $D_T(r; ij)$ . We show there that a tree whose triplet distances  $\{D_T(r; ij) : i, j \in S\}$  are closest to  $L_r$  under  $\ell_\infty$  can be constructed in  $O(n^2)$  time. In this paper we show that it is NP-hard to find an edge-weighted tree  $T$  whose *entire* triplet-distance table  $\{D_T(i; jk) : i, j, k \in S\}$  is closest to a given triplet-dissimilarity table under  $\ell_\infty$ .

The  $\ell_\infty$  norm measures the maximal *difference* between corresponding entries in two triplet-dissimilarity tables:

$$\|\tau_1, \tau_2\|_\infty \triangleq \max_{i,j,k \in S} \{|\tau_1(i; jk) - \tau_2(i; jk)|\}$$

Another distance measure we refer to is *maximal distortion* [2], which is related to the maximal *ratio* between such entries:

$$\text{MaxDist}(\tau_1, \tau_2) \triangleq \max_{i,j,k \in S} \left\{ \frac{\tau_1(i; jk)}{\tau_2(i; jk)} \right\} \cdot \max_{i,j,k \in S} \left\{ \frac{\tau_2(i; jk)}{\tau_1(i; jk)} \right\} \quad (\text{where } 0/0 \triangleq 1)$$

We note that maximal distortion seems to be the most relevant criterion for the evolutionary models assumed in [5, 6] and numerous subsequent works.

Consider the decision version of the 'best-fit to triplet-dissimilarities' problem: given a triplet-dissimilarity table  $\tau$  and a non-negative number  $K$ , is there a tree  $T$  such that  $\|D_T, \tau\|_\infty \leq K$ ? In Section 2 this decision problem is shown to be NP-hard by a polynomial reduction from 3-SAT. In Section 3 we refine the analysis of the reduction to show that it is NP-hard to find a tree whose distance to the input under  $\ell_\infty$  is less than 1.4 times that of the closest tree. In Section 4 we present few other related hardness results implied by our reduction, including the NP-hardness of approximating maximal distortion for triplet dissimilarities by any multiplicative constant. In Section 5 we give an upper bound on the approximation ratio of this problem by showing that a constant-rate approximation for the closest tree to a dissimilarity matrix implies also a constant-rate approximation for the closest tree to a triplet dissimilarity table. We conclude with a short discussion of some relevant open questions.

## 2 A Reduction from 3SAT to the 'Best-Fit to Triplets Under $\ell_\infty$ ' Problem

In this section we present a reduction from 3SAT to the decision version of the 'best-fit to triplets under  $\ell_\infty$ ' problem. This reduction transforms a 3CNF formula  $\varphi$  into a valid triplet-dissimilarity table  $\tau_\varphi$  satisfying three requirements (where  $\Delta$  is a positive constant independent of  $\varphi$ ):

**POLY**  $\tau_\varphi$  can be computed in polynomial time given  $\varphi$ .

**SAT** If  $\varphi$  is satisfiable, then there is a tree  $T$  s.t.  $\|D_T, \tau_\varphi\|_\infty \leq \Delta$ .

**UNSAT** If  $\varphi$  is unsatisfiable, then for every tree  $T$ ,  $\|D_T, \tau_\varphi\|_\infty > \Delta$ ,

Similar to the reduction presented in [1] for the problem of fitting trees to dissimilarity matrices, we first transform the formula  $\varphi$  into a set of upper and lower bounds on *some* triplet distances of a tree  $T$  (see **A1-B3** and Figure 2 below). **UNSAT** is proven by showing that a tree satisfying all these bounds implies a satisfying assignment to  $\varphi$  (Lemma 2.5). These bounds are enforced by the triplet-dissimilarity table  $\tau_\varphi$  in the following way: A bound  $D_T(i; jk) \leq \omega_{ijk}$  is enforced by  $\tau_\varphi(i; jk) = \omega_{ijk} - \Delta$ , and a bound  $D_T(i; jk) \geq \omega_{ijk}$  is enforced by  $\tau_\varphi(i; jk) = \omega_{ijk} + \Delta$ . Clearly, a tree  $T$  satisfying  $\|D_T, \tau_\varphi\|_\infty \leq \Delta$  is guaranteed to obey all bounds<sup>3</sup>.

Requirement **POLY** will be obvious from the description of the transformation. To prove **SAT** we show that a satisfying assignment to  $\varphi$  implies a tree satisfying all bounds (Lemma 2.6). In this tree, triplet-distances corresponding to entries restricted by these bounds are set to satisfy the bounds *with equality*. Other triplet-dissimilarities (not restricted by any bound) are undetermined, however each such dissimilarity falls within one of two intervals:  $[r - \Delta, r + \Delta]$  or  $[s - \Delta, s + \Delta]$ . Entries of  $\tau_\varphi$  corresponding to these dissimilarities are set to the mid-point of the appropriate interval ( $r$  or  $s$ ).

Let us start with some notations: A 3CNF formula  $\varphi$  over a set of variables  $\{x_1, x_2, \dots, x_n\}$  is a conjunction of  $m$  clauses  $\varphi = c_1 \wedge c_2 \wedge \dots \wedge c_m$ , s.t.  $\forall j = 1..m : c_j = (l_1^j \vee l_2^j \vee l_3^j)$ , where  $l_1^j, l_2^j, l_3^j$  are literals (a literal is variable  $x_i$  or its negation  $\bar{x}_i$ ). For such a formula, we define a set of taxa:

$$S_\varphi = \{\mathcal{T}, \mathcal{F}\} \cup \{x_i, \bar{x}_i : i = 1..n\} \cup \{y_1^j, y_2^j, y_3^j : j = 1..m\}$$

We define the following set of bounds on triplet-dissimilarities over  $S_\varphi$  with parameters  $\alpha, \beta > 0$  (Figure 2 can be helpful at this point):

$$\mathbf{A1} \quad D_T(\mathcal{T}, \mathcal{F}) \geq 2\alpha + 2\beta$$

$$\mathbf{A2} \quad \forall i = 1..n : D_T(\mathcal{F}; x_i \bar{x}_i) \leq \alpha ; D_T(\mathcal{T}; x_i \bar{x}_i) \leq \alpha$$

$$\mathbf{B1} \quad \forall j = 1..m : D_T(y_1^j; l_2^j l_3^j) \leq \alpha ; D_T(y_2^j; l_1^j l_3^j) \leq \alpha ; D_T(y_3^j; l_1^j l_2^j) \leq \alpha$$

$$\mathbf{B2} \quad \forall j = 1..m : D_T(y_1^j; \mathcal{T}\mathcal{F}) \geq \alpha ; D_T(y_2^j; \mathcal{T}\mathcal{F}) \geq \alpha ; D_T(y_3^j; \mathcal{T}\mathcal{F}) \geq \alpha$$

$$\mathbf{B3} \quad \forall j = 1..m : D_T(\mathcal{T}; y_1^j y_2^j) \leq \alpha ; D_T(\mathcal{T}; y_1^j y_3^j) \leq \alpha ; D_T(\mathcal{T}; y_2^j y_3^j) \leq \alpha$$

Let  $T$  be a tree satisfying **A1-B3** above. Denote the mid-point of the path connecting  $\mathcal{T}$  and  $\mathcal{F}$  in this tree by  $v_\varphi$ . Note that restriction **A1** implies that  $\mathcal{T}$  and  $\mathcal{F}$  are at distance of at least  $\alpha + \beta$  from  $v_\varphi$ . Denote by  $v_\mathcal{T}$  and  $v_\mathcal{F}$  the points whose distance from  $v_\varphi$ , on the paths leading to  $\mathcal{T}$  and  $\mathcal{F}$  respectively, is *exactly*  $\beta$ . For the sake of the analysis below, we treat the three points  $\varphi, v_\mathcal{T}, v_\mathcal{F}$

<sup>3</sup>Note that in order to keep all entries of  $\tau_\varphi$  nonnegative, we need that  $\Delta \leq \omega_{ijk}$  whenever  $\omega_{ijk}$  is an upper bound on the corresponding entry.

as vertices in the tree (possibly of degree 2), and assume that  $T$  is rooted at  $v_\varphi$  (see Fig. 2).

We now describe and prove the topological restrictions implied by these bounds. Our proof is based on two simple connections between distances and topological properties of quartets (subtrees spanned by four taxa), which we bring next. For vertices  $x, y$  in  $T$ , denote by  $path(x, y)$  the path in  $T$  connecting  $x$  and  $y$ .

**Lemma 2.1.** *For all taxa  $u, v, y$  in  $T$ , we have*

1. *If  $D_T(\mathcal{F}; uv) \leq \alpha$  and  $D_T(\mathcal{T}; uv) \leq \alpha$ , then either  $u$  is a descendant of  $v_{\mathcal{T}}$  and  $v$  is a descendant of  $v_{\mathcal{F}}$  or vice versa.*
2. *If both  $u$  and  $v$  are descendants of  $v_{\mathcal{F}}$ , and  $D_T(y; uv) \leq D_T(y, \mathcal{T}\mathcal{F})$  then  $y$  is also a descendant of  $v_{\mathcal{F}}$ .*

*Proof.*

1. Since  $D_T(\mathcal{F}; uv) + D_T(\mathcal{T}; uv) \leq 2\alpha < 2\alpha + 2\beta \leq D_T(\mathcal{F}, \mathcal{T})$  (by the assumption and bound **A1**), we must have that  $C(u, \mathcal{T}, \mathcal{F})$  and  $C(v, \mathcal{T}, \mathcal{F})$  are distinct vertices on  $path(\mathcal{F}, \mathcal{T})$ . In addition, the assumption also implies that one of them is at distance at most  $\alpha$  from  $\mathcal{T}$  and the other is at distance at most  $\alpha$  from  $\mathcal{F}$ , which proves the claim.
2. Let  $z$  be the father of  $v_{\mathcal{F}}$  (possibly  $z = v_\varphi$ ). Notice that the edge  $(z, v_{\mathcal{F}})$  is in  $path(\mathcal{T}, \mathcal{F})$  (see Fig. 3). Since both  $u$  and  $v$  are descendants of  $v_{\mathcal{F}}$ , we have that if  $y$  is not a descendant of  $v_{\mathcal{F}}$ , then the path from  $y$  to  $path(u, v)$  must contain the edge  $(z, v_{\mathcal{F}})$ , and hence  $D_T(y; uv) \geq D_T(y, \mathcal{T}\mathcal{F}) + w(z, v_{\mathcal{F}}) > D_T(y, \mathcal{T}\mathcal{F})$ , a contradiction.  $\square$

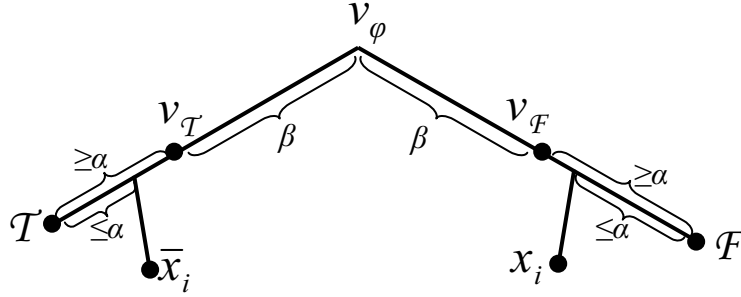


Figure 2: The Topology of a tree satisfying **A1-2**

As a direct consequence of **A2** and Lemma 2.1(1) above, we have the following:

**Corollary 2.2.** *For each  $i = 1..n$ , one of the vertices  $x_i, \bar{x}_i$  is a descendant of  $v_{\mathcal{T}}$  and the other is a descendant of  $v_{\mathcal{F}}$  (see Fig. 2).*

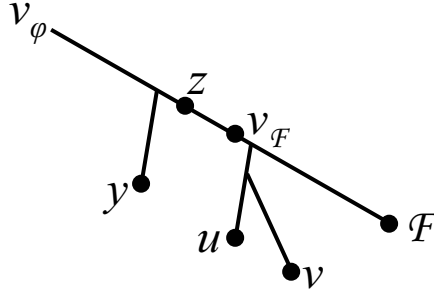


Figure 3: Proof of Lemma 2.1(2)

The above corollary leads to a natural transformation between trees satisfying **A1-2** and truth-assignments to the variables  $x_1..x_n$ :  $\sigma_T(x_i) = TRUE$  if  $x_i$  is a descendant of  $v_T$  and  $\sigma_T(x_i) = FALSE$  otherwise. The consistency of this assignment is guaranteed by Corollary 2.2. Lemma 2.1 (2), and **B1-2** lead to the following corollaries:

**Corollary 2.3.** *Let  $j \in \{1, \dots, m\}$ , and let  $\{a, b, c\} = \{1, 2, 3\}$ . If  $l_a^j$  and  $l_b^j$  are descendants of  $v_F$ , then  $y_c^j$  is also a descendant of  $v_F$ .*

**Corollary 2.4.** *If for some  $j = 1..m$ ,  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_F$ , then the bounds in **B3** cannot hold.*

*Proof.* By Corollary 2.3, if  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_F$ , then so are  $y_1^j, y_2^j, y_3^j$ . This implies, for instance, that  $C(T, y_1^j, y_2^j)$  is a descendant of  $v_F$  as well, so  $D_T(T; y_1^j y_2^j) \geq D_T(T, v_F) \geq 2\beta + \alpha > \alpha$ , contradicting **B3**.  $\square$

Note the slackness (of  $2\beta$ ) we have in the contradiction concluding the proof. This slackness is used to prove hardness of approximation in Section 3. The following lemma concludes the discussion of unsatisfiable formulae:

**Lemma 2.5.** *If  $T$  is an edge-weighted tree over the set of taxa  $S_\varphi$  satisfying all bounds in **A1-B3**, then the assignment  $\sigma_T$  satisfies the formula  $\varphi$ .*

*Proof.* Assume, to the contrary, that  $\sigma_T$  does not satisfy some clause  $c_j$  of  $\varphi$ . Then, by definition of  $\sigma_T$ , the taxa  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_F$ , and so by Corollary 2.4 the bounds in **B3** cannot hold for  $T$ .  $\square$

Lemma 2.5 is later used to ensure requirement **UNSAT**. To show that **SAT** holds we first prove the following lemma:

**Lemma 2.6.** *If the formula  $\varphi$  is satisfiable, then there exists a tree  $T$  over the set of taxa  $S_\varphi$ , satisfying **A1-B3** with equality.*

*Proof.* Let  $\sigma$  be a satisfying assignment of  $\varphi$ . We will construct a tree  $T$  with only two internal vertices  $v_T, v_F$ , and one internal edge of weight  $2\beta$  connecting  $v_T$  and  $v_F$ . All external edges are of weight  $\alpha$ , and all taxa are either connected

to  $v_T$  or  $v_F$ .  $T$  is connected to  $v_T$ ,  $F$  is connected to  $v_F$ , and a literal taxon  $l_a^j$  is connected to  $v_T$  if  $\sigma(l_a^j) = TRUE$ , and to  $v_F$  otherwise. It is easy to see that the bounds in **A1-2** are satisfied by such a tree.

Taxa of the form  $y_a^j$  are connected according to the following scheme (see Fig. 4): if  $l_a^j$  is **the only** literal in clause  $c_j$  satisfied by  $\sigma$ , connect  $y_a^j$  to  $v_F$ ; otherwise connect it to  $v_T$ . **B2** is clearly satisfied by this construction. **B3** is satisfied, since at most one  $y$ -taxon is connected to  $v_F$  for each clause. **B1** is satisfied, since for  $\{a, b, c\} = \{1, 2, 3\}$  the following holds: If  $l_a^j, l_b^j$  are connected to  $v_F$  ( $v_T$ ) then  $y_c^j$  is connected to  $v_F$  ( $v_T$  resp.) as well.  $\square$

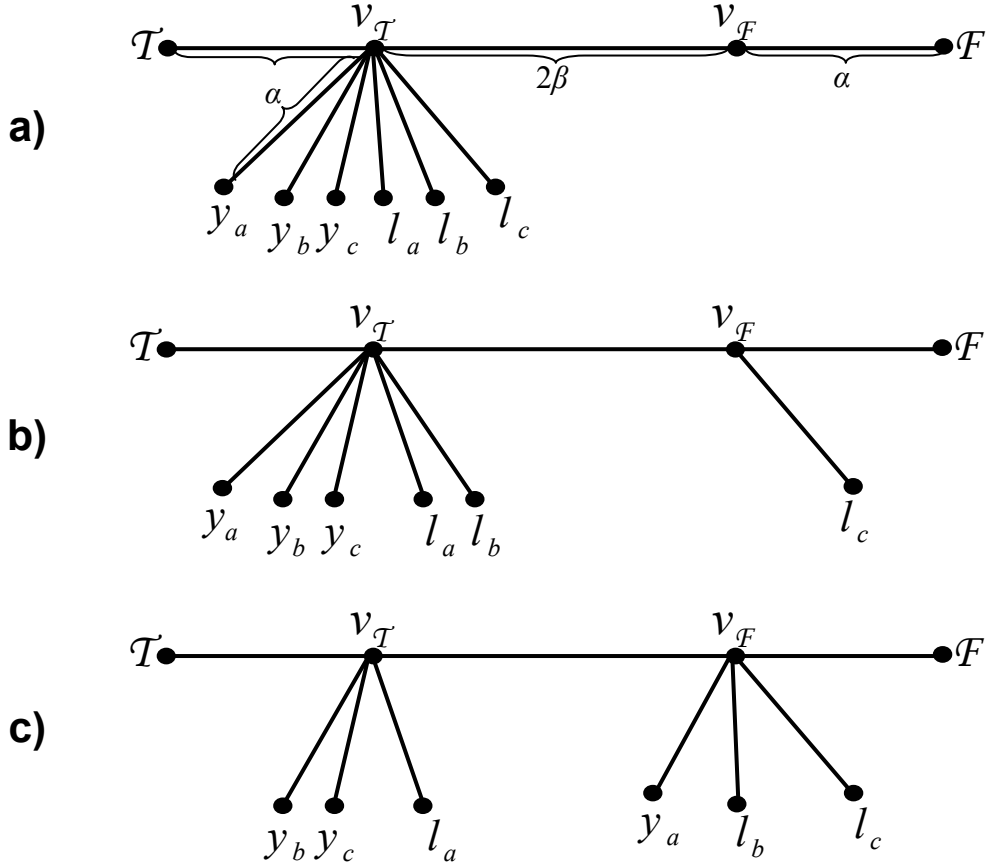


Figure 4: **Construction of a tree given a satisfying assignment.** The figure illustrates how to connect the  $y$ -taxa for each type of satisfied clause:

- a) All literals are satisfied (assigned *TRUE*).
- b) Two literals ( $l_a, l_b$ ) are satisfied.
- c) One literal ( $l_a$ ) is satisfied.

We now describe the reduction of the formula  $\varphi$  to a triplet-dissimilarity table  $\tau_\varphi$ : Entries of  $\tau_\varphi$  corresponding to distances bounded in **A1-B3** are set to enforce the corresponding bounds, as discussed in page 4. The rest of the entries in  $\tau_\varphi$ , and the constant  $\Delta$ , are set so that  $|D_T(i;jk) - \tau_\varphi(i;jk)| \leq \Delta$  will hold for *all* taxon triplets in the tree  $T$  described in the proof of Lemma 2.6, as follows. First, for all distinct  $i, j, k \in S_\varphi$  we have  $D_T(i;jk) \in [\alpha, \alpha + 2\beta]$  (since all external edges are of length  $\alpha$ , and the single internal edge is of length  $2\beta$ ). So we set the corresponding entries of  $\tau_\varphi$  (which do not appear in **A1-B3**) to  $\alpha + \beta$ , and we set  $\Delta = \beta$ . This guarantees that  $|D_T(i;jk) - \tau_\varphi(i;jk)| \leq \Delta$  for the corresponding entries. Similarly, for all distinct  $i, j \in S_\varphi$  we have  $D_T(i;jj) = D_T(i,j) \in [2\alpha, 2\alpha + 2\beta]$ , so we set corresponding entries of  $\tau_\varphi$  to  $2\alpha + \beta$ . Thus, the entries of the triplet-dissimilarity table  $\tau_\varphi$  are defined according to the following rules:

- $\tau_\varphi(\mathcal{T}; \mathcal{FF}) = \tau_\varphi(\mathcal{F}; \mathcal{TT}) = 2\alpha + 3\beta$  (**A1**)
- $\forall i = 1..n : \tau_\varphi(\mathcal{F}; x_i \bar{x}_i) = \tau_\varphi(\mathcal{T}; x_i \bar{x}_i) = \alpha - \beta$  (**A2**)
- $\forall j = 1..m : \tau_\varphi(y_1^j; l_2^j l_3^j) = \tau_\varphi(y_2^j; l_1^j l_3^j) = \tau_\varphi(y_3^j; l_1^j l_2^j) = \alpha - \beta$  (**B1**)
- $\forall j = 1..m : \tau_\varphi(\mathcal{T}; y_1^j y_2^j) = \tau_\varphi(\mathcal{T}; y_1^j y_3^j) = \tau_\varphi(\mathcal{T}; y_2^j y_3^j) = \alpha - \beta$  (**B3**)
- $\forall \{s, t\} (\neq \{\mathcal{T}, \mathcal{F}\}) \subseteq S_\varphi : \tau_\varphi(s; tt) = \tau_\varphi(t; ss) = 2\alpha + \beta$  (arbitrary pairwise-distances)
- For all other entries :  $\tau_\varphi(s; tu) = \alpha + \beta$  (arbitrary triplet-distances and **B2**)

We conclude with the following lemma:

**Lemma 2.7.** *Let  $\varphi$  be a satisfiable formula, and let  $\tau_\varphi$  be the triplet-dissimilarity table as defined above, using any values for  $\alpha, \beta$  s.t.  $\alpha \geq \beta > 0$ . Then there exists a tree  $T$  over the set of taxa  $S_\varphi$ , such that  $\|D_T, \tau_\varphi\|_\infty \leq \beta$ .*

*Proof.* The tree  $T$  corresponding to an assignment  $\sigma$  satisfying  $\varphi$  (as described in the proof of Lemma 2.6) fulfills this requirement. The proof follows directly from the above discussion.  $\square$

**Theorem 2.8.** *The decision version of the problem of finding a tree best fitting a given a triplet-dissimilarity table under the  $\ell_\infty$  norm is NP-Hard.*

*Proof.* By the polynomial-time reduction from 3SAT described above. The reduction  $\varphi \mapsto (\tau_\varphi, \Delta)$  is clearly polynomial (requirement **POLY**). By Lemma 2.7, if  $\varphi$  is satisfiable then there exists a tree  $T$ , s.t.  $\|D_T, \tau_\varphi\|_\infty \leq \Delta$  (**SAT**). If, on the other hand,  $\varphi$  is unsatisfiable, then by Lemma 2.5 there is no tree satisfying **A1-B3**. Due to the construction of  $\tau_\varphi$ , this means there is *no* tree  $T$ , s.t.  $\|D_T, \tau_\varphi\|_\infty \leq \Delta$  (**UNSAT**).  $\square$



### 3 Hardness of Approximation of The ‘*Best-Fit to Triplets Under $\ell_\infty$* ’ Problem

We prove hardness of approximation of this problem by showing that the reduction described in Section 2 satisfies stronger requirements:

**SAT’** If  $\varphi$  is satisfiable, then there is a tree  $T$  s.t.  $\|D_T, \tau_\varphi\|_\infty \leq \beta$ .

**UNSAT’** If  $\varphi$  is unsatisfiable, then for every tree  $T$ ,  $\|D_T, \tau_\varphi\|_\infty \geq 1.4\beta$ .

The first requirement is exactly **SAT** as phrased in the previous section, and so it follows from Lemma 2.7. **UNSAT’** requires proving a stronger version of Lemma 2.5, for a  $\delta$ -relaxed version of inequalities **A1-B3**, for some positive  $\delta$  which will be defined soon.

**A’1**  $D_T(\mathcal{T}, \mathcal{F}) \geq 2\alpha + 2\beta - \delta$

**A’2**  $\forall i = 1..n : D_T(\mathcal{F}; x_i \bar{x}_i) \leq \alpha + \delta ; D_T(\mathcal{T}; x_i \bar{x}_i) \leq \alpha + \delta$

**B’1**  $\forall j = 1..m : D_T(y_1^j; l_2^j l_3^j) \leq \alpha + \delta ; D_T(y_2^j; l_1^j l_3^j) \leq \alpha + \delta ; D_T(y_3^j; l_1^j l_2^j) \leq \alpha + \delta$

**B’2**  $\forall j = 1..m : D_T(y_1^j; \mathcal{T}\mathcal{F}) \geq \alpha - \delta ; D_T(y_2^j; \mathcal{T}\mathcal{F}) \geq \alpha - \delta ; D_T(y_3^j; \mathcal{T}\mathcal{F}) \geq \alpha - \delta$

**B’3**  $\forall j = 1..m : D_T(\mathcal{T}; y_1^j y_2^j) \leq \alpha + \delta ; D_T(\mathcal{T}; y_1^j y_3^j) \leq \alpha + \delta ; D_T(\mathcal{T}; y_2^j y_3^j) \leq \alpha + \delta$

Let  $T$  be a tree satisfying **A’1-B’3** above for some  $\delta < \frac{2\beta}{5}$ , and let  $v_\varphi$  be the mid-point of  $path(\mathcal{F}, \mathcal{T})$ . Let  $v_{\mathcal{T}}$  and  $v_{\mathcal{F}}$  to be the points whose distance from  $v_\varphi$  is **exactly**  $\beta - 1.5\delta$  on the paths to  $\mathcal{T}$  and  $\mathcal{F}$  respectively. Note that by **A’1**,  $D_T(\mathcal{F}, v_{\mathcal{F}}) \geq \alpha + \delta$  and  $D_T(\mathcal{T}, v_{\mathcal{T}}) \geq \alpha + \delta$  (see Fig. 5). Using this, we prove a stronger version of Lemma 2.1:

**Lemma 3.1.** *For all taxa  $u, v, y$  in  $T$ , we have*

1. *If  $D_T(\mathcal{F}; uv) \leq \alpha + \delta$  and  $D_T(\mathcal{T}; uv) \leq \alpha + \delta$ , then either  $u$  is a descendant of  $v_{\mathcal{T}}$  and  $v$  is a descendant of  $v_{\mathcal{F}}$  or vice versa.*
2. *If both  $u$  and  $v$  are descendants of  $v_{\mathcal{F}}$ , and  $D_T(y; uv) < D_T(y, \mathcal{T}\mathcal{F}) + D_T(v_{\mathcal{T}}, v_{\mathcal{F}})$  then  $y$  is not a descendant of  $v_{\mathcal{T}}$ .*

*Proof.*

1. As in the proof of Lemma 2.1(1), since  $D_T(\mathcal{F}; uv) + D_T(\mathcal{T}; uv) < D_T(\mathcal{F}, \mathcal{T})$  we have that  $C(u, \mathcal{T}, \mathcal{F})$  and  $C(v, \mathcal{T}, \mathcal{F})$  are distinct vertices on  $path(\mathcal{F}, \mathcal{T})$ , one at distance at most  $\alpha + \delta$  from  $\mathcal{F}$  and the other at distance at most  $\alpha + \delta$  from  $\mathcal{T}$ .
2. Assume, to the contrary, that  $y$  is a descendant of  $v_{\mathcal{T}}$ . Since both  $u$  and  $v$  are descendants of  $v_{\mathcal{F}}$ , the path from  $y$  to  $path(u, v)$  must contain both  $v_{\mathcal{T}}$  and  $v_{\mathcal{F}}$ , which are both on  $path(\mathcal{F}, \mathcal{T})$  (see Fig. 6). Thus we must have that  $D_T(y, uv) \geq D_T(y, \mathcal{T}\mathcal{F}) + D_T(v_{\mathcal{T}}, v_{\mathcal{F}})$ , contradicting the assumption.

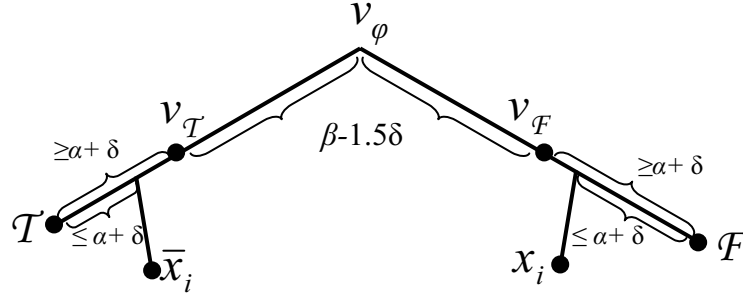


Figure 5: The Topology of a tree satisfying **A'1-2**

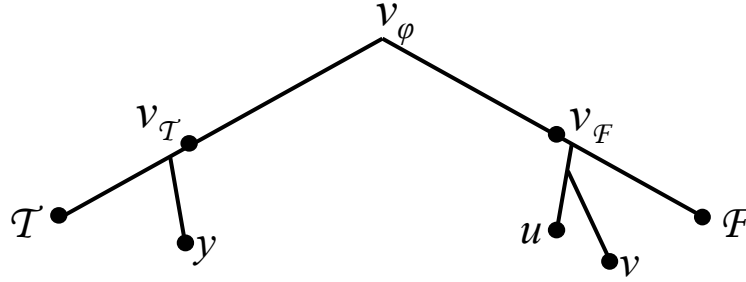


Figure 6: Proof of Lemma 3.1(2)

□

The following corollaries follow from Lemma 3.1 and **A'1-B'3**:

**Corollary 3.2.** Assume that  $\delta < \frac{2\beta}{3}$ . Then for each  $i = 1..n$ , one of the vertices  $x_i, \bar{x}_i$  is a descendant of  $v_T$  and the other is a descendant of  $v_F$ .

**Corollary 3.3.** Assume that  $\delta < \frac{2\beta}{5}$ . Let  $j$  be in  $\{1, .., m\}$ , and let  $\{a, b, c\} = \{1, 2, 3\}$ . If  $l_a^j$  and  $l_b^j$  are descendants of  $v_F$ , then  $y_c^j$  is not a descendant of  $v_T$ .

Notice that the relaxation of the bounds prevents us from proving (as in Corollary 2.3) that the  $y$ -taxa are descendants of  $v_F$ . However, the weaker claim in Corollary 3.3 is sufficient to contradict the bounds in **B'3**, due to the slackness we had in the proof of Corollary 2.4:

**Corollary 3.4.** If for some  $j = 1..m$ ,  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_F$ , then the bounds in **B'3** cannot hold.

*Proof.* By Corollary 3.3, if  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_F$ , then none of  $y_1^j, y_2^j, y_3^j$  are descendants of  $v_T$ . This implies, for instance, that  $C(\mathcal{T}, y_1^j, y_2^j)$  is not a descendant of  $v_T$  as-well, so  $D_T(\mathcal{T}; y_1^j y_2^j) > D_T(\mathcal{T}, v_T) \geq \alpha + \delta$  (by definition of  $v_T$ ), contradicting **B'3**. □

This corollary leads us to the following:

**Lemma 3.5.** *If there is an edge-weighted tree  $T$  over the set of taxa  $S_\varphi$  satisfying  $\|D_T, \tau_\varphi\|_\infty < 1.4\beta$ , then the formula  $\varphi$  is satisfiable.*

*Proof.* A tree satisfying  $\|D_T, \tau_\varphi\|_\infty \leq \beta + \delta$  satisfies the  $\delta$ -relaxed bounds in **A'1-B'3** as well. So if  $\|D_T, \tau_\varphi\|_\infty < 1.4\beta$ , then  $T$  satisfies the  $\delta$ -relaxed bounds for  $\delta = \|D_T, \tau_\varphi\|_\infty - \beta < \frac{2}{5}\beta$ . Now since **A'1-2** hold, the assignment  $\sigma_T$  is well defined (Corollary 3.2). Assume that  $\sigma_T$  does not satisfy some clause  $c_j$  of  $\varphi$ . Then, by definition of  $\sigma_T$ ,  $l_1^j, l_2^j, l_3^j$  are all descendants of  $v_{\mathcal{F}}$ , and by Corollary 3.3 the bounds in **B'3** cannot hold, in contradiction.  $\square$

For a distance table  $\tau$ , let  $OPT(\tau)$  be the minimal value  $k$ , for which there is a tree  $T$  s.t.  $\|D_T, \tau\|_\infty \leq k$ . By Lemma 2.6, if  $\varphi$  is satisfiable then  $OPT(\tau_\varphi) \leq \beta$ , and by Lemma 3.5, if  $\varphi$  is unsatisfiable then  $OPT(\tau_\varphi) \geq 1.4\beta$ . Thus if there was a polynomial time algorithm  $\mathcal{A}$  which is guaranteed to approximate  $OPT(\tau)$  within a factor *smaller* than 1.4, then satisfiability of a formula  $\varphi$  could be determined by executing  $\mathcal{A}$  on  $\tau_\varphi$  and obtaining  $k = \|\tau_\varphi, \mathcal{A}(\tau_\varphi)\|_\infty$ . If  $k < 1.4\beta$  then  $\varphi$  must be satisfiable, and if  $k \geq 1.4\beta$  then  $\varphi$  is unsatisfiable. Hence it is NP-hard to find a tree which approximates the optimal  $\ell_\infty$  distance to a given triplet-dissimilarity table by a ratio smaller than 1.4.

## 4 Hardness of Approximation of Maximal Distortion and Other Implied Results

We now use the reductions presented in the previous sections to obtain several related hardness results. As the constructions are similar to these in previous sections, most proofs in this section are only sketched.

### 4.1 Hardness of Approximation of Maximal Distortion

Recall the *maximal distortion* between two triplet-dissimilarity tables :

$$MaxDist(\tau_1, \tau_2) \triangleq \max_{i,j,k \in S} \left\{ \frac{\tau_1(i;jk)}{\tau_2(i;jk)} \right\} \cdot \max_{i,j,k \in S} \left\{ \frac{\tau_2(i;jk)}{\tau_1(i;jk)} \right\} \quad (\text{where } 0/0 \triangleq 1)$$

We use a reduction similar to the one in Section 2 to prove that  $MaxDist$  of the closest tree cannot be approximated by any multiplicative factor. First, note that scaling a tree by a multiplicative factor does not affect its  $MaxDist$  from a given triplet-dissimilarity table. In other words,  $MaxDist(\tau, D_T) = MaxDist(\tau, D_{[\gamma T]})$ , where  $\gamma T$  is the weighted tree obtained by multiplying edge weights of  $T$  by the positive constant  $\gamma$ . This means that if there is a tree  $T$  s.t.  $MaxDist(\tau, D_T) \leq \rho$ , then there is a tree  $T'$  (obtained by re-scaling  $T$ ) s.t.

$$\max \left\{ \max_{i,j,k \in S} \left\{ \frac{\tau(i;jk)}{D_{T'}(i;jk)} \right\}, \max_{i,j,k \in S} \left\{ \frac{D_{T'}(i;jk)}{\tau(i;jk)} \right\} \right\} \leq \sqrt{\rho}.$$

A CNF formula  $\varphi$  is translated to a triplet-dissimilarity table  $\tilde{\tau}_\varphi$  which enforces the inequalities in **A1-B3** through bounds on maximal distortion as follows: An upper bound  $D_T(i; jk) \leq \omega$  is enforced by setting  $\tilde{\tau}_\varphi(i; jk) = \frac{\omega}{\sqrt{\rho}}$ , and a lower bound  $D_T(i; jk) \geq \omega$  is enforced by setting  $\tilde{\tau}_\varphi(i; jk) = \sqrt{\rho}\omega$ , where  $\rho \geq 1$  will soon be defined. By the argument raised above, a tree whose  $MaxDist$  from  $\tilde{\tau}_\varphi$  is at most  $\rho$  implies a tree satisfying all bounds. We now show how to set  $\rho$  and fill in the rest of the entries of  $\tilde{\tau}_\varphi$ , such that the tree  $T$  described in the proof of Lemma 2.6 satisfies  $MaxDist(\tilde{\tau}_\varphi, D_T) \leq \rho$ : Recall that in such a tree, triplet-dissimilarities not mentioned in **A1-B3** fall within the interval  $[\alpha, \alpha+2\beta]$  for distinct-taxa triplets, and within the interval  $[2\alpha, 2\alpha+2\beta]$  for taxon-pairs (see discussion following Lemma 2.6). In order to allow triplet dissimilarities within these intervals, we set  $\rho = \max\{\frac{\alpha+2\beta}{\alpha}, \frac{2\alpha+2\beta}{2\alpha}\} = 1 + 2\frac{\beta}{\alpha}$ , and set the relevant entries of  $\tilde{\tau}_\varphi$  to  $\sqrt{\rho} \cdot \alpha$  and  $\sqrt{\rho} \cdot 2\alpha$  (corresponding to distinct-taxa triplets and taxon-pairs respectively). The following lemma ensures that  $\tilde{\tau}_\varphi$  and  $\rho$  have the desired properties:

**Lemma 4.1.** *Let  $\alpha, \beta > 0$  be given, and let  $\rho = 1 + 2\frac{\beta}{\alpha}$ . Further, let  $\tilde{\tau}_\varphi$  be the triplet-dissimilarity table defined by  $\alpha, \beta$  and  $\rho$  as described above, then:*

**SAT''** *If  $\varphi$  is satisfiable, then there exists a tree  $T$  s.t.  $MaxDist(D_T, \tilde{\tau}_\varphi) \leq \rho$ .*

**UNSAT''** *If  $\varphi$  is unsatisfiable, then for every tree  $T$ ,  $MaxDist(D_T, \tilde{\tau}_\varphi) \geq \rho \left(1 + \frac{2\beta}{3\alpha}\right)$ .*

*Proof.* (an outline): **SAT''** is guaranteed by the tree construction described in the proof of Lemma 2.6 and by the value we chose for  $\rho$ , as discussed above. **UNSAT''** is proved by adjusting the proof in Section 3. First, we define a set of bounds **A''1-B''3**, obtained by a relaxation of **A1-B3** by a *multiplicative* factor of  $\delta > 1$  as follows:

$$\mathbf{A''1} \quad D_T(\mathcal{T}, \mathcal{F}) \geq (2\alpha + 2\beta)/\delta$$

$$\mathbf{A''2} \quad \forall i = 1..n : D_T(\mathcal{F}; x_i \bar{x}_i) \leq \alpha\delta ; D_T(\mathcal{T}; x_i \bar{x}_i) \leq \alpha\delta$$

$$\mathbf{B''1} \quad \forall j = 1..m : D_T(y_1^j; l_2^j l_3^j) \leq \alpha\delta ; D_T(y_2^j; l_1^j l_3^j) \leq \alpha\delta ; D_T(y_3^j; l_1^j l_2^j) \leq \alpha\delta$$

$$\mathbf{B''2} \quad \forall j = 1..m : D_T(y_1^j; \mathcal{T}\mathcal{F}) \geq \alpha/\delta ; D_T(y_2^j; \mathcal{T}\mathcal{F}) \geq \alpha/\delta ; D_T(y_3^j; \mathcal{T}\mathcal{F}) \geq \alpha/\delta$$

$$\mathbf{B''3} \quad \forall j = 1..m : D_T(\mathcal{T}; y_1^j y_2^j) \leq \alpha\delta ; D_T(\mathcal{T}; y_1^j y_3^j) \leq \alpha\delta ; D_T(\mathcal{T}; y_2^j y_3^j) \leq \alpha\delta$$

Next, we consider a tree satisfying the relaxed bounds, and define the internal points  $v_{\mathcal{T}}, v_{\mathcal{F}}$  to be at distance  $\frac{\alpha+\beta}{\delta} - \delta\alpha$  from  $v_\varphi$ . To prove the analogue of Lemma 3.1(1), it is required that  $D_T(\mathcal{F}, \mathcal{T})$  be strictly larger than  $2\alpha\delta$ , which by **A''1** reduces to  $\frac{\alpha+\beta}{\delta} - \delta\alpha > 0$ , i.e.  $\delta < \sqrt{1 + \frac{\beta}{\alpha}}$ . Corollary 3.3 is proven by the analogue of Lemma 3.1(2). For this we require  $\alpha/\delta + D_T(v_{\mathcal{F}}, v_{\mathcal{T}}) > \alpha\delta$ . Since  $D_T(v_{\mathcal{F}}, v_{\mathcal{T}}) = 2(\frac{\alpha+\beta}{\delta} - \delta\alpha)$ , this is equivalent to  $\delta < \sqrt{1 + \frac{2\beta}{3\alpha}}$ . This latter upper bound on  $\delta$  implies also the previous one, and hence if  $\varphi$  is unsatisfiable,

there is no tree satisfying the  $\delta$ -relaxed bounds in **A”1-B”3** for  $\delta < \sqrt{1 + \frac{2\beta}{3\alpha}}$ . In other words, there is no tree  $T$  satisfying:

$$\max \left\{ \max_{i,j,k \in S_\varphi} \left\{ \frac{\tilde{\tau}_\varphi(i;jk)}{D_T(i;jk)} \right\}, \max_{i,j,k \in S_\varphi} \left\{ \frac{D_T(i;jk)}{\tilde{\tau}_\varphi(i;jk)} \right\} \right\} < \sqrt{\rho \left( 1 + \frac{2\beta}{3\alpha} \right)}.$$

This means that there is no tree whose  $MaxDist$  from  $\tilde{\tau}_\varphi$  is less than  $\rho \left( 1 + \frac{2\beta}{3\alpha} \right)$ , as claimed.  $\square$

Now, assume there was a polynomial-time algorithm  $\mathcal{A}$  which given a triplet-dissimilarity table  $\tau$ , was guaranteed to return a tree whose  $MaxDist$  from  $\tau$  is at most  $K$ -times the  $MaxDist$  of the closest tree to  $\tau$ , for some constant  $K$ . Such an algorithm may be used to efficiently deduce whether a formula  $\varphi$  is satisfiable in the following way: given a formula  $\varphi$ , calculate  $\tilde{\tau}_\varphi$  with parameters  $\alpha, \beta$  s.t.  $K < 1 + \frac{2}{3} \frac{\beta}{\alpha}$ . Execute algorithm  $\mathcal{A}$  on this triplet-dissimilarity table to receive a tree  $T$ , and calculate  $r = MaxDist(\tilde{\tau}_\varphi, D_T)$ . Now, if  $r \leq K\rho$  (where  $\rho = 1 + 2\frac{\beta}{\alpha}$  as previously defined), then  $\varphi$  must be satisfiable due to **UNSAT”**. If, on the other hand,  $r > K\rho$ , then since  $\mathcal{A}$  guarantees a  $K$ -approximation, there is *no* tree whose  $MaxDist$  from  $\tilde{\tau}_\varphi$  is at most  $\rho$ . From **SAT”** follows that  $\varphi$  is unsatisfiable.

## 4.2 Fitting distances of distinct-taxa triplets

Triplet distance tables, as we defined them, contain entries corresponding to distinct-taxa triplets as well as entries corresponding to taxon-pairs (i.e.  $\tau(i;jj)$ ). In some scenarios it is more natural to separately address pairwise dissimilarities and triplet-dissimilarities. Therefore, we are interested in the problem of finding a best-fit tree to a triplet-dissimilarity table  $\tau$ , considering entries corresponding only to *distinct-taxa triplets*. The best-fit analysis can be done under any of the  $\ell_p$  norms or  $MaxDist$ . Results similar to the ones presented above apply in this case as well. The only modification required in order to adapt the reduction to this case is changing the bounds in **A1**, which correspond to the pairwise distance between  $\mathcal{T}$  and  $\mathcal{F}$ . To ensure a similar bound, we introduce an additional taxon into  $S_\varphi$ :  $\mathcal{F}'$ , and replace **A1** by:

$$\overline{\mathbf{A1}} \quad D_T(\mathcal{T}; \mathcal{F}\mathcal{F}') \geq 2\alpha + 2\beta$$

It is easy to see that this new bound implies the desired lower bound on the distance between  $\mathcal{T}$  and  $\mathcal{F}$  (i.e. **A1**). The original set of bounds is, therefore, equivalent to this one, and all claims proven for it apply here as well. The tree described in the proof of Lemma 2.6 is adapted to the introduction of  $\mathcal{F}'$ , by turning the original taxon  $\mathcal{F}$  into an internal vertex, and adding two zero-weight edges from this vertex to  $\mathcal{F}, \mathcal{F}'$ . All triplet-dissimilarities concerning  $\mathcal{F}'$  are set to be equal to their counterparts concerning  $\mathcal{F}$ . The analysis done in previous sections is easily adjusted to accommodate this modification of the reduction.

### 4.3 Best-Fit Ultrametric

It is possible to generalize all hardness results shown in this paper for ultrametrics as well. A weighted tree is called *ultrametric* if it contains a point which is equidistant from all leaves; this point may be an internal vertex, or a degenerate (degree-2) vertex situated on one of the edges. The problem of finding a best-fit ultrametric to a given *dissimilarity matrix* under  $\ell_\infty$  (and *MaxDist*) was shown to have a polynomial-time algorithm in [9, 7].

In the case of triplet-dissimilarities, the same reductions presented in sections 3 and 4.1 imply that it is NP-hard to find (and to approximate) a best-fit ultrametric under the  $\ell_\infty$  norm, as well as *MaxDist*. To see this, observe that if  $\varphi$  is unsatisfiable, then the lower bounds proved for **UNSAT'** in Lemma 3.5 and for **UNSAT''** in Lemma 4.1 (for  $\ell_\infty$  and *MaxDist* resp.) are clearly valid when the trees are restricted to be ultrametrics. We are left to show that if  $\varphi$  is satisfiable then there is an *ultrametric* tree satisfying all bounds. This follows from the fact that the construction described in the proof of Lemma 2.6 yields an ultrametric tree, since the internal point  $v_\varphi$  is at the same distance  $(\alpha + \beta)$  from all taxa.

## 5 A Constant-Rate Approximation Scheme

In this section we present a constant-rate approximation algorithm for the problem of finding a closest tree under  $\ell_\infty$  to a given triplet-dissimilarity table. Our algorithm is based on an approximation algorithm for the corresponding problem concerning pairwise-dissimilarities. The main result is stated in the following theorem.

**Theorem 5.1.** *A polynomial time  $r$ -approximation algorithm for finding a tree closest under  $\ell_\infty$  to a given dissimilarity matrix implies a polynomial time  $(\frac{3}{2}r + 6)$ -approximation algorithm for finding a tree closest under  $\ell_\infty$  to a given triplet-dissimilarity table.*

Our approximation algorithm,  $\mathcal{APP}$ , consists of two stages:

$\mathcal{APP1}$ . Given a triplet-dissimilarity table  $\tau$  over taxon-set  $S$ , calculate a dissimilarity matrix  $D^\tau$  over  $S$  as follows:  $\forall i, j \in S : D^\tau(i, j) = \tau(i, jj)$ .

$\mathcal{APP2}$ . Execute the  $r$ -approximation algorithm on  $D^\tau$  to obtain an edge-weighted tree  $T^{out}$ .

To analyze the approximation ratio of the above algorithm, we start with some notations. For an arbitrary taxon-pair  $i, j \in S$ , denote  $D_{min}^\tau(i, j) = \min_{k \in S} \{\tau(i, jk) + \tau(j, ik)\}$ , and similarly  $D_{max}^\tau(i, j) = \max_{k \in S} \{\tau(i, jk) + \tau(j, ik)\}$ . Furthermore, denote by  $I^\tau = \max_{i, j \in S} \{D_{max}^\tau(i, j) - D_{min}^\tau(i, j)\}$  the maximum difference between  $D_{max}^\tau$  and  $D_{min}^\tau$ . The following lemma contains two basic inequalities required for the proof of our approximation result.

**Lemma 5.2.** *Let  $\tau$  be a triplet-dissimilarity table, and  $D^\tau$  be the corresponding dissimilarity matrix defined in APP1. Let further  $T$  be an edge-weighted tree with corresponding additive distance matrix  $D_T$  and triplet-dissimilarity table  $\tau_T$ . Then we have the following:*

$$\frac{1}{4}I^\tau \leq \|\tau, \tau_T\|_\infty \leq \frac{3}{2}(I^\tau + \|D^\tau, D_T\|_\infty) .$$

*Proof.* First we prove that  $\frac{1}{4}I^\tau \leq \|\tau, \tau_T\|_\infty$ . Let  $i, j$  be a taxon pair s.t.  $D_{max}^\tau(i, j) - D_{min}^\tau(i, j) = I^\tau$ , and let  $k_{max}$  be a taxon s.t.  $D_{max}^\tau(i, j) = \tau(i; jk_{max}) + \tau(j; ik_{max})$ . Since  $\tau_T(i; jk) + \tau_T(j; ik) = D_T(i, j)$ , for all  $k \in S$ , then:

$$\begin{aligned} D_{max}^\tau(i, j) - D_T(i, j) &= [\tau(i; jk_{max}) + \tau(j; ik_{max})] - [\tau_T(i; jk_{max}) + \tau_T(j; ik_{max})] \\ &= [\tau(i; jk_{max}) - \tau_T(i; jk_{max})] + [\tau(j; ik_{max}) - \tau_T(j; ik_{max})] \\ &\leq 2\|\tau, \tau_T\|_\infty \end{aligned} \tag{1}$$

Similarly, if  $k_{min}$  is a taxon s.t.  $D_{min}^\tau(i, j) = \tau(i; jk_{min}) + \tau(j; ik_{min})$ , then:

$$D_{min}^\tau(i, j) - D_T(i, j) \geq -2\|\tau, \tau_T\|_\infty \tag{2}$$

Now, since  $D_{max}^\tau(i, j) - D_{min}^\tau(i, j) = I^\tau$ , then by subtracting (2) from (1) we get  $I^\tau \leq 4\|\tau, \tau_T\|_\infty$ , thus proving the left inequality.

We now turn to prove the right inequality of the lemma. Given an arbitrary taxon-triplet  $i, j, k \in S$ , denote  $\varepsilon(i; jk) = \tau(i; jk) - \tau_T(i; jk)$ . We will show that  $|\varepsilon(i; jk)| \leq \frac{3}{2}(I^\tau + \|D^\tau, D_T\|_\infty)$ . First,

$$\begin{aligned} |\varepsilon(i; jk) + \varepsilon(j; ik)| &= |[\tau(i; jk) - \tau_T(i; jk)] + [\tau(j; ik) - \tau_T(j; ik)]| \\ &= |[\tau(i; jk) + \tau(j; ik)] - [\tau_T(i; jk) + \tau_T(j; ik)]| \\ &= |\tau(i; jk) + \tau(j; ik) - D_T(i, j)| \\ &\leq |\tau(i; jk) + \tau(j; ik) - D^\tau(i, j)| + |D^\tau(i, j) - D_T(i, j)| \\ &\leq |\tau(i; jk) + \tau(j; ik) - D^\tau(i, j)| + \|D^\tau, D_T\|_\infty \\ &\leq I^\tau + \|D^\tau, D_T\|_\infty \end{aligned}$$

The last inequality follows from the fact that  $D^\tau(i, j) = \tau(i; jj) + \tau(j; ij)$ . Using a similar line of argument we get  $|\varepsilon(i; kj) + \varepsilon(k; ij)|, |\varepsilon(j; ki) + \varepsilon(k; ji)| \leq I^\tau + \|D^\tau, D_T\|_\infty$  as well. This is used to obtain the desired bound as follows:

$$\begin{aligned} |\varepsilon(i; jk)| &= \frac{1}{2} |[\varepsilon(i; jk) + \varepsilon(j; ik)] + [\varepsilon(i; kj) + \varepsilon(k; ij)] - [\varepsilon(j; ki) + \varepsilon(k; ji)]| \\ &\leq \frac{1}{2} (|\varepsilon(i; jk) + \varepsilon(j; ik)| + |\varepsilon(i; kj) + \varepsilon(k; ij)| + |\varepsilon(j; ki) + \varepsilon(k; ji)|) \\ &\leq \frac{3}{2} (I^\tau + \|D^\tau, D_T\|_\infty) \end{aligned}$$

□

Our main result (Theorem 5.1) is directly implied by the following lemma:

**Lemma 5.3.** *Given a triplet-dissimilarity table  $\tau$ , denote by  $\tau^{out}$  the triplet-dissimilarity table induced by the output tree  $T^{out}$  returned by the algorithm  $\mathcal{APP}$ . Then for every triplet-dissimilarity table  $\tau_T$  induced by an arbitrary edge-weighted tree  $T$ , we have:*

$$\|\tau, \tau^{out}\|_\infty \leq \left(\frac{3}{2}r + 6\right)\|\tau, \tau_T\|_\infty$$

*Proof.* Denote by  $D^\tau$  the dissimilarity matrix computed in  $\mathcal{APP}1$ , and by  $D^{out}$  and  $D_T$  the metrics induced over the leaves of  $T^{out}$  and  $T$ , respectively. The lemma is proved by the following sequence of inequalities:

$$\|\tau, \tau^{out}\|_\infty \leq \frac{3}{2} (I^\tau + \|D^\tau, D^{out}\|_\infty) \quad (3)$$

$$\leq \frac{3}{2} (4\|\tau, \tau_T\|_\infty + \|D^\tau, D^{out}\|_\infty) \quad (4)$$

$$\leq \frac{3}{2} (4\|\tau, \tau_T\|_\infty + r\|D^\tau, D_T\|_\infty) \quad (5)$$

$$\leq \left(\frac{3}{2}r + 6\right)\|\tau, \tau_T\|_\infty \quad (6)$$

(3) and (4) follow from the right and left inequalities of Lemma 5.2, respectively. The approximation ratio of the algorithm executed during  $\mathcal{APP}2$  implies (5). (6) follows from the fact that  $\|D^\tau, D_T\|_\infty \leq \|\tau, \tau_T\|_\infty$ , which holds since  $D^\tau(i, j) = \tau(i; jj)$  and  $D_T(i, j) = \tau_T(i; jj)$  for every taxon-pair  $i, j \in S$ .  $\square$

By Theorem 5.1, the 3-approximation algorithms for pairwise dissimilarities presented in [1, 8] imply a  $10\frac{1}{2}$  approximation algorithm for triplet dissimilarities. However, the 3-approximation ratio of the algorithms in [1, 8] is proved under the assumption that the input dissimilarity matrix is a *distance metric*. Therefore, this bound (of  $10\frac{1}{2}$ ) is valid only if the matrix  $D^\tau$  computed in  $\mathcal{APP}2$  satisfies the triangle inequality. When the triangle inequality is not assumed, the analysis in [1, 8] can be modified to yield a 6-approximation ratio, rather than the original 3-approximation. This 6-approximation algorithm leads, by Theorem 5.1, to a 15-approximation of the closest tree to an arbitrary triplet-dissimilarity table under  $\ell_\infty$ .

## 6 Discussion

In this paper we discussed the hardness of several problems of fitting a phylogenetic tree to a given triplet-dissimilarity table. This question is motivated by several recent works which reconstruct trees using triplet-dissimilarities [11, 10, 8]. The optimization criteria considered in this paper are the  $\ell_\infty$  norm and *MaxDist*, which measure the maximum discrepancy (difference and ratio resp.) between the input dissimilarities and the ones induced by the desired tree. It is interesting whether similar hardness results apply also for other common distance measures such as the  $\ell_1$  and  $\ell_2$  norms.



The construction in Lemma 2.6 which implies our basic NP-hardness result yields a tree containing two vertices of very high degree. Common models for phylogenetic trees assume a binary tree (meaning that all internal vertices have degree 3). Furthermore, edge-weights are assumed to lie within an interval  $[w_{\min}, w_{\max}]$ , where  $w_{\min}$  and  $w_{\max}$  are strictly positive constants independent of the size of the tree. It is interesting whether our NP-hardness results apply also when introducing these assumptions on the desired tree, and specifically what is the smallest ratio between  $w_{\max}$  and  $w_{\min}$  mentioned above which still gives similar hardness results. Can this ratio be a constant independent on  $n$ ? Does the NP-hardness result apply also for binary trees with *uniform* edge weights?

Another question relates to the approximation ratio given in Section 5. Possibly, a better approximation ratio may be obtained by a closer analysis of the algorithms in [1, 8].

## Acknowledgements

We thanks the anonymous referees for their insightful comments.

## References

- [1] R. Agarwala, V. Bafna, M. Farach, M. Paterson, and M. Thorup. On the approximability of numerical taxonomy (fitting distances by tree metrics). *SIAM Journal on Computing*, 28(3):1073–1085, June 1999.
- [2] Y. Bartal, N. Linial, M. Mendel, and A. Naor. Low dimensional embeddings of ultrametrics. *Eur. J. Comb.*, 25(1):87–92, 2004.
- [3] P. Buneman. The recovery of trees from measures of dissimilarity. *Mathematics in the Archeological and Historical Sciences*, pages 387–395, 1971.
- [4] W. Day. Computational complexity of inferring phylogenies from dissimilarity matrices. *Bulletin of Mathematical Biology*, 49(4):461–467, 1987.
- [5] P. Erdos, M. Steel, L. Szekely, and T. Warnow. A few logs suffice to build (almost) all trees (I). *Random Structures and Algorithms*, 14:153–184, 1999.
- [6] P. Erdos, M. Steel, L. Szekely, and T. Warnow. A few logs suffice to build (almost) all trees (II). *Theoretical Computer Science*, 221:77–118, 1999.
- [7] M. Farach, S. Kannan, and T. Warnow. A robust model for finding optimal evolutionary trees. *Algorithmica*, 13(1/2):155–179, January 1995.
- [8] I. Gronau and S. Moran. Neighbor joining algorithms for inferring phylogenies via LCA-distances. *Journal of Computational Biology*, 14(1):1–15, 2007.

- [9] M. Krivánek. The complexity of ultrametric partitions on graphs. *Inform. Process. Lett.*, 27:265–270, 1988.
- [10] D. Levy, R. Yoshida, and L. Pachter. Beyond pairwise distances: Neighbor-joining with phylogenetic diversity estimates. *Mol Biol Evol*, 23(3):491–498, 2006.
- [11] V. Ranwez and O. Gascuel. Improvement of distance-based phylogenetic methods by a local maximum likelihood approach using triplets. *Mol Biol Evol*, 19(11):1952–1963, 2002.
- [12] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Mol Biol Evol*, 4:406–425, 1987.
- [13] S. Sattath and A. Tversky. Additive similarity trees. *Psychometrika*, 42(3):319–345, 1977.
- [14] M. Waterman, T. Smith, M. Singh, and W. Beyer. Additive evolutionary trees. *J Theor Biol*, 64(2):199–213, January 1977.