

# ROBUSTNESS OF PHYLOGENETIC INFERENCE BASED ON MINIMUM EVOLUTION

FABIO PARDI<sup>1</sup>, SYLVAIN GUILLEMOT<sup>1</sup>, OLIVIER GASCUEL<sup>1</sup>

**ABSTRACT.** Minimum evolution is the guiding principle of an important class of distance-based phylogeny reconstruction methods, including neighbor-joining (NJ), which is the most cited tree inference algorithm to date. The minimum evolution principle involves searching for the tree with minimum length, where the length is estimated using various least-squares criteria. Since evolutionary distances cannot be known precisely but only estimated, it is important to investigate the robustness of phylogenetic reconstruction to imprecise estimates for these distances. The *safety radius* is a measure of this robustness: it consists of the maximum relative deviation that the input distances can have from the correct distances, without compromising the reconstruction of the correct tree structure. Answering some open questions, we here derive the safety radius of two popular minimum evolution criteria: balanced minimum evolution (BME) and minimum evolution based on ordinary least squares (OLS+ME). Whereas BME has a radius of  $\frac{1}{2}$ , which is the best achievable, OLS+ME has a radius tending to 0 as the number of taxa increases. This difference may explain the gap in reconstruction accuracy observed in practice between OLS+ME and BME (which forms the basis of popular programs such as NJ and FastME).

## 1. INTRODUCTION

Minimum evolution methods for reconstructing phylogenetic trees [15, 20] are based upon the following informal idea: given a matrix of distances between each pair of taxa in a set, reconstruct the phylogenetic tree for these taxa that implies the minimum amount of evolution in order to explain the given distances. In practice a method to estimate the length of any possible tree structure is specified and minimum evolution aims to reconstruct the tree with the minimum length estimate.

Several types of minimum evolution can then be defined based on the different methods for estimating the tree length. In this paper we deal with two important versions of minimum evolution: the one based on an ordinary least squares technique for length estimation, and the one based on an increasingly popular “balanced” technique. We call these versions OLS+ME and BME, respectively.

OLS+ME is based on a method for estimating branch lengths which originated from the work by Cavalli-Sforza and Edwards [6, 7] and was the first version of minimum evolution to be

---

*Key words and phrases.* phylogenetics, distance methods, minimum evolution, least squares, safety radius.

<sup>1</sup>Méthodes et Algorithmes pour la Bioinformatique, LIRMM, CNRS - Université de Montpellier, 161 rue Ada, 34392 - Montpellier, France

proposed [15, 20]. It is introduced in more detail in the next section. Interestingly, neighbor-joining (NJ) — one of the most popular algorithms for tree reconstruction — has been shown to have some connections with OLS+ME [23, 18].

However, further mathematical investigation showed that NJ is much more tightly related to the other criterion we study, BME: in fact, NJ turns out to be a greedy agglomerative algorithm aiming to construct the optimal tree with respect to this criterion [10, 14]. The vast popularity of NJ (more than 20,000 citations for the original paper [23] on Google Scholar), the fact that BME length estimates are biologically meaningful variants of the least squares estimates [9], the computational efficiency of the proposed heuristics for BME [8] and, last but not least, the high accuracy of these heuristics observed in simulation studies [8, 26, 9] justify the growing importance of BME as a criterion for reconstructing phylogenetic trees. Again, a more detailed account of BME is given in the next section.

BME and OLS+ME share an important property: if the distances in input perfectly correspond to the distances between leaves in a tree with branch lengths, then these criteria identify this tree as the correct one. Evolutionary distances are estimated using genetic sequences or any other comparative data from the taxa under consideration. Ideal estimation procedures should ensure that, as the amount of compared data increases, the estimated distances converge to those in the correct phylogenetic tree, which means that the property just stated is equivalent to the statistical *consistency* of these tree reconstruction methods. In fact consistency is a desirable property for *any* minimum evolution method [27], and more generally for all tree reconstruction methods.

Note however that the amount of data is usually limited, and the models used to estimate the distances are usually only rough approximations of the reality. As a consequence, the estimated distances  $\delta$  will somehow deviate from the distances  $\mathbf{d}^{\mathcal{T}}$  in the correct tree  $\mathcal{T}$  and the accuracy of tree reconstruction will depend on its robustness to such deviations. Define the  $L_{\infty}$  difference between  $\delta$  and  $\mathbf{d}^{\mathcal{T}}$  as  $\|\delta - \mathbf{d}^{\mathcal{T}}\|_{\infty} = \max_{i,j} |\delta_{ij} - d_{ij}^{\mathcal{T}}|$ , where  $i$  and  $j$  denote two taxa; then, a measure of this robustness, proposed by Atteson [1], is the safety radius defined in the following way, where  $\ell_{\min}^{\mathcal{T}}$  denotes the length of the shortest branch in  $\mathcal{T}$ :

**Definition 1.** A distance-based tree reconstruction method has *safety radius*  $\alpha$  [over  $n$  taxa] if, for every weighted bifurcating tree  $\mathcal{T}$  [over a set of size  $n$ ] and any distances  $\delta$  such that

$$\|\delta - \mathbf{d}^{\mathcal{T}}\|_{\infty} < \alpha \ell_{\min}^{\mathcal{T}},$$

the method reconstructs the topology of  $\mathcal{T}$ .

Note the distinction between the safety radius over a given number of taxa, and safety radius *tout-court*.

An important remark is that there are distances  $\delta$  which lie at  $\frac{1}{2}\ell_{\min}^{\mathcal{T}}$  not only from  $\mathbf{d}^{\mathcal{T}}$ , but also from the distances  $\mathbf{d}^{\mathcal{T}'}$  in a tree  $\mathcal{T}'$  with a different topology from that of  $\mathcal{T}$  [1]. This implies that: (1) robustness must be measured relative to the length of the shortest branch in  $\mathcal{T}$ , as no maximum value for the difference between  $\delta$  and  $\mathbf{d}^{\mathcal{T}}$  can guarantee a correct tree reconstruction

if nothing is assumed regarding  $\ell_{\min}^T$ ; (2) no method reconstructing a unique tree can have a safety radius greater than  $\frac{1}{2}$ .

Atteson [1] also proved that a number of agglomerative algorithms, including NJ, have optimal safety radius  $\frac{1}{2}$ . A related result was recently shown by Bordewich et al. [3], who proved that another heuristic aimed at minimizing BME (based on *subtree pruning and regrafting*) has at least radius  $\frac{1}{3}$ .

Note that these results are about *heuristics* for minimum evolution, not the criterion itself. Here, we identify the safety radius of a tree reconstruction criterion as the safety radius of any algorithm returning the optimal tree with respect to that criterion. Although there is no version of minimum evolution for which we have a practical algorithm for this task, it would still be interesting to know the safety radius that such an algorithm would have, as this may give some insight about the effectiveness of the heuristic algorithms used in practice. (We note that, at least in the case of BME, this optimisation problem is NP-hard: manuscript in preparation.)

What is then known about the safety radii of minimum evolution criteria? For BME, Bordewich et al. [3] have extended their result and proved that also the BME criterion has radius at least  $\frac{1}{3}$ . As for OLS+ME, Willson [28] proved that the safety radius over  $n$  taxa is limited above by a function tending to  $\frac{1}{4}$ , as  $n$  tends to infinity. Another relevant result was obtained by Gascuel and McKenzie [13] for the case where both the estimation of branch lengths and the selection of a tree structure are based on a least squares criterion. They proved that when this criterion is constrained to produce ultrametric trees (i.e., rooted trees in which all leaves have the same distance from the root), its radius tends to 0 as the the number of taxa increases. This result has relevance for the practice of hierarchical clustering, where the ultrametric constraint is usually required, but it does not immediately imply anything regarding the radius of OLS+ME, of relevance in phylogenetics.

The present paper concludes the debate regarding the radii of BME and OLS+ME. We prove that (as conjectured [3]) BME has radius  $\frac{1}{2}$  and that the result by Willson on the relative sensitiveness of OLS+ME can be strengthened: its safety radius actually tends to 0, as  $n$  tends to infinity. There is then a large gap between the safety radius of BME (the maximum possible) and that of OLS+ME (the minimum possible), which is consistent with the greater reconstruction accuracy of BME observed in several simulation studies [22, 16, 11, 8].

The paper is organised as follows: after some necessary notational preliminaries (Sec. 2), we proceed to prove the stated results about BME (Sec. 3) and OLS+ME (Sec. 4). In section 5 we discuss the implications of these results on the heuristics that are used in practice for BME and OLS+ME. Finally, in the Appendix, we discuss the possibility of extending our positive results on the robustness of BME to other ME methods: in particular, we prove that such an extension is not possible for an obvious candidate class generalizing BME.

## 2. PRELIMINARIES

A *phylogenetic tree* over a set  $X$  is a tree (in the graph-theoretic sense) whose leaves are bijectively labeled by the elements of  $X$ , called *taxa*. Here, for convenience of notation, we assume

$X = \{1, 2, \dots, n\}$ . A *weighted phylogenetic tree* is a phylogenetic tree whose branches (edges) are assigned *lengths* (usually non-negative) typically representing the amount of evolution that has occurred between their endpoints. A phylogenetic tree that is not weighted is also called a *topology* and we say that a weighted phylogenetic tree  $\mathcal{T}$  has topology  $T$  if  $\mathcal{T}$  can be obtained from  $T$  by assigning lengths to the branches of  $T$ . Note that, in this paper, calligraphic letters (such as  $\mathcal{T}$  and  $\mathcal{W}$ ) denote weighted phylogenetic trees and their corresponding italic symbols (such as  $T$  and  $W$ ) denote their topologies. The length of branch  $e$  in a tree  $\mathcal{T}$  is denoted by  $\ell_e^{\mathcal{T}}$ . The length of the shortest branch in  $\mathcal{T}$  is denoted by  $\ell_{\min}^{\mathcal{T}}$ .

A phylogenetic tree is *rooted* if one of its nodes is designated as its *root*, representing the ancestor of all the taxa in the tree. Except in one occasion (Lemma 8), we assume that the trees we deal with do not have a root. A phylogenetic tree is said to be *bifurcating* if every internal (i.e., non-leaf) node has degree three, with the exception of the root, if there is one, which is required to have degree two. For simplicity, in this paper we only deal with bifurcating trees.

Every branch  $e$  of a phylogenetic tree over  $X$  induces a bipartition of  $X$  consisting of the two sets of taxa in the two connected components obtained by deleting  $e$ . The bipartitions  $\{A, B\}$  induced in this way by the branches of a phylogenetic tree are called its *splits* and are denoted by  $A|B$ . A *clade* of a phylogenetic tree over  $X$  is a subset  $A \subseteq X$  such that  $A|(X \setminus A)$  is one of the tree's splits. A classic result in phylogenetics is that the topology of a phylogenetic tree is determined, up to isomorphism, by the set of its splits [5], or equivalently by the set of its clades. For this reason, we here identify a topology with the set of its splits, which allows us to write propositions such as  $A|B \in T$  and  $T \neq W$ . We also make no distinction between branches of a tree and the splits they induce, which justifies the use of expressions such as  $\ell_{A|B}^{\mathcal{T}}$  and  $e \in T$  (where  $e$  is a branch).

Given a weighted phylogenetic tree  $\mathcal{T}$ , the *distance in  $\mathcal{T}$*  between taxa  $i$  and  $j$ , denoted  $d_{ij}^{\mathcal{T}}$ , is the sum of the branch lengths in the path between  $i$  and  $j$  in  $\mathcal{T}$ . The matrix containing the distances in  $\mathcal{T}$  between each pair of taxa is denoted by  $\mathbf{d}^{\mathcal{T}} = (d_{ij}^{\mathcal{T}})$ . Distance methods for phylogenetic reconstruction are based upon the assumption that a *distance matrix*  $\boldsymbol{\delta} = (\delta_{ij})$ , somewhat approximating  $\mathbf{d}^{\mathcal{T}}$  in an unknown tree, is given in input. Here, the only assumptions on  $\boldsymbol{\delta}$  will be that, for every  $i$  and  $j$  in  $\{1, 2, \dots, n\}$ ,  $\delta_{ii} = 0$  and  $\delta_{ij} = \delta_{ji} \geq 0$ . (Although the  $\delta_{ij}$  are not, strictly speaking, necessarily distances, the use of this term is standard in phylogenetics.) In the practice of phylogenetic tree reconstruction, the distances  $\boldsymbol{\delta}$  are obtained (from molecular or morphological data) with various techniques aiming to estimate the distances  $\mathbf{d}^{\mathcal{T}}$  in the “true” evolutionary tree  $\mathcal{T}$  that generated the taxa of interest. The aim of distance methods is to reconstruct this unknown “true” tree.

Minimum evolution (ME) is a class of distance methods based upon the two following logical components:

- (1) establish how to assign branch lengths to any given topology  $T$  so that the distances  $d_{ij}^{\hat{\mathcal{T}}}$  in the resulting weighted tree  $\hat{\mathcal{T}}$  “fit” the distances  $\delta_{ij}$  given in input;

- (2) “look for” the bifurcating topology  $T$  that results in the weighted tree  $\hat{T}$  of minimum “length”.

Clearly these two steps are loosely defined. Different ways of defining the terms in quotes lead to different versions of ME.

Although other approaches are possible [15, 25, 12] (differing in the way they deal with negative branch lengths), we here define the *length* of a weighted tree as the sum of all its branches’ lengths [20]; in symbols:

$$\mathcal{L}(T) = \sum_{e \in T} \ell_e^T.$$

Here we study two versions of ME that differ for the chosen approach for the first component, branch length estimation. Possibly the simplest method for this step, called OLS (*ordinary least squares*), is to look for the branch length assignment that minimises the sum  $\sum_{ij} (\delta_{ij} - d_{ij}^{\hat{T}})^2$ . In this case, the optimal branch lengths are linear functions of  $\delta$  (see, e.g., the book chapter by Desper and Gascuel [10]). Therefore, also the total length of the fitted tree  $\hat{T}$  is a linear function of  $\delta$ , which we denote by  $S^T(\delta)$ . We call OLS+ME the version of ME using  $S^T(\delta)$  to estimate the length of  $T$ . Formally, OLS+ME aims to reconstruct  $T^* = \operatorname{argmin}_T S^T(\delta)$ , where  $T$  ranges over all bifurcating topologies over the input taxa.

The second method for branch length estimation that we study here is the *balanced* method. The branch length estimates, in this case, are optimal with respect to a weighted version of the least squares approach sketched above. We refer the reader to other sources for a more detailed introduction [19, 9, 24, 14]. For our purposes it suffices to know that the length of the fitted tree of bifurcating topology  $T$  is given by a simple, again linear, function of  $\delta$  (due to Pauplin [19]):

$$(2.1) \quad B^T(\delta) := \sum_{i < j} 2^{1-t_{ij}} \delta_{ij},$$

where  $t_{ij}$  indicates the topological distance between  $i$  and  $j$ , that is, the number of branches in the path between  $i$  and  $j$  in  $T$ . We call BME the version of ME using  $B^T(\delta)$  to estimate the length of  $T$ . Formally, BME aims to reconstruct  $T^* = \operatorname{argmin}_T B^T(\delta)$ , where  $T$  ranges again over all bifurcating topologies over the input taxa.

**2.1. Characterization of general linear formulae for estimating tree length.** The fact that  $S^T(\delta)$  and  $B^T(\delta)$  are both linear functions of  $\delta$  is not surprising, as it can be shown that optimizing branch lengths with respect to a large class of generalisations of the least squares approach always results in linear functions of the input distances [4]. Furthermore,  $S^T(\delta)$  and  $B^T(\delta)$  have a number of properties that make them suitable for use in combination with minimum evolution [27]. We briefly summarize these properties here, as some of the statements in the rest of this paper are best expressed in terms of the general class of functions satisfying these properties. We start by defining a generalisation of the notion of split metric: [2]

**Definition 2.** Let  $A$  and  $B$  be disjoint subsets of  $\{1, 2, \dots, n\}$ ;  $\sigma^{A|B}$  is defined by

$$\sigma_{ij}^{A|B} = \begin{cases} 1 & \text{if one between } i \text{ and } j \text{ belongs to } A \text{ and the other to } B, \\ 0 & \text{otherwise.} \end{cases}$$

Note that if  $\mathcal{T}$  is a weighted tree with topology  $T$ , then

$$(2.2) \quad \mathbf{d}^T = \sum_{A|B \in \mathcal{T}} \ell_{A|B}^T \sigma^{A|B}.$$

This follows from the fact that  $\sigma_{ij}^{A|B} = 1$  if and only if the branch corresponding to  $A|B$  is on the path joining  $i$  to  $j$  in  $T$ .

The most basic requirement of any linear function  $L^T(\boldsymbol{\delta})$  for estimating the length of a tree  $T$  is that it should give its correct length given perfect data. As remarked below, the functions satisfying this requirement are precisely those in the following class [27].

**Definition 3.** Let  $T$  be a topology over  $\{1, 2, \dots, n\}$ . We denote by  $\mathcal{U}(T)$  the class of linear functions  $L^T$  of any  $n \times n$  distance matrix, such that  $L^T(\sigma^{A|B}) = 1$  for every split  $A|B$  of  $T$ .

*Remark.* By using the decomposition in (2.2) and the linearity of  $L^T$  it is easy to see that the functions in  $\mathcal{U}(T)$  are precisely the linear functions such that  $L^T(\mathbf{d}^T) = \mathcal{L}(T)$ , for any weighted tree  $T$  with topology  $T$ .

Note that if we write  $L^T(\boldsymbol{\delta}) = \sum_{i < j} c_{ij}^T \delta_{ij}$  then the quantity above can simply be expressed as the sum of the coefficients for pairs of taxa in  $A \times B$ :

$$L^T(\sigma^{A|B}) = \sum_{\substack{i \in A \\ j \in B}} c_{ij}^T,$$

where  $c_{ij}^T$  with  $i > j$  is the same as  $c_{ji}^T$ .

The following definition [27] further restricts the set of linear functions suitable to be used in combination with minimum evolution.

**Definition 4.** Let  $T$  be a topology over  $\{1, 2, \dots, n\}$ . We denote by  $\mathcal{U}^+(T)$  the class of linear functions  $L^T \in \mathcal{U}(T)$  such that  $L^T(\sigma^{A|B}) > 1$  for all bipartitions  $A|B$  of  $\{1, 2, \dots, n\}$  that are not splits of  $T$ .

*Remark.* Imagine that for every bifurcating topology  $T$  we have a function from  $\mathcal{U}^+(T)$  that we use to estimate its length. Then, when the input distances  $\mathbf{d}^T$  coincide with those in a weighted bifurcating tree  $\mathcal{T}$  with positive branch lengths, minimum evolution correctly identifies the topology of  $\mathcal{T}$  as the unique optimal one. This is because every “wrong” bifurcating topology  $W$  (i.e., not coinciding with the topology of  $\mathcal{T}$ ) contains at least one split not belonging to the correct topology, and therefore is such that

$$L^W(\mathbf{d}^T) = \sum_{A|B \in \mathcal{T}} \ell_{A|B}^T L^W(\sigma^{A|B}) > \mathcal{L}(T),$$

whereas, as noted before,  $L^T(\mathbf{d}^T) = \mathcal{L}(T)$  when  $T$  is the topology of  $\mathcal{T}$ .

Minimum evolution using formulae from  $\mathcal{U}^+(T)$  is then consistent. As a matter of fact, the proofs of consistency for OLS+ME [21] and BME [9] were essentially equivalent to proving that  $S^T$  and  $B^T$  belong to  $\mathcal{U}^+(T)$ .

### 3. BME HAS SAFETY RADIUS $\frac{1}{2}$

In this section we prove that BME has a safety radius of  $\frac{1}{2}$ , the best achievable by any distance method. The proof proceeds as follows: first (Sec. 3.1, Theorem 5), we derive a sufficient condition on the linear functions  $L^T$  guaranteeing that ME has safety radius (at least)  $\alpha$ ; second (Sec. 3.2, Corollary 9), we show that this condition is satisfied by  $B^T$  with  $\alpha = \frac{1}{2}$ .

**3.1. A general condition for ME to have safety radius  $\alpha$ .** We here show a condition on the formulae used by ME, providing a minimum guaranteed radius. Although here we only prove that this is a sufficient condition, it can be shown that this is in fact also a necessary condition for ME to have at least safety radius  $\alpha$ . This theorem has similarities with results presented by Willson [28], which however were only stated for the OLS tree length estimates  $S^T(\delta)$ .

**Theorem 5.** *Assume that for each bifurcating topology  $T$  we use a linear function  $L^T \in \mathcal{U}^+(T)$  to estimate its length, where  $L^T(\delta) = \sum_{i < j} c_{ij}^T \delta_{ij}$ . Suppose that for any bifurcating topologies  $T$  and  $W$  over the same set of taxa the following holds:*

$$(3.1) \quad \sum_{i < j} (c_{ij}^W (t_{ij} - w_{ij}) - \alpha |c_{ij}^W - c_{ij}^T|) \geq 0,$$

where  $t_{ij}$  and  $w_{ij}$  denote the number of branches in the path between  $i$  and  $j$  in  $T$  and  $W$ , respectively. Then minimum evolution has safety radius (at least)  $\alpha$ .

The proof of Theorem 5 relies on the following two lemmas.

**Lemma 6.** *Let  $L(\delta) = \sum_{i < j} c_{ij} \delta_{ij}$  be a linear function and  $\delta$  and  $\delta'$  two distance matrices whose components differ by less than  $\epsilon > 0$ , that is,  $\|\delta - \delta'\|_\infty < \epsilon$ . Then,*

$$(3.2) \quad |L(\delta) - L(\delta')| < \epsilon \sum_{i < j} |c_{ij}|.$$

*Proof.*

$$|L(\delta) - L(\delta')| = \left| \sum_{i < j} c_{ij} (\delta_{ij} - \delta'_{ij}) \right| \leq \sum_{i < j} |c_{ij}| |\delta_{ij} - \delta'_{ij}| < \epsilon \sum_{i < j} |c_{ij}|.$$

□

**Lemma 7.** *Let  $\mathcal{T}$  be a weighted tree with non-negative branch lengths and a bifurcating topology  $T$  and let  $W$  be another bifurcating topology over the same set of taxa  $\{1, 2, \dots, n\}$ . Assuming  $L^W \in \mathcal{U}^+(W)$ , then*

$$L^W(\mathbf{d}^T) - \mathcal{L}(\mathcal{T}) \geq \ell_{\min}^T L^W(\mathbf{t} - \mathbf{w}),$$

where  $\mathbf{t} = (t_{ij})$  and  $\mathbf{w} = (w_{ij})$  contain the number of branches in the paths between  $i$  and  $j$  in  $T$  and  $W$ , respectively.

*Proof.* Using the decomposition of  $\mathbf{d}^T$  given by (2.2), we obtain

$$L^W(\mathbf{d}^T) - \mathcal{L}(T) = \sum_{A|B \in T} \ell_{A|B}^T \left( L^W(\boldsymbol{\sigma}^{A|B}) - 1 \right).$$

Since  $L^W \in \mathcal{U}^+(W)$ , each  $L^W(\boldsymbol{\sigma}^{A|B}) - 1$  term above is non-negative. This, together with  $\ell_{A|B}^T \geq \ell_{\min}^T \geq 0$ , implies

$$\begin{aligned} L^W(\mathbf{d}^T) - \mathcal{L}(T) &\geq \ell_{\min}^T \sum_{A|B \in T} \left( L^W(\boldsymbol{\sigma}^{A|B}) - 1 \right) \\ &= \ell_{\min}^T \left( L^W \left( \sum_{A|B \in T} \boldsymbol{\sigma}^{A|B} \right) - (2n - 3) \right). \end{aligned}$$

Now note that  $\mathbf{t} = \sum_{A|B \in T} \boldsymbol{\sigma}^{A|B}$  and that  $L^W(\mathbf{w}) = 2n - 3$ . Then

$$\begin{aligned} L^W(\mathbf{d}^T) - \mathcal{L}(T) &\geq \ell_{\min}^T (L^W(\mathbf{t}) - L^W(\mathbf{w})) \\ &= \ell_{\min}^T L^W(\mathbf{t} - \mathbf{w}). \end{aligned}$$

□

*Proof of Theorem 5.* We wish to prove that for every weighted bifurcating tree  $\mathcal{T}$  with topology  $T$ , if  $\|\boldsymbol{\delta} - \mathbf{d}^T\|_\infty < \alpha \ell_{\min}^T$ , then  $\Delta L(\boldsymbol{\delta}) = L^W(\boldsymbol{\delta}) - L^T(\boldsymbol{\delta}) > 0$ , for every bifurcating topology  $W \neq T$ .

First, application of Lemma 6 with  $L = \Delta L$ ,  $\boldsymbol{\delta}' = \mathbf{d}^T$  and  $\epsilon = \alpha \ell_{\min}^T$ , shows that

$$(3.3) \quad \Delta L(\boldsymbol{\delta}) > \Delta L(\mathbf{d}^T) - \alpha \ell_{\min}^T \sum_{i < j} |c_{ij}^W - c_{ij}^T|.$$

Using the fact that  $L^T \in \mathcal{U}^+(T)$  correctly estimates the length of  $\mathcal{T}$  given perfect data,  $\Delta L(\mathbf{d}^T) = L^W(\mathbf{d}^T) - L^T(\mathbf{d}^T) = L^W(\mathbf{d}^T) - \mathcal{L}(T)$ . Lemma 7 shows that a lower bound for this quantity is given by

$$\ell_{\min}^T L^W(\mathbf{t} - \mathbf{w}) = \ell_{\min}^T \sum_{i < j} c_{ij}^W (t_{ij} - w_{ij}).$$

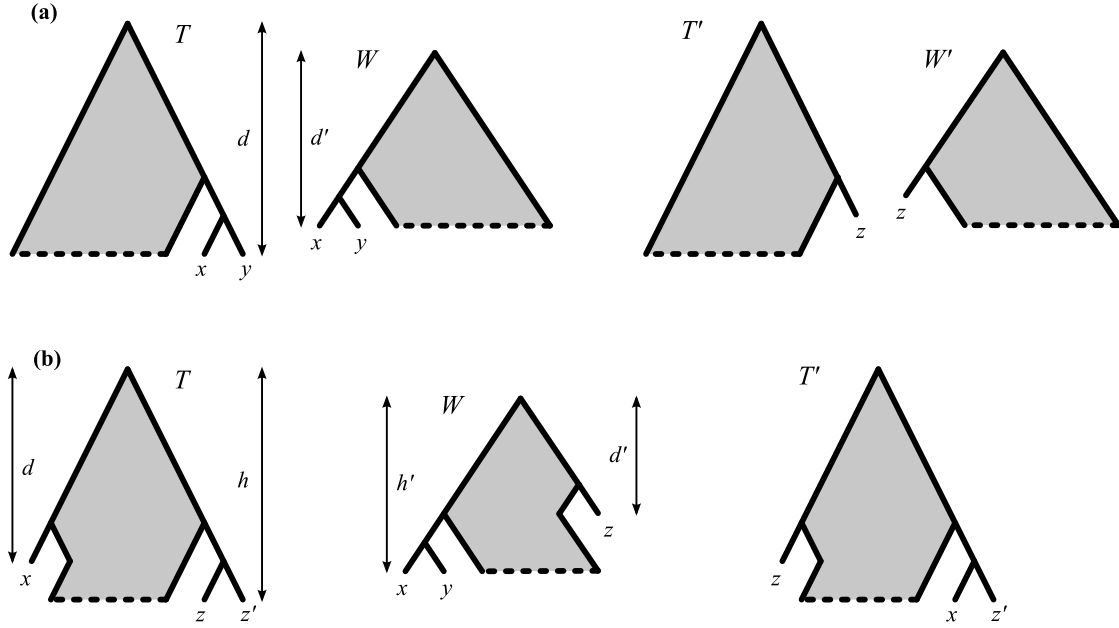
Plugging this expression into (3.3), we obtain

$$\Delta L(\boldsymbol{\delta}) > \ell_{\min}^T \sum_{i < j} (c_{ij}^W (t_{ij} - w_{ij}) - \alpha |c_{ij}^W - c_{ij}^T|).$$

It is then clear that if the sum in the right-hand side above is non-negative for every choice of  $T$  and  $W$ , then ME has safety radius  $\alpha$ . □

**3.2. The length estimates used by BME satisfy Theorem 5.** We now prove that the inequality in Theorem 5 is satisfied for  $L^T = B^T$  and  $\alpha = \frac{1}{2}$ . This will follow as a simple corollary (Corollary 9) of the next lemma, where we need some additional definitions: the *depth* of a taxon  $i$  in a rooted topology  $T$  is the number of branches in the path, in  $T$ , from the root to  $i$ ; also, a *cherry* is a clade of size 2.

FIGURE 3.1. Topologies in the proof of Lemma 8.



**Lemma 8.** Let  $T$  and  $W$  be bifurcating rooted topologies over  $\{1, 2, \dots, n\}$ . Let  $d_i$  and  $d'_i$  denote the depths of  $i \in \{1, 2, \dots, n\}$  in  $T$  and  $W$ , respectively. Then,

$$\sum_{i=1}^n 2^{-d'_i} \left( d_i - d'_i - \frac{1}{2} |2^{d'_i - d_i} - 1| \right) \geq 0.$$

*Proof.* By induction on the number of taxa  $n$ . Denote by  $f(T, W)$  the sum in the lemma's statement. If  $n = 1$  or  $2$ , then  $T$  and  $W$  must be the same topology and therefore  $f(T, W) = 0$ . For the case  $n \geq 3$ , we consider two scenarios: either  $T$  and  $W$  have a common cherry, or they do not.

If  $T$  and  $W$  have a common cherry, say  $\{x, y\}$ , then define  $T'$  and  $W'$  as the trees that are obtained from  $T$  and  $W$ , respectively, by removing the cherry  $\{x, y\}$  and replacing it with a new taxon  $z$  (see fig. 3.1(a)). Since  $T'$  and  $W'$  have  $n - 1$  taxa, we can apply the induction hypothesis and have that  $f(T', W') \geq 0$ . The difference  $f(T, W) - f(T', W')$  is easy to calculate, as the two sums only differ for the terms corresponding to taxa  $x, y$  and  $z$ . Calling  $d$  and  $d'$  the depths of  $x$  (and therefore  $y$ ) in  $T$  and  $W$ , respectively, and noting that therefore the depths of  $z$  in  $T'$  and  $W'$  equal  $d - 1$  and  $d' - 1$  (see fig. 3.1(a)), we have that:

$$\begin{aligned} f(T, W) - f(T', W') &= \\ 2 \cdot 2^{-d} \left( d - d' - \frac{1}{2} |2^{d' - d} - 1| \right) - 2^{-(d' - 1)} \left( d - d' - \frac{1}{2} |2^{d' - d} - 1| \right) &= 0 \end{aligned}$$

Therefore  $f(T, W) = f(T', W') \geq 0$ .

If  $T$  and  $W$  do not have a common cherry, let  $\{x, y\}$  be a cherry with maximum depth in  $W$ . We will show that if we swap taxa in  $T$  so as to form  $\{x, y\}$  in the resulting tree  $T''$ , this tree is such that  $f(T, W) \geq f(T'', W)$ . Since  $T''$  and  $W$  do have a common cherry, for the arguments above, we have that  $f(T'', W)$  – and therefore  $f(T, W)$  – is non-negative.

Let  $\{z, z'\}$  be a cherry with maximum depth in  $T$ . Because  $T$  and  $W$  do not have common cherries, at least one of  $z$  and  $z'$  must not be in  $\{x, y\}$ . Without loss of generality, assume  $z \notin \{x, y\}$ . Swap  $x$  and  $z$  in  $T$  and call the resulting tree  $T'$  (see fig. 3.1(b)). Let  $d$  and  $h$  be the depths in  $T$  of  $x$  and  $z$ , respectively, and note that this implies that, in  $T'$ ,  $x$  and  $z$  have depths  $h$  and  $d$ , respectively. Also, let  $d'$  and  $h'$  be the depths in  $W$  of  $z$  and  $x$ , respectively.

We now show that  $f(T, W) \geq f(T', W)$ . The two sums only differ for the terms corresponding to  $x$  and  $z$ :

$$\begin{aligned} f(T, W) - f(T', W) &= 2^{-h'} \left( d - h' - \frac{1}{2} \left| 2^{h'-d} - 1 \right| \right) - 2^{-h'} \left( h - h' - \frac{1}{2} \left| 2^{h'-h} - 1 \right| \right) \\ &\quad + 2^{-d'} \left( h - d' - \frac{1}{2} \left| 2^{d'-h} - 1 \right| \right) - 2^{-d'} \left( d - d' - \frac{1}{2} \left| 2^{d'-d} - 1 \right| \right) \\ &= \left( 2^{-d'} - 2^{-h'} \right) (h - d) + \frac{1}{2} \left( \left| 2^{-h} - 2^{-h'} \right| - \left| 2^{-d} - 2^{-h'} \right| + \left| 2^{-d} - 2^{-d'} \right| - \left| 2^{-h} - 2^{-d'} \right| \right). \end{aligned}$$

Using the fact that  $|x - y| = x + y - 2 \min\{x, y\}$ , the expression above simplifies into

$$\left( 2^{-d'} - 2^{-h'} \right) (h - d) - \min\{2^{-h}, 2^{-h'}\} + \min\{2^{-d}, 2^{-h'}\} - \min\{2^{-d}, 2^{-d'}\} + \min\{2^{-h}, 2^{-d'}\}.$$

Note that if  $h = d$  then  $f(T, W) - f(T', W) = 0$ . We then assume  $h \geq d + 1$  (remember that  $h$  and  $d$  are integers, as they represent depths in  $T$ ).

Considering the three cases  $d \leq d' \leq h'$ ,  $d' \leq d \leq h'$  and  $d' \leq h' \leq d$ , it is easy to see that

$$\min\{2^{-d}, 2^{-h'}\} - \min\{2^{-d}, 2^{-d'}\} \geq 2^{-h'} - 2^{-d'}.$$

Similarly, considering the three possible positions of  $h$  relative to  $d'$  and  $h'$ , it is easy to see that

$$\min\{2^{-h}, 2^{-d'}\} - \min\{2^{-h}, 2^{-h'}\} \geq 0.$$

Therefore

$$\begin{aligned} f(T, W) - f(T', W) &\geq \left( 2^{-d'} - 2^{-h'} \right) (h - d) + 2^{-h'} - 2^{-d'} \\ &= \left( 2^{-d'} - 2^{-h'} \right) (h - d - 1) \geq 0. \end{aligned}$$

This completes the proof that  $f(T, W) \geq f(T', W)$ . Now  $T'$  has  $\{x, z'\}$  as a cherry. If  $z' = y$ , then define  $T'' = T'$ ; otherwise let  $T''$  be obtained by swapping  $z'$  with  $y$  in  $T'$ . In any case, by the same arguments as above,  $f(T', W) \geq f(T'', W)$ .

Since  $T''$  has a cherry in common with  $W$ , we have that  $f(T'', W) \geq 0$ . But then,

$$f(T, W) \geq f(T', W) \geq f(T'', W) \geq 0.$$

□

*Remark.* The inequality of Lemma 8 holds more generally for any positive sequences  $(d_i)$  and  $(d'_i)$  for which  $\sum_i 2^{-d_i} = \sum_i 2^{-d'_i} = 1$ , because such sequences coincide precisely with the taxon depths in two suitably defined bifurcating rooted trees (not shown).

**Corollary 9.** *Let  $T$  and  $W$  be bifurcating topologies over  $\{1, 2, \dots, n\}$ . Let  $t_{ij}$  and  $w_{ij}$  denote the number of branches in the paths between  $i$  and  $j$  in  $T$  and  $W$ , respectively. Then,*

$$\sum_{i < j} 2^{1-w_{ij}} \left( t_{ij} - w_{ij} - \frac{1}{2} |2^{w_{ij}-t_{ij}} - 1| \right) \geq 0.$$

*Proof.* The above sum can be re-expressed in the following way:

$$\frac{1}{2} \sum_{j=1}^n \sum_{\substack{i=1 \\ i \neq j}}^n 2^{1-w_{ij}} \left( t_{ij} - w_{ij} - \frac{1}{2} |2^{w_{ij}-t_{ij}} - 1| \right).$$

Now define  $T_j$  and  $W_j$  as the rooted topologies that are obtained from  $T$  and  $W$  by rooting in taxon  $j$  and removing the branch at the root. Then,

$$\sum_{\substack{i=1 \\ i \neq j}}^n 2^{1-w_{ij}} \left( t_{ij} - w_{ij} - \frac{1}{2} |2^{w_{ij}-t_{ij}} - 1| \right) = \sum_{\substack{i=1 \\ i \neq j}}^n 2^{-d'_i} \left( d_i - d'_i - \frac{1}{2} |2^{d'_i-d_i} - 1| \right),$$

where  $d_i = t_{ij} - 1$  and  $d'_i = w_{ij} - 1$  are the depths of  $i$  in  $T_j$  and  $W_j$ , respectively. Because of Lemma 8, this sum is non-negative, and therefore the whole sum in the statement — which is a sum of sums with this form — is also non-negative.  $\square$

**3.3. Putting it all together.** We are now ready to complete the proof of the main result of this section.

**Theorem 10.** *BME has safety radius  $\frac{1}{2}$ .*

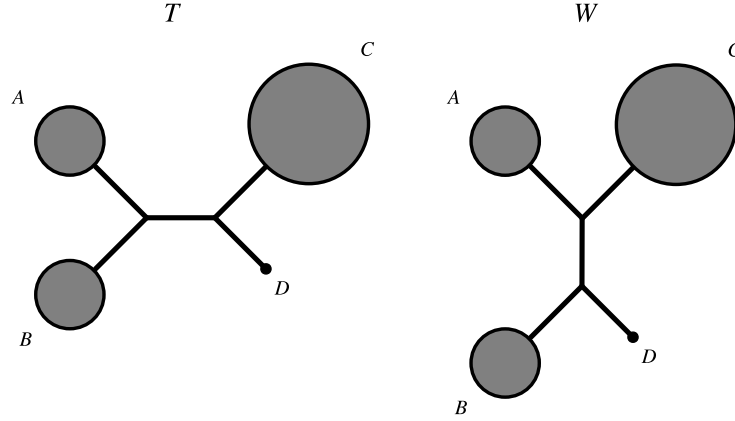
*Proof.* Because  $B^T \in \mathcal{U}^+(T)$  [9], Theorem 5 provides a sufficient condition for BME to have radius  $\frac{1}{2}$ . It is easy to verify that this condition precisely coincides with that proven in Corollary 9.  $\square$

Given the generality of our arguments in section 3.1 — holding for ME using generic formulae in  $\mathcal{U}^+(T)$  — one can wonder whether our proof can be extended to other formulae apart from the  $B^T$  used by BME. However, the next section shows that OLS+ME, which uses  $S^T \in \mathcal{U}^+(T)$  to estimate tree length, has safety radius equal to 0. In order to exclude even more strongly the possibility to extend our arguments beyond BME, in the Appendix we show that even if we only use members of  $\mathcal{U}^+(T)$  that share with  $B^T$  the property of being linear functions with positive coefficients, ME can have an arbitrarily small safety radius.

#### 4. OLS+ME HAS SAFETY RADIUS 0

In the example that follows we construct, for any fixed  $n$  ( $n \geq 4$ ), two bifurcating topologies  $T$  and  $W$  over the same set of  $n$  taxa, a weighted tree  $\mathcal{T}$  with topology  $T$ , and distances  $\delta$  such

FIGURE 4.1. Topologies referred to in Example 11.



that

$$\|\delta - \mathbf{d}^T\|_\infty = \frac{3}{\lfloor \sqrt{n} \rfloor + 1} \ell_{\min}^T \quad \text{and} \quad S^W(\delta) < S^T(\delta).$$

Clearly, this implies that the safety radius of OLS+ME over  $n$  taxa cannot be greater than  $3/(\lfloor \sqrt{n} \rfloor + 1)$ . Because this function converges to 0, then also the safety radius of OLS+ME over  $n$  taxa tends to 0 as the number of taxa grows. This result completes the work by Willson [28], who previously showed that this safety radius is limited above by a function tending to  $\frac{1}{4}$ , but left open the question of how tight this bound was.

The following example makes use of a standard notation that we now introduce. Let  $A$  and  $B$  be two disjoint clades. Given the distances in  $\delta$ , the average distance between  $A$  and  $B$  is defined by:

$$\delta_{AB} = \frac{1}{|A||B|} \sum_{\substack{i \in A \\ j \in B}} \delta_{ij}.$$

**Example 11.** Let  $\mathcal{T}$  be a weighted bifurcating tree whose branch of minimum length separates clades  $A$  and  $B$  on one side from  $C$  and  $D$  on the other side (as in Fig. 4.1, left). Let  $|A| = |B| = \lfloor \sqrt{n} \rfloor$ ,  $|C| = n - 2\lfloor \sqrt{n} \rfloor - 1$  and  $|D| = 1$ . Let  $T$  be the topology of  $\mathcal{T}$ , and let  $W$  be the topology that is obtained by swapping the positions of  $B$  and  $C$  in  $T$  (as in Fig. 4.1, right).

Now define  $\delta$  in the following way:

$$\begin{aligned} \delta_{ij} &= d_{ij}^T + \epsilon & \text{if } (i, j) \in (A \times B) \cup (C \times D) \cup (A \times D) \cup (B \times C), \\ \delta_{ij} &= d_{ij}^T - \epsilon & \text{if } (i, j) \in (A \times C) \cup (B \times D), \\ \delta_{ij} &= d_{ij}^T & \text{if } (i, j) \in A^2 \cup B^2 \cup C^2 \cup D^2. \end{aligned}$$

Clearly,  $\|\delta - \mathbf{d}^T\|_\infty = \epsilon$ . If we express the average distances between the clades  $A$ ,  $B$ ,  $C$  and  $D$ , we then have:

$$\begin{aligned} \delta_{AB} &= d_{AB}^T + \epsilon, & \delta_{CD} &= d_{CD}^T + \epsilon, \\ \delta_{AD} &= d_{AD}^T + \epsilon, & \delta_{BC} &= d_{BC}^T + \epsilon, \\ \delta_{AC} &= d_{AC}^T - \epsilon, & \delta_{BD} &= d_{BD}^T - \epsilon. \end{aligned}$$

Desper and Gascuel [8] have shown how to calculate the difference  $\Delta S(\boldsymbol{\delta}) = S^W(\boldsymbol{\delta}) - S^T(\boldsymbol{\delta})$ , for a generic  $\boldsymbol{\delta}$ , when  $W$  is obtained from  $T$  with a transformation such as that in Figure 4.1 (a *nearest neighbor interchange*, *NNI*). Applying their formula (equation (9) in their paper [8]), we have that:

$$(4.1) \quad \Delta S(\boldsymbol{\delta}) = \frac{1}{2} [(1 - \lambda)(\delta_{AC} + \delta_{BD}) - (1 - \lambda')(\delta_{AB} + \delta_{CD}) + (\lambda - \lambda')(\delta_{AD} + \delta_{BC})],$$

where

$$\lambda = \frac{|A||D| + |B||C|}{(|A| + |B|)(|C| + |D|)} = \frac{1}{2}$$

and

$$\lambda' = \frac{|A||D| + |B||C|}{(|A| + |C|)(|B| + |D|)} = \frac{n \lfloor \sqrt{n} \rfloor - 2 \lfloor \sqrt{n} \rfloor^2}{(n - \lfloor \sqrt{n} \rfloor - 1)(\lfloor \sqrt{n} \rfloor + 1)}.$$

For the calculations that follow it suffices to note that, as  $\lfloor \sqrt{n} \rfloor^2 \leq n$ ,

$$(4.2) \quad \lambda' \geq \frac{n(\lfloor \sqrt{n} \rfloor - 2)}{(n - \lfloor \sqrt{n} \rfloor - 1)(\lfloor \sqrt{n} \rfloor + 1)} > \frac{\lfloor \sqrt{n} \rfloor - 2}{\lfloor \sqrt{n} \rfloor + 1}.$$

Using the expressions for the average distances between  $A$ ,  $B$ ,  $C$  and  $D$  and (4.1), we obtain:

$$\begin{aligned} \Delta S(\boldsymbol{\delta}) &= \frac{1}{2} [(1 - \lambda)(d_{AC}^T + d_{BD}^T - 2\epsilon) - (1 - \lambda')(d_{AB}^T + d_{CD}^T + 2\epsilon) \\ &\quad + (\lambda - \lambda')(d_{AD}^T + d_{BC}^T + 2\epsilon)] \\ &= \frac{1}{2} [(1 - \lambda)(d_{AC}^T + d_{BD}^T) - (1 - \lambda')(d_{AB}^T + d_{CD}^T) + (\lambda - \lambda')(d_{AD}^T + d_{BC}^T)] \\ &\quad + [(\lambda - 1) - (1 - \lambda') + (\lambda - \lambda')] \epsilon \\ &= \Delta S(\mathbf{d}^T) - 2(1 - \lambda)\epsilon \\ &= \Delta S(\mathbf{d}^T) - \epsilon. \end{aligned}$$

We are now going to express  $\Delta S(\mathbf{d}^T)$ . We use the following relationship:

$$d_{AC}^T + d_{BD}^T = d_{AD}^T + d_{BC}^T = d_{AB}^T + d_{CD}^T + 2\ell_{\min}^T,$$

which can be proved by noting that  $d_{ac}^T + d_{bd}^T = d_{ad}^T + d_{bc}^T = d_{ab}^T + d_{cd}^T + 2\ell_{\min}^T$  holds for any  $(a, b, c, d) \in A \times B \times C \times D$ . Applying again (4.1), we then obtain:

$$\begin{aligned} \Delta S(\mathbf{d}^T) &= \frac{1}{2} [(1 - \lambda)(d_{AC}^T + d_{BD}^T) - (1 - \lambda')(d_{AB}^T + d_{CD}^T) + (\lambda - \lambda')(d_{AD}^T + d_{BC}^T)] \\ &= \frac{1}{2} [((1 - \lambda) - (1 - \lambda') + (\lambda - \lambda'))(d_{AB}^T + d_{CD}^T) + 2((1 - \lambda) + (\lambda - \lambda'))\ell_{\min}^T] \\ &= (1 - \lambda')\ell_{\min}^T. \end{aligned}$$

We therefore have:

$$\Delta S(\boldsymbol{\delta}) = (1 - \lambda')\ell_{\min}^T - \epsilon.$$

Using the bound in (4.2), we then have

$$\Delta S(\boldsymbol{\delta}) < \frac{3}{\lfloor \sqrt{n} \rfloor + 1} \ell_{\min}^T - \epsilon$$

from which it becomes clear that if  $\epsilon = \frac{3}{\lfloor \sqrt{n} \rfloor + 1} \ell_{\min}^T$ , then  $S^W(\boldsymbol{\delta}) - S^T(\boldsymbol{\delta}) < 0$ , which is what we wanted to show.

## 5. DISCUSSION

The performance of minimum evolution strongly depends on the chosen method for estimating branch lengths. Not only length estimation determines whether or not ME is consistent [21, 12, 9, 27], but also whether it is robust to noise in the input data. A measure of this robustness is provided by the safety radius studied here.

Previously, it was shown that OLS estimation causes ME to have radius at most  $\frac{1}{4}$  [28]; here we have shown that this version of ME in fact has safety radius 0 (Sec. 4). As for BME, the fact that NJ has optimal radius  $\frac{1}{2}$  [1] led some to conjecture that also BME has radius  $\frac{1}{2}$  [3]; we have proved this here for the first time (Sec. 3).

The result on the safety radius of OLS+ME (0) may at first surprise, as it is very different from the safety radius of NJ ( $\frac{1}{2}$ ), an algorithm that was originally designed to approximate OLS+ME [23]. However it has previously been noted that further optimizing for OLS+ME the trees produced by NJ actually results in a decrease – not an increase – of their accuracy [11, 8]. This observation, together with the lack of robustness of OLS+ME shown here, demonstrates that the merits of NJ lie in the fact that it is a heuristic for BME, and do not come from its relationship with OLS+ME.

In fact, two heuristics precisely aimed at optimizing OLS+ME have been proposed [8] and it is easy to see that also these heuristics have safety radius (converging to) 0. The first of these heuristics (called GME by Desper and Gascuel [8]), is based on a *greedy sequential addition* strategy: starting from the only possible topology on the first three taxa, sequentially add each of the remaining taxa by attaching it (with a new terminal branch) onto the branch that minimises the OLS length  $S^{T'}(\boldsymbol{\delta})$  of the resulting tree  $T'$ . Now imagine applying this greedy algorithm to the distances  $\boldsymbol{\delta}$  constructed in Example 11 where taxon  $n$  (the last to be added by the algorithm) coincides with the only taxon in clade  $D$ . It is then clear that, if  $\epsilon$  is set to the value derived in that example, this algorithm will not reconstruct the correct topology: even if all the taxa up to  $n - 1$  are added in the right places, because  $S^W(\boldsymbol{\delta}) < S^T(\boldsymbol{\delta})$ , taxon  $n$  would be added in the position producing the wrong tree  $W$  rather than that producing  $T$ . The other heuristic for OLS+ME (called FASTNNI [8]), consists of repeatedly swapping the positions of clades separated by three branches (i.e., applying an NNI), again greedily choosing at each step the swap that produces the largest decrease in OLS length, until a local minimum is reached. Again this algorithm does not reconstruct the correct topology when applied to Example 11: because  $W$  is only one swap away from  $T$  and has a smaller OLS length,  $T$  is not a local minimum. Since they do not reconstruct the correct tree when confronted with distances deviating by  $\epsilon$  (with  $\epsilon/\ell_{\min}^T \xrightarrow{n \rightarrow \infty} 0$ ) from the correct distances, both these heuristics have safety radius (converging to) 0.

A related question is that of the safety radius of heuristic algorithms for BME. In addition to NJ, two other such heuristics can simply be obtained by adapting to BME the two

algorithms described above for OLS+ME. A simple inductive reasoning shows that the greedy sequential addition algorithm for BME has radius  $\frac{1}{2}$  (an alternative proof of this result has also been found by R.Mihaescu and M.Bordewich; unpublished, personal communication). Suppose  $\|\delta - \mathbf{d}^{\mathcal{T}}\|_{\infty} < \frac{1}{2}\ell_{\min}^{\mathcal{T}}$ . If  $n = 4$ , then all the three possible bifurcating topologies are taken into account by this algorithm, which then returns the topology with the minimum balanced length. Because of Theorem 10, this is precisely the topology of  $\mathcal{T}$ . If  $n > 4$ , we can assume by inductive hypothesis that the greedy sequential algorithm applied up to taxon  $n - 1$  reconstructs the correct topology of  $\mathcal{T}$  restricted to the first  $n - 1$  taxa. The complete topology of  $\mathcal{T}$  is therefore one of the possible topologies that can be produced when adding the last taxon,  $n$ . Because of Theorem 10, this topology is the one with minimum balanced length and is therefore the one that gets reconstructed.

As for the NNI-based heuristic applied to BME (called BNNI [8]), the question of whether this algorithm is consistent (and therefore that of determining its radius) is still open, although empirical evidence leads us to conjecture that consistency does in fact hold.

In conclusion, there seems to be a clear difference in the robustness of methods based on OLS+ME and those based on BME. Interestingly, although the safety radius is here defined in terms of the maximum norm  $\|\delta - \mathbf{d}^{\mathcal{T}}\|_{\infty}$ , the high sensitivity of OLS+ME and the relative robustness of BME to noise in the input data are independent of the choice of a norm: for instance, whatever the definition of the difference between  $\delta$  and  $\mathbf{d}^{\mathcal{T}}$ , the example in section 4 shows that this difference can become arbitrarily small, with OLS+ME still favoring a wrong tree over the correct one.

BME and OLS+ME therefore have *qualitatively* different properties when confronted with realistic (noisy) data. This may explain the lower reconstruction accuracy of (heuristics for) OLS+ME as compared to (heuristics for) BME, a difference in performance that was observed in several simulation studies [22, 16, 11, 8].

*Acknowledgement.* We wish to thank Mike Steel for helpful suggestions about the original manuscript and Radu Mihaescu for proposing an alternative proof for Lemma 8, based on the observation that the sum in its statement can be seen as the difference between the Kullback-Leibler divergence and the total variation distance between the discrete probability distributions defined by  $(2^{-d_i})_i$  and  $(2^{-d_i})_i$ .

## REFERENCES

- [1] K. Atteson. The performance of neighbor-joining methods of phylogenetic reconstruction. *Algorithmica*, 25:251–278, 1999.
- [2] H.J. Bandelt and A.W. Dress. Split decomposition: a new and useful approach to phylogenetic analysis of distance data. *Molecular Phylogenetics and Evolution*, 1(3):242–252, 1992.
- [3] M. Bordewich, O. Gascuel, K. Huber, and V. Moulton. Consistency of topological moves based on the balanced minimum evolution principle of phylogenetic inference. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 6:110–117, 2009.
- [4] M. Bulmer. Use of the method of generalized least squares in reconstructing phylogenies from sequence data. *Molecular Biology and Evolution*, 8(6):868, 1991.

- [5] P. Buneman. The recovery of trees from measures of dissimilarity. In F. Hodson, editor, *Mathematics in the Archaeological and Historical Sciences*, pages 387–395. Edinburgh University Press, 1971.
- [6] L.L. Cavalli-Sforza and A.W.F. Edwards. Analysis of human evolution. *Proceedings 11th International Congress of Genetics*, 3:923–933, 1964.
- [7] L.L. Cavalli-Sforza and A.W.F. Edwards. Phylogenetic analysis: models and estimation procedures. *American Journal of Human Genetics*, 19(3 part 1):233–257, 1967.
- [8] R. Desper and O. Gascuel. Fast and accurate phylogeny reconstruction algorithms based on the minimum-evolution principle. *Journal of Computational Biology*, 9:687–705, 2002.
- [9] R. Desper and O. Gascuel. Theoretical foundation of the balanced minimum evolution method of phylogenetic inference and its relationship to weighted least-squares tree fitting. *Molecular Biology and Evolution*, 21:587–598, 2004.
- [10] R. Desper and O. Gascuel. The minimum evolution distance-based approach to phylogenetic inference. In O. Gascuel, editor, *Mathematics of Evolution & Phylogeny*, pages 1–32. Oxford University Press, 2005.
- [11] O. Gascuel. On the optimization principle in phylogenetic analysis and the minimum-evolution criterion. *Molecular Biology and Evolution*, 17:401–405, 2000.
- [12] O. Gascuel, D. Bryant, and F. Denis. Strengths and limitations of the minimum evolution principle. *Systematic Biology*, 50:621–627, 2001.
- [13] O. Gascuel and A. McKenzie. Performance analysis of hierarchical clustering algorithms. *Journal of classification*, 21(1):3–18, 2004.
- [14] O. Gascuel and M. Steel. Neighbor-joining revealed. *Molecular Biology and Evolution*, 23:1997–2000, 2006.
- [15] K.K. Kidd and L.A. Sgaramella-Zonta. Phylogenetic analysis: concepts and methods. *American Journal of Human Genetics*, 23:235–252, 1971.
- [16] S. Kumar. A stepwise algorithm for finding minimum evolution trees, 1996.
- [17] R. Mihaescu and L. Pachter. Combinatorics of least squares trees. *Proceedings of the National Academy of Sciences USA*, 105:13206–13211, 2008.
- [18] M. Nei and S. Kumar. *Molecular evolution and phylogenetics*. Oxford University Press, USA, 2000.
- [19] Y. Pauplin. Direct calculation of a tree length using a distance matrix. *Journal of Molecular Evolution*, 51:41–47, 2000.
- [20] A. Rzhetsky and M. Nei. A simple method for estimating and testing minimum-evolution trees. *Molecular Biology and Evolution*, 9:945–967, 1992.
- [21] A. Rzhetsky and M. Nei. Theoretical foundation of the minimum-evolution method of phylogenetic inference. *Molecular Biology and Evolution*, 10:1073–1095, 1993.
- [22] N. Saitou and T. Imanishi. Relative efficiencies of the Fitch-Margoliash, maximum-parsimony, maximum-likelihood, minimum-evolution, and neighbor-joining methods of phylogenetic tree construction in obtaining the correct tree. *Molecular Biology and Evolution*, 6:514–525, 1989.
- [23] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4:406–425, 1987.
- [24] C. Semple and M. Steel. Cyclic permutations and evolutionary trees. *Advances in Applied Mathematics*, 32:669–680, 2004.
- [25] D.L. Swofford, G.J. Olsen, P.J. Waddell, and D.M. Hillis. Phylogenetic inference. In D. Hillis, C. Moritz, and B. Mable, editors, *Molecular Systematics*, pages 407–514. Sinauer, 1996.
- [26] L.S. Vinh and A. von Haeseler. Shortest triplet clustering: reconstructing large phylogenies using representative sets. *BMC Bioinformatics*, 6:92, 2005.
- [27] S.J. Willson. Consistent formulas for estimating the total lengths of trees. *Discrete Applied Mathematics*, 148(3):214–239, 2005.
- [28] S.J. Willson. Minimum evolution using ordinary least-squares is less robust than neighbor-joining. *Bulletin of Mathematical Biology*, 67(2):261–279, 2005.

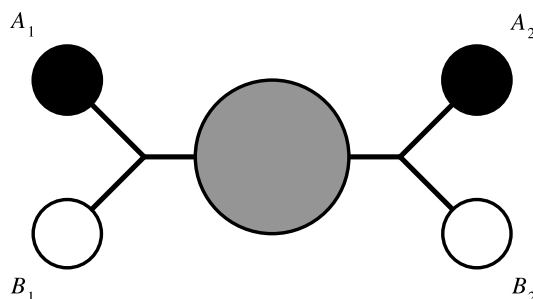
## APPENDIX

Recently, it has been shown that  $S^T$  and  $B^T$  can be seen as special cases of larger classes of length estimators [27, 17].  $\mathcal{U}^+(T)$ , defined in the Preliminaries, is one of these classes (and in fact it can be proven to contain the class by Mihaescu and Pachter [17]; result not shown). Given our results of sections 3 and 4, a reasonable question is whether there is a subclass of  $\mathcal{U}^+(T)$ , containing  $B^T$  but not  $S^T$ , such that ME in combination with the members of this subclass is robust to noisy data.

A good candidate for such a subclass is the set of all the linear functions  $L^T \in \mathcal{U}(T)$  (i.e., those giving the correct length of a tree given perfect data) having positive coefficients: that is, assuming  $L^T(\delta) = \sum_{i < j} c_{ij}^T \delta_{ij}$ , such that  $c_{ij}^T > 0$  for every  $i < j$ . It is easy to verify that  $B^T$  belongs to this set — as its coefficients are powers of two and therefore always positive — whereas  $S^T$ , in general, does not — for example some of the coefficients of  $S^T$  are negative when  $T$  is a caterpillar tree with 8 taxa.

The interest in this class comes from an elegant result which we prove now. We say that two clades  $X$  and  $Y$  of  $T$  are *separated* by a branch  $e$  if any path in  $T$  between a leaf in  $X$  and a leaf in  $Y$  contains  $e$ . Note that clades separated by at least one branch must be disjoint.

FIGURE 5.1. Illustration of the proof of Theorem 12: all the taxa in  $A_1$  and  $A_2$  also belong to  $A$ , which we denote by coloring these clades in black. All the taxa in  $B_1$  and  $B_2$  also belong to  $B$ , which we denote by coloring these clades in white. The rest of the tree is colored in grey to denote that some of its taxa are in  $A$  and some in  $B$ .



**Theorem 12.** *Assume that for each bifurcating topology  $T$  we use a linear function  $L^T \in \mathcal{U}(T)$  to estimate its length, where  $L^T(\delta) = \sum_{i < j} c_{ij}^T \delta_{ij}$  and  $c_{ij}^T > 0$  for every  $T$  and  $i < j$ . Then minimum evolution using the functions in  $\{L^T\}_T$  is consistent.*

*Proof.* Since  $L^T \in \mathcal{U}(T)$ , we just need to prove that  $L^T(\sigma^{A|B}) > 1$  for all bipartitions  $A|B$  of  $\{1, 2, \dots, n\}$  that are not splits of  $T$  to conclude that  $L^T \in \mathcal{U}^+(T)$  and therefore ME is consistent.

Let then  $A|B$  be a bipartition of  $\{1, 2, \dots, n\}$  but not a split of  $T$ . It is easy to realise that there must be in  $T$  four disjoint clades  $A_1$ ,  $B_1$ ,  $A_2$  and  $B_2$ , such that: (1)  $A_1$  and  $B_1$  are separated by exactly two branches, (2)  $A_2$  and  $B_2$  are separated by exactly two branches, (3)

$A_1$  and  $A_2$  are contained in  $A$  and (4)  $B_1$  and  $B_2$  are contained in  $B$  (see fig. 5.1). It is easy to see that

$$\sigma^{A|B} = \sigma^{A_1|B_1} + \sigma^{A_2|B_2} + \mathbf{r},$$

where  $\mathbf{r} = (r_{ij})$  has some entries such that  $r_{ij} = 1$  and the remaining ones such that  $r_{ij} = 0$ .

Because  $L^T$  is linear,

$$L^T(\sigma^{A|B}) = L^T(\sigma^{A_1|B_1}) + L^T(\sigma^{A_2|B_2}) + L^T(\mathbf{r}).$$

Lemma 13 (proved below) shows that  $L^T(\sigma^{A_1|B_1}) = L^T(\sigma^{A_2|B_2}) = \frac{1}{2}$ . Because  $L^T$  only has positive coefficients and  $\mathbf{r}$  is binary with some positive entries, we also have that  $L^T(\mathbf{r}) > 0$ . We can then conclude that  $L^T(\sigma^{A|B}) > 1$ .  $\square$

**Lemma 13.** *Let  $T$  be a bifurcating topology and  $L^T \in \mathcal{U}(T)$ . For every two clades,  $A$  and  $B$ , separated by exactly two branches in  $T$ ,*

$$L^T(\sigma^{A|B}) = \frac{1}{2}.$$

*Proof.* Any internal node of a bifurcating tree  $T$  is attached to three branches, defining three clades  $A$ ,  $B$  and  $C$  each separated from the other two by exactly two branches. Because  $A|B \cup C$ ,  $B|A \cup C$  and  $C|A \cup B$  are all splits of  $T$ , and since  $L^T \in \mathcal{U}(T)$ , the following equalities must hold:

$$\begin{aligned} L^T(\sigma^{A|B \cup C}) &= L^T(\sigma^{A|B}) + L^T(\sigma^{A|C}) = 1, \\ L^T(\sigma^{B|A \cup C}) &= L^T(\sigma^{A|B}) + L^T(\sigma^{B|C}) = 1, \\ L^T(\sigma^{C|A \cup B}) &= L^T(\sigma^{A|C}) + L^T(\sigma^{B|C}) = 1. \end{aligned}$$

Solving the system above for the unknowns  $L^T(\sigma^{A|B})$ ,  $L^T(\sigma^{A|C})$  and  $L^T(\sigma^{B|C})$  leads to the following solution:

$$L^T(\sigma^{A|B}) = L^T(\sigma^{A|C}) = L^T(\sigma^{B|C}) = \frac{1}{2}.$$

Since any two clades separated by exactly two branches must be in the configuration outlined above, our proof is complete.  $\square$

Can then our results on the robustness of BME be extended more generally to all versions of ME using functions from  $\mathcal{U}(T)$  with positive coefficients? Unfortunately, the following example shows that if we estimate the length of each tree using an arbitrary element of this class of functions, ME may have an arbitrarily small safety radius.

**Example 14.** Let  $T$  and  $W$  be defined as in Figure 4.1 except that now we assume that the four corner clades only contain one taxon each, with  $A = \{1\}$ ,  $B = \{2\}$ ,  $C = \{3\}$ ,  $D = \{4\}$ . Give the branches of  $T$  the following (positive) lengths:  $\ell_1, \ell_2, \ell_3, \ell_4$  for the terminal branches incident with 1, 2, 3, 4, respectively and  $\ell_0$  for the internal branch. Call the resulting weighted tree  $\mathcal{T}$  and assume  $\ell_0 = \ell_{\min}^{\mathcal{T}}$ .

Assume ME uses the following formulae to estimate the lengths of  $T$  and  $W$ :

$$L^T(\delta) = \frac{1}{2}(\delta_{12} + \delta_{34}) + \frac{1}{4}(\delta_{13} + \delta_{14} + \delta_{23} + \delta_{24}),$$

$$L^W(\boldsymbol{\delta}) = L_\epsilon^W(\boldsymbol{\delta}) = \frac{1}{2}(\delta_{13} + \delta_{24}) + \left(\frac{1}{2} - \epsilon\right)(\delta_{12} + \delta_{34}) + \epsilon(\delta_{14} + \delta_{23}),$$

with  $\epsilon$  being a small positive number. (It is unimportant what formula ME uses for the third bifurcating topology over  $\{1, 2, 3, 4\}$ .) It is easy to check that  $L^T \in \mathcal{U}(T)$ ,  $L_\epsilon^W \in \mathcal{U}(W)$  and that they both have positive coefficients. We now show that by making  $\epsilon$  arbitrarily small, we can also make the safety radius of ME arbitrarily small.

Let  $\alpha$  be such that  $0 < \alpha \leq \frac{1}{2}$ . Define  $\boldsymbol{\delta}$  in the following way:

$$\begin{aligned}\delta_{12} &= d_{12}^T + \alpha l_0 = l_1 + l_2 + \alpha l_0, \\ \delta_{34} &= d_{34}^T + \alpha l_0 = l_3 + l_4 + \alpha l_0, \\ \delta_{14} &= d_{14}^T + \alpha l_0 = l_1 + l_4 + (1 + \alpha)l_0, \\ \delta_{23} &= d_{23}^T + \alpha l_0 = l_2 + l_3 + (1 + \alpha)l_0, \\ \delta_{13} &= d_{13}^T - \alpha l_0 = l_1 + l_3 + (1 - \alpha)l_0, \\ \delta_{24} &= d_{24}^T - \alpha l_0 = l_2 + l_4 + (1 - \alpha)l_0.\end{aligned}$$

Clearly,  $\|\boldsymbol{\delta} - \mathbf{d}^T\|_\infty = \alpha l_0$ . We also have:

$$L^T(\boldsymbol{\delta}) = l_1 + l_2 + l_3 + l_4 + (1 + \alpha)l_0$$

and

$$L_\epsilon^W(\boldsymbol{\delta}) = l_1 + l_2 + l_3 + l_4 + (1 + 2\epsilon)l_0.$$

These formulae show that, if  $\epsilon < \frac{1}{2}\alpha$ , then  $L_\epsilon^W(\boldsymbol{\delta}) < L^T(\boldsymbol{\delta})$ , that is ME will favour a wrong topology  $W$  over the true topology  $T$ .

Therefore the safety radius of ME using  $L^T(\boldsymbol{\delta})$  and  $L_\epsilon^W(\boldsymbol{\delta})$  (with  $\epsilon$  sufficiently small) is limited above by  $\alpha$  (if the radius were greater than  $\alpha$ , ME should return  $T$  on the  $\boldsymbol{\delta}$  defined above). Since this holds for an arbitrarily small  $\alpha$ , we conclude that using ME with correct linear formulae with positive coefficients does not guarantee any positive safety radius.

The example above suggests that in order to guarantee a certain safety radius, the coefficients of the used linear functions must not only be positive, but also “large enough”. This is an interesting idea for further research, but outside the scope of the present paper.