1    # Information Content of Trees: Three-taxon Statements Inference

2    # Rules and Dependency

3    Valentin Rineau, René Zaragüeta, Jérémie Bardin

4    [1]Valentin Rineau*, [2] René Zaragüeta and [3]Jérémie Bardin

5    [1] *Center for Theoretical Study, Charles University & Czech Academy of Sciences, Jilská 1,*
6    *110 00 Praha 1, Czech Republic.*

7    [2] *Sorbonne Université, CNRS, MNHN, EPHE, Université des Antilles, Institut de*
8    *Systématique, Evolution, Biodiversité, ISYEB, CP 48, 57 rue Cuvier, 75005 Paris, France.*

9    [3] *Sorbonne Université, Centre de Recherche en Paléontologie – Paris (CR2P), 4 place*
10   *Jussieu, barre 46-56 5ème étage, case 104, 75005 Paris, France.*

11   *Corresponding author. Email: valentin.rineau@gmail.com.

Rineau, Zaragüeta and Bardin

12 ABSTRACT

13    The three-taxon statement (also called triplet) is the fundamental unit of rooted trees in

14  phylogenetic systematics. Various supertree and phylogenetic methods use three-taxon

15  statements that are minimal rooted statements of degree of kinship relationships. Because of

16  their fundamental role in phylogenetics, three-taxon statements are present in methodological

17  research of various disciplines in evolutionary biology, as in consensus methods, supertree

18  methods, species-tree methods, distance metrics, phylogenetics, and cladistic biogeography.

19  Three-taxon statements are thus widely used. However, their theoretical properties have been

20  poorly investigated. As a result, three-taxon statements methods are subject to important

21  flaws related to information redundancy. Correcting these biases is essential to improve the

22  efficiency of methods using three-taxon statements. Our aim is to study the behavior of three-

23  taxon statements and the interactions among them in order to enhance their performance in

24  phylogenetic studies. We have identified new types of very specific interactions between

25  three-taxon statements responsible of the emergence of redundancy and dependency in trees.

26  We propose for the first time a classification of three-taxon statements interactions and trace

27  the link between those and the emergence of dependency and redundancy. A new fractional

28  weighting procedure for suppressing redundancy of three-taxon statements is proposed. Our

29  method is subsequently empirically tested in the supertree framework using simulations. We

30  show that three-taxon statements using fractional weights perform drastically better than

31  classical supertree methods such as MRP or methods using unweighted three-taxon

32  statements. Our study shows that appropriate fractional weighting of three taxon statements is

33  an efficient measure of phylogenetic information content for rooted trees. Fractional

34  weighting is of critical importance for removing redundancy in any method using three-taxon

2

35    statements, as in consensus, supertrees, distance metrics, and phylogenetic or biogeographic

36    analyses.

37    Keywords: fractional weighting, information content, three-taxon statements, triplets,

38    cladistics, phylogenetic analysis, supertree methods, consensus trees.

39    Phylogenetic trees are at the heart of most evolutionary biology studies and their

40    reliability is of critical importance. It has long been recognized in the phylogenetic setting the

41    need to understand and manipulate the information contained in phylogenies. Any tree (taxon

42    tree, gene tree, biogeographic tree, etc.) can be fruitfully decomposed into subunits of several

43    kinds (Adams 1986; Vach 1994), as components (rooted subtrees with all taxa and a single

44    informative node; Nelson & Platnick 1981, Wilkinson et al. 2004), three-taxon statements

45    (3ts thereafter; rooted subtrees with three taxa and a single informative node; Adams 1986),

46    splits (bipartitions of a set of taxa; Farris 1970), or quartets (bipartitions of 4 taxa; Strimmer

47    and von Haeseler 1996). These subunits are traditionally considered in the phylogenetic

48    framework as information bricks (e.g. Mickevich & Platnick 1989; Nelson 1979; Nelson &

49    Ladiges 1992; Wilkinson 1994a), and can be used to measure the information content of a

50    phylogenetic tree.

51    Here we focus on rooted trees. Within a rooted tree, the 3ts (see Nelson & Ladiges 1991a,

52    b; Nelson & Platnick 1991) take a special place as they represent the minimal (rooted)

53    phylogenetic information. They have also been called "triad" in Adams (1986) and Vach

54    (1994), "triplet" in Wilkinson (1994a), or still "rooted triple" in (Bryant 2003) in the

55    consensus and supertree literature. A 3ts is a statement in which given three distinct taxa, two

56    are related compared to a third one. For example, given three taxa $a$, $b$ and $c$, $a$ and $b$ are

57    more closely related to each other than either is to $c$: $(a,b),c$.

58    Because of their atomic nature, 3ts are at the heart of many methods. For example, the

59    method of phylogenetic reconstruction commonly referred to as three-taxon analysis (Nelson

60    and Ladiges 1991b; Nelson and Platnick 1991) decomposes hierarchical characters into 3ts

61    matrices and then selects the cladograms that are congruent with the maximal amount of

4

62    those 3ts (Zaragüeta et al. 2012). The same principle is applied in cladistic biogeography and

63    the construction of areagrams (i.e. cladograms of biogeographic areas; Nelson and Ladiges

64    1991a). The use of 3ts in the supertree framework was first proposed by Williams (2004).

65    Currently, several supertree methods (Ranwez et al. 2010; Sevillya et al. 2016; Dannenberg et

66    al. 2019; Mavrodiev et al. 2019) explicitly handle 3ts. Another use of 3ts is also for

67    consensus tree representations (Aho et al. 1981; Adams 1986; Nelson & Ladiges 1994;

68    Wilkinson 1994a; Bryant 2003; Cao et al. 2009). The 3ts permit to precisely quantify the

69    information content of a phylogenetic structure (Mickevich and Platnick 1989; Nelson and

70    Ladiges 1992; Williams and Humphries 2003). The retention index based on 3ts (Kitching et

71    al. 1998) allows measuring the proportion of hierarchical relationships that are found in the

72    optimal cladograms. 3ts-based distance metrics have been used to compare topologies among

73    sets of phylogenetic trees (Grand et al. 2013; Kuhner and Yamato 2015). Finally, 3ts are also

74    explicitly used for coalescent-based species tree estimation from gene tree (Liu et al. 2010;

75    Islam et al. 2020), and for phylogenetic network reconstruction (Poormohammadi et al.

76    2020).

77       The use of 3ts requires their independence. Nelson and Ladiges (1992) showed that the

78    independence among a 3ts set could be achieved by using a fractional weighting to remove

79    redundancy. Wilkinson et al. (2004) criticized this procedure by showing that certain

80    relationships among 3ts generating dependency were not taken into account in the computation

81    of fractional weighting. However, these relationships were neither clearly defined nor

82    formalized in their work. A thorough work on relationships among 3ts in order to propose new

83    operational solutions to remove dependency becomes essential today. The goal of this paper is

84    dual: i) to develop the relationship between dependency and relationships among 3ts and ii) to

85    propose a weighting method to improve the quality and accuracy of phylogeny analyses using

86    3ts.

87        In the first part of this paper, we summarize a historical review of the procedures

88    proposed in the literature to take into account dependency in trees. We highlight main

89    challenges and discuss the solutions that have been proposed in the past. In the second

90    section, we demonstrate the link between tree-shape and dependency within phylogenetic

91    trees. We recognize different types of relationships between internal nodes that lead to

92    dependency. In the third part, we exhaustively describe, formalize, and classify the different

93    relationships among couples of 3ts, and how to accordingly combine 3ts in trees or

94    decompose trees in 3ts. Our work on tree-shape and on 3ts is then used in a fourth part to

95    propose a new procedure called *corrected fractional weighting* that removes all types of

96    dependency between 3ts. This procedure is needed to fulfil the requirement of independence

97    on which are based most of the analyses using phylogenetic information. Finally, we used

98    simulated datasets to compare different supertree methods that use different representation of

99    phylogenetic information and show that methods using fractional weighting performs

100    significantly better.

101    DEPENDENCY AND WEIGHTING 3TS: AN HISTORICAL REVIEW

102        The use of 3ts as a unit of measurement for the information content of phylogenies

103    requires that they be independent of each other (i.e. absolute independence *sensu* Wilkinson

104    et al. 2004). Performing an analysis based on non-independent data leads to biases that can

105    have important consequences on the quality of the result. Independence between two

106    relationships implies that the truth of one does not involve the truth or falsity of the other.

107    The 3ts $a(bc)$ and $a(bd)$ are independent because the truth of one tells nothing the truth

6

108   of the other. On the other hand, the 3ts $a(bc)$, $a(bd)$, and $a(cd)$ are not independent relative to

109   each other, because the truth of two of these 3ts necessarily implies the truth of the third

110   (Nelson and Ladiges 1992):

111       $$a(bc) + a(bd) + a(cd) \rightarrow a(bcd)$$

112       $$a(bc) + a(bd) \rightarrow a(bcd)$$

113       $$a(bc) + a(cd) \rightarrow a(bcd)$$

114       $$a(bd) + a(cd) \rightarrow a(bcd)$$

115   Relationships among 3ts may lead to dependency and therefore to redundancy. In the given

116   example, the information content of these three 3ts cannot be 3; since the information given

117   by the three 3ts is identical to that given by two of them, there is a repetition of information

118   that leads to the overvaluation of 3ts.

119       Fractional weighting (FW; Nelson and Ladiges 1992) was the first attempt to correct the

120   overweighting of 3ts. Since any two 3ts are needed among the three, the amount of

121   information carried by each 3ts is not one but 2/3. The corrected value of the set is 2. A

122   component containing $t$ taxa of which $n$ are included in the informative node contains $n(t-$

123   $n)(n-1)/2$ 3ts, but only requires $(n-1)(t-n)$ 3ts to be reconstructed. This quantity $(n-1)(t-n)$

124   corresponds to the phylogenetic value of a component sensu Nelson and Ladiges (1992). In

125   this case, the value of a 3ts is $2/n$, i.e. the number of independent 3ts divided by the total

126   number of 3ts.

127       The presence of several nested internal nodes implies another type of dependency. The

128   first attempt of fractional weighting by Nelson and Ladiges (1992) was applied component by

129   component, herein *fractional weighting per component* ($FW_{comp.}$) did not take into account

7

Rineau, Zaragüeta and Bardin

130    the dependence. Indeed, in FW$_{comp.}$, the weight of a 3ts present in several components equals

131    the sum of its weights for each component in which it appears. The tree $a(b(cd))$ can be

132    decomposed into two components and then into 3ts:

133        $a(b(cd) \rightarrow a(bcd) + ab(cd)$

134        $a(bcd) \rightarrow a(bc) + a(bd) + a(cd)$ (weight for each 3ts: 2/3)

135        $ab(cd) \rightarrow a(cd) + b(cd)$ (weight for each 3ts: 1)

136    The information content of the tree is 4. 3ts weights are subject to redundancy: 3ts $a(cd)$ is

137    overweighted (5/3) because it is present in both components $a(bcd)$ and $ab(cd)$ while it is

138    present only once in $a(b(cd))$. For this reason, Nelson and Ladiges (1992: 491) describe

139    FW$_{comp.}$ and *Uniform Weighting* (sum of 3ts per component, see table I) as "*both misleading*".

140    A 3ts cannot be present more than once in a tree, and therefore must not have a weight

141    greater than 1. The table I shows another problematic example with the tree $a(bc(de))$. To

142    overcome this difficulty, Nelson and Ladiges (1992) propose a different way to calculate the

143    absolute value of a 3ts called herein *total fractional weighting* (FW$_{total}$). It is a reappraisal of

144    the idea that the weight of a 3ts is the number of independent 3ts divided by the total number

145    of 3ts, but at the scale of the whole tree, by summing the weights of 3ts for each component:

146    "*For each of their cladograms, we list the number of independent three-taxon statements and*

147    *the absolute value of each of all possible statements, stated as the ratio of independent*

148    *statements to all possible statements*" (Nelson and Ladiges 1992: 491). The FW$_{total}$ becomes

149    for a 3ts the sum of all independent 3ts for each component / the sum of all 3ts for each

150    component. For C components deducible from a phylogenetic tree:

8

INFORMATION CONTENT OF TREES

151
$$3is_{total\ FW} = \frac{\sum_{i=1}^{C}(n_i-1)(t-n_i)}{\sum_{i=1}^{C}\dfrac{(t-n_i)(n_i-1)n_i}{2}} = \frac{2}{\sum_{i=1}^{C}n_i}$$

152    As a result, all the 3ts of a tree have the same weight. The four 3ts of the tree $a(b(cd))$ have a

153    weight of 1. Nelson and Ladiges conclude by considering this calculation of the absolute

154    value of a 3ts as the most relevant (Nelson and Ladiges 1992: 494).

155        Wilkinson et al. (2004) gave a new attempt to describe the dependencies between 3ts using

156    Dekker's (1986) dyadic inference rules. According to Dekker, only two inferences are

157    possible from pairs of quartets to deduce additional quartets. By transposing these inference

158    rules to 3ts, Wilkinson et al. (2004: 994) propose three rules for inferring 3ts exemplified as

159    follows:

160        $b(cd) + b(ce) \rightarrow b(de)$

161        $b(cd) + c(de) \rightarrow b(ce)$

162        $d(ab) + a(de) \rightarrow e(ab)$

163    The original Dekker rules imply that any tree containing two quartets involved by the dyadic

164    inference rules also contains an additional quartet. Wilkinson et al. (2004) point out that for

165    3ts, several rules can intervene from a single couple. For example:

166        $b(cd) + c(de) \rightarrow b(ce)*$

167        $b(cd) + b(ce)* \rightarrow b(de)$

168    Which results in:

169        $b(cd) + c(de) \rightarrow b(c(de)) \rightarrow b(cd) + c(de) + b(ce) + b(de)$ (fig. 4b in Wilkinson et al.
170    2004)

9

171    This implies that the presence of some couples of 3ts can imply the truth of two secondary

172    3ts. For Wilkinson et al. (2004), this type of dependence is not taken into account by Nelson

173    and Ladiges' fractional weighting. Since the tree $b(c(de))$ needs only two 3ts to be rebuilt, its

174    weight should be 2. Consequently, the FW$_{total}$ can only take into account one type of

175    dependence corresponding to the first rule presented above. These rules lead to the problems

176    of dependency inside a set of 3ts; the purpose of weighting being to suppress them.

177        However, Wilkinson et al. do not give all the possible combinations that lead to additional

178    3ts. This is justified by the fact that any fully resolved tree with $t$ taxa can be fully entailed by

179    specific sets of $t$-2 3ts (Steel, 1992). Consequently, any dichotomous tree includes only $t$-2

180    independent 3ts, each triplet having a weight of 6 / ($N^2$ - $N$). Wilkinson et al. named this

181    procedure *Minimal Weighting* (MW), which aimed to take into account all possible inference

182    rules without having to define them. The arguments put forward in favor of the MW are: (i)

183    an identical weight for all 3ts because there is no argument to weight differentially the 3ts

184    from a tree, (ii) a linear evolution of the total weight of a tree with the number of taxa

185    (polynomial with FW$_{comp.}$ and FW$_{total}$), and (iii) an evolution of the total weight of a tree that

186    is not influenced by tree shape.

187    DEPENDENCY AND TREE SHAPE

188    Our aim being to find a way of deleting dependency issues in phylogenetic trees, we must

189    define in a first place where the dependency can arise. Here we show for the first time how

190    the tree shape can modify the dependency between nodes, and by corollary altering the

191    information content of phylogenies. To fulfil this aim, we need to precisely define the various

192    types of nodes of phylogenetic trees. There are three types of internal nodes : apical nodes,

193    orthologous nodes, and paralogous nodes.

10

INFORMATION CONTENT OF TREES

194    **Definition 1.** *Apical nodes* are nodes that only include leaves. The tree $b(c(de))$ contains only

195    one apical node, $(de)$.

196    **Definition 2.** *Orthologous nodes* (sensu Zaragüeta et al. 2004, also called asymmetric or

197    unbalanced nodes) include only one internal node as direct descendant. In the tree $b(c(de))$,

198    the node $(cde)$ is orthologous because it leads only to the apical node $(de)$, and the root node

199    $(bcde)$ because it immediately leads only to $(cde)$.

200    **Definition 3.** *Paralogous nodes* (sensu Zaragüeta et al. 2004, also called symmetric or

201    balanced nodes) are nodes that include two or more internal nodes as direct descendants. The

202    tree $(ab)(cd)$ contains one paralogous node, the root $(abcd)$, which contains two internal

203    nodes $(ab)$ and $(cd)$.

204       The relationships between internal nodes in the tree can be linked to that of dependence.

205    Orthologous nodes show dependency where a node is differentiated in another node, while

206    paralogous nodes show occurrences of independence, where a node is divided into several

207    independent and separate lineages. Dependency, i.e. inclusion, is a kind of information that

208    may be taken into account, in contrast to independency, i.e. disjunction. Therefore, a

209    pectinate tree contains more information that a symmetric tree containing independent nodes.

210    As we believe as Williams and Ebach (2008: 213) that in a phylogenetic analysis, a multi-

211    state character is more informative than the pair of corresponding binary characters (because

212    the relationships between states are explicitly showed only in the multi-state character), we

213    can say more broadly that a tree carries more information than the sum of all its components.

214    The dependence or independence between nodes should have an impact on the information

215    content, and thus on 3ts weights. Consequences for 3ts are explained in the next section.

11

Rineau, Zaragüeta and Bardin

216 CLASSIFICATION OF 3TS STATEMENTS RELATIONSHIPS

217 Relationships among 3ts can be classified relatively to dependency. We propose here a

218 first complete classification of the types of relationships between pairs of 3ts – dyadic

219 relationships. We show that the specifics of relationships among 3ts depends on their

220 common leaves and positions. Two 3ts can be identical or not, compatible or not, combinable

221 or not. Combinability can itself be decomposed in four possible relationships. All definitions

222 are given below.

223 **Definition 1.** Two distinct 3ts are *compatible* if they can be true at the same time, i.e. if it

224 exists *at least* one tree containing both.

225 **Corollary.** Two 3ts are incompatible if they cannot be true at the same time, i.e. the truth

226 of one necessarily implies the falsity of the other. Two 3ts are incompatible if they have

227 exactly the same taxa but state different relationships (e.g. $a(bc)$ and $b(ac)$).

228 **Definition 2.** Two compatible 3ts are *combinable* if they can be unambiguously gathered

229 in a single tree. In other words, two 3ts are combinable if and only if they are present in the

230 strict consensus of all possible trees containing them. Formally, two 3ts with their respective

231 taxa sets D = $\{u, v, w\}$ and E = $\{x, y, z\}$ are combinable if it exists an informative strict

232 consensus which taxa set $F$ satisfies $|F| = |D \cup E| = 4$. We note that 3ts are combinable if and

233 only if they have two taxa in common.

234 **Corollary.** If two 3ts are not combinable, the above-mentioned strict consensus is a bush

235 (a tree without any informative node). Combinability logically implies compatibility, but the

236 reverse is not true: compatible couples of 3ts sharing less than two taxa are non-combinable.

237 The couple $a(bc)$ and $a(de)$ is non-combinable because $a$ is the only common taxon, and

12

238    entails several possibilities, as $(de)(a(bc))$ or $a(b(d(ce)))$ and which strict consensus is a bush

239    (no unambiguous common statement between all possible trees).

240        Four relationships among combinable 3ts exist, with specific properties, depending on the

241    location of their common taxa (Figure 1).

242        **Definition 3.** Two combinable 3ts are in *in-in nts relationship* if both common taxa are

243    located inside the apical node of both 3ts, and the combination of the 3ts results in a single n-

244    taxon statement (i.e. a tree with a single informative node; nts thereafter; Wilkinson 1994).

245        The example of $a(cd) + b(cd)$ allows us to build the tree $ab(cd)$. The tree can be

246    decomposed into the two 3ts that allowed its construction:

247            $a(cd) + b(cd) \rightarrow ab(cd) \rightarrow a(cd) + b(cd)$

248    Information conveyed by each 3ts allows an additional taxon to be added as a sister group to

249    the apical node ($cd$). The fact that no additional 3ts is generated implies that no dependency is

250    generated: that is the reason two-leaf apical nodes never cause fractional weight, as Nelson

251    and Ladiges (1992) noticed.

252        **Definition 4.** Two combinable 3ts are in *in-out nts relationship* when one of the common

253    taxa is connected directly to the root and the other is connected to the apical node in both 3ts.

254    The combination of two 3ts in an in-out nts relationship leads to build a nts with a single

255    informative node of three leaves.

256        The couple $a(bc)$ and $a(cd)$ illustrates the in-out nts relationship, whose analysis highlights

257    the implication of a third 3ts:

258            $a(bc) + a(cd) \rightarrow a(bcd) \rightarrow a(bc) + a(cd) + a(bd)$

13

259    The in-out nts relationship between two 3ts generates a third 3ts with two taxa in common

260    with the other two 3ts, at the same positions. The truth of these two 3ts implies the truth of a

261    third. Among the three 3ts $a(bc)$, $a(cd)$, and $a(bd)$, it can be seen that whatever the two 3ts

262    chosen, they will be in the in-out nts relationship and will involve the third one. It is precisely

263    this kind of relationships which is taken into account by Nelson and Ladiges' FW to correct

264    the redundancy. This type of dependency is the only one present within the components. If a

265    new 3ts is added and forms additional in-out nts relationships, then all 3ts can be combined,

266    as in-out relationship is transitive:

267          $a(bc) + a(cd) + a(de) \rightarrow a(bcde) \rightarrow a(bc) + a(cd) + a(de) + a(bd) + a(be) + a(ce)$

268    Three 3ts are linked by two in-out relationships. Each in-out relationship generates an

269    additional 3ts:

270          $a(bc) + a(cd) \rightarrow a(bd)$

271          $a(cd) + a(de) \rightarrow a(ce)$

272    The generation of the secondary 3ts leads to the appearance of two new in-out dependency

273    relationships, leading to the same 3ts:

274          $a(bd) + a(de) \rightarrow a(be)$

275          $a(bc) + a(ce) \rightarrow a(be)$

276    We call the initial 3ts as primary 3ts, $a(bc)$, $a(cd)$, and $a(de)$ in the example. The 3ts

277    generated from the combination of primary 3ts are called secondary 3ts, here $a(bd)$ and $a(ce)$.

278    The 3ts generated by the combination of secondary 3ts are called tertiary 3ts, here $a(be)$ and

279    so on (to $n$-ary 3ts for the $n$th level of combination). For the combination of the set of three

280    primary 3ts, a total of six $n$-ary 3ts are entailed. This is why coding a set of six 3ts $a(bc)$,

14

281    $a(cd)$, $a(de)$, $a(bd)$, $a(be)$, $a(ce)$, requires a weight of ½ for each 3ts since only three 3ts

282    among the six are needed to build the node (given that all terminals are represented among

283    the set). From the three primary 3ts we can deduce all the secondary and tertiary 3ts. The

284    corollary is that 3ts may not satisfy the pairwise compatibility theorem (Estabrook et al.

285    1976; Wilkinson 1994b). Indeed, the fact that all primary 3ts are pairwise compatible does

286    not imply that all n-ary 3ts are compatible.

287    **Definition 5.** Two combinable 3ts are in *orthologous relationship* when both have a

288    common taxon connected to the apical node, and when the other common taxon is connected

289    to the apical node in one of the 3ts and to the root in the other. Unlike with the in-in and in-

290    out nts relationships, the trees constructed from two 3is linked by an orthologous relationship

291    have two informative nodes, one of which is orthologous (see fig. 4b in Wilkinson et al.

292    2004). If the tree $a(b(cd))$ is decomposed into its two components $ab(cd)$ and $a(bcd)$, the first

293    generates only one in-in nts relationship and the other three in-out nts relationships

294    (Appendix 1a). The combination of these two components implies the existence of two new

295    couples of 3ts that were not present in the independent components:

296        $b(cd)$ (from $ab(cd)$) + $a(bd)$ (from $a(bcd)$) → $a(b(cd))$ → $b(cd)$ + $a(bd)$ + $a(cd)$ +

297    $a(bc)$

298        $b(cd)$ (from $ab(cd)$) + $a(bc)$ (from $a(bcd)$) → $a(b(cd))$ → $b(cd)$ + $a(bd)$ + $a(cd)$ +

299    $a(bc)$

300    These two couples linked by an orthologous relationship produce the complete tree

301    $a(b(cd))$. If combined, 3is linked by orthologous relationships, therefore, lead to inter-node

302    dependency. This shows that the quantification of the information contained in a tree is not

303    equal to the sum of the information of each component. Methods that decompose trees in

15

304   their components without taking into account the relationships between these components,

305   such as the supertree method MRP (Baum 1992; Ragan 1992) or the biogeographical method

306   BPA (Wiley 1986; Brooks 1990), are thus flawed. The information conveyed by the inclusion

307   (dependency) of one node in another must be taken into account.

308   **Definition 6.** Two combinable 3ts are in *paralog relationship* when among the two taxa

309   they have in common, one is connected to the root in the first 3ts and the apical node in the

310   second 3ts, and the other is connected to the apical node in the first 3ts and to the root in the

311   second. The paralogous relationship appears when at least two informative nodes are disjoint

312   (i.e. if there is a paralogous node in the tree where the 3ts come from).

313   The following example illustrates a paralogous relationship between two 3ts that generates

314   a symmetric tree (Appendix 1b):

315   $$c(ab) + a(cd) \rightarrow (ab)(cd) \rightarrow c(ab) + a(cd) + d(ab) + b(cd)$$

316   This combination allows two 3ts, $d(ab)$ and $b(cd)$, to be deduced.

317   Some general corollaries may be listed from the set of above described relationships:

318   — Any hierarchical tree with at least one informative node can be analyzed into a set

319       of 3ts with no incompatible relationships.

320   — The four dyadic relationships that are in-in, in-out, orthologous, and paralogous are

321       necessary and sufficient to reconstruct any phylogenetic tree that can be deduced

322       from a set of 3ts, including Bryant and Steel's "irreducible high order rules" (1995:

323       445–452) that are not formalized but only exemplified, their finding being the

324       result of an incomplete axiomatization.

16

325        — Any tree containing an informative node and more than one leaf outside this

326            informative node contains in-in relationships (Fig. 1a).

327        — Any tree containing an internal informative node and more than two leaves

328            connected to that node contains in-out relationships (Fig. 1b).

329        — Compatible but not combinable 3ts (zero to one taxon in common) are present in a

330            tree as long as there is an informative node comprising more than three leaves. As

331            soon as there are at least two informative nodes in a tree, paralogous (fig. 1c)

332            and/or orthologous relationships (fig. 1d) appear. It is, therefore, possible to find

333            all these types of relationships within the same tree (Appendix 1c).


334    WEIGHTING PROCEDURES

335    Formalized types of relationships among 3ts that lead to dependency are the only way to

336    remove redundancy in 3ts decomposition. Firstly, we will show that MW fails to remove all

337    instances of dependency. Secondly, we will propose a new procedure, grounded on the FW

338    rationale, that succeeds at removing every dependency from any 3ts set.


339    *Minimal weighting*

340    Wilkinson et al. (2004) designed the MW in order to remove all dependencies issues using

341    a simple formula ($t$-2). However, there are two arguments for rejecting the MW.

342    The first concerns the identical weight given to each triplet, as in Nelson and Ladiges'

343    $FW_{total}$ computation. This weight implies that among the four 3ts of ($b(c(de))$), only two are

344    necessary to reconstruct the tree. The total weight of the tree (its information content) is 2,

345    and each 3ts has an information content of ½. Wilkinson et al. (2004) rationale implies that

346    each of the 3ts has the same information content. If this was true, any set of two 3ts among

17

347    $b(cd)$, $c(de)$, $b(ce)$, and $b(de)$, would allow to reconstitute the tree, which is not correct. Here

348    is the result of the different combinations of 3ts:

349        (1) $b(cd) + c(de) \rightarrow b(c(de))$

350        (2) $c(de) + b(ce) \rightarrow b(c(de))$

351        (3) $b(cd) + b(ce) \rightarrow b(cde)$

352        (4) $b(cd) + b(de) \rightarrow b(cde)$

353        (5) $b(ce) + b(de) \rightarrow b(cde)$

354        (6) $c(de) + b(de) \rightarrow bc(de)$

355    (1) and (2) are the only couples that allow the tree to be fully reconstituted. The other couples

356    allow only one of the two components to be rebuilt: the couples (3) to (5) allow rebuilding

357    $b(cde)$, and the couple (6) allows rebuilding $bc(de)$. In order to keep only the minimum

358    number of 3ts, the statement "any 2 among 4" is incorrect. Therefore, MW is inconsistent

359    with Wilkinson et al.'s (2004) rationale.

360      The second argument against MW is a strong limitation in its conception: it can only

361    handle fully dichotomous trees. The total weight of a tree of $t$ taxa is $t$-2 for a totally

362    bifurcating tree and, to our knowledge, no analytical solution has so far been proposed to

363    include multifurcations. Using the total weight formula of $t$-2 for polytomous trees implies

364    that the information content of a given tree is identical to that of any of its subtrees as long as

365    they have the same leaves. The trees $(ab)(cd)$, $ab(cd)$ and $(ab)cd$ all have an information

366    content of 2 according to this formula. This conclusion is clearly illogical.

367    *Corrected fractional weighting*

368      To successfully eliminate dependency among a 3ts set, the different types of 3ts

369    relationships and the scale at which they operate (intra-node or inter-node) must be taken into

18

370  account. MW fails to recognize in detail the different types of relationships. Both $FW_{total}$ and

371  $FW_{comp.}$ takes into account a single relationship (in-out nts), which is insufficient. Our

372  strategy to suppress the dependency is set up in sequential steps, each step corresponding to a

373  type of 3ts relationship (table II). We will explain here a corrected version of the FW ($FW_{cor.}$)

374  that take into account all types of dependency defined above.

375  The first step corresponds to the procedure set out by Nelson and Ladiges in 1992. It

376  consists in removing the dependency linked to in-out relationships within each component

377  itself. The formula for the number of independent 3ts for a component is $(t-n)(n-1)$, $t$ being

378  the total number of taxa and $n$ being the number of taxa connected to the informative node

379  (Nelson and Ladiges, 1992). The factor $(t-n)$ corresponds to the number of leaves connected

380  directly to the root. Increasing their number only generates in-in relationships. The factor $(n-$

381  $1)$ is linked to in-out relationships. For $a(bcde)$, the factor $(n-1)$ is 3. Only three 3ts are

382  needed to rebuild the tree. The factor $(t-n)$ is 1 and acts as a multiplier according to the

383  number of leaves at the root. For example, by adding a leaf $f$ to the root, $(n-1)$ does not

384  change. $(t-n)$ on the other hand, equals 2:

385      $a(bcde)$ [3] + $f(bcde)$ [3] → $af(bcde)$ [6]

386  The numbers in square brackets are the weights of the corresponding trees. The minimum

387  number of pairs of leaves pairs (e.g. $a,b$ and $c,d$) necessary for each leaf to be in a pair and

388  each pair to have a leaf in common with at least one another pair is $(n-1)$. $(n-1)$ is then

389  multiplied by the number of leaves connected to the root $(t-n)$ to give the formula $(t-n)(n-1)$.

390  $(t-n)$ does not change between the formula of the number of independent 3ts and the formula

391  of the total number of 3ts. On the other hand, the value $(n-1)$ changes: to have the total

392  number of 3ts, we need to search for all the combinations of two leaves among $n$, i.e.:

19

393
$$\binom{n}{2} = \frac{n!}{2(n-2)!} = \frac{n(n-1)}{2}$$

394     By replacing (n-1) by $\binom{n}{2}$ in (t-n)(n-1), we obtain the total number of 3ts:

395
$$\text{Number of 3ts for a component} = \frac{n(n-1)(t-n)}{2}$$

396     Nelson and Ladiges' $FW_{comp.}$ is, therefore, $FW_{cor.}$'s first step in eliminating redundancy.

397     This step of removing intra-node redundancy requires the analysis of each component

398     independently. For each component, all its 3ts are generated, and the weight of 2/$n$ (number

399     of independent 3ts / total number of 3ts) is assigned to each 3ts. This weight eliminates

400     overweight due to in-out relationships.

401     For example, component $a(bcd)$ can be decomposed into three 3ts: $a(bc)$, $a(cd)$, $a(bd)$.

402     Each of the three pairs of 3ts forms an in-out nts relationship which logically implies the third

403     3ts. Each 3ts is therefore overweighed because it exists both as a primary and secondary 3ts.

404     More generally, a 3ts present in a component is overweighed if it can be generated from other

405     3ts. This first step makes it possible to define the phylogenetic weight of each node.

406     $FW_{cor.}$ is obtained by removing the dependency linked to in-out nts relationships, as in

407     $FW_{comp}$, and then to orthologous relationships among 3ts of different nodes (inter-node

408     redundancy). Since 3ts have already been weighted in the previous step, the correction of the

409     orthologous dependency (i) concerns inclusions between informative nodes are present in the

410     tree, and (ii) revises the weight of the 3ts assigned in the first step.

411     Nelson and Platnick (1991) showed that every 3ts generated from a fully pectinate tree has

412     a weight of 1, applying FW or not because there is no redundancy between the 3ts of a fully

413     pectinate tree. The correction of inter-node redundancy requires checking the 3ts of each pair

20

414    of nodes involving direct inclusion. Two 3ts in an orthologous relationship will generate

415    secondary 3ts that already exist as primary ones. Using the example of the tree $a(b(cd))$

416    analyzed in Figure 2, Nelson and Ladiges (1992: p. 491) showed that the statement $a(cd)$

417    occurs in both series. This 3ts is the only one having a weight greater than 1 5/3 (2/3 + 1): it

418    is the overweighed 3ts. The component $a(bcd)$ bears the redundant 3ts because it is the only

419    component of the tree that contains this particular 3ts corrected by $FW_{comp}$. The other

420    component consists only of in-in relationships: its 3ts are therefore all essential to the

421    reconstruction of the tree and have a weight of 1. The analysis of the two orthologous

422    relationships of the tree highlights the overweight of $a(cd)$:

423        $a(bc) + b(cd) \rightarrow a(bd) + a(cd)$

424        $a(bd) + b(cd) \rightarrow a(bc) + a(cd)$

425    The secondary 3ts generated are $a(bd)$, $a(bc)$ and $a(cd)$ (twice). All the 3ts resulting from

426    orthologous relationships of both components allow generating only one component: $a(bcd)$.

427    However, among all the 3ts, $a(cd)$ appears twice while the others appear only once because

428    the orthologous relationship that overweighed this particular 3ts produces also a

429    supplementary secondary 3ts, the only 3ts present in both components. The uncertainty of

430    'two 3ts out of three' from the component $a(bcd)$ can then be resolved: the two necessary 3ts

431    are $a(bd)$ and $a(bc)$. $a(cd)$ is deleted from the component $a(bcd)$, and its weight (2/3) is

432    equally distributed among the others — if its weight was not distributed between the

433    remaining 3ts, the result would be that knowledge about orthologous relationship between

434    nodes would lead to a loss of phylogenetic information. In $(a(b(cd)))$, the weight of each of

435    the four 3ts becomes 1. The ambiguity is resolved. The analysis of inter-node relationships

21

436    allows us to correct the weights of 3ts from intra-node relationships because the inclusion of

437    one node in another carries specific phylogenetic information.

438    To exemplify how polytomies are handled by this procedure, we apply these two steps to

439    the tree $a(bc(de))$ (Figure 3 and Table I). The six 3ts generated by the first component are

440    weighted 1/2 (in-out relationships) and the three 3ts generated by the second component are

441    weighed 1 (in-in relationships). The total information content of the tree is 6. It is noted that

442    $a(de)$ is present in both components, and can, therefore, be removed from the component

443    $a(bcde)$. After correction, the total weight of the tree is still 6, but the weights of the 3ts have

444    been corrected. This tree shows a polytomy. Unlike the weights of the 3ts derived from the

445    tree $(a(b(cd)))$, which all end up balanced at 1, the weights of some of the 3ts derived from

446    $(a(bc(de)))$ are still fractional, even if the correction makes them tend towards 1. The

447    dichotomous pectinate structure is the most informative: it contains the maximum number of

448    components, and these components are all dependent on each other. Uncertainty increases

449    with the number of polytomies, and with the number of branches involved in polytomies.

450    The last type of dependency is related to paralogous relationships. Since paralogous

451    relationships only concern independent, i.e. disjoint, groups, there is no dependency to

452    manage at the level of 3ts for correcting 3ts weights. Indeed, paralogous relationships never

453    provide redundancy of any particular 3ts.

454    An example with the analysis of the tree $((ab)(cd))$ will illustrate this point (Appendix

455    1b):

456    $((ab)(cd))) \rightarrow (ab)cd + ab(cd)$

457    $(ab)cd \rightarrow c(ab) + d(ab)$

22

458      $ab(cd) \rightarrow a(cd) + b(cd)$

459      The weighting schemes proposed in the literature relay on the number of nodes. (i) The tree

460      has a weight of 2 (corresponding to MW). Taken *independently*, each component of the tree

461      also has a weight of 2; the choice of this weighting leads to accepting that the tree $(ab)(cd)$)

462      contains the same amount of phylogenetic information as $(ab)cd$ or $ab(cd)$. (ii) The tree has a

463      weight of 4 and can be divided into two components without any loss of information

464      (corresponding to the FW). However, since these two components are disjoint, they are

465      independent. Here, the tree information is the sum of the information of its two independent

466      components. Both components have a weight of 2. Paralogous relationships do not create

467      differential weights between the 3ts because no 3ts is overrepresented by n-ary statements.

468      Therefore, there is no justification for correcting in any way the weight of the 3ts intra-node.

469      Following Nelson and Ladiges (1992: 492): '*For cladogram 11, there is one internal*

470      *node, ((AB)(CD)), with four statements, each with an absolute value of ¾*'. The authors

471      provided no justification of their result (their result cannot be reached nor with the $FW_{comp.}$,

472      nor with the $FW_{total}$). However, as already pointed out by Nelson and Ladiges, if two totally

473      dichotomous trees with the same number of leaves (more than four) are compared, the tree

474      with more paralogous nodes will have less information content. The information content does

475      need to reflect that. The information is maximal in pectinate trees because they are devoid of

476      independent couples of nodes. For example, the weight of $a(b(c(de)))$ is 10 whereas the

477      weight of $a((bc)(de))$ is 9. In conclusion, our amended version $FW_{cor.}$ is the unique weighting

478      scheme that takes into account dependency emerging both from intra-node (in-out

479      relationships) and inter-node redundancy (ortholog relationships). Our proposal of $FW_{cor.}$ has

480      been implemented in the software LisBeth (Zaragüeta et al. 2012).

23

481    SIMULATIONS

482    *MATERIAL AND METHODS*

483    Following theoretical foundations of the relationship between dependency and weighting

484    developed in the previous parts, our aim here is to simulate datasets and analyses to

485    empirically test the efficiency of $FW_{cor.}$ versus current alternative methods. Each simulation

486    consists in 1) generate a tree (named hereafter 'the original tree'), 2) produce two subtrees, 3)

487    generate noise in the subtrees, 4) use the subtrees with different weighting schemes to 5)

488    reconstruct an optimal tree (and produce a strict consensus if several optimal trees exist) and

489    6) compare this optimal tree (or the strict consensus of all optimal trees) with the original one

490    to assess the relative efficiency of the different weighing schemes (Fig. 4). Each original tree

491    is randomly generated using the python package ete3 (Huerta-Cepas et al. 2016) and contains

492    11 terminal taxa. Two subtrees are generated. Each is a copy of the original tree from which a

493    defined number *m* of terminal taxa has been removed (*m* goes from 0 to 8; the choice of taxa

494    is random). For each generation, 1000 simulations are run, leading to a total of 9000 runs. All

495    the runs were repeated by introducing a varying amount of random noise by taxa

496    permutations to simulate homoplasy and/or distortion. The different amount of distortion is

497    included by pruning and regrafting a number *n* of taxa (*n* goes from 0 to 5). The choice of the

498    pruned taxa and the regrafting locations are both random. The 9,000 runs were thus repeated

499    6 times leading to 54,000 runs. To combine both subtrees in the optimal tree, we used two

500    supertree methods to reconstruct the optimal tree, one based on a parsimony approach (i.e.

501    Matrix Representation with Parsimony) using PAUP 4.0a (Swofford 2003), and the other

502    based on a 3ta approach (Three-taxon analysis; Nelson and Platnick 1991) with three

503    different weighing schemes: $FW_{cor.}$, $FW_{comp.}$, and MW (in our simulations MW and $FW_{total.}$

24

504    give always the exact same results because the trees of each analysis always have the same

505    number of terminal taxa), using Lisbeth (Zaragüeta et al. 2012). Note that when all trees of an

506    analysis have the same number of taxa, the total fractional weight would lead to the same

507    results as the minimal weight because the weights of all 3ts are the same for a given tree. The

508    total number of runs is increased to 216,000. The ability of the methods to reconstruct the

509    original tree is assessed by three complementary metrics. Triplet distance metrics (Tavares

510    2018) were preferred to Robinson-Foulds, because the latter is strongly influenced by

511    wildcard taxa. This is a common problem in supertree simulation analyses (Penny et al.

512    1982). Finally, the triplet distance metrics can be decomposed for our purpose in true

513    resolutions (*TR*) and false resolutions (FR; Grand et al. 2013; Rineau et al. 2018). The first

514    one, named *TR*, is the number of 3ts founded in both the original and optimal trees divided by

515    the number of 3ts of the original tree. This can be viewed as the proportion of the 3ts from the

516    original tree that is present in the optimal one. The second one, named *FR*, is the number of

517    3ts founded in the optimal tree but not in the original tree, divided by the number of 3ts of the

518    optimal tree. This can be viewed as the proportion of the 3ts from the optimal tree that were

519    not present in the original one. At last, the efficiency (*TR − FR*; range: [-100; 100]), allow to

520    synthesize the results. -100 is the case where all phylogenetic relationships are false and the

521    tree is completely resolved (i.e. dichotomous); 100 where all phylogenetic relationships are

522    true and the tree is completely resolved. A null result means that there are as many true and

523    false phylogenetic relationships, regardless of how well resolved the optimal tree is. In order

524    to characterize the quantity of information brought by optimal trees, we counted internal

525    nodes and 3ts from optimal tree for maximizing precision, hereafter named *constrict*

526    *resolution* and *constrict information content* respectively.

527    The results were then analyzed using R software (R Development Core Team 2008). We

528    first explored the simulations by performing a PCA to provide a global overview of the

529    covariations between the variables involved in the analyses. These variables consist in the

530    starting parameters (i.e. number of permutations and number of taxa in subtrees), the number

531    of optimal trees as well as their number of taxa, both the resolution and the information

532    content of the strict consensus and finally, the 3 measures of the results quality (i.e. true

533    resolutions, false resolutions and efficiency). Then, we describe in more details the

534    relationships between the variables of interest. As those relationships show an important

535    degree of complexity, we were not able to model them in a satisfactory manner. As a

536    consequence, we used linear models to describe and test the relative efficiency of the

537    methods and weighting schemes as well as other variables that can modify their relative

538    efficiencies.

539    *RESULTS AND DISCUSSION*

540    The two first axes of the PCA represent 73 % of the dataset variance. The first axis

541    represents mostly the contribution of TR and four positively correlated variables that are: the

542    constrict resolution, the constrict information content, the number *n* of taxa in subtrees, and

543    the number of taxa in the optimal tree (related to the more or less important overlapping of

544    the taxa sets of the subtrees; Fig. 5). For a given number of permutations, the amount of TR

545    compared to the number of taxa in subtrees has an S-shaped distribution (see SM 1a),

546    consequently TR tends to a horizontal asymptote (e.g. 100% of TR for 0 permutations when

547    the number of taxa in subtrees tends to the one in the original tree). As the number of

548    permutations increases, the asymptote decreases. The second axis of the PCA is mostly

549    structured by FR (Fig. 5). As expected, the number of permutations (i.e. the quantity of noise)

26

550    correlates mainly with this axis. In details, for a given number of permutations, the number of

551    FR firstly increases with the number of taxa in subtrees and then decreases (SM 1b).

552    Obviously, the efficiency (TR – FR) positively correlates to the first axis (thus to TR) and

553    negatively to the second axis (thus to FR) (Fig. 5, SM 1c). Finally, the number of optimal

554    trees does not correlate much with the other variables.

555    The relationships between all variables were thoroughly explored. The relative amount of

556    TR and FR depends also greatly on the number of taxa in the optimal tree (SM 2). For a given

557    number of taxa in the optimal tree, most of the runs fall on a line going from 0% TR and

558    100% FR to 0% FR and a percentage $p$ of TR. This percentage $p$ reaches 100% when the

559    number of taxa in the optimal tree equals the one in the original tree. These lines represent the

560    maximal resolution reachable with fully bifurcating optimal trees.

561    As concluded from the PCA, the variance of the dataset can be mostly summed up by TR,

562    FR, the number of taxa in subtrees, and the number of permutations. The relationships

563    between those variables are mostly monotonous. Thus, we computed two linear models with

564    respectively TR and FR as dependent variables and as independent variables for both models:

565    the number of taxa in subtrees, the number of permutations, and the different types of

566    analyses (MRP, 3ta with $FW_{cor.}$, $FW_{comp., or}$ MW/ $FW_{total}$). These models thus have three

567    variables and we included all the interactions terms. The homoscedasticity is not verified. For

568    example the variance of TR drastically increases in relation to the number of taxa in subtrees.

569    As a consequence, we do not expect from these models to test for the effects but rather to

570    investigate main tendencies and to observe specific interactions between the different types of

571    analyses and the other variables (summaries of linear models are provided in SM 3). We use

572    the function *interact_plot* from the R package *jtools* to visualize to three interactions as well

573 as main effects from isolated variables at the same time (Fig. 6). Linearity checks are
574 provided in the supplementary information (SM 4). The relationships are consistent with
575 what inferred from the PCA. The increase of TR with the number of taxa in subtrees depends
576 on the number of permutations. The higher the number of permutations is, the lower the
577 increase of TR in relation to the number of taxa in subtrees. The main results with the
578 substantial effect are that $FW_{cor.}$ and $FW_{comp.}$ are indistinct from each other and better
579 compared to MRP and MW as the number of taxa in subtrees increases. The increase of FR
580 with the number of permutations also depends on the number of taxa in subtrees and the type
581 of analysis. For a few taxa in subtrees, the analyses are indistinct no matter the number of
582 permutations. With both the increase of the number of taxa and permutations, differences
583 between analyses increase with the order of FR amount: MW, MRP, and indistinctively
584 $FW_{cor.}$ and $FW_{comp.}$ (Fig. 6). The efficiency decreases with the number of permutations. It also
585 increases with the number of taxa in subtrees, especially at a low number of taxa. The
586 numerous taxa in subtrees are, the higher the slope of this relationship is. In other terms, the
587 efficiency does not really increase with the number of taxa in subtrees for too noisy datasets.
588 This being said, $FW_{cor.}$ and $FW_{comp.}$ have better efficiencies compared to MRP and MW as
589 the number of taxa in subtrees increases (Fig. 6).

590 In conclusion, $FW_{cor.}$ and $FW_{comp.}$ provide more resolved trees in comparison to MRP and
591 MW. As the number of true resolutions (TR) is generally higher than the number of false
592 resolutions (FR), efficiency is a bit better for $FW_{cor.}$ and $FW_{comp}$. This leads us to recommend
593 these weighting methods for the purpose of supertree reconstruction methods, cladistics
594 methods, or metrics using triplets. These are very preliminary results and they need to be
595 confirmed in other conditions such as the type of noise, the tree shape. Future works will
596 have to study more precisely the differential effect of $FW_{cor.}$ versus $FW_{comp.}$, which was

597    indistinguishable in our simulations. The absence of significant differences between $FW_{cor.}$

598    and $FW_{comp.}$ may be related to the very low number of source trees. The higher the number of

599    source trees is, the more incongruent relationships there is. It is the incongruence that will

600    lead to the choice of this or that 3ts, and the choice should be biased with the $FW_{comp.}$ which

601    will weight some 3ts to more than 1, which is theoretically flawed.


602    CONCLUSION

603        Our paper focuses on the different types of relationships within 3ts sets, some of them

604    being directly responsible for the emergence of dependency. We also proposed a method for

605    evolutionary biologists whose aim is to completely eliminate these different types of

606    dependency from which all methods based on three-taxon statements currently suffer.

607        In-out nts, orthologous and paralogous relationships lead to the redundancy of some 3ts.

608    We propose for the first time a method of 3ts weighting based on these different types of

609    redundancy. Ideally, any phylogenetic or supertree analysis should only include datasets

610    whose elements are independent of each other, to eliminate any possibility of bias. We have

611    shown that the simple decomposition of a tree into its components leads to a loss of

612    information and that our method of corrected fractional weighting is the only procedure able

613    to fully delete all the redundancies. We have also shown using computer simulations that the

614    absence of weighting can significantly reduce the effectiveness of any method based on the

615    3ts. We also proposed the first empirical comparison between supertree component methods

616    (MRP) and 3ts (three-taxon analysis), confirming the efficiency of methods which

617    decompose trees into 3ts. More precisely, the most efficient methods are those which weight

618    3ts through a decomposition step into components, as Fractional Weighting. Finally,

619    understanding the dependency relationships between the 3ts will permit to improve the

29

620    accuracy and efficiency of supertree methods, consensus methods and phylogenetic analysis

621    methods based on 3ts. Triplet distance metrics should expect the same handling. To conclude,

622    we hope to open new debates and to fuel methodological discussions on how phylogenetic

623    information should be weighted in order to improve the accuracy of analyses.

634    LITERATURE CITED

635    Adams E.N. 1986. N-trees as nestings: Complexity, similarity, and consensus. Journal of

636        Classification. 3:299–317.

637    Aho A.V., Sagiv Y., Szymanski T.G., Ullman J.D. 1981. Inferring a tree from lowest

638        common ancestors with an application to the optimization of relational expressions.

639        SIAM Journal on Computing. 10:405–421.

640     Baum B. 1992. Combining Trees as a Way of Combining Data Sets for Phylogenetic

641             Inference, and the Desirability of Combining Gene Trees. Taxon. 41:3–10.

642     Brooks D.R. 1990. Parsimony Analysis in Historical Biogeography and Coevolution:

643             Methodological and Theoretical Update. Systematic Zoology. 39:14–30.

644     Bryant D. 2003. A classification of consensus methods for phylogenetics. DIMACS series in

645             discrete mathematics and theoretical computer science. 61:163–184.

646     Bryant D., Steel, M. 1995. Extension Operations on Sets of Leaf-Labelled Trees. Advances in

647             applied mathematics. 16:425–453.

648     Cao N., Bourdon E., El Azawi M., Zaragüeta R. 2009. Three-item analysis and parsimony,

649             intersection tree and strict consensus: a biogeographical example. Bulletins de la

650             société géologique de France. 180:13–15.

651     Dannenberg K., Jansson J., Lingas A., Lundell E.M. 2019. The approximability of maximum

652             rooted triplets consistency with fan triplets and forbidden triplets. Discrete Applied

653             Mathematics. 257:101–114.

654     Dekker M.C.H. 1986. Reconstruction methods for derivation trees. Unpublished Master

655             Thesis. Department of Mathematics and Computer Science, Vrije Universiteit,

656             Amsterdam.

657     Estabrook G.F., Johnson Jr.C.S., McMorris F.R. 1976. A mathematical foundation for the

658             analysis of cladistic character compatibility. Mathematical Biosciences. 29:181–187.

659     Farris J.S. 1970. Methods for Computing Wagner Trees. Systematic Zoology. 19:83–92.

31

660    Fitch W.M. 1970. Distinguishing homologous from analogous proteins. Systematic Zoology.

661        19:99–113.

662    Grand A., Corvez A., Duque Velez L.M., Laurin M. 2013. Phylogenetic inference using

663        discrete characters: performance of ordered and unordered parsimony and of three-

664        item statements. Biological Journal of the Linnean Society. 110:914–930.

665    Huerta-Cepas J., Serra F., Bork P. 2016. ETE 3: Reconstruction, Analysis, and Visualization

666        of Phylogenomic Data. Molecular Biology and Evolution. 33:1635–1638.

667    Islam M., Sarker K., Das T., Reaz R., Bayzid M. S. 2020. STELAR: A statistically consistent

668        coalescent-based species tree estimation method by maximizing triplet consistency.

669        BMC genomics, 21: 1–13.

670    Kitching I.J., Forey P.L., Humphries C.J., Williams D.M. 1998. Cladistics: The Theory and

671        Practice of Parsimony Analysis. Second edition. Oxford: Oxford University Press.

672    Kuhner M.K., Yamato J. 2015. Practical Performance of Tree Comparison Metrics.

673        Systematic Biology. 64:205–214.

674    Liu L., Yu L., Edwards S.V. 2010. A maximum pseudo-likelihood approach for estimating

675        species trees under the coalescent model. BMC Evolutionary Biology. 302:1–18.

676    Mavrodiev E.V., Williams D.M., Ebach M.C. 2019. On the Typology of Relations.

677        Evolutionary Biology. 46:71–89.

678    Mickevich M.F., Platnick N.I. 1989. On the information content of classifications. Cladistics.

679        5:33–47.

680    Nelson G. 1979. Cladistic analysis and synthesis: principles and definitions, with a historical

681        note on Adanson's Familles des Plantes (1763–1764). Systematic Biology, 28:1–21.

682    Nelson G., Ladiges P.Y. 1991a. Standard assumptions for biogeographic analysis. Australian

683        Systematic Botany. 4:41–58.

684    Nelson G., Ladiges P.Y. 1991b. Three-area statements: standard assumptions for

685        biogeographic analysis. Systematic Biology. 40:470–485.

686    Nelson G., Ladiges P.Y. 1992. Information content and fractional weight of three-item

687        statements. Systematic biology. 41:490–494.

688    Nelson G., Ladiges P.Y. 1994. Three-item consensus empirical test of fractional weighting.

689        Systematics Association Special Volume. 52: 193–209.

690    Nelson G., Platnick N.I. 1981. Systematics and Biogeography: Cladistics and Vicariance.

691        Columbia University Press. New York.

692    Nelson G., Platnick N.I. 1991. Three-Taxon Statements: A More Precise Use of Parsimony?

693        Cladistics. 7:351–366.

694    Penny D., Foulds L.R., Hendy M.D. 1982. Testing the theory of evolution by comparing

695        phylogenetic trees constructed from five different protein sequences. Nature.

696        297:197–200.

697    Poormohammadi H., Zarchi M. S., Ghaneai H. 2020. NCHB: A Method for Constructing

698        Rooted Phylogenetic Networks from Rooted Triplets based on Height Function and

699        Binarization. Journal of Theoretical Biology, 2–35.

33

700 Prin S. 2012. Structure mathématique des hypothèses cladistiques et conséquences pour la

701      phylogénie et l'évolution. Avec une perspective sur l'analyse cladistique.

702      Unpublished PhD Thesis.

703 R Development Core Team. 2008. R: A language and environment for statistical computing.

704      Vienna, Austria. R Foundation for Statistical Computing.

705 Ragan M. 1992. Phylogenetic Inference Based on Matrix Representation of Trees. Molecular

706      Phylogenetics and Evolution. 1:53–58.

707 Ranwez V., Criscuolo A., Douzery E.J.P. 2010. SUPERTRIPLETS: a triplet-based supertree

708      approach to phylogenomics. Bioinformatics. 26:i115–i123.

709 Rineau V., Zaragüeta R., Laurin M. 2018. Impact of errors on cladistic inference: simulation-

710      based comparison between parsimony and three-taxon analysis. Contributions to

711      Zoology. 87:25–40.

712 Sevillya G., Frenkel Z., Snir S. 2016. Triplet MaxCut: a new toolkit for rooted supertree.

713      Methods in Ecology and Evolution. 7:1359–1365.

714 Strimmer K., von Haeseler A. 1996. Quartet Puzzling: A Quartet Maximum-Likelihood

715      Method for Reconstructing Tree Topologies. Molecular Biology and Evolution.

716      13:964–969.

717 Swofford D.L. 2003. PAUP*. Phylogenetic Analysis Using Parsimony (*and Other

718      Methods). Version 4. Sunderland, Massachusetts: Sinauer Associates.

719   Tavares B.L. 2018. A synopsis of comparative metrics for classifications. arXiv preprint

720        arXiv:1804.03929.:37.

721   Vach W. 1994. Preserving consensus hierarchies. Journal of Classification. 11:59–77.

722   Wiley E.O. 1986. Methods in vicariance biogeography. In: Hovenkamp P., editor.

723        Systematics and Evolution: a matter of Diversity. Utrecht. p. 283–306.

724   Wilkinson M. 1994a. Common Cladistic Information and its Consensus Representation:

725        Reduced Adams and Reduced Cladistic Consensus Trees and Profiles. Systematic

726        Biology. 43:343–368.

727   Wilkinson M. 1994b. Three-taxon statements: when is a parsimony analysis also a clique

728        analysis? Cladistics. 10:221–223.

729   Wilkinson M., Cotton J., Thorley J. 2004. The Information Content of Trees and Their Matrix

730        Representations. Systematic Biology. 53:989–1001.

731   Williams D.M. 2004. Supertrees, components and three-item data. In: Bininda-Emonds O.,

732        editor. Phylogenetic supertrees: combining information to reveal the tree of life.

733        Springer Science & Business Media. p. 389–408.

734   Williams D.M., Ebach M.C. 2008. Foundations of systematics and biogeography. Springer

735        Science & Business Media.

736   Williams D.M., Humphries C.J. 2003. Component Coding, Three-item Coding, and

737        Consensus Methods. Sys. Biol. 52:255–259.

738 Zaragüeta R., Lelièvre H., Tassy P. 2004. Temporal paralogy, cladograms, and the quality of

739         the fossil record. Geodiversitas. 26:381–389.

740 Zaragüeta R., Ung V., Grand A., Vignes-Lebbe R., Cao N., Ducasse J. 2012. LisBeth: New

741         cladistics for phylogenetics and biogeography. Comptes Rendus Palevol. 11:563–566.

742 CAPTIONS



744 Figure 1.

745 Four possible types of relationships between pairs of 3ts can be combined. Two 3ts

746 sharing two taxa in common can be combined into a 4-taxon tree from which it is possible to

747 deduct 0, 1 or 2 3ts. The black dots correspond to the leaves common to both 3ts. a. in-in nts

748 relationship. b. in-out nts relationship. c. ortholog relationship. d. paralog relationship.

749

Rineau, Zaragüeta and Bardin



Figure 2.

Fractional weighting procedure with correction of inter-node redundancy on $(a(b(cd)))$.

The dichotomous pectinate tree only produces 3ts of weight 1.

INFORMATION CONTENT OF TREES



Figure 3.

Fractional weighting procedure on $(a(bc(de)))$. The tree is not completely resolved, which

leads to uncertainty amongst 3ts and the persistence of fractional weightings for five of them.

759

760    Figure 4.

761    Flow-chart describing the steps for each run of the simulations: generation of an original

762    tree, sampling of two subtrees with $n$ taxa, pruning and regrafting of taxa, supertree analysis,

763    comparison of the optimal tree(s) with the original tree in order to calculate $TR$, $FR$, and

764    efficiency.

Figure 5.

Variables factor map of the PCA for the two first principal components (*PCA* function from *FactoMineR* R package). The first axis represents 51.34% of the variance and correlates positively with *TR* (true resolutions), numbers of taxa in subtrees and optimal tree, strict consensus (sc) information content and resolution. The second axis represents 21.02% of the variance and correlates positively with FR (false resolutions) and the number of permutations. The variance of the efficiency is split between the first and second axis.

41

Rineau, Zaragüeta and Bardin



773

774

42

775     Figure 6.

776     Interaction plots illustrating three linear models all including as independent variables the

777     number of taxa in subtrees, the number of permutations and the types of analysis (function

778     *interact_plot* from *jtools* R package. a. TR (number of true resolutions) as dependent

779     variable. b. FR (number of false resolutions) as dependent variable. c. Efficiency as

780     dependent variable. Types of lines represent the types of analysis, solid line: $FW_{cor.}$, dashed

781     line: $FW_{comp.}$, pointed line: MW, dashed and pointed line: MRP.

782     Table I.

783     Decomposition of the tree ($a(bc(de))$) into 3ts and calculation of their weights using

784     different methods.

| 3ts | Fractionnal weighting proposed herein | Nelson and Ladiges (1992) total fractionnal weighting | Nelson and Ladiges (1992) fractionnal weighting per component | Uniform weighting | Minimal weighting (Wilkinson et al. 2004) |
|---|---|---|---|---|---|
| (d,e),c | 1 | 6/8 | 1 | 1 | 3/8 |
| (d,e),b | 1 | 6/8 | 1 | 1 | 3/8 |
| (d,e),a | 1 | 6/8 | 3/2 (1+1/2) | 2 | 3/8 |
| (b,c),a | 6/10 (=1/2+1/10) | 6/8 | 1/2 | 1 | 3/8 |
| (b,d),a | 6/10 (=1/2+1/10) | 6/8 | 1/2 | 1 | 3/8 |
| (b,e),a | 6/10 (=1/2+1/10) | 6/8 | 1/2 | 1 | 3/8 |
| (c,d),a | 6/10 (=1/2+1/10) | 6/8 | 1/2 | 1 | 3/8 |
| (c,e),a | 6/10 (=1/2+1/10) | 6/8 | 1/2 | 1 | 3/8 |
| Total | 6 | 6 | 6 | 9 | 3 |

785

786     Table II.

787     Pseudocode detailing the steps of removing the dependency of a set of 3ts from a tree by

788     fractional weighting.

789

43

Pseudocode. Corrected_Fractional_weighting(tree)

| | | |
|---|---|---|
| 1 | **if** tree as informative nodes **then:** | |
| 2 | **for** each informative node: | Nelson and Ladiges (1992) fractionnal weighting per component |
| 3 | generate the component associated to the informative node | |
| 4 | calculate the number of leaves (n) branched to the informative node of the component | |
| 5 | decompose the component in all its implied 3ts | |
| 6 | attach a weight of 2/n to each 3ts | |
| 7 | **if** number of informative nodes > 1: | |
| 8 | **for** each couple of informative nodes directly included into each other from least inclusive to most inclusive nodes: | Weighting correction proposed herein |
| 9 | calculate the sum of weights (t) of the 3ts from the least inclusive node that are common with the most inclusive node | |
| 10 | suppress the 3ts from the least inclusive node that are common with the most inclusive node | |
| 11 | calculate the number of remaining 3ts (i) | |
| 12 | **for** each remaining 3ts: | |
| 13 | Add t/i to the weight of the 3ts | |
| 14 | **return** the computed 3is with their associated fractionnal weighting | |
| 15 | **else** | |
| 16 | **return** | |

790

791    Appendix 1

792       a. Analysis of the different types of relationships between 3ts from a pectinate tree

793    ($a$($b$($cd$))) and its components ($ab$($cd$)) and ($a$($b$($cd$)) taken independently. b. Analysis of the

794    different types of relationships between 3ts from a symmetric tree (($ab$)($cd$)) and its

795    components ($ab$($cd$)) and (($ab$)$cd$) taken independently. c. Analysis of the different types of

796    relationships between 3ts from a mixed tree ((($ab$)($cd$))$e$) and its components (($ab$)$cde$),

797    ($ab$($cd$)$e$), and (($abcd$)$e$) taken independently. The different types of relationships are in-in,

798    in-out, orthologous (o), paralogous (p). An empty cell means 3ts are compatible but not

799    combinable.

INFORMATION CONTENT OF TREES

800        a.

Component 1:
(a,b,(c,d))

|  | b(c,d) | a(c,d) |
|---|---|---|
| b(c,d) | . | in-in |
| a(c,d) |  | . |

801

Component 2: (a,(b,c,d))

|  | a(b,d) | a(c,d) | a(b,c) |
|---|---|---|---|
| a(b,d) | . | in-out | in-out |
| a(c,d) |  | . | in-out |
| a(b,c) |  |  | . |

802

Tree: (a,(b,(c,d))) (component 1 + 2)

|  | b(c,d) | a(b,d) | a(c,d) | a(b,c) |
|---|---|---|---|---|
| b(c,d) | . | o | in-in | o |
| a(b,d) |  | . | in-out | in-out |
| a(c,d) |  |  | . | in-out |
| a(b,c) |  |  |  | . |

803

804

45

805      b.

Component 1:
(a,b,(c,d))

|        | a(c,d) | b(c,d) |
|--------|--------|--------|
| a(c,d) | .      | in-in  |
| b(c,d) |        | .      |

806

Component 2:
((a,b),c,d)

|        | c(a,b) | d(a,b) |
|--------|--------|--------|
| c(a,b) | .      | in-in  |
| d(a,b) |        | .      |

807

Tree: ((a,b),(c,d)) (component 1 + 2)

|        | c(a,b) | a(c,d) | d(a,b) | b(c,d) |
|--------|--------|--------|--------|--------|
| c(a,b) | .      | p      | in-in  | p      |
| a(c,d) |        | .      | p      | in-in  |
| d(a,b) |        |        | .      | p      |
| b(c,d) |        |        |        | .      |

808

809

46

INFORMATION CONTENT OF TREES

810      c.

Component 1: ((a,b),c,d,e)

|        | c(a,b) | d(a,b) | e(a,b) |
|--------|--------|--------|--------|
| c(a,b) | .      | in-in  | in-in  |
| d(a,b) |        | .      | in-in  |
| e(a,b) |        |        | .      |

811

Component 2: (a,b,(c,d),e)

|        | a(c,d) | b(c,d) | e(c,d) |
|--------|--------|--------|--------|
| a(c,d) | .      | in-in  | in-in  |
| b(c,d) |        | .      | in-in  |
| e(c,d) |        |        | .      |

812

Component 3: ((a,b,c,d),e)

|        | a(b,c) | a(c,d) | a(d,e) | a,(b,d) | a,(b,e) | a,(c,e) |
|--------|--------|--------|--------|---------|---------|---------|
| a(b,c) | .      | in-out |        | in-out  | in-out  | in-out  |
| a(c,d) |        | .      | in-out | in-out  |         | in-out  |
| a(d,e) |        |        | .      | in-out  | in-out  | in-out  |
| a,(b,d)|        |        |        | .       | in-out  |         |
| a,(b,e)|        |        |        |         | .       | in-out  |
| a,(c,e)|        |        |        |         |         | .       |

813

Tree: (((a,b),(c,d)),e) (component 1 + 2 + 3)

|        | c(a,b) | d(a,b) | e(a,b) | a(c,d) | b(c,d) | e(c,d) | e(a,c)  | e(a,d)  | e(b,c)  | e(b,d)  |
|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|---------|
| c(a,b) | .      | in-in  | in-in  | p      | p      |        | o       |         | o       |         |
| d(a,b) |        | .      | in-in  | p      | p      |        |         | o       |         | o       |
| e(a,b) |        |        | .      |        |        |        | in-out  | in-out  | in-out  | in-out  |
| a(c,d) |        |        |        | .      | in-in  | in-in  | o       | o       |         |         |
| b(c,d) |        |        |        |        | .      | in-in  |         |         | o       | o       |
| e(c,d) |        |        |        |        |        | .      | in-out  | in-out  | in-out  | in-out  |
| e(a,c) |        |        |        |        |        |        | .       | in-out  | in-out  |         |
| e(a,d) |        |        |        |        |        |        |         | .       |         | in-out  |
| e(b,c) |        |        |        |        |        |        |         |         | .       | in-out  |
| e(b,d) |        |        |        |        |        |        |         |         |         | .       |

814

47

815    SUPPLEMENTARY MATERIAL

816      SM 1.

817      Boxplots of true resolutions (TR), false resolutions (FR) and efficiency in regard to the

818    number of taxa in subtrees, split by the type of analysis (MRP: matrix representation with

819    parsimony, $FW_{cor.}$: 3ta with corrected fractional weighting, MW: 3ta with minimal weighting,

820    $FW_{comp.}$: 3ta with fractional weighting per component) and the number of permutations.

821      SM 2.

822      Movie of a three dimensions block diagram with the percentages of true and false

823    resolutions (TR and FR), the number of taxa in the optimal tree (also highlighted by colors).

824    Darkening of each color depends on the resolutions of the optimal tree, dark corresponds to

825    poorly resolved trees and bright corresponds to well-resolved tree.

826      SM 3.

827      Summaries of the three linear models containing formulas of the models (bold), the

828    estimates of each coefficient as well as standard error, t-values, and p-values.

829      SM 4.

830      Loess smoothed lines indicating a deviation from linearity effect from the three linear

831    models.