1     **Significant phylogenetic signal is not enough to trust phylogenetic predictions**

2     Rafael Molina-Venegas[1,2*], Ignacio Morales-Castilla[2] and Miguel Á. Rodríguez[2]

3

4     1. *Department of Ecology, Faculty of Sciences, Universidad Autónoma de Madrid, Madrid, Spain.*

5     2. *Universidad de Alcalá, Department of Life Sciences, Global Change Ecology and Evolution*

6     *Group, Alcalá de Henares, Madrid, Spain.*

7     * Correspondence: rafmolven@gmail.com

8

10

11     **Abstract**

12     In a recent study, Cantwell-Jones et al. (2022) proposed a list of 1044 species as

13     promising key sources of B vitamins based primarily on phylogenetic predictions. To

14     identify candidate plants, they fitted lambda models of evolution to edible species with

15     known values in each of six B vitamins (232 to 280 species) and used the estimated

16     parameters to predict B-vitamin profiles of edible plants lacking nutritional data (6460

17     to 6508 species). The latter species were defined as potential sources of a B vitamin if

18     the predicted vitamin content was ≥15% towards recommended dietary allowances for

19     active females between 31-50 years per 100 g of fresh edible plant material consumed.

20     Unfortunately, the reliability of the predictions that informed the list of candidate

21     species is questionable due to insufficient phylogenetic signal in the data (Pagel's $\lambda$

22     between 0.171 and 0.665) and a high incidence of species with missing values (over

23     95% of all the species analyzed in the study). We found that of the 1044 species

24     proposed as promising B-vitamin sources, 626 to 993 species showed accuracies that

25     were indistinguishable from those obtained under a white noise model of evolution (i.e.

26     random predictions conducted in absence of any phylogenetic structure) in at least one

27    of the vitamins, which proves the weakness of the inference drawn from imputed

28    information in the original study. We hope this commentary serves as a cautionary note

29    for future phylogenetic imputation exercises to carefully assess whether the data meet

30    the requirements for the predictions to be valuable, or at least more accurate than

31    expected by chance.

32

33    **Main**

34    In a recent study, Cantwell-Jones et al.[1] proposed a list of 1044 species as promising

35    key sources of B vitamins based primarily on phylogenetic predictions. To identify

36    candidate plants, they fitted lambda models of evolution to edible species with known

37    values in each of six B vitamins and used the estimated parameters to predict B-vitamin

38    profiles of edible plants lacking nutritional data. The latter species were defined as

39    'sources' of a B vitamin if the predicted vitamin content was $\geq 15\%$ towards

40    recommended dietary allowances for active females between 31–50 years per 100 g of

41    fresh edible plant material consumed. This phylogenetic imputation exercise is exciting

42    and inspiring, more so as it would have the potential to be replicated for further goals

43    such as, for example, identifying species of pharmaceutical interest. Unfortunately, the

44    predictions from this study are questionable. The predictive capability of lambda

45    models is primarily determined by the amount of phylogenetic signal in the known data

46    (i.e. the extent to which closely related species share similar values in the trait of

47    interest), which means that predictions based on low phylogenetic signals, even if

48    statistically significant, are typically valueless[2-4]. Further, simulations have shown that

49    even under 'strong' phylogenetic signal (i.e. Pagel's $\lambda = 1$), acceptable predictions can

50    only be expected when the most recent common ancestor (MRCA) of a target species

51    and its closest relative with known trait value (hereafter 'predictive MRCA') is

52    relatively recent[4]. As such, the older the predictive MRCAs, the lesser the difference

53    between phylogenetic imputations and random predictions (Fig. 1). Based on 232 to 280

54    nutritionally known species (observed data), Cantwell-Jones et al.[1] predicted B-vitamin

55    profiles for 6460 to 6508 nutritionally unknown species (over 95% of all the species

56    analyzed in the study), and they did so relying on very weak-to-moderate phylogenetic

57    signals in the observed data (Pagel's $\lambda$ between 0.171 and 0.665). These figures suggest

58    that the predictions conducted by the authors are unreliable due to both insufficient

59    phylogenetic signal and a high incidence of relatively old predictive MRCAs (Fig. 1).

60    Moreover, the unreliability of their predictions calls into question the proposed list of

61    1044 species as promising sources of B vitamins.

62         Cantwell-Jones et al.[1] conducted leave-one-species-out cross-validation trials on

63    the observed values of the B vitamins they analyzed (they referred to this procedure as

64    "jackknifing") and estimated 95% confidence intervals (CI) around each re-estimated

65    (predicted) value. They found that ≥91.4% of nutritionally known species had measured

66    (observed) values within the 95% CI of the values predicted in the trials. Here, we

67    conducted the same analysis but randomizing the authors' dataset in a two-step

68    procedure. First, the observed values of each vitamin were reshuffled across the species

69    with known value in the vitamins 500 times, and then each reshuffled set (n = 500 sets

70    per vitamin) was checked to show complete lack of phylogenetic signal (i.e. $\lambda < 0.001$

71    and $p > 0.05$). Otherwise, the values were reshuffled iteratively until the conditions

72    were met. We found that ≥92.8% of the observed values in the randomized sets sit

73    within 95% CIs of their corresponding predicted values (Supplementary Data 1-5). This

74    demonstrates that the authors simply found the null expectation of the analysis, which

75    deems the CI-based evidence of their study as invalid.

76    On the other hand, the authors used generalised least squares models (with a

77    variance structure in the error term) to test the strength of the relationship between

78    predicted and observed values, and they found significant relationships for all the

79    vitamins. Here, we used a simple and more intuitive prediction coefficient ($P^2$) to assess

80    the overall predictive power of the lambda models they employed:

81    $$P^2 = 1 - \frac{\sum_{i=1}^{n}(\hat{y}_i - y_i)^2/n}{s_y^2}$$

82    where $\hat{y}_i$ and $y_i$ are respectively the predicted and observed value for the nutritionally

83    known species $i$ and $s_y^2$ is the sample variance of the observed values of the trait[4-6]. $P^2$

84    can be interpreted similarly as Ezekiels' adjusted coefficient of determination[7], so that

85    $P^2 = 1$ when all predictions perfectly match the observations and $P^2 = 0$ when the model

86    is no better in predicting values than simply taking the mean of the observed values

87    (note that $P^2$ has no negative boundary, as there is no theoretical limit to how badly a

88    model can predict observed trait values). We found that $P^2$ was close to zero in all cases

89    except for folate, which showed a slightly higher (yet weak) score (Table 1). This is not

90    surprising, as previous simulations have shown that phylogenetic predictions perform

91    very poorly when $\lambda$ is lower than ~0.6[4]. There are numerous studies that caution against

92    using phylogenies alone to predict missing data under low phylogenetic signal[2-4,8,9], yet

93    none of them was acknowledged in Cantwell-Jones et al.[1]

94    In line with the authors' observation that "median differences between predicted

95    and observed values for each nutrient were <33% of the standard deviation across

96    species", it could be argued that some of the predictions may still be useful despite the

97    discouraging results from the leave-one-species-out cross-validation trials (Table 1).

98    Thus, we used the method described in Molina-Venegas et al.[4] to assess the expected

99    species-level accuracy of the predictions that informed the list of candidate plants (see

100   Supplementary Information for details). Depending on whether *p*-values were

101 Bonferroni-corrected ($n = 1616$, two-sample Wilcoxon tests) and the nominal alpha

102 criterion (5% or 0.1%), we found that of the 1044 species proposed as promising B-

103 vitamin sources, 626 to 993 species showed accuracies that were indistinguishable from

104 those obtained under a 'white noise' model of evolution (i.e. random predictions

105 conducted in absence of any phylogenetic structure) in at least one of the vitamins

106 (Supplementary Data 6). Moreover, regarding the species that showed statistically

107 significant accuracies, it should not be concluded that their predictions are necessarily

108 valuable, only that they are better than drawing values at random independently of the

109 phylogeny. This is illustrated by the rather modest maximum accuracy that was

110 recorded across all candidate species and B vitamins analyzed, which corresponded to

111 *Hordeum bulbosum* and niacin with only 54.4% of $P^2_{sim}$ values $\geq 0.75$ (Supplementary

112 Data 7). These results are not surprising either. The extremely high incidence of missing

113 values in the dataset makes predictive MRCAs prone to be relatively old (Fig. 1), which

114 has a proven negative impact on prediction accuracy[4,8]. As such, simulations show that

115 only species whose predictive MRCA is younger than ~10% of the height of the

116 phylogeny are expected to show consistently accurate predictions (i.e. at least 75% of

117 $P^2_{sim}$ values $\geq 0.75$) under a scenario of 'strong' phylogenetic signal (i.e. $\lambda = 1$)[4].

118 It is very exciting to see a burgeoning interest in connecting phylogenetic

119 information with human well-being, but researchers should be clear on the limitations of

120 phylogenetic predictive methods and utilize imputed information with caution and

121 restraint, especially if the goal is employing individual predictions separately (e.g. to

122 evaluate if a species has potential to be 'source' of a B vitamin, as in Cantwell-Jones et

123 al.[1]). The dataset gathered by Cantwell-Jones et al.[1] is impressive, and it would not be

124 surprising if the 1044 species they proposed as candidate sources attract the attention of

125 professionals willing to invest resources in measuring their B-vitamin content.

5

126     Unfortunately, the reliability of the predictions that informed the list of candidate

127     species ranges between weak to very weak. While the authors did not intend to predict

128     the exact nutrient content of species but to evaluate if the predicted values were above

129     or below a certain threshold, the fact that most of their predictions are indistinguishable

130     from pure randomness proves the weakness of the inference drawn from imputed

131     information in their study. We hope this commentary serves as a cautionary note for

132     future phylogenetic imputation exercises to carefully assess whether the data meet the

133     requirements for the predictions to be valuable, or at least more accurate than expected

134     by chance.

135

136     **Acknowledgements**

142

143     **Author contributions**

144     R.M.-V. conceived and discussed the idea with M.A.R. and I.M.-C., performed the

145     calculations and led the writing, I.M.-C. drew the figure, and all the authors contributed

146     to the writing.

147

148     **Competing Interests**

149     The authors declare no competing interests

150

151    **Additional information**

152    **Supplementary information.** The online version contains supplementary material

153    available at XXX

154

155    **Data availability.** All the datasets generated for this study are available from the

156    Figshare Digital Repository at https://doi.org/10.6084/m9.figshare.19682880.v1.

157

158    **Correspondence** should be addressed to R.M.-V.

159

160    **Table 1.** Phylogenetic signal λ of each B vitamin analyzed across $N$ nutritionally known

161    species and prediction coefficients $P^2$ derived from the leave-one-species-out cross-

162    validation trials.

163

| B vitamin | $N$ | λ | $P^2$ |
|---|---|---|---|
| Thiamine | 280 | 0.293 | -0.010 |
| Riboflavin | 277 | 0.192 | 0.035 |
| Niacin | 278 | 0.665 | 0.086 |
| Pantothenic acid | 232 | 0.171 | 0.023 |
| Folate | 256 | 0.302 | 0.183 |

164

165    **Figure 1.** Conceptual model on the expected accuracy of phylogenetic predictions as a

166    function of phylogenetic signal and distance to predictive MRCAs, this is, the MRCA

167    of each target species and its closest relative with known trait value. Accuracy is

168    expected to be acceptable only under strong phylogenetic signal and relatively recent

169    predictive MRCAs (short distances), and predictions will be uncertain if any of these

170    conditions are not met. The probability of finding relatively old predictive MRCAs

171    (long distances) increases with the amount of missing data, which may lead to the

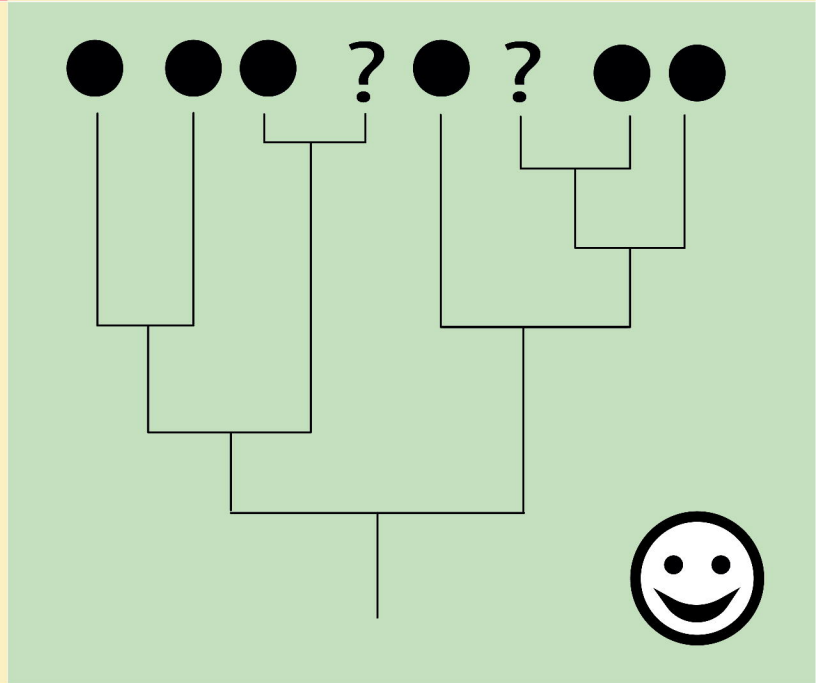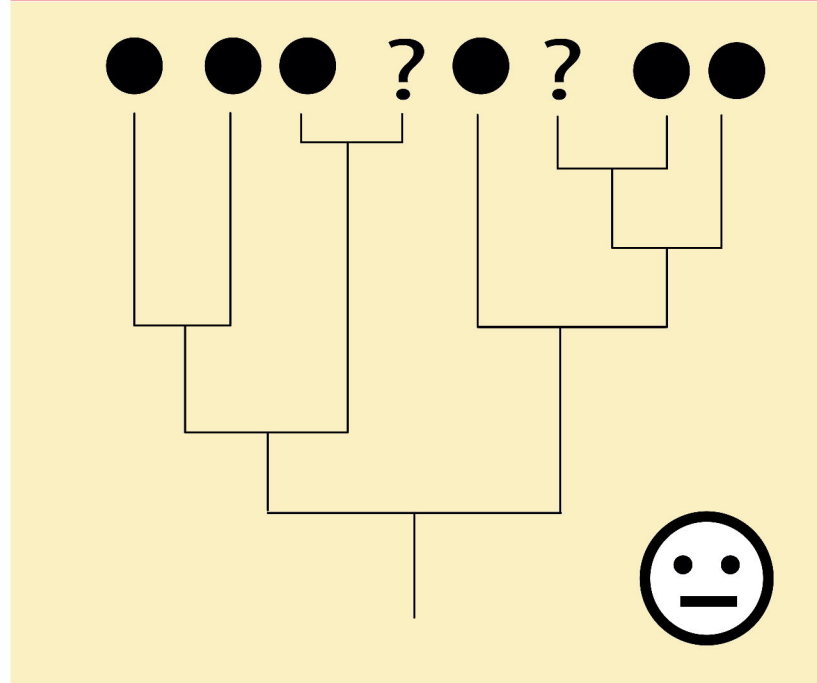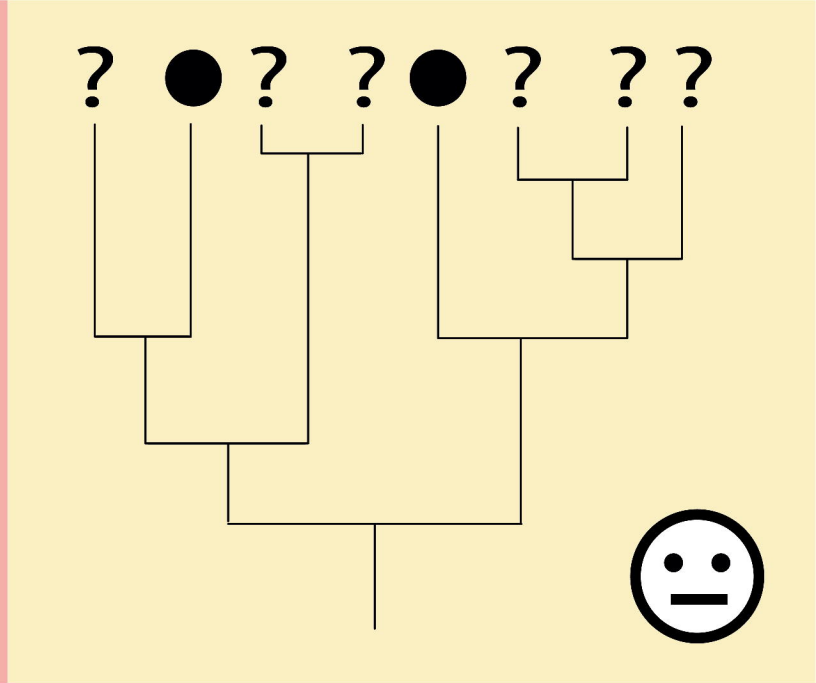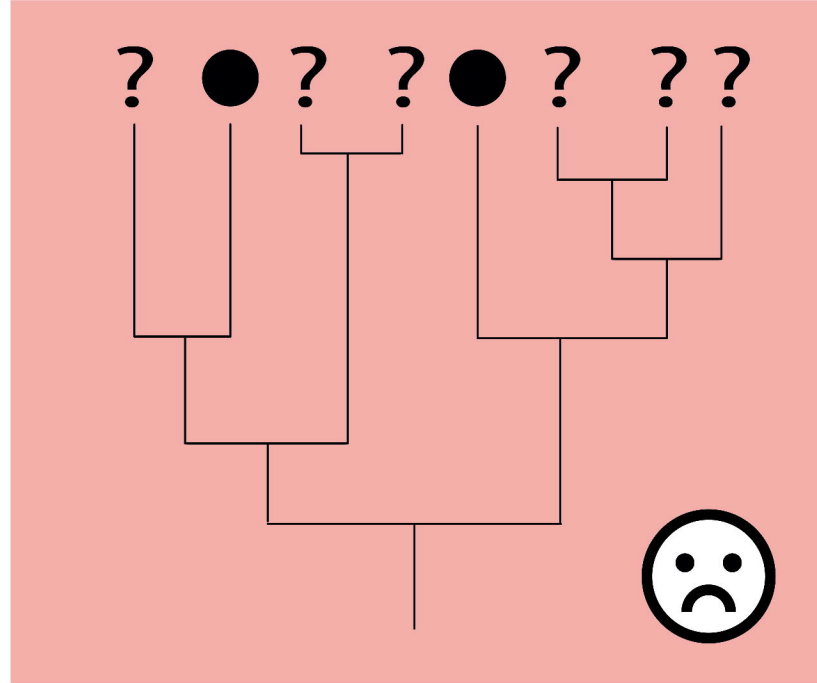172    worst-case scenario if phylogenetic signal is weak.

173

174    **References**

175    1. Cantwell-Jones, A. *et al.* Global plant diversity as a reservoir of micronutrients for

176       humanity. *Nature Plants* 1–8 (2022).

177    2. Swenson, N. G. Phylogenetic imputation of plant functional trait databases.

178       *Ecography* **37**, 105–110 (2014).

179    3. Swenson, N. G. *et al.* Phylogeny and the prediction of tree functional diversity across

180       novel continental settings. *Glob Ecol Biogeogr* **26**, 553–562 (2017).

181    4. Molina-Venegas, R. *et al.* Assessing among-lineage variability in phylogenetic

182       imputation of functional trait datasets. *Ecography* **41**, 1740–749 (2018).

183    5. Guénard, G., Legendre, P. & Peres-Neto, P. Phylogenetic eigenvector maps: a

184       framework to model and predict species traits. *Methods Ecol Evol* **4**, 1120–1131

185       (2013).

186    6. Guénard, G., Ohe, P. C. von der, Walker, S. C., Lek, S. & Legendre, P. Using

187       phylogenetic information and chemical properties to predict species tolerances to

188       pesticides. *Proc R Soc B Biol Sci* **281**, 20133239 (2014).

189    7. Ezekiel, M. *Methods of correlation analysis*. John Wiley and Sons, New York, USA

190       (1930).

191    8. Johnson, T. F., Isaac, N. J. B., Paviolo, A. & González-Suárez, M. Handling missing

192       values in trait data. *Glob Ecol Biogeogr* **30**, 51–62 (2021).

193    9. Debastiani, V. J., Bastazini, V. A. G. & Pillar, V. D. Using phylogenetic information

194       to impute missing functional trait values in ecological databases. *Ecol Inform* **63**,

195       101315 (2021).