

21 Abstract

22 The problem of pattern and scale is a central challenge in ecology. The problem
23 of scale is central to community ecology, where functional ecological groups are
24 aggregated and treated as a unit underlying an ecological pattern, such as ag-
25 gregation of “nitrogen fixing trees” into a total abundance of a trait underlying
26 ecosystem physiology. With the emergence of massive community ecological
27 datasets, from microbiomes to breeding bird surveys, there is a need to objec-
28 tively identify the scales of organization pertaining to well-defined patterns in
29 community ecological data.

30 The phylogeny is a scaffold for identifying key phylogenetic scales associ-
31 ated with macroscopic patterns. Phylofactorization was developed to objec-
32 tively identify phylogenetic scales underlying patterns in relative abundance
33 data. However, many ecological data, such as presence-absences and counts,
34 are not relative abundances, yet it is still desirable and informative to identify
35 phylogenetic scales underlying a pattern of interest. Here, we generalize phylo-
36 factorization beyond relative abundances to a graph-partitioning algorithm for
37 any community ecological data.

38 Generalizing phylofactorization connects many tools from data analysis to
39 phylogenetically-informed analysis of community ecological data. Two-sample
40 tests identify three phylogenetic factors of mammalian body mass which arose
41 during the K-Pg extinction event, consistent with other analyses of mammalian
42 body mass evolution. Projection of data onto coordinates defined by the phy-
43 logeny yield a phylogenetic principal components analysis which refines our un-
44 derstanding of the major sources of variation in the human gut microbiome.
45 These same coordinates allow generalized additive modeling of microbes in Cen-
46 tral Park soils and confirm that a large clade of Acidobacteria thrive in neutral
47 soils. Generalized linear and additive modeling of exponential family random

48 variables can be performed by phylogenetically-constrained reduced-rank regres-
49 sion or stepwise factor contrasts. We finish with a discussion of how phylofac-
50 torization produces an ecological species concept with a phylogenetic constraint.
51 All of these tools can be implemented with a new R package available online.

52 Keywords

53 Phylofactorization, phylogeny, microbiome, ecological data, big data, graph par-
54 titioning, dimensionality reduction

55 Introduction

56 The problem of pattern and scale is a central problem in ecology [27]. Ecological
57 patterns of interest, such as ecosystem physiology, species abundance distribu-
58 tions, epidemics, ecosystem services of animal-associated microbial communi-
59 ties, and more, are often the result of processes that operate at multiple scales.
60 Traditionally, the “scales” of interest are space, time, and levels of ecological
61 organization ranging from individuals to populations to ecosystems. Predic-
62 tion of spatial variation over different scales, millimeters, meters, or kilometers,
63 requires incorporation of different processes driving patterns observed. Pre-
64 dicting climatic and weather patterns over days, years, or millennia requires
65 different data, processes and models. Predicting the collective behavior of a
66 school of fish requires interfacing individual behavior with interaction networks
67 of those individuals [25] whereas predicting the ability of a forest to act as a car-
68 bon sink requires interfacing weather, nutrient cycles, and competition between
69 trees with different traits, such as nitrogen fixation [11]. Understanding emer-

70 gent infectious diseases requires interfacing processes over scales ranging from
71 animal population dynamics, reservoir epizootiology, and human epidemiology
72 [37]. Ecological theory requires interfacing phenomena across scales believed
73 to be important, and continually updating our beliefs about which scales are
74 important to interface.

75 For a novel or unfamiliar pattern, such as a change in microbial community
76 composition along environmental gradients, how can one objectively identify
77 the appropriate scales of ecological organization? In macroscopic systems, a
78 researcher will use intuition derived from natural history knowledge to determine
79 scales of interest. Models of how the presumably important natural history traits
80 affect the pattern will be constructed, and the goodness of fit to the pattern of
81 interest will be used as a metric for the successful identification of ecological
82 scales/traits. However, for some patterns, such as the ecosystem physiology of
83 the human microbiome, there is limited natural history knowledge to draw on to
84 assist the decision of the appropriate scales of interest. There is a need for rules,
85 algorithms and laws for the simplification, aggregation, and scaling of ecological
86 phenomena.

87 A central feature of biological systems is the existence of a hierarchical as-
88 semblage of entities, from genes to species, whose relationships and evolutionary
89 history can be estimated and organized into a hierarchical tree. The estimated
90 phylogeny contains edges along which mutations occur and new traits arise.
91 When the phylogeny correctly captures the evolution of discrete, functional eco-
92 logical traits underlying a pattern of interest, the phylogeny is a natural scaffold
93 for simplification, aggregation, and scaling in ecological systems. Patterns such
94 as the change of bacterial abundances following antibiotic exposure, whose func-
95 tional ecological traits of antibiotic resistance are laterally transferred, can still
96 be simplified by constructing a phylogeny of the laterally transferred genes, such

97 as the beta-lactamases[18], as a natural scaffold for defining the entities with
98 different responses to antibiotics.

99 The phylogeny contains a hierarchy of possible scales for aggregation. Gra-
100 ham et al. [17] develop the term “phylogenetic scale” to refer to the depth of the
101 tree over which we aggregate information from a clade. Functional ecological
102 traits often arise at different depths of the tree. Many ecological phenomena
103 may be driven by traits not properly summarized or aggregated by moving the
104 phylogeny along a constant depth. Instead, there may be multiple phylogenetic
105 scales, or grains, underlying an ecological pattern of interest. For example, the
106 patterns of vertebrate abundances on land and water are simplified by nested
107 clades: Tetrapods, Cetaceans, Pinnipeds, etc. There is a need for general sta-
108 tistical methods to partition the phylogeny into the grains with significantly
109 different associations or contributions to the ecological pattern. Such a method
110 can objectively identify the phylogenetic scales underlying an ecological pattern
111 of interest.

112 Phylofactorization [51] was developed to identify the phylogenetic scales in
113 compositional (relative abundance) data by iteratively constructing variables
114 corresponding to edges in the phylogeny and selecting variables which maxi-
115 mize an objective function. The variables used were a common transform from
116 compositional data analysis [1], referred to as the isometric log-ratio transform
117 [10, 9], which contrast the relative abundances of species separated by an edge
118 in the phylogeny. A coordinate in an isometric log-ratio transform aggregates
119 relative abundances within clades by a geometric mean and contrasts clades
120 through log-ratios of the clades’ geometric mean relative abundances. The
121 isometric log-ratio transform also allows the construction of non-overlapping
122 contrasts, thereby reducing an obvious source of dependence in phylogenetic
123 variables. The isometric log-ratio transform is used to identify phylogenetic

124 scales capturing large blocks of variation in relative-abundance data and con-
125 struct coordinates that correspond to edges along which hypothesized functional
126 ecological traits arose.

127 However, many ecological data are more appropriately viewed as counts, not
128 compositions. For example, the presence/absence of bird species across conti-
129 nents are best modelled as Bernoulli random variables, not compositional data.
130 In this paper, we extend phylofactorization to broader classes of data types
131 by generalizing the logic of phylofactorization and to a set of three operations:
132 aggregation, contrast, and an objective function defined by the pattern of in-
133 terest. The nested dependence of clades within clades is avoided by defining
134 phylofactorization as a graph-partitioning algorithm that contrasts species sep-
135 arated by edges and iteratively partitioning the phylogeny along edges that best
136 differentiate species.

137 After defining phylofactorization as a graph-partitioning algorithm, we il-
138 lustrate the generality of the algorithm through several examples. First, we
139 show that two-sample tests, such as t-tests and Fisher's exact test, are natural
140 operations for phylofactorization - they first aggregate data from two groups
141 through means, contrast the aggregates via a difference of means, and have nat-
142 ural objective functions defined by their test-statistics. We illustrate the use of
143 two-sample tests by performing phylofactorization of a dataset of mammalian
144 body mass.

145 Then, we show how the phylogeny serves as a scaffold for changing variables
146 in biological data through a contrast basis - the same basis used in the isomet-
147 ric log-ratio transform - which can be used to identify the phylogenetic scales
148 providing low-rank, phylogenetically-interpretable representations of a dataset.
149 Defining the contrast basis allows us to introduce a phylogenetic analog of prin-
150 cipal components analysis - phylogenetic components analysis - which identifies

151 the dominant, phylogenetic scales capturing variance in a dataset. We perform
152 phylogenetic components analysis on the American Gut microbiome dataset
153 (www.americangut.org) and reveal that some of the dominant clades explaining
154 variation in the American Gut correspond to clades within Bacteroides and Fir-
155 micutes, thereby providing finer, phylogenetic resolution of a known, major axis
156 of variation in human gut microbiomes found to be associated with obesity [47],
157 age [31] and more. Another phylogenetic factor of variance in the American
158 Gut is a clade of Gammaproteobacteria strongly associated with IBD, corrobo-
159 rating a recent study's use of phylofactorization to diagnose patients with IBD
160 [49]. The contrast basis can also be used with regression if the data assumed
161 to be approximately normal, log-normal, logistic-normal or otherwise related
162 to the normal distribution through a monotonic transformation. We illustrate
163 regression-phylofactorization through a generalized additive model analysis of
164 how microbial abundances change across a range of pH, Nitrogen, and Carbon
165 concentrations in soils. The resulting contrast basis and its fitted values from
166 generalized additive modeling yield a low-rank representation of biological big-
167 data and translates to clear biological hypotheses aiming to identify the traits
168 driving observed non-linear patterns of abundance across pH [39].

169 Datasets comprised of non-Gaussian, exponential family random variables
170 can still be analyzed through regression-phylofactorization. We present four
171 methods for generalized regression-phylofactorization in exponential family data.
172 The first method is to use the contrast basis for constrained, reduced-rank re-
173 gression to obtain a low-rank approximation of coefficient matrices in multivari-
174 ate generalized linear models. The second uses a two-level factor, a surrogate
175 variable `phylo` indicating which side of an edge a species is found, to define
176 objective functions based on the deviance or the magnitude of the coefficients
177 for the factor-contrast. The third method aggregates exponential family data

178 within clades to marginally stable distributions within the exponential family,
179 and then performs `phylo` factor contrasts described above. The fourth method
180 is a mix of the first and second, developed to have the accuracy of the second
181 method while reducing the computational costs. The mixed method considers
182 `phylo` factors for only a subset of the best edges obtained from reduced-rank
183 approximation of the coefficient matrix.

184 We finish with a discussion of the challenges, and opportunities, for future
185 development of phylofactorization, and provide an R package - `phylofactor` -
186 available at <https://github.com/reptalex/phylofactor>.

187 Phylofactorization

188 Which vertebrates live on land, and which vertebrates live in the sea (Figure
189 1a)? Most children have enough natural history knowledge to say “fish live in the
190 sea”, thus correctly identifying one of the most important phylogenetic factors of
191 land/sea associations in vertebrates. The statement “fish live in the sea” can be
192 mathematically formalized by noting that one edge in the vertebrate phylogeny
193 separates “fish” from “non-fish” (Figure 1b). Partitioning the phylogeny along
194 the edge basal to tetrapods can separate vertebrates fairly well by land/sea asso-
195 ciations. An algorithm for identifying that edge by land/sea associations alone,
196 without requiring detailed knowledge of macroscopic life and morphological and
197 physiological traits, can correctly identify an edge along which functional ecolog-
198 ical traits and life-history traits arose. Controlling for the previously identified
199 edge, one might be able to identify the edges basal to Cetaceans and Pinnipeds,
200 tetrapods which live in the sea (Figure 1b). Three edges can capture most of
201 the variation in land/sea associations across thousands of vertebrate species.

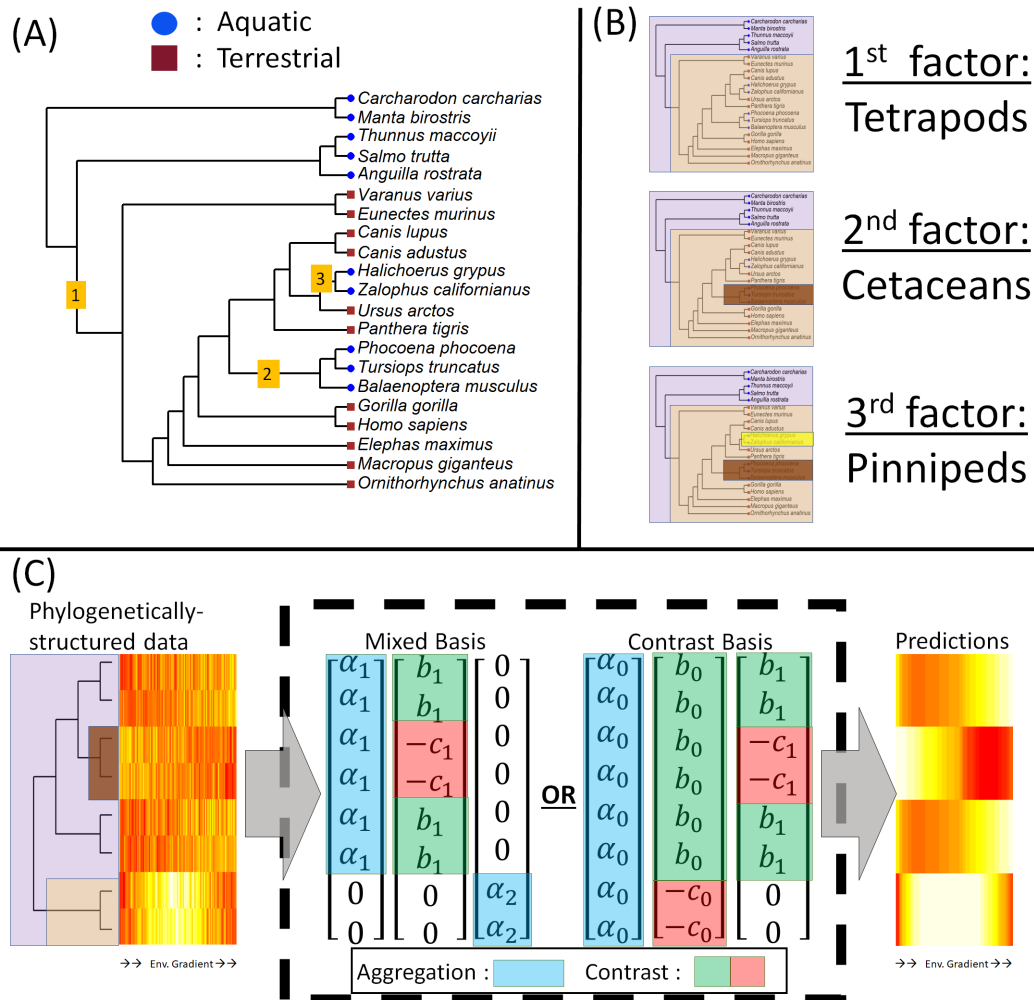


Figure 1: Phylofactorization generalizes the logic of how to simplify phylogenetically-structured datasets. (A) Vertebrate land/water associations can be simplified by partitioning the tree into the edges along which major traits arose. (B) The first phylogenetic factor of vertebrate land/water associations is the edge along which tetrapods arose - an edge along which lungs and limbs evolved that allowed colonization of land. Downstream factors can refine the original partitioning, and include the Cetaceans and Pinnipeds, among other edges along which adaptation to aquatic life arose among tetrapods. (C) Phylogenetic factorization generalizes this same logic for phylogenetically-structured data in which traits might not be known or their evolution easily modeled, including traits like a non-linear relationship between abundance and an environmental gradient. Phylogenetically-structured data can be partitioned through operations of aggregation and contrast. Pure aggregations (blue) sum data within a clade, whereas contrasts (green/red) are differences between two clades. Low-rank, phylogenetically-interpretable predictions of our data can be obtained through a mixed basis of a series of aggregations and contrasts, or a “contrast basis” in which there is a global aggregate partitioned in subsequent contrasts.

202 Ancestral state reconstruction of habitat association provides a well-known
203 means of making such inferences. However, sometimes the desired traits and
204 ecological patterns of interest are more complicated and their ancestral state re-
205 construction dubious. For instance, how can we identify the phylogenetic scales
206 of changes in microbial community composition along a pH gradient, allow-
207 ing possible non-linear associations that could be detected through generalized
208 additive modeling? Answering such a question through ancestral state recon-
209 struction requires conceiving and analyzing an evolutionary model of how the
210 generalized additive models of pH association evolve along a tree. Phylofactor-
211 ization aims to generalize the phylogenetic logic used for land/sea associations
212 in order to identify phylogenetic scales for more complicated functional traits
213 and ecological patterns, for which an evolutionary model would be dubious.
214 Phylogenetic factorization generalizes the logic of land/sea associations through
215 a graph partitioning algorithm iteratively identifying edges in the phylogeny
216 along which meaningful differences arise (Figure 1c).

217 **General Algorithm**

218 Phylofactorization requires a set of disjoint phylogenies spanning the set of
219 species considered in the data. The phylogenies are rooted or unrooted graphs
220 with no cycles, containing and connecting the units of interest in our data (the
221 units can be species, genes other evolving units of interest). Phylofactoriza-
222 tion can be implemented with disjoint sub-graphs, such as viral phylogenies for
223 which there are not clear common ancestors, and the sub-phylogenies can either
224 be kept separate or joined at a polytomous root. The phylogeny may have an
225 arbitrary number and degree of polytomies.

226 Let $[x]_{i,j}$ be the data for species $i = 1, \dots, m$ in sample $j = 1, \dots, n$. Let
227 $\mathbf{x}_{R,j}$ be the vector of a subset of species, R , in sample j . Let \mathbf{Z} be the $n \times p$

228 matrix containing p additional meta-data variables for each sample. Let \mathcal{T} be
229 the phylogenetic tree and let edge e partition the phylogeny into disjoint groups
230 R and S . Phylofactorization requires:

- 231 • An aggregation function, $A(\mathbf{x}_{R,j}, \mathcal{T})$ which aggregates any subset, R , of
232 species
- 233 • A contrast function, $C(A(\mathbf{x}_{R,j}, \mathcal{T}), A(\mathbf{x}_{S,j}, \mathcal{T}), \mathcal{T}, e)$ which contrasts the
234 aggregates of two disjoint subsets of species, R and S , possibly using
235 information from the tree \mathcal{T} and edge, e .
- 236 • An objective function, $\omega(C, \mathbf{Z})$.

237 With these operations, phylofactorization is defined iteratively as a special case
238 of a graph partitioning algorithm (Figure 2). The steps of phylofactorization
239 are:

- 240 1. For each edge, e , separating disjoint groups of species R_e and S_e within the
241 sub-tree containing e , compute $C_e = C(A(\mathbf{x}_{R_e,j}, \mathcal{T}), A(\mathbf{x}_{S_e,j}, \mathcal{T}), \mathcal{T}, e)$
- 242 2. compute edge objective $\omega_e = \omega(C_e, \mathbf{Z})$ for each edge, e
- 243 3. Select winning edge $e^* = \underset{e}{\operatorname{argmax}}(\omega_e)$
- 244 4. Partition the sub-tree containing e^* along e^* , forming two disjoint sub-
245 trees.
- 246 5. Repeat 1-5 until a stopping criterion is met.

247 Unlike more general graph-partitioning algorithms, phylofactorization does not
248 impose a balance constraint - it does not require that the partitions have a simi-
249 lar size or weight. Furthermore, phylofactorization, by working with phylogenies
250 or graphs without cycles is centered around aggregation and contrast as princi-
251 ple operations for defining scales and units of organization. Phylofactorization is

252 limited to contrasts of non-overlapping groups, and the constraint of contrasting
253 aggregates is used to formalize the process of aggregation prior to contrasting
254 groups - such formalization ensures one can subsequently aggregate the bins of
255 species partitioned in phylofactorization according to the method of aggregation
256 by which the bins were discovered to be different. The incorporation of the tree,
257 \mathcal{T} , in the contrast function encompasses a class of ancestral state reconstruction
258 reconstruction methods. Ancestral state reconstruction with non-overlapping
259 contrasts can be done with time-reversible models of evolution; in this case,
260 phylofactorization contrasts the root ancestral states obtained in which the two
261 nodes adjacent an edge are considered roots of the subtrees separated by an
262 edge.

263 The edges, e^* and their contrasts, C_e , are interchangeably referred to as
264 the “phylogenetic factors” due to their correspondence to hypothesized latent
265 variables (traits) and their ability to construct basis elements that allow ma-
266 trix factorization [51]. It’s possible to define objective functions through pure
267 aggregation, but we limit our focus to contrast-based phylofactorizations which
268 identify edges along which meaningful differences arose for reasons discussed
269 later in the section on the “contrast basis”.

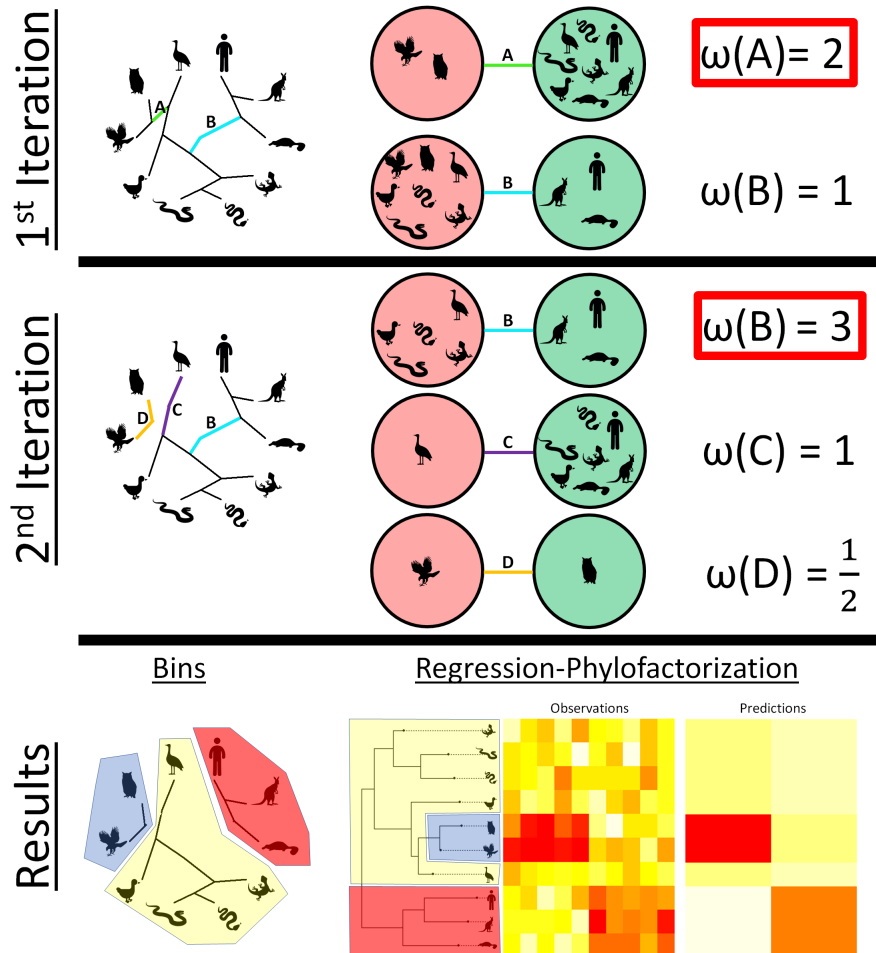


Figure 2: Phylofactorization is a graph partitioning algorithm. An objective function, ω , of a contrast of species separated by an edge allows one to iteratively partition the phylogeny along edges maximizing the objective function (1st iteration). After partitioning the phylogeny, the objective functions are re-computed to contrast species in the same sub-tree separated by an edge. Edge B in the first iteration contrasted mammals from non-mammals, but in the second iteration it contrasts mammals from non-mammals, excluding raptors (partitioned in the first iteration). The result of k iterations of phylofactorization is a set of $k + 1$ bins of species with similar within-group behavior. A particularly useful case is “regression-phylofactorization”. Regression-phylofactorization is implemented by defining contrasts through the contrast basis (Figure 1c) and defining an objective function through regression on the component scores of each candidate contrast basis element. Regression-phylofactorization is a flexible way to search for clades with similar patterns of association with environmental meta-data while also obtaining low-rank, phylogenetically-interpretable representations of a data matrix.

270 The result of phylofactorization after t iterations is a set of t inferences on
271 edges or links of edges. Links of edges occur following a previous partition,
272 when two adjoining edges separate the same two groups in the resultant sub-
273 tree. Partitioning the phylogeny along t edges results in $t + 1$ bins of species,
274 referred to as “binned phylogenetic units”. In general, the problem of maximizing
275 some global objective function, $\omega(e_1^*, \dots, e_t^*)$, for a set of t edges, $\{e_1^*, \dots, e_t^*\}$, is
276 NP hard [6]. However, stochastic searches of the space of possible partitions,
277 via a stochastic computation of ω_e in step 2 or a weighted draw of e^* in step 3,
278 may yield better approximations of a global maximum [32, 20, 23].

279 Aggregation, contrast, and objective functions are several junctures to define
280 and interpret meaningful quantities and outcomes from data analysis. Explicit
281 decisions about aggregation formalize how a researcher would summarize data
282 from an arbitrary set of species. Explicit decisions about contrast formalize how
283 a researcher differentiates two arbitrary, disjoint groups of species - these com-
284 mon operations form an organizational framework for ecologists studying phy-
285 logenetic scales. Aggregation can be done through many operations, including
286 but not limited to addition, multiplication, generalized means, and maximum
287 likelihood estimation of ancestral states under models of trait diffusion away
288 from the focal node. Likewise, examples of contrasts are differences, ratios, var-
289 ious two-sample tests, and more complicated metrics of dissimilarity such as the
290 deviance of a factor contrast in a generalized additive model. Researchers must
291 decide for themselves how best to aggregate information in groups of species,
292 contrast two groups, and decide which group maximizes the objective for a
293 research goal pertaining to a particular ecological pattern. Doing so allows ob-
294 jective, a priori definitions of what makes an informative phylogenetic scale,
295 and the operations chosen are integrated into a broader theoretical framework
296 of phylofactorization.

297 Below, we develop the generality and illustrate the results from phylofac-
298 torization. These examples were run using the R package “phylofactor”, using
299 relevant functions for analyzing and visualizing phylogenies from the R packages
300 ape [36], phangorn [43], phytools [40], and ggtree [53]. Scripts and datasets for
301 every analysis are available in the supplemental materials.

302 **Example 1: two-sample tests and mammalian body-mass** 303 **phylofactorization**

304 If the data are a single vector of observations, \mathbf{x} , similar to the land/sea associ-
305 ations of vertebrates, phylofactorization can be implemented through standard-
306 ized tests for differences of means or rate parameters in the two sets of species,
307 R and S .

308 To illustrate, we phylofactorize a dataset of mammalian body mass from
309 PanTHERIA [24] and the open tree of life using the R package “rotl” [33]. A
310 single vector of data assumed to be log-normal can be factored based on a two-
311 sample t-test (Figure 3a). In this case, $A(\mathbf{x}_R) = \overline{\log(\mathbf{x}_R)}$ is the arithmetic mean
312 of the log-body-mass; we use the contrast operation

$$C = \frac{|A(\mathbf{x}_R) - A(\mathbf{x}_S)|}{\sqrt{\frac{1}{r} + \frac{1}{s}}} \quad (1)$$

313 and the objective function $\omega_e = C_e$. Equation (1) defines the test-statistic for
314 a two-sample t-test with the assumption of constant variance. Maximization of
315 the objective function yields edges with the most significant difference in body
316 mass of organisms on different sizes of the tree.

317 The first five phylogenetic factors of mammalian body mass in these data are
318 Euungulata, Ferae, Laurasiatheria (excluding Euungulata and Ferae), a clade
319 of rodent sub-orders Myodonta, Anomaluromorpha, and Castorimorpha, and

320 the simian parvorder Catarrhini. Five factors produce six binned phylogenetic
321 units of species with different average body mass (Figure 3a). The most sig-
322 nificant phylogenetic partition of mammalian body mass occurs along the edge
323 basal to Euungulata, containing 296 species with significantly larger body mass
324 than other mammals. The second partition corresponds to Ferae, containing 242
325 species which have body masses larger than other mammals, excluding Euungu-
326 lata. The third partition corresponds to 864 remaining species in Laurasiathe-
327 ria, excluding Euungulata and Ferae, which contains Chiroptera, Erinaceomor-
328 pha, and Soricomorpha. These mammals have lower body mass than non-
329 Laurasiatherian mammals. The fourth partition identifies three rodent sub-
330 orders comprising 926 species with lower body mass than non-Laurasiatherian
331 mammals. Finally, 106 species comprising the Simian parvorder Catarrhini
332 are factored as having higher body mass than the remaining mammals. These
333 factors are fairly robust: 3000 replicates of stochastic Metropolis-Hasting phylo-
334 factorization, drawing edges in proportion to C^λ with $\lambda = 6$ (producing a 1/4
335 probability of drawing the most dominant edge) could not improve upon these
336 5 factors.

337 The first two phylogenetic factors of mammalian body size partition the
338 mammalian tree at deep edges with ancestors near the K-Pg extinction event,
339 corroborating evidence of ecological release [2, 3] and the exponential growth
340 of maximum body sizes following the K-Pg extinction event [46] for these two
341 dominant clades. The crown group of modern Euungulata are thought to have
342 originated in the late Cretaceous [54] and its representatives may have expanded
343 into previously dinosaur-occupied niches during the rapid evolution of body
344 size in mammals immediately after the K-Pg extinction event at the Creta-
345 ceous/Paleogene boundary [45]. Cope's rule posits that lineages tend to in-
346 crease in body size over time, and a recent study [4] confirms Cope's rule and

347 found that mammals have, along all branch lengths in their phylogeny, tended
348 to increase in size. The phylogenetic factors of mammalian body size discovered
349 here illustrate an important feature of phylofactorization: correlated evolution
350 within a clade, such as a consistently high body-size increase among lineages in
351 a clade, can cause the edge basal to a clade to be an important partition for
352 capturing variance in a trait. A more robust phylofactorization may be done
353 through iterative ancestral-state reconstruction of the roots of subtrees parti-
354 tioned by each edge (where the subtrees are re-rooted at the nodes adjacent
355 the edge), but this unsupervised phylogenetic factorization body masses in 3374
356 mammals takes 15 seconds on a laptops and yields partitions which simplify the
357 story of mammalian body-mass variation to a set of 5 edges forming 6 binned
358 phylogenetic units.

359 Two-sample tests can be used for phylogenetic factorization of any vector of
360 trait data. For another example, Bernoulli trait data, such as presence/absence
361 of a trait, can be factored using Fisher's exact test that there is the same
362 proportion of presences in two groups, R and S . In this case, the aggregation
363 operation $A(\mathbf{x}_R) = \sum_{i \in R} x_i$ counts the number of successes in group R , the
364 contrast operation is the computation of the P-value using Fisher's exact test
365 with the contingency table

Successes	Failures	Total
$A(\mathbf{x}_R)$	$r - A(\mathbf{x}_r)$	r
$A(\mathbf{x}_S)$	$s - A(\mathbf{x}_S)$	s
$A(\mathbf{x}_R) + A(\mathbf{x}_S)$	$r + s - (A(\mathbf{x}_r) + A(\mathbf{x}_S))$	$r + s$

366 An objective function can be defined as the inverse of the P-value from Fisher's
367 exact test, $\omega_e = |C_e^{-1}|$. The phylofactorization of vertebrates by land/water
368 association in Figure 1, using an ad-hoc selection of vertebrates for illustration,
369 was performed using Fisher's exact test, and the factors obtained correspond to
370 Tetrapods, Cetaceans, and Pinnipeds. Unlike the phylofactorization of mam-

371 malian body mass, all three factors obtained from phylofactorization of verte-
372 brate land/water association correspond to a set of traits. Tetrapods evolved
373 lungs and limbs which allowed them to live on land. Cetaceans evolved fins and
374 blowholes, and Pinnipeds evolved fins, all traits adaptive to life in the water.

375 Two-sample tests are used when partitioning a vector of traits and not con-
376 trolling for additional meta-data such as sampling effort or other confounding
377 effects. Phylofactorization of body mass and land/water associations illustrate
378 two potential evolutionary models under which edges are important: correlated
379 evolution of members of a clade and punctuated equilibria. Edges identified from
380 more complicated methods of phylofactorization may correspond to traits, or
381 they may correspond to directional evolutionary processes shared among mem-
382 bers of a clade or their ancestors, such as ecological release or niche partitioning.
383 When the objective function from two-sample tests has a well-defined null dis-
384 tribution, as is the case for the two-sample *t*-test and Fisher’s exact test, the
385 uniformity of the distribution of P-values can be used to define a stopping criteria
386 as discussed later (see: “stopping criteria”).

387 **Example 2: Contrast basis and phylogenetic components** 388 **analysis**

389 The phylogeny provides a natural scaffold for low-rank, phylogenetically in-
390 terpretable approximations of the data. As a sphere defines a natural set of
391 coordinates for GPS data, the phylogeny defines a natural set of coordinates
392 that can be used for a variety of data analyses. One example of a natural coor-
393 dinate in the phylogeny is aggregation: the sum of abundances of species within
394 a clade. Another natural coordinate is a contrast: the differences of abundance
395 between two clades, either sister clades or a monophyletic clade and its comple-
396 ment. Together, these operations allow one to construct natural coordinates for

397 more sophisticated analyses of phylogenetically-structured ecological data.

398 Phylogenetically-interpretable, low-rank approximations of data can be ob-
399 tained by constructing basis elements through aggregation and contrast vectors
400 (Figure 1c). An aggregation basis element for a group $Q = R \cup S$ can be
401 constructed through a vector whose i th element is

$$\mathbf{v}_{A_Q,i} = \begin{cases} a & i \in Q \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

402 and such aggregation basis elements can be subsequently partitioned with a
403 contrast vector

$$\mathbf{v}_{C_{R|S},i} = \begin{cases} b & i \in R \\ -c & i \in S \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where $b > 0$ and $c > 0$. By meeting the criteria

$$rb - sc = 0 \quad (4)$$

$$rb^2 + sc^2 = 1 \quad (5)$$

, one can ensure that \mathbf{v}_{A_Q} and \mathbf{v}_{C_Q} are orthogonal and with unit norm. These criteria are satisfied by

$$b = \sqrt{\frac{s}{r(r+s)}} \quad (6)$$

$$c = \sqrt{\frac{r}{s(r+s)}}. \quad (7)$$

When projecting data from sample j , \mathbf{x}_j , onto a contrast vector, the aggregation

and contrast operations are

$$\begin{aligned} A(\mathbf{x}_{R,j}) &= \bar{\mathbf{x}}_{R,j} \\ C(A(\mathbf{x}_{R,j}), A(\mathbf{x}_{S,j})) &= \sqrt{\frac{rs}{r+s}} (\bar{\mathbf{x}}_{R,j} - \bar{\mathbf{x}}_{S,j}). \end{aligned} \quad (8)$$

404 where $\bar{\mathbf{x}}_{R,j}$ is the sample mean of species in group R and sample j . Projecting
405 a dataset onto $\mathbf{v}_{C_{R|S}}$ yields coordinates which are a standardized difference of
406 means similar to equation (1). The contrast vector is comprised of two sub-
407 aggregations of opposite sign, one for group R and the other for group S . By
408 ensuring criterion (4), the groups aggregated within a contrast vector can be sub-
409 sequently partitioned with additional, orthogonal contrast vectors splitting each
410 group R and S . Maintaining criterion (5), the aggregation and contrast vectors
411 defined here can be used to construct an orthonormal basis for describing data
412 containing our species, $\mathbf{x}_j \in \mathbb{R}^m$, by defining a set of $q \leq m$ orthogonal aggrega-
413 tion vectors corresponding to disjoint sets of species Q_l such that the entire set
414 of aggregations, $\bigcup_{l=1}^q Q_l = \{1, \dots, n\}$, covers the entire set of m species. Then,
415 $m - q$ contrast vectors partitioning the aggregations and the sub-aggregations
416 within contrast vectors can complete the basis (Figure 1c). Of note is that, as
417 defined in equations (2) and (3), the span of any aggregate and its contrast is
418 equal to the span of the contrasts' sub-aggregates, i.e. for $R \cup S = Q$,

$$\text{span}(\mathbf{v}_{A_Q}, \mathbf{v}_{C_{R|S}}) = \text{span}(\mathbf{v}_{A_R}, \mathbf{v}_{A_S}) \quad (9)$$

419 (Figure 1c) and the two natural ways of changing variables with the phylogeny,
420 an aggregate of species and its orthogonal contrast (grouping species and parti-
421 tioning the group) or two orthogonal aggregates (two disjoint groups of species),
422 are rotations of one-another. Aggregation and contrast vectors translate the no-
423 tion of phylogenetic scale and group-differences into a basis that can be used to

424 analyze community ecological data.

425 Pure aggregation vectors as defined in equation (2) can be defined a priori
426 based on traits or clades of species thought to be important for the question
427 at hand (e.g. aggregate “terrestrial” and “aquatic” animals), or defined by the
428 data through myriad clustering algorithms or phylofactorization based purely on
429 aggregation by converting steps (1) and (2) in the phylofactorization algorithm
430 into a single step: maximizing an objective function of the aggregate of a clade.
431 A special case occurs when data are compositional [1], in which case the sum
432 of any sample across all species in the community will equal 1 and thus the
433 data are constrained by an aggregation element - the aggregate of all species
434 - which can only be subsequently contrasted. Phylofactorization via contrasts
435 of log-relative abundance data allows one to construct an isometric log-ratio
436 transform, a commonly used and well-behaved transform for the analysis of
437 compositional data [10, 9, 44]. Since the span of an aggregate and its contrast
438 is equal to the span of the contrasts’ two aggregates (equation 9), we simplify
439 construction of the basis by considering, from here on out, only the “contrast
440 basis” in which the an initial aggregate of all species is then partitioned with a
441 series of contrasts.

442 An orthonormal basis, including one constructed via aggregation and con-
443 trast vectors, enables researchers to partition the variance along each of a set
444 of orthogonal directions corresponding to discrete, identifiable features in the
445 phylogeny. Using the phylofactorization algorithm, a dataset $\mathbf{X} = [x]_{i,j}$ can be
446 summarized by defining the objective function

$$\omega_e = \text{Var} [\mathbf{v}_{C_e}^T \mathbf{X}] \quad (10)$$

447 where \mathbf{v}_{C_e} is the contrast vector from (3) corresponding to the sets of species, R
448 and S , split by edge e . The objective function in equation (10) yields a phyloge-

449 netic decomposition of variance we define as “phylogenetic components analysis”
450 or PhyCA. PhyCA is a constrained version of principal components analysis,
451 allowing researchers to identify the dominant axes of variation, constrained to
452 axes which contrast species separated by an edge.

453 The variance of component scores, $\mathbf{y}_e = \mathbf{v}_{C_e}^T \mathbf{X}$, are easiest to understand
454 when the data $[x_{i,j}]$ are assumed to be standard Gaussian. The component
455 score for sample j , $\mathbf{y}_{e,j}$, can be written as

$$\mathbf{y}_{e,j} = \sqrt{\frac{rs}{r+s}} (\bar{\mathbf{x}}_{R,j} - \bar{\mathbf{x}}_{S,j}) \quad (11)$$

456 where $\bar{\mathbf{x}}_{R,j}$ is the sample mean of $x_{i,j}$ for $i \in R$ and $\bar{\mathbf{x}}_{S,j}$ is the sample mean of
457 $x_{i,j}$ for $i \in S$. The variance of the component score across all samples $j = 1, \dots, n$
458 is

$$\text{Var}[\mathbf{y}_e] = \frac{rs}{r+s} (\text{Var}[\bar{\mathbf{x}}_R] + \text{Var}[\bar{\mathbf{x}}_S] - 2\text{Cov}[\bar{\mathbf{x}}_R, \bar{\mathbf{x}}_S]). \quad (12)$$

459 The variance of \mathbf{y}_e increases through a combination of variances in aggregations
460 of groups R and S across samples ($\bar{\mathbf{x}}_R$ and $\bar{\mathbf{x}}_S$, respectively) and a high negative
461 covariance between aggregations for groups R and S across samples. Species
462 with a negative covariance may be competitively excluding one-another or may
463 be differentiated due to a trait which arose along edge e which causes different
464 habitat associations or responses to treatments. Edges extracted from PhyCA
465 are edges along which putative functional ecological traits arose differentiating
466 the species in R and S in the dataset of interest.

467 **Phylogenetic Components of the American Gut** To illustrate, we per-
468 form PhyCA to identify 10 factors from a sub-sample of the American Gut
469 dataset and the greengenes phylogeny [8] containing $m = 1991$ species and
470 $n = 788$ samples from human feces (Figure 3b). The American Gut dataset
471 was filtered to only fecal samples with over 50,000 sequence counts and, for

472 those samples, otus with an average of more than one sequence count per
473 sample. After performing PhyCA, each identified resulting component score,
474 y_{e^*} , was assessed for a linear association with seven explanatory variables:
475 types_of_plants (a question asking participants how many types of plants
476 they've eaten in the past week), age, bmi, alcohol consumption frequency, sex,
477 antibiotic use (ABX), and inflammatory bowel disease (subset_ibd) (Figure
478 3b). The raw P-values are presented below, but for a reference, the P-value
479 threshold for a 5% family-wise error rate is 7.1×10^{-4} .

480 The first factor splits 1229 species of Firmicutes from the remainder of mi-
481 crobes. The component score for the first factor, $y_{e_1^*}$, is strongly associated with
482 antibiotic use ($P=3.6 \times 10^{-4}$), with dramatic decreases in relative abundance
483 in patients who have taken antibiotics in the past week or month. The second
484 factor identifies 217 species of several genera of Lachnospiraceae, a clade con-
485 tained within the Firmicutes, strongly associated with age ($P=1.2 \times 10^{-15}$) and
486 bmi ($P=3.2 \times 10^{-6}$) and alcohol ($P=6.4 \times 10^{-3}$). The third factor is a clade of
487 81 Bacteroides most strongly associated with types_of_plants ($P=2 \times 10^{-9}$).
488 By identifying a clade of Bacteroides as a major axis of variation, factors 1
489 and 3 refine the Firmicutes to Bacteroidetes ratio commonly used to describe
490 variation in the gut microbiome and found associated with obesity and other
491 disease states [28, 7]. It's been found that the Firmicutes/Bacteroidetes ratio
492 changes with age [31], but the picture from phylofactorization is more nuanced:
493 the large clade of Firmicutes in the first factor does not change with age, but
494 the Lachnospiraceae within that clade decrease strongly with age relative to
495 the remaining Firmicutes, while the Bacteroides show only a moderate decrease
496 with age. The strong decrease with age in Lachnospiraceae is found in a few
497 other clades within the Firmicutes: the 4th factor identified a clade of Firmi-
498 cutes of the family Ruminococcaceae strongly associated with types of plants

499 ($P=3.6 \times 10^{-5}$), sex ($P=5.9 \times 10^{-4}$) and decreasing with age ($P=9.2 \times 10^{-4}$),
500 and the 5th factor identified a group of Firmicutes of the family Tissierellaceae
501 that decrease strongly with age ($P=1.9 \times 10^{-5}$).

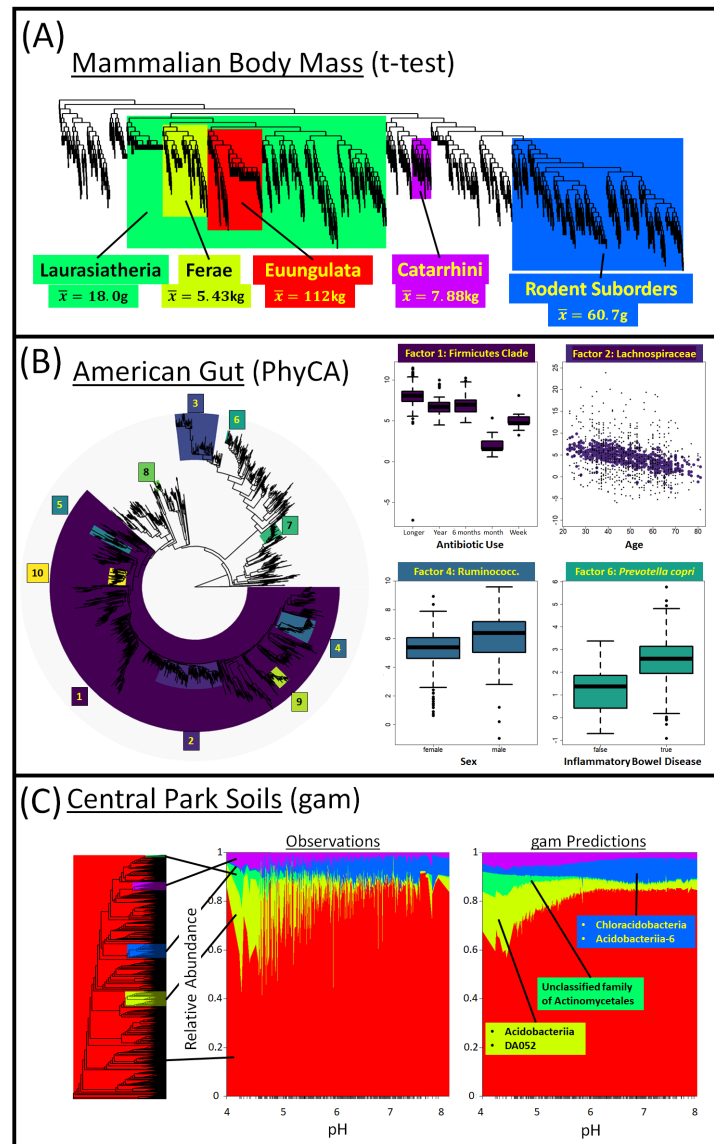


Figure 3: Phylofactorization with contrast basis. (A) The contrast basis defines variables similar to t-statistics, and maximizing the projection of data onto the contrast basis can identify phylogenetic factors. Five iterations of phylofactorization on a dataset of mammalian log-body mass yields five clades with very different body masses. (B) Maximizing the variance of component scores, y_e , of log-relative abundance data produces a “phylogenetic components analysis” (PhyCA) of the American Gut dataset. The most variable clades cover a range of phylogenetic scales. Downstream analysis of component scores tested associations with meta-data - plotted are linear predictors against relevant meta-data; the plot of Lachnospiraceae includes the raw data as black dots. (C) More complicated methods can be used, such as generalized additive modeling with y_e . Using the central park soils dataset, y_e of log-relative abundances, the model $y_e \sim s(\log(\text{Carbon})) + s(\log(\text{Nitrogen})) + s(\text{pH})$, and the objective of maximizing the explained variance, we obtained the same 4 factors obtained using generalized linear modeling in the original data, including the misnomer group of Chloracidobacteria that don't thrive in low pH environments. The relative importance of pH in the generalized additive models and exact clades with a high amount of variance explained by pH allows a projection of 3000 species into 5 BPUs for clear visualization of a dominant feature of how soil bacterial communities change along a key environmental gradient.

502 The sixth factor is a small group of 5 OTUs of *Prevotella copri* strongly as-
503 sociated with types_of_plants ($P=2.8 \times 10^{-4}$) and inflammatory bowel disease
504 ($P=2.5 \times 10^{-3}$). Previous studies have found that *Prevotella copri* abundances
505 are correlated with rheumatoid arthritis in people and inoculation of *Prevotella*
506 *copri* exacerbates colitis in mice. Consequently, *Prevotella copri* is hypothesized
507 to increase inflammation in the mammalian gut [42], and the discovery of *Pre-*
508 *vetella copri* as one of the dominant phylogenetic factors of the American Gut, as
509 well as the discovery of its association with IBD, corroborates the hypothesized
510 relationship between *Prevotella copri* and inflammation. Likewise, the seventh
511 factor is a clade of 41 Gammaproteobacteria of the order Enterobacteriales also
512 associated with types_of_plants ($P=6.7 \times 10^{-8}$) and weakly associated with
513 inflammatory bowel disease ($P=0.022$). Gammaproteobacteria were used as
514 biomarkers of Crohn’s disease in a recent study [49] and their associations with
515 IBD in the American Gut project corroborates the possible use of Gammapro-
516 teobacterial abundances for detection of IBD from stool samples. Summaries of
517 the models for all factors’ component scores are in the supplemental information.

518 **Example 3: Compositional, log-normal and Gaussian regression-** 519 **phylofactorization**

520 The contrast basis can be used to define more complicated objective functions
521 for data assumed to be Gaussian or easily mapped to Gaussian, such as logistic-
522 normal compositional data or log-normal data. Conversion of the data to an
523 assumed-Gaussian form can then allow one to perform least-squares regression
524 using \mathbf{y}_e as either an independent or dependent variable. Rather than per-
525 forming PhyCA and subsequent regression, one can choose phylogenetic factors
526 based on their associations with meta-data of interest.

527 Maximizing the explained variance from regression identifies clades through

528 the product of a high contrast-variance from equation (10) and the percent
529 of explained-variance from regression - such clades can capture large blocks
530 of explained variance in the dataset. Another common objective function is
531 the deviance or F -statistic from regression which identifies clades with more
532 predictable responses - such clades can be seen as bioindicators or particularly
533 sensitive clades, even if they are not particularly large or variable clades in
534 the data. Regression-phylofactorization can use the component scores as an
535 independent variable, as was used in the phylofactorization-based classification
536 of Crohn's disease [49]. For multiple regression, one can use the explanatory
537 power of the entire model, or a more nuanced objective function of a subset of
538 the model. More complicated regression models can be considered, including
539 generalized additive models, regularized regression, and more.

540 To illustrate the flexibility of regression phylofactorization to identify phy-
541 logenetic scales corresponding to nonlinear patterns of abundance-habitat asso-
542 ciations, we perform a generalized additive model analysis of the Central Park
543 soils dataset [39] analyzed previously using a generalized linear model. To iden-
544 tify non-linear associations between clades and pH, Carbon and Nitrogen, we
545 perform a generalized additive model of the form

$$\mathbf{y}_e \sim s(\text{pH}) + s(\text{Carbon}) + s(\text{Nitrogen}) \quad (13)$$

546 and maximize the explained variance (Figure 3c). The resultant phylofactor-
547 izations identifies the same 4 factors as the generalized linear model, but allows
548 nonlinear and multivariate analysis of how community composition changes over
549 environmental meta-data. Combining the high relative-importance of pH with
550 the identified 4 factors, splitting over 3,000 species 5 binned phylogenetic units,
551 allows clear and simple visualization of otherwise complex behavior of how a
552 community of several thousand microbes changes across several hundred soil

553 samples. As with the original analysis, the generalized additive modeling phylo-
554 factorization identifies a clade of Acidobacteria - the Chloracidobacteria - which
555 have highest relative abundances in more neutral soils.

556 **Example 4: Phylofactorization through generalized linear** 557 **models**

558 Many ecological data are not Gaussian. Presence-absence data or count data
559 with many zeros cannot be easily transformed to yield approximately Gaussian
560 random variables. Data assumed to be exponential family random variables can
561 be analyzed with regression-phylofactorization by adapting concepts in gener-
562 alized linear models.

563 We present four options for phylofactorization through generalized linear
564 models. These options correspond to the contrast basis, either explicitly us-
565 ing the contrast basis to approximate the coefficient matrix in multivariate
566 generalized linear models, or implicitly using a form of the contrast basis in
567 the likelihood function when performing shared-coefficient or factor-contrasts in
568 generalized linear modeling.

569 **Coefficient Contrast** The first method, related to reduced rank regression
570 for vector generalized linear models [52], uses the contrast basis to provide a
571 reduced-rank approximation of the coefficient matrix from multivariate general-
572 ized linear models. Multivariate (vector) generalized linear models assume the
573 data \mathbf{X} are drawn from an exponential family distribution with canonical pa-
574 rameters for each species, $\boldsymbol{\eta} \in \mathbb{R}^m$, related to the meta-data \mathbf{Z} through a linear
575 model

$$\boldsymbol{\eta} \sim \mathbf{B}\mathbf{Z} \tag{14}$$

576 where $\mathbf{B} \in \mathbb{R}^{m \times p}$ is the coefficient matrix and $\mathbf{Z} \in \mathbb{R}^{p \times n}$ is the matrix of meta-
577 data. Instead of using $m \times p$ coefficients, one can represent the coefficient matrix
578 \mathbf{B} through contrast basis elements and their component scores

$$\mathbf{B} = \mathbf{1}\mathbf{w}_0^T + \mathbf{V}\mathbf{W} + \epsilon \quad (15)$$

579 where $\mathbf{1} \in \mathbb{R}^m$ is the one vector, $\mathbf{w}_0 \in \mathbb{R}^p$ contains the sum of the regression
580 coefficients for each of the p predictors, $\mathbf{V} \in \mathbb{R}^{m \times t}$ is a matrix whose columns
581 are contrast basis elements obtained from t iterations of phylofactorization and
582 $\mathbf{W} \in \mathbb{R}^{t \times p}$ is a matrix whose rows are the component scores for each contrast
583 basis element. If one is interested in partitioning species based on a subset, P ,
584 of the explanatory variables, one can implement equation (15) for the matrix
585 \mathbf{B}_P containing only the partitioning variables for phylofactorization.

586 To put multiple independent meta-data from multiple species on the same
587 scale, it's important to standardize the coefficients $\beta_{i,j}$ by dividing them by
588 their standard error. We refer to these standard coefficients as $\beta_{i,j}^0$ and the ma-
589 trix of such standard coefficients for partitioning variables as the “standardized
590 coefficient matrix”, \mathbf{B}_P^0 .

591 A useful objective function for approximating the coefficient matrix with
592 the contrast basis is the Euclidean norm of the projection of the standardized
593 coefficient matrix onto contrast basis elements,

$$\omega_e = \|\mathbf{v}_{C_e}^T \mathbf{B}_P^0\| \quad (16)$$

594 which captures the extent to which coefficients in \mathbf{B}_P^0 differ between the sets
595 of species partitioned by the edge e . Coefficient contrasts are fast and easy
596 to compute, but the algorithm described here minimizes the distance between
597 $\mathbf{V}\mathbf{W}$ and \mathbf{B}_P^0 . Other algorithms described below can more robustly identify the

598 edge, e , whose reduced-rank approximation maximizes the likelihood.

599 **Stepwise phylo factor contrasts** Other options for aggregation and con-
600 trast exploit the factor-contrasts built into generalized linear and additive mod-
601 eling machinery. Factor contrasts using a variable `phylo` $\in \{R, S\}$, indicating
602 which group a species is in, can capture the assumption of shared coefficients
603 within-groups and contrast the coefficients between-groups in multivariate gen-
604 eralized linear modeling across all species. Stepwise, maximum-likelihood selec-
605 tion of `phylo` factor contrasts are a more accurate, yet computationally intensive,
606 algorithm for partitioning exponential family random variables.

607 For example, a data frame contrasting how the counts of “birds” from “non-
608 birds” react to meta-data z_2 while controlling for z_1 can be constructed as
609 follows

Site	Species	Abundance	z_1	z_2	<code>phylo</code>
1	Pigeon	10	1	.5	<i>R</i>
1	Dove	8	1	.5	<i>R</i>
1	Lizard	1	1	.5	<i>S</i>
1	Mouse	3	1	.5	<i>S</i>
1	Cat	1	1	.5	<i>S</i>
2	Pigeon	2	0	-2	<i>R</i>
2	Dove	1	0	-2	<i>R</i>
2	Lizard	10	0	-2	<i>S</i>
2	Mouse	4	0	-2	<i>S</i>
2	Cat	3	0	-2	<i>S</i>
...

610 Phylofactorization can be implemented through a generalized linear model for

611 a count family (e.g. Poisson, binomial, or negative binomial) using the formula

$$\text{Abundance} \sim z_1 + \text{phylo} \times z_2. \quad (17)$$

612 The `phylo` factor contrasts birds from non-birds and using its deviance as the
613 objective function will find the edge e^* whose `phylo` factor maximizes the like-
614 lihood of the data.

In stepwise `phylo` factor contrasts, aggregation occurs within the likelihood function. The likelihood $\mathcal{L}(\mathbf{x}_j; \boldsymbol{\eta})$ for a vector of binomial random variables \mathbf{x}_j can be written in exponential family form

$$\mathcal{L}(\mathbf{x}_j; \boldsymbol{\eta}) = h(\mathbf{x}_j) \exp \{ \boldsymbol{\eta}' \mathbf{x} - \mathcal{A}(\boldsymbol{\eta}) \}. \quad (18)$$

615 A two-factor model, such as $\mathbf{x} \sim \text{phylo}$, will reduce the likelihood function from
616 s parameters in $\boldsymbol{\eta}$ to two parameters, $\boldsymbol{\eta} \in (\eta_R, \eta_S)$, yielding

$$\mathcal{L}(\mathbf{x}_j; \text{phylo}) = h(\mathbf{x}_j) \exp \left\{ \eta_R \sum_{i \in R} x_{i,j} + \eta_S \sum_{i \in S} x_{i,j} - \mathcal{A}(\boldsymbol{\eta}) \right\}.$$

617 Aggregation, within the likelihood function above, is summation of data within
618 groups. Obtaining the maximum likelihood estimates, $\hat{\eta}_R$ and $\hat{\eta}_S$, a contrast
619 function can be defined as a difference of η_R and η_S , or test-statistic from the
620 null hypothesis that $\eta_R = \eta_S$. For general purposes, the deviance of the `phylo`
621 term in generalized linear or additive models serves as a useful contrast allowing
622 one to identify the edge e^* whose `phyloe` factor that maximizes the likelihood
623 for the regression model containing the `phylo` factor.

624 Stepwise selection of maximum-likelihood `phylo` factor contrasts is a very
625 accurate method for regression-phylofactorization of exponential family ran-
626 dom variables. However, unlisting an entire dataset, computing a glm, and

627 re-computing the glm for each edge in the phylogeny is computationally inten-
628 sive.

629 **Marginally Stable (mStable) Aggregation** Another option, aimed to al-
630 low maximum-likelihood estimation of phylo factor contrasts while reducing the
631 computational difficulty, is to aggregate the data \mathbf{X} prior to maximizing the
632 likelihood in the generalized linear model. The method we present is to assume
633 within-group homogeneity and aggregate exponential family random variables
634 to a “marginally stable” exponential family random variable that can be used
635 for downstream analysis. Marginal stability, to the best of our knowledge, has
636 not been explicitly defined elsewhere, and thus we introduce the term here by
637 loosening the definition of stable distributions [41].

638 **Stable distribution** A distribution with parameters θ , $\mathcal{F}(\theta)$, is said to be
639 stable if a linear combination of two independent random variables from $\mathcal{F}(\theta)$
640 is also in $\mathcal{F}(\theta)$, up to location and scale parameters.

641 **Marginally stable distribution** A distribution with parameters $\{\theta_1, \theta_2\}$,
642 $\mathcal{F}(\theta_1, \theta_2)$, is said to be marginally stable on θ_1 if $\mathcal{F}(\theta_1, \theta_2)$ is it is stable condi-
643 tioned on θ_1 being fixed.

644

645 For example, the Gaussian distribution is stable: the sum of two Gaus-
646 sian random variables is also Gaussian. Meanwhile, binomial random variables
647 $Binom(\rho, N)$ are marginally stable on ρ ; random variables $x_i \sim Binom(\rho, N_i)$
648 can be summed to yield $A(\mathbf{x}) \sim Binom(\rho, \sum N_i)$. The marginal stability can
649 also be used with transformations that connect the assumed distribution of the
650 data to a marginally stable distribution. Log-normal random variables can be
651 converted to Gaussians through exponentiation; chi random variables can be

652 converted to chi-squared through squaring - random variables from many dis-
653 tributions may be analyzed by transformation to a stable or marginally stable
654 family of distributions. Such transformation-based analyses implicitly define
655 aggregation through a generalized f -mean

$$A_f(\mathbf{x}_R) = f^{-1} \left(\sum_{i \in R} f(x_i) \right) \quad (19)$$

656 where $f(x) = \log(x)$ for log-normal random variables, $f(x) = x^2$ for Chi ran-
657 dom variables, etc. The goal of such aggregation, whether through exploiting
658 marginal stability or generalized f -means or other group operations in the ex-
659 ponential family, is to produce summary statistics for each group, R and S , in a
660 manner that permits generalized linear modeling of the summary statistics. By
661 ensuring summary statistics are also exponential-family random variables, one
662 can perform a factor-contrast style analysis as described above but only on the
663 two summary statistics and not on all s species. Doing so can greatly reduce
664 the computational load of phylofactorizing large datasets and, as we show be-
665 low, can increase the power of edge-identification even when the within-group
666 homogeneity assumption does not hold. Marginal stability, for the purposes of
667 phylofactorization, must be on the parameter of interest in generalized linear
668 modeling (Figure 3a).

669 Marginal stability opens up more distributions to stable aggregation. Pres-
670 ence absence data, for instance, can be assumed to be Bernoulli random vari-
671 ables. The assumption of within-group homogeneity for the probability of pres-
672 ence, ρ , allows addition of Bernoulli random variables within each group, R and
673 S , to yield a respective binomial random variable, x_R and x_S . Likewise, the ad-
674 dition of a set of binomial random variables with the same probability of success,
675 ρ , yields an aggregate binomial random variable. A set of exponential random
676 variables with the same rate parameter, λ , can be added to form a gamma ran-

677 dom variable. Gamma random variables, $x_i \sim \text{Gamma}(\kappa_i, \theta)$, parameterized by
678 their shape, κ_i , and scale, θ , are marginally stable on θ . Addition of geometric
679 random variables with the same rate parameter forms a negative binomial, and
680 the addition of a set of negative binomial random variables, $x_i \sim \text{NB}(\pi_i, \rho)$,
681 with the same probability of success ρ but different numbers of failures, π_i , can
682 be aggregated into $x_R = \sum_{i \in R} x_i$ where $x_R \sim \text{NB}(\sum_{i \in R} \pi_i, \rho)$. All of these
683 distributions are not stable, but they are marginally stable.

684 Marginally stable aggregation can be made efficient by matrix multiplication
685 onto one-vectors $\mathbf{1}_R$ and $\mathbf{1}_S$ whose i th entries are 1 for all $i \in R, S$, respectively,
686 and 0 otherwise. Assuming a Poisson or negative binomial count model for the
687 bird/non-bird data frame above, the data frame is reduced to

Site	Species	Abundance	z_1	z_2	phylo
1	Bird	18	1	.5	R
1	Non-Bird	5	1	.5	S
2	Bird	3	0	-2	R
2	Non-Bird	17	0	-2	S
...

688 and the same equation (17) can be used for phylofactorization through phylo
689 factor-contrasts.

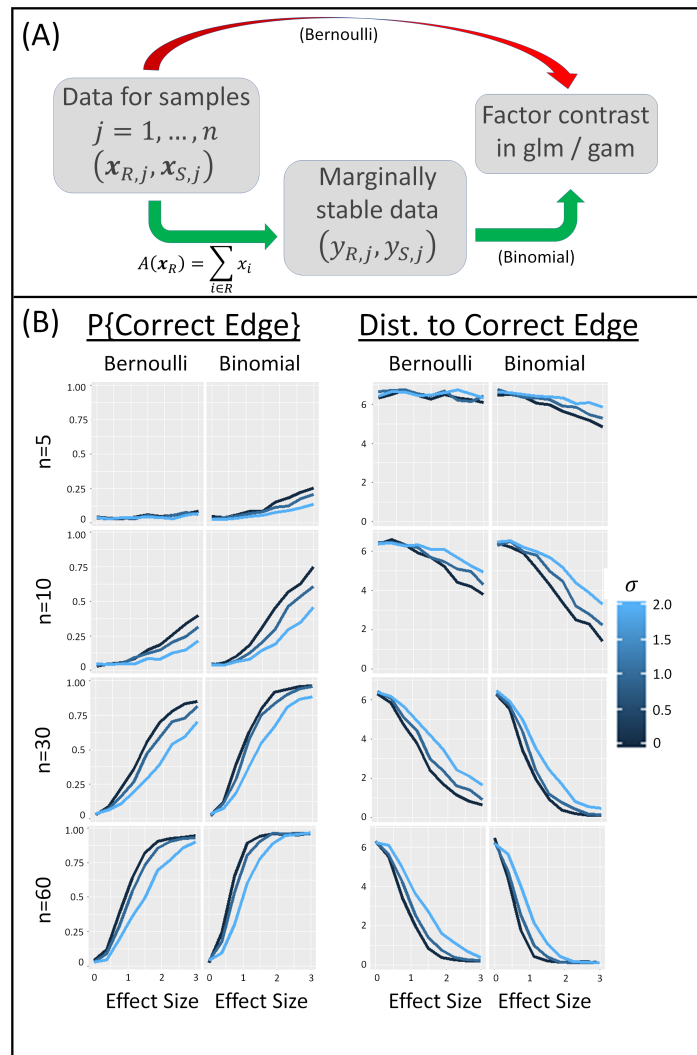


Figure 4: phylo factor contrasts can allow phylofactorization of exponential family random variables. (A) Each edge separates the species in a sample into two groups. These groups can be used as factors directly in a generalized linear model as in equations 17 and 21. Alternatively, a within-group homogeneity assumption can be used to aggregate data of many exponential family random variables to a marginally stable distribution, such as addition of Bernoulli random variables with the same probability of success to a binomial random variable. Regression on marginally stable random variables may dramatically reduce computational costs and, if within-group heterogeneity is low, improve accuracy. (B) Simulations of Bernoulli presence/absence data of 30 species with a random phylogeny suggest that aggregation to binomial improves power across a range of effect sizes, δ (x-axis), sample sizes, n (rows), and within-group heterogeneity, σ . Here, aggregation of presence-absence data to binomial random variables for subsequent factor-contrasts outperformed the raw factor contrast of Bernoulli presence/absence data, suggesting it is at least a viable tool for large datasets. The generality of improved power of regression on surrogate, marginally stable aggregates remains to be seen.

Aggregation to a marginally stable distribution is computationally efficient but will only outperform maximum-likelihood estimation if the within-group heterogeneity is small. For 700 replicates for each combination of sample size $n \in \{5, 10, 30, 60\}$, effect size $\delta \in \{0, 0.375, 0.75, 1.125, 1.5, 1.875, 2.25, 2.625, 3\}$, and within-group variance $\sigma \in \{0, 1, 2\}$, we simulated three explanatory variables $\{z_1, z_2, z_3\}$ as independent, identically distributed n -vectors of standard normal random variables. The log-odds of presence for individual i in group R or group S was modeled as

$$\begin{aligned}\eta_{R,i} &= z_1 + z_2 + \left(0.1 + \frac{\delta}{2}\right) z_3 + z_{4,i} \\ \eta_{S,i} &= z_1 - z_2 + \left(0.1 - \frac{\delta}{2}\right) z_3 + z_{4,i}\end{aligned}\tag{20}$$

690 where $z_{4,i} \stackrel{i.i.d.}{\sim} Gsn(0, \sigma^2)$ are independent Gaussian random variables particular
691 to the individual and sample. The data were either kept as Bernoulli random
692 variables or aggregated via summation to binomial random variables and then
693 analyzed using factor contrasts in a generalized linear model of the form

$$\eta = z_1 + \text{phylo} \times z_2 + \text{phylo} \times z_3.\tag{21}$$

694 The objective function was the deviance from the final term, $\text{phylo} \times z_3$. The
695 probability of identifying the correct edge and the distance between the iden-
696 tified and correct edge (in the number of nodes separating the two edges) are
697 plotted in Figure 4b. The method of factor-contrasts for glm-phylofactorization
698 asymptotically approaches perfect edge-identification, both in the probability of
699 detecting the correct edge and in distance from the correct edge, as the sample
700 sizes and effect sizes increase. Aggregation to binomial and subsequent factor-
701 contrast of the aggregates slightly improved the power of edge-identification in
702 these simulations. The improved accuracy of marginally-stable aggregation de-

703 creases with differences in within-group means, as opposed to an addition of
704 individual within-group variance through $z_{4,i}$, as illustrated below. However,
705 marginally-stable aggregation performs reasonably well and, crucially, scales
706 well with increasing numbers of species and sample size. Consequently, if the
707 datasets are large and the within-group homogeneity across samples is small,
708 marginally-stable aggregation and stepwise construction of factor contrasts may
709 be a useful tool for regression-phylofactorization of exponential family random
710 variables.

711 **Mixed Algorithm** Coefficient contrasts are computationally easy yet inac-
712 curate, whereas stepwise `phylo` factor selection (without marginally-stable ag-
713 gregation) is accurate yet computationally demanding (Figure 5). It's possible
714 to develop mixed algorithms with accuracy similar to stepwise `phylo` factor se-
715 lection and reduced computational costs more similar to coefficient contrasts
716 or marginally-stable aggregation. We present one example. In the first stages
717 of the algorithm, multivariate generalized linear modelling is performed as for
718 coefficient contrasts. For each iteration, coefficient contrasts (equation 16) are
719 used to narrow down the set of possible edges, $\{e\}_{top}$, to a set of edges with high
720 objective functions from standardized coefficient contrasts. We use the top 20%
721 of edges based on ω_e in equation 16, resulting in an approximately 80% speed-
722 up compared to the brute-force `phylo` factor contrast algorithm. For only these
723 edges, `phylo` factors are considered and the winning edge is the top-quantile
724 edge which maximizes the deviance of its `phylo` factor contrast.

725 **Algorithm comparison** We compare the performance of the four algorithms
726 listed above by testing how well they can correctly identify the affected edges,
727 $\{e_1^*, e_2^*\}$, and how long they take to extract a variable number of factors. The
728 four algorithms tested are:

- 729 • “B” : standardized coefficient contrasts
- 730 • “phylo”: Unaggregated phylo factor contrasts
- 731 • “mStable”: marginally-stable aggregation followed by phylo factor con-
732 trasts
- 733 • “mix”: Use of the “phylo” algorithm on only the top 20% of edges.

734 For edge identification, presence/absence data, $x_{i,j}$, were simulated for a set of
735 $s = 50$ species and $n = 40$ samples. The logit probability of all species was
736 modelled as

$$\eta_i \sim \beta_{i,0} + 0.1z_1 + 0.1z_2 \quad (22)$$

737 where $\beta_{i,0} \stackrel{i.i.d.}{\sim} N(0, 1)$ broke the within-group homogeneity in mean-probability
738 of presence/absence. For comparison, the case with $\beta_{0,i} = 0$ for all species
739 i is also considered. The other two explanatory variables, z_1 and z_2 , were
740 the partitioning variables differentiating species separated by edges. Two non-
741 nested clades, one containing 21 species and the other containing 5 species, had
742 a different association with the meta-data:

$$\eta_i \sim z_{0,i} - 0.2z_1 + 0.6z_2$$

743 for i in either of the two affected clades. To add an additional level of complexity,
744 the two meta-data variables were given multicollinearity by simulating $z_1 \sim$
745 $Gsn(0, 1)$ and $z_2 \sim Gsn(z_1, 1)$. The algorithms were run for two factors and the
746 number of correctly identified edges (out of 2) was tallied across 1000 replicates
747 (e.g. an algorithm that was 80% correct identified 1600 correct edges over 1000
748 replicates). The times for each of these algorithms to compute two factors
749 was also recorded. To compare the scaling of the algorithms, null data were
750 simulated across a range of species richness $m \in \{50, 100, 150, 200, 250, 300\}$

751 and across a range of factors $t \in \{1, 2, 3\}$.

752 The stepwise `phylo` factor contrasts by maximizing the total deviance of
753 `phylo * (z1 + z2)` had the greatest accuracy but also the slowest computation
754 time (Figure 5). The time required to compute `phylo` factor contrasts scale
755 quadratically with the number species, m , whereas coefficient contrasts and
756 marginally stable (`mStable`) aggregation scale linearly. Marginally stable ag-
757 gregation only performs well when $\beta_{i,0} = 0$ for all species, i , and when the
758 within-group heterogeneity is small. The accuracy of `phylo` factor contrasts
759 can be preserved and the computation time reduced by selecting the top 20% of
760 edges based on coefficient contrasts. The computation of multiple generalized
761 linear models across edges can be parallelized to reduce computation time, and
762 such parallelization is built into the R package `phylofactor`.

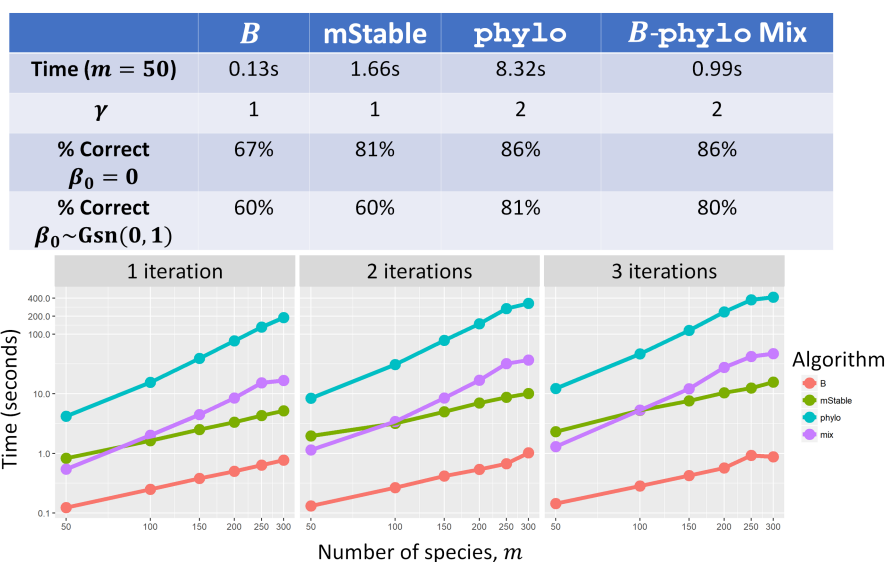


Figure 5: The accuracy, computation time and scaling of four algorithms for generalized phylofactorization. Algorithms are compared via the baseline time for two factors with $m = 50$ species, the scaling coefficient γ in $time \propto m^\gamma$, and percent of correctly identified edges in simulated data with $m = 50$ species and 2 affected clades. Stepwise phylo factor contrasts have high accuracy but are computationally costly and scale quadratically with the number of species. Marginally stable (mStable) aggregation scales linearly with m but only performs well when $\beta_0 = 0$. Computation time can be reduced and accuracy preserved if coefficient contrasts in equation 16 are used to narrow the set of edges considered for rigorous phylo factor contrasts.

763 **Summary of generalized phylofactorization** We have presented algo-
 764 rithms to perform regression-phylofactorization for non-Gaussian data. These
 765 algorithms can be called within the function `gpf()`. The stepwise selection of
 766 phylo factor contrasts is best able to correctly identify edges and is easily par-
 767 allelizable. The computation time of stepwise phylo factor contrasts can be
 768 reduced by narrowing the set of considered edges to those with high coefficient
 769 contrasts. Marginally stable aggregation may be a promising alternative for
 770 faster algorithms as it scales linearly with the number of species, but marginally
 771 stable aggregation only performs well when there is little systematic difference
 772 in the mean, $\beta_{i,0}$, across species, i .

773 There are fruitful avenues for future research to refine the algorithms for
774 phylofactorization of big-data consisting of non-Gaussian exponential family
775 random variables. These algorithms are intimately related to reduced rank
776 regression and generalized linear modelling with shared coefficients. Reduced-
777 rank regression considers a compact set of possible basis vectors and, conse-
778 quently, can use gradient-descent methods to find maximum-likelihood esti-
779 mates. The constrained set of allowable contrasts in the phylogeny precludes
780 gradient-descent and produces problems directly analogous to those in phylo-
781 genetic components analysis and thus we have focused on explicit testing of all
782 possible allowable contrasts in the phylogeny or, in the case of the mixed algo-
783 rithm, testing a subset of contrasts believed to contain the winning edge, e^* .
784 These methods can extend to generalized additive models and, as we discuss
785 below, spatial and time-series data as well.

786 **Phylogenetic factors of space and time**

787 Phylofactorization can also be used in analyses of spatial and temporal patterns.
788 We've demonstrated phylofactorization through examples of cross-sectional data
789 through two-sample tests, analyses of contrast-basis projections, and use of
790 `phylo` factor contrasts in communities sampled across a range of meta-data.
791 These same tools can be used for phylofactorization-based analysis of spatial
792 and temporal ecological data. Samples of a community over space and time can
793 be projected onto contrast basis elements and the resulting component scores,
794 y_e , can be analyzed much like PhyCA to identify the phylogenetic partitions
795 of community composition over space and time. Spatial samples can be an-
796 alyzed using `phylo` factor contrasts as defined for generalized linear models.
797 Multivariate Autoregressive Integrated Moving Average (ARIMA) models can
798 be constructed either as ARIMA models of the component scores, y_e , or as

799 multivariate ARIMA models with phylo factor contrasts as used in general-
800 ized linear models perform phylogenetic partitions based on differences in drift,
801 volatility, and other features of interest. Coefficient matrices, including spatial
802 and temporal autocorrelation matrices or coefficients for extrinsic meta-data \mathbf{Z} ,
803 can be approximated with phylogenetic contrast-bases as in equation (15).

804 Marginally stable aggregation in spatial and temporal data requires a more
805 complex consideration of the marginal stability of spatially explicit random vari-
806 able and stochastic processes. “Stability”, for spatially and temporally explicit
807 random variables, must preserve the underlying model for the spatial or tem-
808 poral process being used for analysis. An example of a less obvious marginally
809 stable aggregation of time-series data is the stability of neutral drift (sensu
810 Hubbell [22]) to grouping.

811 Neutral communities fluctuate, and those fluctuations have a drift and volatil-
812 ity unique to neutral drift. Neutral drift can also be defined either by discrete,
813 finite-community size urn processes or stochastic differential equations for the
814 continuous approximations of finite but large communities. Recently, Wash-
815 burne et al. [50] articulated the importance of a feature of neutral drift which
816 enables time-series neutrality tests: its invariance to grouping of species. If a
817 stochastic process of relative abundances, \mathbf{X}_t , obeys the probability law defined
818 by neutral drift, then any complete, disjoint groupings of \mathbf{X}_t also obeys the
819 probability law for a lower-dimensional neutral drift. Thus, neutral processes
820 are stable to aggregation by grouping or summation of relative abundances. Col-
821 lapsing all species into two disjoint groups, R and S , yields a two-dimensional
822 neutral drift with a well-defined neutrality test for time-series data. Specifically,
823 if \mathbf{X}_t is a Wright Fisher process and R and S are disjoint groups whose union

824 is the entire community, the quantity

$$\nu_t = \arcsin \left(\left(\sum_{i \in R} X_{i,t} \right) - \left(\sum_{j \in S} X_{j,t} \right) \right) \quad (23)$$

has a constant volatility which can be used to define a neutrality test for time-series data. Thus, phylofactorization can be done to partition edges across which the dynamics appear to be the least neutral. For the test developed by Washburne et al., the aggregation operation is the L_1 norm and the contrast operation is subtraction:

$$\begin{aligned} A(\mathbf{x}_R) &= |\mathbf{x}_R| \\ C(A(\mathbf{x}_R), A(\mathbf{x}_S)) &= \arcsin(A(\mathbf{x}_R) - A(\mathbf{x}_S)) \end{aligned} \quad (24)$$

825 and the objective function, ω , for edge e is the test-statistic of a homoskedasticity
826 test of C_e . Neutrality is a relative measure - biological units are neutral relative
827 to one-another - and thus the use of aggregation of species into a unit and a
828 contrast of two units is a natural connection between the theory and operations
829 of phylofactorization and the concept of neutrality.

830 **Statistical Challenges**

831 We present a unifying algorithm which partition organisms into functional groups
832 by identifying meaningful differences or contrasts along edges in the phylogeny.
833 Phylofactorization is formally defined as a graph-partitioning algorithm. How-
834 ever, maximizing the variance of the data projected onto contrast basis elements
835 corresponding to edges in the phylogeny is a constrained principal components
836 analysis. The use of regression-based objective functions and the iterative con-
837 struction of a low-rank approximation of a data matrix is similar to factor
838 analysis. The discovery of a sequence of orthogonal factor contrasts in general-

839 ized linear models is a form of stepwise/hierarchical regression and partitioning
840 a coefficient matrix \mathbf{B} is a reduced-rank regression method. The maximization
841 of the objective function at each iteration is a greedy algorithm. Each of these
842 connections between phylofactorization and other classes of methods produces a
843 body of literature from related methods which could inform phylofactorization
844 and facilitate rapid development of this exploratory tool into a more robust,
845 inferential one.

846 There are statistical challenges common across many methods for phylofac-
847 torization. In this section, we enumerate some of the statistical challenges and
848 discuss work that has been done so far. First, as with any method using the phy-
849 logeny as a scaffold for creating variables or making inferences, the uncertainty
850 of the phylogeny and the common use of multiple equally likely phylogenies war-
851 rant consideration and further method development. Other challenges discussed
852 here are: understanding the propagation of error; development of Metropolis al-
853 gorithms to better arrive at global maxima; the appropriateness, and error rates,
854 of phylofactorization under various evolutionary models underlying the effects
855 (e.g. trait differences, habitat associations, etc.) and residuals in our data;
856 understanding graph-topological biases and confidence regions; cross-validating
857 the partitions and inferences from phylofactorization; determining the appropri-
858 ate number of factors and stopping criteria to stop a running phylofactorization
859 algorithm; and understanding the null distribution of test-statistics when objec-
860 tive functions being maximized are themselves test-statistics from a well-known
861 distribution. Any exploratory data analysis tool can be made into an inferential
862 tool with appropriate understanding of its behavior under a null hypothesis,
863 and the connections of phylofactorization to related methods can accelerate the
864 development of well-calibrated statistical tests for phylogenetic factors.

865 **Phylogenetic inference** So far we have assumed that the phylogeny is known
866 and error free, but the true evolutionary history is not known - it is estimated.
867 Consequently, phylofactorizations are making inferences on an uncertain scaf-
868 fold; the more certain the scaffold, the more certain our inferences about a
869 clade. Two challenges remain for dealing with phylofactorization on an uncer-
870 tain phylogeny. For a consensus tree, there is the question of what statistics of
871 the consensus are most easily integrated for precise statements of uncertainty
872 in phylofactorization inferences. Bootstrapped confidence limits for monophyly
873 [12] are the most commonly used statement of uncertainty for a consensus tree,
874 but there may be others as well. Different organisms will have different leverages
875 in regression or two-sample test phylofactorization, and thus monophyly is only
876 part of the picture: leverage is another. For a set of equally likely bootstrapped
877 trees, there is a need to integrate phylofactorization across trees. Phylofactor-
878 ization of sets of equally likely phylogenies has not yet been done, but may be
879 a fruitful avenue for future research. One last option for researchers with trees
880 containing clades with low bootstrap monophyly is to lower the resolution of the
881 tree. Phylofactorization can still be performed on a tree with polytomies - the
882 mammalian phylogeny used above contained many - and reducing the number
883 of edges considered at each iteration can focus statistical effort (and chances of
884 false-discovery) on clades about which the researcher is more certain.

885 **Propagation of error** Phylofactorization is a greedy algorithm. Like any
886 greedy algorithm, the deterministic application of phylofactorization is non-
887 recoverable. Choosing the incorrect edge at one iteration can cause error to
888 propagate, potentially leading to decreased reliability of downstream edges. Lit-
889 tle research has been done towards managing the propagation of error in phylo-
890 factorization, but recognizing the method as a greedy algorithm suggests options
891 for improving performance. Stochastic-optimization schemes, such as replicate

892 phylofactorizations using Metropolis algorithms and stochastic sampling as im-
893 plemented in the mammalian tree phylofactorization (sampling of edges with
894 probabilities increasing monotonically with ω_e and picking the phylofactor ob-
895 ject which maximizes a global objective function), may reduce the risk of error
896 cascades in phylofactorization [20].

897 **Behavior under various evolutionary models** Phylofactorization is hy-
898 pothesized to work well under a punctuated-equilibrium model of evolution or
899 jump-diffusion processes [15, 26] in which jumps are infrequent and large, such
900 as the evolution of vertebrates to land or water. If few edges have large changes
901 in functional ecological traits underlying the pattern of interest, phylofactor-
902 ization is hypothesized to work well. Phylofactorization may also work well
903 when infrequent life-history traits arise or evolutionary events occur (such as
904 ecological release) along edges and don't yield an obvious trait but instead yield
905 a correlated, directional evolution among descendants. Phylofactorization of
906 mammalian body sizes yielded a scenario hypothesized to be in this category.
907 In this case the exact trait may not have arisen along the edge identified, but a
908 precursor trait, or a chance event such as extinctions or the emergence of novel
909 niches, may precipitate downstream evolution of the traits underlying phylofac-
910 torization. Both aggregation and contrast functions can incorporate phyloge-
911 netic structure and edge lengths to partition the tree based on likelihoods of
912 such evolutionary models. The sensitivity of phylofactorization to alternative
913 models, such as continuous Brownian motion and Ornstein-Uhlenbeck models
914 commonly used in phylogenetic comparative methods [13, 19], remains to be
915 tested and will likely vary depending on the particular method used.

916 **Basal/distal biases** Researchers may be interested in the distribution of fac-
917 tored edges in the tree. If a dataset of microbial abundances in response to

918 antibiotics is analyzed by regression-phylofactorization and results in many tips
919 being selected, a researcher may be interested in quantifying the probability of
920 drawing a certain number of tips given t iterations of phylofactorization. Alter-
921 natively, if several edges are drawn in close proximity researchers may wonder
922 the probability of drawing such clustered edges under a null model of phylofac-
923 torization. For another example, researchers may wonder if the number of im-
924 portant functional ecological traits arose in a particular historical time window
925 (e.g. due to some hypothesis of important evolutionary event or environmental
926 change), and thus want to test the probability of drawing as many or more
927 edges than observed under a null model of phylofactorization. All of these tests
928 would require an accurate understanding of the probability of drawing edges in
929 different locations of the tree.

930 All methods described here, save the Fisher exact test, have a bias for tips in
931 the phylogeny (Figure 6). Such biases affect the calibration of statistical tests
932 of the location of phylogenetic factors, such as a test of whether/not there is
933 an unusually large number of differentiating edges in mammalian body mass
934 during or after the K-Pg extinction event.

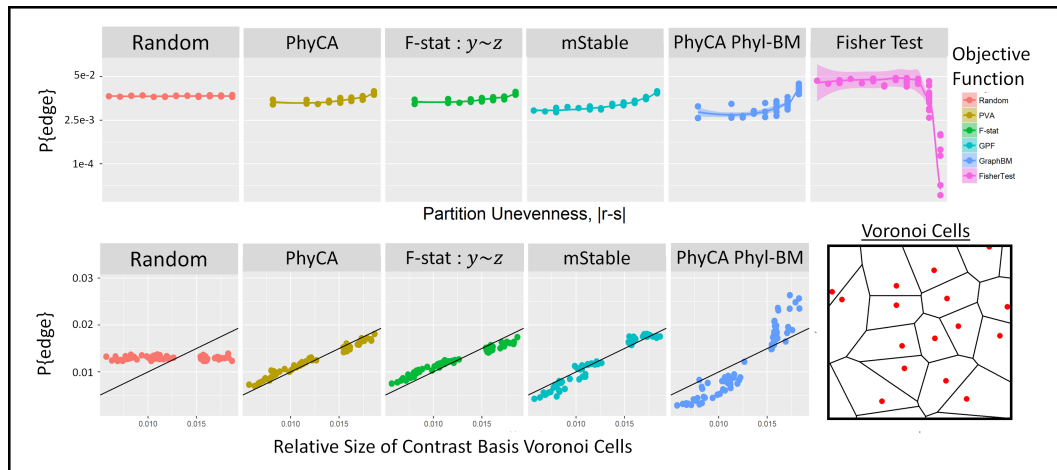


Figure 6: Graph topological bias in null data and the relative size of Voronoi cells of contrast basis elements. The method and the null distribution of the data determine graph-topological bias of phylofactorization. A random draw of edges does not discriminate against edges based on the relative sizes of two groups contrasted by the edge, but 16,000 replicate phylofactorizations of null data reveal that contrast-basis methods are slightly biased towards uneven splits (e.g. tips of the phylogeny). Standard Gaussian null data were used for PhyCA, F-statistics from regression on contrast basis elements ($y_e \sim z$), and binomial null data was used for generalized phylofactorization (gpf) through marginally-stable aggregation. Other methods, such as Fisher’s exact test of a vector of Bernoulli random variables, have opposite biases. The tip-bias of contrast-basis analysis is amplified for marginal-stable aggregation in generalized phylofactorization, and amplified even more if the null data have residual structure from a Brownian motion diffusion along the phylogeny (Phyl-BM). The common bias when using contrast bases across a range of objective functions is related to the uneven relative sizes of Voronoi cells produced by the bases, simulated here by equation (25).

935 Phylofactorization using the contrast basis is biased towards the tips of
 936 the tree. Some progress can be made towards understanding the source of
 937 basal/distal biases in phylofactorization via the contrast-basis. The biases from
 938 analyses of contrast basis coordinates, \mathbf{y}_e , stem from a common feature of the
 939 set of K_t candidate basis elements $\{\mathbf{v}_{C_e}\}_{e=1}^{K_t}$ considered at iteration t of phylo-
 940 factorization. For the example of the t-test phylofactorization of a vector of
 941 data, \mathbf{x} , the winning edge e^* is

$$e^* = \operatorname{argmax}_e |\mathbf{v}_{C_e}^T \mathbf{x}|. \quad (25)$$

942 If all basis elements have unit norm, which they do under equation (5), then
943 each basis element being considered corresponds to a point on an m -dimensional
944 unit hypersphere. If the data, \mathbf{x} , are drawn at random, such that no direction
945 is favored over another, the probability that a particular edge e is the winning
946 edge is proportional to the relative size of its Voronoi cell on the surface of the
947 unit m -hypersphere. Thus, the basal/distal biases for contrast-basis analyses
948 with null data assumed to be drawn from a random direction can be boiled
949 down to calculating or computing the relative sizes of Voronoi cells. For our
950 simulation, we estimated the size of Voronoi cells through matrix multiplication

$$\mathbf{Y}_{null} = \mathbf{V}^T \mathbf{X}_{null} \quad (26)$$

951 where \mathbf{V} is a matrix whose columns j is the contrast basis elements for edge
952 e_j being considered and \mathbf{X}_{null} is the dataset simulated under the null model
953 of choice whose columns are independent samples \mathbf{x}_j . Each column of \mathbf{Y}_{null}
954 contains the projections of a single random vector - the element of each column
955 with the largest absolute value is the edge closest to that random vector.

956 **Graph-topology and confidence regions** As a graph-partitioning algo-
957 rithm, phylofactorization invites a novel description of confidence regions over
958 the phylogeny. The graph-topology of our inferences - edges, and their proximity
959 to other edges, both on the phylogeny and in the m -dimensional hypersphere
960 discussed above - can be used to refine our statements of uncertainty. 95%
961 Confidence intervals for an estimate, e.g. the sample mean, give bounds within
962 which the true value is likely to fall 95% of the time in random draws of the
963 estimate. Confidence regions are multi-dimensional extensions of confidence
964 intervals. Conceptually, it's possible to make similar statements regarding phy-
965 logenetic factors - confidence regions on a graph indicating the regions in which

966 the true, differentiating edge is likely to be.

967 Extending the concept of confidence regions to the graph-topological infer-
968 ences from phylofactorization requires useful notions of distance and “regions”
969 in graphs. One example of such a distance between two edges is a walking
970 distance: the number of nodes one crosses along the geodesic path between
971 two edges. Alternatively, one could define regions in terms of years or branch-
972 lengths. Defining confidence regions in phylofactorization must combine the
973 uneven Voronoi cell sizes as well as the geometry of the contrast basis. For
974 low effect sizes, confidence regions extend to distant edges on the graph whose
975 contrast basis have a large relative Voronoi cell size (e.g. the tips). As the effect
976 sizes increase, confidence regions over the graph are better described in terms of
977 angular distances between the contrast basis elements and that of the winning
978 edge, e^* (Figure 7).

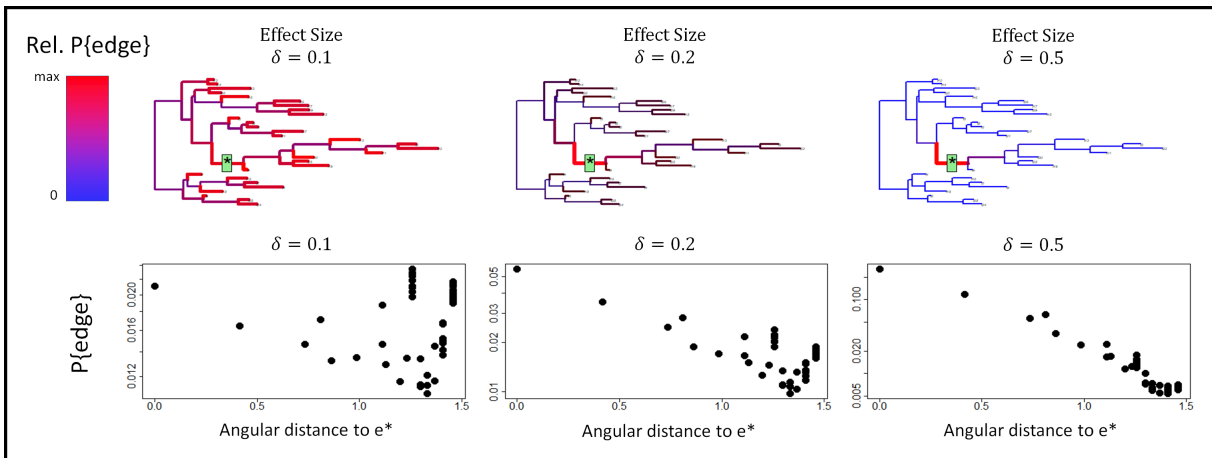


Figure 7: **Graph-topological confidence regions for phylofactorization.** Confidence regions around inferred edges must use distances relevant to the method and graph topology. A tree with 30 species was given a fixed effect about edge e^* in their mean values as a function of meta-data $z \sim Gsn(\pm\delta/2, 1)$. 7×10^5 iterations of phylofactorization were run and the relative probability of drawing each edge was visualized through both the color and width of the edge. The relationship between the angular distance of an edge's contrast basis element to that of e^* and the probability of drawing the edge indicate that for low effects, confidence intervals must incorporate a mix of tip-bias and angular distance, but larger effect sizes, in which the edge drawn is reliably in the neighborhood of e^* , the angular distance of contrast basis elements capture confidence regions around the location of inferred phylogenetic factors.

979 **Cross-validation** How do we compare phylofactorization across datasets to
980 cross-validate our results? If a researcher observes a pattern in the ratio of
981 squamates to mammalian abundances in North America, say a decrease in the
982 ratio of lizard and snake to mammal abundance with increasing altitude, they
983 may wish to cross-validate their findings in other regions, including regions with
984 few or none of the same species in the original study. Researchers replicating
985 the study in Australia and New Zealand would have to grapple with whether
986 or not to include monotremes in their grouping of “mammals” and whether or
987 not to include the tuatara, a close relative of squamates, in their grouping of
988 “squamates” - such branches were basal to the squamate & mammalian clades
989 contrasted in the hypothetical North American study.

990 Phylofactorization formalizes the issues arising with such phylogenetic cross-
991 validation (Figure 8). If all species in the training/testing datasets can be
992 located on a universal phylogeny, phylofactorization of a training set of species
993 and data identifies edges or links of edges in the training phylogeny which are
994 guaranteed to correspond to edges or links of edges in the universal phylogeny.
995 The testing set of species may introduce new edges to the phylogeny which
996 interrupt the links of edges in the universal phylogeny along which training
997 contrasts were conducted. In the example above, the tuatara and monotremes
998 all interrupt the link of edges separating North American mammals from North
999 American reptiles on the universal phylogeny.

1000 Robust cross-validation for phylofactorization requires directly addressing
1001 the issues arising from the interruptions of edges produced by novel species.
1002 Interruptions may be either ignored, or used to refine the inference. Returning
1003 to the previous example, one can use the presence of monotremes and tuatara to
1004 refine the definition of North American mammals to mean “all mammals” and
1005 “all placental and marsupial mammals”, and likewise one can optionally refine
1006 the definition of “squamates” to the broader “Lepidosauria” clade.

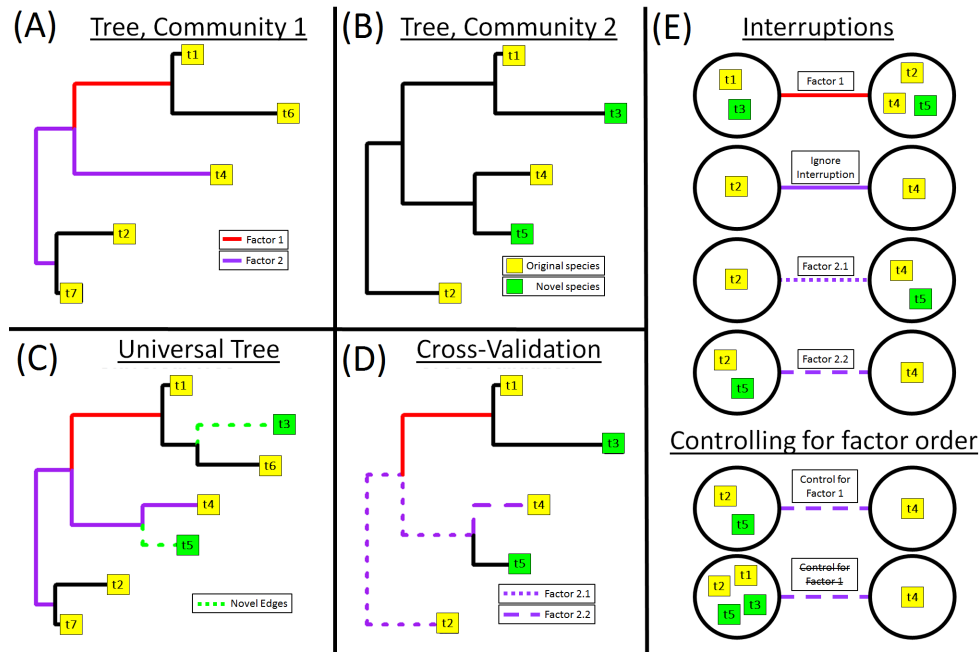


Figure 8: Graph-topological considerations with cross-validation. (A) The training community has 5 species (yellow boxes) split into two factors. The second factor forms a partition separating t_4 from $\{t_2, t_7\}$. The second factor does not correspond to a single edge, but instead a chain of two edges. (B) A second, testing community is missing species t_6 and t_7 and contains novel species t_3 and t_5 (green boxes). (C) All factors can be mapped to chains of edges on a universal phylogeny. Novel species “interrupt” edges in the original tree; cross-validation requires deciding what to do with novel species and interrupted edges. Species t_3 does not interrupt a factored edge, and so t_3 can be reliably grouped with t_1 in factor 1. However, species t_5 interrupts one of the edges in the edge-path of factor 2. (D-E) Interruptions can be ignored, or they can be used to refine the location of important edges (illustrated in Factor 2.1 and Factor 2.2). Another topological and statistical question is whether/not to control for factor order. For instance, controlling for factor order with Factor 2.2 would partition t_4 from $\{t_2, t_5\}$. Not controlling for factor order would partition t_4 from $\{t_1, t_2, t_3, t_5\}$.

1007 **Stopping Criteria** Often, it’s desirable to obtain a minimal set of partitions
 1008 to prioritize findings, simplify high-dimensional data, and focus effort on more
 1009 certain inferences. Doing so requires a method for stopping phylofactorization.
 1010 There are two broad options for stopping phylofactorization: a stopping func-
 1011 tion demonstrated to be sufficiently conservative, and null simulations allowing
 1012 quantile-based cutoffs (e.g. stop phylofactorization when the percent variance

1013 explained by PhyCA is within the 95% quantile of null phylofactorizations).
1014 Null simulations may allow statistical statements stemming from a clear null
1015 model, but stopping criteria can be far more computationally efficient and can
1016 be constructed to be conservative.

1017 Washburne et al. [51] proposed a stopping criterion for regression phylofac-
1018 torization which extends to all methods of phylofactorization using an objective
1019 function that is a test-statistic whose null-distribution is known. The original
1020 stopping criterion is based on the fact that, if the null hypothesis is true, the
1021 distribution of P-values from multiple hypothesis tests is uniform. Phylofactor-
1022 ization performs multiple hypothesis tests at each iteration. At each iteration,
1023 one can perform a one-tailed KS test on the uniformity of the distribution of the
1024 P-values from the test-statistics on each edge; if the KS-test is non-significant,
1025 stop phylofactorization. KS-test stopping criteria can conservatively stop simu-
1026 lations at the appropriate number of factors when there is a discrete subset of
1027 edges with effects. Such a method performs similarly to Horn's stopping crite-
1028 rion for factor analysis [21], whereby one stops factorization when the scree plot
1029 from the data crosses that expected from null data (figure 9). It's also possible
1030 to first use a stopping criterion and subsequently run null simulations to under-
1031 stand the likelihood of observed results under a null model of the researcher's
1032 choice (figure 9).

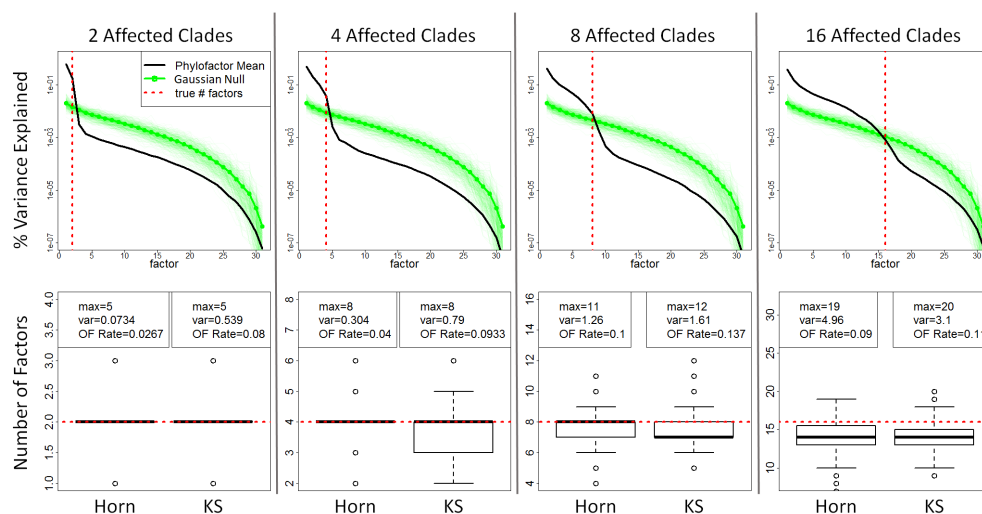


Figure 9: Null simulations and stopping criteria. A challenge of phylofactorization is determining the number of factors, K , to include in an analysis. Null simulations allow quantile-based cutoffs such as those in Horn’s parallel analysis from factor analysis. Stopping criteria stop phylofactorization using features available during phylofactorization of the observed data. Abundances of $m = 32$ species across $n = 10$ samples were simulated as i.i.d. standard Gaussian random variables. A set of u clades were associated with environmental meta-data, \mathbf{z} , where $z_j \stackrel{i.i.d.}{\sim} Gsn(0, 1)$. Regression-phylofactorization on the contrast-basis scores y_e was performed on 300 datasets for each $u \in \{2, 4, 8, 16\}$ and on data with and without effects. The objective function was the total variance explained by regression $y_e \sim z$. (**top row**) The percent of the variance in the dataset explained at each factor (EV) decreases with factor, t , and the mean EV curve for data with u affected clades intersects the mean EV curve for null data near where $t = u$, motivating a stopping criterion (Horn) based on phylofactorization of null datasets. (**bottom row**) The Horn stopping criterion has a slightly lower over-factorization (OF) rate than the standard KS stopping criterion (where OF rate is the fraction of the 300 phylofactorizations of data with simulated effects in which $t > u$). However, the algorithms were not extremely different and both criteria can be modified to be made more conservative. The KS stopping criterion is far less computationally intensive for large datasets as it requires running phylofactorization only once. Null simulations, however, can allow inferential statistical statements regarding the null distribution of test statistics in phylofactorization.

Calibrating Statistical Tests for ω_{e^*} Often, the objective function for the winning edge in phylofactorization, ω_{e^*} , corresponds directly to a common test-statistic. Applying a standard test for the resultant test-statistic, however, will lead to a high false-positive rate and an over-estimation of the significance of an effect, as the statistic was drawn as the best of many. Even when using

a test-statistic not equal to the objective function, researchers should be cautious of dependence between their test-statistic and the objective function as a possible source of high false-positive rates. Two methods for calibrating, or making conservative, statistical tests of ω_{e^*} are multiple-comparisons corrections to control a family-wise error rate (or other multiple-hypothesis-test methods) or conservative bounds on the distribution of the maximum of many independent, identically distributed statistics. For example, if each edge of K_t edges considered at iteration t resulted in an independent F -statistic, F_e , then the distribution of the maximum F -statistics, F_{e^*} , is

$$\begin{aligned} P\{F_{e^*} > F\} &= P\{F_{e_1} > F \cap F_{e_2} > F \cap \dots \cap F_{e_K}\} \\ &= P\{F_e > F\}^{K_t}. \end{aligned} \tag{27}$$

1033 Such an approximation may be used to yield conservative estimates, but the
1034 F -statistics are not independent and thus more nuanced analyses are needed for
1035 well-calibrated statistical tests. More research is needed to obtain conservative
1036 bounds on test-statistics in phylofactorization.

1037 **Summary of limitations** Phylofactorization can be a reliable statistical tool
1038 with a careful understanding of the statistical challenges inherent in the method
1039 and shared with related methods such as graph-partitioning, greedy algorithms,
1040 factor analysis, and the use of a constrained, biased basis for matrix factoriza-
1041 tion. Phylofactorization can first and easiest be an exploratory tool, but all
1042 exploratory tools can be made inferential with suitable understanding of their
1043 behavior under an appropriate null model. For example, principal components
1044 analysis was and still is primarily an exploratory tool, but the discovery of the
1045 Marcenko-Pastur distribution [30] has improved the calibration of statistical
1046 tests on principal components for standardized, mean-centered data. Improved

1047 understanding of how uncertainties in phylogenetic inference translate to uncer-
1048 tainties in phylofactorization, conservative stopping criteria, null distributions
1049 of test-statistics for winning edges, propagation of error and stochastic sampling
1050 algorithms to avoid deterministic ruts, graph-topological biases and confidence
1051 regions on a graph, can all improve the reliability of phylofactorization as an
1052 inferential tool.

1053 While phylofactorization was built with an evolutionary model of punctu-
1054 ated equilibria in mind, it may also work well under other evolutionary models
1055 such as correlated evolution among descendants of an edge. There are also many
1056 evolutionary models under which phylofactorization does not perform well. For
1057 instance the graph-topological biases of PhyCA are increased under a Brownian
1058 motion model of evolution. All statistical tools operate well under appropriate
1059 assumptions, and understanding the assumptions, as well as the known limita-
1060 tions, are necessary for responsible and academically fruitful use of statistical
1061 tools like phylofactorization.

1062 Discussion

1063 Functional ecological traits underlie many observed patterns in ecology, includ-
1064 ing species abundances, presence/absence of species, and responses of traits
1065 or abundances to experimental conditions or along environmental gradients.
1066 Where the ecological pattern of interest is associated with heritable traits, the
1067 phylogeny provides a scaffold for the discovery of functional groupings of clades
1068 underlying the ecological pattern of interest. Traits arise along edges, and con-
1069 trasting taxa on opposing sides of an edge allows one to uncover edges best
1070 separating species with different functional associations or links to the ecologi-
1071 cal pattern. By noting that each edge partitions the phylogeny into two disjoint

1072 sets of species, by generalizing the operations of “grouping” - aggregating and
1073 contrasting disjoint sets of species - and by defining the objective function of
1074 interest (the pattern), we have proposed a universal method for identifying rel-
1075 evant phylogenetic scales in ecological datasets.

1076 Phylofactorization is a graph-partitioning algorithm intended to separate the
1077 phylogeny into binned phylogenetic units with a combination of high within-
1078 group similarity and high between-group differences. Two-sample tests are a
1079 natural method for making such partitions in vectors of data; such partitions
1080 can also be made with ancestral state reconstruction. The quantities used in
1081 two-sample tests can be extended to larger, real-valued datasets by analysis
1082 of a contrast basis. Objective functions for choosing the appropriate contrast
1083 basis include maximizing variance - a phylogenetic analog of principal com-
1084 ponents analysis - maximizing explained variance from regression, maximizing
1085 F-statistics from regression, and more. By partitioning coefficient matrices and
1086 using `phylo` factor contrasts, phylofactorization can be extended to generalized
1087 linear models, generalized additive models, and analyses of spatial and temporal
1088 patterns in ecological data.

1089 We’ve illustrated that two-sample tests can partition a dataset of mam-
1090 malian body mass into groups with very different average body masses. We’ve
1091 demonstrated that maximizing variance of data projected onto a contrast basis
1092 can identify major clades of bacteria in human feces that have been known, at
1093 a coarser resolution, to be highly variable and determined that one of the top
1094 phylogenetic factors in the American Gut dataset is a clade of Gammaproteobac-
1095 teria associated with IBD and used recently in an effort to diagnose patients
1096 with Crohn’s disease. We’ve shown that analyses of contrast bases can use non-
1097 linear regression, and within minutes of analysis on a laptop found a natural
1098 way put over 3,000 species into 5 binned phylogenetic units, sort them along

1099 an axis of the dominant explanatory variable, and produce a simplified story of
1100 how community composition changes in Central Park soils.

1101 One can also perform phylofactorization when doing maximum-likelihood
1102 regression of exponential family random variables. The coefficient matrix can be
1103 approximated using the contrast basis, resulting in a phylogenetically-interpretable
1104 reduced-rank regression. Alternativley, it's possible to use `phylo` factor con-
1105 trasts for a shared-coefficients model and maximum-likelihood based selection
1106 of edges for partitioning. One can either perform the factor contrasts on the
1107 raw data, or, for many exponential family random variables, one can aggregate
1108 the data from each group to a marginally stable distribution for more compu-
1109 tationally efficient factor contrasts. These methods can be extended to spatial
1110 and temporal data. All methods discussed here can be implemented with the
1111 R package “`phylofactor`”, and scripts for running all analyses in this paper are
1112 available in the supplemental materials.

1113 As with any method, there are limitations to be aware of. First, the general
1114 problem of separating species into k bins that maximize a global objective func-
1115 tion is an NP hard problem. Second, like any greedy algorithm, purely deter-
1116 ministic phylofactorization may fall into ruts and errors in one step might prop-
1117 agate into downstream inferences. Third, the null distribution of test-statistics
1118 resulting from phylofactorization is not known; the resultant test statistics are
1119 biased towards extreme values. Null simulations, conservative stopping func-
1120 tions, and/or extremely stringent multiple comparisons corrections can be used
1121 to make inferences through phylofactorization while maintaining conservative
1122 bounds in family-wise error or false-discovery rate. When the objective func-
1123 tion being maximized is also a test-statistic with a well-defined null distribution,
1124 one-sided KS-tests of the P-values from the test-statistic can serve as a computa-
1125 tionally efficient and conservative stopping function. Fourth, common objective

1126 functions using the contrast basis will be biased due to the unequal relative
1127 sizes of the Voronoi cells of the contrast basis elements in the unit hypersphere
1128 in which they lie, with contrast basis elements corresponding to tips of the
1129 phylogeny tending to have larger relative Voronoi cell size than contrast basis
1130 elements corresponding to interior edges. Understanding the graph-topology
1131 of errors can assist the description of graph-topological confidence regions for
1132 each inference. Finally, phylofactorization formalizes the logic and challenges of
1133 cross-validating ecological comparisons even when the training and testing sets
1134 of species are completely disjoint. Many of these limitations may be resolved
1135 with future work, allowing the general algorithm and its common implementa-
1136 tions to become a reliable, well-calibrated inferential tool.

1137 Phylofactorization can objectively identify phylogenetic scales for ecologi-
1138 cal big-data and instantly produce avenues for future natural history research.
1139 By iteratively identifying clades, phylofactorization provides a sequence of low-
1140 rank approximations of a dataset, such as that visualized in figure 3c, which
1141 correspond to groups of species with a shared evolutionary history. What traits
1142 characterize the Chloracidobacteria which don't like acidic soils? What traits
1143 characterize the monophyletic clade of Gammaproteobacteria that are associ-
1144 ated with IBD? What traits underlie the Clostridia/Erysipelotrichi being such
1145 variable species in the American gut? The low-rank approximations of eco-
1146 logical data obtained by phylofactorization motivate subsequent questions best
1147 answered by life history comparisons, comparative genomics, microbial phys-
1148 iological studies, and other avenues of future research contrasting the species
1149 partitioned.

1150 **Relation to other phylogenetic methods** Phylofactorization is proposed
1151 amidst an explosion of literature in phylogenetic comparative methods and var-
1152 ious other phylogenetic methods for analyzing ecological datasets [29, 38, 14],

1153 and some careful thinking is beneficial to clarify the distinctions between the
1154 myriad methods.

1155 Phylogenetic generalized least squares [16] aims to control for residual struc-
1156 ture in the response variable expected under a model of trait evolution, and
1157 is thus used when performing regression on a trait, whereas phylofactorization
1158 aims to partition observed trait values or abundances into groups, separated by
1159 edges, with different means or associations with meta-data. Thus, while meth-
1160 ods of phylogenetic signal, such as Pagel's λ [35] or Blomberg's κ [5], summarize
1161 global patterns of phylogenetic signal by parameterizing the extent to which a
1162 particular model of evolution can be assumed to underlie the residual structure
1163 of observed traits (often for downstream use in PGLS), phylofactorization it-
1164 eratively identifies precise locations of putative changes and precise locations
1165 partitioning phylogenetic signal or structure.

1166 Phylofactorization can be implemented by a contrast of ancestral state re-
1167 constructions of nodes separated by edges, for example by looking for edges with
1168 nodes whose reconstructed ancestral states are most different, but is limited by
1169 disallowing the descendant clade of an edge to impact the ancestral state of the
1170 edge's basal node - a proper non-overlapping contrast would separate the groups
1171 of species being used to reconstruct each node, and thus phylofactorization can
1172 be implemented with ancestral state reconstruction under the assumption of
1173 time-reversible evolutionary models.

1174 Phylogenetically independent contrasts [13] produces variables correspond-
1175 ing to contrasts of descendants from each node, whereas phylofactorization uses
1176 contrasts of species separated by an edge, picks out the best edge, splits the tree,
1177 and repeats. Phylofactorization develops a set of variables and an orthonormal
1178 basis to describe ecological data, but limits itself to bases interpretable as non-
1179 overlapping contrasts along edges; eigenvectors of phylogenetic distances matri-

1180 ces or covariance matrices under diffusion models of traits [35], are not encom-
1181 passed in phylofactorization as they do not construct non-overlapping contrasts
1182 along edges. Such eigenvector methods construct quantities whose evolutionary
1183 interpretation is less clear. Unlike many modern methods for re-defining dis-
1184 tances, such as UniFrac distances [29] or phylogenetically-defined inner prod-
1185 ucts [38], phylofactorization is principally about discovering phylogenetically-
1186 interpretable directions - vectors which characterize primary axes of variation
1187 in the community and represented through the contrast basis, a multilevel-factor
1188 developed from stepwise selection of factor contrasts, or a basis made of aggre-
1189 gations of the binned phylogenetic units.

1190 **Phylofactorization as a species concept** There is great debate about what
1191 constitutes a species in microbes, let alone all organisms. There is a need for
1192 objectivity and universality in the definition of “species” and other units in
1193 ecology and evolution. The biological species concept is complicated by asex-
1194 ual reproduction. Genetic species concepts are limited by the subjectivity of a
1195 sequence-similarity cutoff, such as the 97% sequence similarity commonly used
1196 in defining operational taxonomic units or OTUs, which is additionally compli-
1197 cated by the fact that functional ecological similarity may not be uniform at
1198 a given sequence-similarity cutoff. Ecological species concepts are often useful
1199 once researchers have a clear sense of the functional ecological groups, but it is
1200 difficult to objectively define what constitutes an important functional ecologi-
1201 cal group, especially for taxa whose life histories are unknown. Species concepts
1202 coarse-grain the diversity of life in a way that connects our coarse-grained units
1203 to biological, ecological, and evolutionary theory. To that end, phylofactoriza-
1204 tion can be seen as defining a species concept.

1205 Species concepts are fundamental to biology as they partition the diversity of
1206 life into units between which we define ecological interactions and within which

1207 we define evolution and natural selection. At the heart of species concepts are
1208 the operations fundamental to phylofactorization: aggregation, contrast, and an
1209 objective function. Species are aggregations of finer units of diversity: individual
1210 subpopulations of individual organisms and their individual cells and the cells'
1211 individual genes are all aggregated to define a "population". Aggregation in a
1212 species concept defines a clear partition for later "within-species" contrasts (evo-
1213 lution) and "between-species" interactions (competition & ecological interactions
1214 among populations or aggregates of species). A species concept must meaning-
1215 fully contrast the units of diversity - the biological species concept contrasts
1216 species based on reproductive isolation, the genetic species concept contrasts
1217 species based on genetic dissimilarity, and ecological species concepts contrast
1218 species based on distinct functional ecological traits. The objective function in
1219 phylofactorization is the theoretical placeholder for a researcher's "meaningful
1220 contrast". The units for aggregation and contrast must be done in light of some
1221 objective, such as a common fitness or pattern of relative abundance within units
1222 over time, space, across environmental gradients and/or between experimental
1223 treatments. A full theoretical consideration of phylofactorization as a species
1224 concept, as it relates to evolutionary and ecological theory, is saved for future
1225 research. For the time being, we note that phylofactorization partitions diver-
1226 sity and yields notions of a "species" which can be aggregated and contrasted
1227 with other "species".

1228 Phylofactorization is a flexible species concept, a hybrid of the phylogeny-
1229 based phylogenetic species concept [34] and the character-based ecological species
1230 concept [48]. After k iterations of phylofactorization, the phylogeny is par-
1231 titioned into $k + 1$ bins of species referred to as "binned phylogenetic units"
1232 (BPUs). BPUs are aggregations of the phylogeny which, up to a certain level
1233 of partitioning, are more similar to one-another with respect to the aggrega-

1234 tion, contrast and objective function, than they are to other groups. BPU are
1235 a coarse-grained way to cluster entities into “units” of organization with com-
1236 mon behavior with respect to the ecological pattern defined in the objective
1237 function. Phylofactorization defines functional groups based on phylogenetic
1238 partitions and a similar association with some ecological pattern of interest.
1239 Consequently, phylofactorization can be seen as an ecological species concept
1240 constrained to a phylogenetic scaffold. Whereas the phylogenetic species con-
1241 cept is character-based and pattern oriented, phylofactorization is pattern-based
1242 and phylogenetically-constrained. A textbook example of a phylofactorization-
1243 derived species are “land-dwelling tetrapods”, a group which can be obtained
1244 objectively through phylofactorization and which defines a scale for aggregating
1245 and summarizing the pattern of vertebrate species-abundances across land/water
1246 habitats.

1247 Phylofactorization permits optional fine-graining and coarse-graining of our
1248 patterns of diversity. Phylofactorization provides an algorithm for identifying
1249 relevant units, and those units may be at different taxonomic or phylogenetic
1250 depths but species within those units will have shared evolutionary history and
1251 similar associations with the ecological pattern of interest. For microorganisms,
1252 for which the biological species concept doesn’t apply, the genetic species con-
1253 cept appears too detached from ecology, and the ecological species concept is
1254 unavailable due to lack of life history detail, phylofactorization serves as a way
1255 to organize diversity for focused between-species interactions and within-species
1256 comparisons.

1257 **R package: phylofactor** An R package is in development and, prior to its
1258 stable release to CRAN, publicly available at <https://github.com/reptalex/phylofactor>.
1259 The R package contains detailed help functions and supports flexible definition
1260 of two-sample tests (the function `twoSampleFactor`), contrast-basis analyses with

1261 the function `PhyloFactor`, and generalized phylofactorization of exponential fam-
1262 ily random variables with the function `gpf`. Phylofactorization is highly par-
1263 allelizable, and the R package functions have built-in parallelization. The R
1264 package in development also works with phylogenies containing polytomies, al-
1265 lowing researchers to collapse clades with low bootstrap support to make more
1266 robust inferences. The output from phylofactorization is a “phylofactor” object
1267 containing the contrast basis, the BPU, and other details allowing one to input
1268 the object into various functions which summarize, plot, cross-validate, run null
1269 simulations, and parse out the information from phylofactorization. Researchers
1270 are invited to contact the corresponding author for assistance with the package
1271 and how to produce their own customized phylofactorizations - such feedback
1272 will be essential for a user-friendly stable release to CRAN.

1273 Until then, the supplemental information contains the data and scripts used
1274 for all analyses done in this manuscript in an effort to accelerate method devel-
1275 opment in this field.

1276 **“Everything makes sense in light of evolution”** Phylogenetic factoriza-
1277 tion is a new paradigm for analyzing a large class of biological data. Ecological
1278 big-data, as Thomas Dhobzansky noted about biology in general, makes sense
1279 “in light of evolution”. Phylofactorization extends a broad category of data anal-
1280 yses - two sample tests, generalized linear modelling, factor analysis and PCA,
1281 and analysis of spatial and temporal patterns - to incorporate a natural set of
1282 variables and operations defined by the phylogeny. Phylofactorization localizes
1283 inferences in big data to particular edges or chains of edges on the phylogeny
1284 and, in so doing, accelerates our understanding of the phylogenetic scales under-
1285 lying ecological patterns of interest. The problem of pattern and scale is central
1286 to biology, and phylofactorization uses the pattern to objectively uncover the
1287 relevant phylogenetic scales in ecological datasets.

1288 Acknowledgments

1289 This work is published in loving memory of Diana Nemergut. This research was
1290 developed with funding from the Defense Advanced Research Projects Agency
1291 (DARPA; D16AP00113).

1292 **Table of mathematical notation**

Term	Description
$A(\cdot)$	Aggregation operator
$C(\cdot, \cdot)$	Contrast operator
$\mathcal{F}(\theta)$	Distribution parameterized by θ
F_e	F-statistic for edge e
K_t	Number of edges considered in iteration t of phylofactorization
N	Size of a binomial random variable
Q	A group $Q = R \cup S$ aggregated at a current or previous iteration
R, S	Two groups contrasted containing r and s species, respectively
U, B, P	Meta-data subsets for phylofactorization
\mathcal{T}	Phylogenetic tree
B	$m \times p$ coefficient matrix
W	matrix of component scores corresponding to V
V	m matrix of contrast basis elements
X	$m \times n$ data matrix used for phylofactorization
Y	$K \times n$ matrix of component scores, one for each edge considered
Z	$n \times p$ matrix of meta-data used in regression-phylofactorization
a	Coefficient in aggregation vector
b, c	Coefficients in a contrast vector
e_k	Edge k
e^*	Winning edge
e_t^*	Winning edge at iteration t
$f(\cdot)$	Transformation in generalized f -mean
g	Factor containing two levels, $\{R, S\}$
i, j, k, l	Indexes. Often, i is the index for species and j for samples.
m	Number of species
n	Number of samples
p	Number of meta-data types for each sample

Terms	Description
q	Number of pure aggregates in a basis for \mathbb{R}^m
r, s	Numbers of species in groups R, S respectively
$s(\cdot)$	Smoothing spline notation for term in generalized additive model
t	Iteration of phylofactorization
$x_{i,j}$	The i, j th element of data matrix \mathbf{X}
$x_{R,j}$	Aggregate, $A(\mathbf{x}_j)$ of group R for sample j . If j is missing then sample is arbitrary.
$x_{S,j}$	See $x_{R,j}$ above.
x_i	A random variable (assumed to be a single species i for arbitrary sample)
$[x]_{i,j}$	i, j th entry of data matrix, \mathbf{X}
z_i	Column of meta-data matrix, \mathbf{Z}
$v_{Q,i}$	i th element of aggregation basis for set Q
$v_{C_{R S}}$	Contrast vector splitting groups R and S
v_{C_e}	Contrast vector for edge e (which splits sub-tree into two disjoint groups)
$\mathbf{x}_{R,j}$	r -vector containing only the species in group R for sample j
$\mathbf{x}_{S,j}$	See $\mathbf{x}_{R,j}$ above.
\mathbf{x}	m -vector of species' data for an arbitrary sample
$\bar{\mathbf{x}}$	Sample mean of vector \mathbf{x}
\mathbf{y}_e	n -vector of component scores for edge e
\mathbf{z}_k	Vector of meta-data of type k .
β_i	Coefficients for linear model
η	Natural parameter for exponential-family random variable
κ	Scale parameter for Gamma distribution
π	Number of failures parameter for Negative Binomial distribution.
ρ	Probability of success for Bernoulli, Binomial, Negative Binomial distributions
σ	Standard deviation for Gaussian random variable
θ	Arbitrary parameters for probability distribution

1293 References

- 1294 [1] J. AITCHISON, *The statistical analysis of compositional data*, (1986).
- 1295 [2] J. ALROY, *Cope's rule and the dynamics of body mass evolution in north*
1296 *american fossil mammals*, *Science*, 280 (1998), pp. 731–734.
- 1297 [3] ———, *The fossil record of north american mammals: evidence for a pale-*
1298 *ocene evolutionary radiation*, *Systematic Biology*, 48 (1999), pp. 107–118.
- 1299 [4] J. BAKER, A. MEADE, M. PAGEL, AND C. VENDITTI, *Adaptive evolution*
1300 *toward larger size in mammals*, *Proceedings of the National Academy of*
1301 *Sciences*, 112 (2015), pp. 5093–5098.
- 1302 [5] S. P. BLOMBERG, T. GARLAND JR, A. R. IVES, AND B. CRESPI, *Test-*
1303 *ing for phylogenetic signal in comparative data: behavioral traits are more*
1304 *labile*, *Evolution*, 57 (2003), pp. 717–745.
- 1305 [6] A. BULUÇ, H. MEYERHENKE, I. SAFRO, P. SANDERS, AND C. SCHULZ,
1306 *Recent advances in graph partitioning*, in *Algorithm Engineering*, Springer,
1307 2016, pp. 117–158.
- 1308 [7] J. C. CLEMENTE, L. K. URSELL, L. W. PARFREY, AND R. KNIGHT, *The*
1309 *impact of the gut microbiota on human health: an integrative view*, *Cell*,
1310 148 (2012), pp. 1258–1270.
- 1311 [8] T. Z. DESANTIS, P. HUGENHOLTZ, N. LARSEN, M. ROJAS, E. L.
1312 BRODIE, K. KELLER, T. HUBER, D. DALEVI, P. HU, AND G. L. ANDER-
1313 SEN, *Greengenes, a chimera-checked 16s rrna gene database and workbench*
1314 *compatible with arb*, *Applied and environmental microbiology*, 72 (2006),
1315 pp. 5069–5072.

- 1316 [9] J. J. EGOZCUE AND V. PAWLOWSKY-GLAHN, *Groups of parts and their*
1317 *balances in compositional data analysis*, *Mathematical Geology*, 37 (2005),
1318 pp. 795–828.
- 1319 [10] J. J. EGOZCUE, V. PAWLOWSKY-GLAHN, G. MATEU-FIGUERAS, AND
1320 C. BARCELO-VIDAL, *Isometric logratio transformations for compositional*
1321 *data analysis*, *Mathematical Geology*, 35 (2003), pp. 279–300.
- 1322 [11] C. E. FARRIOR, R. DYBZINSKI, S. A. LEVIN, AND S. W. PACALA, *Com-*
1323 *petition for water and light in closed-canopy forests: a tractable model of*
1324 *carbon allocation with implications for carbon sinks*, *The American Natu-*
1325 *ralist*, 181 (2013), pp. 314–330.
- 1326 [12] J. FELSENSTEIN, *Confidence limits on phylogenies: an approach using the*
1327 *bootstrap*, *Evolution*, (1985), pp. 783–791.
- 1328 [13] ———, *Phylogenies and the comparative method*, *The American Naturalist*,
1329 125 (1985), pp. 1–15.
- 1330 [14] L. Z. GARAMSZEGI, *Modern phylogenetic comparative methods and their*
1331 *application in evolutionary biology*, *Concepts and Practice*. London, UK:
1332 Springer, (2014).
- 1333 [15] N. E.-S. J. GOULD, *Punctuated equilibria: an alternative to phyletic grad-*
1334 *ualism*, (1972).
- 1335 [16] A. GRAFEN, *The phylogenetic regression*, *Philosophical Transactions of the*
1336 *Royal Society of London. Series B, Biological Sciences*, 326 (1989), pp. 119–
1337 157.
- 1338 [17] C. H. GRAHAM, D. STORCH, AND A. MACHAC, *Phylogenetic scale in*
1339 *ecology and evolution*, *bioRxiv*, (2017).

- 1340 [18] B. G. HALL AND M. BARLOW, *Evolution of the serine β -lactamases: past,*
1341 *present and future*, Drug Resistance Updates, 7 (2004), pp. 111–123.
- 1342 [19] T. F. HANSEN, *Stabilizing selection and the comparative analysis of adap-*
1343 *tation*, Evolution, 51 (1997), pp. 1341–1351.
- 1344 [20] W. K. HASTINGS, *Monte carlo sampling methods using markov chains and*
1345 *their applications*, Biometrika, 57 (1970), pp. 97–109.
- 1346 [21] J. L. HORN, *A rationale and test for the number of factors in factor anal-*
1347 *ysis*, Psychometrika, 30 (1965), pp. 179–185.
- 1348 [22] S. P. HUBBELL, *The Unified Neutral Theory of Biodiversity and Bio-*
1349 *geography (MPB-32)*, Princeton University Press, 2001.
- 1350 [23] M. JERRUM AND G. B. SORKIN, *The metropolis algorithm for graph bi-*
1351 *section*, Discrete Applied Mathematics, 82 (1998), pp. 155–175.
- 1352 [24] K. E. JONES, J. BIELBY, M. CARDILLO, S. A. FRITZ, J. O'DELL,
1353 C. D. L. ORME, K. SAFI, W. SECHREST, E. H. BOAKES, C. CAR-
1354 BONE, ET AL., *Pantheria: a species-level database of life history, ecology,*
1355 *and geography of extant and recently extinct mammals*, Ecology, 90 (2009),
1356 pp. 2648–2648.
- 1357 [25] Y. KATZ, K. TUNSTRØM, C. C. IOANNOU, C. HUEPE, AND I. D. COUZIN,
1358 *Inferring the structure and dynamics of interactions in schooling fish*, Pro-
1359 ceedings of the National Academy of Sciences, 108 (2011), pp. 18720–18725.
- 1360 [26] M. J. LANDIS, J. G. SCHRAIBER, AND M. LIANG, *Phylogenetic analysis*
1361 *using lévy processes: finding jumps in the evolution of continuous traits*,
1362 Systematic biology, 62 (2012), pp. 193–204.
- 1363 [27] S. A. LEVIN, *The problem of pattern and scale in ecology: the robert h.*
1364 *macarthur award lecture*, Ecology, 73 (1992), pp. 1943–1967.

- 1365 [28] R. E. LEY, P. J. TURNBAUGH, S. KLEIN, AND J. I. GORDON, *Microbial*
1366 *ecology: human gut microbes associated with obesity*, *Nature*, 444 (2006),
1367 pp. 1022–1023.
- 1368 [29] C. LOZUPONE AND R. KNIGHT, *Unifrac: a new phylogenetic method for*
1369 *comparing microbial communities*, *Applied and environmental microbiol-*
1370 *ogy*, 71 (2005), pp. 8228–8235.
- 1371 [30] V. A. MARČENKO AND L. A. PASTUR, *Distribution of eigenvalues for*
1372 *some sets of random matrices*, *Mathematics of the USSR-Sbornik*, 1 (1967),
1373 p. 457.
- 1374 [31] D. MARIAT, O. FIRMESSE, F. LEVENEZ, V. GUIMARÃES, H. SOKOL,
1375 J. DORÉ, G. CORTIER, AND J. FURET, *The firmicutes/bacteroidetes ra-*
1376 *tio of the human microbiota changes with age*, *BMC microbiology*, 9 (2009),
1377 p. 123.
- 1378 [32] N. METROPOLIS, A. W. ROSENBLUTH, M. N. ROSENBLUTH, A. H.
1379 TELLER, AND E. TELLER, *Equation of state calculations by fast computing*
1380 *machines*, *The journal of chemical physics*, 21 (1953), pp. 1087–1092.
- 1381 [33] F. MICHONNEAU, J. W. BROWN, AND D. J. WINTER, *rotl: an r package to*
1382 *interact with the open tree of life data*, *Methods in Ecology and Evolution*,
1383 7 (2016), pp. 1476–1481.
- 1384 [34] K. C. NIXON AND Q. D. WHEELER, *An amplification of the phylogenetic*
1385 *species concept*, *Cladistics*, 6 (1990), pp. 211–223.
- 1386 [35] M. PAGEL, *Inferring the historical patterns of biological evolution*, *Nature*,
1387 401 (1999), pp. 877–884.
- 1388 [36] E. PARADIS, J. CLAUDE, AND K. STRIMMER, *Ape: analyses of phyloge-*
1389 *netics and evolution in r language*, *Bioinformatics*, 20 (2004), pp. 289–290.

- 1390 [37] R. K. PLOWRIGHT, C. R. PARRISH, H. MCCALLUM, P. J. HUDSON,
1391 A. I. KO, A. L. GRAHAM, AND J. O. LLOYD-SMITH, *Pathways to zoonotic*
1392 *spillover*, Nature Reviews Microbiology, (2017).
- 1393 [38] E. PURDOM, *Analysis of a data matrix and a graph: Metagenomic data and*
1394 *the phylogenetic tree*, The Annals of Applied Statistics, (2011), pp. 2326–
1395 2358.
- 1396 [39] K. S. RAMIREZ, J. W. LEFF, A. BARBERÁN, S. T. BATES, J. BET-
1397 LEY, T. W. CROWTHER, E. F. KELLY, E. E. OLDFIELD, E. A. SHAW,
1398 C. STEENBOCK, ET AL., *Biogeographic patterns in below-ground diversity*
1399 *in new york city’s central park are similar to those observed globally*, in
1400 Proc. R. Soc. B, vol. 281, The Royal Society, 2014, p. 20141988.
- 1401 [40] L. J. REVELL, *phytools: an r package for phylogenetic comparative biology*
1402 *(and other things)*, Methods in Ecology and Evolution, 3 (2012), pp. 217–
1403 223.
- 1404 [41] K.-I. SATO, *Lévy processes and infinitely divisible distributions*, Cambridge
1405 university press, 1999.
- 1406 [42] J. U. SCHER, A. SZESNAK, R. S. LONGMAN, N. SEGATA, C. UBEDA,
1407 C. BIELSKI, T. ROSTRON, V. CERUNDOLO, E. G. PAMER, S. B. ABRAM-
1408 SON, ET AL., *Expansion of intestinal prevotella copri correlates with en-*
1409 *hanced susceptibility to arthritis*, Elife, 2 (2013), p. e01202.
- 1410 [43] K. P. SCHLIEP, *phangorn: phylogenetic analysis in r*, Bioinformatics, 27
1411 (2011), pp. 592–593.
- 1412 [44] J. D. SILVERMAN, A. D. WASHBURNE, S. MUKHERJEE, AND L. A.
1413 DAVID, *A phylogenetic transform enhances analysis of compositional mi-*
1414 *crobiota data*, Elife, 6 (2017), p. e21887.

- 1415 [45] F. A. SMITH, A. G. BOYER, J. H. BROWN, D. P. COSTA, T. DAYAN,
1416 S. M. ERNEST, A. R. EVANS, M. FORTELIUS, J. L. GITTLEMAN, M. J.
1417 HAMILTON, ET AL., *The evolution of maximum body size of terrestrial*
1418 *mammals*, *science*, 330 (2010), pp. 1216–1219.
- 1419 [46] F. A. SMITH AND S. K. LYONS, *How big should a mammal be? a macroe-*
1420 *cological look at mammalian body size over space and time*, *Philosophical*
1421 *Transactions of the Royal Society of London B: Biological Sciences*, 366
1422 (2011), pp. 2364–2378.
- 1423 [47] P. J. TURNBAUGH, R. E. LEY, M. A. MAHOWALD, V. MAGRINI, E. R.
1424 MARDIS, AND J. I. GORDON, *An obesity-associated gut microbiome with*
1425 *increased capacity for energy harvest*, *nature*, 444 (2006), pp. 1027–131.
- 1426 [48] L. VAN VALEN, *Ecological species, multispecies, and oaks*, *Taxon*, (1976),
1427 pp. 233–239.
- 1428 [49] Y. VÁZQUEZ-BAEZA, A. GONZALEZ, Z. Z. XU, A. WASHBURNE, H. H.
1429 HERFARTH, R. B. SARTOR, AND R. KNIGHT, *Guiding longitudinal sam-*
1430 *pling in ibd cohorts*, *Gut*, (2017), pp. gutjnl–2017.
- 1431 [50] A. D. WASHBURNE, J. W. BURBY, AND D. LACKER, *Novel covariance-*
1432 *based neutrality test of time-series data reveals asymmetries in ecological*
1433 *and economic systems*, *PLoS computational biology*, 12 (2016), p. e1005124.
- 1434 [51] A. D. WASHBURNE, J. D. SILVERMAN, J. W. LEFF, D. J. BENNETT,
1435 J. L. DARCY, S. MUKHERJEE, N. FIERER, AND L. A. DAVID, *Phyloge-*
1436 *netic factorization of compositional data yields lineage-level associations in*
1437 *microbiome datasets*, *PeerJ*, 5 (2017), p. e2969.
- 1438 [52] T. W. YEE AND T. J. HASTIE, *Reduced-rank vector generalized linear*
1439 *models*, *Statistical modelling*, 3 (2003), pp. 15–41.

- 1440 [53] G. YU, D. K. SMITH, H. ZHU, Y. GUAN, AND T. T.-Y. LAM, *ggtree: an*
1441 *r package for visualization and annotation of phylogenetic trees with their*
1442 *covariates and other associated data*, *Methods in Ecology and Evolution*, 8
1443 (2017), pp. 28–36.
- 1444 [54] X. ZHOU, S. XU, J. XU, B. CHEN, K. ZHOU, AND G. YANG, *Phy-*
1445 *logenomic analysis resolves the interordinal relationships and rapid diver-*
1446 *sification of the laurasiatherian mammals*, *Systematic biology*, 61 (2011),
1447 pp. 150–164.