

Predicting environmental gradients with fern species composition in
Brazilian Amazonia

**Gabriela Zuquim, Hanna Tuomisto, Mirkka Jones, Jefferson Prado, Fernando F.O.G.
Figueiredo, Gabriel M. Moulatlet, Flavia R.C. Costa, Carlos A. Quesada & Thaise Emilio**

Zuquim, G (corresponding author, gabizuquim@gmail.com), Tuomisto, H
(hanna.tuomisto@utu.fi), Jones, M (mirkka.jones@gmail.com) & Moulatlet, G.M
(mandaprogabriel@gmail.com): Department of Biology, University of Turku, FI-20014 Turku,
Finland.

Jones, M: Ecoinformatics and Biodiversity Group, Department of Bioscience, Aarhus University,
DK-8000, Aarhus C., Denmark.

Prado, J. (jp Prado.01@uol.com.br): Instituto de Botânica, Herbário SP, C.P. 68041, CEP 04045-
972, São Paulo, SP, Brazil.

Figueiredo, F.O.G (nandoeco06@gmail.com) & Emilio, T (thaise.emilio@gmail.com): Programa
de Pós-Graduação em Ecologia, Instituto Nacional de Pesquisas da Amazônia (INPA/MCTI).
Manaus, Brazil.

Moulatlet, G.M., Costa, F.R.C (flaviacosta001@gmail.com) & Quesada, C.A
(quesada.beto@gmail.com): Instituto Nacional de Pesquisas da Amazônia, CP 478, 69011-970,
Manaus, AM, Brazil.

25 **Abstract**

26 **Aim:** A major problem for conservation in Amazonia is that species distribution maps are
27 inaccurate. Consequently, conservation planning needs to be based on other information sources
28 such as vegetation and soil maps which are inaccurate as well. We propose and test the use of
29 biotic data on a common and relatively easily inventoried group of plants to infer environmental
30 conditions that can be used to improve maps of floristic patterns for plants in general.

31 **Location:** Brazilian Amazonia.

32 **Methods:** We sampled 326 plots of 250 m x 2 m separated by distances of 1 to 1800 km.
33 Terrestrial fern individuals were identified and counted. Edaphic data were obtained from soil
34 samples and analyzed for cation concentration and texture. Climatic data were obtained from
35 Worldclim. We performed multivariate regression tree to evaluate the hierarchical importance of
36 soils and climate for fern communities and identified significant indicator species for the
37 resultant classification. We then tested how well the edaphic properties of the plots could be
38 predicted on the basis of their floristic composition using two calibration methods, weighted
39 averaging and k-nearest neighbour estimation.

40 **Results:** Soil cation concentration emerged as the most important variable in the regression tree,
41 whereas soil textural and climatic variation played secondary roles. Almost all the plot classes
42 had several fern species with high indicator values for that class. Soil cation concentration was
43 also the variable most accurately predicted on the basis of fern community composition ($R^2 =$
44 0.65-0.75 for log-transformed data). Predictive accuracy varied little among the calibration
45 methods, and was not improved by the use of abundance data instead of presence-absence data.

46 **Main conclusions:** Fern species composition can be used as an indicator of soil cation
47 concentration, which can be expected to be relevant also for other components of rain forests.

48 Presence-absence data are adequate for this purpose, which makes the collecting of additional
49 data potentially very rapid. Comparison with earlier studies suggests that edaphic preferences of
50 fern species have good transferability across geographical regions within lowland Amazonia.
51 Therefore, species and environmental datasets already available in the Amazon region represent
52 a good starting point for generating better environmental and floristic maps for conservation
53 planning.

54 **Keywords:** pteridophytes; tropical forest; edaphic characteristics; floristic composition;
55 vegetation maps; *k*-NN; weighted averaging; calibration methods; indicator species.

56 **Nomenclature:** The International Plant Name Index (IPNI) (www.ipni.org; accessed 22 July
57 2013)

58 **Abbreviations:** db-MRT = distance-based Multivariate Tree Regression; *k*-NN = *k* nearest-
59 neighbours; RMSE = Root Mean Squared Error; WA = Weighted Averaging

60 **Running head:** Predicting soil fertility using Amazonian ferns

61

62

63 Introduction

64 Understanding the spatial heterogeneity of environmental conditions and species
65 distributions in Amazonia is a major challenge for conservation planning. A generally accepted
66 principle is that the network of conservation units should contain adequate representation of
67 different habitats, so as to collectively provide living space for species adapted to different
68 habitats. Currently, sufficiently detailed maps that would allow assessing whether this aim has
69 been fulfilled do not exist for Amazonia. The available soil and species distributions maps are
70 inaccurate and give an incomplete representation of the known Amazonian heterogeneity.

71 Several soil maps are available for Amazonia ([RADAMBRASIL 1978](#); SOTERLAC -
72 [Dijkshoorn et al. 2005](#); [Quesada et al. 2011](#)), but all of them are coarse-grained because there is a
73 general paucity of ground data. While information on broad-scale variation in soil properties can
74 be extracted from such maps, this is not sufficient to take into account the documented effects of
75 soil variation on biotic heterogeneity at local to landscape scales ([Phillips et al. 2003](#); [Tuomisto](#)
76 [et al. 2003a, b, c](#); [Costa et al. 2005](#); [Kinupp & Magnusson 2005](#); [Jones et al. 2006](#); [Ruokolainen](#)
77 [et al. 2007](#); [Zuquim et al. 2009a](#); [Higgins et al. 2011](#)). Consequently, there is a general lack of
78 knowledge of the distribution of Amazonian habitat types ([Emilio et al. 2010](#)) and species
79 ([Schulman et al. 2007a](#)), which forces conservation planning in Amazonia to be based on the use
80 of more or less unreliable surrogates ([Schulman et al. 2007b](#)).

81 When information on environmental gradients is needed but measurements of
82 environmental variables cannot be made, biotic communities have been used as predictors of the
83 environmental conditions. For example, paleo-environmental reconstructions ([Birks et al. 2010](#))
84 use modern species-climatic relationships to infer past climatic conditions according to the
85 analogue fossil record ([ter Braak & van Dam 1989](#); [Birks et al. 1990](#)). The same approach was
86 used by Sirén et al. (2013) to generate predictive maps of soil fertility based on fern and

87 lycophyte species composition in a lowland rainforest area in Ecuadorian Amazonia. The authors
88 used floristic and soil data from other parts of western Amazonia ([Tuomisto et al. 2003a](#) and
89 unpublished data) to determine fern and lycophyte species' optima on a soil cation concentration
90 gradient. Then they used those optima to estimate soil cation concentrations in their study area,
91 where fern and lycophyte species lists were available but direct measurements of soil properties
92 were not. Suominen et al. (2013) recently evaluated the application of similar estimation
93 techniques for predicting chemical soil properties in western Amazonia using species occurrence
94 data of the plant family Melastomataceae.

95 Specific taxa can also be used as indicators of particular environments or habitat types
96 ([Ruokolainen et al. 1997](#); [Ruokolainen et al. 2007](#); [Margules et al. 2002](#); [Tuomisto et al. 2003a](#);
97 [Salovaara et al. 2004](#)). The use of indicator species ([Noss 1990](#)) is an important method in
98 conservation biology because it is flexible ([Dufrière & Legendre 1997](#)) and conceptually
99 straightforward ([McGeoch 1998](#)). Well-chosen indicator taxa can contribute significantly to a
100 conservation strategy by facilitating the recognition and mapping of habitats ([Noss 1990](#);
101 [Howard et al. 1998](#)).

102 Ferns have been proposed as a suitable indicator group in Amazonia because they are
103 easy to observe and identify. Several studies have documented edaphic affinities of selected fern
104 species in the western Amazon region in relation to either a simple classification of soil types
105 ([Tuomisto & Poulsen 1996](#); [Salovaara et al. 2004](#); [Cárdenas et al. 2007](#)), or quantitative soil
106 gradients ([Tuomisto et al. 1998, 2002](#); [Tuomisto 2006](#)). Some of these studies have only reported
107 results for a few species within selected genera, and none has explicitly assessed the accuracy of
108 soil property estimates when these are based on indicator values of the species.

109 In this study, we investigate the use of ferns as environmental indicators in central and
110 northern Amazonian lowlands. First, we clarify the main environmental drivers of fern

111 community composition and define the environmental optima and tolerances for each species
112 along each of these gradients. Then we use species optima to predict environmental variable
113 values and test the accuracy of these predictions. Finally, we assess whether species abundance
114 data are needed to obtain useful predictions, or whether the more easily obtainable presence-
115 absence data are adequate.

116 **Methods**

117 **Study area and sampling design**

118 A total of 326 plots were sampled (Fig. 1). Plots were located in Brazilian
119 Amazonian lowlands in the states of Acre (7 plots), Amazonas (129 plots), Pará (101 plots),
120 Rondônia (30 plots) and Roraima (59 plots). All study sites are part of the Brazilian Biodiversity
121 Research Program (PPBio, <http://ppbio.inpa.gov.br/>). Minimum distance between plots was 1 km
122 and maximum ca. 1800 km. Plots were established in private lands or in conservation units along
123 the highways BR-163, BR-230 (Transamazônica) and BR-319 and in the protected areas of
124 ReBio Uatumã, ESEC Maracá, PN Viruá, BDFFP and PE Chandless. In every location, 5 to 30
125 plots were established according to the RAPELD methodology ([Magnusson et al. 2005](#)). The
126 plots were 250 x 2 m in size and placed so that the longer axis followed the topographic contour
127 in order to minimize internal heterogeneity in soil properties and drainage, which often correlate
128 with topographic position ([Chauvel et al. 1987](#); [Mertens 2004](#)). Vegetation structure in the plots
129 varied from tall and dense rainforests to white sand forests with a more simple canopy structure
130 (campinaranas) and in extreme cases edaphic savannas (IBGE 2004). According to the Soil and
131 Terrain Database for Latin America and the Caribbean (SOTERLAC - [Dijkshoorn et al. 2005](#)),
132 six main soil classes dominated the areas where the plots were situated: Ferralsols (157 plots),
133 Podzols (29 plots), Plinthosol (91 plots), Acrisols (37 plots), Leptosols (5 plots), and Cambisols
134 (7 plots). Because local-scale soil variation does not appear in broad-scale maps, it is possible
135 that some of the plots were in fact situated in a different soil type than the one dominating the

136 region. Average annual rainfall in the plots ranged from 1,633 to 2,655 mm and annual mean
137 temperature from 25 to 27°C. General characteristics of the study sites can be found in Table 1,
138 and a more detailed description of each region in appendix S1.

139 **Data collection**

140 FLORISTIC DATA - In each plot, all terrestrial fern individuals with at least one leaf
141 longer than 10 cm were counted and identified to species. Inventories were done between 2004
142 and 2011. Voucher specimens were collected to verify species identities. Full sets of the
143 vouchers are deposited in Herbaria at the Instituto de Botânica, São Paulo (SP) and privately
144 with the first author. Duplicates of fertile specimens are also deposited in the nearest regional
145 herbarium either at Instituto Nacional de Pesquisas da Amazônia (INPA), Herbário Rondoniense
146 (RON) or Universidade Federal do Acre (UFACPZ).

147 ENVIRONMENTAL DATA - Surface soil samples (topmost layer of the mineral soil
148 sampled down to 5-10 cm depth) were taken every 50 m along the long axis of each plot. The six
149 soil samples from the same plot were either bulked into a single composite sample before
150 laboratory analyses or analyzed separately. In the latter case, the obtained values were averaged
151 to obtain a single value for each edaphic variable for each plot. Before laboratory analyses, the
152 soil samples were air-dried, cleaned of roots and other detritus and sieved through a 2 mm mesh.
153 Analyses included soil texture (percentage of clay, silt and sand, by the pipette method) and
154 exchangeable bases (Ca, Mg by KCl 1 M and K by Mehlich 1 standard methods for
155 exchangeable cations). All soil samples were analyzed in the Thematic Laboratory of Soils and
156 Plants at INPA. Floristic data, soil data and geographical coordinates of the plots are publicly
157 available at <http://ppbio.inpa.gov.br/knb/style/skins/ppbio/>. The plots were georeferenced in the
158 field using a hand-held GPS (Garmin 12XL or Garmin 60X).

159 Climatic data were derived from monthly temperature and rainfall values available
160 in Bioclim ([Hijmans et al. 2005](#)). The variables used were annual temperature range, annual

161 precipitation, precipitation seasonality and precipitation of the wettest quarter (Bioclim variables
162 7, 12, 15 and 16, respectively). The data were downloaded from WorldClim database
163 (<http://www.worldclim.org/bioclim>) in 2.5 arc-minutes resolution (about 4.7 km). The remaining
164 15 climatic variables available in Bioclim were not included either because they were strongly
165 correlated with an already selected variable and hence provided little additional information, or
166 because they varied so little within our study region that it seemed unlikely that it would result in
167 floristic response. Amazonia has few climatic stations, so the real resolution of the data is
168 probably much poorer than the nominal pixel size, and there are known problems of data
169 uncertainty (Hijmans et al. 2005). Nevertheless, this is currently the best available source of
170 temperature and rainfall data for the area. The climatic values for each plot were extracted using
171 the free software DIVA-GIS ([Hijmans et al. 2012](#)).

172 **Data analysis**

173 Fern species that occurred in less than five plots were excluded from all analyses,
174 as species optima based on so few data points were considered too unreliable. Twenty-one of the
175 plots had no fern species with the minimum frequency determined. These plots were excluded
176 from the analyses, which therefore were run on 305 plots. The sum of exchangeable bases
177 (concentration of Ca+Mg+K, all in $\text{cmol}_+ \text{kg}^{-1}$) was logarithmically transformed (base 10) before
178 numerical analyses. This was done because it is reasonable to assume that plants react to relative
179 rather than absolute differences in the availability of soil nutrients, i.e. small differences in soil
180 cation concentration are ecologically important if the overall cation concentration is low but
181 inconsequential if the overall cation concentration is high.

182 REGRESSION TREES AND INDICATOR SPECIES - To evaluate the hierarchical
183 importance of edaphic and climatic conditions in structuring fern communities, we carried out a
184 distance-based multivariate regression tree analysis (db-MRT; [De'ath 2002](#)). MRT is based on
185 repeatedly splitting the plots into two groups that are separated by a single value along one of the

186 environmental gradients. At each split, the gradient and the threshold value are selected so as to
187 minimize the between-plot compositional dissimilarities within each group. As a measure of
188 compositional dissimilarity, we used the extended Bray-Curtis dissimilarity index ([De'ath 1999](#))
189 based on species proportional abundances (number of individuals as a proportion of all
190 individuals in the plot). The extended rather than classical Bray-Curtis index was used because
191 our data covered long environmental gradients, so a large proportion of the plots shared no
192 species. This leads to poor model fit if not corrected for ([De'ath 1999](#); [Tuomisto et al. 2012](#);
193 [Zuquim et al. 2012](#)). To find the best db-MRT classification, we used cross-validation and
194 selected the db-MRT with the smallest error, given by the sum of squares ([De'ath 2002](#)). We then
195 assessed whether any species were significantly associated with the groups of plots obtained
196 from the db-MRT by calculating the indicator value of each species for each group. A high
197 indicator value is obtained for species that combine high specificity (most individuals of the
198 species are within the group) and high fidelity (most sites of the group contain the species). The
199 IndVal index was used for this purpose ([Dufrière & Legendre 1997](#); Legendre & Legendre
200 1998).

201 ENVIRONMENTAL PREDICTIONS BASED ON *k*-NN AND WA - Next we asked how
202 accurately it is possible to estimate the values of environmental variables for a plot on the basis
203 of its floristic composition. Each variable was estimated for each plot using the species-
204 environment relationships as deduced from the remaining plots. We applied two methods that are
205 commonly used in paleoecology: the *k*-Nearest Neighbours (*k*-NN) and Weighted Averaging
206 calibration (WA) with inverse deshrinking.

207 *K*-NN is a non-parametric method that estimates the value of an environmental
208 variable in a focal plot on the basis of the average value of the variable in the *k* nearest
209 neighbouring plots. We used similarity in species composition as the measure of nearness, and
210 calculated it with either the Bray-Curtis index (for proportional abundance data) or the Sørensen

211 index (for presence-absence data). Each of the 305 plots was used as the focal plot in turn. The
212 results will depend on the value of k : when $k = 1$, the predicted value of the variable depends on
213 its value in a single plot, which may lead to noisy results, but when k increases, the predicted
214 value will tend towards its overall mean in the dataset. Different values of k may work best for
215 different kinds of data, so we run the analyses with $k=1$ to $k=20$ in order to find the value of k
216 that gives the most accurate predictions for this dataset.

217 WA estimates the value of an environmental variable in a focal plot as the weighted
218 average of the indicator values (optima) of the species occurring in the plot. We calculated the
219 optimum of a species along an environmental gradient as the weighted average of the
220 environmental variable values in those plots where the species had been observed, with species
221 abundance in a plot being used as the weight (eq. 4 in [ter Braak & van Dam 1989](#)). We ran these
222 analyses both using the number of individuals as the abundance measure, and using presence-
223 absence data (i.e. abundance was set to unity if the species was present and to zero if it was
224 absent). The optimum value carries no information on how broad the species' distribution is, so
225 in a second set of analyses we weighted each species' optimum value by the inverse of its
226 tolerance. Tolerance is a measure of the variability in species occurrences around the optimum,
227 and is obtained as the root mean squared error (RMSE) calculated between the species optimum
228 and the observed environmental variable value for each individual (eq. 7 in [ter Braak & van Dam](#)
229 [1989](#)). Because the WA computation involves the taking of averages twice, the range of the
230 estimated values tends to shrink, i.e. to become smaller than the range of the original
231 observations. We used inverse linear deshrinking to restore the original range of the variable ([ter](#)
232 [Braak & Juggins 1993](#)). WA is based on the idea of unimodal species response curves along the
233 environmental gradients, which we considered appropriate because our dataset is highly
234 heterogeneous ([Zuquim et al. 2012](#)).

235 Prediction accuracy was quantified with cross-validation for each environmental
236 variable separately using root mean squared error (RMSE) and the coefficient of determination
237 (R^2) between the measured and predicted values. Cross-validation was done using the leave-one-
238 out method for WA and by bootstrapping for k -NN. In our sampling design, the plots were
239 placed in 37 locations spread across eight regions (Fig. 1). Each location had 5 - 30 plots with
240 distances from 1 to 5 km between each other and in a regular arrangement within a few square
241 kilometers, so spatial autocorrelation might cause the predictive power of the calibration
242 methods to appear unrealistically high. For this reason, more stringent cross-validations were
243 also done by leaving out all plots that were in the same location as the focal plot when
244 calculating the predicted values.

245 Both k -NN and WA analyses were carried out separately using abundance and
246 presence-absence data. This was done because collecting abundance data is much more time-
247 consuming than collecting presence-absence data, so it is of interest to test if this is justified by
248 more accurate predictions.

249 All statistical analyses were carried out using the RStudio (v. 0.97.173; RStudio,
250 Inc., Boston, USA) interface to R (R Foundation for Statistical Computing, Vienna, AT).
251 Multivariate Regression Trees were made using the R package *mvpart* (v. 1.6-0) and Indicator
252 Species analysis with *indicspecies* (v. 1.6.5; de Caceres & Legendre 2009). K -NN, WA and
253 associated calculations of species optima and tolerances were done using the R package *Rioja*
254 (07-3).

255 RESULTS

256 GENERAL – After excluding species occurring in less than 5 plots, the 326 plots contained a total
257 of 29 202 individuals of ferns representing 54 species. Twenty-one plots contained no ferns at
258 all, or were left empty after the exclusion of the rare species. Twenty of the excluded plots were
259 in Roraima in the northern part of the study area, and one was in Pará. The most species-rich
260 genera were *Adiantum* (17 species), *Trichomanes* (7 species), *Lindsaea* (5 species), and
261 *Triplophyllum* (5 species). The most abundant species were *Trichomanes pinnatum* Hedw. (8512
262 individuals), *Adiantum argutum* Splitg. (8560 individuals), and *A. pulverulentum* L. (1593
263 individuals). The most frequent species were *T. pinnatum* (205 plots), *Lindsaea lancea* (L.)
264 Bedd. (132 plots) and *A. cajennense* Willd. (115 plots).

265 FERN COMMUNITY STRUCTURE AND INDICATOR SPECIES - The first division in the
266 multivariate regression tree (Fig. 2) was determined by the community response to the sum of
267 bases in the soil. One branch contained 79 plots with soil cation concentrations exceeding 0.68
268 $\text{cmol}_+ \text{kg}^{-1}$, and the other contained 226 plots with lower-cation soils. The second division was
269 defined by soil clay content in the richer-soils branch and by annual rainfall in the poorer-soils
270 branch. Textural components of soils determined two more hierarchical divisions within the plot
271 groups characterized by low-cation soils and low annual rainfall (Fig. 2). The other climatic
272 variables did not define any divisions in the regression tree. In preliminary analyses, we also
273 included latitude and longitude, because the climatic variables show clear spatial gradients across
274 Amazonia. However, neither latitude nor longitude substituted any of the climatic variables in
275 the regression tree, and since they are not direct environmental variables, they were left out of
276 the final analyses.

277 Most of the statistically significant indicator species were associated with the
278 branch containing the high-cation sites (Fig. 2). Nine out of seventeen species of *Adiantum* were

279 significant indicators of this branch and only two *Adiantum* species were significantly associated
280 with the poorer-soils branch, although the genus as a whole was represented over the entire
281 gradient. Both *Pteris* species were also associated with the richer-soils branch. Almost all of the
282 18 richer-soils indicator species were also significantly associated with the rich soils-high clay
283 content branch in the second level division.

284 Five out of seven *Trichomanes* species were indicators of some secondary or
285 tertiary division within the poorer-soils branch, and the very frequent *Trichomanes pinnatum*
286 indicated poor soils generally. Three out of five *Lindsaea* species were indicators of the poorer-
287 soils branch and none was significantly associated with the richer soils. The majority of poor soil
288 indicator species were associated with sites with relatively high total annual rainfall (≥ 2163 mm).
289 Only a few species were indicators of habitats with both poor soils and low rainfall.

290 There was a gradual turnover of species optimum values along the soil cation
291 concentration gradient, although most species optima were concentrated towards the low-cation
292 end (Fig. 3). In agreement with the results of the indicator value analysis, all species of the
293 genera *Lindsaea* and *Trichomanes* had low cation optima, whereas those of *Thelypteris* and
294 *Pteris* had high optima. *Adiantum phyllitidis* and *Cyclopeltis semicordata* were the two species
295 with the highest optima. Most *Adiantum* species optima were positioned in the intermediate part
296 of the gradient, but the genus had representatives along the whole gradient.

297 PREDICTING ENVIRONMENTAL VARIABLES FROM FERN INVENTORIES - The edaphic
298 variable that could be best predicted by fern species composition was the sum of bases. All
299 methods of calibration produced R^2 values that were between 0.64 and 0.75 when the focal plot
300 was excluded in cross-validation. When all plots from the same locality as the focal plot were
301 excluded in leave-group-out cross-validation, R^2 values decreased to between 0.46 and 0.64
302 (Table 2). There was variation among the regions in the slope of the regression line between

303 predicted and observed soil cation concentration, with the predictions for the Acre region
304 becoming especially inaccurate when leave-group-out cross-validation was used (Fig. 4). The R^2
305 values of the predictions for soil clay, sand and silt contents were never higher than 0.48 (Table
306 2). This is in accordance with the regression tree results, which suggested that ferns respond
307 more strongly to soil cation concentrations than to soil textural properties.

308 The best results (smallest RMSEs) for predictions using k -NN were achieved with
309 between four and seven neighbouring plots ($k=4$ to $k=7$). The differences in prediction accuracy
310 between k values in this range were generally small, so for simplicity we report the results for
311 $k=4$ in all cases. There were slight variations in prediction accuracy among methods, but none of
312 them was consistently better than the others for all the edaphic variables. Weighted Averaging
313 achieved lower RMSEs and higher R^2 values than k -NN when abundance data were used (Table
314 2), but with presence-absence data, k -NN gave similar or higher R^2 values.

315 Weighting species by the inverse of their tolerance improved the predictions in
316 some cases but not universally. When leave-group-out cross-validation was used, the differences
317 in accuracy between weighted and non-weighted estimations (R^2 and RMSE) were small. In
318 general, the availability of abundance data did not improve model performance. In fact, k -NN
319 always performed better with presence-absence data than with abundance data, and even WA did
320 so in most cases (Table 2).

321 **Discussion**

322 Earlier studies that have been carried out mostly in western Amazonia have
323 proposed that ferns and lycophytes are good indicators of environmental conditions, especially
324 soil cation concentration and particle size distribution ([Ruokolainen et al. 1997](#); [Ruokolainen et](#)

325 [al. 2007](#); [Tuomisto et al. 2003a, c](#); [Higgins et al. 2011](#)). Here we tested this proposal in central
326 Amazonia by making explicit predictions of soil properties and climatic variables on the basis of
327 information about fern species composition.

328 Our results supported the conclusions of earlier studies. The sum of bases emerged
329 as the most important variable in the regression tree, and was also the variable for which the
330 most accurate predictions could be obtained on the basis of fern community composition. Soil
331 textural and climatic variation played secondary roles in the regression tree, and soil texture was
332 predicted less accurately than soil base cation concentration. Soil texture is not a physiologically
333 important edaphic factor, but it correlates with other relevant environmental characteristics, such
334 as nutrient retention and water holding capacity. Climate is also relevant in structuring fern
335 communities at broad scales ([Zuquim et al. 2012](#); [Jones et al. 2013](#)), but in the present study its
336 role was minor. This is in agreement with the findings of Tuomisto and Poulsen (1996), who
337 found that even in a dataset where annual rainfall varied more than in ours, the main floristic
338 gradient still seemed to correspond to soil properties more than to rainfall.

339 PREDICTING EDAPHIC CONDITIONS FROM FERN INVENTORIES - We
340 found that sum of bases in the soil can be well predicted based on fern species composition. Our
341 analyses were carried out with log-transformed data, which means that prediction errors related
342 with large values of the variable of interest are downweighted. In other words, whether a
343 prediction is considered accurate or not depends more on how large the error is in relation to the
344 actual value of the variable of interest, rather than on the absolute error value. This is an
345 appropriate model in the present context, given that the final aim is to use the predicted soil
346 values to infer habitat characteristics and occurrence patterns for such plant groups that have not
347 been directly observed in the field.

348 Another result that has practical implications is that prediction accuracy for a
349 particular environmental variable was rather consistent among calibration methods. This
350 parallels the observations of Suominen et al. (2013), who tested the *k*-NN and WA methods in
351 western Amazonian transects using the family Melastomataceae as a model group. In a
352 theoretical sense, both methods have their strengths and weaknesses ([Birks et al. 2010](#)), but in
353 practical applications both seem to perform equally well. As could be expected, prediction
354 accuracy appeared generally higher when only the focal plot was left out of the training set than
355 when all the plots from the same site were left out (R^2 between 0.64–0.75 vs. 0.46–0.64). Fig. 4
356 shows that the decrease in prediction accuracy was most notable for the plots situated in Acre
357 state, for which the predictions fell dramatically below the observed values in the leave-group-
358 out cross-validation. This reflects the fact that the plots in Acre had the highest observed cation
359 concentrations in the entire data set, so when all of them were excluded from the training set, no
360 accurate analogue remained for the Acre plots. As with other modelling methods, attempts to
361 extrapolate predictions of WA calibration and *k*-NN estimation beyond the observed range of the
362 input variables can lead to seriously inaccurate results.

363 A third interesting result is that the prediction accuracies for the edaphic variables
364 were very similar whether species presence-absence data or abundance data were used. Even
365 though we expected abundance data to provide better estimates of species optima, and that this
366 would lead to more accurate predictions, this was not the case. One possible reason is that the
367 species abundances are so symmetrically distributed along the relevant environmental gradients
368 that the optimum is in practice at the midpoint of the species range, and can hence be identified
369 equally well with presence-absence and abundance data. Another possibility is that species
370 abundances depend on many different factors that are not necessarily linked to the factor being
371 evaluated. For example, fertility may limit the range of species, which is captured by presence
372 absence data, but may not be the main driver of local abundances, which may be controlled by
373 biotic interactions or more local factors such as light. These unmeasured factors may cause a

374 species to be relatively abundant far away from its optimum for a given variable, or not so
375 abundant close to its optimum, which then biases the estimate for that variable.

376 Earlier studies have obtained mixed results on whether using abundance data
377 increases or decreases the correlations between species turnover and edaphic differences
378 ([Tuomisto et al. 2003a](#); [Ruokolainen et al. 2007](#)). Our results support the suggestion that at least
379 when the observed soil gradients are relatively long, presence-absence data are adequate for
380 many purposes ([Tuomisto et al. 2002, 2003a](#); [Higgins & Ruokolainen 2004](#); [Higgins et al. 2011](#)).
381 This is good news, because collecting only presence-absence data speeds up the fieldwork
382 considerably. Moreover, these results suggests that it is feasible to tap edaphic information from
383 non-quantitative species lists and floras (e.g., [Tuomisto & Poulsen 1996](#); [Edwards 1998](#); [Costa et](#)
384 [al. 1999](#); [Freitas & Prado 2005](#); [Costa et al. 2006](#); [Costa & Pietrobon 2007](#); [Maciel et al. 2007](#);
385 [Prado & Moran 2009](#); [Zuquim et al. 2009b](#)), and perhaps even from herbarium records through
386 online databases such as GBIF. For example, linking species lists with the species' environmental
387 optima and tolerances enables inferences about site environmental conditions. This opens up new
388 and unexplored possibilities for assessing representativeness of conservation area networks
389 based on the use of readily available biotic data as indicators of habitat types.

390 SPECIES OPTIMA, TOLERANCES AND INDICATOR VALUES - In our data
391 set, the species optimum values were distributed along the entire gradient of soil cation
392 concentration (Fig. 3), but most of them were in the low end. This contrasts with the results of
393 earlier studies, which have found more fern species in high-cation soils than in low-cation soils
394 ([Tuomisto & Poulsen 1996](#); [Tuomisto et al. 2002, 2003b](#)). The difference is likely due to biases
395 in sampling. Our data set contained many more plots with low cation concentration than high
396 cation concentration, and most of the plots that in our data represent the high end of the gradient
397 were relatively cation-poor compared to the cation-rich soils in the western Amazonian data.
398 This probably explains why most of the genera that in earlier studies have been thought to

399 indicate cation-rich soils (e.g. *Diplazium*, *Tectaria* and *Thelypteris*) were absent or rare in our
400 data.

401 For those genera that were well represented in both geographical areas, our results
402 agreed with the earlier ones from western Amazonia. The genus *Adiantum* was found throughout
403 the soil nutrient gradient, but most *Adiantum* species occurred in intermediate to richer soils, in
404 agreement with the results of Tuomisto et al. (1998). They observed that *A. tomentosum* and *A.*
405 *pulverulentum* occur at opposite ends of the soil cation gradient and never co-occur, and this was
406 the case also in our data.

407 Species differed in how accurate they seem to be as indicators of environmental
408 variables. For example, *Trichomanes pinnatum* had a high indicator value for cation-poor soils in
409 general, and some other species of the same genus appeared as significant indicators for the finer
410 clusters within that group of sites. Although our sampling is relatively extensive, it still covers
411 only a small part of the environmental variation within Amazonia. Therefore, the optima and
412 tolerances of species shown in Fig. 3 are still preliminary, and should not be taken at face value.
413 A veiled gradient will push optimum values towards the mean of the gradient for those species
414 whose ranges extend beyond the part of the gradient sampled, so the values we obtained for the
415 species at the cation-rich end of the gradient can be expected to be especially inaccurate.
416 However, the high congruence between our results and those from western Amazonia suggest
417 that the positions of the species optima in relation to each other, and the degrees of overlap in
418 tolerance ranges, are probably rather reliable.

419 In spite of this limitation in the extreme of the soil gradient, it is noteworthy how
420 well our results on species optima agree with the suggestions made in earlier studies, although
421 the earlier datasets were much smaller, less quantitative and represented a different geographical
422 region (e.g., [Tuomisto & Poulsen 1996](#); [Tuomisto et al. 1998](#), [2002](#); [Tuomisto et al. 2003b](#);

423 [Cárdenas et al. 2007](#)). Such congruence indicates that the inferences on the edaphic preferences
424 of ferns have a good transferability across geographical regions.

425 The methods we used are based on general ecological principles and can therefore
426 be applied to any biogeographical area. The prerequisite is that the training dataset is suitable for
427 the task at hand: it needs to cover the relevant environmental gradients sufficiently well and to
428 contain an adequate number of species from the area of interest. Our present data can be used as
429 the training set for other studies in central Amazonia, but studies focusing on western or eastern
430 Amazonia should complement the training set locally. Failure to do so would compromise the
431 accuracy of the predictions, as illustrated with the relatively low prediction accuracy for the Acre
432 sites in the leave-group-out cross-validation. At least one study in Ecuadorian Amazonia (Sirén
433 et al. 2013) has produced a map of estimated soil cation concentrations without having had
434 access to direct soil data from the area of interest. Instead, they made fern inventories and used
435 data from existing inventories from other parts of NW Amazonia as the training set to estimate
436 soil cation concentrations through calibration. Then they used satellite imagery to generate
437 extrapolated soil fertility maps. These kinds of maps can be used to identify areas with different
438 site conditions, and thereafter to assess whether all the recognised habitat variation is adequately
439 represented in conservation area networks.

440 Additional data with a more complete geographical coverage will make it possible
441 to select a limited number of good indicator species that combine high environmental specificity
442 with sufficient frequency in suitable conditions ([Diekmann 2003](#)). Indicator plants reflect
443 environmental conditions as integrated over extended time periods, whereas soil samples give
444 snapshot information of the measured variables. Therefore, the species composition of an
445 indicator plant group can be expected to provide information that is relevant for plants in
446 general. The same approach could also be tested in other relatively well inventoried plant groups
447 such as palms ([Vormisto et al. 2000](#); [Costa et al. 2009](#); [Svenning 1999](#)), trees ([Pitman et al.](#)

448 [2001](#); [Castilho et al. 2006](#); [Stropp et al. 2009](#)) and gingers ([Figueiredo et al. 2013](#)). Our results
449 demonstrate that the species and environmental datasets already available in the Amazon region
450 are a good starting-point towards better tools and maps for conservation planning.

451 **Acknowledgements**

452 We are thankful to several field assistants that made this work possible. ICMBIO provided
453 permits and infrastructure facilities. Financial support to field work was provided by Biological
454 Dynamics of Forest Fragments (BDFFP), MCT/CNPq/PPG7 no. 48/2005 (led by William E.
455 Magnusson), Brazilian Program of Biodiversity Research - PPBio, CNPq/FAPEAM/PRONEX
456 project no. 673/2010 (led by William E. Magnusson, INPA), FINEP/Projeto Integrado MCT-
457 EMBRAPA (led by Ana L. K. M. Albernaz), Hidroveg Project -FAPEAM/FAPESP no.
458 1428/2010 (led by Flávia R. C. Costa and Javier Tomasella). Gabriela Zuquim was supported by
459 CNPq, CAPES and Academy of Finland (research grant to Hanna Tuomisto). We thank Lassi
460 Suominen for helpful comments on the manuscript. Many people are acknowledged for their
461 efforts to make data and analytical tools freely available. This is publication number 627 ST of
462 the BDFFP technical series.

463 **References**

- 464 Birks, H.J.B., Heiri, O., Seppä, H., & Bjune, A.E. 2010. Strengths and Weaknesses of
465 Quantitative Climate Reconstructions Based on Late-Quaternary Biological Proxies. *The open*
466 *Ecology Journal* 3: 68-110.
- 467 Birks, H.J.B., Line, J.M., Juggins, S., Stevenson, A.C. & ter Braak, C.J.F. 1990. Diatoms and pH
468 reconstruction. *Philosophical Transactions of the Royal Society of London Series B-*
469 *Biological Sciences* 327: 263-278.
- 470 Cárdenas, G.G., Halme, K.J. & Tuomisto, H. 2007. Riqueza y Distribución Ecológica de
471 Especies de Pteridofitas en la Zona del Río Yavarí-Mirín, Amazonía Peruana. *Biotropica* 39:
472 637-646.
- 473 Castilho, C.V., Magnusson, W.E., Araújo, R.N.O., Luizão, R.C.C., Luizão, F.J., Lima, A.P. &
474 Higuchi, N. 2006. Variation in aboveground tree live biomass in a central Amazonian Forest:
475 Effects of soil and topography. *Forest Ecology and Management* 234: 85-96.

- 476 Chauvel, A., Lucas, Y. & Boulet, R. 1987. On the genesis of the soil mantle of the region of
477 Manaus, Central Amazonia, Brasil. *Experientia* 43: 234-240.
- 478 Costa, F.R., Guillaumet, J.-L., Lima, A.P. & Pereira, O.S. 2009. Gradients within gradients: The
479 mesoscale distribution patterns of palms in a central Amazonian forest. *Journal of Vegetation*
480 *Science* 20: 69-78.
- 481 Costa, F.R.C., Magnusson, W.E. & Luizão, R.C. 2005. Mesoscale distribution patterns of
482 Amazonian understorey herbs in relation to topography, soil and watersheds. *Journal of*
483 *Ecology* 93: 863-878.
- 484 Costa, J.M. & Pietrobon, M.R. 2007. Pteridófitas (Lycophyta e Monilophyta) da Ilha de
485 Mosqueiro, município de Belém, estado do Pará, Brasil. *Boletim do Museu Paraense Emílio*
486 *Goeldi. Ciências Naturais* 2: 45-56.
- 487 Costa, J.M., Souza, M.G.C. & Pietrobon, M.R. 2006. Levantamento florístico das Pteridófitas
488 (Lycophyta e Monilophyta) do Parque Ambiental de Belém (Belém, Pará, Brasil). *Revista de*
489 *Biologia Neotropical* 3: 4-12.
- 490 Costa, M.A.S., Prado, J., Windisch, P.G., Freitas, C.A.A. & Labiak, P. 1999. Pteridophyta. In:
491 Ribeiro, J.E.L.S., Hopkins, M.J.G., Vicentini, A., Sothers, C.A., Costa, M.A.S., Brito, J.M.,
492 Souza, M.A., Martins, L.H., Lohmann, L.G., Assunção, P.A.C.L., Pereira, E.C., Silva, C.F. &
493 Procópio, L.C. 1999. *Flora da Reserva Ducke - Guia de identificação das plantas vasculares*
494 *de uma floresta de terra firme na Amazônia Central*. Editora INPA, Manaus, BR.
- 495 De Caceres, M. & Legendre, P. 2009. Associations between species and groups of sites: indices
496 and statistical inference. *Ecology* 90: 3566-3574.
- 497 De'ath, G. 1999. Extended dissimilarity: a method of robust estimation of ecological distances
498 from high beta diversity data. *Plant Ecology* 144: 191-199.
- 499 De'ath, G. 2002. Multivariate regression trees: a new technique for modeling species-
500 environment relationships. *Ecology* 83: 1105-1117.
- 501 Diekmann, M. 2003. Species indicator values as an important tool in applied plant ecology - a
502 review. *Basic and Applied Ecology* 4: 493-506.
- 503 Dijkshoorn J.A., Huting, J.R.M. & Tempel, P. 2005. *Update of the 1:5 million Soil and Terrain*
504 *Database for Latin America and the Caribbean (SOTERLAC; version 2.0)*. URL:
505 http://www.isric.org/sites/default/files/ISRIC_Report_2005_01.pdf
- 506 Dufrene, M. & Legendre, P. 1997. Species assemblages and indicator species: the need for a
507 flexible asymmetrical approach. *Ecological Monographs* 67: 345-366.
- 508 Edwards, P.J. 1998. The pteridophytes of the Ilha de Maracá. In: Milliken, W. & Ratter, J.A.
509 (eds.) *Maracá: the biodiversity and environment of an Amazonian rainforest*. John Wiley &
510 Sons, Chichester, UK.
- 511 Emilio, T., Nelson, B.W., Schiatti, J., Desmoulière, S.J.-M., Espírito Santo, H.M.V. & Costa,
512 F.R.C. 2010. Assessing the relationship between forest types and canopy tree beta diversity in
513 Amazonia. *Ecography* 33: 738-747.
- 514 Figueiredo, F.O.G., Costa, F.R.C., Nelson, B.W., Pimentel, T.P. 2013. Validating forest types
515 based on geological and land-form features in central Amazonia. *Journal of Vegetation*
516 *Science*. Doi:10.1111/jvs.12078
- 517 Freitas, C.A.A. & Prado, J. 2005. Lista anotada das pteridófitas de florestas inundáveis do alto
518 Rio Negro, Município de Santa Isabel do Rio Negro, AM, Brasil. *Acta Botanica Brasilica* 19:
519 399-406.
- 520 geomorfologia, pedologia, vegetação e uso potencial da terra. Departamento Nacional de
521 Produção Mineral. Rio de Janeiro.
- 522 Higgins, M.A. & Ruokolainen, K. 2004. Rapid tropical forest inventory: a comparison of
523 techniques based on inventory data from western Amazonia. *Conservation Biology* 18: 799-
524 811.
- 525 Higgins, M.A., Ruokolainen, K., Tuomisto, H., Llerena, N., Cardenas, G., Phillips, O.L.,
526 Vasquez, R. & Räsänen, M. 2011. Geological control of floristic composition in Amazonian
527 forests. *Journal of Biogeography* 38: 2136-2149.

528 Hijmans, R.J., Cameron, S.E., Parra, J.L., Jones, P.G. & Jarvis, A. 2005. Very high resolution
529 interpolated climate surfaces for global land areas. *International Journal of Climatology* 25:
530 1965-1978.

531 Hijmans, R.J., Guarino, L., Mathur, P. 2012. *DIVA-GIS Version 7.5. Manual*. URL:
532 http://www.diva-gis.org/docs/DIVA-GIS_manual_7.pdf

533 Howard, P.C., Viskanic, P., Davenport, T.R.B., Kigenyi, F.W., Baltzer, M., Dickinson, C.J.,
534 Lwanga, J.S., Matthews, R.A. & Balmford, A. 1998. Complementarity and the use of
535 indicator groups for reserve selection in Uganda. *Nature* 394: 472-475.

536 IBGE - Instituto Brasileiro de Geografia e Estatística. 2004. *Mapa de Vegetação do Brasil*. 3rd
537 Edition. URL: ftp://ftp.ibge.gov.br/Cartas_e_Mapas/Mapas_Murais/

538 Jones, M. M., Tuomisto, H., Clark, D. B. & Olivas, P. 2006. Effects of mesoscale environmental
539 heterogeneity and dispersal limitation on floristic variation in rain forest ferns. *Journal of*
540 *Ecology* 94: 181-195.

541 Jones, M.M., Ferrier, S., Condit, R., Manion, G., Aguilar, S. & Pérez, R. 2013. Strong
542 congruence in tree and fern community turnover in response to soils and climate in central
543 Panama. *Journal of Ecology* 101: 506-516.

544 Kinupp, V.F. & Magnusson, W.E. 2005. Spatial patterns in the understorey shrub genus
545 *Psychotria* in central Amazonia: effects of distance and topography. *Journal of Tropical*
546 *Ecology* 21: 363-374.

547 Legendre P. & Legendre L. 1998. *Numerical Ecology*. 2nd Edition. Elsevier, Amsterdam, NL.

548 Maciel, S., Souza, M.G.C. & Pietrobon, M.R. 2007. Licófitas e monilófitas do Bosque
549 Rodrigues Alves Jardim Botânico da Amazônia, município de Belém, estado do Pará, Brasil.
550 *Boletim do Museu Paraense Emílio Goeldi. Ciências Naturais* 2: 69-83.

551 Magnusson, W.E., Lima, A.P., Luizão, R.C.C., Luizão, F., Costa, F.R.C., Castilho, C.V. &
552 Kinupp, V.P. 2005. RAPELD: a modification of the Gentry method for biodiversity surveys
553 in long-term ecological research sites. *Biota Neotropica* 5(2). URL:
554 [http://www.biotaneotropica.org.br/v5n2/pt/download?point-of-](http://www.biotaneotropica.org.br/v5n2/pt/download?point-of-view+bn01005022005+abstract)
555 [view+bn01005022005+abstract](http://www.biotaneotropica.org.br/v5n2/pt/download?point-of-view+bn01005022005+abstract)

556 Margules, C.R., Pressey, R.L., Williams, P.H. 2002. Representing biodiversity: data and
557 procedures for identifying priority areas for conservation. *Journal of Biosciences* 27: 309-
558 326.

559 McGeoch, M.A. 1998. The selection, testing and application of terrestrial insects as
560 bioindicators. *Biological Reviews* 73: 181-201.

561 Mertens, J. 2004. *The characterization of selected physical and chemical soil properties of the*
562 *surface soil layer in the 'Reserva Ducke', Manaus, Brazil, with emphasis on their spatial*
563 *distribution*. Bachelor thesis, Humboldt University, Berlin, GE.

564 Noss, R.F. 1990. Indicators for Monitoring Biodiversity: A Hierarchical Approach. *Conservation*
565 *Biology* 4: 355-364.

566 Phillips, O. L., Vargas, P. N., Monteagudo, A. L., Cruz, A. P., Zans, M. -E. C., Sánchez, W. G.,
567 Yli-Halla, & Rose, S. 2003. Habitat association among Amazonian tree species: a landscape-
568 scale approach. *Journal of Ecology* 91: 757-775.

569 Pitman, N.C.A., Terborgh, J.W., Miles, R., Silman, P.N.V., Neill, D.A., Cerón, C.E., Palacios,
570 W.A. & Aulestia, M. 2001. Dominance and distribution of tree species in upper Amazonian
571 terra firme forests. *Ecology* 83: 2101-2117.

572 Prado, J. & Moran, R.C. 2009. Checklist of the ferns and lycophytes of Acre state, Brazil. *Fern*
573 *Gazette* 18: 230-263.

574 Quesada, C.A., Lloyd, J., Anderson, L.O., Fyllas, N.M., Schwarz, M. & Czimeczik, C.I. 2011.
575 Soils of Amazonia with particular reference to the RAINFOR sites. *Biogeosciences* 8: 1415-
576 1440.

577 RADAMBRASIL 1978. Projeto RADAMBRASIL. Vol. (1:34). Geologia,

- 578 Ruokolainen, K., Linna, A., & Tuomisto, H. 1997. Use of Melastomataceae and pteridophytes
579 for revealing phytogeographical patterns in Amazonian rain forests. *Journal of Tropical*
580 *Ecology* 13: 243-256.
- 581 Ruokolainen, K., Tuomisto, H., Macía, M. J., Higgins, M.A. & Yli-Halla, M. 2007. Are floristic
582 and edaphic patterns in Amazonian rain forests congruent for trees, pteridophytes and
583 Melastomataceae? *Journal of Tropical Ecology* 23: 13-25.
- 584 Salovaara, K.J., Cárdenas, G.G. & Tuomisto, H. 2004. Forest classification in an Amazonian
585 rainforest landscape using pteridophytes as indicator species. *Ecography* 27: 689-700.
- 586 Schulman, L., Ruokolainen, K., Junikka, L., Sääksjärvi, I.E., Salo, M., Juvonen, S., Salo, J. &
587 Higgins, M. 2007b. Amazonian biodiversity and protected areas: do they meet? *Biodiversity*
588 *and Conservation* 16: 3011-3051.
- 589 Schulman, L., Toivonen, T. & Ruokolainen, K. 2007a. Analysing botanical collecting effort in
590 Amazonia and correcting for it in species range estimation. *Journal of Biogeography* 34:
591 1388-1399.
- 592 Sirén, A., Tuomisto, H., & Navarrete, H. 2013. Mapping environmental variation in lowland
593 Amazonian rainforests using remote sensing and floristic data. *International Journal of*
594 *Remote Sensing* 34: 1561-1575.
- 595 Sombroek, W. G. 2001. Spatial and temporal patterns of Amazon rainfall. Consequences for the
596 planning of agricultural occupation and the protection of primary forests. *Ambio* 30: 388-396.
- 597 Stropp, J., ter Steege, H., Malhi, Y., ATDN & RAINFOR. 2009. Disentangling regional and
598 local tree diversity in the Amazon. *Ecography* 32: 46-54.
- 599 Suominen, L., Ruokolainen, K., Tuomisto, H., Llerena, N. & Higgins, M.A. 2013. Predicting soil
600 properties from floristic composition in western Amazonian rainforests: performance of k-
601 nearest neighbour estimation and weighted averaging calibration. *Journal of Applied Ecology*.
602 Doi: 10.1111/1365-2664.12131.
- 603 Svenning, J.-C. 1999. Microhabitat specialization in a speciesrich palm community in
604 Amazonian Ecuador. *Journal of Ecology*. 87: 55-65.
- 605 ter Braak, C.J.F. & Juggins, S. 1993. Weighted averaging partial least squares regression (WA-
606 PLS): an improved method for reconstructing environmental variables from species
607 assemblages. *Hydrobiologia* 269/270: 485-502.
- 608 ter Braak, C.J.F. & van Dam, H. 1989. Inferring pH from diatoms: a comparison of old and new
609 calibration methods. *Hydrobiologia* 178: 209-223.
- 610 Tuomisto, H. 2006. Edaphic niche differentiation among Polybotrya ferns in Western Amazonia:
611 implications for coexistence and speciation. *Ecography* 29: 273-284.
- 612 Tuomisto, H., & Poulsen, A.D. 1996. Influence of edaphic specialization on the distribution of
613 pteridophyte in neotropical forests. *Journal of Biogeography* 23: 283-293.
- 614 Tuomisto, H., Poulsen, A.D. & Moran, R.C. 1998. Edaphic distribution of some species of the
615 fern genus Adiantum in Western Amazonia. *Biotropica* 30: 392-399.
- 616 Tuomisto, H., Poulsen, A.D., Ruokolainen, K., Moran, R.C., Quintana, C., Celi, J. & Cañas, G.
617 2003c. Linking floristic patterns with soil heterogeneity and satellite imagery in Ecuadorian
618 Amazonia. *Ecological Applications* 13: 352-371.
- 619 Tuomisto, H., Ruokolainen, K. & Yli-Halla, M. 2003a. Dispersal, environmental, and floristic
620 variation of Western Amazonian forests. *Science* 299: 241-244.
- 621 Tuomisto, H., Ruokolainen, K., Aguilar, M. & Sarmiento, A. 2003b. Floristic patterns along a
622 43-km long transect in an Amazonian rain forest. *Journal of Ecology* 91: 743-756.
- 623 Tuomisto, H., Ruokolainen, K., Poulsen, A.D., Moran, R.C., Quintana, C., Cañas, G. & Celi, J.
624 2002. Distribution and diversity of pteridophytes and Melastomataceae along edaphic
625 gradients in Yasuní national park, Ecuadorian amazonia. *Biotropica* 34: 516-533.
- 626 Tuomisto, H., Ruokolainen, L. & Ruokolainen, K. 2012. Modelling niche and neutral dynamics:
627 on the ecological interpretation of variation partitioning results. *Ecography* 35: 961-971.

- 628 Vormisto, J., Phillips, O., Ruokolainen, K., Tuomisto, H. & Vásquez, R. 2000. A comparison of
629 fine-scale distribution patterns of four plant groups in an Amazonian rainforest. *Ecography*
630 23: 349-359.
- 631 Zuquim, G., Costa, F.R.C., Prado, J. & Braga-Neto, R. 2009a. Distribution of pteridophyte
632 communities along environmental gradients in Central Amazonia, Brazil. *Biodiversity and*
633 *Conservation* 18: 151-166.
- 634 Zuquim, G., Prado, J. & Costa, F.R.C. 2009b. An annotated checklist of ferns and lycophytes
635 from the Biological Reserve of Uatumã, an area with patches of rich-soils in central
636 Amazonia, Brazil. *Fern Gazette* 18: 286-306.
- 637 Zuquim, G., Tuomisto, H., Costa, F.R.C., Prado, J., Magnusson, W.E., Pimentel, T., Braga-Neto,
638 R. & Figueiredo, F.O.G. 2012. Broad Scale Distribution of Ferns and Lycophytes along
639 Environmental Gradients in Central and Northern Amazonia, Brazil. *Biotropica* 44: 752-762.
640

641 Table 1. Mean, standard deviation (\pm) and range (in parentheses) of species richness and
 642 environmental variables per plot for 326 plots in Brazilian Amazon. Climatic data was obtained
 643 from WorldClim database in 2.5 arc-minutes resolution (ca. of 4.7 km).

	Values
Species richness	4.9 \pm 3.6 (0-20)
Species abundance (individuals)	90 \pm 153 (0-1131)
Sum of Bases (cmol _c kg ⁻¹)	1.34 \pm 4.16 (0.08-38.11)
Clay (%)	29 \pm 22 (0.5-90)
Silt (%)	25 \pm 18 (0.5-76)
Sand (%)	47 \pm 25 (1.7-90)
Temperature annual range (°C)	12.4 \pm 2 (10.2-19.4)
Annual precipitation (mm)	2177 \pm 270 (1633-2655)
Precipitation seasonality (coefficient of variation)	57 \pm 13 (33-80)
Precipitation of the wettest quarter	925 \pm 57 (815-1082)

644

645 Table 2. Prediction accuracy given by the Root Mean Squared Error (RMSE) and coefficient of determination (R^2) of the regressions between predicted
646 and observed edaphic properties in 305 plots in Brazilian Amazonia. The accuracy of the predictions for the k -Nearest-Neighbours (k -NN) method
647 reported here is based on $k=4$ neighbours. The deshrinking method applied in Weighted Averaging (WA) was inverse deshrinking. Down-weighting in
648 WA was done by inversely-weighting species optima by their tolerances along the environmental gradient when generating the predicted values. In k -
649 NN, down-weighting was done by inversely-weighting the selected neighbouring plots by their floristic similarity to the focal. Cross-validation
650 methods were bootstrap (k -NN) and Leave-one-out (WA) except when mentioned. "Crossval=lgo" refers to Leave-group-out cross-validation method
651 and "Pres.-Abs." refers to presence-absence input species data.

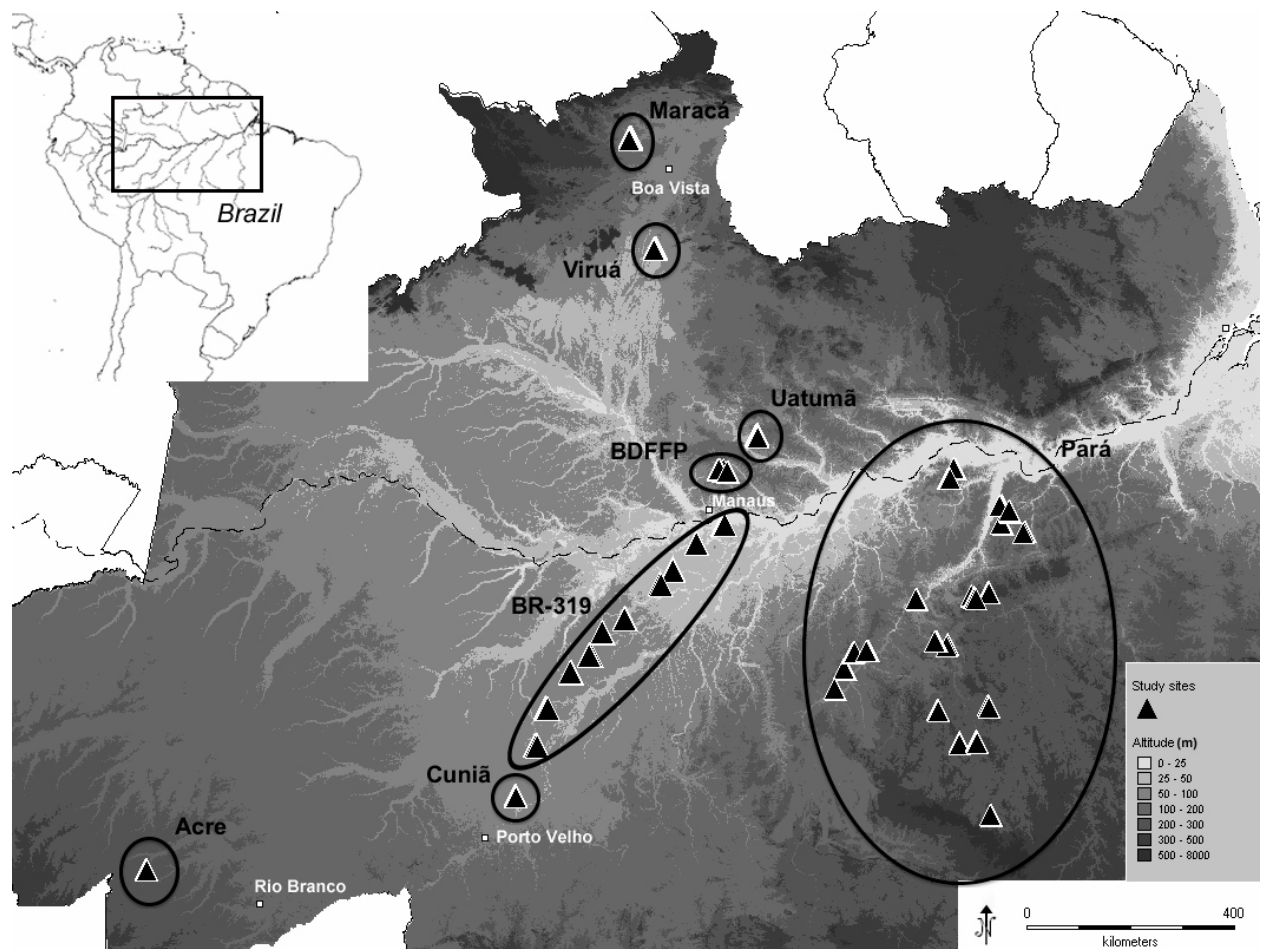
species input data	downweighting	Log (Sum of Bases) crossval=lgo				Clay		Silte		Sand		
		RMSE	R^2	RMSE	R^2	RMSE	R^2	RMSE	R^2	RMSE	R^2	
K-nn	Abundance	no	0.31	0.68	0.31	0.59	20.13	0.35	16.28	0.39	24.53	0.21
		similarity	0.33	0.64	0.32	0.59	20.10	0.35	16.82	0.35	25.15	0.14
	Pres.-Abs.	no	0.28	0.74	0.30	0.62	18.76	0.46	14.80	0.48	23.98	0.24
		similarity	0.28	0.75	0.31	0.64	19.54	0.41	15.65	0.43	24.24	0.19
WA	Abundance	no	0.29	0.65	0.33	0.55	18.67	0.30	14.53	0.37	22.15	0.18
		tolerance	0.27	0.70	0.36	0.46	18.10	0.34	14.00	0.41	21.79	0.20
	Pres.-Abs.	no	0.29	0.65	0.32	0.55	17.58	0.38	13.90	0.42	21.75	0.21
		tolerance	0.27	0.68	0.33	0.54	17.92	0.35	14.45	0.38	21.93	0.20

652

653

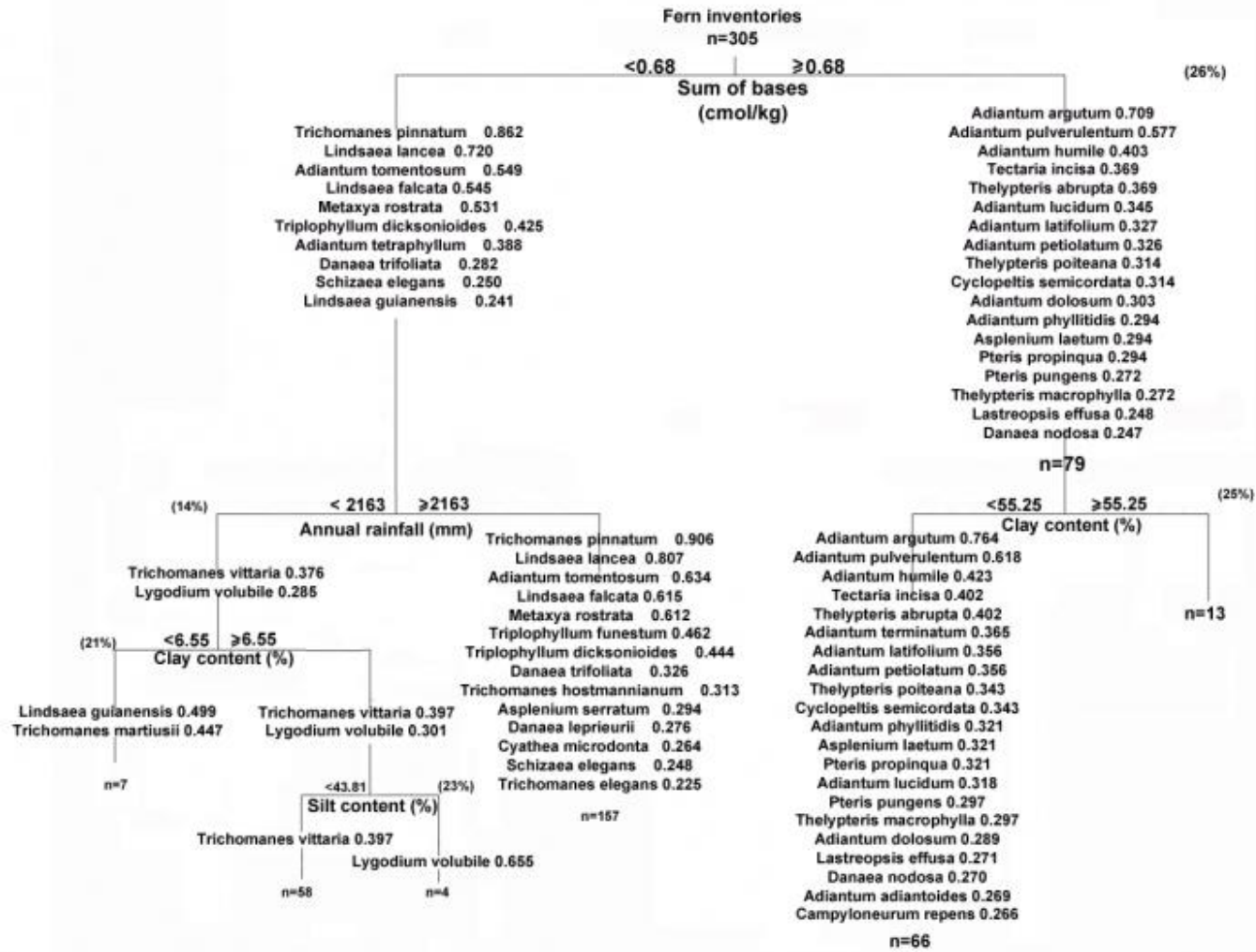
654

655

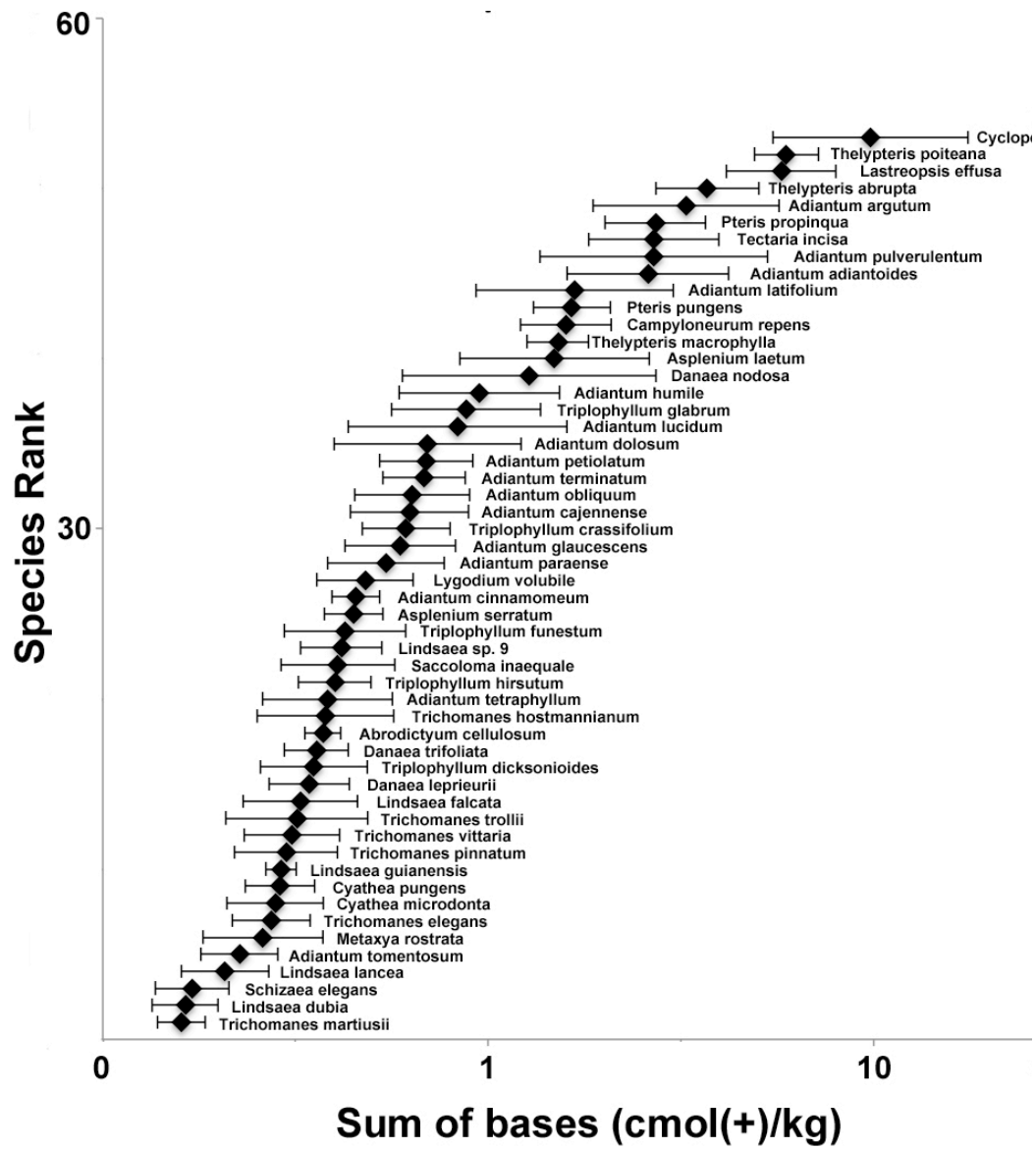


656

657 Figure 1. Location of 326 plots established in Brazilian Amazonia (black triangles) divided in
658 eight regions. Black lines are country boundaries and the dashed line is the main channel of the
659 Amazonas River. Gray scale represents altitude according to SRTM. More detailed description
660 of the circled regions in S1.



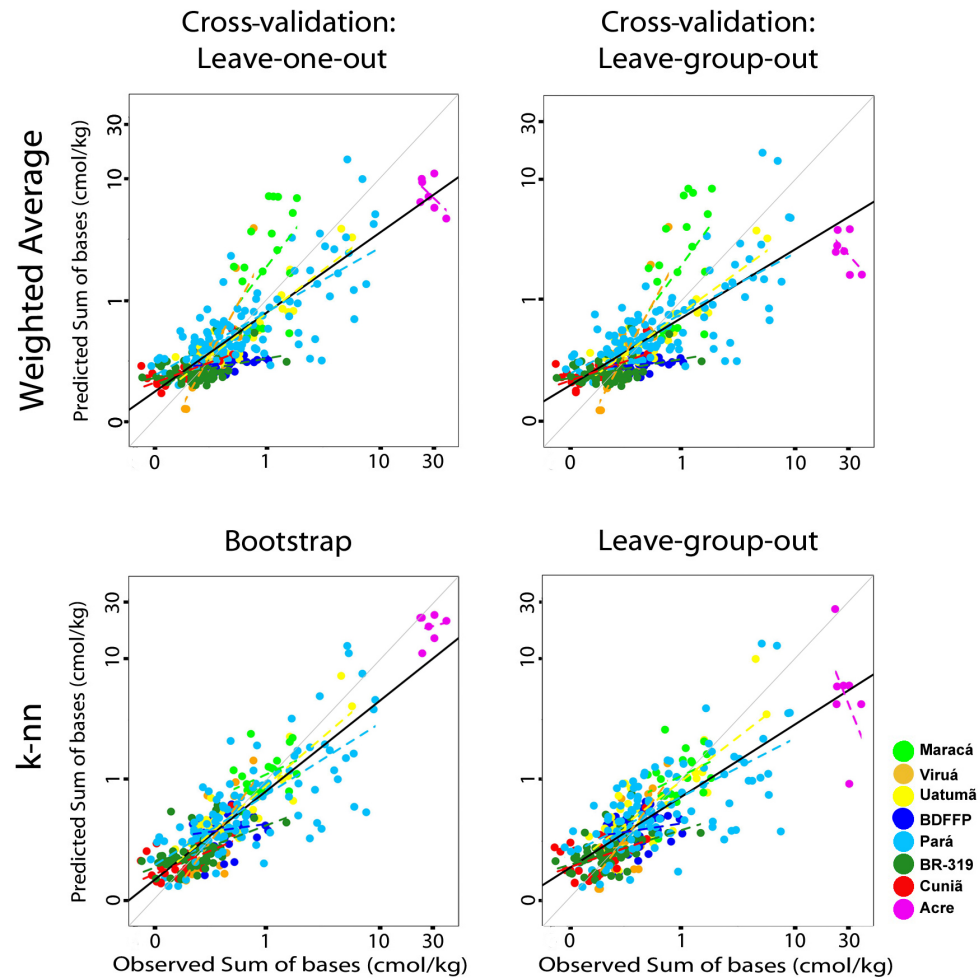
662 Figure 2. Results of the distance-based Multivariate Regression Tree (db-MRT) of fern inventories in 305 plots in Brazilian Amazonia. A list of
663 significant ($\alpha \leq 0.05$) indicator species followed by their indicator values is presented for each branch. The percentage of improvement in model
664 performance given by each division is in parentheses.



666

667 Figure 3. Estimated optima and tolerances of fern species along the sum of bases
 668 305 plots in Brazilian Amazonia based on abundance data. Values on the x-axis
 669 a logarithmic scale.

670



671

672 Figure 4. Predicted *vs.* observed sum of bases in 305 plots in Brazilian Amazonia. The solid black lines corresponds to the regression line for all the
 673 predicted *vs.* observed values. Dashed lines corresponds to the regression lines based on the same predictions but shown for each regional subset of the

674 plots to illustrate the variation among regions. The 1:1 line used in accuracy assessment to calculate the root mean squared errors (RMSEs) is shown in
675 gray. The deshrinking method applied in Weighted Averaging (WA) was inverse deshrinking. Both axes are on a logarithmic scale.