



HAL
open science

Identification of ecological thresholds from variations in phytoplankton communities among lakes: contribution to the definition of environmental standards

V. Roubex, P.A. Danis, T. Feret, J.M. Baudoin

► **To cite this version:**

V. Roubex, P.A. Danis, T. Feret, J.M. Baudoin. Identification of ecological thresholds from variations in phytoplankton communities among lakes: contribution to the definition of environmental standards. *Environmental Monitoring and Assessment*, 2016, 188 (4), pp.246. <10.1007/s10661-016-5238-y>. <hal-01529258>

HAL Id: hal-01529258

<https://hal.science/hal-01529258v1>

Submitted on 30 May 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

1 **Identification of ecological thresholds from variations in phytoplankton**
2 **communities among lakes: contribution to the definition of environmental**
3 **standards**

4 Vincent Roubeix¹, Pierre-Alain Danis², Thibaut Feret³ and Jean-Marc Baudoin²

5 ¹Irstea, UR RECOVER, Pôle Onema-Irstea hydroécologie plans d'eau, centre d'Aix-en-Provence, 3275 route
6 Cézanne, F-13612 Le Tholonet Aix-en-Provence, France

7 ²Onema, Pôle Onema-Irstea hydroécologie plans d'eau, F-13612 Le Tholonet Aix-en-Provence, France

8 ³Irstea, UR EABX, centre de Bordeaux, F-33612 Gazinet Cestas, France

9 Corresponding author: vincent.roubeix@irstea.fr +33 4 42 66 79 35

1 Abstract

2 In aquatic ecosystems, the identification of ecological thresholds may be useful for managers as it
3 can help to diagnose ecosystem health and to identify key levers to enable the success of
4 preservation and restoration measures. A recent statistical method, gradient forest, based on
5 random forests, was used to detect thresholds of phytoplankton community change in lakes along
6 different environmental gradients. It performs exploratory analyses of multivariate biological and
7 environmental data to estimate the location and importance of community thresholds along
8 gradients. The method was applied to a dataset of 224 French lakes which were characterized by 29
9 environmental variables and the mean abundances of 196 phytoplankton species. Results showed
10 the high importance of geographic variables for the prediction of species abundances at the scale of
11 the study. A second analysis was performed on a subset of lakes defined by geographic thresholds
12 and presenting a higher biological homogeneity. Community thresholds were identified for the most
13 important physico-chemical variables including water transparency, total phosphorus, ammonia,
14 nitrates and dissolved organic carbon. Gradient forest appeared as a powerful method at a first
15 exploratory step, to detect ecological thresholds at large spatial scale. The thresholds that were
16 identified here must be reinforced by the separate analysis of other aquatic communities and may be
17 used then to set protective environmental standards after consideration of natural variability among
18 lakes.

19

20 Keywords : community thresholds; environmental standards; gradient forest; lakes; phytoplankton

21

22

1 Introduction

2 The response of an ecosystem to a gradual change in environmental conditions may be smooth, in
3 proportion to the change, or abrupt if a critical level is reached, with potentially an hysteresis effect
4 when the change is reversed (Scheffer and Carpenter 2003). The hypothesis of multiple stable states
5 of ecosystems and communities (May 1977) argues in favor of non-linear responses with ecological
6 thresholds. This vision of ecosystems dynamics is supported by simple mathematical models and
7 empirical observations of sudden shifts from one state to another. In the cases of shallow lakes, the
8 existence of several states characterized by phytoplankton and macrophyte species has been
9 demonstrated, with transitions possibly due to changes in nutrient concentrations (Scheffer et al.
10 1997; Scheffer et al. 2003). Following the idea of discontinuous variations triggered by controlling
11 variables, numerous scientific studies have focused in the last 15 years on the research of thresholds
12 from the response of aquatic communities to environmental parameters, such as water
13 transparency, phosphorus and nitrogen concentrations or pH (Graham et al. 2004; Holt et al. 2003;
14 Richardson et al. 2007). The identification of thresholds in the gradients of human-influenced
15 variables is particularly interesting for environmental management, in order to anticipate dramatic
16 ecological changes in case of degradation or to set targets for restoration (Chambers et al. 2012a;
17 Soranno et al. 2008; Vollenweider 1975). In the European Water Framework Directive (WFD)
18 (European Commission 2000) concerning all types of water masses including lakes, member states
19 have to set environmental standards for general physico-chemical parameters which support the
20 ecological assessment made from biological quality indices. These standards may be simply derived
21 from the distribution of measured values by the use of quantiles or by the division of observed
22 gradients (e.g. in Chambers et al. 2012a), without any link to biological elements. However, it is more
23 relevant for lake management to propose standards based on ecological thresholds, through the
24 analysis of community responses to water quality gradients (Poikane et al. 2014; Solheim et al. 2008;
25 Penning et al. 2008; Free et al. 2006).

26 An ecological threshold may be generally defined as a critical point where moderate variations of an
27 environmental parameter produce large responses in ecosystem state (Groffman et al. 2006).
28 Thresholds can be *a posteriori* identified from time-series analysis (Andersen et al. 2009), but they
29 can be also detected from the spatial comparisons of systems in environmental gradients (Catalan et
30 al. 2009; Holt et al. 2003; Soranno et al. 2008). Amongst the response variables for lakes, those
31 related to the composition of aquatic communities can be used to demonstrate the existence of
32 thresholds.

1 Ecological thresholds imply non-linear relationships between environmental drivers and some
2 biological variables. A large set of statistical methods has been proposed for threshold identification
3 in environmental gradients (Brenden et al. 2008; Dodds et al. 2010; Qian 2014) with applications in
4 various ecological contexts, mostly for streams (Black et al. 2011; Chambers et al. 2012b; Evans-
5 White et al. 2009; Richardson et al. 2007; Smith and Tran 2010). The different approaches found in
6 the literature can be classified according to their exploratory character and *a priori* knowledge of the
7 relationships between environment and communities. Most of them focus on one environmental
8 factor which is known to influence communities directly in water (e.g. phosphorus, pH) or from the
9 watershed (e.g. urbanization, agriculture). Concerning the biological variables, some studies use the
10 abundance of one sensitive species or a selection of species according to their indicator value (e.g.
11 King et al. 2011; Richardson et al. 2007); in others the response variables are aggregate or synthetic
12 metrics (e.g. Black et al. 2011; Evans-White et al. 2009). In some cases, the search for thresholds is
13 carried out separately in different groups of sites, making the hypothesis that some typological
14 features may influence the results (Utz et al. 2009).

15 In this study, an exploratory approach was followed to find thresholds through an analysis of
16 multivariate biological and environmental data, minimizing preliminary variable selection. The
17 phytoplankton communities of various types of lakes from different regions of France were
18 considered together. Phytoplankton can be found in every lake and might be particularly sensitive to
19 the commonly measured water physico-chemical parameters. It is therefore well-suited for an
20 analysis of thresholds in environmental gradients on a large spatial scale. The abundances of all
21 phytoplankton species were taken into account, assuming that every species can be an indicator of
22 an environmental parameter. Thus, the thresholds that were sought out here, rather refer to the
23 concept of community threshold that can be defined as a zone in an environmental gradient in which
24 the rate of change in community structure is enhanced relative to the rest of the gradient, as a result
25 of sharp increases or decreases in the abundances of several species (Baker and King 2010; Catalan et
26 al. 2009).

27 With the method used in the present study, called gradient forest (Ellis et al. 2012; Pitcher et al.
28 2012), it is possible to include in the analysis a large number of potentially interacting environmental
29 factors, without any assumption on their real influence on the biological communities. Phytoplankton
30 data could therefore be simultaneously related either to human-influenced or natural factors.
31 Gradient forest is an extension to the community level of random forest (Breiman 2001) which
32 estimates variable importance for one species and allows to detect at which levels in an
33 environmental gradient the main changes in abundance occur. By aggregating the results of random
34 forests for all species reported in a survey, gradient forest orders environmental variables by

1 importance for communities and indicates along gradients the compositional turnover (cumulative
2 rate of change for all species), whose potential peaks may indicate community thresholds.

3 The objectives of this work were to apply a recent statistical method, gradient forest, to a broad scale
4 phytoplankton database, in order to identify potential thresholds in the gradients of the most
5 important physico-chemical parameters, especially those related to eutrophication. The final goal is
6 to give ecological justifications for the setting of environmental standards used in lake management.

7 **Materials and methods**

8 **Database**

9 **Environmental description**

10 The analysis was initially performed on 224 lakes located in France (Fig. 1). All types of lentic water
11 masses were included (natural lakes, reservoirs, aquaculture ponds, gravel pits) provided that lake
12 surface area was above 0.5 km². Data were provided by French water basin agencies that organize
13 the monitoring of French lakes according to standard protocols and as required by the European
14 Water Framework Directive (European Commission 2000). The 29 environmental variables which
15 were taken into account in the analysis are presented in Table 1. This selection results from a trade-
16 off between sufficient ecological description and the maximization of lake-sample size given the
17 incompleteness of the national data set. The diversity of lakes was characterized by some basic
18 geographic or physical variables such as latitude, longitude, altitude and lake maximal depth. Field
19 physico-chemical measurements were done in each lake at the point of greatest depth. Water
20 transparency was assessed using a Secchi disk (NF EN ISO 7027) and the depth of the euphotic zone
21 (Zeu) was assumed to be 2.5 times Secchi depth. Average values of temperature, pH, conductivity
22 and O₂ saturation were derived from the integration of vertical profiles over Zeu (1 m-measurement
23 interval). Nutrients, dissolved organic carbon and alkalinity analyses were performed in an integrated
24 water sample collected in the euphotic zone, following national and European standards (NF EN ISO
25 10304, NF EN ISO 6878, NF EN 1484, NF EN ISO 9963). Nutrient concentrations below the
26 quantification limit (LQ) were given an arbitrary value of LQ/2 in order to keep the most oligotrophic
27 lakes in the analysis.

28 **Phytoplankton communities**

29 Phytoplankton was analyzed from the integrated water sample used for nutrients (Laplace-Treytoure
30 et al. 2009). The composition of communities was determined at a specific level under inverted
31 microscope according to Utermöhl's method (NF EN 15204). Cell counts of each species were

1 converted into biovolumes (in $\text{mm}^3 \cdot \text{L}^{-1}$) using standard specific cell volumes. The occurrence of more
2 than 600 species was reported in the lakes. Rare species which were absent in more than 95 % of
3 sites were excluded from the analysis, so that the final species number was 196.

4 **Data temporal aggregation**

5 Data were organized in campaigns along annual cycles between 2006 and 2012. There were 4
6 campaigns per annual cycle and 3 of them during the most productive months (from April to
7 October). Physico-chemical measurements and phytoplankton samples were realized during the
8 same campaigns. For 80 % of the lakes, data were available only for one annual cycle, whereas the
9 rest of the lakes had data for 2 (14 %), 3 (4.5 %) or 4 (1 %) annual cycles. As the approach in this study
10 is based on an inter-lake comparison, seasonal and inter-annual variations were not taken into
11 account in the analysis. Thus, arithmetic means of phytoplankton biovolumes from all campaigns in
12 each lake were considered. All environmental data except longitude, latitude, water temperature, pH
13 and O_2 saturation level were \log_{10} -transformed to facilitate graphical visualization of the results.
14 Then, annual medians and maximal or minimal values (when relevant) of water physico-chemical
15 variables were computed to describe abiotic conditions for phytoplankton in the lakes (Table 1).
16 According to Ellis et al (2012), all variables were included in the GF analysis without selection, even if
17 some of them were highly correlated. A sensitivity analysis with and without some correlated
18 variables showed very little effect on the main results of gradient forest.

19 **Data analysis**

20 **Method description**

21 Gradient forest (Ellis et al. 2012) is a computer intensive method based on classification and
22 regression tree analysis (De'ath and Fabricius 2000; Breiman et al. 1984). Regression trees repeatedly
23 partitions the values of a single response variable (e.g. one species' abundance) into two mutually
24 exclusive groups. These two groups correspond to the values of an explanatory variable which are
25 below and above a split value. The explanatory variable splitting the data and the split value are
26 determined so that the homogeneity of the groups is maximized as regards the response variable.
27 Each split results in two branches and the recursive partitioning of sub-groups gives rise to a tree
28 (e.g. in Fig. 2.2). At this step, it is important to notice that a split value can be interpreted as a
29 threshold in an environmental gradient, from which the response variable changes substantially. The
30 importance of the threshold can be measured by the fit improvement, i.e. the deviance reduction
31 resulting from the split.

32 Regression trees can be used to predict the abundance of a species given the values of
33 environmental variables. However, the results may depend on the observations and predictors

1 considered for the analysis. To gain more stability, it is preferable to repeat a high number of times
2 the construction of the tree, with random selection of observations and predictors to be considered,
3 and finally average the predictions of all trees (making a forest). In random forest (Prasad et al. 2006;
4 Breiman 2001), each tree uses a random sample of the observations and each split is determined
5 from a random subset of predictors. The observations which are not taken into account for the
6 construction of a tree are used to cross-validate the performance of the tree. The predictive
7 performance of the forest is the mean cross-validated performance of all trees (R^2). The importance
8 of each predictor is estimated by the increase in mean square prediction error when the values of
9 the predictor are randomly permuted.

10 Random forest can be applied to each species occurring in a survey using the same environmental
11 variables. Then gradient forest combines the results of all random forests to derive information at
12 the community level. Thus, the community level importance of an environmental variable is the
13 average of all species level importances weighted by species R^2 (Fig. 2.4b). This overall importance in
14 the survey can be fractionated along the variable gradient into quanta of community change,
15 considering the splits due to the variable in all trees of all forests. The location of the quanta in the
16 gradient is given by the split values and the magnitude, by the aggregation of the associated fit
17 improvements. The aggregation of fit improvements takes into account the variable importance and
18 the species R^2 in each random forest (see Ellis et al 2012 for details). The results can be represented
19 in the form of a barplot of aggregated fit improvements in the environmental gradient (Fig. 2.4a). To
20 help interpretation, the data can be smoothed with a density curve (density of splits). The density of
21 splits, whose peaks may mark a community threshold, can be biased by a non-uniform distribution of
22 data along the gradients of environmental variables. Therefore the ratio of split-over-data densities
23 must be considered to better identify and characterize community thresholds (Fig. 2.4a).

24 **Computer implementation**

25 Gradient forest was performed using the two R packages 'extendedForest' and 'gradientForest'
26 (R_Core_Team 2013; Ellis et al. 2012). A total of 500 trees were generated for each random forest.
27 The split criterion was the sum of square deviations about the mean. In order to stabilize variance,
28 biovolume data were log-transformed after the addition of the minimal strictly positive value for
29 each species. For the estimation of variable importance, a conditional approach was followed in
30 order to limit the importance of correlated variables (Strobl et al. 2008; Ellis et al. 2012). For each
31 explanatory environmental variable, constrained permutations were carried out in each tree within
32 partitions of correlated variables (>0.5 Pearson) obtained after up to 5 splits. For density estimation,
33 the environmental gradients were divided into 201 bins. Local regression was applied to log-
34 transformed biovolume data using R 'lowess' function. The probability of presence along a gradient

1 was calculated in the following way. The numbers of measurements and occurrences in each bin of
2 the gradient were first determined. Then, the variations of these numbers were smoothed using R
3 density function with a gaussian kernel and the same bandwidth as used for the density of data in
4 the split density plot. Finally, the density of occurrences was divided by the density of measurements
5 to get the probability of presence.

6 **Results**

7 **National dataset**

8 Gradient forest was first applied to the whole data set. Using biovolume data, 147 species had a
9 positive R^2 and the mean value was 0.20. Indeed in gradient forest, species R^2 may be 0 or even
10 negative (Ellis et al. 2012). Seven species had R^2 above 0.5 and were particularly well predicted by the
11 29 environmental variables considered (Fig. 3a). Surprisingly, latitude and longitude were clearly the
12 most important variables for phytoplankton abundance (Fig. 3b). The next variables in order of
13 decreasing importance, i.e Alk, PO₄, NH₄, Secchi, TP and DOC, had close values that were
14 approximately 3-fold lower than latitude and longitude. The distribution of splits along the
15 geographic gradients revealed two major thresholds in the dataset, a broad one around latitude 44°N
16 and a more precise one at longitude 5°E (Fig. 4). The integration of the ratio of densities for each
17 species involved in gradient forest (specific cumulative importance, Fig. 4) confirms the existence of
18 spatial community thresholds. The concentrations of vertical sections of the species curves just
19 above 44°N and at 5°E shows that several species exhibited breakpoints in their relations to latitude
20 and longitude.

21 **Reduced dataset and physico-chemical thresholds**

22 Given the importance of geographic variables and since the focus was rather on water physico-
23 chemical characteristics, a second analysis was performed on a reduced data set with higher
24 biological homogeneity. In order to include the highest number of lakes as possible, the lakes located
25 between the two thresholds, with latitude higher than 44.5°N and longitude lower than 4.5°E, were
26 selected (Fig. 1). These regionalisation of the study excluded the lakes of the east and south parts of
27 the country and corresponds approximately to the Atlantic biogeographic area. Thus, 129 lakes
28 remained in the analysis with 147 species associated. Considering water quality, the reduced data set
29 was characterized by lower water transparency (median = 1.1 m) and alkalinity (0.74 meq.L⁻¹) and
30 higher DOC (7.1 mg C.L⁻¹) and TP (44.7 µg.L⁻¹) concentrations. Following GF analysis, 77 species had a
31 positive R^2 and the mean value was 0.15, slightly lower than that of the whole data set (Fig. 5a). The
32 seven species with the highest R^2 (> 0.3) were *Phacotus lenticularis*, *Neodesmus danubialis*,

1 *Monoraphidium arcuatum*, *Pediastrum duplex*, *Nitzschia acicularis*, *Acutodesmus obliquus* and
2 *Pediastrum tetras*. The order of variable importance was notably modified: although geographic
3 variables remained important, several water physico-chemical parameters had the same levels of
4 importance, such as Alk, TP, NH₄, NO₃max and DOC. Secchi became the most important variable
5 (Fig. 5b). Detailed results of this second analysis are given for the following variables which are
6 amongst the most important and which are largely influenced by human activities, particularly in a
7 context of eutrophication: Secchi, TP, NH₄, NO₃max and DOC. For each variable, in addition to global
8 GF results concerning all species, the response of one species to the gradient was also provided as an
9 example with a simple plot of its log-transformed biovolume. This species was selected according to
10 two criteria: (1) the variable of interest was the most important for the species, otherwise it would
11 be difficult to see a threshold from the plot of all biovolume data (without any preliminary split) and
12 (2) the species contributed to the main thresholds detected by gradient forest.

13 In the Secchi gradient, splits are clearly restricted to the zone below 1.3 m (Fig. 6). It means that
14 above this value, water transparency does not influence species abundance. The zone where the
15 most important changes occur is around 0.4 m, where the ratio of densities peaks and several
16 species exhibit a sudden rise in specific cumulative importance. The Chlorophyceae *Monoraphidium*
17 *arcuatum* presents a threshold in its response to the Secchi gradient around 1 m. This can be seen
18 through the raw abundance data of this species since Secchi is by far the most important variable
19 determining its biovolume. The threshold can be identified on the graph by a decrease in the slope of
20 the regression curve and by a more frequent absence of the species resulting in a drop in the
21 probability of presence.

22 The density curves for TP reveal a threshold at 50 µg.L⁻¹ and a larger zone of change from 100 µg.L⁻¹
23 up to the end of the gradient (Fig. 6). The threshold results from the change in abundance of some
24 species at the middle of the gradient, as illustrated by the specific cumulative importances. Among
25 these species, the diatom *Aulacoseira granulata* shows the most important change. Its biovolume is
26 mostly determined by TP. The regression curve shows a faster increase in biovolume at 50 µg.L⁻¹ and
27 the species is present in almost all sites located further in the gradient.

28 Two main peaks of split density appear in the NH₄ gradient at 50 and 125 µg.L⁻¹ (Fig. 7). Taking into
29 account the density of data, the most important is rather the second one which corresponds to steep
30 rises in specific cumulative importance for several species. From the beginning of the gradient,
31 crossing this threshold implies for the Chlorophyceae *Monoraphidium tortile*, higher biovolumes and
32 a sharp increase in the presence rate.

1 For NO₃max, the most important splits are at the end of the gradient after 10 mg.L⁻¹ (Fig. 7), giving
2 rise to a single large peak of the ratio of densities. The specific cumulative importance graph shows
3 that most species respond only to large concentrations of NO₃max. The diatom *Stephanodiscus*
4 *hantzschii* exhibits one of the most important response between 20 and 30 mg.L⁻¹, characterized by
5 higher abundance and a more frequent presence in the lakes.

6 In the DOC gradient, important splits only occur from 10 mg.L⁻¹ onwards (Fig. 8). Considering the
7 density of data, it appears that most of the lakes are below the community thresholds which can be
8 located at 11 and 18 mg.L⁻¹. The Chlorophyceae *Scenedesmus ecornis* shows a first breakpoint in its
9 relation to DOC at 11 mg.L⁻¹, defined by marked increases in abundance and presence rate. A second
10 step in specific cumulative importance corresponds to the second threshold and is due to even
11 higher abundance and no absence.

12 Discussion

13 Statistical approach

14 The basic method used in gradient forest to find thresholds is the non-parametric deviance reduction
15 (NDR) (Qian et al. 2003; Brenden et al. 2008) which is used in regression trees to determine splits.
16 NDR has been used in many studies on ecological thresholds with univariate data (Chambers et al.
17 2012b; Evans-White et al. 2009; Holt et al. 2003; Smith and Tran 2010; Soranno et al. 2008). It is well
18 suited for data following step function models but it is less effective in finding thresholds when
19 biological response patterns are smoother with less abrupt thresholds (Brenden et al. 2008).
20 Especially, the method may improperly detect thresholds in the case of linear relations (Fig. 2) (Daily
21 et al. 2012). Nevertheless, when the abundance of individual species is considered, the data are
22 often sparse and discontinuous with many zeros, and they generally exhibit stair-step patterns in
23 environmental gradients rather than linear trends. Gradual responses are more often observed when
24 species data are aggregated, e.g. into higher level taxonomic groups or functional metrics (King and
25 Baker 2010; Utz et al. 2009). It is then important to examine graphically each relationship and to
26 question the relevancy of a threshold. When response models differ largely from a step-function, the
27 concept of threshold corresponds more generally to a change point (or change zone), rather to the
28 classic definition given above, i.e. a small change in an explanatory variable giving rise to large
29 variations in the response variable.

30 Another limit of the deviance reduction approach is its sensitivity to the distribution of data in the
31 gradient. It was numerically demonstrated that such an analysis do not find the same thresholds
32 whether sampling is uniform or not (Cuffney and Qian 2013; Daily et al. 2012). In gradient forest, the

1 use of the ratio of split-over-data densities adjusts for the bias created by data skewness and leads to
2 a better identification of ecological thresholds (e.g. in Fig. 2).

3 Unlike most other threshold detection methods, gradient forest takes into account all species
4 reported in a survey and assumes that each taxon can be an indicator of one of the drivers entering
5 the analysis. In this point of view, it is similar to TITAN method (Baker and King 2010) which aims at
6 finding a community threshold along a gradient; i.e. a zone in the gradient where there are
7 'synchronous' changes in the abundances of many species. Congruence in species' responses to a
8 gradient may be unlikely given the diversity of ecological traits among taxa (Luck 2005). However, the
9 existence of community thresholds may find a justification in an evolutionary point of view. Indeed,
10 concordant declines of many species may be expected when an environmental variable (especially
11 human-influenced) goes outside the range of variations in which these species have co-evolved (King
12 and Baker 2010). In gradient forest, the absence of threshold in a gradient may mean either that all
13 species respond randomly to the environmental variable which is not important (Fig. 2), or that the
14 species have different thresholds along the gradient.

15 The changes in abundance due to limits in the biogeographic extension area of species, may interfere
16 in the search for environmental thresholds (false zeros) when analysis is conducted at species level
17 and at large spatial scale (Cuffney and Qian 2013; Utz et al. 2009). Aggregating species into metrics or
18 higher taxonomic levels decreases spatial dependency but implies a loss of information. At large
19 spatial scales, phytoplankton data are often analyzed at low taxonomic resolution level to cope with
20 unharmonized species names and differences in taxonomic resolution among regional databases
21 (Maileht et al. 2013). However, it may be always preferable to address the issue of thresholds at
22 species level because species aggregation implies more difficult threshold identification due to the
23 linearization of biological response.

24 The problem of spatial dependence of species abundance can be extended to any environmental
25 variable interacting with the variable under study. For example, the identification of thresholds for
26 zooplankton communities along acidity gradients in lakes can be confounded by morphometric
27 factors such as lake depth or area (Holt et al. 2003). Regression tree analysis which is the basis of
28 gradient forest, can deal with many explanatory variables and account for complex interactions
29 (De'ath and Fabricius 2000). As a result of recursive partitioning, the effect of an environmental
30 variable on biotic communities is assessed among sites which are the most homogenous as possible
31 considering other more important variables. If natural variables are included (such as morphometric
32 or geographic), a typology is implicitly made in the analysis. However, the community responses
33 (thresholds) in the implicit types are aggregated in the outputs of gradient forest, so that it is not
34 possible to analyze the interactions in details. Multivariate Regression Tree (De'Ath 2002) which also

1 deals with community change in a multivariate environment, is similar to gradient forest and more
2 transparent as regards interactions. Nonetheless, gradient forest adds the performance of an
3 ensemble method (random forest) and can indicate a rate of community change along
4 environmental gradients, which is useful to identify thresholds.

5 If a natural factor is suspected by its importance, to modulate a community response to a variable of
6 interest, the effect of the interaction may be investigated by analyzing separately different groups of
7 sites, defined by the thresholds associated with this variable. In this study, the reduction of the data
8 set according to longitude and latitude, assumed that there might be different thresholds in regions
9 with distinct phytoplankton flora. The difference may be due to regional extirpation or adaptation of
10 sensitive species (Utz et al. 2009). Even if gradient forest automatically accounts for regional effects
11 on the abundance of each species, focusing on a more floristically homogeneous region limits the
12 number of thresholds and facilitates the interpretation of the results. Lake typologies are commonly
13 used in the context of the European Water Framework Directive for the application of biological
14 indicators (European Commission 2000). The main natural criteria used to define lake types are
15 altitude, depth and alkalinity. In the GF analyses, these variables had not a prominent importance
16 compared to human-influenced variables. Thus, it did not appear necessary to carry out further
17 typological data splits which would reduce the number of observations and decrease the robustness
18 of the detected thresholds. Moreover, the biological phytoplankton index developed for French lakes
19 does not consider any lake typology (Ferret and Laplace-Treyture 2013). The thresholds found in this
20 study are relevant for the concerned biogeographic region (NW France) and should not be used for
21 the management of lakes in other regions without more investigations.

22 **Drivers of phytoplankton communities**

23 A striking result of the gradient forest analysis on phytoplankton species was the overwhelming
24 importance of the geographic variables (latitude and longitude) (Fig. 3). This can be explained by a
25 difference in spatial scales between species distribution areas and survey area. Generally, species
26 distribution areas are determined by a combination of migration processes and a selection by
27 environmental factors. As these distributions were not governed by a single dominant factor among
28 those that were used in the analysis, geographic coordinates were the best predictors of species
29 abundance. At the scale of the United States of America, Stomp et al. (2011) also demonstrated a
30 large influence of geographic variables on phytoplankton diversity.

31 After exclusion of the eastern and southern regions, the relative importance of geographic variables
32 was reduced in favor of structuring environmental variables, such as Secchi or TP (Fig. 5).

33 Nevertheless, latitude remained an important variable since a secondary threshold was also detected

1 around 47°N (Fig. 4). Most of the important variables identified by gradient forest in the second
2 analysis are related to lake eutrophication, as shown by the high correlation coefficients with the
3 mean annual log-transformed chlorophyll-*a* concentration: -0.77, 0.68 and 0.59 for Secchi, TP and
4 DOC respectively (Pearson correlation). A Secchi depth measures water transparency and determines
5 the amount of light available for microalgal growth. Considering phytoplankton as a whole, Secchi
6 does not constitute a limiting resource (as phosphorus does) but rather a consequence of algal
7 biomass which reduces the penetration of light in the water column. However, at the species level,
8 Secchi may select species according to their light requirements or their ability to adapt to low or high
9 irradiance levels. Its prominent role in algal physiology and interspecific competition results in a
10 leading position of Secchi among the other explanatory variables.

11 The importance of TP and dissolved inorganic nitrogen (DIN) is consistent with the well-known
12 trophic control on phytoplankton production. An increase in nutrient concentrations generally
13 stimulates phytoplankton biomass but its effect is not the same on all species or algal groups
14 (Watson et al. 1997). Some species show preferences for low nutrient concentrations whereas others
15 develop mostly in nutrient-rich waters. The classification of species according to their affinity for
16 nutrients has been the basis for the development of many biological indices of eutrophication
17 (Carvalho et al. 2013).

18 As lakes are generally limited by phosphorus, the importance of DIN might not be expected (Fig. 5).
19 However, eutrophicated lakes with a high phosphorus concentration can become limited by nitrogen
20 (Solheim et al. 2008; Donald et al. 2013). Changes in DIN concentration modify the N/P ratio and
21 influence the composition of phytoplankton communities, particularly the abundance of nitrogen-
22 fixing species (Schindler 1977). Several reasons may explain the higher importance of NH₄ over NO₃.
23 Most ammonia comes from the decomposition of organic matter. Its concentration is thus closely
24 related to lake productivity and eutrophication, as shown by the correlation between NH₄ and TP
25 ($\rho=0.66$), whereas NO₃ was independent of TP ($\rho=-0.01$). High concentration of ammonia can also be
26 related to organic pollution from sewage effluents (Beklioglu et al. 1999; García-Ferrer et al. 2003)
27 which may differentially impact phytoplankton species (Katsiapi et al. 2013; Villena and Romo 2003).
28 Finally, NH₄⁺ can be toxic to phytoplankton when it turns into its unionized form (NH₃) as pH and
29 temperature increase (Camargo and Alonso 2006; König et al. 1987). Unlike TP or NH₄, the effect of
30 nitrates was better expressed when the annual maximum was considered instead of the annual
31 median. In average, the maximum concentration of nitrates was measured at the year's first
32 campaign at the end of winter.

1 Another interesting result was the relatively high importance of DOC to predict phytoplankton
2 abundance. The role of DOC in lake ecosystem functioning has been neglected in the last decades,
3 whereas research in limnology has focused on the nutrient-productivity relationship (Williamson et
4 al. 1999; Carpenter et al. 1998). DOC can affect phytoplankton communities through the attenuation
5 of solar radiation and interactions with nutrients and contaminants (Jansson et al. 2000; Wall and
6 Briand 1979). High DOC concentrations might result either from input of colored organic matter from
7 the watershed, limiting light penetration and inhibiting primary production, or from the production
8 of aquatic plants in a eutrophicated system (transparent, labile DOC) (Bade et al. 2007). A better
9 assessment of DOC effect on phytoplankton would require a distinction between the allochthonous
10 and autochthonous forms.

11 Due to its interactions with pH and nutrients (CO₂ concentration and phosphate bioavailability) and
12 its link with water conductivity, alkalinity is a general factor influencing the distribution of microalgal
13 species. In another study at broad spatial scale, it was shown that alkalinity could explain a large part
14 of variations in phytoplankton communities (Maileht et al. 2013) but maximal depth was the most
15 determining factor. Here, the morphological variable Z_{max} was considerably less important than
16 Secchi and TP. This could be explained by the conditional approach used here to estimate variable
17 importance, since lake depth is correlated with factors which might influence more directly
18 phytoplankton communities, like TP or Secchi ($\rho < -0.7$). Indeed, when lake depth decreases,
19 resuspension of sediment is more likely and inputs from the watershed (phosphorus) become
20 concentrated in a lower lake water volume.

21 **Thresholds and their use for lake management**

22 The community thresholds that were identified in this study can be used to derive ecologically sound
23 environmental standards. Unlike 'chemical' thresholds like those derived from percentile analysis
24 (Chambers et al. 2011; Smith and Tran 2010), community thresholds take into account the response
25 of biological compartments to environmental stressors and are in line with the preservation of
26 biodiversity and ecosystem functioning (Brenden et al. 2008). The necessity to link environmental
27 standards with the response of communities has been strengthened by the European Water
28 Framework Directive (European Commission 2000), which states that physico-chemical quality
29 elements must support the achievement of good status for the biological quality elements, such as
30 phytoplankton, fish, macrophytes and macro-invertebrates. Recently, Poikane et al. (2014) have
31 proposed ecological boundaries for chl-*a* on the basis of the response of lake ecosystems to
32 eutrophication, considering phytoplankton (Cyanobacteria) and macrophytes.

1 The gradient forest analysis provided interesting thresholds for several environmental variables.
2 Geographic thresholds isolated the Southern and especially the Eastern part of the country, which is
3 characterized by several chains of mountains (mainly the Alps). These limits form a biogeographic
4 typology which can be used to address the issue of thresholds separately for each region. Focusing in
5 this study on the largest identified group without Mediterranean or alpine influence, interesting
6 values were detected by gradient forest for the human-influenced variables (Fig. 6, 7 and 8). For TP
7 and Secchi, the thresholds identified here can be compared to the trophic class limits proposed by
8 OECD (1982): $35 \mu\text{g TP.L}^{-1}$ and $100 \mu\text{g TP.L}^{-1}$ or 1 m Secchi depth for the mesotrophic/eutrophic and
9 eutrophic/hypereutrophic limits, respectively.

10 Ecological threshold values reported in the literature concern mostly total phosphorus and have the
11 same order of magnitude as the threshold at $50 \mu\text{g.L}^{-1}$ found in the present study. Concerning lakes,
12 TP thresholds have been reported at $18 \mu\text{g.L}^{-1}$ (Soranno et al. 2008), 10, 25 and $70 \mu\text{g.L}^{-1}$ (Free et al.
13 2006), 20 and $50 \mu\text{g.L}^{-1}$ (Penning et al. 2008) using phytoplankton or macrophyte data. Most other
14 references of phosphorus thresholds come from periphyton or macro-invertebrates in rivers. The
15 values are the following: 30-60 (Dodds et al. 2002), 12-15 (Richardson et al. 2007) , 60-90 (Evans-
16 White et al. 2009), 9-70 (Smith and Tran 2010), 30-280 (Black et al. 2011) and $21\text{-}63 \mu\text{g TP.L}^{-1}$
17 (Chambers et al. 2012b). Thresholds can be more or less broad in a gradient depending on their
18 sharpness, the variability in response among species or biological compartments and the uncertainty
19 associated with detection methods. A broad threshold, as observed here for TP, may also reflect
20 hysteresis in the response of a biological compartment to the gradient, especially if some lakes under
21 restoration are included in the analysis.

22 Since the variables presented in Fig. 6, 7 and 8 (Secchi, TP, NH_4 , NO_3max , DOC) may respond to
23 human disturbance, they can be used as physico-chemical indicators to evaluate the lake alteration
24 level and to protect ecosystems. All of them might be sensitive to eutrophication or organic
25 pollution, but also to hydromorphological alterations. Thresholds identified in this study can be
26 translated into environmental standards which should not be exceeded to prevent ecosystems from
27 important changes in their communities. However, variations in these physico-chemical parameters
28 among lakes may be partly natural and independent of anthropogenic pressures. For example, a lake
29 can present high values of colored DOC due to a large proportion of wetlands or coniferous forests in
30 its watershed, resulting in a low Secchi depth (humic lakes). Similarly, naturally eutrophic lakes may
31 exhibit high phosphorus concentrations which are not linked to any pollution (Borics et al. 2013),
32 even if most European lakes in Europe have lower TP concentrations than the threshold of $50 \mu\text{g.L}^{-1}$
33 (Cardoso et al. 2007). Therefore, reference values must be defined in the evaluation process in order
34 to take into account non-anthropogenic variations. Soranno et al. (2008) proposed a framework to

1 derive site-specific nutrient criteria from the knowledge of biological thresholds, reference values,
2 and measured values. The lowest biological threshold that is above the modeled reference value,
3 may be retained as an operational standard. When several thresholds are identified above a
4 reference value, they may account for the increasing levels of alteration, delimiting the different
5 classes of quality (Solheim et al. 2008; Free et al. 2006). According to the European Water Framework
6 Directive (European Commission 2000), there must be four boundaries delimiting five quality classes
7 for each physico-chemical quality element. The contribution of this study consists in the detection of
8 some boundaries in some gradients. The research of thresholds should be extended to other lake
9 communities (macrophytes, fish, macro-invertebrates, phytobenthos), which could reveal new
10 critical zones because of their different sensitivities to pollution or on the contrary, reinforce the
11 position of some thresholds affecting several biological groups (Richardson et al. 2007). Nevertheless,
12 one critical question remains for lake management: which limit between quality classes should be
13 attributed to a given ecological threshold? The main threshold following the reference conditions
14 may be the good-moderate boundary, but if this threshold appears too distant from the reference, it
15 may be assigned a lower quality limit. The choice can be guided by the levels of boundaries set in
16 other regions through other methods (Claussen et al. 2012) or by the observed ranges of parameters
17 in groups of lakes of each quality class, as determined by intercalibrated biological indices (Phillips et
18 al. 2013). Missing limits can eventually be derived from simple interval divisions.

19 **Conclusion**

20 Gradient forest is an exploratory method which deals with multivariate biological and environmental
21 data. It is well suited for the research of community thresholds on a large spatial scale. The main
22 drivers of phytoplankton variations among the French lakes were identified and community
23 thresholds were detected for spatial variables (longitude and latitude) and human-influenced
24 variables in a chosen biogeographic region. These thresholds can contribute to the definition of
25 environmental standards for lake management. Further investigations may include the definition of
26 appropriate thresholds for the other biogeographic regions derived from the first analysis and the
27 questions of criteria for other pressures on lakes such as acidification or climate change (Cardoso et
28 al. 2009). Thresholds obtained from other lake communities (e.g. macrophytes or fish) could also be
29 compared to those inferred here from phytoplankton. Finally, models that include
30 hydromorphological explanatory variables should be developed to (1) identify the main causes of
31 variations among lakes of the most important physico-chemical variables identified in this study, and
32 (2) to predict natural reference conditions which are necessary for the evaluation of human
33 alteration of ecosystems.

1 **Acknowledgments**

2 This research was funded by The French National Agency for Water and Aquatic Environments
3 (ONEMA). The authors are grateful to the French Water Basin Agencies and their partners who
4 contributed to the lake data acquisition and to the Onema-Irstea consortium for Lake Hydroecology
5 (Pôle Onema-Irstea d'études et de recherche "hydroécologie des plans d'eau", Aix-en-Provence,
6 France), who currently maintains the national biological and physico-chemical database for French
7 lakes; special thanks to Nathalie Reynaud and Thierry Point (database) and Milena Borissova (text
8 edition).

9 **References**

- 10 **Andersen, T., Carstensen, J., Hernández-García, E., & Duarte, C. M. (2009). Ecological thresholds**
11 **and regime shifts: approaches to identification. *Trends in Ecology & Evolution*, 24(1), 49-57,**
12 **doi:<http://dx.doi.org/10.1016/j.tree.2008.07.014>.**
- 13 **Bade, D. L., Carpenter, S. R., Cole, J. J., Pace, M. L., Kritzberg, E., Van de Bogert, M. C., et al. (2007).**
14 **Sources and fates of dissolved organic carbon in lakes as determined by whole-lake carbon**
15 **isotope additions. *Biogeochemistry*, 84(2), 115-129, doi:DOI 10.1007/s10533-006-9013-y.**
- 16 **Baker, M. E., & King, R. S. (2010). A new method for detecting and interpreting biodiversity and**
17 **ecological community thresholds. *Methods in Ecology and Evolution*, 1(1), 25-37,**
18 **doi:10.1111/j.2041-210X.2009.00007.x.**
- 19 **Beklioglu, M., Carvalho, L., & Moss, B. (1999). Rapid recovery of a shallow hypertrophic lake**
20 **following sewage effluent diversion: Lack of chemical resilience. *Hydrobiologia*, 412, 5-15.**
- 21 **Black, R. W., Moran, P. W., & Frankforter, J. D. (2011). Response of algal metrics to nutrients and**
22 **physical factors and identification of nutrient thresholds in agricultural streams.**
23 ***Environmental Monitoring and Assessment*, 175(1-4), 397-417, doi:10.1007/s10661-010-**
24 **1539-8.**
- 25 **Borics, G., Nagy, L., Miron, S., Grigorszky, I., László-Nagy, Z., Lukács, B. A., et al. (2013). Which**
26 **factors affect phytoplankton biomass in shallow eutrophic lakes? *Hydrobiologia*, 714(1),**
27 **93-104.**
- 28 **Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5-32, doi:10.1023/A:1010933404324.**
- 29 **Breiman, L., Friedman, J., Stone, C. J., & Olshen, R. A. (1984). *Classification and Regression Trees:***
30 **Chapman and Hall/CRC**
- 31 **Brenden, T. O., Wang, L., & Su, Z. (2008). Quantitative Identification of Disturbance Thresholds in**
32 **Support of Aquatic Resource Management. *Environmental Management*, 42(5), 821-832,**
33 **doi:10.1007/s00267-008-9150-2.**
- 34 **Camargo, J. A., & Alonso, A. (2006). Ecological and toxicological effects of inorganic nitrogen**
35 **pollution in aquatic ecosystems: A global assessment. *Environment International*, 32(6),**
36 **831-849, doi:DOI 10.1016/j.envint.2006.05.002.**
- 37 **Cardoso, A. C., Free, G., Nöges, P., Kaste, Ø., Poikane, S., & Solheim, A. L. (2009). Lake**
38 **Management, Criteria. In G. E. Likens (Ed.), *Encyclopedia of Inland Waters* (Vol. 1, pp. 310-**
39 **331). Oxford.**
- 40 **Cardoso, A. C., Solimini, A., Premazzi, G., Carvalho, L., Lyche, A., & Rekolainen, S. (2007).**
41 **Phosphorus reference concentrations in European lakes. *Hydrobiologia*, 584, 3-12,**
42 **doi:10.1007/s10750-007-0584-y.**
- 43 **Carpenter, S. R., Cole, J. J., Kitchell, J. F., & Pace, M. L. (1998). Impact of dissolved organic carbon,**
44 **phosphorus, and grazing on phytoplankton biomass and production in experimental lakes.**
45 ***Limnology and Oceanography*, 43(1), 73-80.**

- 1 **Carvalho, L., Poikane, S., Lyche Solheim, A., Phillips, G., Borics, G., Catalan, J., et al. (2013). Strength**
2 **and uncertainty of phytoplankton metrics for assessing eutrophication impacts in lakes.**
3 ***Hydrobiologia*, 704(1), 127-140, doi:10.1007/s10750-012-1344-1.**
- 4 **Catalan, J., Barbieri, M. G., Bartumeus, F., Bitusik, P., Botev, I., Brancelj, A., et al. (2009). Ecological**
5 **thresholds in European alpine lakes. *Freshwater Biology*, 54(12), 2494-2517,**
6 **doi:10.1111/j.1365-2427.2009.02286.x.**
- 7 **Chambers, P. A., Benoy, G. A., Brua, R. B., & Culp, J. M. (2011). Application of nitrogen and**
8 **phosphorus criteria for streams in agricultural landscapes. *Water Science and Technology*,**
9 **64(11), 2185-2191, doi:Doi 10.2166/Wst.2011.760.**
- 10 **Chambers, P. A., Culp, J. M., Roberts, E. S., & Bowerman, M. (2012a). Development of**
11 **Environmental Thresholds for Streams in Agricultural Watersheds. *Journal of***
12 ***Environmental Quality*, 41(1), 1-6, doi:10.2134/jeq2011.0338.**
- 13 **Chambers, P. A., McGoldrick, D. J., Brua, R. B., Vis, C., Culp, J. M., & Benoy, G. A. (2012b).**
14 **Development of Environmental Thresholds for Nitrogen and Phosphorus in Streams.**
15 ***Journal of Environmental Quality*, 41(1), 7-20, doi:10.2134/jeq2010.0273.**
- 16 **Claussen, U., Müller, P., & Arle, J. (2012). Comparison of Environmental Quality Objectives,**
17 **Threshold Values or Water Quality Targets Set for the Demands of the European Water**
18 **Framework Directive. *WFD CIS ECOSTAT WG A Report* (pp. 27). CIRCABC**
- 19 **Cuffney, T. F., & Qian, S. S. (2013). A critique of the use of indicator-species scores for identifying**
20 **thresholds in species responses. *Freshwater Science*, 32(2), 471-488, doi:10.1899/12-056.1.**
- 21 **Daily, J. P., Hitt, N. P., Smith, D. R., & Snyder, C. D. (2012). Experimental and environmental factors**
22 **affect spurious detection of ecological thresholds. *Ecology*, 93(1), 17-23.**
- 23 **De'Ath, G. (2002). Multivariate regression trees: a new technique for modeling species-**
24 **environment relationships. *Ecology*, 83(4), 1105-1117, doi:Doi 10.2307/3071917.**
- 25 **De'ath, G., & Fabricius, K. E. (2000). Classification and regression trees: A powerful yet simple**
26 **technique for ecological data analysis. *Ecology*, 81(11), 3178-3192, doi:Doi**
27 **10.2307/177409.**
- 28 **Dodds, W. K., Clements, W. H., Gido, K., Hilderbrand, R. H., & King, R. S. (2010). Thresholds,**
29 **breakpoints, and nonlinearity in freshwaters as related to management. *Journal of the***
30 ***North American Benthological Society*, 29(3), 988-997, doi:Doi 10.1899/09-148.1.**
- 31 **Dodds, W. K., Smith, V. H., & Lohman, K. (2002). Nitrogen and phosphorus relationships to benthic**
32 **algal biomass in temperate streams. *Canadian Journal of Fisheries and Aquatic Sciences*,**
33 **59(5), 865-874, doi:Doi 10.1139/F02-063.**
- 34 **Donald, D. B., Bogard, M. J., Finlay, K., Bunting, L., & Leavitt, P. R. (2013). Phytoplankton-Specific**
35 **Response to Enrichment of Phosphorus-Rich Surface Waters with Ammonium, Nitrate, and**
36 **Urea. *Plos One*, 8(1), doi:DOI 10.1371/journal.pone.0053277.**
- 37 **Ellis, N., Smith, S. J., & Pitcher, C. R. (2012). Gradient forests: calculating importance gradients on**
38 **physical predictors. *Ecology*, 93(1), 156-168.**
- 39 **European Commission (2000). Directive 2000/60/EC of the European Parliament and of the Council**
40 **of 23 October 2000 establishing a Framework for Community Action in the Field of Water**
41 **Policy. (pp. 72): The European Parliament and Council.**
- 42 **Evans-White, M. A., Dodds, W. K., Huggins, D. G., & Baker, D. S. (2009). Thresholds in**
43 **macroinvertebrate biodiversity and stoichiometry across water-quality gradients in Central**
44 **Plains (USA) streams. *Journal of the North American Benthological Society*, 28(4), 855-868,**
45 **doi:10.1899/08-113.1.**
- 46 **Feret, T., & Laplace-Treyture, C. (2013). IPLAC : l'indice Phytoplancton Lacustre : Méthode de**
47 **développement, description et application nationale 2012. *Rapport convention***
48 ***Onema/Irstea 2012* (pp. 69). Bordeaux: Irstea, UR REBX.**
- 49 **Free, G., Little, R., Tierney, D., Donnelly, K., & Caroni, R. (2006). A reference based typology and**
50 **ecological assessment system for Irish lakes. Preliminary investigations. (pp. 266). Wexford,**
51 **Ireland: Environmental Protection Agency.**

- 1 **García-Ferrer, I., Camacho, A., Armengol, X., Miracle, M. R., & Vicente, E. (2003). Seasonal and**
 2 **spatial heterogeneity in the water chemistry of two sewage-affected saline shallow lakes**
 3 **from central Spain. *Hydrobiologia*, 506-509, 101-110.**
- 4 **Graham, J. L., Jones, J. R., Jones, S. B., Downing, J. A., & Clevenger, T. E. (2004). Environmental**
 5 **factors influencing microcystin distribution and concentration in the Midwestern United**
 6 **States. *Water Research*, 38(20), 4395-4404, doi:DOI 10.1016/j.watres.2004.08.004.**
- 7 **Groffman, P., Baron, J., Blett, T., Gold, A., Goodman, I., Gunderson, L., et al. (2006). Ecological**
 8 **thresholds: The key to successful environmental management or an important concept**
 9 **with no practical application? *Ecosystems*, 9(1), 1-13, doi:DOI 10.1007/s10021-003-0142-z.**
- 10 **Holt, C. A., Yan, N. D., & Somers, K. M. (2003). pH 6 as the threshold to use in critical load modeling**
 11 **for zooplankton community change with acidification in lakes of south-central Ontario:**
 12 **accounting for morphometry and geography. *Canadian Journal of Fisheries and Aquatic***
 13 ***Sciences*, 60(2), 151-158, doi:doi:10.1139/f03-008.**
- 14 **Jansson, M., Bergström, A.-K., Blomqvist, P., & Drakare, S. (2000). allochthonous organic carbon**
 15 **and phytoplankton/bacterioplankton production relationships in lakes. *Ecology*, 81(11),**
 16 **3250-3255, doi:10.1890/0012-9658(2000)081[3250:AOCAPB]2.0.CO;2.**
- 17 **Katsiapi, M., Moustaka-Gouni, M., Vardaka, E., & Kormas, K. A. (2013). Different phytoplankton**
 18 **descriptors show asynchronous changes in a shallow urban lake (L. Kastoria, Greece) after**
 19 **sewage diversion. *Fundamental and Applied Limnology*, 182(3), 219-230.**
- 20 **King, R. S., & Baker, M. E. (2010). Considerations for analyzing ecological community thresholds in**
 21 **response to anthropogenic environmental gradients. *Journal of the North American***
 22 ***Benthological Society*, 29(3), 998-1008, doi:10.1899/09-144.1.**
- 23 **King, R. S., Baker, M. E., Kazyak, P. F., & Weller, D. E. (2011). How novel is too novel? Stream**
 24 **community thresholds at exceptionally low levels of catchment urbanization. *Ecological***
 25 ***Applications*, 21(5), 1659-1678, doi:10.1890/10-1357.1.**
- 26 **Konig, A., Pearson, H. W., & Silva, S. A. (1987). Ammonia Toxicity to Algal Growth in Waste**
 27 **Stabilization Ponds. *Water Science and Technology*, 19(12), 115-122.**
- 28 **Laplace-Treytore, C., Barbe, J., Dutartre, A., Druart, J.-C., Rimet, F., & Anneville, O. (2009). Standard**
 29 **protocol for sampling, conservation, observation and counting of lake phytoplankton for**
 30 **application of the WFD. Version 3.3.1. (pp. 42). Cestas, France: Cemagref UR REBX.**
- 31 **Luck, G. W. (2005). An introduction to ecological thresholds. *Biological Conservation*, 124(3), 299-**
 32 **300, doi:10.1016/j.biocon.2005.01.042.**
- 33 **Maileht, K., Nöges, T., Nöges, P., Ott, I., Mischke, U., Carvalho, L., et al. (2013). Water colour,**
 34 **phosphorus and alkalinity are the major determinants of the dominant phytoplankton**
 35 **species in European lakes. *Hydrobiologia*, 704(1), 115-126, doi:10.1007/s10750-012-1348-x.**
- 36 **May, R. M. (1977). Thresholds and Breakpoints in Ecosystems with a Multiplicity of Stable States.**
 37 ***Nature*, 269(5628), 471-477, doi:Doi 10.1038/269471a0.**
- 38 **OECD (1982). Eutrophication of waters -monitoring, assessment and control. Paris: Organisation**
 39 **for Economic Co-operation and Development.**
- 40 **Penning, W. E., Dudley, B., Mjelde, M., Hellsten, S., Hanganu, J., Kolada, A., et al. (2008). Using**
 41 **aquatic macrophyte community indices to define the ecological status of European lakes.**
 42 ***Aquatic Ecology*, 42(2), 253-264, doi:10.1007/s10452-008-9183-x.**
- 43 **Phillips, G., Lyche-Solheim, A., Skjelbred, B., Mischke, U., Drakare, S., Free, G., et al. (2013). A**
 44 **phytoplankton trophic index to assess the status of lakes for the Water Framework**
 45 **Directive. *Hydrobiologia*, 704(1), 75-95, doi:10.1007/s10750-012-1390-8.**
- 46 **Pitcher, C. R., Lawton, P., Ellis, N., Smith, S. J., Incze, L. S., Wei, C. L., et al. (2012). Exploring the role**
 47 **of environmental variables in shaping patterns of seabed biodiversity composition in**
 48 **regional-scale ecosystems. *Journal of Applied Ecology*, 49(3), 670-679, doi:DOI**
 49 **10.1111/j.1365-2664.2012.02148.x.**
- 50 **Poikane, S., Portielje, R., van den Berg, M., Phillips, G., Brucet, S., Carvalho, L., et al. (2014).**
 51 **Defining ecologically relevant water quality targets for lakes in Europe. *Journal of Applied***
 52 ***Ecology*, 51(3), 592-602, doi:10.1111/1365-2664.12228.**

- 1 Prasad, A. M., Iverson, L. R., & Liaw, A. (2006). Newer classification and regression tree techniques:
 2 Bagging and random forests for ecological prediction. *Ecosystems*, 9(2), 181-199, doi:DOI
 3 10.1007/s10021-005-0054-1.
- 4 Qian, S. S. (2014). Ecological threshold and environmental management: A note on statistical
 5 methods for detecting thresholds. *Ecological Indicators*, 38, 192-197,
 6 doi:10.1016/j.ecolind.2013.11.008.
- 7 Qian, S. S., King, R. S., & Richardson, C. J. (2003). Two statistical methods for the detection of
 8 environmental thresholds. *Ecological Modelling*, 166(1–2), 87-97,
 9 doi:http://dx.doi.org/10.1016/S0304-3800(03)00097-8.
- 10 R_Core_Team (2013, R: A language and environment for statistical computing.
- 11 Richardson, C. J., King, R. S., Qian, S. S., Vaithyanathan, P., Qualls, R. G., & Stow, C. A. (2007).
 12 Estimating ecological thresholds for phosphorus in the Everglades. *Environmental Science
 13 & Technology*, 41(23), 8084-8091, doi:10.1021/es062624w.
- 14 Scheffer, M., & Carpenter, S. R. (2003). Catastrophic regime shifts in ecosystems: linking theory to
 15 observation. *Trends in Ecology & Evolution*, 18(12), 648-656,
 16 doi:10.1016/j.tree.2003.09.002.
- 17 Scheffer, M., Rinaldi, S., Gragnani, A., Mur, L. R., & vanNes, E. H. (1997). On the dominance of
 18 filamentous cyanobacteria in shallow, turbid lakes. *Ecology*, 78(1), 272-282.
- 19 Scheffer, M., Szabo, S., Gragnani, A., van Nes, E. H., Rinaldi, S., Kautsky, N., et al. (2003). Floating
 20 plant dominance as a stable state. *Proceedings of the National Academy of Sciences of the
 21 United States of America*, 100(7), 4040-4045, doi:10.1073/pnas.0737918100.
- 22 Schindler, D. W. (1977). Evolution of Phosphorus Limitation in Lakes. *Science*, 195(4275), 260-262,
 23 doi:DOI 10.1126/science.195.4275.260.
- 24 Smith, A. J., & Tran, C. P. (2010). A weight-of-evidence approach to define nutrient criteria
 25 protective of aquatic life in large rivers. *Journal of the North American Benthological
 26 Society*, 29(3), 875-891, doi:Doi 10.1899/09-076.1.
- 27 Solheim, A. L., Rekolainen, S., Moe, S. J., Carvalho, L., Phillips, G., Ptacnik, R., et al. (2008).
 28 Ecological threshold responses in European lakes and their applicability for the Water
 29 Framework Directive (WFD) implementation: synthesis of lakes results from the REBECCA
 30 project. *Aquatic Ecology*, 42(2), 317-334, doi:DOI 10.1007/s10452-008-9188-5.
- 31 Soranno, P. A., Cheruvilil, K. S., Stevenson, R. J., Rollins, S. L., Holden, S. W., Heaton, S., et al.
 32 (2008). A framework for developing ecosystem-specific nutrient criteria: Integrating
 33 biological thresholds with predictive modeling. *Limnology and Oceanography*, 53(2), 773-
 34 787, doi:10.4319/lo.2008.53.2.0773.
- 35 Stomp, M., Huisman, J., Mittelbach, G. G., Litchman, E., & Klausmeier, C. A. (2011). Large-scale
 36 biodiversity patterns in freshwater phytoplankton. *Ecology*, 92(11), 2096-2107.
- 37 Strobl, C., Boulesteix, A. L., Kneib, T., Augustin, T., & Zeileis, A. (2008). Conditional variable
 38 importance for random forests. *Bmc Bioinformatics*, 9, doi:Doi 10.1186/1471-2105-9-307.
- 39 Utz, R. M., Hilderbrand, R. H., & Boward, D. M. (2009). Identifying regional differences in threshold
 40 responses of aquatic invertebrates to land cover gradients. *Ecological Indicators*, 9(3), 556-
 41 567, doi:10.1016/j.ecolind.2008.08.008.
- 42 Villena, M. J., & Romo, S. (2003). Phytoplankton changes in a shallow Mediterranean lake
 43 (Albufera of Valencia, Spain) after sewage diversion. *Hydrobiologia*, 506-509, 281-287.
- 44 Vollenweider, R. A. (1975). Input-output models - With special reference to the phosphorus loading
 45 concept in limnology. *Source of the Document Schweizerische Zeitschrift für Hydrologie*,
 46 37(1), 53-84.
- 47 Wall, D., & Briand, F. (1979). Response of lake phytoplankton communities to in situ manipulations
 48 of light intensity and colour. *Journal of Plankton Research*, 1(1), 103-112,
 49 doi:10.1093/plankt/1.1.103.
- 50 Watson, S. B., McCauley, E., & Downing, J. A. (1997). Patterns in phytoplankton taxonomic
 51 composition across temperate lakes of differing nutrient status. *Limnology and
 52 Oceanography*, 42(3), 487-495.

1 **Williamson, C. E., Morris, D. P., Pace, M. L., & Olson, A. G. (1999). Dissolved organic carbon and**
2 **nutrients as regulators of lake ecosystems: Resurrection of a more integrated paradigm.**
3 ***Limnology and Oceanography*, 44(3), 795-803.**

4

5

1 **Table 1** List of environmental variables used in the analysis with basic distribution statistics
 2 (minimum, median and maximum) among lakes. The codes are used to denominate the variables in
 3 the following. For variables having annual variations (lines below Zmax), the first code corresponds to
 4 the annual median and the codes with 'min' and 'max' endings are the minimal and maximal annual
 5 values, respectively. CV is the mean annual coefficient of variation.

Parameter	code	min	med	max	CV (%)
Latitude (°N)	lat	41.47	46.42	50.87	-
Longitude (°E)	lon	-4.01	2.26	9.48	-
Altitude (m)	Alt	0	213	2082	-
Maximal depth (m)	Zmax	0.8	15	309.7	-
Water temperature (°C)	Temp	7.1	17	22.8	40
	Tempmin	0.3	7	14.5	-
	Tempmax	8.9	21.5	29.2	-
Secchi depth (m)	Secchi	0.1	1.6	18.7	38
	Secchimmin	0.05	0.95	12	-
O ₂ saturation (%)	SatO2	53.9	94.8	126.9	17
	SatO2min	15.2	78.3	114.3	-
Alkalinity (meq.L ⁻¹)	Alk	0.1	1.0	4.6	-
pH	pH	5.7	8.1	9.9	6
	pHmin	5.1	7.5	9.1	-
	pHmax	6.9	8.6	12.5	-
Conductivity (µS.cm ⁻¹)	Cond	14	239	1307	12
Dissolved organic carbon (mg.L ⁻¹)	DOC	0.3	4.9	36.7	19
	DOCmax	0.4	6.0	56	-
Nitrates (µg.L ⁻¹)	NO3	120	1220	39260	81
	NO3max	250	3950	72000	-
Ammonia (µg.L ⁻¹)	NH4	9	50	389	71
	NH4max	10	130	3300	-
Nitrites (µg.L ⁻¹)	NO2	3.5	18	229	67
	NO2max	6	50	690	-
Phosphates (µg.L ⁻¹)	PO4	6	16	972	72
	PO4max	6	40	2116	-
Total phosphorus (µg.L ⁻¹)	TP	6	30	599	49
	TPmax	7	46	1490	-
Dissolved silicon (mg SiO ₂ .L ⁻¹)	SiO2	0.18	3.3	20.78	54

6

7

1 **Table 2** List of taxa used in the gradient forest analysis restricted to the North-West part of France (n =number
2 of records).

Taxon	n	Taxon	n	Taxon	n
Achnanthydium minutissimum	8	Dinobryon sociale	17	Peridinium inconspicuum	8
Acanthoceras zachariasii	23	Dinobryon suecicum	16	Phacus tortus	8
Actinastrum hantzschii	28	Discostella pseudostelligera	26	Phacotus lenticularis	24
Acutodesmus acuminatus	16	Discostella stelligera	13	Planktothrix agardhii	73
Acutodesmus obliquus	55	Dolichospermum flos-aquae	8	Plagioselmis nannoplanctica	83
Ankistrodesmus falcatus	10	Elakatothrix gelatinosa	37	Planctonema lauterbornii	11
Ankyra ancora	16	Erkenia subaequiciliata	34	Pseudodidymocystis fina	33
Ankyra judayi	47	Fragilaria capucina	10	Pseudodidymocystis planctonica	7
Aphanocapsa delicatissima	15	Fragilaria crotonensis	43	Pseudanabaena catenata	9
Aphanocapsa elachista	17	Goniochloris mutica	14	Pseudanabaena limnetica	16
Aphanocapsa holsatica	7	Granulocystopsis coronata	8	Pseudanabaena mucicola	13
Aphanizomenon flos-aquae	15	Hariotina reticulata	18	Puncticulata radiosia	32
Aphanothece clathrata	7	Kephyrion rubri-claustri	7	Rhodomonas lacustris	65
Asterionella formosa	97	Kirchneriella obesa	7	Scenedesmus aculeolatus	15
Aulacoseira ambigua	19	Koliella longiseta	64	Scenedesmus bicaudatus	32
Aulacoseira distans	44	Koliella planctonica	21	Scenedesmus ecornis	37
Aulacoseira granulata	98	Lagerheimia ciliata	19	Scenedesmus ellipticus	31
Aulacoseira subarctica	17	Lagerheimia genevensis	59	Scenedesmus obtusus	14
Choricystis minor	12	Lanceola spatulifera	71	Scenedesmus pulloideus	13
Ceratium hirundinella	40	Limnothrix redekei	31	Scenedesmus quadricauda var. longispinus	17
Chroococcus minutus	11	Mallomonas akrokomos	30	Scenedesmus verrucosus	12
Chrysococcus rufescens	9	Melosira varians	24	Schroederia setigera	11
Closterium aciculare	11	Merismopedia tenuissima	31	Sphaerocystis schroeteri	31
Closterium acutum	16	Microcystis aeruginosa	17	Staurastrum cingulum	8
Closterium acutum var. variabile	18	Micractinium pusillum	10	Staurastrum pingue	17
Cocconeis placentula	22	Monoraphidium arcuatum	95	Stephanodiscus hantzschii	66
Coelastrum astroideum	10	Monoraphidium circinale	41	Stephanodiscus minutulus	7
Coelastrum microporum	31	Monoraphidium contortum	111	Stephanodiscus parvus	29
Crucigenia tetrapedia	93	Monoraphidium convolutum	10	Staurosira construens	11
Cryptomonas curvata	11	Monoraphidium griffithii	55	Stelaxomonas dichotoma	45
Cryptomonas marssonii	27	Monoraphidium komarkovae	36	Spermatozopsis exsultans	15
Cryptomonas ovata	22	Monoraphidium minutum	50	Tabellaria flocculosa	12
Cyclotella meneghiniana	44	Monoraphidium nanum	8	Tetrachlorella alternans	30
Cyclotella ocellata	16	Monoraphidium tortile	37	Tetraedron caudatum	66
Cyclostephanos dubius	77	Monactinus simplex	50	Tetraedron incus	27
Cyclostephanos invisitatus	16	Mucidosphaerium pulchellum	20	Tetraedron minimum	104
Chrysolykos planctonicus	8	Navicula gregaria	13	Tetraedron triangulare	39
Desmodesmus abundans	21	Navicula lanceolata	18	Tetraëdriella regularis	19
Desmodesmus armatus	64	Nephrochlamys rostrata	13	Tetrastrum staurogeniaeforme	21
Desmodesmus communis	82	Neodesmus danubialis	11	Tetrastrum triangulare	7
Desmodesmus intermedius	45	Nitzschia acicularis	92	Trachelomonas volvocina	80
Desmodesmus opoliensis	61	Nitzschia palea	11	Trachelomonas volvocina var. punctata	19
Desmodesmus spinosus	21	Oocystis lacustris	17	Treubaria planctonica	33
Diatoma tenuis	17	Oocystis parva	11	Treubaria setigera	10
Dictyosphaerium subsolitarium	25	Pandorina morum	11	Ulnaria delicatissima var. angustissima	25
Dinobryon bavaricum	20	Pediastrum boryanum	52	Ulnaria ulna var. acus	42
Dinobryon crenulatum	11	Pediastrum duplex	83	Ulnaria ulna	70
Dinobryon divergens	40	Pediastrum boryanum var. longicorne	10	Urosolenia longiseta	77
Dinobryon sertularia	10	Pediastrum tetras	73	Woronichinia naegeliana	38

3

4

1 **Figure captions**

2

3 **Fig. 1** Location of the 224 lakes composing the data set for analysis. The dashed lines indicate latitude
4 +44.5°N and longitude +4.5°E.

5

6 **Fig. 2** Gradient forest analysis applied to a simulated data set of 4 species' abundances and 5
7 environmental variables. (1) Specific simulated relationships for 5 species-variable pairs (following
8 simple functions: logistic, Gaussian, step-wise and linear, with Poisson error), independence was
9 assumed for the other pairs. (2) Example of a regression tree for one species; splitting factors and
10 values are indicated at the nodes of the branches. Leaves give the mean values of species abundance
11 in the corresponding groups. (3) For each species, n trees are constructed, making a random forest
12 from which a species R^2 and a specific variable importance can be derived. (4) Gradient forest
13 combines the results of the random forests to give a community level variable importance and to
14 locate the splits and their importance in the environmental gradients (black curve for density of
15 splits, dashed curve for density of data and thick grey curve for the ratio of densities; the horizontal
16 dashed line indicates where the ratio is 1, each curve integrates to the variable importance). There is
17 one community threshold in the gradient of A and two in the gradient of B. The variable E shows no
18 threshold (independence). There are 2 equally important thresholds in the gradient of D after
19 accounting for data distribution (grey curve), and variable C exhibits a spurious threshold (see
20 discussion).

21

22 **Fig. 3** Species R^2 derived from random forests (a) and overall variable importance as determined by
23 gradient forest (b). The sum of variable importances equals the average specific R^2 . Codes refer to
24 Table 1.

25

26 **Fig. 4** Density of splits (top) and specific cumulative importance (bottom) in the gradients of the most
27 important variables identified by gradient forest : latitude (lat) and longitude (lon). The specific
28 cumulative importance is the integration of the ratio of densities for each species.

29

30 **Fig. 5** Species R^2 derived from random forests with the reduced data set (a) and overall variable
31 importance as determined by gradient forest with biovolume data (b). The sum of variable
32 importances equals the average specific R^2 . Codes refer to Table 1.

33

1 **Fig. 6** Results of gradient forest for the variables Secchi and TP (one column per variable). (a) Graph
2 of densities: the dashed curve indicates the density of data, the thin black curve, the density of splits,
3 the thick grey curve, the ratio of densities and the horizontal dashed line shows where the ratio is
4 one. Each curve integrates to the variable importance. (b) Curves of the specific cumulative
5 importance which is the integration of the ratio of densities for each species. The black curve
6 corresponds to the species shown as an example below. (c) Plots of the mean biovolume of
7 *Monoraphidium arcuatum* (Chlorophyceae) and *Aulacoseira granulata* (Bacillariophyceae) in the
8 gradients of Secchi and TP, respectively. The smoothing curve is continuous and the dashed curve
9 indicates the probability of presence of the species among the sampled sites. (d) Barplot of the
10 relative impurity importance of the variables for the species shown for example.

11

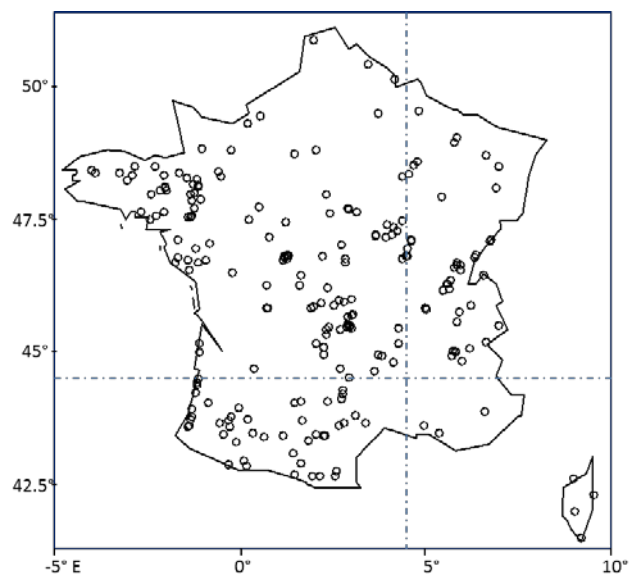
12 **Fig. 7** Results of gradient forest for the variables NH₄ and NO₃max (one column per variable). (a)
13 Graph of densities: the dashed curve indicates the density of data, the thin black curve, the density of
14 splits, the thick grey curve, the ratio of densities and the horizontal dashed line shows where the
15 ratio is one. Each curve integrates to the variable importance. (b) Curves of the specific cumulative
16 importance which is the integration of the ratio of densities for each species. The black curve
17 corresponds to the species shown as an example below. (c) Plots of the mean biovolume of
18 *Monoraphidium tortile* (Chlorophyceae) and *Cyclotella choctawhatcheeana* (Bacillariophyceae) in the
19 gradients of NH₄ and NO₃max, respectively. The smoothing curve is continuous and the dashed
20 curve indicates the probability of presence of the species among the sampled sites. (d) Barplot of the
21 relative impurity importance of the variables for the species shown for example.

22

23 **Fig. 8** Results of gradient forest for the variable DOC. (a) Graph of densities: the dashed curve
24 indicates the density of data, the thin black curve, the density of splits, the thick grey curve, the ratio
25 of densities and the horizontal dashed line shows where the ratio is one. Each curve integrates to the
26 variable importance. (b) Curves of the specific cumulative importance which is the integration of the
27 ratio of densities for each species. The black curve corresponds to the species shown as an example
28 below. (c) Plots of the mean biovolume of *Scenedesmus ecornis* (Chlorophyceae) in the gradients of
29 DOC. The smoothing curve is continuous and the dashed curve indicates the probability of presence
30 of the species among the sampled sites. (d) Barplot of the relative impurity importance of the
31 variables for the species shown for example.

32

1 Fig 1

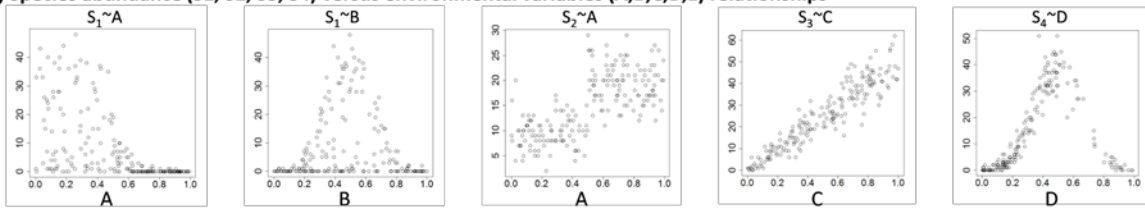


2

3

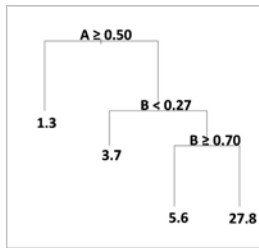
1 Fig 2

1) Species abundance (S1, S2, S3, S4) versus environmental variables (A,B,C,D,E) relationships



2) Regression tree per species

➤ Ex : $S_1 = f(A,B,C,D,E)$



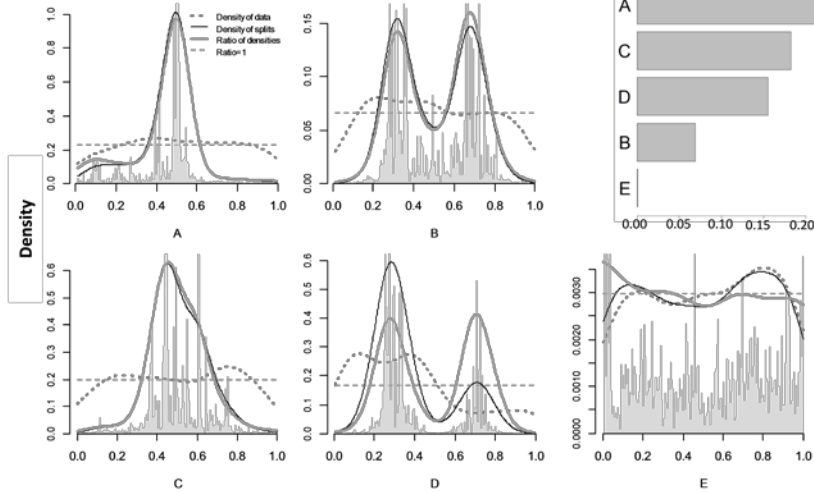
3) $\times n$ trees => random forest

- Cross-validated species predictive performance (R^2)
- Variable importance

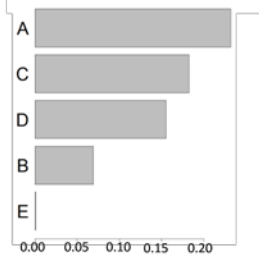
4) $\times p$ species => gradient forest

➤ Aggregation of RF results over all p species

4-a) Density of splits along gradients (from $n \times p$ trees)



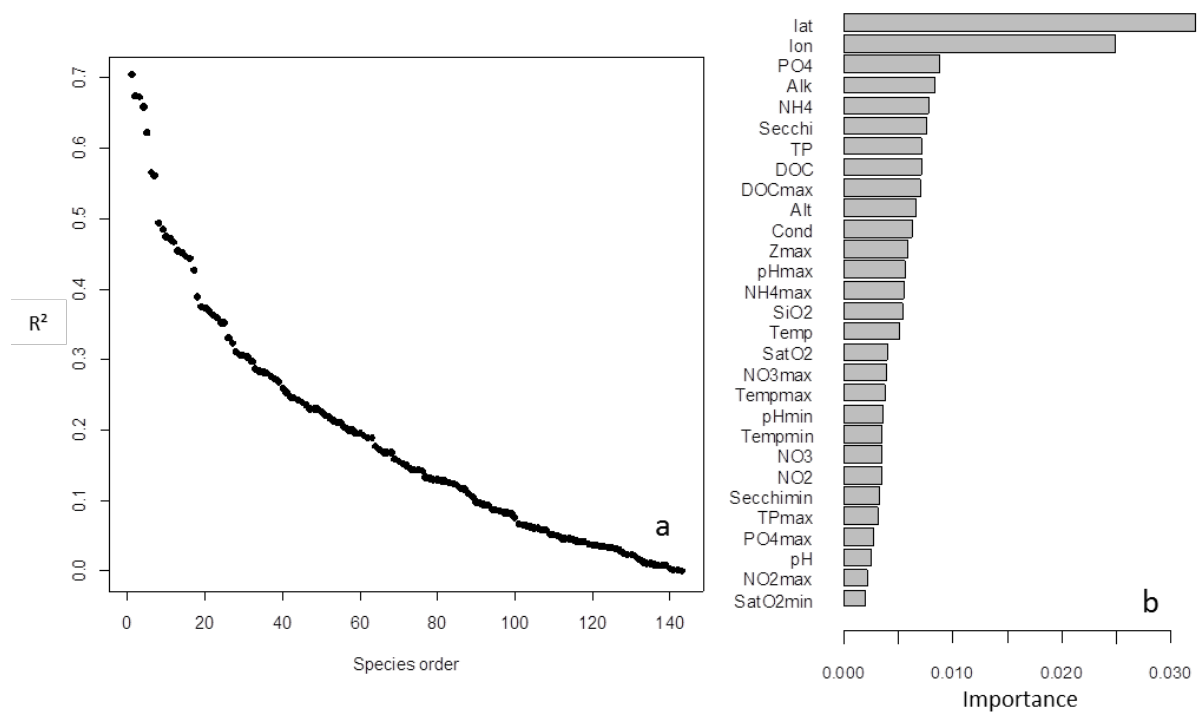
4-b) Overall variable importance (R^2 weighted)



2

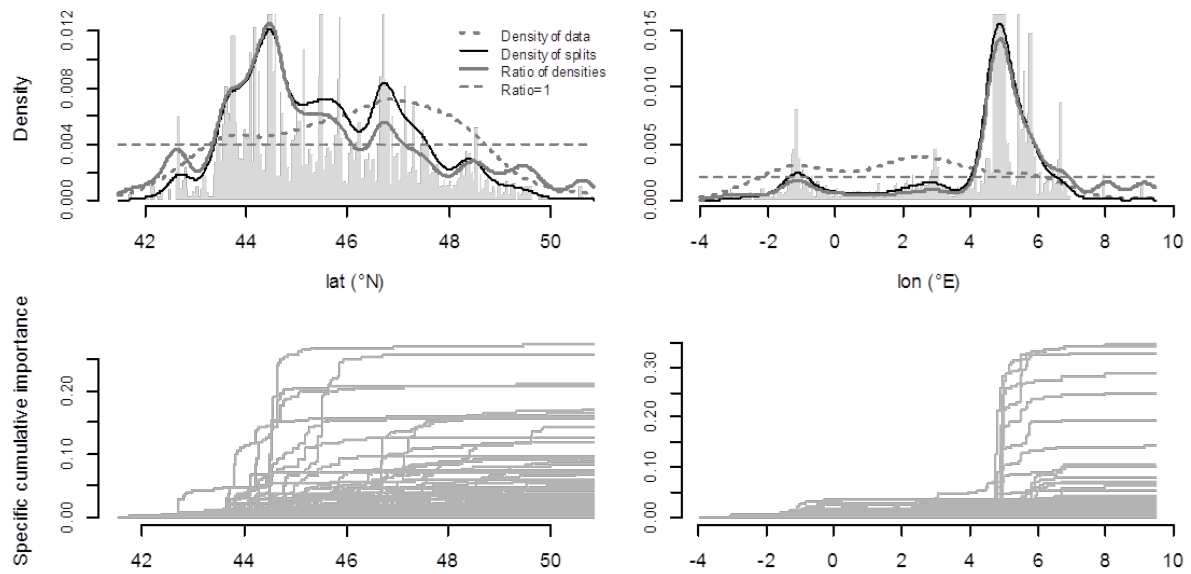
3

1 Fig 3



2

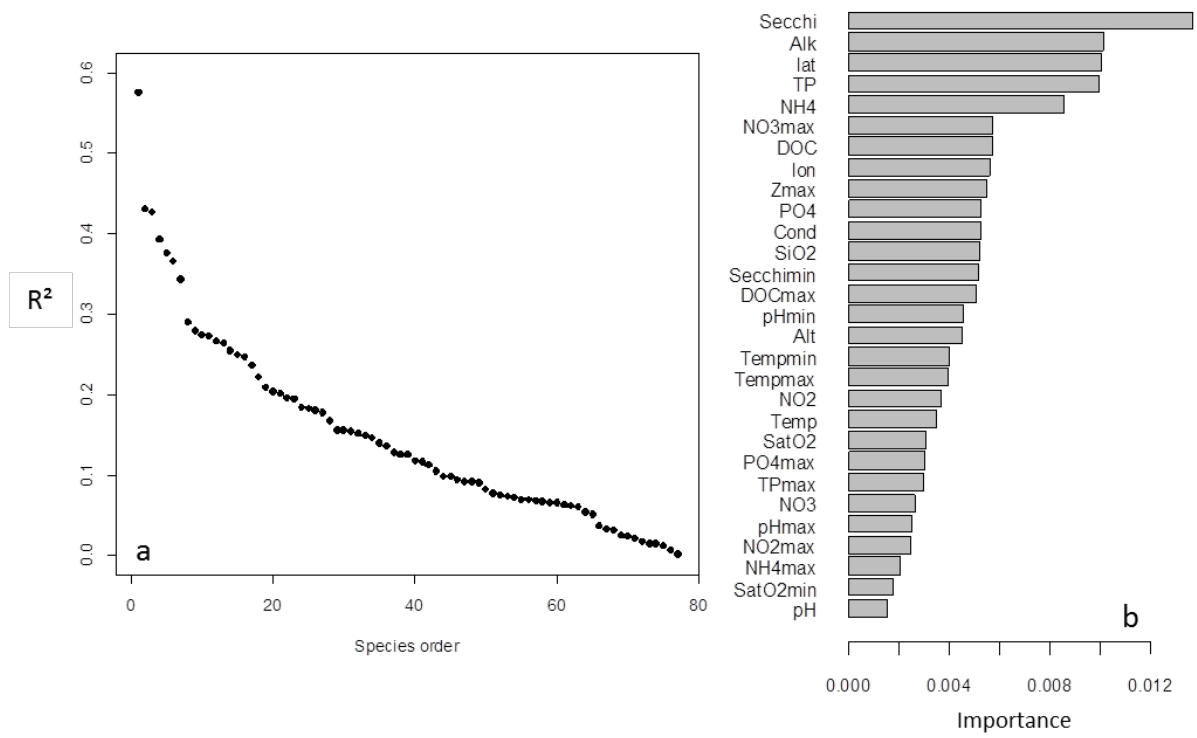
1 Fig 4



2

3

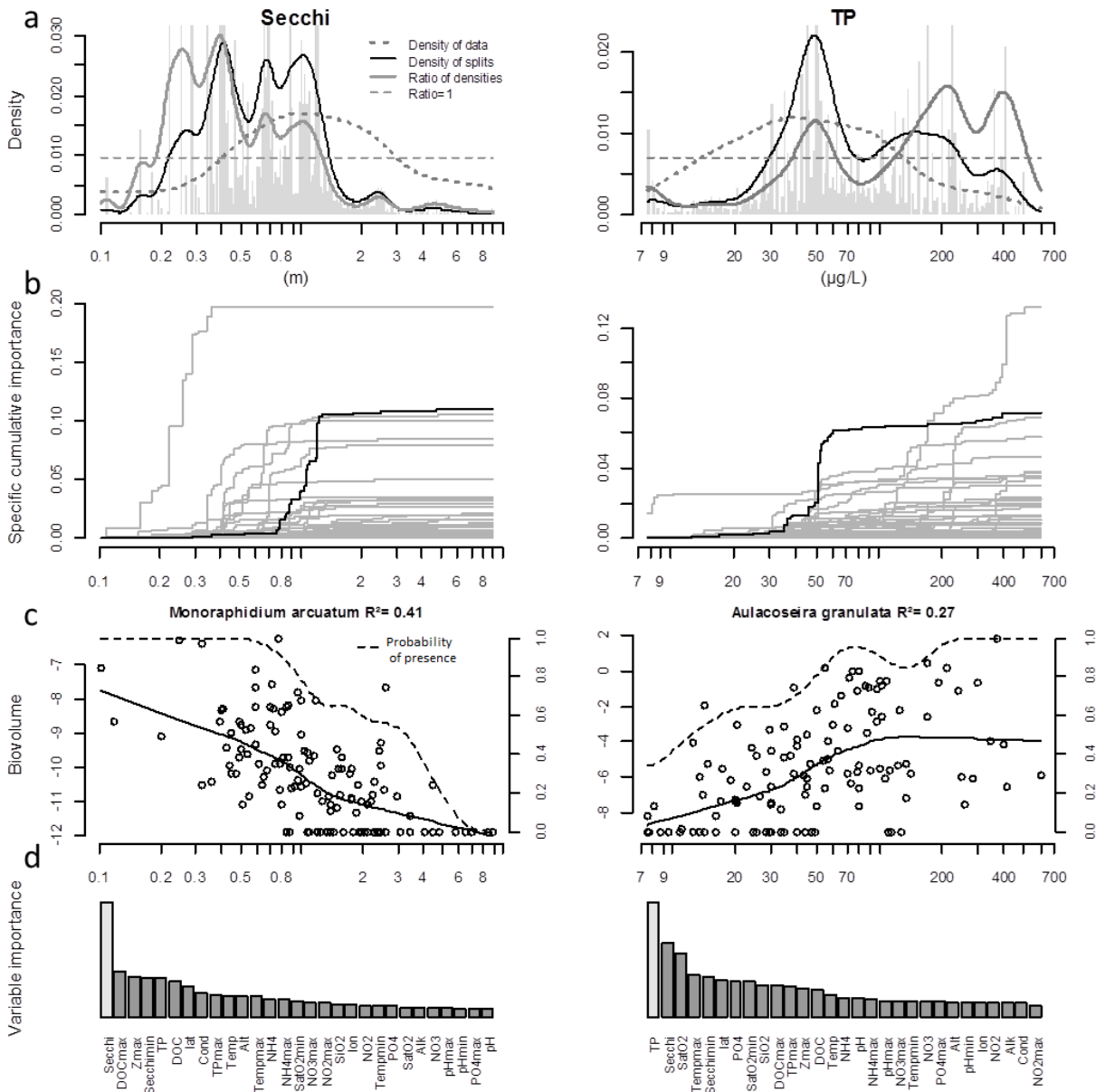
1 Fig 5



2

3

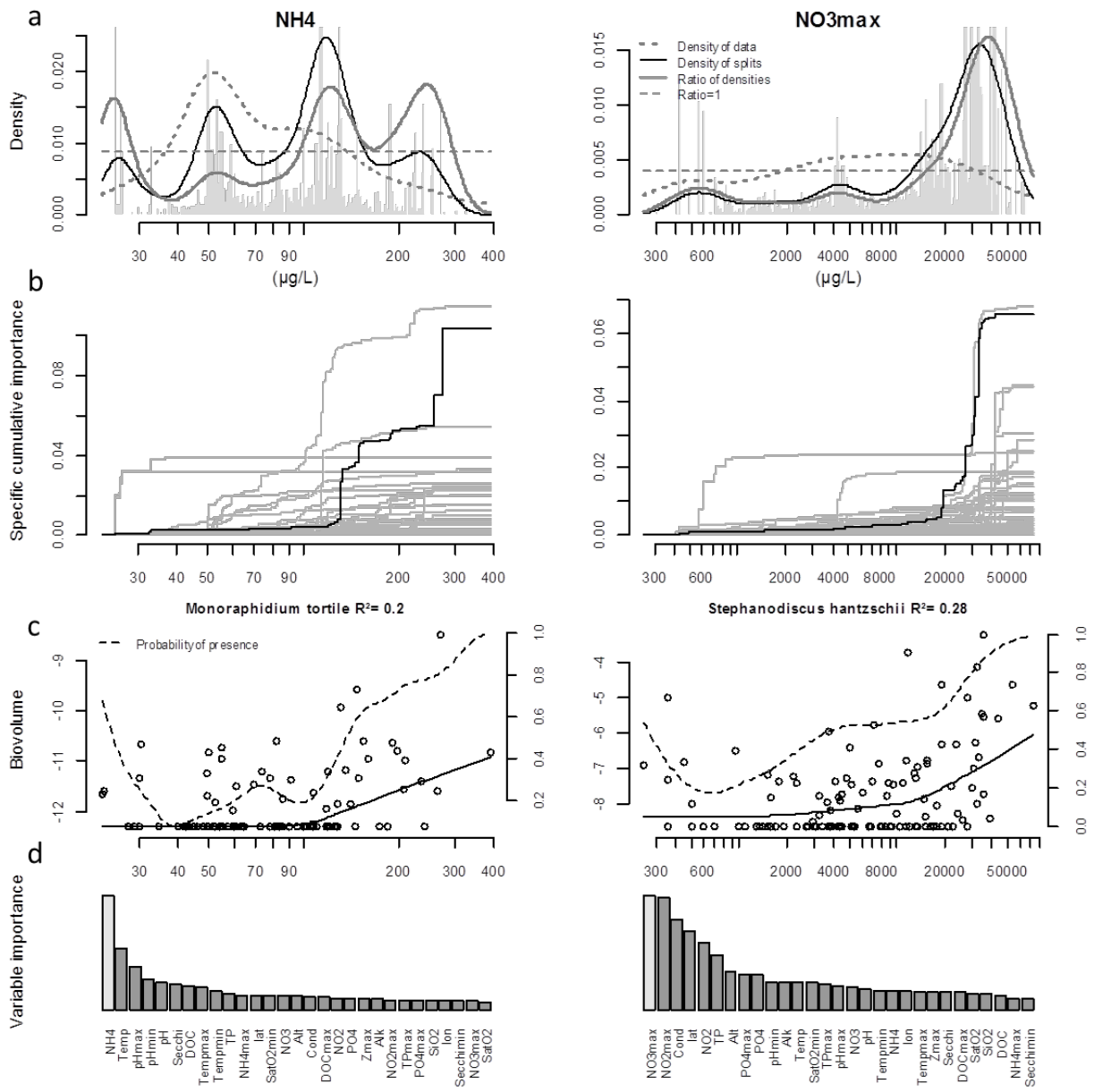
1 Fig 6



2

3

1 Fig 7



2

3

1 Fig 8

