# WILEY

Once you have Acrobat Reader open on your computer, click on the Comment tab at the right of the toolbar:

| 🗎 | 💾 | 🖨 | ✉ | ⬆ | ⬇ | 1 | / 27 | ⊖ | ⊕ | 70.4% | ▾ | 💾 | ⤢ ▾ | Tools | Comment | Share |

This will open up a panel down the right side of the document. The majority of tools you will use for annotating your proof will be in the Annotations section, pictured opposite. We've picked out some of these tools below:

▾ **Annotations**

💬 📝 📎 🔊 👤▾

T̄ₐ 🔁 ┿ T T̄ₚ
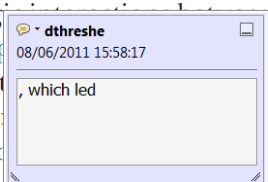
---

## 1. Replace (Ins) Tool – for replacing text.

🔁 Strikes a line through text and opens up a text box where replacement text can be entered.

**How to use it**

- Highlight a word or sentence.
- Click on the Replace (Ins) icon in the Annotations section.
- Type the replacement text into the blue box that appears.

ndard framework for the analysis of m
icy. Nevertheless, it also led to exoge
ole of strategi interaction between fi
nber of comp     o
 is that the st     of
nain compo     be
level, are exc     nc
important works on entry by Ghiro
M henceforth)[1] we open the 'black h

> 💬 ▾ **dthreshe** ⬜
> 08/06/2011 15:58:17
>
> , which led

---

## 2. Strikethrough (Del) Tool – for deleting text.

┿ Strikes a red line through text that is to be deleted.

**How to use it**

- Highlight a word or sentence.
- Click on the Strikethrough (Del) icon in the Annotations section.

there is no room for extra profits a
ups are zero and the number of
et) values are not determined by
Blanchard ~~and Kiyotaki~~ (1987),
rfect competition in general equili
ts of aggregate demand and supply
lassical framework assuming mono
een an exogenous number of firms

---

## 3. Add note to text Tool – for highlighting a section to be changed to bold or italic.

T̄ₚ Highlights text in yellow and opens up a text box where comments can be entered.

**How to use it**

- Highlight the relevant section of text.
- Click on the Add note to text icon in the Annotations section.
- Type instruction on what should be changed regarding the text into the yellow box that appears.

namic responses of mark ups
ent with the VAR evidence

sation     ith
y Ma     ell
 and     led
on n     ber
to a     on
stent also with the demand-

> 💬 ▾ **dthreshe**
> 08/06/2011 15:31:38

---

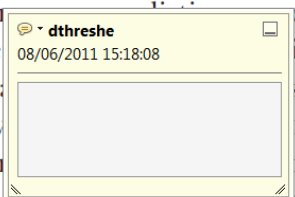## 4. Add sticky note Tool – for making notes at specific points in the text.

💬 Marks a point in the proof where a comment needs to be highlighted.

**How to use it**

- Click on the Add sticky note icon in the Annotations section.
- Click at the point in the proof where the comment should be inserted.
- Type the comment into the yellow box that appears.

land and supply shocks. Most of
a💬mi     eti
numbe     iff
dard fra     sis
icy. Nev     o
ole of st     we
ber of competitors and the imp
is that the structure of the secto

> 💬 ▾ **dthreshe** ⬜
> 08/06/2011 15:18:08

# WILEY

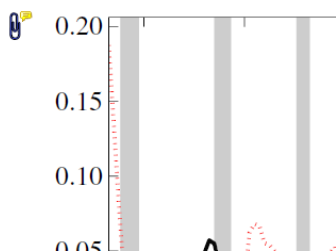5. Attach File Tool – for inserting large amounts of text or replacement figures.

Inserts an icon linking to the attached file in the appropriate place in the text.
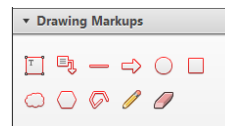
How to use it

- Click on the Attach File icon in the Annotations section.
- Click on the proof to where you'd like the attached file to be linked.
- Select the file to be attached from your computer or network.
- Select the colour and type of icon that will appear in the proof. Click OK.

END



6. Drawing Markups Tools – for drawing shapes, lines and freeform annotations on proofs and commenting on these marks. Allows shapes, lines and freeform annotations to be drawn on proofs and for comment to be made on these marks.



How to use it

- Click on one of the shapes in the Drawing Markups section.
- Click on the proof at the relevant point and draw the selected shape with the cursor.
- To add a comment to the drawn shape, move the cursor over the shape until an arrowhead appears.
- Double click on the shape and type any text in the red box that appears.

# Modeling lake trophic state: a random forest approach

Jeffrey W. Hollister,† W. Bryan Milstead, and Betty J. Kreakie

*Atlantic Ecology Division, Office of Research and Development, National Health and Environmental Effects Research Laboratory, United States Environmental Protection Agency, 27 Tarzwell Drive, Narragansett, Rhode Island 02882 USA*

**Abstract.**   Productivity of lentic ecosystems is well studied, and it is widely accepted that as nutrient inputs increase, productivity increases and lakes transition from lower trophic state (e.g., oligotrophic) to higher trophic states (e.g., eutrophic). These broad trophic state classifications are good predictors of ecosystem condition, services (e.g., recreation and esthetics), and disservices (e.g., harmful algal blooms). While the relationship between nutrients and trophic state provides reliable predictions, it requires in situ water quality data to parameterize the model. This limits the application of these models to lakes with existing and, more importantly, available water quality data. To address this, we take advantage of the availability of a large national lakes water quality database (i.e., the National Lakes Assessment), land-use/land-cover data, lake morphometry data, and other universally available data, and we apply data-mining approaches to predict trophic state. Using these data and random forests, we first model chlorophyll *a* and then classify the resultant predictions into trophic states. The full model estimates chlorophyll *a* with both in situ and universally available data. The mean-squared error and adjusted $R^2$ of this model was 0.09 and 0.8, respectively. The second model uses universally available GIS data only. The mean-squared error was 0.22, and the adjusted $R^2$ was 0.48. The Kappa coefficients of the trophic state classifications derived from the chlorophyll *a* predictions were 0.57 for the full model and 0.29 for the "GIS-only" model. Random forests extend the usefulness of the class predictions by providing prediction probabilities for each lake. This allows us to make trophic state predictions and also indicate the level of uncertainty around those predictions. For the full model, these predicted class probabilities ranged from 0.42 to 1. For the GIS-only model, they ranged from 0.33 to 0.96. It is our conclusion that in situ data are required for better predictions, yet GIS and universally available data provide trophic state predictions, with estimated uncertainty, that still have the potential for a broad array of applications. The source code and data for this manuscript are available from https://github.com/USEPA/LakeTrophicModelling.

**Key words:**   cyanobacteria; harmful algal blooms; National Lakes Assessment; nutrients; open science.

**Received** 30 November 2015; accepted 8 December 2015. Corresponding Editor: D. P. C. Peters.
† **E-mail:** hollister.jeff@epa.gov

# Introduction

Productivity in lentic systems is often categorized across a range of trophic states (e.g., the trophic continuum) from early successional (i.e., oligotrophic) to late successional lakes (i.e., hypereutrophic) with lakes naturally occurring across this range (Carlson 1977). Oligotrophic lakes occur in nutrient poor areas or have a more recent geologic history, are often found in higher elevations, have clear water, and are usually favored for drinking water or direct contact recreation (e.g., swimming). Lakes with higher productivity (e.g., mesotrophic and eutrophic lakes) have greater nutrient loads, tend to be less clear, have greater density of aquatic plants, and often support more diverse and abundant fish communities. Higher primary productivity is not

necessarily a predictor of poor ecological condition as it is natural for lakes to shift from lower to higher trophic states, but this is a slow process (Rodhe 1969). However, at the highest productivity levels (hypereutrophic lakes), biological integrity is compromised (Hasler 1969, Smith et al. 1999, Schindler and Vallentyne 2008).

Monitoring trophic state allows for rapid assessment of a lakes biological productivity and identification of lakes with unusually high productivity (e.g., hypereutrophic). These cases are indicative of lakes under greater anthropogenic nutrient loads, also known as cultural eutrophication, and are more likely to be at risk of fish kills, beach fouling, and harmful algal blooms (Smith 1998, Smith et al. 1999, 2006). Given the association between trophic state and many ecosystem services and disservices, being able to accurately model trophic state could provide a first cut at identifying lakes with the potential for harmful algal blooms (i.e., from cyanobacteria) or other problems associated with cultural eutrophication. This type of information could be used for setting priorities for management and allow for more efficient use of limited resources.

As trophic state and related indices can be best defined by a number of in situ water-quality parameters (modeled or measured), most models have used this information as predictors (Imboden and Gächter 1978, Salas and Martino 1991, Carvalho et al. 2011, Milstead et al. 2013). This leads to accurate models, but these data are often sparse and not always available, thus limiting the population of lakes for which we can make predictions. A possible solution for this issue is to build models that use widely available data that are correlated to many of the in situ variables. For instance, landscape metrics of forests, agriculture, wetlands, and urban land in contributing watersheds have all been shown to explain a significant proportion of the variation (ranging from 50% to 86%, depending on study) of nutrients in receiving waters (Jones et al. 2001, 2004, Seilheimer et al. 2013). Building on these previously identified associations might allow us to use only landscape and other universally available data to build models. Identifying predictors using this type of ubiquitous data would allow for estimating trophic state in both monitored and unmonitored lakes. Furthermore, being able to classify a large number of lakes would

have implications for the management of lakes. A broader discussion of ecological classification and resource management is beyond the scope of this article, but see Carpenter et al. (1999) for [5] more information on this topic.

Many published models of nutrients and trophic state in freshwater systems are based on linear modeling methods such as standard least squares regression or linear mixed models (Jones et al. 2001, 2004). While these methods have proven to be reliable, they have limitations (e.g., independence, distribution assumptions, and outlier sensitivity). Using data-mining approaches, such as random forests, can avoid many of the limitations, may reduce bias, and will often provide better predictions (Breiman 2001, Cutler et al. 2007, Peters et al. 2007, Fernández-Delgado et al. 2014). For instance, random forests are nonparametric, and thus, the data do not need to come from a specific distribution (e.g., Gaussian) and can contain collinear variables (Cutler et al. 2007). Second, random forests work well with very large numbers of predictors (Cutler et al. 2007). Finally, random forests can deal with model selection uncertainty as predictions are based on a consensus of many models and not just a single model selected with some measure of goodness of fit.

The research presented here builds on past work in three areas. First, we built, assessed, and compared two random forest models of chlorophyll *a* with (1) in situ and universally available GIS data and then (2) universally available GIS data only. Second, we converted the chlorophyll *a* estimates, for both models, to trophic state and assessed prediction accuracy and uncertainty. Third, we examined the important predictors for both models. Finally, to promote transparency in our work, the analysis code and data are available as an R package from https://github.com/USEPA/LakeTrophicModelling.

## METHODS

### Data and study area

We utilized three primary sources of data for this study, the National Lakes Assessment (NLA), the National Land Cover Dataset (NLCD), and lake morphometry modeled from the NHDPlus and National Elevation Data Set (Homer et al. 2004, USEPA 2009, Xian et al. 2009, Hollister and Milstead 2010, Hollister et al.

2011, Hollister 2014). All datasets are national in extent and provide a unique snapshot view of the condition of lakes in the conterminous United States during the summer of 2007.

The NLA dataset was collected during the summer of 2007, and the final datasets were released in 2009 (USEPA 2009). With consistent methods and metrics collected at over 1000 locations across the conterminous United States (Fig. 1), the NLA provides a unique opportunity to examine broad scale patterns in lake productivity. The NLA collected data on biophysical measures of lake water quality and habitat as well as an assessment of the phytoplankton community. For this analysis, we only use the various water quality measurements from the NLA (USEPA 2009). Additionally, the NLA included ecological regions as defined in the Wadeable Streams Assessment (Fig. 2) (Omernik 1987, USEPA 2006).
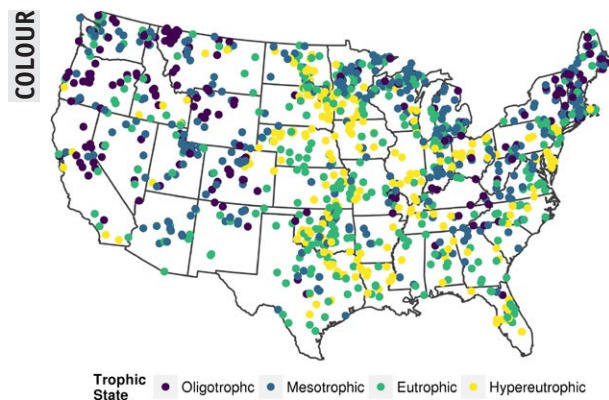


Fig. 1. Map of the distribution of National Lakes Assessment sampling locations.
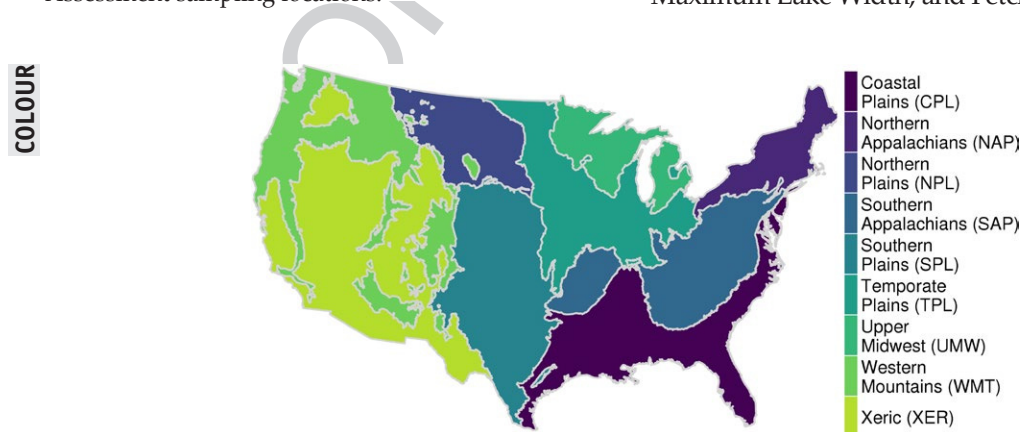
Adding to the monitoring data collected via the NLA, we used the 2006 NLCD data to examine landscape-level drivers of trophic status in lakes. The NLCD is a national land-use/land-cover dataset that also provides estimates of impervious surface. We calculated total proportion of each NLCD land-use land-cover class and total percent impervious surface within a 3-km buffer surrounding each lake (Homer et al. 2004, Xian et al. 2009). We chose this buffer distance for several reasons. First, in some preliminary efforts, we tried a variety of scales (300 m, 1.5 km, and 3 km), and they had little impact on prediction accuracy. Second, since we also include local lake-specific variables (see below) as well as the broader scale ecoregions, we chose the 3-km buffer as it made intuitive sense as representative of land-use impacts that would not be accounted for these other variables. While many regional classifications and scales have been shown to be effective (e.g., Cheruvelil et al. 2013), we chose a 3-km buffer as it represented an intermediate scale that is greater than immediate parcels but smaller than regional and whole-basin measures.

Local lake-specific characteristics have been shown to be important (Read et al. 2015). Thus to account for this, we used measures of lake morphometry (depth, volume, fetch, etc.). As these data are difficult to obtain for large numbers of lakes over broad regions, we used modeled estimates of lake morphometry (Hollister and Milstead 2010, Hollister et al. 2011, Hollister 2014). These included surface area, shoreline length, Shoreline Development, Maximum Depth, Mean Depth, Lake Volume, Maximum Lake Length, Mean Lake Width, Maximum Lake Width, and Fetch.



Fig. 2. Wadeable Streams Assessment ecoregions.

*Predicting trophic state with random forests*

Random forest is a machine-learning algorithm that aggregates numerous decision trees to obtain a consensus prediction of the response categories (Breiman 2001). Bootstrapped sample data are recursively partitioned according to a given random subset of predictor variables, and a predetermined number of decision trees are developed. With each new tree, the sample data subset is randomly selected, and with each new split, the subset of predictor variables is randomly selected. For a more detail description of random forests, see Breiman (2001) and Cutler et al. (2007).

Random forests are able to handle numerous correlated variables without a decrease in prediction accuracy; however, one possible shortcoming of this approach is that the resulting model may be difficult to interpret, thus selecting the most important variables is an important first step. Several methods have been proposed to do this with random forest. For instance, this is a problem often faced in genomics and in that field, a variable selection method based on random forest has been successfully applied and implemented in the R Language as the varSelRF package (Díaz-Uriarte and De Andres 2006), but this is limited to classification problems. Additionally, others have suggested alternative variable importance measures, but this is only needed with a large number of categorical variables which are selected against with traditional random forest approach (Strobl et al. 2007).

In our case, we predicted a continuous variable, chlorophyll *a*, directly thus varSelRF, does not apply, and nearly all of our variables are continuous **6** so the approach suggested by Strobl (2007) is not necessary. Thus, we developed an approach, similar to varSelRF but applied to random forest with regression trees. With this approach, we fit a full random forest model that includes all variables and a large number of trees. We then rank the variables using the increase in mean square error, which has been shown to be a less biased metric of importance than the mean decrease in the Gini coefficient (Strobl et al. 2007). Using this ranking, we then iterate through the variables and create a random forest with the top two variables and record mean square error and adjusted *R* of the resultant random forest. We then repeat this process by adding the next most important variable in order of importance. With this information, we

identify both the top variables and the point at which adding variables does not improve the fit of the overall model. These variables are selected and used as the "reduced model." With this method, a minimum set of variables that maximizes model accuracy is provided. This allows us to start with a full suite of predictor variables from which to select a minimum easier to interpret set of variables.

*Model details*

We used the randomForest package in R to build predictive models of chlorophyll *a* with two sets of predictors (Liaw and Wiener 2002). The first included in situ and universally available GIS predictors. We refer to this as the "all-variables" model. For the second model, we used just the universally available data (i.e., no in situ information). This is referred to as the GIS-only model. A list of all considered variables is in Appendix S1: Table S1. Our separation of predictors was chosen so that we could highlight the additional predictive performance provided by adding the in situ water quality variables on top of the GIS-only variables. Finally, we used only complete cases (i.e., missing data were removed) so the total number of observations varied among models.

Our modeling work flow was as follows:

1. Identify a minimal set of variables from the full suite of variables (Appendix S1: Table S1) that maximize accuracy of the random forest algorithm. This minimal set of variables, the reduced model, is calculated for each of the models.
2. Using R's randomForest package, we develop two random forest models with 5000 trees (all variables and GIS only).
3. Assess model performance for both the predicted chlorophyll *a* and categorical trophic state classifications. Trophic state was defined using the NLA chlorophyll *a* trophic state cut-offs (Table 1).
4. Examine importance and partial dependence of the most important variables.

*Measures of model performance and variable importance*

We assessed the performance of the random forest two ways. First, we compared the root mean square error and the adjusted *R* of the

Table 1.   Chlorophyll a-based trophic state cut-offs.

| Trophic state (4 class) | Trophic state (2 class) | μg/L Cut-offs |
| --- | --- | --- |
| Oligotrophic | Oligotrophic/mesotrophic | ≤2 |
| Mesotrophic | Oligotrophic/mesotrophic | >2 and ≤7 |
| Eutrophic | Eutrophic/hypereutrophic | >7 and ≤30 |
| Hypereutrophic | Eutrophic/hypereutrophic | >30 |

Table 2.   Random forest confusion matrix for all-variables model converted to four trophic states. Columns show predicted values and rows show observed values. Agreement indicated on diagonal and accuracy for each trophic state indicated in "Class Accuracy" column.

|  | Oligotrophic | Mesotrophic | Eutrophic | Hypereutrophic | Class accuracy (%) |
| --- | --- | --- | --- | --- | --- |
| Oligotrophic | 115 | 31 | 0 | 0 | 78.77 |
| Mesotrophic | 67 | 251 | 63 | 0 | 65.88 |
| Eutrophic | 7 | 61 | 217 | 75 | 60.28 |
| Hypereutrophic | 0 | 5 | 29 | 159 | 82.38 |

Table 3.   Random forest confusion matrix for GIS-only model converted to four tropic states. Columns show predicted values and rows show observed values. Agreement indicated on diagonal and accuracy for each trophic state indicated in "Class Accuracy" column.

|  | Oligotrophic | Mesotrophic | Eutrophic | Hypereutrophic | Class accuracy (%) |
| --- | --- | --- | --- | --- | --- |
| Oligotrophic | 65 | 14 | 6 | 0 | 76.47 |
| Mesotrophic | 101 | 213 | 98 | 18 | 49.53 |
| Eutrophic | 29 | 126 | 193 | 141 | 39.47 |
| Hypereutrophic | 1 | 8 | 38 | 87 | 64.93 |

models. Second, we examined the accuracy of the model predictions when converted to trophic states classes via a confusion matrix (Tables 2 and 3). A confusion matrix shows agreement and disagreement in a tabular form with predicted values forming the columns of the matrix and observed values, the rows. From this tabulated information, we calculated the total accuracy (i.e., percent correctly predicted) and the Kappa coefficient, which takes into account the error expected by chance alone (i.e., the off-diagonal values of the matrix) (Cohen 1960, Hubert and Arabie 1985). The Kappa coefficient can range from −1 to 1 with 0 equaling the agreement expected by chance alone. Values >0 represent agreement greater than would be expected by chance. A Kappa coefficient greater than approximately 0.6 is considered "substantial" agreement (Landis and Koch 1977). Negative values are rare and would indicate no agreement between the predicted and observed values. We use Kappa as a means of comparison across models as well as within subsets of a given model. Additionally, random forest builds each tree on bootstrapped, random subsets of the original data, and thus, a separate independent validation dataset is not required and random forest error estimates are expected to be unbiased (Breiman 2001).

Random forests explicitly measure variable importance with two metrics: mean decrease in Gini and percent increase in mean-squared error. These measure the impact on the overall model when a particular variable is included and thus can be used to assess importance (Breiman 2001). The Gini Index has been shown to have a bias (Strobl et al. 2007), and thus, we used percent increase in mean-squared error to assess variable importance. Finally, partial dependence plots provide a mechanism to examine the partial relationship between individual

variables and the response variable (Jones and Linder 2015). We examined these plots for the top variables as assigned by percent increase in mean-squared error for each the reduced models.

### Trophic state probabilities

One of the powerful features of random forests is the ability to aggregate a very large number of competing models or trees. Each tree provides an independent prediction or vote for a possible outcome. In the context of our chlorophyll *a* models, we have 5000 estimates of chlorophyll *a* for each lake. We convert these values to trophic states (Table 1) and then sum total votes for each class and divide by total possible votes to get an estimate of the probability that a lake is in a given trophic state. For instance, for a single lake (National Lake Assessment ID = NLA06608-0005), the vote probabilities for the all-variables model were 95% for oligotrophic, 5% for mesotrophic, 0% for eutrophic, and 0% for hypereutrophic. The maximum probability provides the predicted class, in this case oligotrophic, and suggests little uncertainty in this prediction. We refer to this value as the "prediction probability."

Further, we might expect higher total accuracy for lakes that have more certain predictions. This should be evident by looking at the Kappa coefficient of lakes given their prediction probability is at or above a certain probability. To test this, we use an approach similar to one outlined by Paul and McDonald (2005) and implemented by Hollister et al. (2008) and examine the change in Kappa coefficient as a function of the prediction probability for both models.

## RESULTS

Our complete dataset included 1148 lakes; however, five lakes did not have chlorophyll *a* data. Thus, the base dataset for our modeling was conducted on data for 1143 lakes. The lakes were well distributed across the four trophic state categories (Table 1) and spatially throughout the United States (Fig. 1).

### Models: "All variables"

The model built with all predictors used 1080 total observations, had a mean-squared error
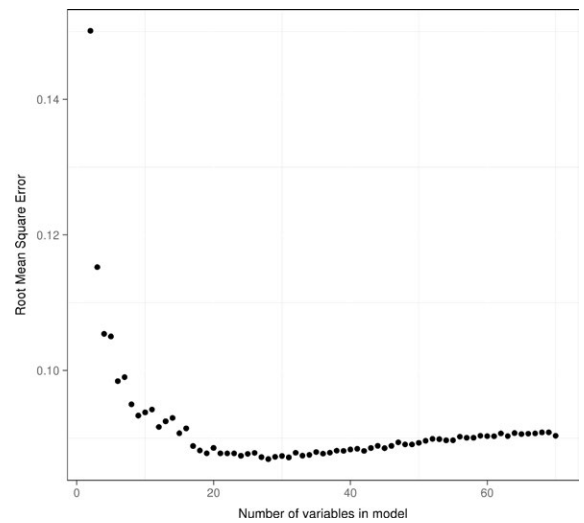


Fig. 3. Variable selection plot for all variables. Shows percent increase in mean-squared error as a function of the number of variables.

of 0.09 and *R* of 0.8. The accuracy of the four trophic states was 68.7%, the Kappa coefficient was 0.57, and class accuracy ranged from 60% to 82% (Table 2). The variable selection process identified a reduced model with 20 variables (Fig. 3). The six most important variables were turbidity, total phosphorus, total nitrogen, elevation, total organic carbon, and N:P ratio (Figs. 4 and 5). The role that each played in predicting chlorophyll *a* varied (Fig. 5).

### Models: GIS-only variables

The GIS-only model was built using 1138 total observations, had a mean-squared error of 0.22 and *R* 0.48. Four trophic states were predicted with a total accuracy of 49%, had a Kappa coefficient of 0.29, and class accuracy ranged from 39% to 76% (Table 3). The variable selection process for this model produced a reduced model with 15 variables (Fig. 6). The six most important variables were ecoregion, percent cropland, elevation, latitude, percent evergreen forest, and mean lake depth (Figs. 7 and 8).

### Trophic state probabilities

The all-variables model provides more certain model predictions with a median prediction probability of 0.81 vs. 0.72 for the GIS-only model (Fig. 9). Additionally, the Kappa
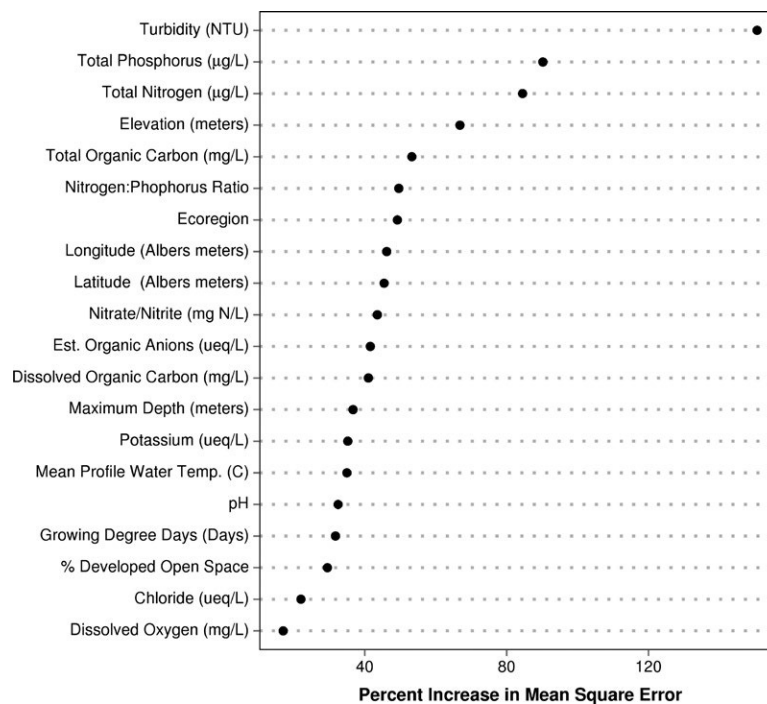
Fig. 4. Importance plot for all variables. Shows percent increase in mean square error. Higher values of percent increase in mean-squared error indicate higher importance.

coefficient of the predictions is a function of this uncertainty. Lakes with more certain predictions were more accurately classified and had higher Kappa coefficients (Fig. 10). For both models, when prediction probabilities are approximately 0.8 or higher, the models had a Kappa coefficient of ~1. This represents 55% of the lakes for the all-variables model and 22% of the lakes for the GIS-only model. A Kappa coefficient of 0.6 or higher is considered "substantial" agreement (Landis and Koch 1977). For the GIS-only model, this is seen with 52% of the lakes. Finally, as prediction probabilities increased, the difference in Kappa coefficient between the two models decreased (Fig. 10 and Tables 4 and 5).

## DISCUSSION

### Trophic state probabilities

Not surprisingly, lakes with more certain predictions (i.e., higher prediction probabilities) were more accurately predicted (Fig. 10). The fact that the difference in accuracy (as measured by the Kappa coefficient) between the two

models decreased as certainty in the prediction increased suggests that models with lower overall accuracy, such as the GIS-only model, may have acceptable accuracy for many individual cases (Tables 4 and 5). Additionally, the prediction probabilities may be mapped for each of the four classes (Fig. 11). The spatial patterns show little variability between the all-variables and GIS-only models; thus, we only show the results from the more broadly applicable GIS-only model (Fig. 11).

This map provides several insights. First, since low uncertainty is associated with high accuracy, this map shows the broad spatial patterns of lake trophic state across the United States (i.e., darker colors more likely to be correctly predicted). Hypereutrophic lakes are much more commonly predicted in the Midwest and southeastern United States. Clear, oligotrophic lakes are in the northwestern United States, through the western mountains, and in the northeastern United States. The middle trophic states are more evenly distributed across the country. Finally, this particular map is very similar to simply mapping the raw data. However, it highlights what could
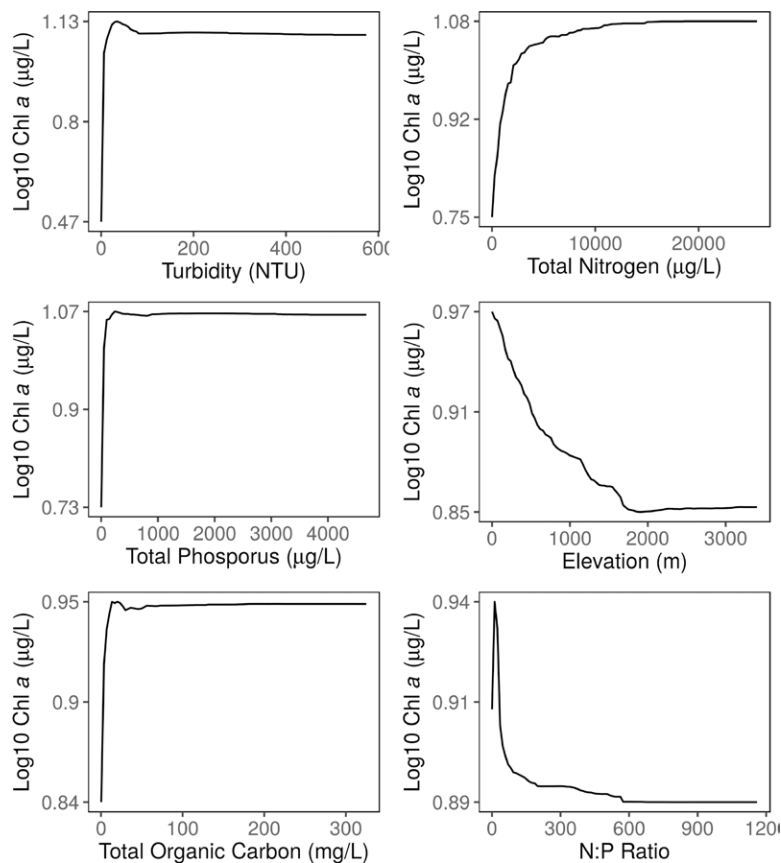
Fig. 5. All-variables partial dependence plots for the six most important variables.
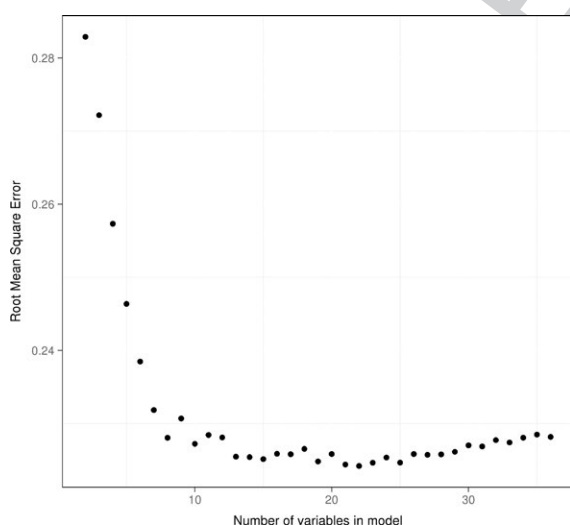


Fig. 6. Variable selection plot for GIS-only variables. Shows percent increase in mean-squared error as a function of the number of variables.

be done if the GIS-only model were used to map data without measured chlorophyll *a* values which would provide probabilities of given trophic states for all lakes in the United States.

*Partial dependencies of explanatory variables*

In line with past predictive modeling of chlorophyll *a* concentrations, the all-variables model selected the water quality variables (turbidity, total organic carbon, total nitrogen, total phosphorus, and N:P ratios) as important variables (Downing et al. 2001). While there is variation in the response of chlorophyll *a* to changes in nutrient concentrations, the general pattern suggests that limiting nutrients have predictable impacts. If we examine the partial dependencies of these variables, we see a general linear increase in log chlorophyll *a* with nitrogen, phosphorus, and organic carbon concentrations (Fig. 5). This relationship holds until nutrient
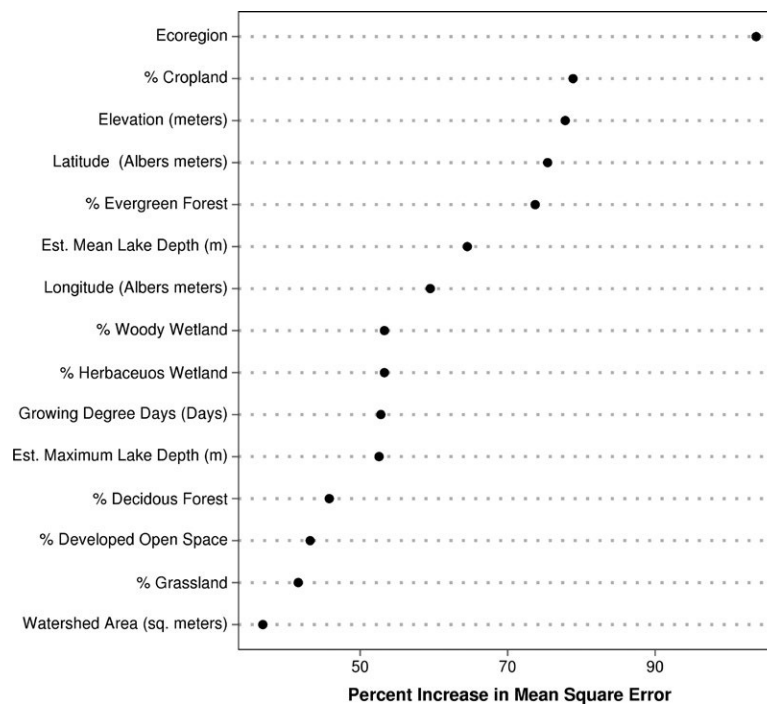
Fig. 7. Importance plot for GIS-only variables. Shows percent increase in mean square error. Higher values of percent increase in mean-squared error indicate higher importance.

concentrations become saturated. The partial dependency plots (Fig. 5) for the nitrogen:phosphorus ratio is more complicated, indicating that for ratios less than ~14 chlorophyll *a* increases but after ~14 there is marked decrease. The effect of the nitrogen:phosphorus ratio on chlorophyll has been the subject of considerable research, and our results are consistent with the majority of the findings suggesting that at low ratio values nitrogen is limiting (Downing and McCauley 1992, Smith and Schindler 2009). Conversely, at higher ratios, the phosphorus levels may be limiting. This would be a cause for concern with linear models; however, linearity is not an assumption of tree-based modeling approaches such as random forest.

Turbidity was selected as the most important variable in the all-variables model. The partial dependency analysis shows that, similar to the nutrients discussed above, log chlorophyll *a* increases with increased turbidity. At first this may seem counter intuitive, since we might expect productivity to decrease as turbidity increases, and therefore light availability decreases (Tilzer 1988, Bilotta and Brazier 2008). However, algal

biomass can contribute heavily to measures of turbidity and we expect greater productivity to lead to increased turbidity (Hansson 1992). We interpret this pattern as indicating that as chlorophyll *a* concentrations increase we see a concomitant increase in turbidity due to increased algal cell densities.

Elevation was selected as an important predictive variable in both the all-variables and the GIS-only models; the partial dependencies (Figs. 5 and 8) indicate a negative relationship between elevation and chlorophyll *a* concentration that is probably due to fact that the location of mountains in the United States is the spatial inverse of the distribution of agricultural and urban lands. As elevation increases, we expect decreased loads due to smaller watershed contributing areas. In contrast, lower elevation sites will have larger drainage areas and greater potential for increased nutrient loads from urban and agricultural sources.

The variables in the GIS-only model captured the large-scale spatial pattern of the trophic status gradient of lakes across the United States. In addition to elevation, mentioned above, the model
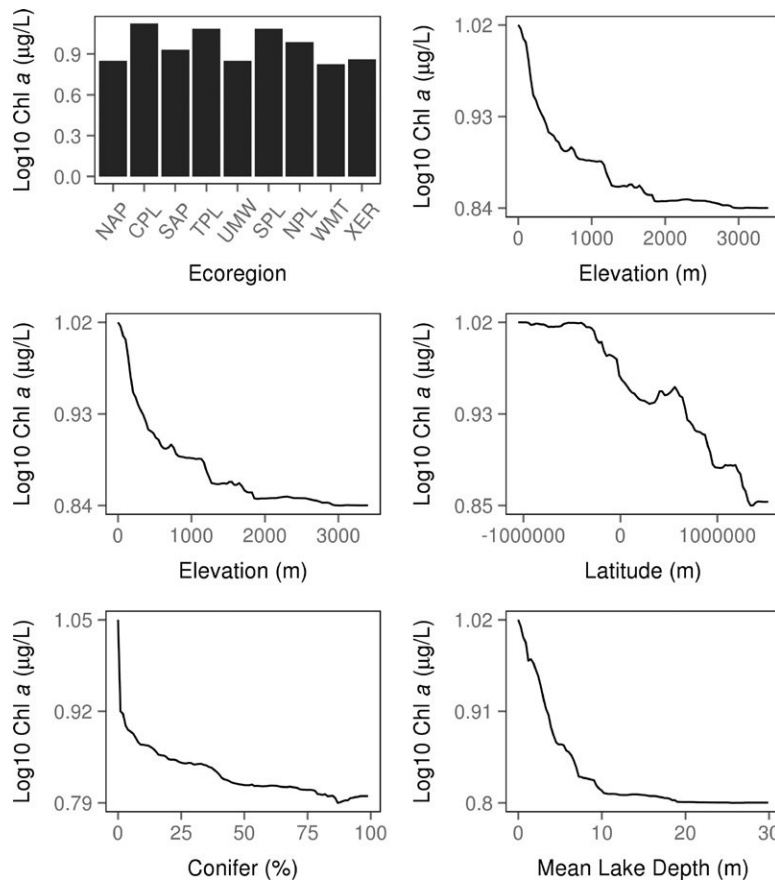
Fig. 8.  GIS-only variables partial dependence plots for the six most important variables.
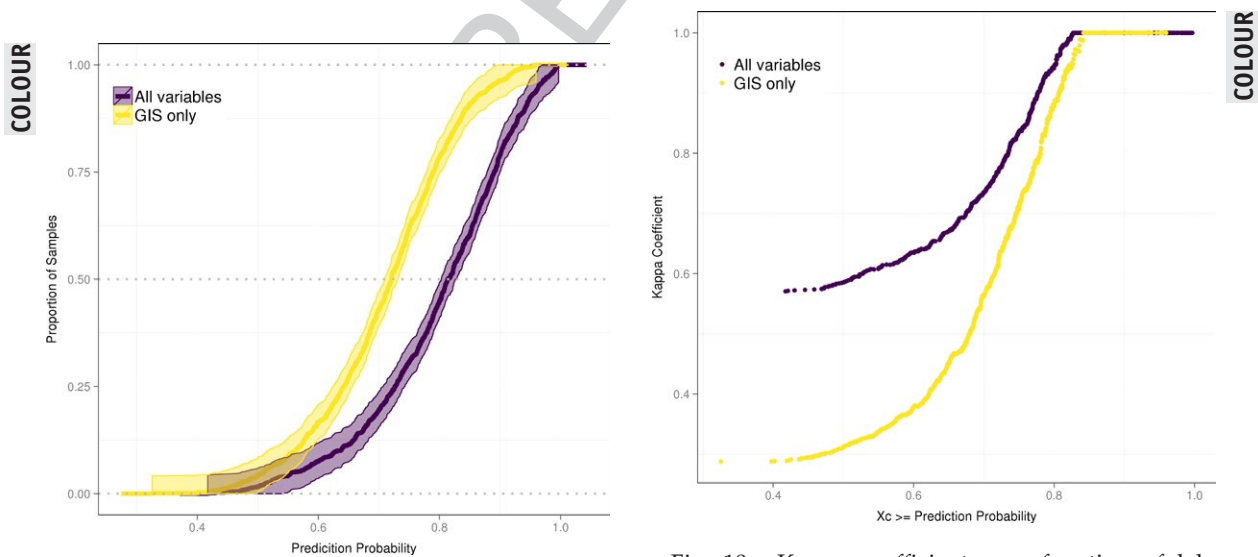


Fig. 9.  Prediction probabilities for the all-variables and GIS-only models.



Fig. 10.  Kappa coefficient as a function of lake prediction probability. The *x*-axis represents lakes with a prediction probability at a given level or higher.

Table 4. Summary of relationship between prediction probabilities, Kappa coefficient, and number of lakes for the all-variables model.

| Prediction prob. | Kappa coefficient | Percent of sample | Number of samples |
|---|---|---|---|
| All | 57 | 100 | 1080 |
| 0.50 | 59 | 98 | 1063 |
| 0.60 | 63 | 92 | 999 |
| 0.70 | 73 | 81 | 870 |
| 0.80 | 95 | 55 | 596 |
| 0.90 | 100 | 21 | 227 |

Table 5. Summary of relationship between prediction probabilities, Kappa coefficient, and number of lakes for the GIS-only model.

| Prediction prob. | Kappa coefficient | Percent of sample | Number of samples |
|---|---|---|---|
| All | 29 | 100 | 1138 |
| 0.50 | 31 | 96 | 1091 |
| 0.60 | 38 | 83 | 949 |
| 0.70 | 56 | 57 | 651 |
| 0.80 | 88 | 22 | 247 |
| 0.90 | 100 | 4 | 43 |

**COLOUR**



Fig. 11. Maps of prediction probabilities for each of the four chlorophyll *a* trophic states.

was most sensitive to latitude and ecoregion. In general, chlorophyll *a* concentrations are highest in the Southern portions of the study area where temperatures can be higher (a known driver of productivity), elevations lower, and agricultural impacts more pronounced. Likewise, ecoregion (Figs. 2 and 8) has a pronounced affect indicating continental scale effects of land use and geogra-

phy. Agriculturally dominated landscapes such as the Temperate Plains, Southern Plains, and Coastal Plains show the highest levels of chlorophyll *a*, whereas high elevation zones (Western Mountains), arid lands (Xeric), and Northern habitats (Upper Midwest) have lower concentrations.

Further evidence for the role of land-use/land-cover variables, and similar to results from Read et al. (2015), is shown by the selection of the percent cropland and percent evergreen forest variables. As indicated by the partial dependency plots (Fig. 8), chlorophyll *a* increases with cropland and decreases with evergreen cover. It is not surprising that croplands were selected given the overwhelming impact of agriculture on the eutrophication process. Evergreens and chlorophyll *a* concentrations show a negative association (Fig. 8). As the percent of evergreens increases, we are likely to see increased elevation and soil differences that limit agriculture.

Lastly, morphometry (e.g., depth) also proved to be important in the prediction of lake trophic state (Genkai-Kato and Carpenter 2005). As morphometry shows little to no broad scale spatial pattern and is unique to a given lake, these data are likely illuminating the local, lake-scale drivers such as in-lake nutrient processing and residence time.

## Conclusions

Our research goals were to explore the utility of a widely used data-mining algorithm, random forests, in the modeling of chlorophyll *a* and lake trophic state. Further, we hoped to examine the utility of these models when built with only ubiquitous GIS data, which allows estimation of trophic state for all lakes in the United States. The all-variables model had an [7] RMSE of 0.09 and an adjusted $R$ of 0.8, whereas the GIS-only models had an RMSE of 0.22 and the adjusted $R$ was 0.48. Our total accuracy in predicting chlorophyll *a*-based trophic states was 69% for the all-variables model and 49% for the GIS-only model.

While the GIS-only model showed lower prediction accuracies than the all-variables model, the association between the uncertainty of prediction and the Kappa coefficient (Fig. 10 and Tables 4 and 5) suggest that the GIS-only model will provide reasonable estimates of trophic state for many lakes

across the United States. Furthermore, we can map the uncertainty of the predictions, and thus, we know the spatial patterns and location of the lakes for which we are certain, or not, of their predicted trophic state. With this, plus the fact that these models may be applied to any lake in the United States, we can recommend using this model.

Future iterations of this modeling effort may be able to utilize modeled predictions of nutrients to improve accuracy and also maintain broad applicability (Milstead et al. 2013). Changes such as these have several advantages. First, this would allow for estimating changes to chlorophyll *a* and trophic state as a function of changing nutrient loads, which are expected due to climate change (Adrian et al. 2009, Jeppesen et al. 2011, Moss et al. 2011, Jones and Brett 2014). Second, with the ability to make predictions for most lakes in the United States, the GIS-only models could be used as a source of information on national scale phenomena. For example, predictions of chlorophyll *a*, with measures of uncertainty, could be used in efforts to scale up the contributions from lakes to broad scale estimates of gross primary production.

For the all-variables model, the in situ water quality variables drove the predictions. This is not surprising. For the GIS-only model, the results were more nuanced. Three broad categories were routinely being selected as important: broad scale spatial patterns in trophic state, land-use/land-cover controls of trophic state, and local, lake-scale control driven by lake morphometry.

Our results raise three important considerations related to managing eutrophication. First, the broad scale patterning, indicated by ecoregion as an important variable, suggests regional trends. This is noteworthy because it suggests that efforts to monitor, model, and manage eutrophication and cyanobacteria should be undertaken at both national and regional levels. This corroborates past findings that regional drivers are important for water quality (Cheruvelil et al. 2013). Second, while direct control of water quality in lakes would have a large impact, the land-use/land-cover drivers (i.e., nonpoint sources) of water quality are also important, and better management of the spatial distribution of important classes such as forest and agriculture can provide some level of control on trophic state and amount of cyanobacteria present. Third, in-lake processes (residence time, nutrient cycling, etc.)

are, as expected, important and need to be part of any management strategy. Building on these efforts through updated models, direct prediction of cyanobacteria, and additional information on the regional differences will help us get a better handle on the broad scale dynamics of productivity in lakes and the potential risk to human health from cyanobacteria blooms.

## ACKNOWLEDGMENTS

## LITERATURE CITED

Adrian, R., et al. 2009. Lakes as sentinels of climate change. Limnology and Oceanography 54:2283–2297.

Bilotta, G., and R. Brazier. 2008. Understanding the influence of suspended solids on water quality and aquatic biota. Water Research 42:2849–2861.

Breiman, L. 2001. Random forests. Machine Learning 45:5–32.

Carlson, R. E. 1977. A trophic state index for lakes. Limnology and Oceanography 22:361–369.

Carpenter, C. A., W. N. Busch, D. T. Cleland, J. Gallegos, R. Harris, R. Holm, C. Topik, and A. Williamson. 1999. The use of ecological classification in management. Pages 395–430 in R. C. Szaro, N. C. Johnson, W. T. Sexton, and A. J. Malk, editors. Ecological stewardship: a common reference for ecosystem management, Vol. 2. Elsevier Science, Oxford, UK.

Carvalho, L., C. A. Miller, E. M. Scott, G. A. Codd, P. S. Davies, and A. N. Tyler. 2011. Cyanobacterial blooms: statistical models describing risk factors for national-scale lake assessment and lake management. Science of the Total Environment 409:5353–5358.

Cheruvelil, K., P. Soranno, K. Webster, and M. Bremigan. 2013. Multi-scaled drivers of ecosystem state: quantifying the importance of the regional spatial scale. Ecological Applications 23:1603–1618.

Cohen, J. 1960. A coefficient of agreement for nominal scales. Educational and Psychological Measurement 20:37–46.

Cutler, D. R., T. C. Jr Edwards, K. H. Beard, A. Cutler, K. T. Hess, J. Gibson, and J. J. Lawler. 2007. Random forests for classification in ecology. Ecology 88:2783–2792.

Díaz-Uriarte, R., and S. A. De Andres. 2006. Gene selection and classification of microarray data using random forest. BMC Bioinformatics 7:3.

Downing, J. A., and E. McCauley. 1992. The nitrogen:phosphorus relationship in lakes. Limnology and Oceanography 37:936–945.

Downing, J. A., S. B. Watson, and E. McCauley. 2001. Predicting cyanobacteria dominance in lakes. Canadian Journal of Fisheries and Aquatic Sciences 58:1905–1908.

Fernández-Delgado, M., E. Cernadas, S. Barro, and D. Amorim. 2014. Do we need hundreds of classifiers to solve real world classification problems? Journal of Machine Learning Research 15:3133–3181.

Genkai-Kato, M., and S. R. Carpenter. 2005. Eutrophication due to phosphorus recycling in relation to lake morphometry, temperature, and macrophytes. Ecology 86:210–219.

Hansson, L.-A. 1992. Factors regulating periphytic algal biomass. Limnology and Oceanography 37:322–328.

Hasler, A. D. 1969. Cultural eutrophication is reversible. BioScience 19:425–431.

Hollister, J. W. 2014. Lakemorpho: Lake morphometry in R. R package version 1.0. http://CRAN.R-project.org/package=lakemorpho

Hollister, J., and W. B. Milstead. 2010. Using GIS to estimate lake volume from limited data. Lake and Reservoir Management 26:194–199.

Hollister, J. W., H. A. Walker, and J. F. Paul. 2008. CProb: a computational tool for conducting conditional probability analysis. Journal of Environmental Quality 37:2392–2396.

Hollister, J. W., W. B. Milstead, and M. A. Urrutia. 2011. Predicting maximum lake depth from surrounding topography. PLoS ONE 6:e25764.

Homer, C., C. Huang, L. Yang, B. Wylie, and M. Coan. 2004. Development of a 2001 national land-cover database for the United States. Photogrammetric Engineering and Remote Sensing 70:829–840.

Hubert, L., and P. Arabie. 1985. Comparing partitions. Journal of Classification 2:193–218.

Imboden, D., and R. Gächter. 1978. A dynamic lake model for trophic state prediction. Ecological Modelling 4:77–98.

Jeppesen, E., et al. 2011. Climate change effects on nitrogen loading from cultivated catchments in Europe: implications for nitrogen retention, ecological state of lakes and adaptation. Hydrobiologia 663:1–21.

Jones, J., and M. T. Brett. 2014. Lake nutrients, eutrophication, and climate change. Pages 273–279 *in* B. Freedman. Global environmental change. Springer, ?????, the Netherlands.

Jones, Z., and F. Linder. 2015. Exploratory data analysis using random forests. http://zmjones.com/static/papers/rfss_manuscript.pdf

Jones, K. B., A. C. Neale, M. S. Nash, R. D. Van Remortel, J. D. Wickham, K. H. Riitters, and R. V. O'Neill. 2001. Predicting nutrient and sediment loadings to streams from landscape metrics: a multiple watershed study from the United States mid-Atlantic region. Landscape Ecology 16:301–312.

Jones, J., M. Knowlton, D. Obrecht, and E. Cook. 2004. Importance of landscape variables and morphology on nutrients in Missouri reservoirs. Canadian Journal of Fisheries and Aquatic Sciences 61:1503–1512.

Landis, J. R., and G. G. Koch. 1977. The measurement of observer agreement for categorical data. Biometrics 33:159–174.

Liaw, A., and M. Wiener. 2002. Classification and regression by randomForest. R News 2:18–22.

Milstead, W. B., J. W. Hollister, R. B. Moore, and H. A. Walker. 2013. Estimating summer nutrient concentrations in northeastern lakes from SPARROW load predictions and modeled lake depth and volume. PLoS ONE 8:e81457.

Moss, B., et al. 2011. Allied attack: climate change and eutrophication. Inland Waters 1:101–105.

Omernik, J. M. 1987. Ecoregions of the conterminous united states. Annals of the Association of American Geographers 77:118–125.

Paul, J. F., and M. E. McDonald. 2005. Development of empirical, geographically specific water quality criteria: a conditional probability analysis approach. ????? 41:1211–1223.

Peters, J., B. D. Baets, N. E. Verhoest, R. Samson, S. Degroeve, P. D. Becker, and W. Huybrechts. 2007. Random forests as a tool for ecohydrological distribution modelling. Ecological Modelling 207:304–318.

Read, E. K., et al. 2015. The importance of lake-specific characteristics for water quality across the continental united states. Ecological Applications 25:943–955.

Rodhe, W. 1969. Crystallization of eutrophication concepts in northern Europe. Pages 50–64 *in* ?????. ?????, editors. Proceedings of Symposium on Eutrophication: causes, consequences, correctives. National Academy of Sciences, Washington, D.C.

Salas, H. J., and P. Martino. 1991. A simplified phosphorus trophic state model for warm-water tropical lakes. Water Research 25:341–350.

Schindler, D. W., and J. R. Vallentyne. 2008. The algal bowl: overfertilization of the world's freshwaters and estuaries. University of Alberta Press, Edmonton, California.

Seilheimer, T. S., P. L. Zimmerman, K. M. Stueve, and C. H. Perry. 2013. Landscape-scale modeling of water quality in Lake Superior and Lake Michigan watersheds: how useful are forest-based indicators? Journal of Great Lakes Research 39:211–223.

Smith, V. H. 1998. Cultural eutrophication of inland, estuarine, and coastal waters. Pages 7–49 *in* M. L. Pace, and M. L. Groffman, editors. Successes, limitations, and frontiers in ecosystem science. Springer, New York, New York, USA.

Smith, V. H., and D. W. Schindler. 2009. Eutrophication science: where do we go from here? Trends in Ecology and Evolution 24:201–207.

Smith, V. H., G. D. Tilman, and J. C. Nekola. 1999. Eutrophication: impacts of excess nutrient inputs on freshwater, marine, and terrestrial ecosystems. Environmental Pollution 100:179–196.

Smith, V. H., S. B. Joye, and R. W. Howarth. 2006. Eutrophication of freshwater and marine ecosystems. Limnology and Oceanography 51:351–355.

Strobl, C., A. Boulesteix, A. Zeileis, and T. Hothorn. 2007. Bias in random forest variable importance measures: illustrations, sources and a solution. BMC Bioinformatics 8:25.

Tilzer, M. M. 1988. Secchi disk – chlorophyll relationships in a lake with highly variable phytoplankton biomass. Hydrobiologia 162:163–171.

USEPA. 2006. Wadeable streams assessment: a collaborative survey of the nation's streams. EPA 841-b-06-002. Office of Water; Office of Research; Development, United States Environmental Protection Agency, Washington, D.C.

USEPA. 2009. National lakes assessment: a collaborative survey of the nation's lakes. EPA 841-r-09-001. Office of Water; Office of Research; Development, United States Environmental Protection Agency Washington, D.C.

Xian, G., C. Homer, and J. Fry. 2009. Updating the 2001 national land cover database land cover classification to 2006 by using Landsat imagery change detection methods. Remote Sensing of Environment 113:1133–1147.

## Supporting Information

Additional supporting information may be found online at: http://onlinelibrary.wiley.com/doi/10.1002/ecs2.1321/supinfo

# Author Query Form

Dear Author,

During the copy-editing of your paper, the following queries arose. Please respond to these by marking up your proofs with the necessary changes/additions. Please write your answers on the query sheet if there is insufficient space on the page proofs. Please write clearly and follow the conventions shown on the attached corrections sheet. If returning the proof by fax do not write too close to the paper's edge. Please remember that illegible mark-ups may delay publication.

Many thanks for your assistance.

| Query reference | Query | Remarks |
|---|---|---|
| 1 | **WILEY: Please supply date of revision.** | |
| 2 | **AUTHOR: Please confirm that given names (red) and surnames/family names (green) have been identified correctly.** | |
| 3 | **AUTHOR: Please define GIS.** | |
| 4 | **AUTHOR: Please check the hierarchy of heading levels.** | |
| 5 | **AUTHOR: Carpenter 1999 has been changed to Carpenter et al. 1999 so that this citation matches the Reference List. Please confirm that this is correct.** | |
| 6 | **AUTHOR: Strobl (2007) has not been included in the Reference List, please supply full publication details.** | |
| 7 | **AUTHOR: Please define RMSE.** | |
| 8 | **AUTHOR: Please provide the publisher location for reference Jones and Brett (2014).** | |
| 9 | **AUTHOR: Please provide the journal title for reference Paul and McDonald (2005).** | |
| 10 | **AUTHOR: Please provide the editor for reference Rodhe (1969).** | |

# MARKED PROOF

## Please correct and return this set

Please use the proof correction marks shown below for all alterations and corrections. If you wish to return your proof by fax you should ensure that all amendments are written clearly in dark ink and are made well within the page margins.

| Instruction to printer | Textual mark | Marginal mark |
|---|---|---|
| Leave unchanged | • • • under matter to remain | (√) |
| Insert in text the matter indicated in the margin | ⅄ | New matter followed by ⅄ or ⅄② |
| Delete | / through single character, rule or underline _or_ ⊢——⊣ through all characters to be deleted | ⌀ or ⌀② |
| Substitute character or substitute part of one or more word(s) | / through letter  or ⊢——⊣ through characters | new character / or new characters / |
| Change to italics | — under matter to be changed | ⌣ |
| Change to capitals | ≡ under matter to be changed | ≡ |
| Change to small capitals | = under matter to be changed | = |
| Change to bold type | ∿ under matter to be changed | ∿ |
| Change to bold italic | ≈ under matter to be changed | ≈ |
| Change to lower case | Encircle matter to be changed | ≢ |
| Change italic to upright type | (As above) | ⅄ |
| Change bold to non-bold type | (As above) | ⅄ |
| Insert 'superior' character | / through character   or ⅄ where required | ⅄ or ⅄ under character e.g. ⅄ or ⅄ |
| Insert 'inferior' character | (As above) | ⅄ over character e.g. ⅄ |
| Insert full stop | (As above) | ⊙ |
| Insert comma | (As above) | , |
| Insert single quotation marks | (As above) | ⅄ or ⅄ and/or ⅄ or ⅄ |
| Insert double quotation marks | (As above) | ⅄ or ⅄ and/or ⅄ or ⅄ |
| Insert hyphen | (As above) | ⊢⊣ |
| Start new paragraph | ⌐ | ⌐ |
| No new paragraph | ↪ | ↪ |
| Transpose | ⊔⊓ | ⊔⊓ |
| Close up | linking ⌢ characters | ⌢ |
| Insert or substitute space between characters or words | / through character   or ⅄ where required | ⅄ |
| Reduce space between characters or words | \| between characters or words affected | ↑ |