

# Modelling lake trophic state: A random forest approach

Jeffrey W. Hollister<sup>\*</sup> <sup>1</sup> W. Bryan Milstead <sup>1</sup> Betty J. Kreakie <sup>1</sup>

<sup>1</sup>US Environmental Protection Agency, Office of Research and Development, National Health and Environmental Effects Research Laboratory, Atlantic Ecology Division, 27 Tarzwell Drive Narragansett, RI, 02882, USA

<sup>\*</sup> corresponding author: [hollister.jeff@epa.gov](mailto:hollister.jeff@epa.gov)

## 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 aesthetics), 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 in order 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, other universally available data, and apply data mining approaches to predict trophic state. Using these data and random forests, we first model chlorophyll *a*, 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 accuracy of the trophic state classifications derived from the chlorophyll *a* predictions were 69% for the full model and 49% 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>.

**Keywords:** Harmful Algal Blooms; Cyanobacteria; Open Science; Nutrients; National Lakes Assessment

# 1 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

61 (Imboden and Gächter 1978, Salas and Martino 1991, Carvalho et al. 2011, Milstead et al. 2013).  
 62 This leads to accurate models, but these data are often sparse and not always available, thus  
 63 limiting the population of lakes for which we can make predictions. A possible solution for this  
 64 issue is to build models that use widely available data that are correlated to many of the *in situ*  
 65 variables. For instance, landscape metrics of forests, agriculture, wetlands, and urban land in  
 66 contributing watersheds have all been shown to explain a significant proportion of the variation  
 67 (ranging from 50-86%, depending on study) in nutrients in receiving waters (Jones et al. 2001,  
 68 2004, Seilheimer et al. 2013). Building on these previously identified associations might allow us  
 69 to use only landscape and other universally available data to build models. Identifying predictors  
 70 using this type of ubiquitous data would allow for estimating trophic state in both monitored  
 71 and unmonitored lakes.

72 Many published models of nutrients and trophic state in freshwater systems are based on linear  
 73 modelling methods such as standard least squares regression or linear mixed models (Jones et  
 74 al. 2001, 2004). While these methods have proven to be reliable, they have limitations (e.g.,  
 75 independence, distribution assumptions, and outlier sensitivity). Using data mining approaches,  
 76 such as random forests, avoids many of the limitations, may reduce bias, and often provides  
 77 better predictions (Breiman 2001, Cutler et al. 2007, Peters et al. 2007, Fernández-Delgado et  
 78 al. 2014). For instance, random forests are non-parametric and thus the data do not need to  
 79 come from a specific distribution (e.g., Gaussian) and can contain collinear variables (Cutler et  
 80 al. 2007). Second, random forests work well with very large numbers of predictors (Cutler et al.  
 81 2007). Lastly, random forests can deal with model selection uncertainty as predictions are based  
 82 upon a consensus of many models and not just a single model selected with some measure of  
 83 goodness of fit.

84 The research presented here builds on past work in three areas. First, we built, assessed,  
 85 and compared two random forest models of chlorophyll *a* with 1) *in situ* and universally  
 86 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. Lastly, to promote transparency in our work, the analysis code and data are available as an R package from <https://github.com/USEPA/LakeTrophicModelling>.

## 2 Methods

### 2.1 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 (Figure 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 National Lakes Assessment (USEPA 2009). Additionally, the NLA included ecological regions as defined in the Wadeable Streams Assessment (Figure 2) (Omernik 1987, USEPA 2006).

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 kilometer buffer surrounding each lake (Homer et al. 2004, Xian et al. 2009). A three kilometer buffer was selected to represent 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 (i.e., 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.

## 2.2 Predicting Trophic State with Random Forests

Random forest is a machine learning algorithm that aggregates numerous decision trees in order 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 are 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 gene selection 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 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  $R^2$  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.

## 2.3 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 the full suite of variables tested is in Appendix 1. 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. Lastly, we used only complete cases (i.e., missing data were removed) so the total number of observations varied among models.

Our modelling work flow was as follows:

1. Identify a minimal set of variables 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 for 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.

## 2.4 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  $R^2$  of the models. Second, we examined the accuracy of the model predictions when converted to trophic states classes via a confusion matrix (Table 1). 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 greater than 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. Additionally, random forest builds each tree on bootstrapped, random subsets of the original data, 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), thus, we used percent increase in mean squared error to assess variable importance. Lastly, 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.

## 2.5 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 5,000 estimates of chlorophyll *a* for each lake. We convert these values to trophic states (Table 1) then count up 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



205 provides the predicted class, in this case oligotrophic, and suggests little uncertainty in this  
206 prediction. We refer to this value as the “prediction probability.”

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

## 212 3 Results

213 Our complete dataset included 1148 lakes; however 5 lakes did not have chlorophyll *a* data. Thus,  
214 the base dataset for our modelling was conducted on data for 1143 lakes. The lakes were well  
215 distributed across the four trophic state categories (Table 1) and spatially throughout the United  
216 States (Figure 1).

### 217 3.1 Models: All Variables

218 The model built with all predictors used 1080 total observations, had a mean squared error  
219 of 0.09 and  $R^2$  of 0.8. The accuracy of the four trophic states was 68.7% and the kappa  
220 coefficient was 0.57 (Table 2). The variable selection process identified a reduced model with 20  
221 variables (Figure 3). The six most important variables were turbidity, total phosphorus, total  
222 nitrogen, elevation, total organic carbon, and N:P ratio (Figures 4). The role that each played in  
223 predicting chlorophyll *a* varied (Figure 5).

## 3.2 Models: GIS Only Variables

The GIS only model was built using 1138 total observations, had a mean squared error of 0.22 and  $R^2$  0.48. Four trophic states were predicted with a total accuracy of 49% and had a kappa coefficient of 0.29 (Table 3). The variable selection process for this model produced a reduced model with 15 variables (Figure 6). The six most important variables were ecoregion, percent cropland, elevation, latitude, percent evergreen forest, and mean lake depth (Figures 7 & 5).

## 3.3 Trophic State Probabilities

The “All variables” model provides more certain model predictions with a median prediction probability of 0.81 versus 0.72 for the “GIS only” model (Figure 9). Additionally, total accuracy of the predictions is a function of this uncertainty. Lakes with more certain predictions were more accurately classified (Figure 10). For both models, when prediction probabilities are approximately 0.8 or higher, the models had an accuracy of ~100%. This represents 55% of the lakes for the “All variables” model and 22% of the lakes for the “GIS only” model. Lastly, as prediction probabilities increased, the difference in total accuracy between the two models decreased (Figure 10 & Table 4).

# 4 Discussion

## 4.1 Trophic State Probabilities

Not surprisingly, lakes with more certain predictions (i.e., higher prediction probabilities) were more accurately predicted (Figure 10). The fact that the difference in accuracy 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 (Table 4). Additionally, the prediction probabilities may be mapped for each of the four classes (Figure 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 (Figure 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. Lastly, this particular map is very similar to simply mapping the raw data. However, it highlights what could 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

## 4.2 Partial dependencies of explanatory variables

In line with past predictive modelling 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 (Figure 5). This relationship holds until nutrient concentrations become saturated. The partial dependency plots (Figure 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 modelling 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 (Figures 5 & 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 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 (see Figures 2 & 8) has a pronounced affect indicting continental scale effects of land use and geography. 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), 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 (Figure 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 (Figure 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.

## 5 Conclusions

Our research goals were to explore the utility of a widely used data mining algorithm, random forests, in the modelling 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 RMSE of 0.09 and an adjusted  $R^2$  of 0.8 whereas, the GIS only models had an RMSE of 0.22 and the adjusted

317  $R^2$  was 0.48. Our total accuracy in predicting chlorophyll *a* based trophic states was 69% for the  
318 “All variables” model and 49% for the “GIS only” model.

319 While the “GIS only” model showed lower prediction accuracies than the “All variables” model,  
320 the association between the uncertainty of prediction and total accuracy (Figure 10 and Table 4)  
321 suggest that the “GIS only” model will provide reasonable estimates of trophic state for many  
322 lakes across the United States. Furthermore, we can map the uncertainty of the predictions,  
323 thus, we know the spatial patterns and location of the lakes for which we are certain, or not, of  
324 their predicted trophic state. Given this and that these models may be applied to any lake in  
325 the United States we can recommend using this model. Future iterations of this modelling effort  
326 may be able to utilize modeled predictions of nutrients to improve accuracy and also maintain  
327 broad applicability (Milstead et al. 2013).

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

333 A potentially useful benefit of models of trophic state and chlorophyll *a* are their use in assessing  
334 risk due to cyanobacteria. Cyanobacteria biomass should be closely associated with chlorophyll *a*  
335 and trophic state as cyanobacteria contribute to the chlorophyll concentration in a lake. If these  
336 associations are strong enough we may be able to expand models such as those reported here to  
337 also predict probability of cyanobacteria blooms and other indices related to cyanobacteria (e.g.,  
338 toxin presence). Others have seen these associations. For instance, Kasinak et al. (2015) used  
339 bench top fluorimeters and showed a strong correlation between chlorophyll *a* and phycocyanin.  
340 Using the NLA data, we see a positive trend between chlorophyll *a* and cyanobacteria abundance  
341 (Figure 12). Both of these suggest that trophic state may be an acceptable proxy for cyanobacteria.

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. Second, while direct control of water quality in lakes would have a large impact, the land use/land cover drivers (i.e., non-point 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 (i.e., 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.

## 6 Acknowledgements

We would like to thank Farnaz Nojavan, Nathan Schmucker, John Kiddon, Joe LiVolsi, Tim Gleason, and Wayne Munns for constructive reviews of this paper. This paper has not been subjected to Agency review. Therefore, it does not necessary reflect the views of the Agency. Mention of trade names or commercial products does not constitute endorsement or recommendation for use. This contribution is identified by the tracking number ORD-011075 of the Atlantic Ecology Division, Office of Research and Development, National Health and Environmental Effects Research Laboratory, US Environmental Protection Agency.

363 **7 Figures**

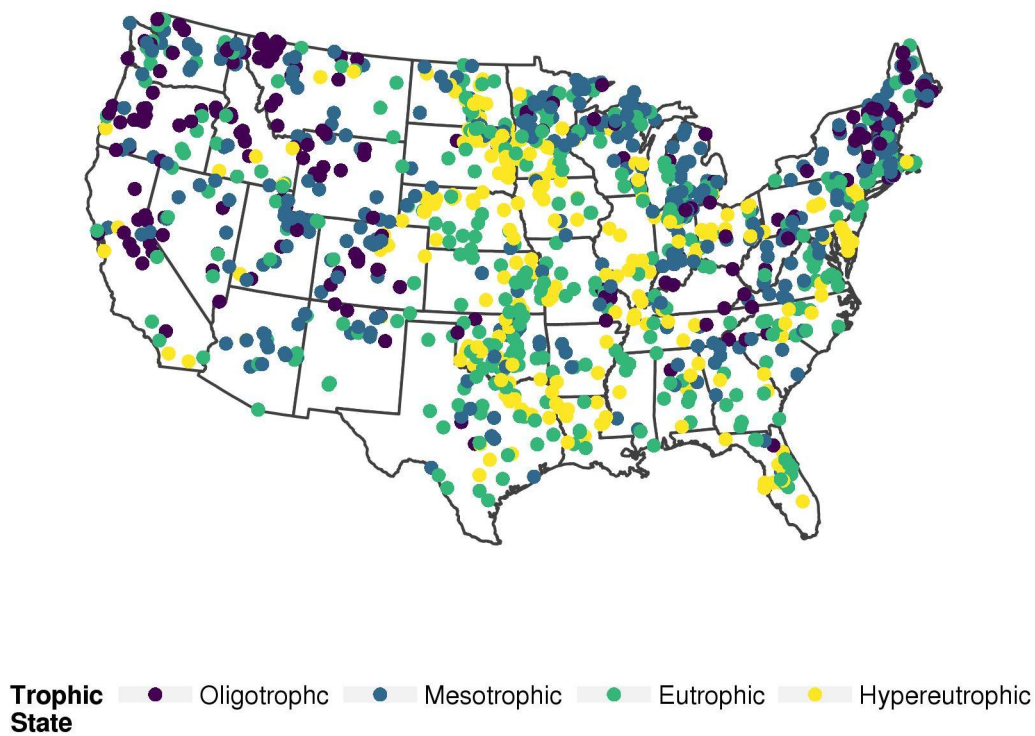


Figure 1: Map of the distribution of National Lakes Assessment Sampling locations



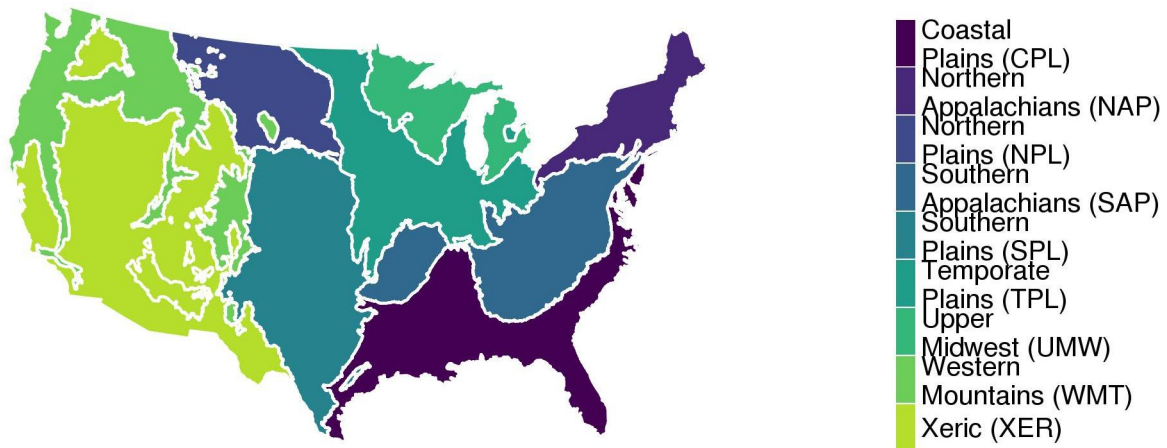


Figure 2: Wadeable Streams Assesment ecoregions

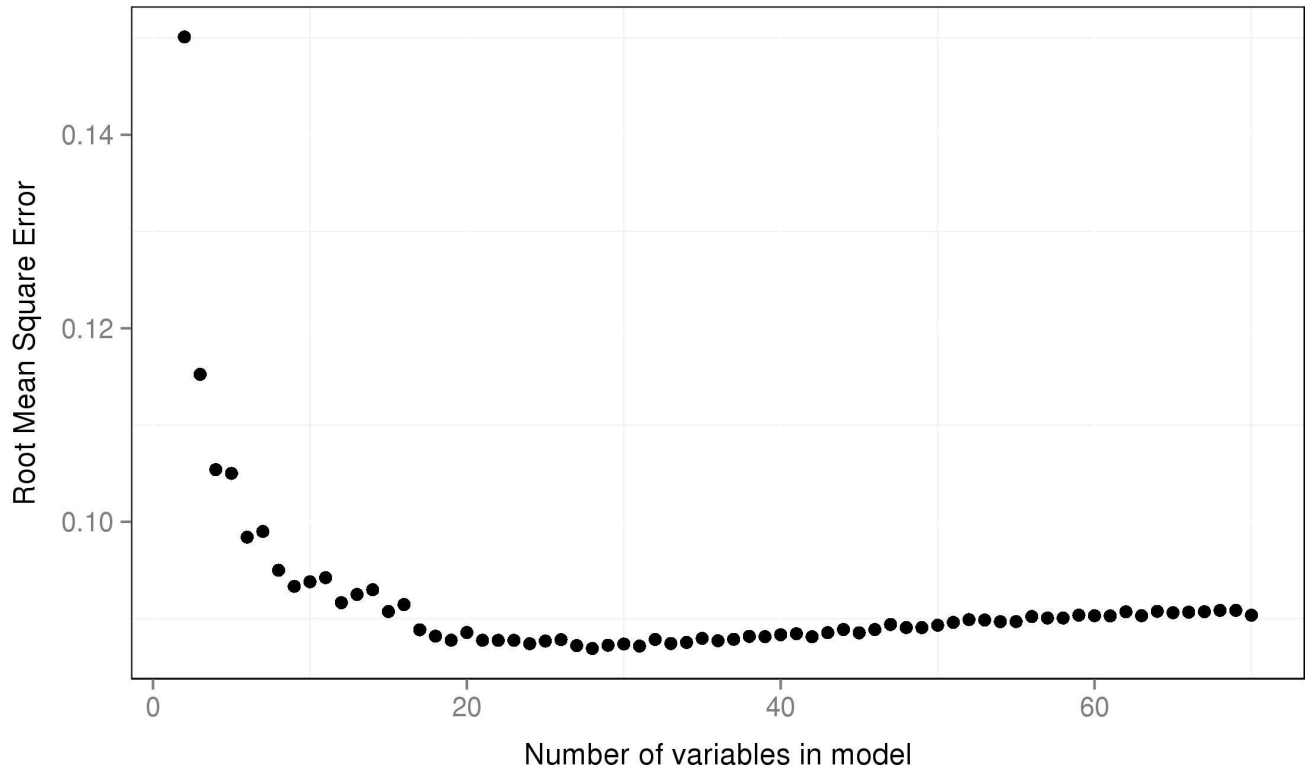


Figure 3: Variable selection plot for all variables. Shows percent increase in mean squared error as a function of the number of variables.



Figure 4: Importance plot for All Variables., shows percent increase in mean square error. Higher values of percent increase in mean squared error indicates higher importance.

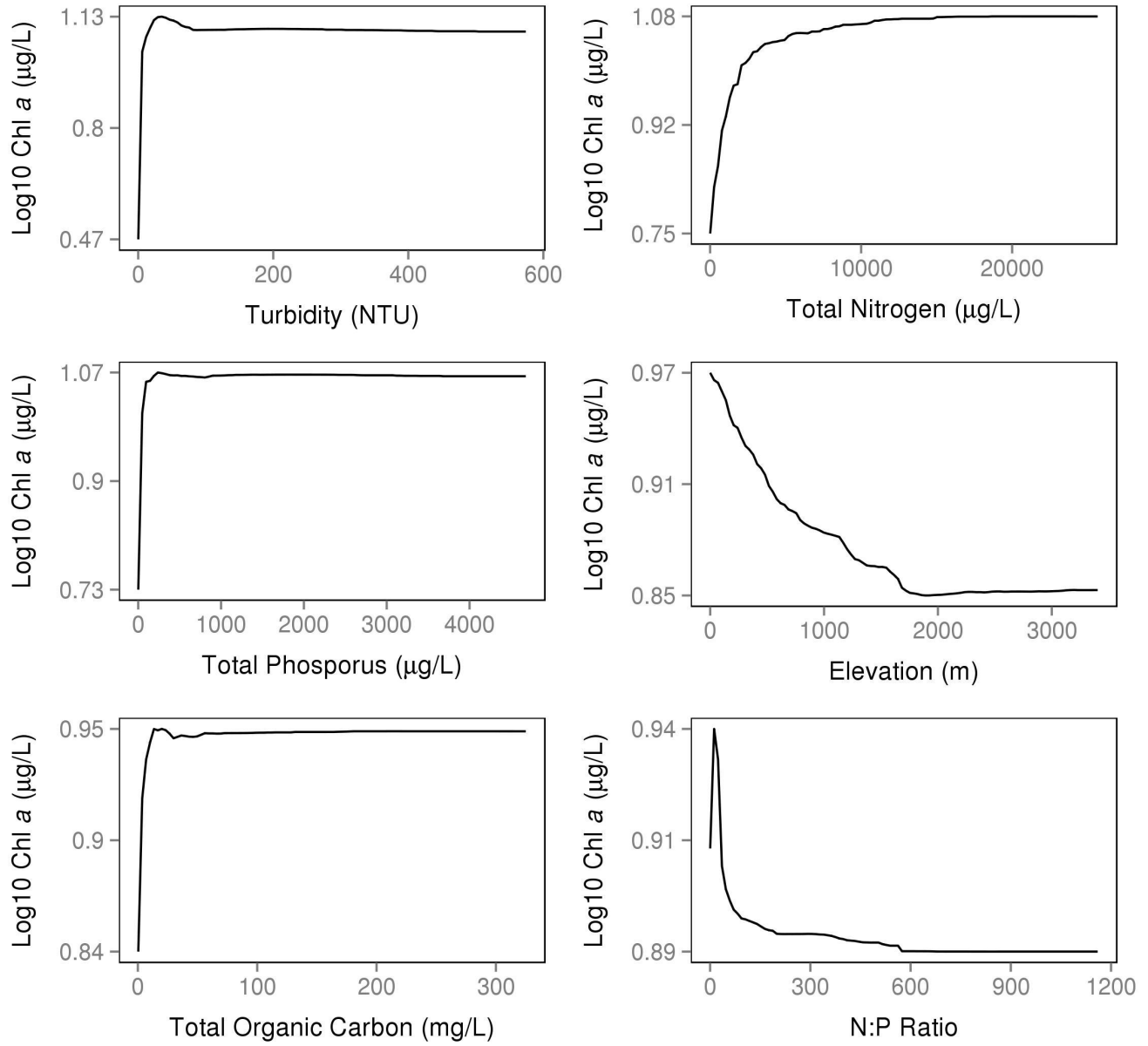


Figure 5: All Variables partial dependence plots for the top 5 most important variables.

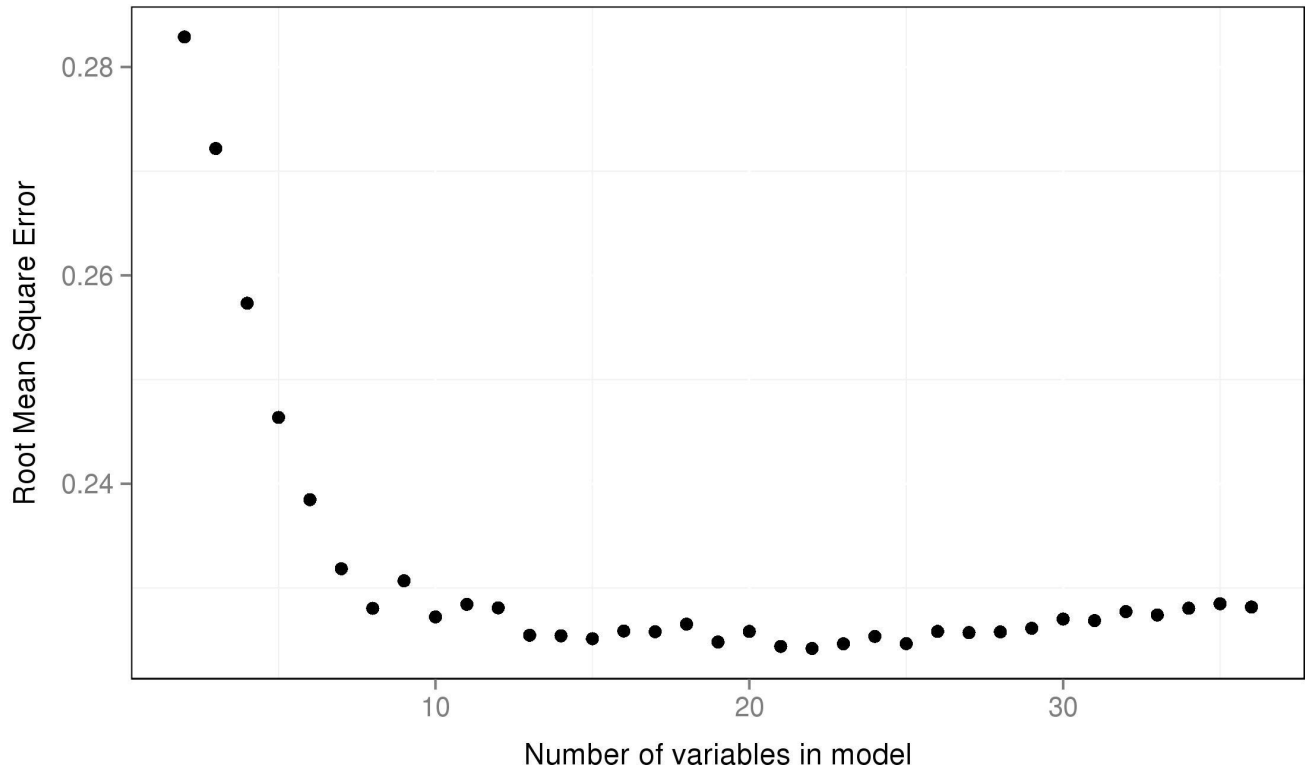


Figure 6: Variable selection plot for GIS only variables. Shows percent increase in mean squared error as a function of the number of variables.

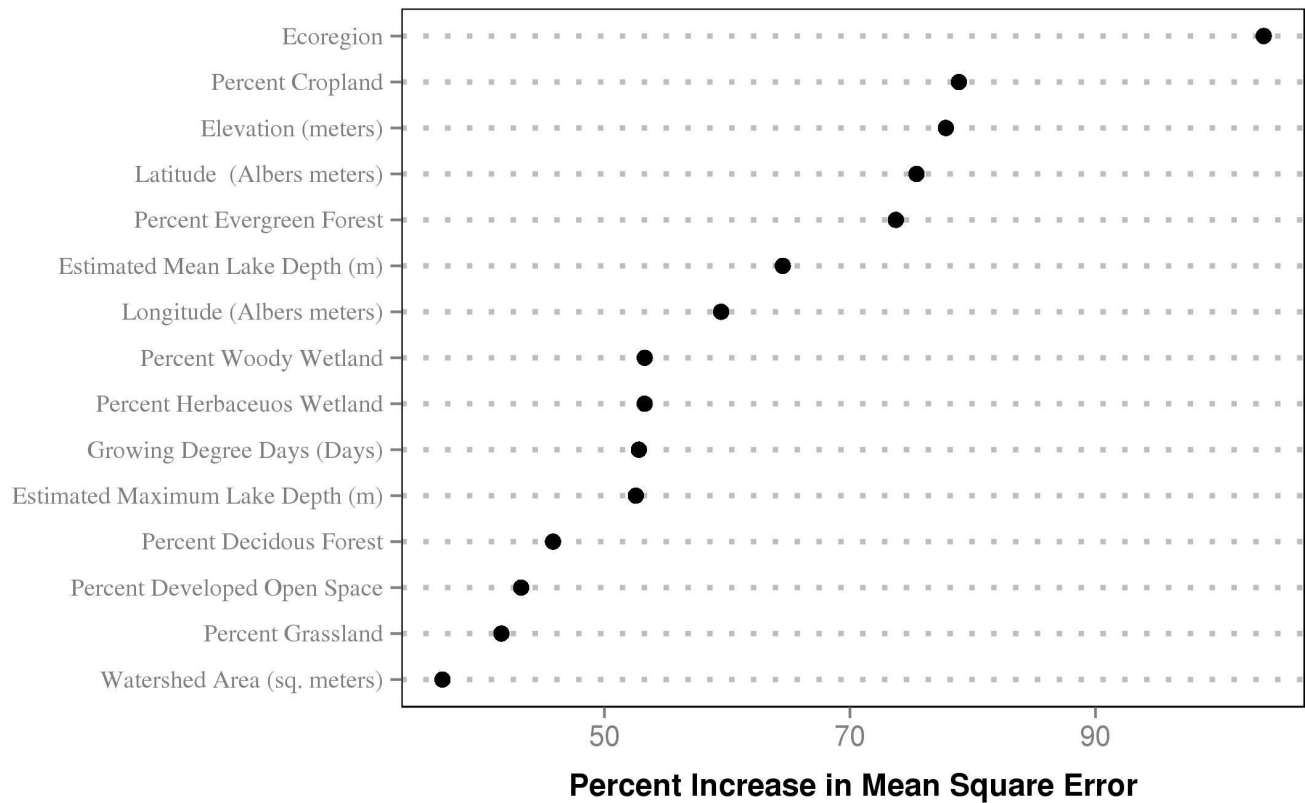


Figure 7: Importance plot for GIS Only Variables., shows percent increase in mean square error. Higher values of percent increase in mean squared error indicates higher importance.

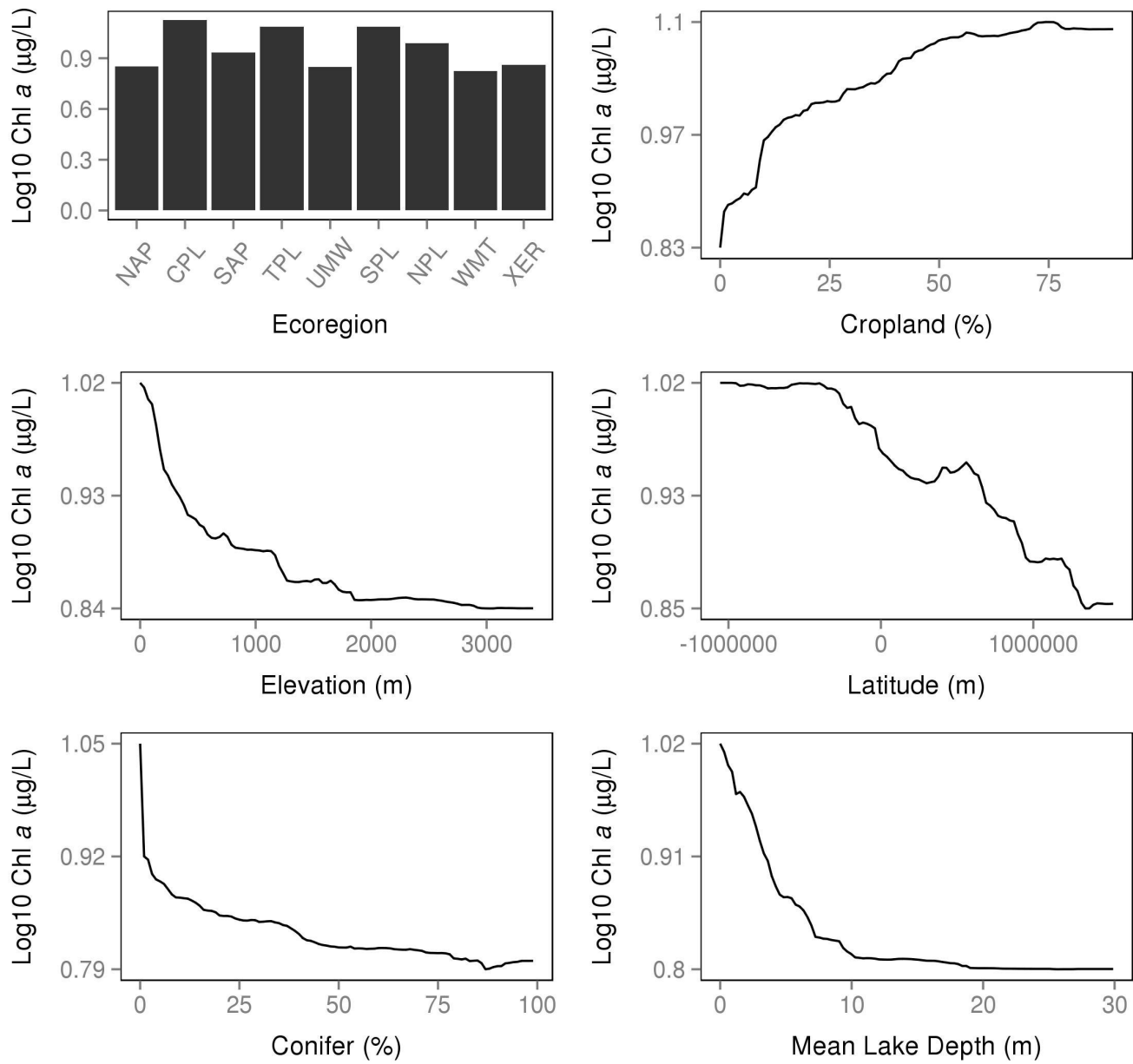


Figure 8: GIS Only Variables partial dependence plots for the top 5 most important variables.

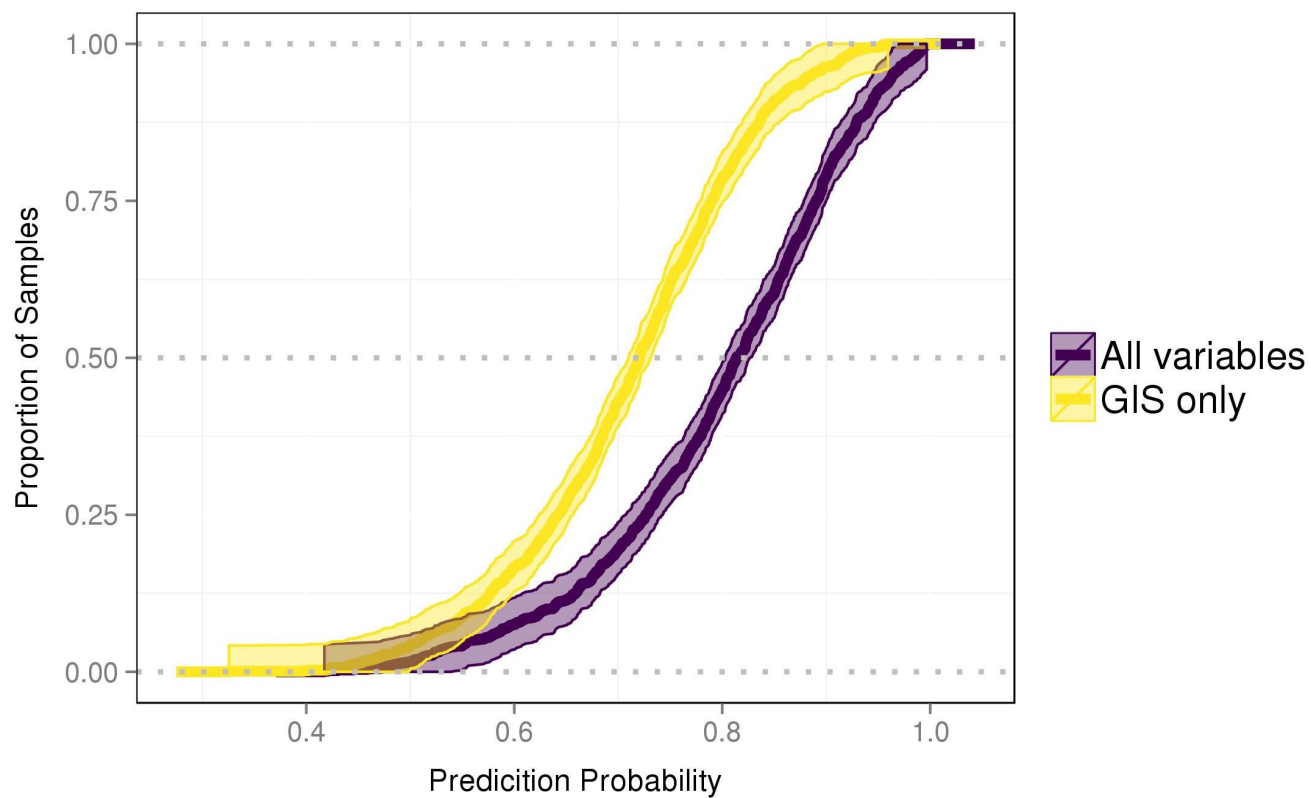


Figure 9: Prediction probabilities for the All Variables and GIS Only models.



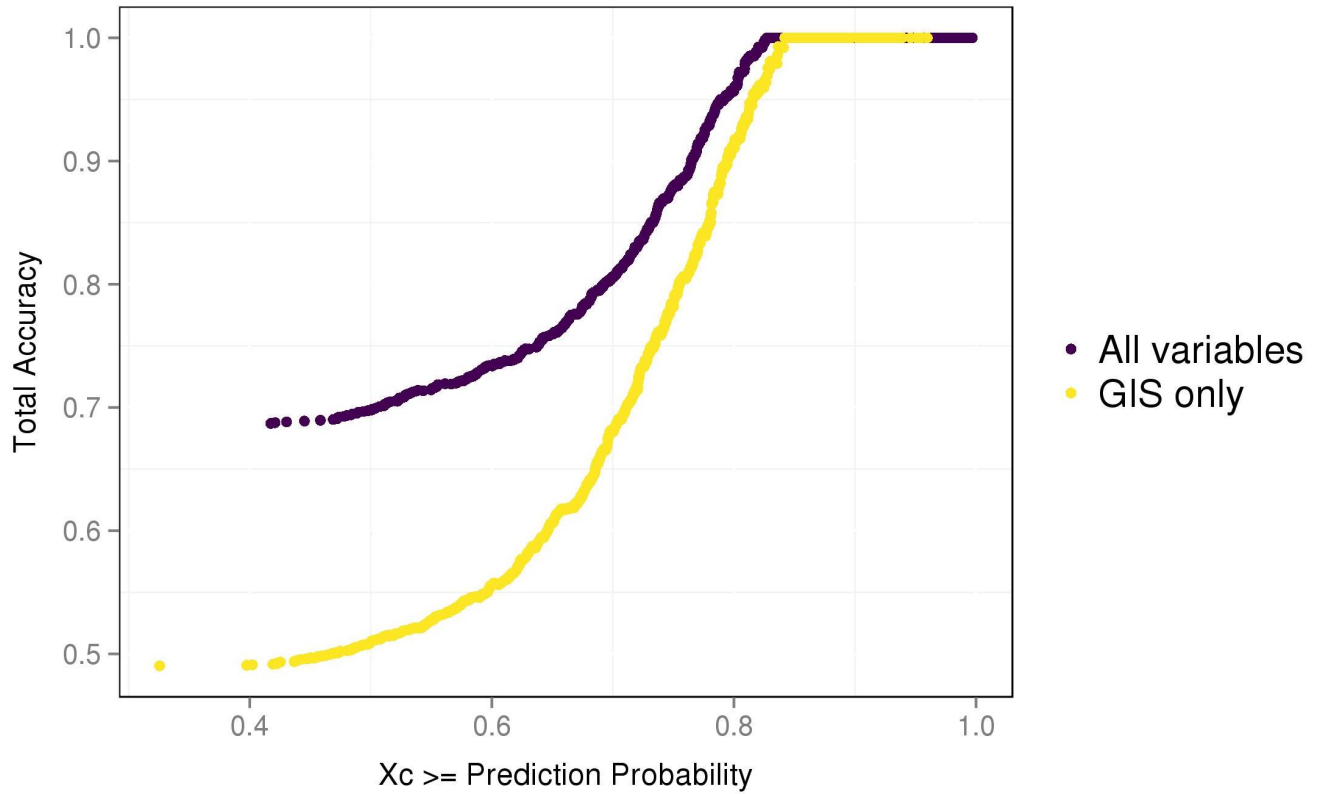


Figure 10: Accuracy of predictions as a function of lake prediction probability. The x-axis represents lakes with a prediction probability at a given level or higher.

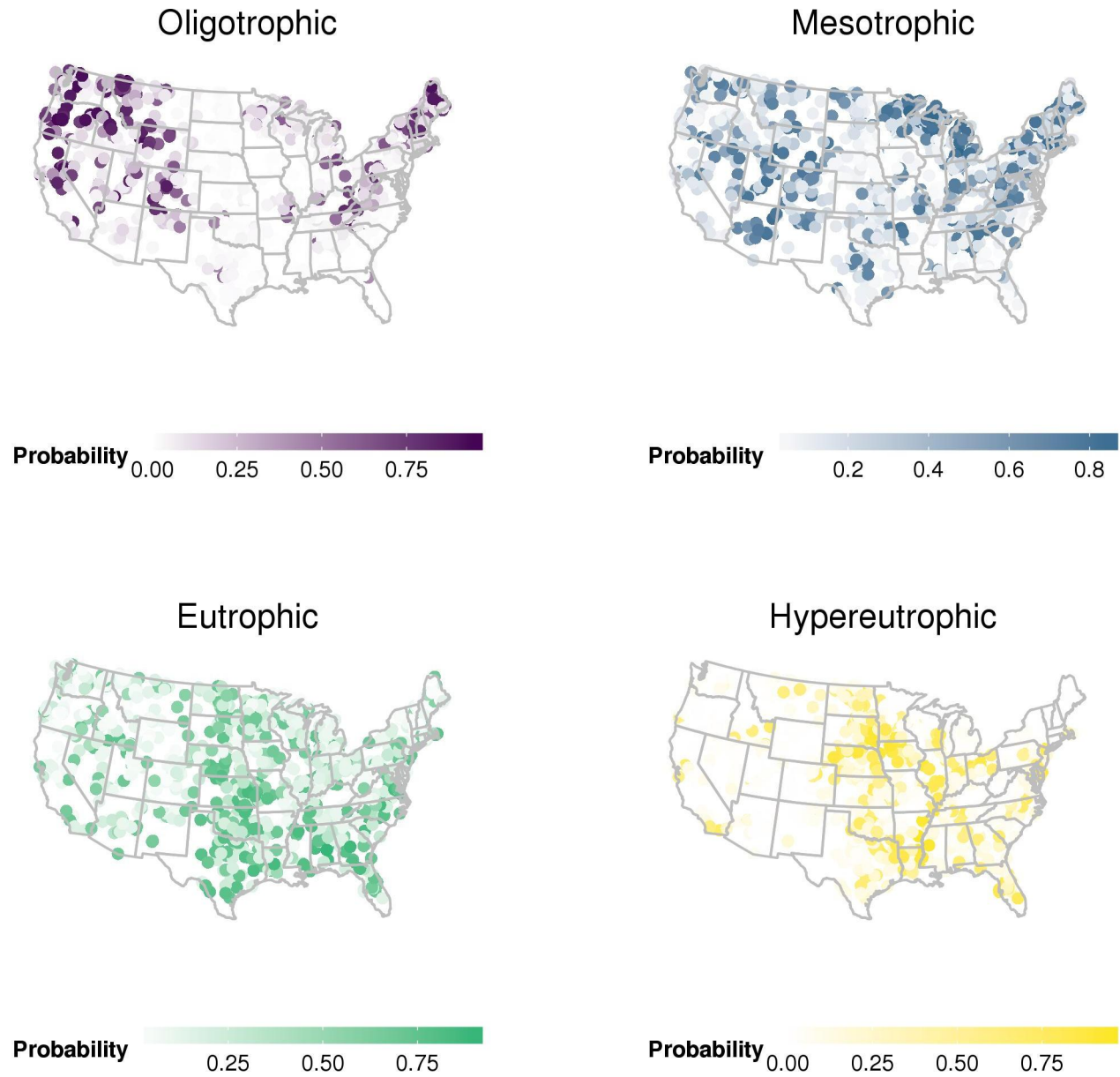


Figure 11: Maps of prediction probabilities for each of the four chlorophyll *a* trophic states

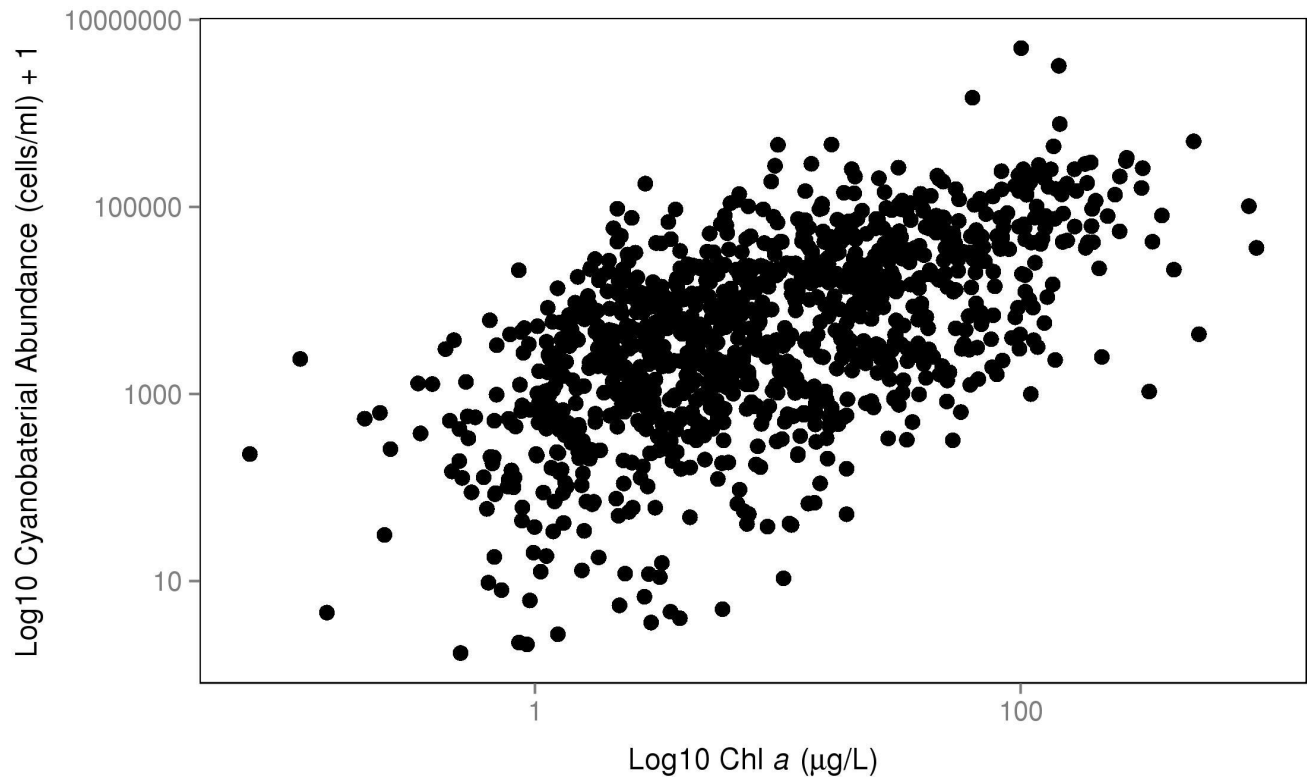


Figure 12: Cholorphyll *a* and cyanobacteria abundance scatterplot

364 **8 Tables**

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

Trophic State (4 class)	Trophic State (2 class)	$\mu\text{g/L}$ Cut-off
oligotrophic	oligotrophic/mesotrophic	$\leq 2$
mesotrophic	oligotrophic/mesotrophic	$>2-7$
eutrophic	eutrophic/hypereutrophic	$>7-30$
hypereutrophic	eutrophic/hypereutrophic	$>30$

Table 2: Random Forest confusion matrix for All Variables model converted to 4 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.

	oligo	meso	eu	hyper	Class Accuracy (%)
oligo	115	31	0	0	78.77
meso	67	251	63	0	65.88
eu	7	61	217	75	60.28
hyper	0	5	29	159	82.38

Table 3: Random Forest confusion matrix for GIS Only model converted to 4 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.

	oligo	meso	eu	hyper	Class Accuracy (%)
oligo	65	14	6	0	76.47
meso	101	213	98	18	49.53
eu	29	126	193	141	39.47
hyper	1	8	38	87	64.93

Table 4: Summary of relationship between prediction probabilities, total accuracy, and number of lakes.

Prediction Prob.	“All Var.”			“GIS Only”		
	Total Accuracy	Percent of Sample	Number of Samples	Total Accuracy	Percent of Sample	Number of Samples
All	69	100	846	49	100	878
0.50	70	98	829	51	95	834
0.60	73	91	770	56	81	711
0.70	81	77	654	68	56	490
0.80	96	51	434	91	24	212
0.90	100	20	173	100	5	41

## 9 Appendix 1. Variable Definitions

Variable Names	Description	Source	Mean	Std. Error
AlbersX	Longitude (Albers meters)	GIS	126757.1	34305.5
AlbersY	Latitude (Albers meters)	GIS	436908.1	17367.2
BASINAREA	Watershed Area (sq. meters)	GIS	3208.5	788.1
BarrenPer_3000m	% Barren	GIS	0.7	0.1
CropsPer_3000m	% Cropland	GIS	13.3	0.6
DDs45	Growing Degree Days (Days)	GIS	2750.0	41.0
DeciduousPer_3000m	% Deciduous Forest	GIS	17.1	0.6
DevHighPer_3000m	% High Intensity Development	GIS	0.4	0.0
DevLowPer_3000m	% Low Intensity Development	GIS	3.0	0.2
DevMedPer_3000m	% Medium Intensity Development	GIS	1.4	0.1
DevOpenPer_3000m	% Developed Open Space	GIS	5.4	0.2
ELEV_PT	Elevation (meters)	GIS	607.6	20.1
EvergreenPer_3000m	% Evergreen Forest	GIS	12.2	0.6
FetchE	Fetch from East (m)	GIS	1652.8	80.3
FetchN	Fetch from North (m)	GIS	2009.6	106.9
FetchNE	Fetch from Northeast (m)	GIS	1645.0	80.9
FetchSE	Fetch from Southeast (m)	GIS	1642.0	80.5
GrassPer_3000m	% Grassland	GIS	13.8	0.7
HerbWetPer_3000m	% Herbaceous Wetland	GIS	1.7	0.1
IceSnowPer_3000m	% Ice/Snow	GIS	0.0	0.0
LakeArea	Lake Surface Area (sq. meters)	GIS	12.2	2.3
LakePerim	Lake Perimeter (meters)	GIS	33.6	4.5
MaxDepthCorrect	Est. Maximum Lake Depth (m)	GIS	8.4	0.3
MaxLength	Maximum Lake Length (m)	GIS	2972.1	137.2



Variable Names	Description	Source	Mean	Std. Error
MaxWidth	Maximum Lake Width (m)	GIS	1567.5	76.0
MeanDepthCorrect	Est. Mean Lake Depth (m)	GIS	2.9	0.1
MeanWidth	Mean Lake Width (m)	GIS	1370.1	122.6
MixedForPer_3000m	% Mixed Forest	GIS	3.8	0.3
PasturePer_3000m	% Pasture	GIS	7.7	0.3
PercentImperv_3000m	% Impervious	GIS	2.6	0.2
ShoreDevel	Shoreline Development Index	GIS	2.7	0.1
ShrubPer_3000m	% Shrub/Scrub	GIS	10.4	0.6
VolumeCorrect	Est. Lake Volume (cubic meters)	GIS	101211909.9	27438696.4
WSA_ECO9	Ecoregion	GIS	NA	NA
WaterPer_3000m	% Water	GIS	4.1	0.2
WoodyWetPer_3000m	% Woody Wetland	GIS	5.2	0.3
ANC	Acid Neutralizing Capacity (ueq/L)	NLA	2584.2	171.7
ANDEF2	Anion Deficit (ueq/L)	NLA	-506.4	143.2
ANSUM2	Sum of Anions using ANC (ueq/L)	NLA	8043.1	1197.9
BALANCE2	Ion Balance (%)	NLA	-0.7	0.1
CA	Calcium (ueq/L)	NLA	1388.3	54.0
CATSUM	Sum of Cations (ueq/L)	NLA	7536.7	1105.0
CL	Chloride (ueq/L)	NLA	1600.3	438.2
COLOR	Color (PCU)	NLA	16.1	0.5
CONCAL2	Calculated Conductivity (uS/cm)	NLA	949.0	148.1
COND	Conductivity (uS/cm)	NLA	656.0	72.6
CONDHO2	D-H-O Calculated Conductivity (uS/cm)	NLA	618.6	55.1
DATE_COL	Date Samples Collected	NLA	NA	NA
DEPTHMAX	Maximum Depth (meters)	NLA	9.6	0.3

Variable Names	Description	Source	Mean	Std. Error
DO2_2M	Dissolved Oxygen (mg/L)	NLA	7.9	0.1
DOC	Dissolved Organic Carbon (mg/L)	NLA	8.6	0.5
H	Hydrogen Ions (ueq/L)	NLA	0.2	0.1
K	Potassium (ueq/L)	NLA	245.6	40.6
MG	Magnesium (ueq/L)	NLA	2190.4	282.2
NH4	Ammonium (mg/L)	NLA	2.9	0.2
NH4ION	Calculated Ammonium (ueq/L)	NLA	2.5	0.2
NO3	Nitrate (ueq/L)	NLA	5.4	0.7
NO3_NO2	Nitrate/Nitrite (mg N/L)	NLA	0.1	0.0
NPratio	Nitrogen:Phophorus Ratio	NLA	34.5	1.8
NTL	Total Nitrogen ( $\mu$ g/L)	NLA	1109.9	56.4
Na	Sodium (ueq/L)	NLA	3709.7	816.3
OH	Hydroxide (ueq/L)	NLA	3.1	0.2
ORGION	Est. Organic Anions (ueq/L)	NLA	85.9	4.8
PH_FIELD	pH	NLA	8.1	0.0
PTL	Total Phosphorus ( $\mu$ g/L)	NLA	103.1	7.8
SIO2	Silica (mg/L)	NLA	8.6	0.3
SO4	Sulfate (ueq/L)	NLA	3853.4	935.7
SOBC	Sum of Base Cation (ueq/L)	NLA	7534.1	1105.0
TOC	Total Organic Carbon (mg/L)	NLA	9.6	0.6
TURB	Turbidity (NTU)	NLA	12.3	1.0
TmeanW	Mean Profile Water Temp. (C)	NLA	24.1	0.1

```

366 ## R version 3.2.1 (2015-06-18)
367 ## Platform: x86_64-redhat-linux-gnu (64-bit)
368 ## Running under: Red Hat Enterprise Linux Server release 6.7 (Santiago)
369 ##
370 ## locale:
371 ##   [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
372 ##   [3] LC_TIME=en_US.UTF-8      LC_COLLATE=en_US.UTF-8
373 ##   [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
374 ##   [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
375 ##   [9] LC_ADDRESS=C             LC_TELEPHONE=C
376 ##  [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
377 ##
378 ## attached base packages:
379 ##   [1] parallel  stats4    grid      stats     graphics  grDevices  utils
380 ##   [8] datasets  methods   base
381 ##
382 ## other attached packages:
383 ##   [1] condprob2_2.0             viridis_0.1
384 ##   [3] maptools_0.8-36          sfsmisc_1.0-27
385 ##   [5] mapproj_1.2-3            maps_2.3-10
386 ##   [7] rmarkdown_0.6.1          caret_6.0-52
387 ##   [9] lattice_0.20-31          dplyr_0.4.2
388 ##  [11] e1071_1.6-4              rgdal_1.0-4
389 ##  [13] sp_1.1-1                 ggplot2_1.0.1
390 ##  [15] knitr_1.10.5             doParallel_1.0.8
391 ##  [17] iterators_1.0.7          foreach_1.4.2
392 ##  [19] interpretR_0.2.3         randomForest_4.6-11

```

```

393 ## [21] broom_0.3.7           party_1.0-21
394 ## [23] strucchange_1.5-1     sandwich_2.3-3
395 ## [25] zoo_1.7-12           modeltools_0.2-21
396 ## [27] mvtnorm_1.0-2         tidyr_0.2.0
397 ## [29] pander_0.5.2         edarf_0.1
398 ## [31] wesanderson_0.3.2     LakeTrophicModelling_0.1
399 ##
400 ## loaded via a namespace (and not attached):
401 ## [1] splines_3.2.1         AUC_0.3.0             gtools_3.5.0
402 ## [4] assertthat_0.1       highr_0.5             coin_1.0-24
403 ## [7] yaml_2.1.13          quantreg_5.11         digest_0.6.8
404 ## [10] minqa_1.2.4          colorspace_1.2-6     htmltools_0.2.6
405 ## [13] Matrix_1.2-1         plyr_1.8.3            psych_1.5.4
406 ## [16] varSelRF_0.7-5       BradleyTerry2_1.0-6   SparseM_1.6
407 ## [19] scales_0.2.5         brglm_0.5-9          lme4_1.1-8
408 ## [22] mgcv_1.8-6           car_2.0-25           lazyeval_0.1.10
409 ## [25] nnet_7.3-9           pbkrtest_0.4-2        mnormt_1.5-3
410 ## [28] proto_0.3-10         survival_2.38-1      magrittr_1.5
411 ## [31] evaluate_0.7         nlme_3.1-120         MASS_7.3-40
412 ## [34] foreign_0.8-63       class_7.3-12         tools_3.2.1
413 ## [37] formatR_1.2          stringr_1.0.0        munsell_0.4.2
414 ## [40] snowfall_1.84-6      nloptr_1.0.4         labeling_0.3
415 ## [43] gtable_0.1.2         codetools_0.2-11     DBI_0.3.1
416 ## [46] reshape2_1.4.1       R6_2.0.1             rgeos_0.3-11
417 ## [49] stringi_0.5-5        Rcpp_0.11.6

```

## References

- 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.
- 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.
- Cohen, J. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* 20:37–46.
- Cutler, D. R., T. C. Edwards Jr, 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.

- 439 Genkai-Kato, M., and S. R. Carpenter. 2005. Eutrophication due to phosphorus recycling in  
440 relation to lake morphometry, temperature, and macrophytes. *Ecology* 86:210–219.
- 441 Hansson, L.-A. 1992. Factors regulating periphytic algal biomass. *Limnology and Oceanography*  
442 37:322–328.
- 443 Hasler, A. D. 1969. Cultural eutrophication is reversible. *BioScience* 19:425–431.
- 444 Hollister, J. W. 2014. Lakemorpho: Lake morphometry in R. R package version 1.0.  
445 <http://CRAN.R-project.org/package=lakemorpho>.
- 446 Hollister, J. W., W. B. Milstead, and M. A. Urrutia. 2011. Predicting maximum lake depth  
447 from surrounding topography. *PLoS ONE* 6:e25764.
- 448 Hollister, J. W., H. A. Walker, and J. F. Paul. 2008. CProb: A computational tool for conducting  
449 conditional probability analysis. *Journal of environmental quality* 37:2392–2396.
- 450 Hollister, J., and W. B. Milstead. 2010. Using GIS to estimate lake volume from limited data.  
451 *Lake and Reservoir Management* 26:194–199.
- 452 Homer, C., C. Huang, L. Yang, B. Wylie, and M. Coan. 2004. Development of a 2001 national  
453 land-cover database for the united states. *Photogrammetric Engineering & Remote Sensing*  
454 70:829–840.
- 455 Hubert, L., and P. Arabie. 1985. Comparing partitions. *Journal of classification* 2:193–218.
- 456 Imboden, D., and R. Gächter. 1978. A dynamic lake model for trophic state prediction. *Ecological*  
457 *modelling* 4:77–98.
- 458 Jones, J., M. Knowlton, D. Obrecht, and E. Cook. 2004. Importance of landscape variables  
459 and morphology on nutrients in missouri reservoirs. *Canadian Journal of Fisheries and Aquatic*  
460 *Sciences* 61:1503–1512.

- 461 Jones, K. B., A. C. Neale, M. S. Nash, R. D. Van Remortel, J. D. Wickham, K. H. Riitters,  
462 and R. V. O'Neill. 2001. Predicting nutrient and sediment loadings to streams from landscape  
463 metrics: A multiple watershed study from the united states mid-atlantic region. *Landscape*  
464 *Ecology* 16:301–312.
- 465 Jones, Z., and F. Linder. 2015. Exploratory data analysis using random forests. *in* The 73rd  
466 annual mPSA conference. MPSA.
- 467 Kasinak, J.-M. E., B. M. Holt, M. F. Chislock, and A. E. Wilson. 2015. Benchtop fluorometry of  
468 phycocyanin as a rapid approach for estimating cyanobacterial biovolume. *Journal of Plankton*  
469 *Research* 37:248–257.
- 470 Landis, J. R., and G. G. Koch. 1977. The measurement of observer agreement for categorical  
471 data. *biometrics* 33:159–174.
- 472 Liaw, A., and M. Wiener. 2002. Classification and regression by randomForest. *R News* 2:18–22.
- 473 Milstead, W. B., J. W. Hollister, R. B. Moore, and H. A. Walker. 2013. Estimating summer  
474 nutrient concentrations in northeastern lakes from SPARROW load predictions and modeled  
475 lake depth and volume. *PloS one* 8:e81457.
- 476 Omernik, J. M. 1987. Ecoregions of the conterminous united states. *Annals of the Association*  
477 *of American geographers* 77:118–125.
- 478 Paul, J. F., and M. E. McDonald. 2005. Development of empirical, geographically specific water  
479 quality criteria: A conditional probability analysis approach 41:1211–1223.
- 480 Peters, J., B. D. Baets, N. E. Verhoest, R. Samson, S. Degroeve, P. D. Becker, and W. Huybrechts.  
481 2007. Random forests as a tool for ecohydrological distribution modelling. *Ecological Modelling*  
482 207:304–318.
- 483 Read, E. K., V. P. Patil, S. K. Oliver, A. L. Hetherington, J. A. Brentrup, J. A. Zwart, K. M.

- Winters, J. R. Corman, E. R. Nodine, R. I. Woolway, and others. 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.
- 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. Page 334. University of Alberta Press Edmonton.
- 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* Successes, limitations, and frontiers in ecosystem science. Springer.
- Smith, V. H., and D. W. Schindler. 2009. Eutrophication science: Where do we go from here? *Trends in Ecology & Evolution* 24:201–207.
- Smith, V. H., S. B. Joye, R. W. Howarth, and others. 2006. Eutrophication of freshwater and marine ecosystems. *Limnology and Oceanography* 51:351–355.
- 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.
- Strobl, C., A.-L. 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



506 phytoplankton biomass. *Hydrobiologia* 162:163–171.

507 USEPA. 2006. Wadeable streams assessment: A collaborative survey of the nation's streams. ePA  
508 841-b-06-002. Office of Water; Office of Research; Development, US Environmental Protection  
509 Agency Washington, DC.

510 USEPA. 2009. National lakes assessment: A collaborative survey of the nation's lakes. ePA  
511 841-r-09-001. Office of Water; Office of Research; Development, US Environmental Protection  
512 Agency Washington, DC.

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