

Modeling Lake Trophic State: A Data Mining Approach

Jeffrey W. Hollister^{*} ¹ W. Bryan Milstead¹ Betty J. Kreakie¹

¹US 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

^{*} corresponding author: hollister.jeff@epa.gov

Abstract

Productivity of lentic ecosystems has been well studied and it is widely accepted that as nutrient inputs increase, productivity increases and lakes transition from low trophic state (e.g. oligotrophic) to higher trophic states (e.g. eutrophic). These broad trophic state classifications are good predictors of ecosystem health and ecosystem services and disservices (e.g. recreation, aesthetics, fisheries, and 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 expand our ability to predict trophic state in lakes without water quality data, we take advantage of the availability of a large national lakes water quality database, land use/land cover data, lake morphometry data, other universally available data, and modern data mining approaches to build and assess models of lake trophic state that may be more universally applied. We use random forests and random forest variable selection to identify variables to be used for predicting trophic state and we compare the performance of two sets of models of trophic state (as determined by chlorophyll *a* concentration). The first set of models estimates trophic state with *in situ* as well as universally available data and the second set of models uses universally available data only. For each of these models we used three separate trophic state categories, for a total of six models. Overall accuracy for models built from *in situ* and universal data ranged from 0.669% to 0.867%. For the universal data only models, overall accuracy ranged from 0.489% to 0.757%. Lastly, it is believed that the presence and abundance of cyanobacteria is strongly associated with trophic state. To test this we examine the association between estimates of cyanobacteria abundance and measured chlorophyll *a* and find a positive relationship. Expanding these preliminary results to include cyanobacteria taxa indicates that cyanobacteria are significantly more likely to be found in highly eutrophic lakes. These results suggest that predictive models of lake trophic state may be improved with additional information on the landscape surrounding lakes and that those models provide additional information on the presence of potentially harmful cyanobacteria taxa.

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

39 productivity (e.g. mesotrophic and eutrophic lakes) have greater nutrient loads, tend to be less clear,
40 have greater density of aquatic plants, and often support more diverse and abundant fish communities.
41 Higher primary productivity is not necessarily a predictor of poor ecological condition as it is natural
42 for lakes to shift from lower to higher trophic states but this is a slow process. However, at the highest
43 productivity levels (hypereutrophic lakes) biological integrity is compromised (Hasler 1969, Smith et al.
44 1999, Schindler and Vallentyne 2008).

45 Monitoring trophic state allows the identification of rapid shifts in trophic state or locating lakes with
46 unusually high productivity (e.g. hypereutrophic). These cases are indicative of lakes under greater
47 anthropogenic nutrient loads, also known as cultural eutrophication, and are more likely to be at risk of
48 fish kills, fouling, and harmful algal blooms (Smith 1998, Smith et al. 1999, 2006). Given the association
49 between trophic state and many ecosystem services and disservices, being able to accurately model
50 trophic state could provide a first cut at identifying lakes with the potential for harmful algal blooms or
51 other problems associated with cultural eutrophication.

52 As trophic state and related indices can be best defined by a number of *in situ* water quality parameters
53 (modeled or measured), most models have used this information as predictors (Imboden and Gächter
54 1978, Salas and Martino 1991, e.g., Carvalho et al. 2011, Milstead et al. 2013). This leads to accurate
55 models, but also requires data that are often sparse and not always available, thus limiting the population
56 of lakes for which we can make predictions. A possible solution for this is to build models that use widely
57 available data that are correlated to many of the *in situ* variables. For instance, landscape metrics of
58 forests, agriculture, wetlands, and urban land in contributing watersheds have all been shown to explain
59 a significant proportion of the variation (ranging from 50-86%, depending on study) in nutrients in
60 receiving waters (Jones et al. 2001, 2004, Seilheimer et al. 2013). Building on these previously identified
61 associations might allow us to use only landscape and other universally available data to build models.
62 Identifying predictors using this type of ubiquitous data would allow for estimating trophic state in
63 both monitored and unmonitored lakes.

64 Many published models of nutrients and trophic state in freshwater systems are based on linear modelling
65 methods such as standard least squares methods or linear mixed models (Jones et al. 2001, e.g., 2004).
66 While these methods have proven to be reliable, they have limitations (e.g. independence and distribution

assumptions, and outlier sensitivity). Using data mining approaches, such as random forests, avoids many of the limitations, may reduce bias and often provides better predictions (Breiman 2001, Cutler et al. 2007, Peters et al. 2007). For instance, random forests are non-parametric 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). Lastly, random forests can deal with model selection uncertainty as predictions are based upon a consensus of many models and not just a single model selected with some measure of goodness of fit.

To build on past work we have identified four goals for this research. First, update trophic state modelling efforts with the use of random forests. Second, assess the accuracy of predicted trophic state in lakes with the full suite of data and then with the universally available data only. Third, identify important variables for describing lake trophic state and lastly, explore associations between trophic state and cyanobacteria to begin to understand how changes in trophic state may be linked to an important ecosystem disservice.

2 Methods

2.1 Data and Study Area

We utilize four primary sources of data for this study, the National Lakes Assessment (NLA), the National Land Cover Dataset (NLCD), modeled lake morphometry, and cyanobacteria abundance (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 scale and provide a unique snapshot view of the condition of lakes in the conterminous United States' during the summer of 2007.

The NLA data were collected during the summer of 2007 and the final data were released in 2009 (USEPA 2009). With consistent methods and metrics collected at 1056 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 examined the water quality

92 measurements and total cyanobacteria abundance from the National Lakes Assessment (USEPA 2009).
93 Adding to the monitoring data collected via the NLA, we use the 2006 NLCD data to examine landscape-
94 level drivers of trophic status in lakes. The NLCD is a national land use/land cover dataset that also
95 provides estimates of impervious surface. We calculated total proportion of each NLCD land use land
96 cover class and total percent impervious surface within a 3 kilometer buffer surrounding the lake (Homer
97 et al. 2004, Xian et al. 2009). A three kilometer buffer was selected as an intermediate measure of
98 adjacent neighborhood; the three kilometer buffer size is greater than the immediate parcel but smaller
99 than regional measures.

100 To account for unique aspects of each lake and characterize lake productivity, we also used measures
101 of lake morphometry (i.e. depth, volume, fetch, etc.). As these data are difficult to obtain for large
102 numbers of lakes over broad regions, we used modeled estimates of lake morphometry (Hollister and
103 Milstead 2010, Hollister et al. 2011, Hollister 2014). From these prior efforts we included, Surface Area,
104 Shoreline Length, Shoreline Development, Maximum Depth, Mean Depth, Lake Volume, Maximum
105 Lake Length, Mean Lake Width, Maximum Lake Width, and Fetch.

106 **2.2 Predicting Trophic State with Random Forests**

107 Random forest is a machine learning algorithm that aggregates numerous decision trees in order to
108 obtain a consensus prediction of the response categories (Breiman 2001). Bootstrapped sample data
109 are recursively partitioned according to a given random subset of predictor variables and trees are
110 completely grown without pruning. With each new tree, the sample data subset is randomly selected
111 and with each new split, the subset of predictor variables are randomly selected.

112 While random forests are able to handle numerous correlated variables without a decrease in prediction
113 accuracy, one possible downfall to this approach is that the resulting model may be difficult to interpret.
114 This is a problem often faced in gene selection and in that field, a variable selection method based
115 on random forest has been successfully applied and implemented in the R Language as `varSelRF`
116 (Díaz-Uriarte and De Andres 2006). With this method, a minimum set of variables that maximizes
117 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. One issue with the approach in `varSelRF` is that because of the randomization inherent in random forests it is possible to get variation in the minimum selected set of variables. To account for this we repeated `varSelRF` 100 times. In our case, repeating the procedure 100 times quickly converged on a set of all possible important variables.

2.3 Model Details

Using both `varSelRF` and `randomForest` we ran models for six sets of variables and trophic state classifications. These included three different combinations of the Chlorophyll *a* trophic states as the dependent variables and using all variables (*in situ* and GIS variables) or the GIS only variables (i.e. no *in situ* information) as the independent variables in the random forest. A listing of all considered variables is in Appendix 1. Trophic state was defined using the NLA chlorophyll *a* trophic state cut offs and the three combinations of trophic state were used to highlight the possible error caused by misclassification of adjacent classes, such as mesotrophic and eutrophic (Table 1). The six model combinations were:

- **Model 1:** Chlorophyll *a* trophic state - 4 class = All variables (*in situ* water quality, lake morphometry, and landscape)
- **Model 2:** Chlorophyll *a* trophic state - 3 class = All variables (*in situ* water quality, lake morphometry, and landscape)
- **Model 3:** Chlorophyll *a* trophic state - 2 class = All variables (*in situ* water quality, lake morphometry, and landscape)
- **Model 4:** Chlorophyll *a* trophic state - 4 class = All variables (lake morphometry, and landscape)
- **Model 5:** Chlorophyll *a* trophic state - 3 class = All variables (lake morphometry, and landscape)
- **Model 6:** Chlorophyll *a* trophic state - 2 class = All variables (lake morphometry, and landscape)

Our modelling work flow was as follows:

1. Use `iterVarSelRF` in the `LakeTrophicModelling` R package to identify a minimal set of variables that maximize accuracy of the random forest algorithm (Diaz-Uriarte 2010, Hollister et al. 2014). This subset of variables, the reduced model, is calculated for each of our 6 models.

2. Using R's `randomForest` package, we pass the reduced models selected with `iterVarSelRF` and assess model performance (Liaw and Wiener 2002).

2.4 Measures of Model Performance and Variable Importance

We assess the performance of the random forest models by comparing the total prediction accuracy and the kappa coefficient of the final confusion matrix. As random forest builds each tree on bootstrapped, random subsets of the original data, a separate independent validation dataset is not required and random forest error estimates expected to be unbiased [breiman2001random]. For each of the models, the final predictions were compared to the original data via a confusion matrix. The total accuracy (i.e. percent correctly predicted) was calculated. Since some agreement can be expected by chance alone, it is also useful to take this type of error into account. For this we calculated the kappa coefficient for each model as well (Cohen 1960, Hubert and Arabie 1985). The kappa coefficient can range from -1 to 1 with 0 equalling the agreement expected by chance alone. Values greater than 0 represent agreement greater than would be expected by chance, with values greater than 0.61 considered “substantial” agreement (Landis and Koch 1977).

Lastly, the random forest algorithm explicitly measures variable importance as mean decrease in Gini. The Gini Index is a measure of how well the data are classified into homogeneous groups. For every node, the splitting variables are permuted and the change in actual Gini and permuted Gini is recorded. The mean decrease Gini is a summed and standardized value for each variable (Breiman 2001). Higher values of mean decrease Gini suggest a higher importance for that variable.

3 Results and Discussion

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

168 3.1 Models

169 Accuracy for the models built with all predictors ranged from 0.669 to 0.867 and the kappa coefficient
170 had a minimum value of 0.549 and maximum of 0.734. The GIS only models had a total accuracy
171 between 0.489 and 0.757 and kappa coefficient between 0.302 and 0.515. The importance of variables for
172 the models including the *in situ* data were fairly stable while There was considerably more variation in
173 variable importance for the three different GIS only models. Details for each model are discussed below.

174 *Model 1: 4 Trophic States ~ All Variables* The reduced model for Model 1 included potassium,
175 nitrogen:phosphorus, total nitrogen, total phosphorus, total organic carbon, turbidity, ecoregion, organic
176 ions, dissolved organic carbon, and maximum depth and of these, turbidity, total phosphorus, total
177 nitrogen, and total organic carbon were the most four most important predictors of the four classes of
178 trophic state (Figure 2). Total accuracy for Model 1 was 0.669% and the Cohen's Kappa was 0.549
179 (Table 2).

180 *Model 2: 3 Trophic States ~ All Variables* For Model 2, the reduced model included turbidity, total
181 phosphorus, total nitrogen, total organic carbon, nitrogen:phosphorus, longitude, pH, estimated organic
182 anions, elevation, maximum depth, dissolved organic carbon, potassium, latitude, ecoregion, chloride,
183 ammonium and percent cropland (Figure 3). The top predictors for 3 trophic state classes were turbidity,
184 total phosphorus, total nitrogen, and total organic carbon (Figure 3). Model 2 accuracy was 0.795%
185 and the Cohen's Kappa was 0.613 (Table 3).

186 *Model 3: 2 Trophic States ~ All Variables* The reduced model for Model 3 was similar to Model 1
187 and Model 2 and included turbidity, total phosphorus, total nitrogen, nitrogen:phosphorus, potassium,
188 ecoregion, elevation, total organic carbon, growing degree days, longitude, sodium, maximum depth,
189 estimated organic anions, latitude, and dissolved organic carbon (Figure 4). The top three predictors
190 were the same for Model 3; however, elevation and growing degree days had a higher importance than
191 total organic carbon. (Figure 4). Total accuracy for Model 3 was 0.867% and the Cohen's Kappa was
192 0.734 (Table 4).

193 *Model 4: 4 Trophic States ~ GIS Only Variables* The selected variables for the Model 4 were longitude,
194 latitude, elevation, estimated mean lake depth, percent evergreen forest, estimated maximum lake depth,

percent cropland, and ecoregion (Figure 5). The most important variables were percent evergreen forest, ecoregion, percent cropland, and longitude (Figure 5). Total accuracy for Model 4 is 0.489% and the Cohen's Kappa is 0.302 (Table 5).

Model 5: 3 Trophic States ~ GIS Only Variables The reduced model for Model 5 included estimated mean lake depth, percent cropland, longitude, latitude, percent evergreen forest, elevation, estimated maximum lake depth, estimated lake volume, percent deciduous forest, percent developed open space, ecoregion, percent woody wetland, and percent shrub/scrub (Figure 6). The most important variables for model 5 were ecoregion, percent evergreen forest, percent cropland, and estimated mean depth. (Figure 6). Total accuracy for Model 5 is 0.676% and the Cohen's Kappa is 0.347 (Table 6).

Model 6: 2 Trophic States ~ GIS Only Variables The variable selection process for Model 6 produced a reduced model with ecoregion, growing degree days, percent evergreen forest, percent cropland, elevation, estimated mean lake depth, longitude, latitude, watershed area, estimated maximum lake depth, percent developed open space, percent deciduous forest, and estimated lake volume (Figure 7). Similar to models 4 and 5, the four most important variables were ecoregion, percent evergreen forest, percent cropland, and elevation (Figure 7). Total accuracy for Model 6 0.757% and the Cohen's Kappa is 0.515 (Table 7).

3.2 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 trophic state models, we have 10,000 votes for each lake. These values may be interpreted as 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 Model 1 were 0.81 for oligotrophic, 0.19 for mesotrophic, 0 for eutrophic, and 0 for hypereutrophic. This suggests little uncertainty in the predicted oligotrophic state.

Further, the maximum probability for each lake can be used as a measure of how certain the random forest model was of the prediction. We would expect higher total accuracy for lakes that had more

221 certain predictions. To test this we can examine the accuracy of trophic state predictions across the full
222 range of trophic state probabilities, similar to an approach outlined by Paul and MacDonald (2005)
223 and implemented by Hollister *et al.* (2008). We utilize this approach and examine the change in total
224 accuracy as a function of the maximum probability for each lake. As expected, lakes with higher
225 maximum vote probabilities are more accurately predicted (Figure 12). This suggest that even for
226 models with low overall accuracy there will also be a large number of cases that are predicted with high
227 accuracy.

228 3.3 Variable Selection and Importance

229 There was a great deal of agreement on the important variables for each set of models. In line with
230 past predictive modeling of cyanobacteria abundance, the *in situ* models consistently select the water
231 quality variables (turbidity, total nitrogen, total phosphorus, and N:P ratios) as important variables
232 (Downing et al. 2001). While there is variation in the response of cyanobacteria to changes in relative
233 nutrient concentrations, the general pattern suggests that limiting nutrients have considerable impact
234 once amounts increase beyond expected levels.

235 The mechanistic role of turbidity on lake trophic state is more complex. Light availability in turbid
236 waters is lower than in clear waters. This would suggest a negative relationship between turbidity
237 and chlorophyll *a*. Second, chlorophyll *a* can also be a component of turbidity and lakes with higher
238 chlorophyll *a* concentrations will also be more turbid. Last, chlorophyll *a* is not the only component of
239 of turbidity and turbid waters can be caused by, for example, increased sediment loads or tannin. This
240 would be a cause for concern with linear models; however, linearity is not an assumption of tree-based
241 modelling approaches such as random forest [need cite].

242 Our GIS models are capturing the large scale spatial pattern of trophic status gradient of lakes across
243 the United States. We reliably see latitude and longitude and ecoregion selected as important variables.
244 It is also possible that other variables selected as important are also capturing a portion of this trend.
245 For instance, elevation and growing degree days both have obvious spatial components, but may also be
246 accounting for variation in temperature.

247 The land use/land cover variables are also important variables in describing trophic state patterns.
248 Like elevation and growing degree days, there are broad scale spatial patterns inherent in the data.
249 For instance, the relative continental position of mountains in the United States is the spatial inverse
250 of the distribution of agricultural lands. However, it is known that forests are positively associated
251 with lower nutrient loads where as agricultural land shows a negative association. These more local
252 scale relationships with land use/land cover are likely providing additional predictive power to the
253 information in the broader scale data.

254 Lastly, morphometry (e.g. depth and volume) also proved to be important in the prediction of lake
255 trophic state. As morphometry shows little to no broad scale spatial pattern and is unique to a given
256 lake, these data are likely illuminating the local, lake scale drivers of trophic state. As only depth and
257 volume were selected this probably shows the importance of in lake nutrient processing and residence
258 time.

259 3.4 Associating Trophic State and Cyanobacteria

260 Cyanobacteria biomass should be closely related to trophic state as they contribute to the chlorophyll
261 concentration in a lake. These associations have been seen by others. If these associations are strong
262 enough we may be able to expand models such as those reported here to also predict probability of
263 cyanobacteria blooms. To test if trophic state may be used to differentiate cyanobacteria abundance
264 we examine distribution of cyanobacteria abundance for each trophic state and we also explore linear
265 associations between Chlorophyll *a* and cyanobacteria abundance.

266 The distribution of cyanobacteria abundance shows separation between all of the trophic state classifica-
267 tions (Figures 8, 9, and 10). Furthermore, there is a significant linear relationship ($r^2=0.33$) between
268 chlorophyll *a* and cyanobacteria abundance (Figure 11). Further, Yuan et al. (2014) used the 2007 NLA
269 to demonstrate that total nitrogen and chlorophyll *a* concentrations were good predictors of World
270 Health Organization microcystin (a toxin produced by some cyanobacteris) criteria exceedences. These
271 results suggest that trophic state is indeed an acceptable proxy for cyanobacteria abundance and that
272 in lakes with higher trophic state it is also reasonable to expect higher cyanobacteria.

273 4 Conclusions

274 Our research goals were to explore the utility of a widely used data mining algorithm, random forests,
275 in the modelling of lake trophic state. Further, we hoped to examine the utility of these models when
276 built with only ubiquitous GIS data, which would allow for making trophic state estimates for all lakes
277 in the United States. We were able to successfully predict a variety of trophic state classes. With the
278 GIS only data models our total accuracy ranged from 0.4894552 to 0.7574692 and with the full suite of
279 data our model accuracy had a minimum accuracy of 0.6690018 and maximum accuracy of 0.8669002.

280 While some of the models (i.e. Model 4) show relatively low prediction accuracies, another feature of
281 the random forest, votes, can provide additional information. In addition to providing a single estimate
282 of trophic state for each lake, our models also indicate the probability that a lake was classified in
283 any of the categories. These probabilities may be mapped directly to show the uncertainty of a given
284 predicted class. Furthermore as the certainty of prediction increases so to does the overall trophic
285 state classification accuracy (Figure 12). These results suggest that our models will provide reasonable
286 estimates of trophic state across the United States.

287 There was great deal of agreement on the important variables for each set of models. For the combined
288 *in situ* and GIS Models, the *in situ* water quality variables drove the predictions. This is expected. For
289 the GIS only models, the results were more nuanced with three broad categories routinely being selected
290 as important: broad scale spatial patterns in trophic state, land use/land cover controls of trophic state,
291 and local, lake-scale control driven by lake morphometry. Lastly, associations between trophic state
292 and cyanobacteria show that at the broad scale of the 2007 NLA there is a linear relationship between
293 chlorophyll *a* and cyanobacteria abundance and that using trophic state as a proxy for cyanobacteria
294 has potential.

295 These broad categories and the association between trophic state and total cyanobacteria abundance raise
296 three important considerations related to managing eutrophication. First, the broad scale patterning
297 suggests regional trends. This is important because it suggests that efforts to monitor, model and
298 manage eutrophication and cyanobacteria should be undertaken at both national and regional levels.
299 Second, while direct control of water quality in lakes would have a large impact, the land use/land

300 cover drivers (i.e. non-point sources) of water quality are also important and better management of
301 the spatial distribution of important classes such as forest and agriculture can provide some level of
302 control on trophic state and amount of cyanobacteria present. Third, in-lake processes (i.e. residence
303 time, nutrient cycling, etc.) are, as expected, very important and need to be part of any management
304 strategy. Building on these efforts through updated models, direct prediction of cyanobacteria, and
305 additional information on the regional differences will help us get a better handle on the broad scale
306 dynamics of productivity in lakes and the potential risk to human health from cyanobacteria blooms.

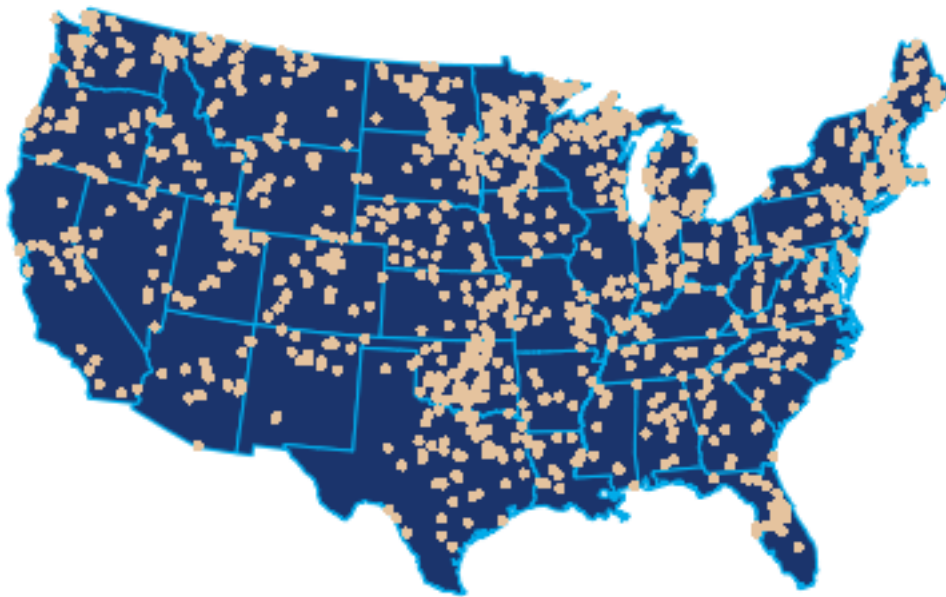


Figure 1: Map of the distribution of National Lakes Assessment Sampling locations

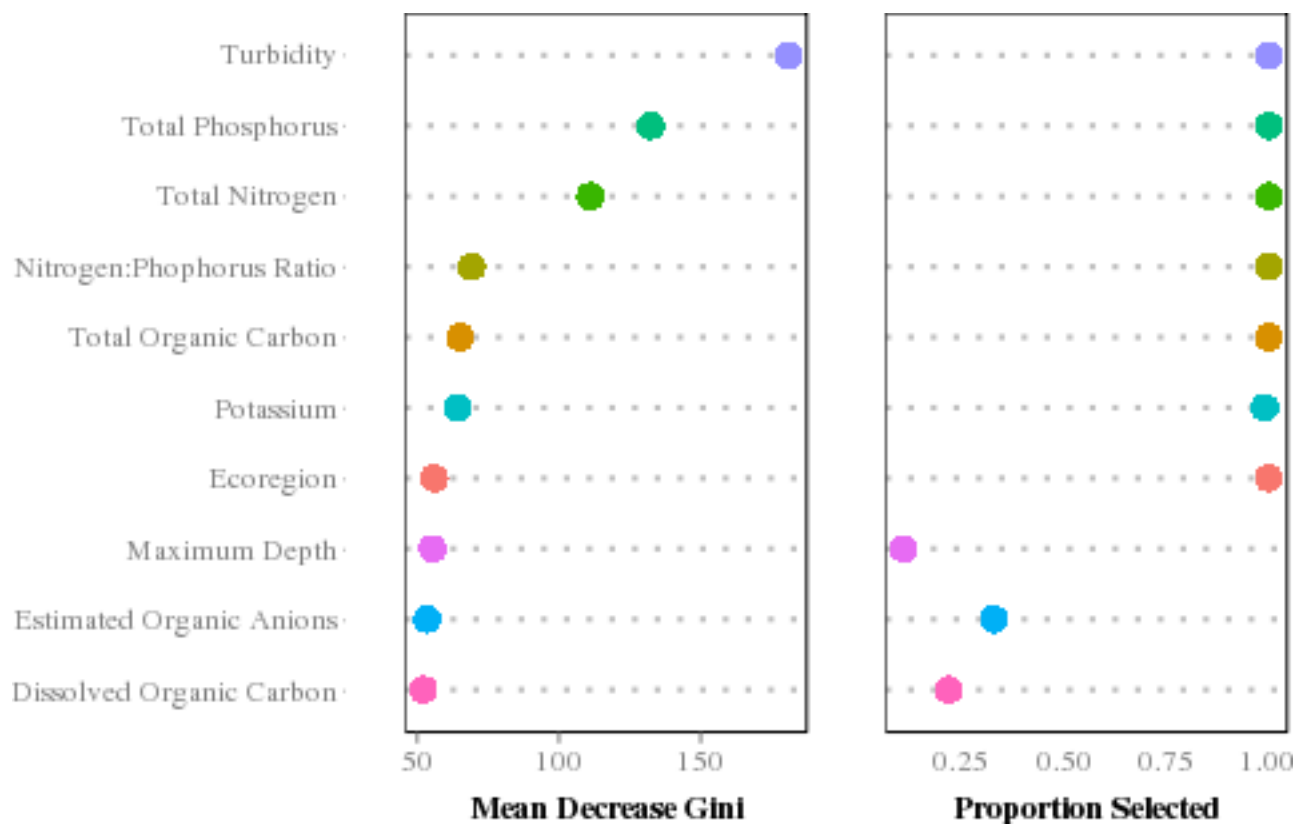


Figure 2: Importance plot for Model 1

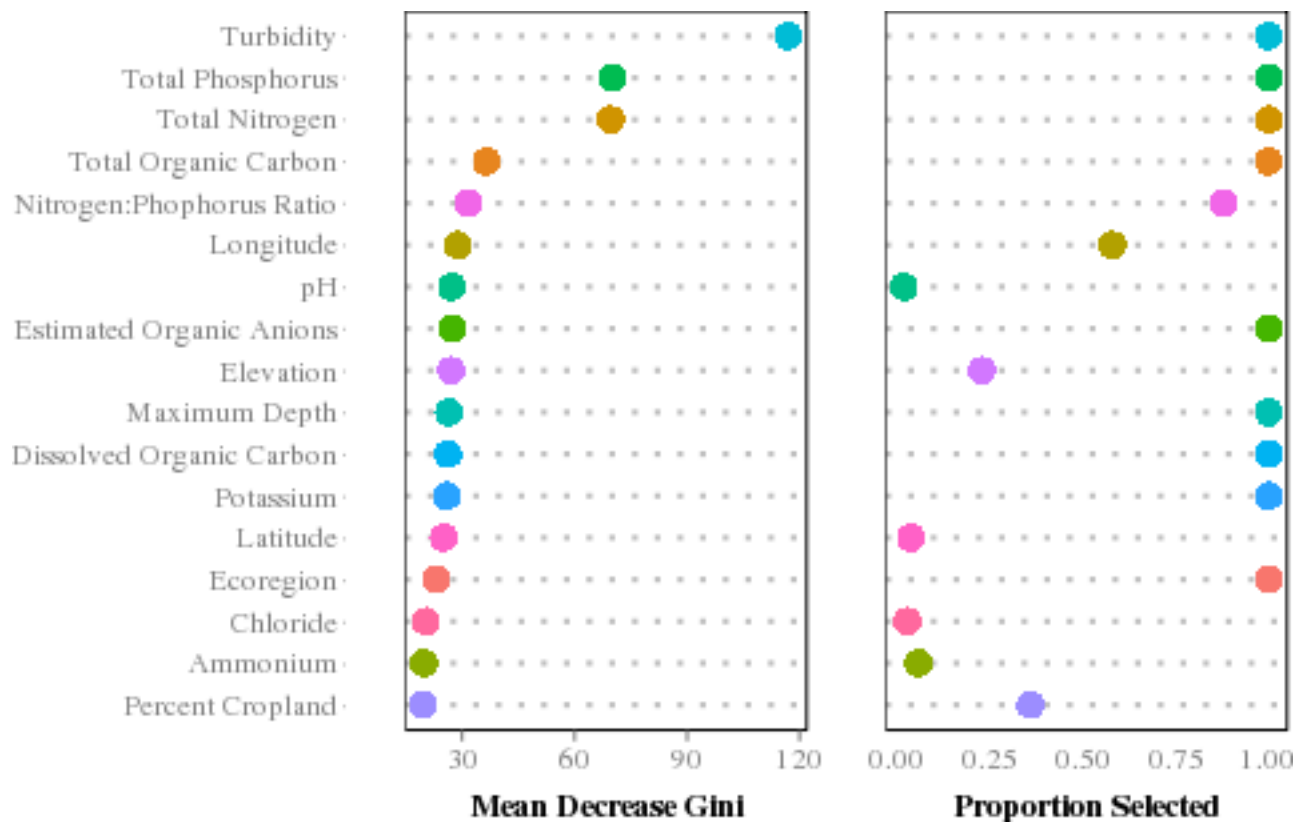


Figure 3: Importance plot for Model 2

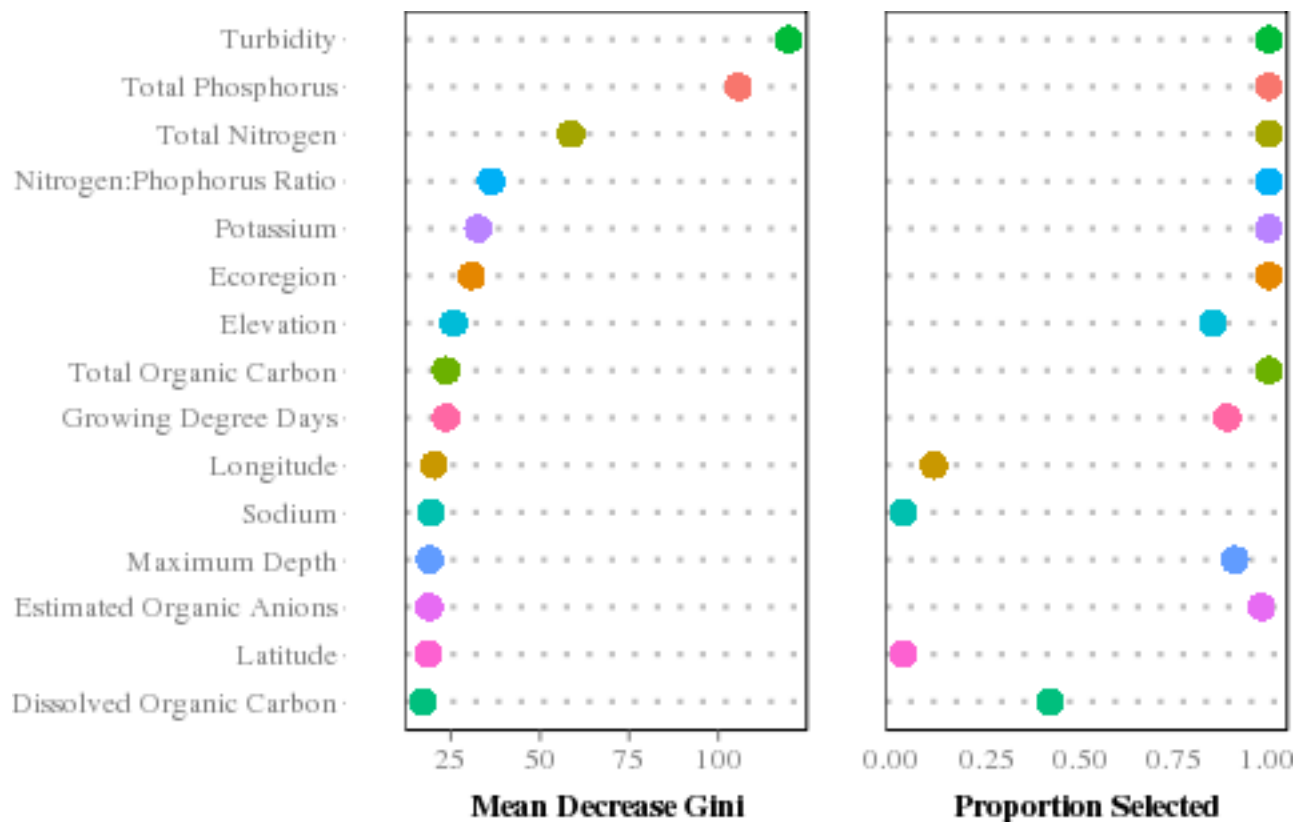


Figure 4: Importance plot for Model 3

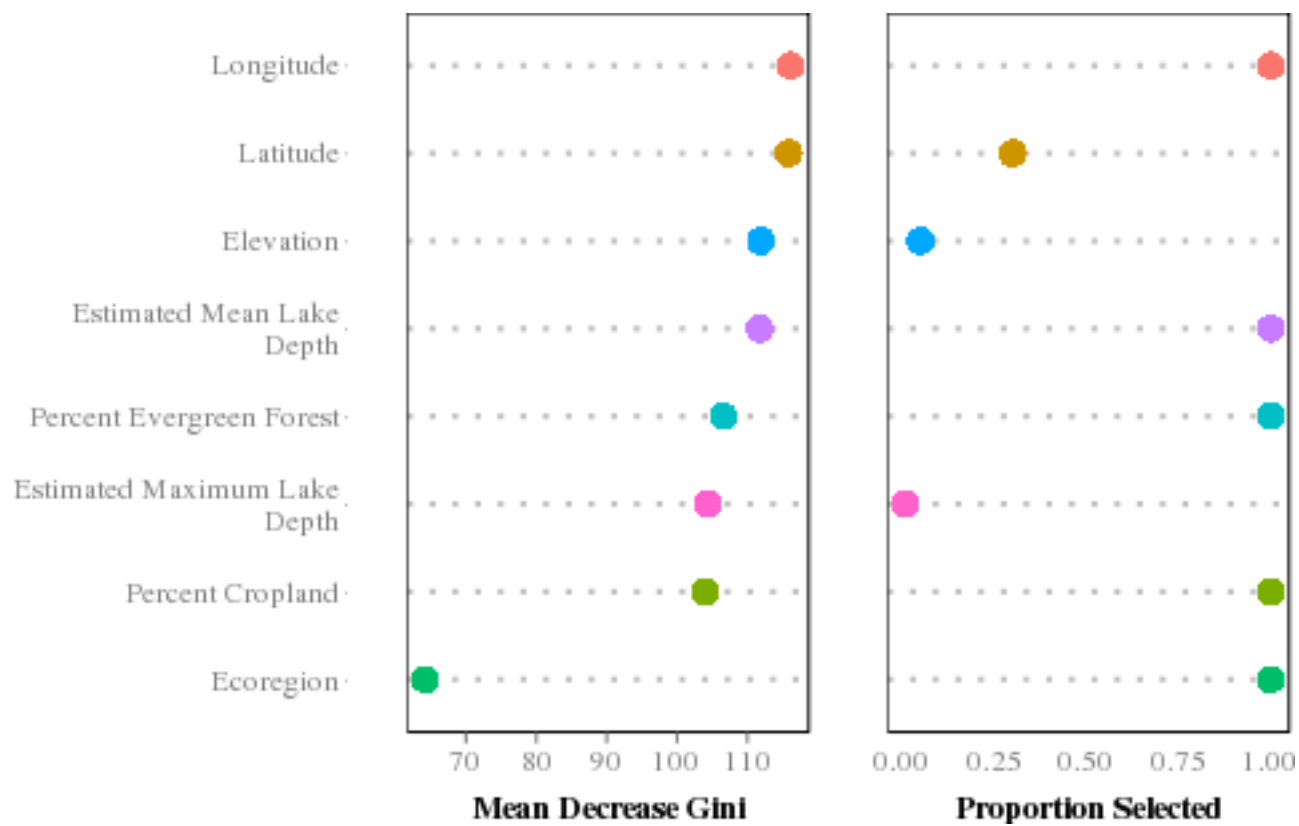


Figure 5: Importance plot for Model 4

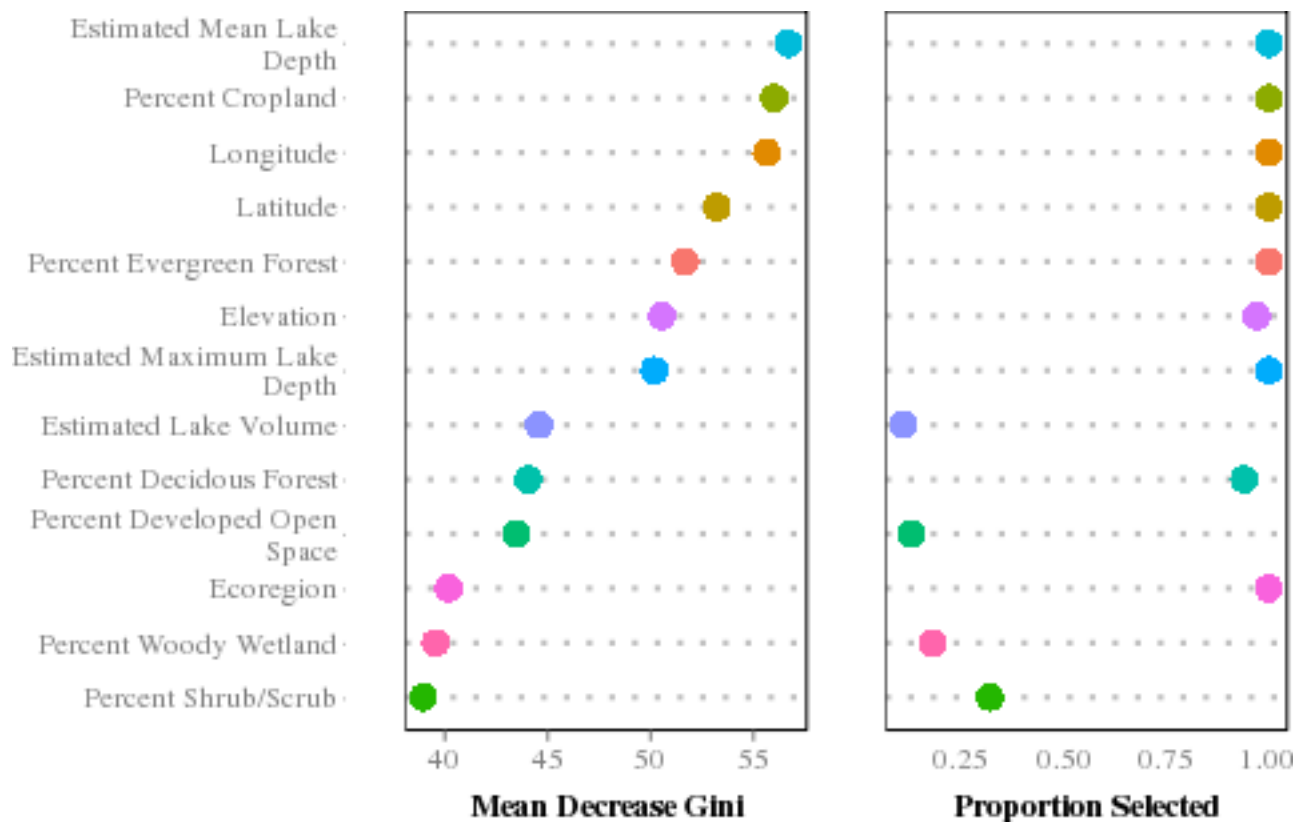


Figure 6: Importance plot for Model 5

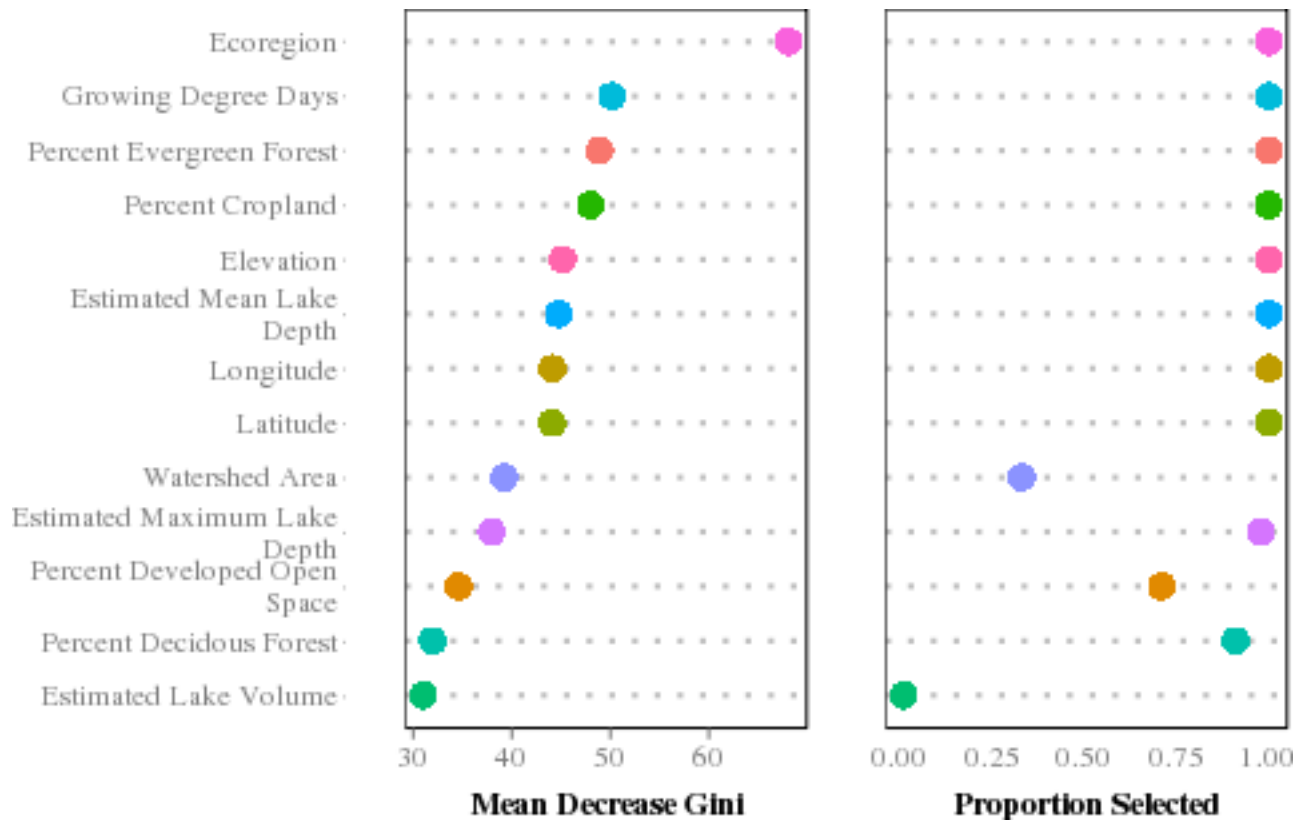


Figure 7: Importance plot for Model 6

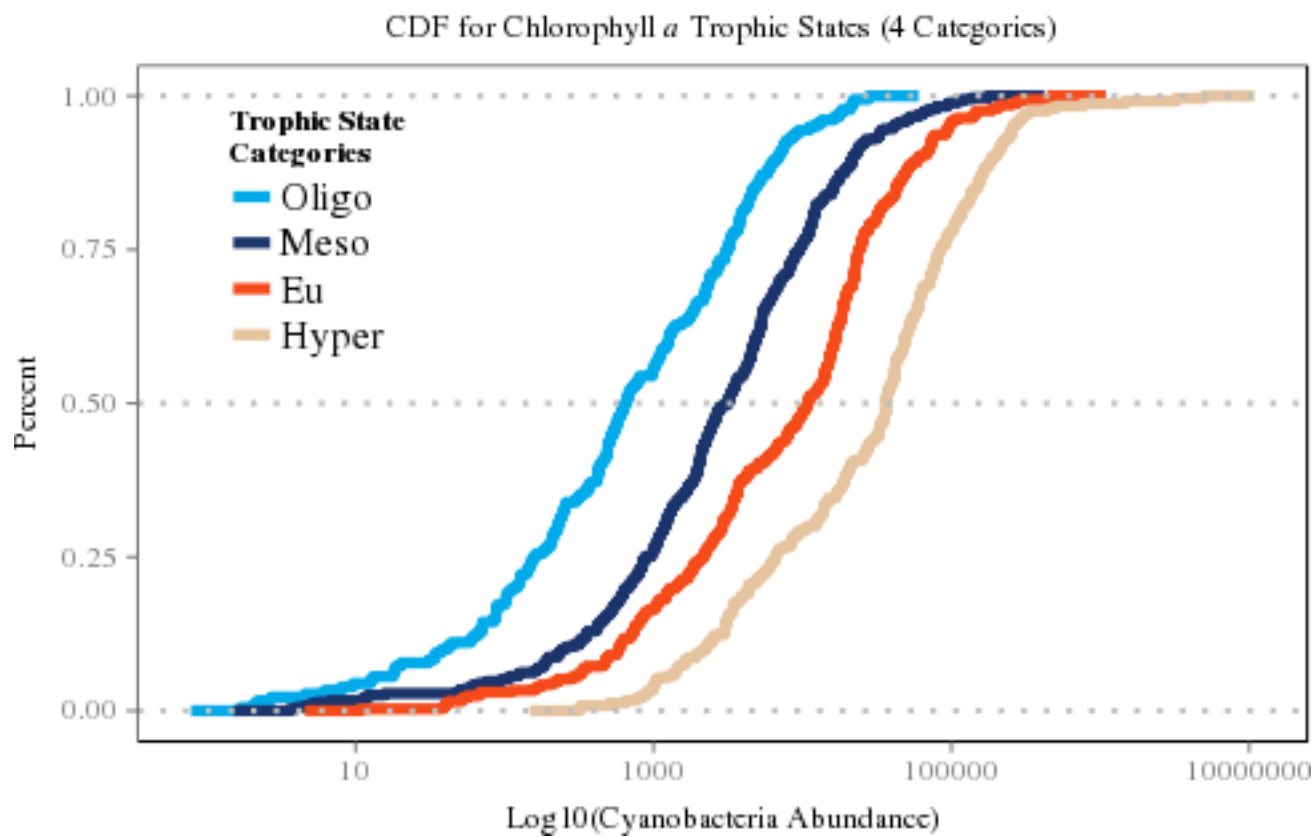


Figure 8: Cumulative distribution function of cyanobacteria abundance for 4 trophic state classes

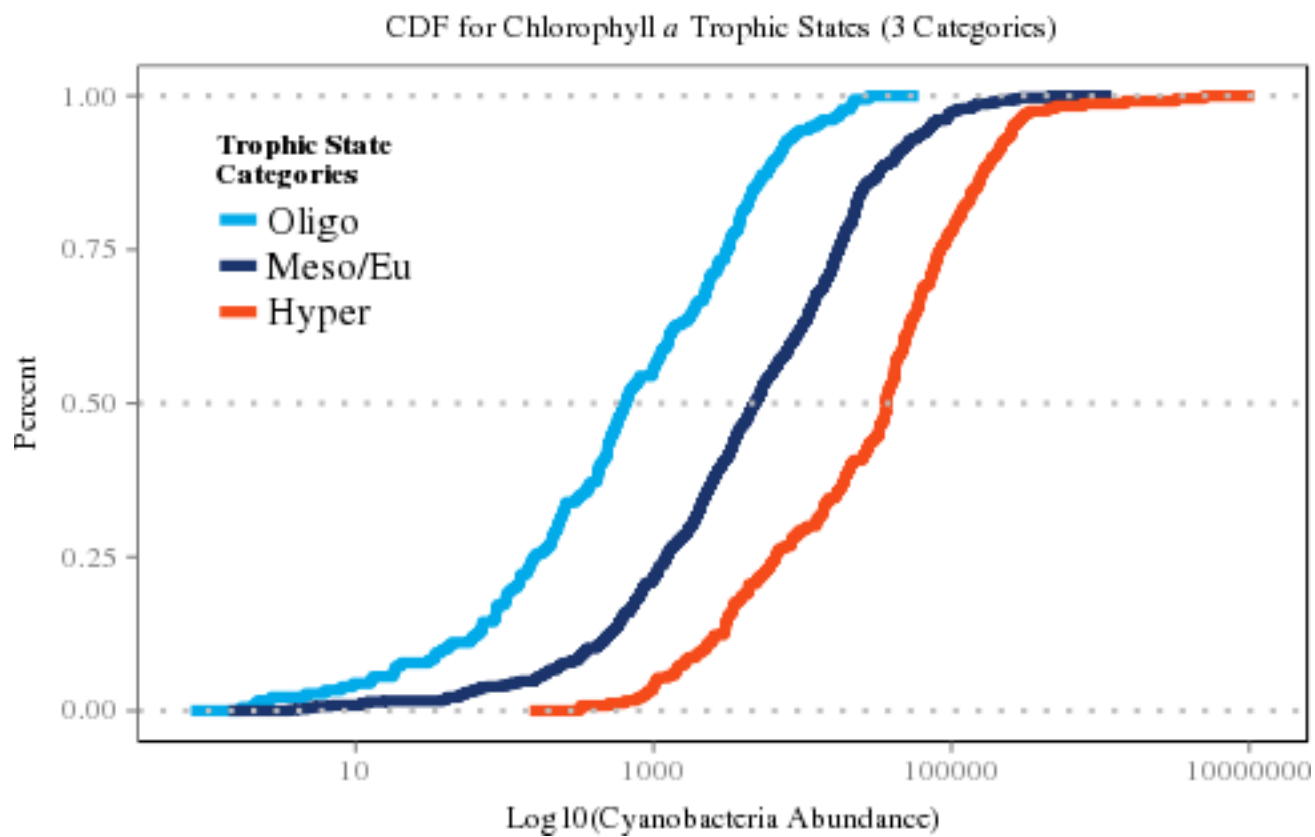


Figure 9: Cumulative distribution function of cyanobacteria abundance for 3 trophic state classes

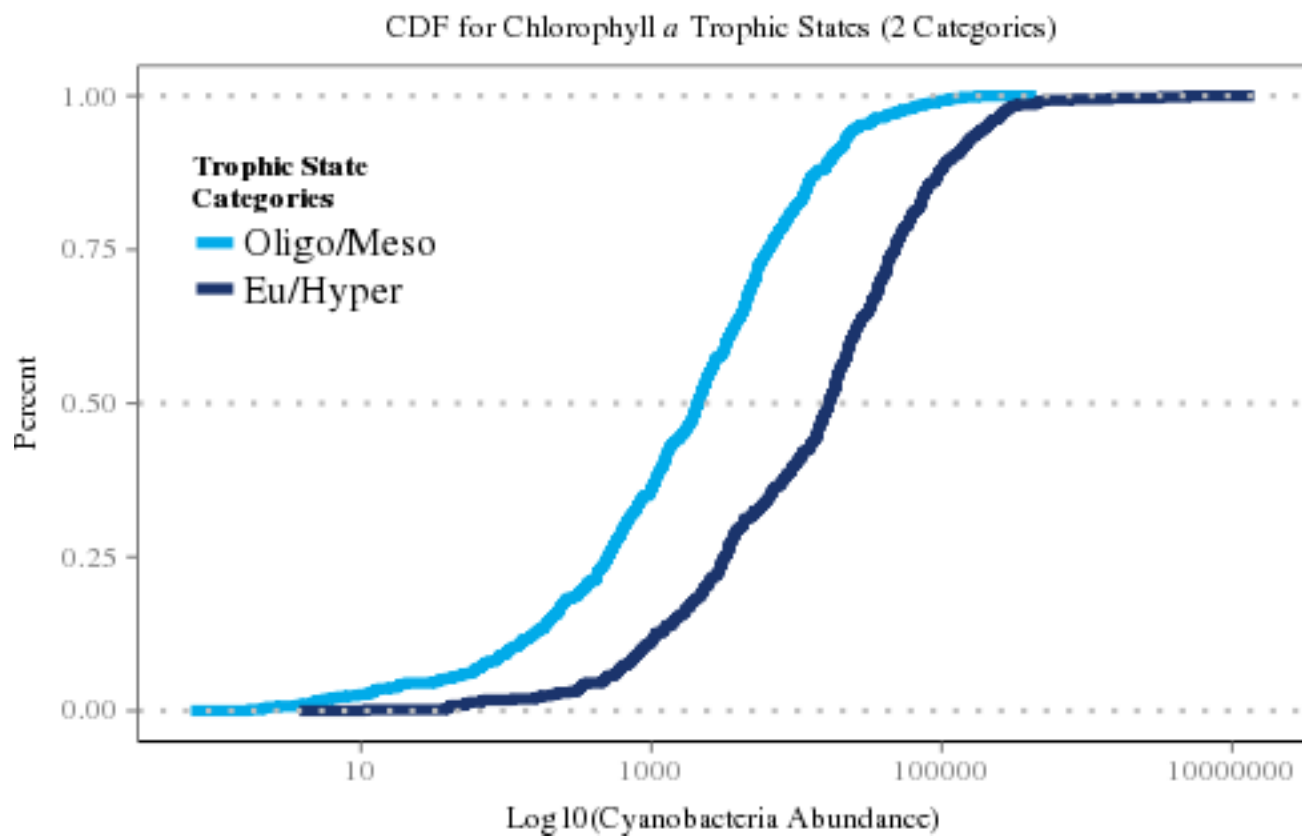


Figure 10: Cumulative distribution function of cyanobacteria abundance for 2 trophic state classes

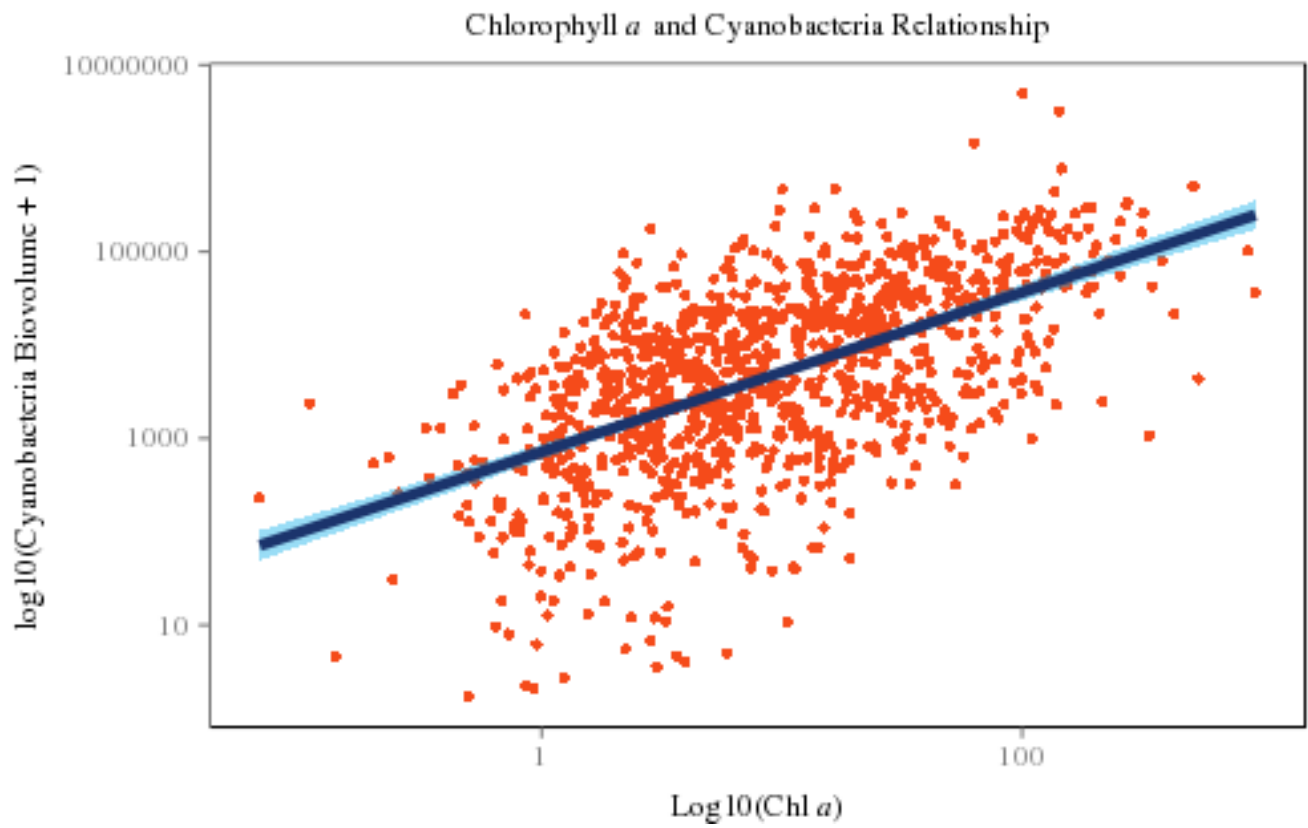


Figure 11: Cholorphyll *a* and cyanobacteria abundance scatterplot

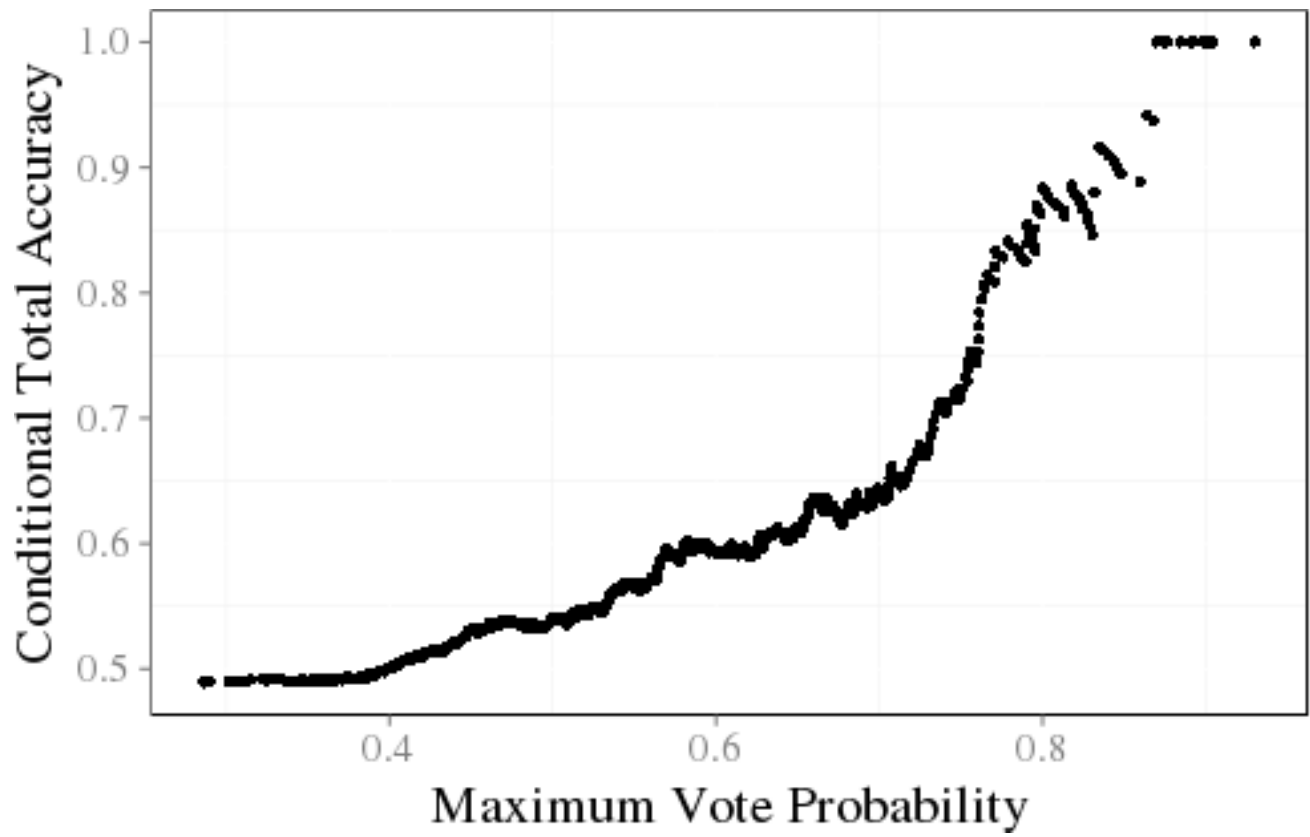


Figure 12: Comparison of certainty of trophic state prediction and total accuracy

Trophic State (4)	Trophic State (3)	Trophic State (2)	Cut-off	n
oligo	oligo	oligo/meso	≤ 0.2	198
meso	meso/eu	oligo/meso	$>2-7$	362
eu	meso/eu	eu/hyper	$>7-30$	337
hyper	hyper	eu/hyper	>30	246

Table 1: Chlorophyll a based trophic state cut-offs

Oligo	Meso	Eu	Hyper	class.error
135	58	4	1	0.32
42	233	77	10	0.36
2	66	222	46	0.34
0	3	69	174	0.29

Table 2: Random Forest confusion matrix for Model 1

Oligo	Meso/Eu	Hyper	class.error
122	74	0	0.38
43	604	42	0.12
0	72	173	0.29

Table 3: Random Forest confusion matrix for Model 2

Oligo/Meso	Eu/Hyper	class.error
485	75	0.13
77	505	0.13

Table 4: Random Forest confusion matrix for Model 3

Oligo	Meso	Eu	Hyper	class.error
94	72	28	2	0.52
50	201	80	30	0.44
21	110	131	73	0.61
1	34	80	131	0.47

Table 5: Random Forest confusion matrix for Model 4

Oligo	Meso/Eu	Hyper	class.error
80	115	1	0.59
50	585	61	0.16
0	142	104	0.58

Table 6: Random Forest confusion matrix for Model 5

Oligo/Meso	Eu/Hyper	class.error
428	129	0.23
147	434	0.25

Table 7: Random forest confusion matrix for Model 6

variable_names	description	type
PercentImperv_3000m	Percent Impervious	GIS
WaterPer_3000m	Percent Water	GIS
IceSnowPer_3000m	Percent Ice/Snow	GIS
DevOpenPer_3000m	Percent Developed Open Space	GIS
DevLowPer_3000m	Percent Low Intensity Development	GIS
DevMedPer_3000m	Percent Medium Intensity Development	GIS
DevHighPer_3000m	Percent High Intensity Development	GIS
BarrenPer_3000m	Percent Barren	GIS
DeciduousPer_3000m	Percent Deciduous Forest	GIS
EvergreenPer_3000m	Percent Evergreen Forest	GIS
MixedForPer_3000m	Percent Mixed Forest	GIS
ShrubPer_3000m	Percent Shrub/Scrub	GIS
GrassPer_3000m	Percent Grassland	GIS
PasturePer_3000m	Percent Pasture	GIS
CropsPer_3000m	Percent Cropland	GIS
WoodyWetPer_3000m	Percent Woody Wetland	GIS
HerbWetPer_3000m	Percent Herbaceous Wetland	GIS
AlbersX	Longitude	GIS
AlbersY	Latitude	GIS
LakeArea	Lake Surface Area	GIS
LakePerim	Lake Perimeter	GIS

variable_names	description	type
ShoreDevel	Shoreline Development Index	GIS
DATE_COL	Date Samples Collected	Water Quality
WSA_ECO9	Ecoregion	GIS
BASINAREA	Watershed Area	GIS
DEPTHMAX	Maximum Depth	Water Quality
ELEV_PT	Elevation	GIS
DO2_2M	Dissolved Oxygen	Water Quality
PH_FIELD	pH	Water Quality
COND	Conductivity	Water Quality
ANC	Acid Neutralizing Capacity	Water Quality
TURB	Turbidity	Water Quality
TOC	Total Organic Carbon	Water Quality
DOC	Dissolved Organic Carbon	Water Quality
NH4	Ammonium	Water Quality
NO3_NO2	Nitrate/Nitrite	Water Quality
NTL	Total Nitrogen	Water Quality
PTL	Total Phosphorus	Water Quality
CL	Chloride	Water Quality
NO3	Nitrate	Water Quality
SO4	Sulfate	Water Quality
CA	Calcium	Water Quality
MG	Magnesium	Water Quality

variable_names	description	type
Na	Sodium	Water Quality
K	Potassium	Water Quality
COLOR	Color	Water Quality
SIO2	Silica	Water Quality
H	Hydrogen Ions	Water Quality
OH	Hydroxide	Water Quality
NH4ION	Calculate Ammonium	Water Quality
CATSUM	Cation Sum	Water Quality
ANSUM2	Anion Sum	Water Quality
ANDEF2	Anion Deficit	Water Quality
SOBC	Base Cation Sum	Water Quality
BALANCE2	Ion Balance	Water Quality
ORGION	Estimated Organic Anions	Water Quality
CONCAL2	Calculated Conductivity	Water Quality
CONDHO2	D-H-O Calculated Conductivity	Water Quality
TmeanW	Mean Profile Water Temperature	Water Quality
DDs45	Growing Degree Days	GIS
MaxLength	Maximum Lake Length	GIS
MaxWidth	Maximum Lake Width	GIS
MeanWidth	Mean Lake Width	GIS
FetchN	Fetch from North	GIS
FetchNE	Fetch form Northeast	GIS

variable_names	description	type
FetchE	Fetch from East	GIS
FetchSE	Fetch from Southeast	GIS
MaxDepthCorrect	Estimated Maximum Lake Depth	GIS
VolumeCorrect	Estimated Lake Volume	GIS
MeanDepthCorrect	Estimated Mean Lake Depth	GIS
NPratio	Nitrogen:Phophorus Ratio	Water Quality

References

- 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 (nee Ferguson), 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.
- Diaz-Uriarte, R. 2010. varSelRF: Variable selection using random forests.
- 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., S. B. Watson, and E. McCauley. 2001. Predicting cyanobacteria dominance in lakes. *Canadian journal of fisheries and aquatic sciences* 58:1905–1908.
- Hasler, A. D. 1969. Cultural eutrophication is reversible. *BioScience* 19:425–431.
- Hollister, J. W. 2014. lakemorpho: Lake morphometry in r.
- Hollister, J. W., W. B. Milstead, and B. J. Kreakie. 2014. LakeTrophicModelling: Package to reproduce hollister et al. (2014) modeling lake trophic state: A data mining approach.
- Hollister, J. W., W. B. Milstead, and M. A. Urrutia. 2011. Predicting maximum lake depth from surrounding topography. *PLoS ONE* 6:e25764.
- 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.

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

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

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

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

340 Jones, J., M. Knowlton, D. Obrecht, and E. Cook. 2004. Importance of landscape variables and
341 morphology on nutrients in missouri reservoirs. *Canadian Journal of Fisheries and Aquatic Sciences*
342 61:1503–1512.

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

346 Landis, J. R., and G. G. Koch. 1977. The measurement of observer agreement for categorical data.
347 *biometrics*:159–174.

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

349 Milstead, W. B., J. W. Hollister, R. B. Moore, and H. A. Walker. 2013. Estimating summer nutrient
350 concentrations in northeastern lakes from sPARROW load predictions and modeled lake depth and
351 volume. *PloS one* 8:e81457.

352 Paul, J. F., and M. E. McDonald. 2005. Development of empirical, geographically specific water quality
353 criteria: a conditional probability analysis approach. *Wiley Online Library*.

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

356 Salas, H. J., and P. Martino. 1991. A simplified phosphorus trophic state model for warm-water tropical

357 lakes. *Water research* 25:341–350.

358 Schindler, D., and J. Vallentyne. 2008. Algal bowl: overfertilization of the world’s freshwaters and
359 estuaries.

360 Seilheimer, T. S., P. L. Zimmerman, K. M. Stueve, and C. H. Perry. 2013. Landscape-scale modeling of
361 water quality in lake superior and lake michigan watersheds: How useful are forest-based indicators?
362 *Journal of Great Lakes Research* 39:211–223.

363 Smith, V. H. 1998. Cultural eutrophication of inland, estuarine, and coastal waters. Pages 7–49 *in*
364 *Successes, limitations, and frontiers in ecosystem science*. Springer.

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

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

369 USEPA. 2009. National lakes assessment: a collaborative survey of the nation’s lakes. ePA 841-r-09-001.
370 Office of Water; Office of Research; Development, US Environmental Protection Agency Washington,
371 DC.

372 Xian, G., C. Homer, and J. Fry. 2009. Updating the 2001 national land cover database land cover clas-
373 sification to 2006 by using landsat imagery change detection methods. *Remote Sensing of Environment*
374 113:1133–1147.

375 Yuan, L. L., A. I. Pollard, S. Pather, J. L. Oliver, and L. D’Anglada. 2014. Managing microcystin:
376 identifying national-scale thresholds for total nitrogen and chlorophyll a. *Freshwater Biology* 59:1970–
377 1981.