

Modeling Lake Trophic State: A Data Mining Approach

Jeffrey W. Hollister^{*} ¹ W. Bryan Milstead ¹ Betty J. Kreakie ¹

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/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 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 models of trophic state (as determined by chlorophyll *a* concentration). The first model estimates trophic state with *in situ* as well as universally available data and the second model 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 biovolume and the measured and predicted trophic state 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) (Carlson 1977) and lakes naturally occur across this range. Naturally oligotrophic lakes occur in nutrient poor areas or have a more recent geologic history. These lakes are often found in higher elevations, have clear water, and are often favored for drinking water or direct contact recreation (e.g. swimming).

38 Lakes with higher productivity (e.g. eutrophic lakes) have greater nutrient loads, tend to be less clear,
39 have greater density of aquatic plants, and often support more diverse and abundant fish communities.
40 Higher primary productivity is not necessarily a predictor of poor ecological condition.

41 It is natural for lakes to shift from lower to higher trophic states but this is a slow process. Given this
42 fact, monitoring trophic state allows the identification of rapid shifts in trophic state or locating lakes
43 with unusually high productivity (e.g. hypereutrophic). These cases are indicative of lakes under greater
44 anthropogenic nutrient loads, also known as cultural eutrophication, and are more likely to be at risk of
45 fish kills, fouling, and harmful algal blooms (Smith 1998, Smith et al. 1999, 2006). Given the association
46 between trophic state and many ecosystem services and disservices, being able to model trophic state
47 could allow for estimating trophic state in unmonitored lakes and provide a first cut at identifying lakes
48 with the potential for harmful algal blooms and other problems associated with cultural eutrophication.

49 As trophic state can be defined by a number of *in situ* water quality measurements, most models have
50 used this information as predictors. This leads to accurate models, but also requires data that is sparse
51 and not always available, limiting the population of lakes that can be modeled. Landscape data is
52 ubiquitous . . .

53 We have three goals for this preliminary research. First, we build and assess multiple models of lake
54 trophic state using a full suite of data including *in situ* water quality and universally available data
55 (e.g. landscape data). Second, we assess the accuracy of predicted trophic state in lakes with only the
56 universally available data. Lastly, we explore associations between trophic state and cyanobacteria to
57 explore.

58 **2 Methods**

59 **2.1 Data and Study Area**

60 We utilize four primary sources of data for this study, the National Lakes Assessment (NLA), the National
61 Lake Cover Dataset (NLCD), modeled lake morphometry, and estimated cyanobacteria biovolumes
62 (Homer et al. 2004, USEPA 2009, Xian et al. 2009, Hollister and Milstead 2010, Hollister et al. 2011,

63 Beaulieu et al. 2013, Hollister 2014). All datasets are national in scale and provide a unique snapshot
64 view of the condition of lakes in the United States’.

65 The NLA data were collected during the summer of 2007 and the final data were released in 2009.
66 With consistent methods and metrics collected at 1056 locations across the conterminous United States
67 (Figure 1), the NLA provides a unique opportunity to examine broad scale patterns in lake productivity.
68 The NLA collected data on biophysical measures of lake water quality and habitat. For this analysis
69 we primarily examined the water quality measurements from the NLA (USEPA 2009). Adding to
70 the monitoring data collected via the NLA, we use the 2006 NLCD data to examine the possible
71 landscape-level drivers of trophic status in lakes. The NLCD is a nationally collected land use land
72 cover dataset that also provides estimates of impervious surface. We collected total land use land
73 cover and total percent impervious surface within a 3 kilometer buffer surrounding the lake to examine
74 larger landscape-level effect (Homer et al. 2004, Xian et al. 2009). We also used various measures
75 of lake morphometry (i.e. depth, volume, fetch, etc.) as they are important in understanding lake
76 productivity, yet many of these data are difficult to obtain for large numbers of lakes over broad regions.
77 To add this information we modeled lake morphometry (Hollister and Milstead 2010, ???, Hollister
78 et al. 2011, Hollister 2014). Lastly, to explore associations between trophic state and cyanobacteria,
79 we used estimates of cyanobacterial biovolume calculated by Beaulieu *et al.* (2013). Cyanobacteria
80 biovolumes are a truer measure of cyanobacteria dominance than abundance as there is great variability
81 in the size within and between species. We have consolidated the taxa level estimates from Beaulieu *et*
82 *al.* (2013) and summed that information on a per-lake basis.

83 2.2 Predicting Trophic State with Random Forests

84 Random forest is a machine learning algorithm that aggregates numerous decision trees in order to
85 obtain a consensus prediction of the response categories (Breiman 2001). Bootstrapped sample data is
86 recursively partitioned according to a given random subset of predictor variables and completely grown
87 without pruning. With each new tree, both the sample data and predictor variable subset is randomly
88 selected.

89 While random forests are able to handle numerous correlated variables without a decrease in prediction

90 accuracy, unusually large numbers of related variables can reduce accuracy and increase the chances
 91 of over-fitting the model. This is a problem often faced in gene selection and in that field, a variable
 92 selection method based on random forest has been succesfully applied (Díaz-Uriarte and De Andres
 93 2006). We use varselRF in R to initially examine the importance of the water quality and GIS derived
 94 variables and select a subset, the reduced model, to then pass to random forest(Diaz-Uriarte 2010).

95 Using R's randomForest package, we pass the reduced models selected with varSelRF and calculate
 96 confusion matrices, overall accuracy and kappa coeffecient (Liaw and Wiener 2002). From the reduced
 97 model random forests we collect a consensus prediction and calculate a confusion matrix and summary
 98 stats.

99 2.3 Model Details

100 Using a combination of the `varSelRF` and `randomForest` we ran models for six combinations of variables
 101 and trophic state classifications. These combinations included different combinations of the Chlorophyll *a*
 102 trophic states (Table 1) along with all variables and the GIS only variables (i.e. no *in situ* information).
 103 The six model combinations were:

- 104 1. Chlorophyll *a* trophic state - 4 class = All variables (*in situ* water quality, lake morphometry, and
 105 landscape)
- 106 2. Chlorophyll *a* trophic state - 3 class = All variables (*in situ* water quality, lake morphometry, and
 107 landscape)
- 108 3. Chlorophyll *a* trophic state - 2 class = All variables (*in situ* water quality, lake morphometry, and
 109 landscape)
- 110 4. Chlorophyll *a* trophic state - 4 class = All variables (lake morphometry, and landscape)
- 111 5. Chlorophyll *a* trophic state - 3 class = All variables (lake morphometry, and landscape)
- 112 6. Chlorophyll *a* trophic state - 2 class = All variables (lake morphometry, and landscape)

113 **3 Results and Discussion**

114 **3.1 Model 1: 4 Trophic States ~ All Variables**

115 The selected variables that made up Model 1 were Potassium, Nitrogen:Phosphorus, Total Nitrogen,
116 Total Phosphorus, Total Organic Carbon, Turbidity, Ecoregion, Organic Ions, Dissolved Organic Carbon,
117 and Maximum Lake Depth (Table 2). Total accuracy for Model 1 is 0.669% and the Cohen's Kappa is
118 0.549 (Table 3).

119 Lastly, turbidity, total phosphorus, total nitrogen, and total organic carbon were the most important
120 predictors of the 4 classes of trophic state (Figure 2).

121 **3.2 Model 2: 3 Trophic States ~ All Variables**

122 Total accuracy for Model 2 is 0.796% and the Cohen's Kappa is 0.613.

123 **3.3 Model 3: 2 Trophic States ~ All Variables**

124 Total accuracy for Model 3 is 0.867% and the Cohen's Kappa is 0.734.

125 **3.4 Model 4: 4 Trophic States ~ GIS Only Variables**

126 Total accuracy for Model 4 is 0.489% and the Cohen's Kappa is 0.302.

127 **3.5 Model 5: 3 Trophic States ~ GIS Only Variables**

128 Total accuracy for Model 5 is 0.676% and the Cohen's Kappa is 0.347.

129 **3.6 Model 6: 2 Trophic States ~ GIS Only Variables**

130 Total accuracy for Model 6 0.757% and the Cohen's Kappa is 0.515.

131 **3.7 Associating Trophic State and Cyanobacteria**

132 Arm waving goes here.

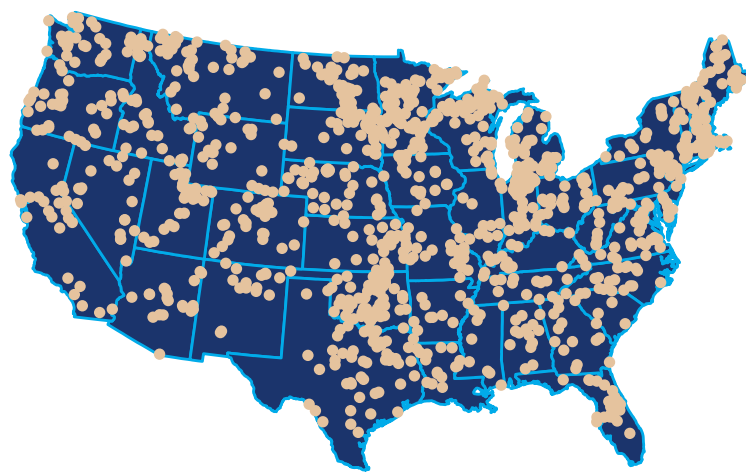


Figure 1: Map of the distribution of National Lakes Assessment Sampling locations

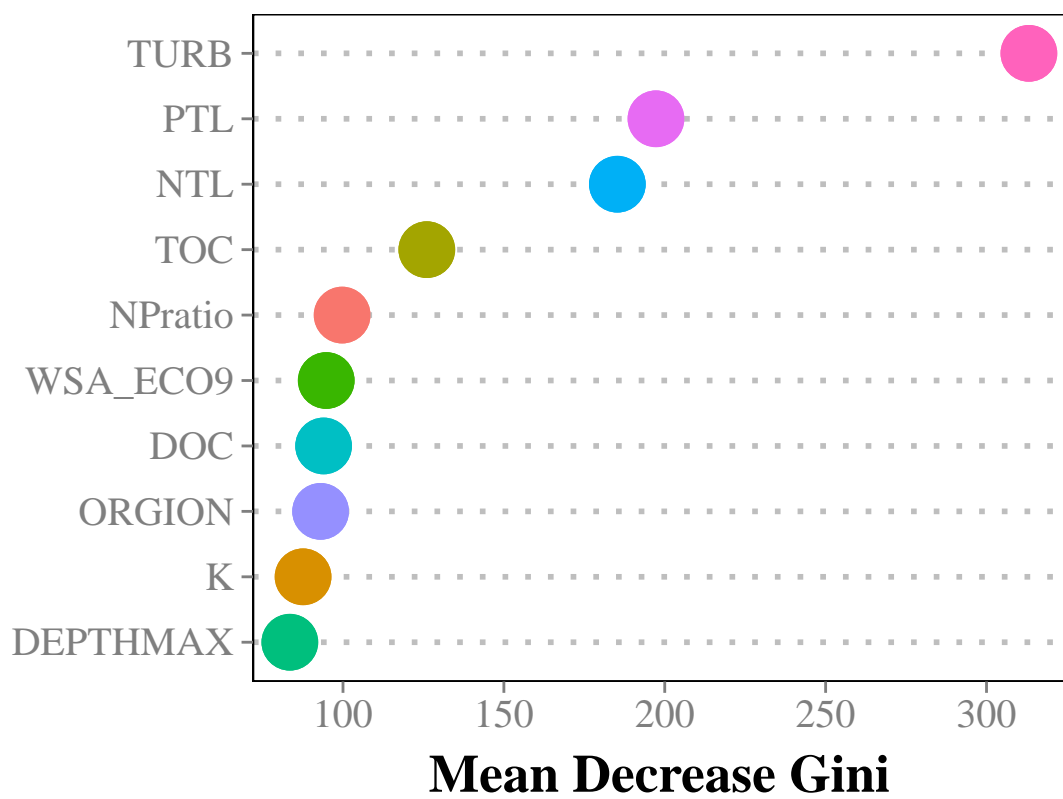


Figure 2: Importance plot for Model 1

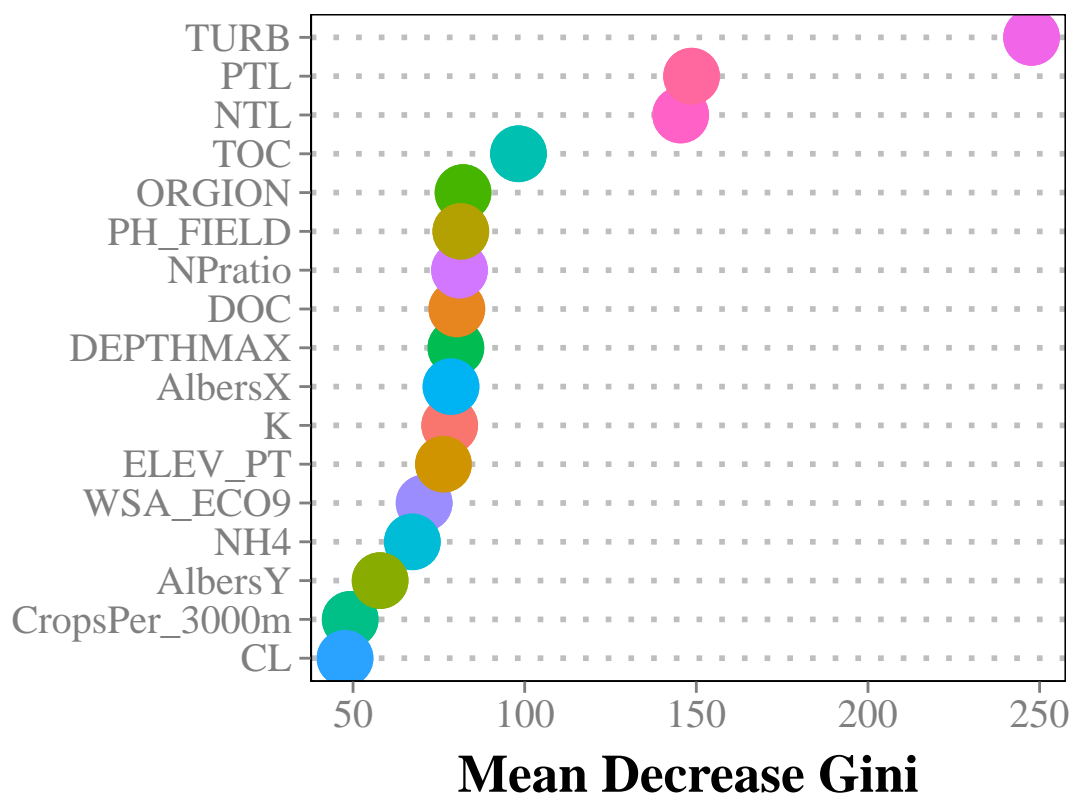


Figure 3: Importance plot for Model 2

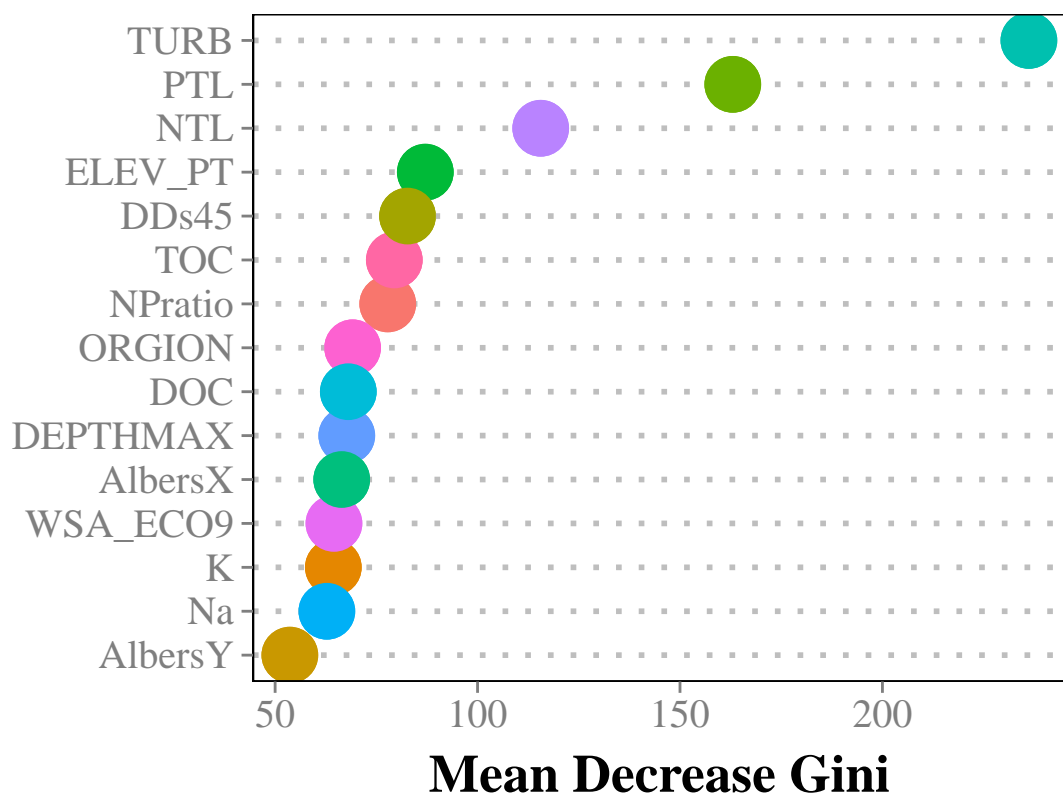


Figure 4: Importance plot for Model 3

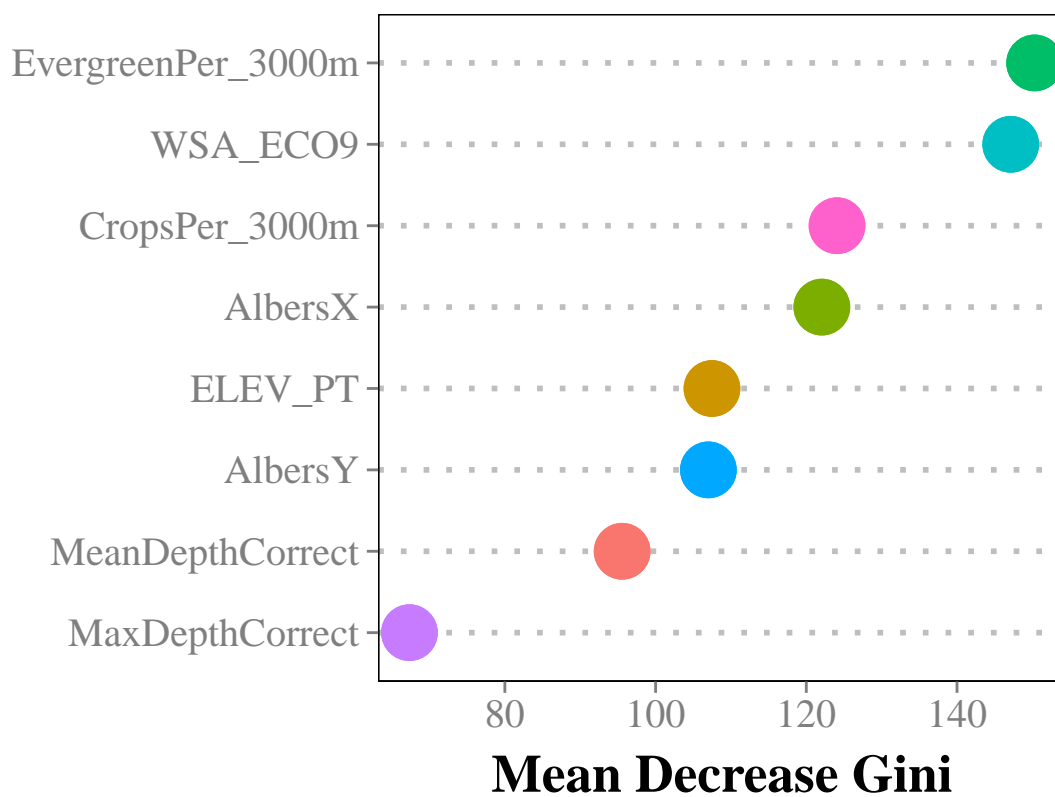


Figure 5: plot of chunk Importance_Model4

134 :Importance plot for Model 4

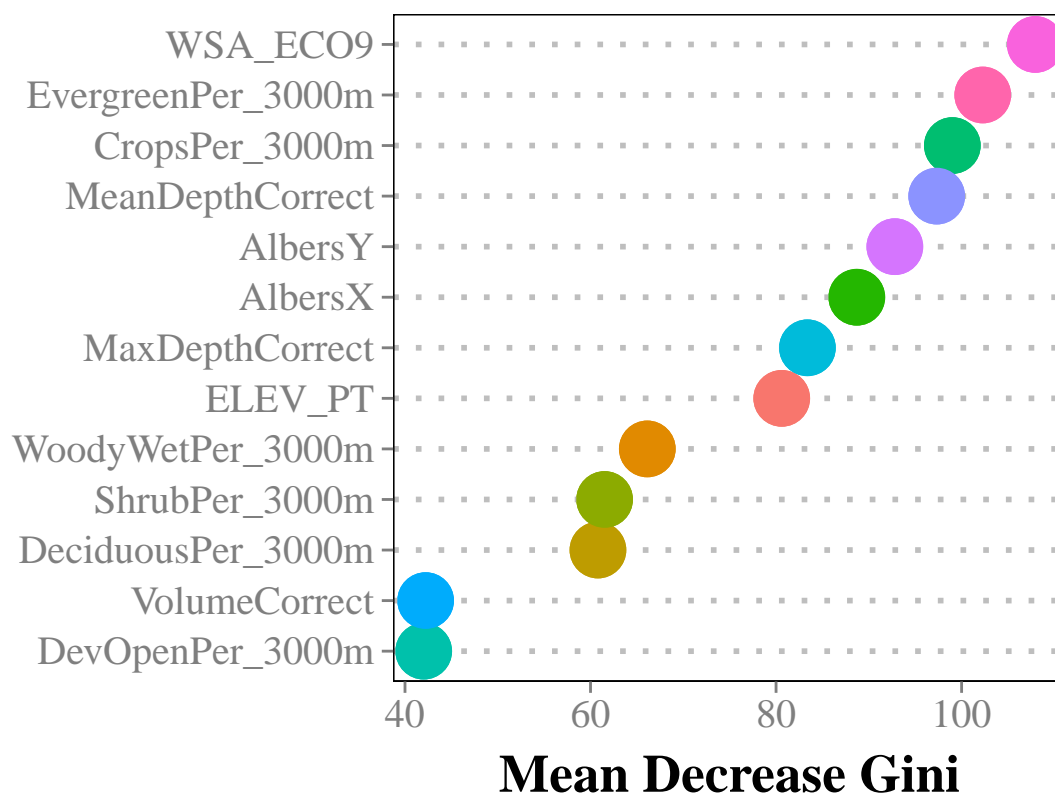


Figure 6: Importance plot for Model 5

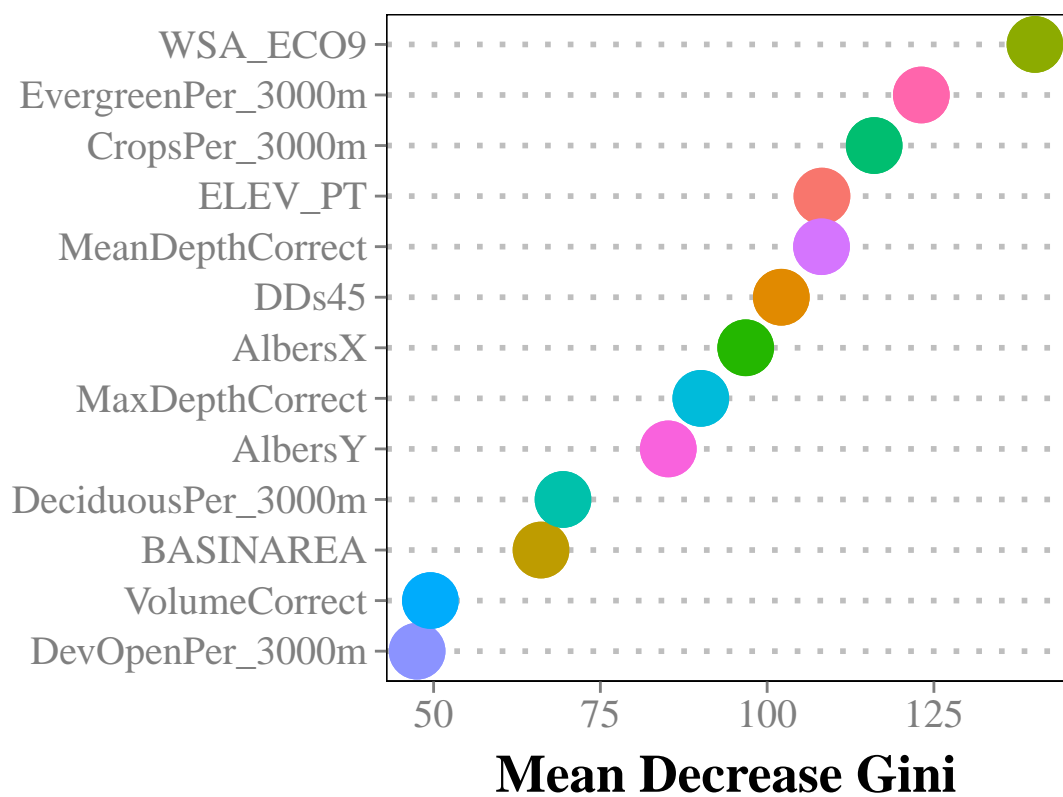


Figure 7: Importance plot for Model 6

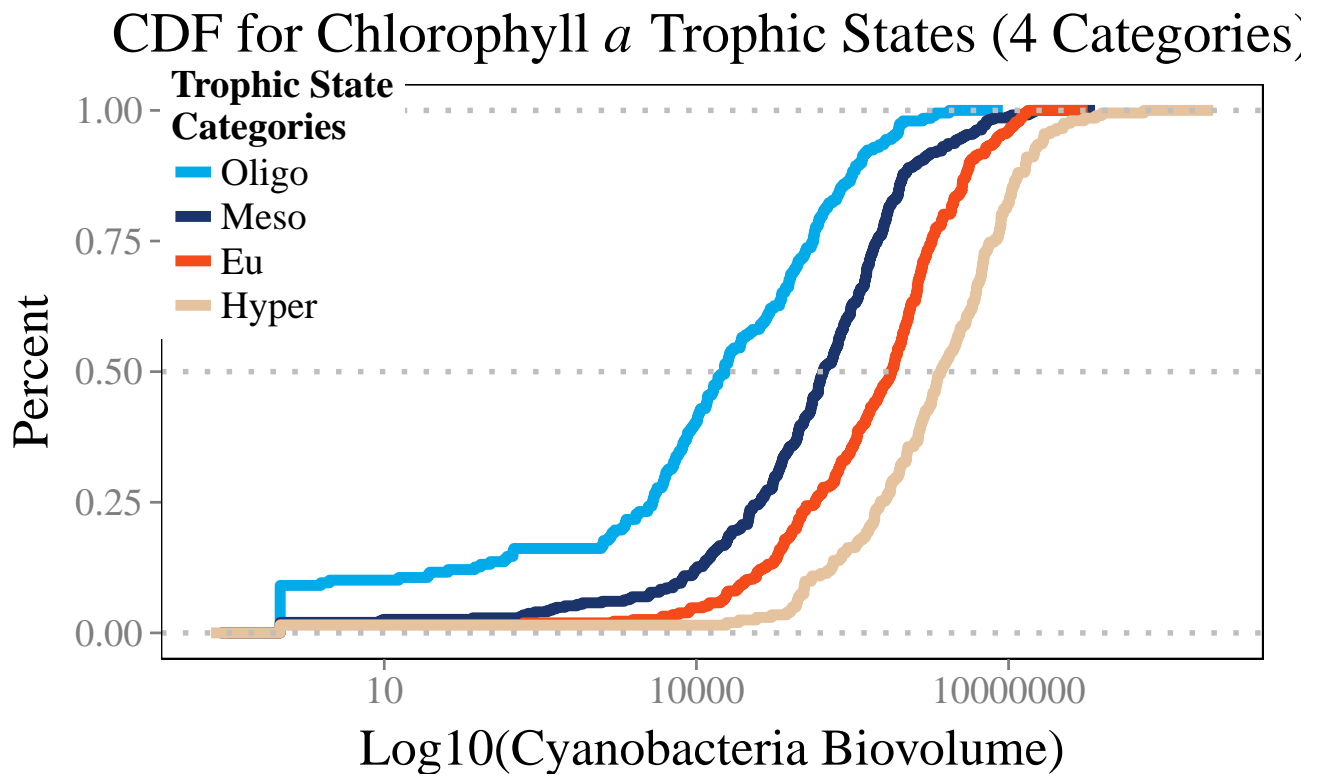


Figure 8: Cumulative distribution function of cyanobacteria biovolume for 4 trophic state classes

CDF for Chlorophyll *a* Trophic States (3 Categories)

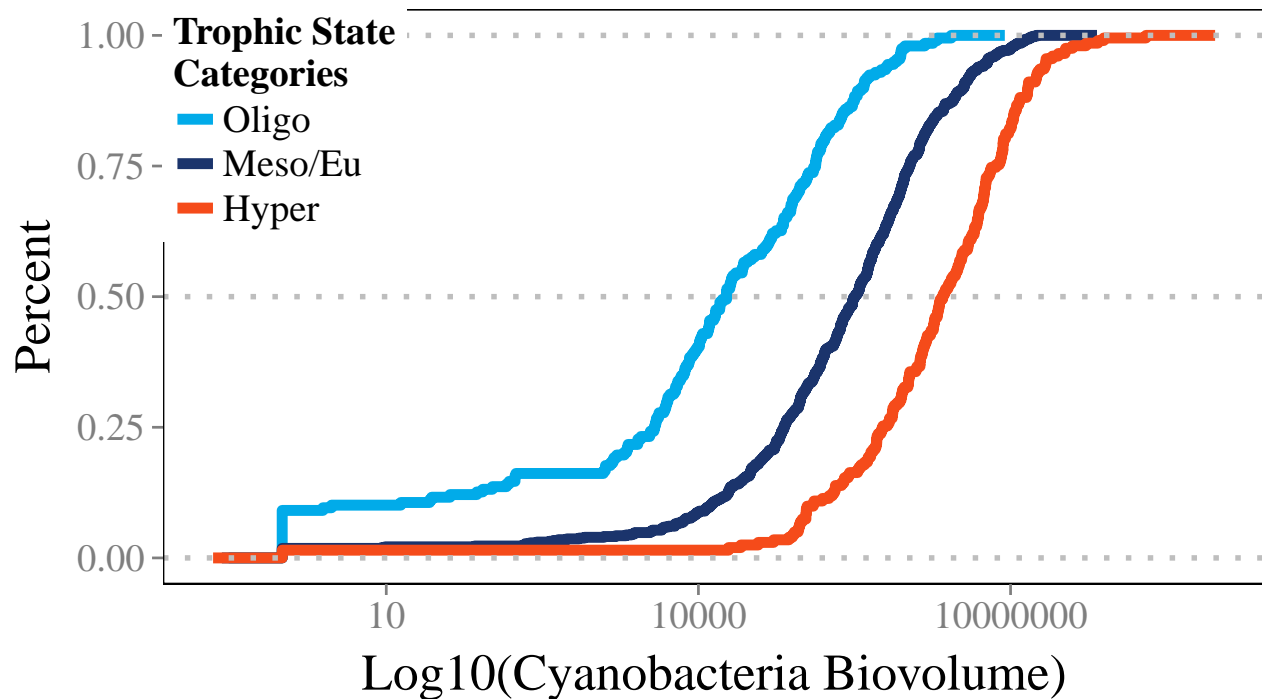


Figure 9: Cumulative distribution function of cyanobacteria biovolume for 3 trophic state classes

CDF for Chlorophyll *a* Trophic States (2 Categories)

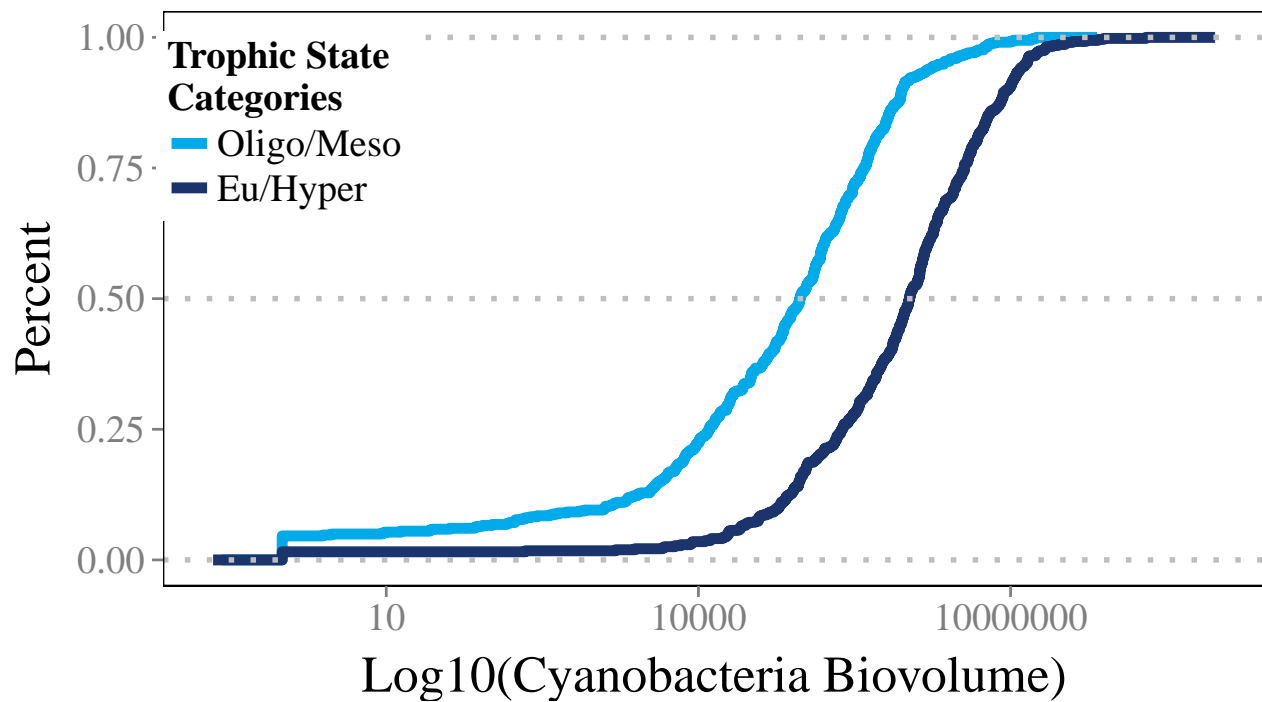


Figure 10: Cumulative distribution function of cyanobacteria biovolume for 2 trophic state classes

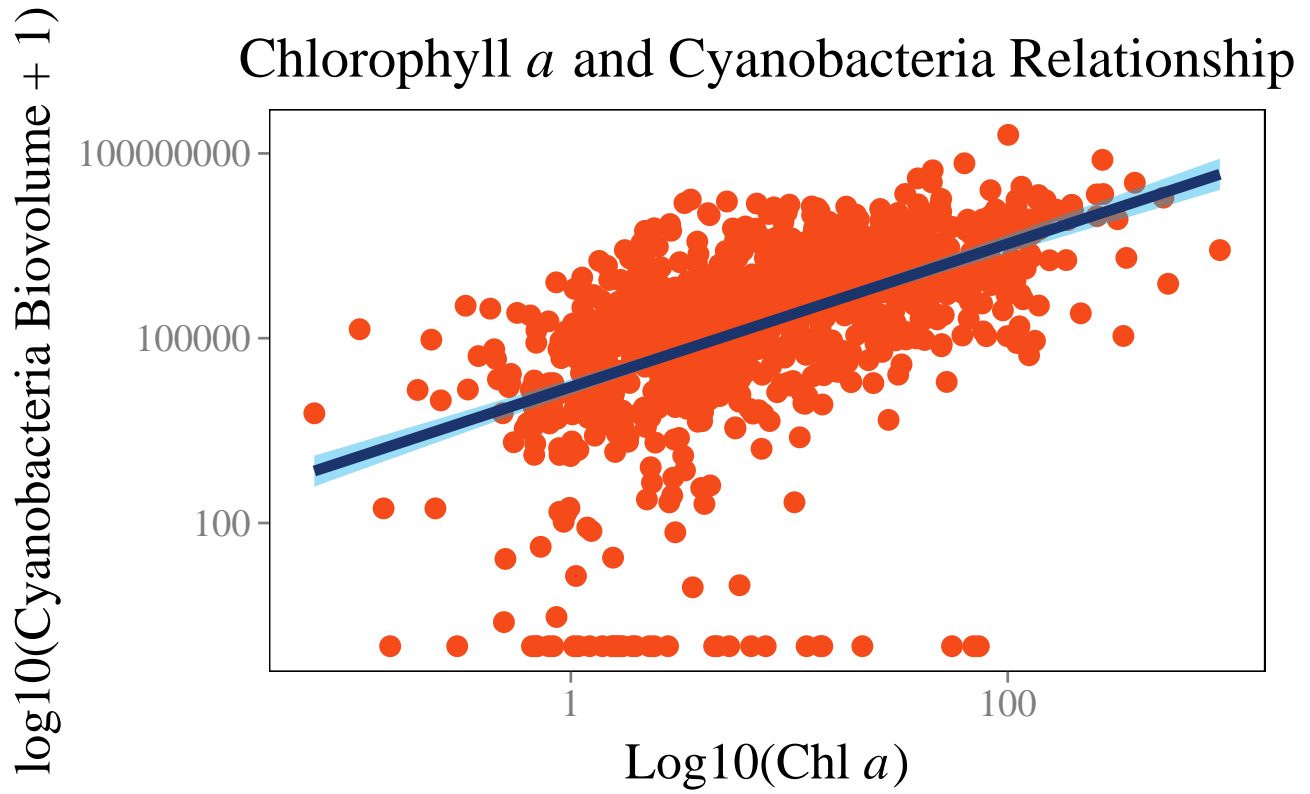


Figure 11: Cholorphyll *a* and cyanobacteria biovolume scatterplot

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

Table 1: Chlorophyll a based trophic state cut-offs

Variable	Percent
NPratio	1.00
NTL	1.00
PTL	1.00
TOC	1.00
TURB	1.00
WSA_ECO9	1.00
K	0.99
ORGION	0.33
DOC	0.22
DEPTHMAX	0.11

Table 2: Variable selection results for Model 1

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 3: Random Forest confusion matrix for Model 1

Variable	Percent
DEPTHMAX	1.00
DOC	1.00
K	1.00
NTL	1.00
ORGION	1.00
PTL	1.00
TOC	1.00
TURB	1.00
WSA_ECO9	1.00
NPratio	0.88
AlbersX	0.58
CropsPer_3000m	0.36
ELEV_PT	0.23
NH4	0.06
AlbersY	0.04
CL	0.03
PH_FIELD	0.02

Table 4: Variable selection results for Model 2

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

Table 5: Random Forest confusion matrix for Model 2

Variable	Percent
K	1.00
NPratio	1.00
NTL	1.00
PTL	1.00
TOC	1.00
TURB	1.00
WSA_ECO9	1.00
ORGION	0.98
DEPTHMAX	0.91
DDs45	0.89
ELEV_PT	0.85
DOC	0.42
AlbersX	0.11
AlbersY	0.03
Na	0.03

Table 6: Variable selection results for Model 3

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

Table 7: Random Forest confusion matrix for Model 3

Variable	Percent
AlbersX	1.00
CropsPer_3000m	1.00
EvergreenPer_3000m	1.00
MeanDepthCorrect	1.00
WSA_ECO9	1.00
AlbersY	0.30
ELEV_PT	0.05
MaxDepthCorrect	0.01

Table 8: Variable selection results for Model 4

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 9: Random Forest confusion matrix for Model 4

Variable	Percent
AlbersX	1.00
AlbersY	1.00
CropsPer_3000m	1.00
EvergreenPer_3000m	1.00
MaxDepthCorrect	1.00
MeanDepthCorrect	1.00
WSA_ECO9	1.00
ELEV_PT	0.97
DeciduousPer_3000m	0.94
ShrubPer_3000m	0.32
WoodyWetPer_3000m	0.18
DevOpenPer_3000m	0.13
VolumeCorrect	0.11

Table 10: Variable selection results for Model 5

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

Table 11: Random Forest confusion matrix for Model 5

Variable	Percent
AlbersX	1.00
AlbersY	1.00
CropsPer_3000m	1.00
DDs45	1.00
ELEV_PT	1.00
EvergreenPer_3000m	1.00
MeanDepthCorrect	1.00
WSA_ECO9	1.00
MaxDepthCorrect	0.98
DeciduousPer_3000m	0.91
DevOpenPer_3000m	0.71
BASINAREA	0.33
VolumeCorrect	0.01

Table 12: Variable selection results for Model 6

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

Table 13: Random forest confusion matrix for Model 6

6 References

- Beaulieu, M., F. Pick, and I. Gregory-Eaves. 2013. Nutrients and water temperature are significant predictors of cyanobacterial biomass in a 1147 lakes data set. *Limnol. Oceanogr* 58:1736–1746.
- 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.
- 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.
- Hollister, J. W. 2014. lakemorpho: Lake morphometry in r.
- Hollister, J. W., and W. B. Milstead. In Preparation. National lake morphometry dataset v1.0.
- Hollister, J. W., W. B. Milstead, and M. A. Urrutia. 2011. Predicting maximum lake depth from surrounding topography. *PLoS ONE* 6:e25764.
- Hollister, J., and W. B. Milstead. 2010. Using gIS to estimate lake volume from limited data. *Lake and Reservoir Management* 26:194–199.
- 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 & Remote Sensing* 70:829–840.
- Liaw, A., and M. Wiener. 2002. Classification and regression by randomForest. *R News* 2:18–22.
- 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., 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

158 on freshwater, marine, and terrestrial ecosystems. *Environmental pollution* 100:179–196.

159 USEPA. 2009. National lakes assessment: a collaborative survey of the nation’s lakes. ePA 841-r-09-001.

160 Office of Water; Office of Research; Development, US Environmental Protection Agency Washington,

161 DC.

162 Xian, G., C. Homer, and J. Fry. 2009. Updating the 2001 national land cover database land cover clas-

163 sification to 2006 by using landsat imagery change detection methods. *Remote Sensing of Environment*

164 113:1133–1147.