

Introduction

Cyanobacteria are an important taxonomic group associated with harmful algal blooms in lakes. Understanding the drivers of cyanobacteria presence has important implications for lake management and for the protection of human and ecosystem health. Chlorophyll *a* concentration, a measure of the biological productivity of a lake, is one such driver that is largely determined by nutrient inputs. As nutrient inputs increase, productivity increases and lakes transition from low trophic state (e.g. oligotrophic) to higher trophic states (e.g. hypereutrophic). These broad trophic state classifications are associated with ecosystem health and ecosystem services and disservices. Thus, models of trophic state might be used to predict things like cyanobacteria. In the preliminary work reported here, we:

- 1. Build and assess models of lake trophic state predictions
- 2. Assess ability to predict trophic state in lakes without available *in situ* water quality data
- 3. Explore association between cyanobacteria and trophic state.

Methods

Data

We utilize four primary sources of data for this study.

1. National Lakes Assessment (NLA) 2007: Using consistent methods and metrics, the NLA collected data from ~1150 lakes across the conterminous United States, on biophysical measures of lake water quality and habitat (Map 1). For this analysis we primarily examined the water quality measurements from the NLA (USEPA 2009).

2. National Land Cover Dataset (NLCD) 2006: The NLCD is a national land use/land cover dataset. We calculated total land use/land cover and total percent impervious surface within a 3 kilometer buffer of each lake to examine larger landscapelevel effects (Homer et al. 2004; Xian, Homer & Fry 2009).

3. Modeled lake morphometry: Various measures of lake morphometry (i.e. depth, volume, fetch, etc.) are important in understanding lake productivity, yet are often difficult to obtain for large numbers of lakes. Modeled estimates solved this problem. (Hollister & Milstead 2010; Hollister, Milstead, & Urrutia 2011; Hollister 2013; Hollister & Milstead in prep).

4. Estimated Cyanobacteria Biovolumes: Measuring of cyanobacteria dominance is best done with biovolume as there is great size variability within and among taxa. Beaulieu et al. (2013) used literature values to estimate biovolumes for the taxa in the NLA. They shared these data with our colleagues and we have summed that information on a per-lake basis.

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 is recursively partitioned according to a given random subset of predictor variables and completely grown without pruning. With each new tree, both the sample data and predictor variable subset is randomly selected.

Random forests are able to handle numerous correlated variables without a decrease in prediction accuracy. But large numbers of related variables can reduce accuracy and lead to over-fitting. This problem, faced in gene selection, has been addressed with a variable selection method based on random forest (Díaz-Uriarte & De Andres 2006). Using 100 iterations of varSelRF in R, we determine how often our variables are included in a final model (Diaz-Uriarte 2010). The most commonly selected variables (i.e. the reduced model) are used to develop a final random forest model (Liaw & Wiener 2002). From these random forests we collect a consensus prediction and calculate a confusion matrix and summary stats.

Model Details

Using a combination of the varSelRF and randomForest we ran models for six combinations of variables and trophic state classifications. These combinations included different combinations of the Chlorophyll *a* trophic states (Table 2) with all variables or with the GIS variables (i.e. no *in situ* information). The six model combinations were:

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

Trophic State (4)	Trophic State (3)	Trophic State (2)	Cut-
oligo	oligo	oligo/meso	<= 0.
meso	meso/eu	oligo/meso	>2-7
eu	meso/eu	eu/hyper	>7-30
hyper	hyper	eu/hyper	>30

Table 1. Chlorophyll *a* (µg/l) trophic state categories used as dependent variable

References

1.) Beaulieu, Marieke, Frances Pick, and Irene Gregory-Eaves. 2013. "Nutrients and Water Temperature Are Significant Predictors of Cyanobacterial Biomass in a 1147 Lakes Data Set." Limnol. Oceanogr 58 (5): 1736–1746.

2.) Breiman, Leo. 2001. "Random Forests." Machine Learning 45 (1): 5–32.

3.) Diaz-Uriarte, Ramon. 2010. VarSelRF: Variable Selection Using Random Forests. http://CRAN.R-project.org/package=varSelRF. 4.) Díaz-Uriarte, Ramón, and Sara Alvarez De Andres. 2006. "Gene Selection and Classification of Microarray Data Using Random Forest." BMC Bioinformatics 7 (1): 3.

5.) Hollister, Jeffrey. 2013. Lakemorpho: Lake Morphometry in R. http://www.github.com/USEPA/lakemorpho.

6.) Hollister, Jeffrey W, and W Bryan Milstead. *in prep*. "National Lake Morphometry Dataset V1.0."

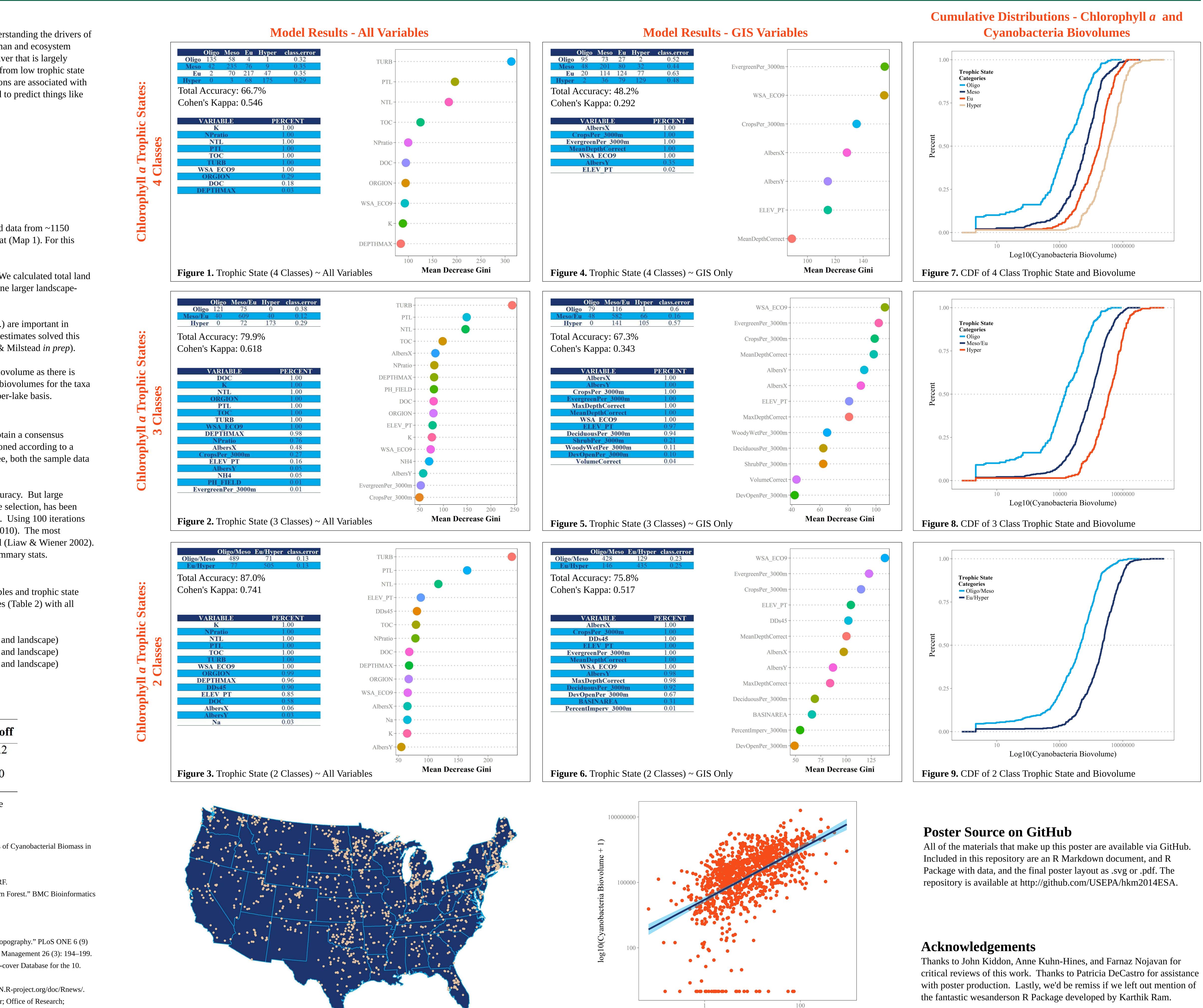
7.) Hollister, Jeffrey W., W. Bryan Milstead, and M. Andrea Urrutia. 2011. "Predicting Maximum Lake Depth from Surrounding Topography." PLoS ONE 6 (9) 8.) Hollister, Jeffrey, and W Bryan Milstead. 2010. "Using GIS to Estimate Lake Volume from Limited Data." Lake and Reservoir Management 26 (3): 194–199. 9.) Homer, Collin, Chengquan Huang, Limin Yang, Bruce Wylie, and Michael Coan. 2004. "Development of a 2001National Land-cover Database for the 10. United States." Photogrammetric Engineering & Remote Sensing 70 (7):829–840.

10.) Liaw, Andy, and Matthew Wiener. 2002. "Classification and Regression by RandomForest." R News 2 (3): 18–22.http://CRAN.R-project.org/doc/Rnews/. 11.) USEPA. 2009. "National Lakes Assessment: a Collaborative Survey of the Nation's Lakes. EPA 841-r-09-001." Office of Water; Office of Research; Development, US Environmental Protection Agency Washington, DC.

12.) Xian, George, Collin Homer, and Joyce Fry. 2009. "Updating the 2001 National Land Cover Database Land Cover Classification to 2006 by Using Landsat Imagery Change Detection Methods." Remote Sensing of Environment113 (6): 1133–1147.

Expanding Models of Lake Trophic State to Predict Cyanobacteria in Lakes: A Data Mining Approach Jeffrey W. Hollister, W. Bryan Milstead, and Betty J. Kreakie

U.S. Environmental Protection Agency, Office of Research and Development, National Health and Environmental Effects Research Laboratory, Atlantic Ecology Division, Narragansett, RI 02882



Log10(Chl a)**Figure 10.** Relationship between Chlorophyll *a* and Cyanobacteria