# Application of a continuous lake trophic state index on lakes with limited data

**Farnaz Nojavan A.[1], Betty J. Kreakie[2], Jeffrey W. Hollister[2], and Song S. Qian[1]**

[1]**Farnaz's new affiliation address**

[2]**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**

[3]**Department of Environmental Sciences, The University of Toledo, Toledo, OH, United States**

Corresponding author:

Farnaz Nojavan A.[1]

Email address: `f.nojavan@gmail.com`

## ABSTRACT

Lake trophic state indices have long been used to provide a measure of the trophic state of lakes. Over time it has been determined that these indices perform better when they utilize mutliple metrics and provide a continuous measurement of trophic state. We utilize such a method for trophic state that is based upon a Proportional Odds Logistic Regression (POLR) model and extend this model with a Bayesian multilevel model that predicts nutrient concentrations from universally available GIS data. This Bayesian multilevel model provides relatively accurate measures of trophic state and has an overall accuraccy of 60%. The approach illustrates a method for estimating a continuous, mutli-metric trophic state index for any lake in the United States. Future improvements to the model will focus on improving overall accuracy and use variables that are more sensitive to change over time.

## INTRODUCTION

In this brief research note, we extend a model for estimating a continuous mulit-metric trophic state index described by Nojavan et al. (n.d.). This model uses lake elevation and *in situ* measurements of total nitrogen, total phosphorus, and secchi depth to provide a continuous index of trophic state. The drawback of the developed POLR model is the cost of monitoring multiple predictor variables (e.g., nutrients). This is addressed in the extended application by linking nitrogen and phosphorus to universally available GIS variables. The goal of the extended POLR model is to allow prediction of the trophic state of all lakes (i.e. lakes with limited field data) in the United States.

## METHODS

We present the extended application of the developed POLR model using a Bayesian multilevel model. Our modeling work flow is as follows:

1. Develop a random forest model, using R's `randomForest` package, with 5000 trees using only GIS variables to identify the best predictor variables for nitrogen and phosphorus.
2. Develop the extended application model (the Bayesian multilevel model) using R's `rjags` package to run Just Another Gibbs Sampler (JAGS) from inside of R. JAGS is a program for simulation and analysis of Bayesian hierarchical models using Markov Chain Monte Carlo (MCMC).
3. Assess the performance of the extended application model using a hold-out validation method (90% training set, 10% evaluation set).

We link nitrogen and phosphorus in the POLR model to a separate nutrient model built from universally available GIS data, thereby, avoiding the need for nitrogen and phosphorus data, costly variables to

44 measure for all lakes. The number of variables for each response variable, nitrogen or phosphorus, was
45 decided using random forest model's variable selection plots(Hollister, Milstead, and Kreakie 2016).

## RESULTS AND DISCUSSION

47 Selected GIS variables for nitrogen and phosphporus were initially screened with variable selection plots
48 (Figures~1 &~2).The figures show model mean squared error as a function of the number of variables.
49 The best representation of nitrogen and phosphorus could be achieved using three variables, adding more
50 than three variables had incremental ($< 0.1$) impact on root mean square error. The three most important
51 variables were ecoregion, % evergreen forest, and latitude. The random forest models provided estimates
52 of variable importance for nitrogen and phosphorus and the results are reported in figures~3 &4.

53 Figure~5 represents the regression models. The extended POLR model is grouped into two blocks
54 (gray shaded rectangles). The trophic state classification regression, the POLR model in the lower block,
55 includes nitrogen, phosphorus, secchi disk, and elevation as predictors. The nutrient model, in the upper
56 block, estimates the means of nitrogen and phosphorus based on ecoregion, % evergreen forest, and
57 latitude. The two blocks are connected through the estimated means of nitrogen ($\mu_{Nitrogen}$) and phosphorus
58 ($\mu_{Phosphorus}$) to form the combined model which enables trophic state classification for all lakes without
59 the costly sampling requirement. The relationship between nitrogen, phosphorus, and their predictors was
60 examined using multilevel linear regression models. The standard deviation of the normal distribution, as
61 well as each parameter in the regression model, were then assigned non-informative prior distributions
62 (uniform, or nearly so, to allow the information from the likelihood to be interpreted probabilistically).

63 The three selected variables, latitude, eco-region, and % evergreen forest, appear to be capturing
64 patterns of total nutrient concentration at three different spatial scales. Figures~6 &~7 depict the partial
65 dependency plot for latitude, the marginal effect of latitude on the predicted outcome of nitrogen or
66 phosphorus in the random forest model. For example for predicted total nitrogen, high concentrations
67 in the northern and southern extremes of the continental US and the lowest predicted concentrations
68 correspond to the mid-latitudes. The ecoregion variable represents an intermediate scale among these
69 three variables and represents the variation between the regions. Finally, the % evergreen variable was
70 summarized within a 3 kilometer buffer around each lake and is presumably summarizing more local land
71 use decisions that are adjacent to lakes.

72 As mentioned, the extension of the developed POLR model uses eco-region, latitude, and watershed
73 level % evergreen forest as predictors for nitrogen and phosphorus. This contrasts with prior trophic
74 state classification models that are applied to all lakes, regardless of the differences across scale. Lake
75 trophic index, and hence lake trophic classes, should be calculated differently in different eco-regions to
76 accommodate variation in landform and climate characteristics and our proposed model and extension
77 bares this out by identifying and including and eco-regional approach to quantifying trophic state.
78 Furthermore, the developed multilevel model structure can be further expanded to lake-specific trophic
79 state index, upon availability of multiple measurements for each lake.

80 Mathematically, the models were set up as follows:

$$\text{Nitrogen}_{ij} \sim \mathcal{N}(\mu_{Nitrogen_{ij}}, \sigma^2_{Nitrogen}) \tag{1}$$

81 where $\mu_{Nitrogen_{ij}} = X_{Nitrogen}B$, $X_{Nitrogen}$ is the matrix of predictors, and $B$ is the vector of coefficients.
82 *Nitrogen$_{ij}$ is the $\{i\}$th nitrogen observation in the $\{j\}$th ecoregion.*

$$\text{Phosphorus}_{ij} \sim \mathcal{N}(\mu_{Phosphorus_{ij}}, \sigma^2_{Phosphorus}) \tag{2}$$

83 *where $\mu_{Phosphorus_{ij}} = X_{Phosphorus}\Gamma$, $X_{Phosphorus}$ is the matrix of predictors, and $\Gamma$ is the vector of coefficients.*
84 *Phosphorus$_{ij}$ is the $\{i\}$th phosphorus observation in the $\{j\}$th ecoregion.\*

85 *The overall accuracy of the extended POLR model was 0.6 and the balanced accuracies were*
86 *0.78, 0.77, 0.69, 0.68 for oligotrophic, mesotrophic, eutrophic, and hypereutrophic classes, respectively*
87 *(Table~1). Table~2 shows the confusion matrix for the extended POLR model.*

88 *The extended POLR model calculates lake trophic index and classes differently for different eco-*
89 *regions. Please refer to Table~1 for varying coefficients in different eco-regions. For example, eco-regions*
90 *3, 6, and 5, corresponding to Northern Plains, Temperate Plains, and Southern Plains, have the highest*

*positive coefficients for nitrogen. Hence, nitrogen plays a significant role in moving the trophic state index*
*and class toward the eutrophic/hypereutrophic side of the trophic continuum. Further Table~1 shows the*
*coefficients for latitude and % evergreen. We included these predictors as they were selected as important*
*variables by the random forest model. They may not help predictions dramatically but they do not hurt the*
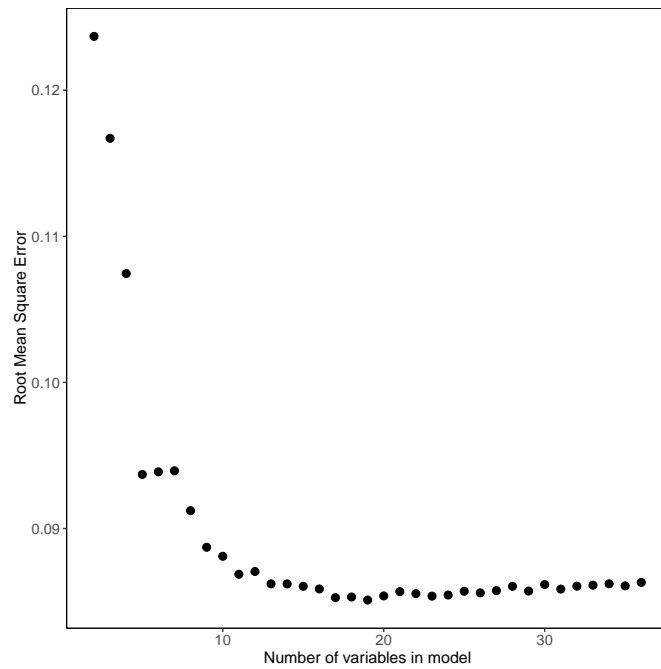*results.*

96 ## *TABLES*

**Table 1.** Coefficients for the extended POLR model.

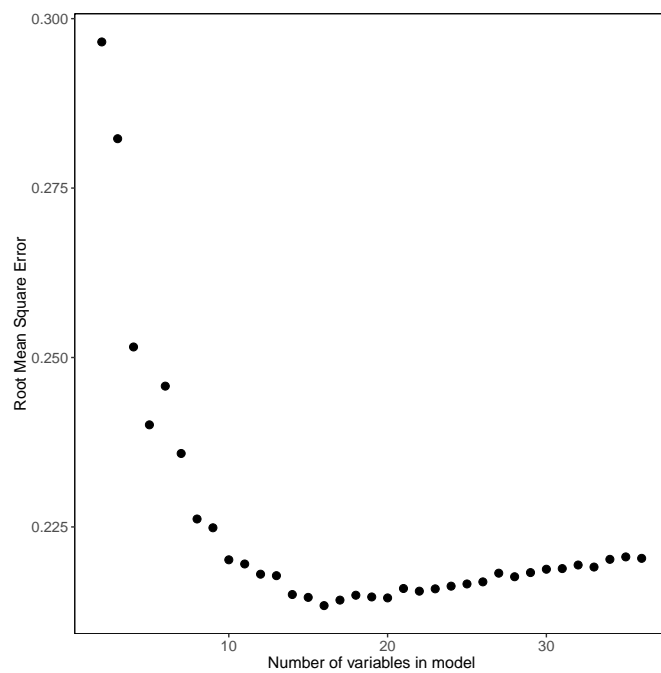|  |  | Mean | Standard Deviation |
|---|---|---|---|
| *Cutoff points/Thresholds* | $C_{Oligo|Meso}$ | -156.60 | 44.04 |
|  | $C_{Meso|Eu}$ | -6.18 | 8.29 |
|  | $C_{Eu|Hyper}$ | 121.32 | 35.04 |
| *POLR model coefficients* | $\alpha_{Elevation}$ | -40.20 | 12.86 |
|  | $\alpha_{Nitrogen}$ | -44.33 | 29.29 |
|  | $\alpha_{Phosphorus}$ | 165.90 | 46.96 |
|  | $\alpha_{\text{Secchi Disk Depth}}$ | 0.18 | 5.23 |
| *Multilevel model coefficients for nitrogen* | $\beta_{\%Evergreen}$ | 0.00 | 0.01 |
|  | $\beta_{Ecoregion_1}$ | 0.34 | 0.13 |
|  | $\beta_{Ecoregion_2}$ | -0.78 | 0.12 |
|  | $\beta_{Ecoregion_3}$ | 0.96 | 0.15 |
|  | $\beta_{Ecoregion_4}$ | -0.37 | 0.10 |
|  | $\beta_{Ecoregion_5}$ | 0.59 | 0.10 |
|  | $\beta_{Ecoregion_6}$ | 0.68 | 0.09 |
|  | $\beta_{Ecoregion_7}$ | -0.01 | 0.10 |
|  | $\beta_{Ecoregion_8}$ | -1.00 | 0.10 |
|  | $\beta_{Ecoregion_9}$ | 0.11 | 0.12 |
|  | $\beta_{Latitude}$ | 0.11 | 0.05 |
| *Multilevel model coefficients for phosphorus* | $\gamma_{\%Evergreen}$ | -0.00 | 0.01 |
|  | $\gamma_{Ecoregion_1}$ | 0.40 | 0.09 |
|  | $\gamma_{Ecoregion_2}$ | -0.90 | 0.09 |
|  | $\gamma_{Ecoregion_3}$ | 0.73 | 0.11 |
|  | $\gamma_{Ecoregion_4}$ | -0.38 | 0.08 |
|  | $\gamma_{Ecoregion_5}$ | 0.53 | 0.08 |
|  | $\gamma_{Ecoregion_6}$ | 0.71 | 0.07 |
|  | $\gamma_{Ecoregion_7}$ | -0.32 | 0.08 |
|  | $\gamma_{Ecoregion_8}$ | -0.69 | 0.08 |
|  | $\gamma_{Ecoregion_9}$ | 0.07 | 0.09 |
|  | $\gamma_{Latitude}$ | -0.03 | 0.03 |
| *Logistic distribution's scale parameter* | $\sigma$ | 75.64 | 21.27 |

**Table 2.** Confusion matrix for multilevel POLR model. Each element of the matrix is the number of cases for which the actual state is the row and the predicted state is the column.

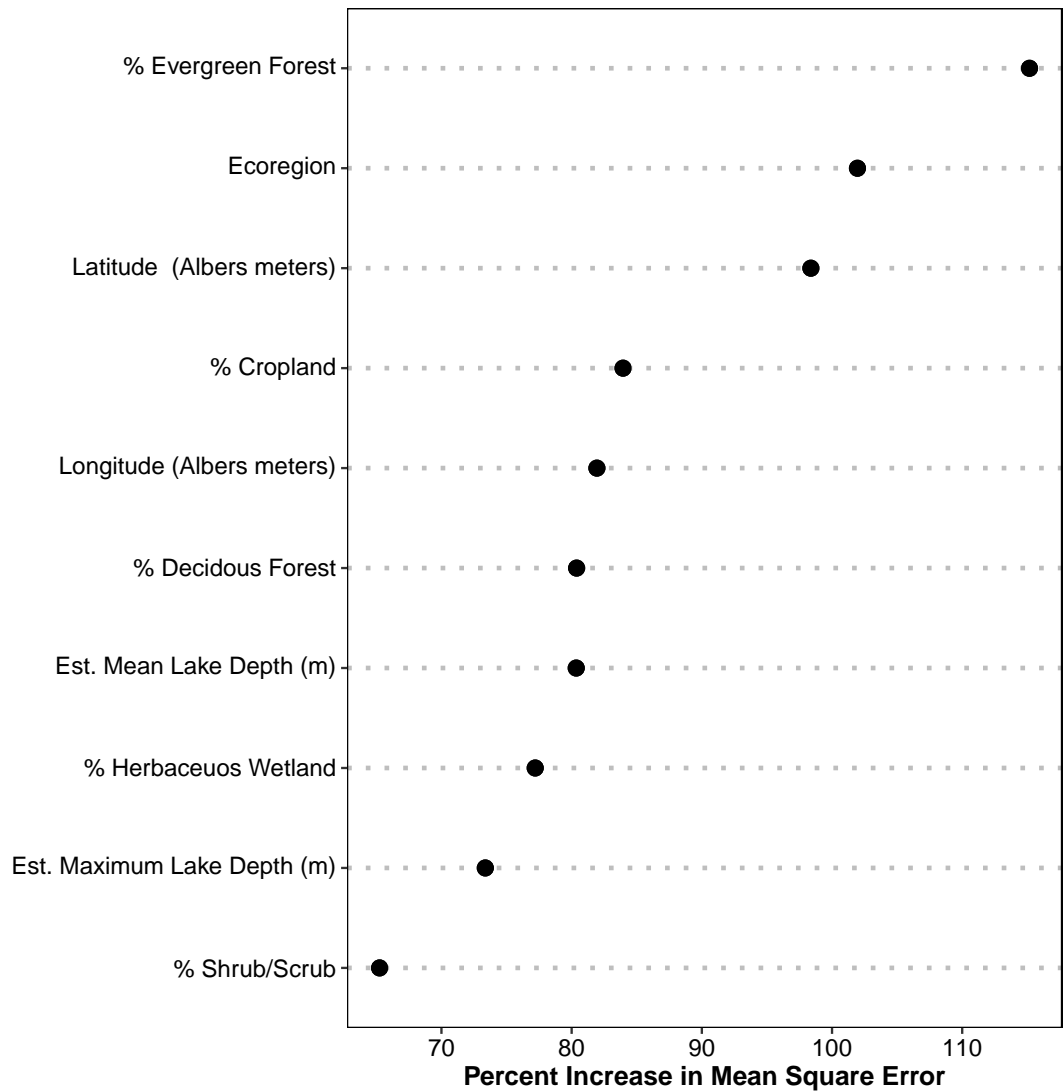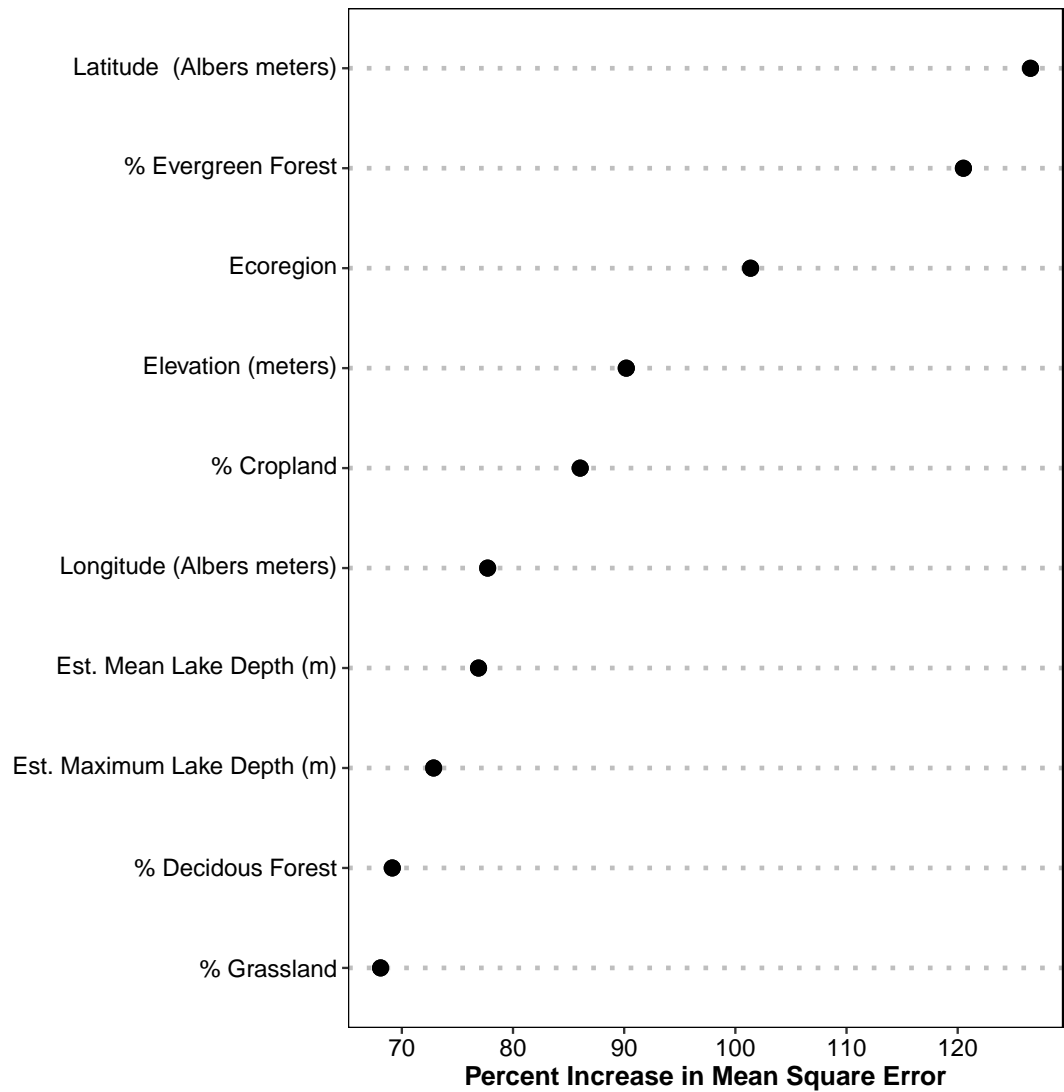|  | Oligo | Meso | Eu | Hyper |
|---|---|---|---|---|
| Oligo | 5 | 3 | 0 | 0 |
| Meso | 3 | 12 | 7 | 1 |
| Eu | 0 | 0 | 16 | 10 |
| Hyper | 0 | 1 | 3 | 9 |

97 ## *FIGURES*

**Figure 1.** Random Forest model's output for nitrogen with GIS only variables as predictors. Shows model mean squared error as a function of the number of variables.
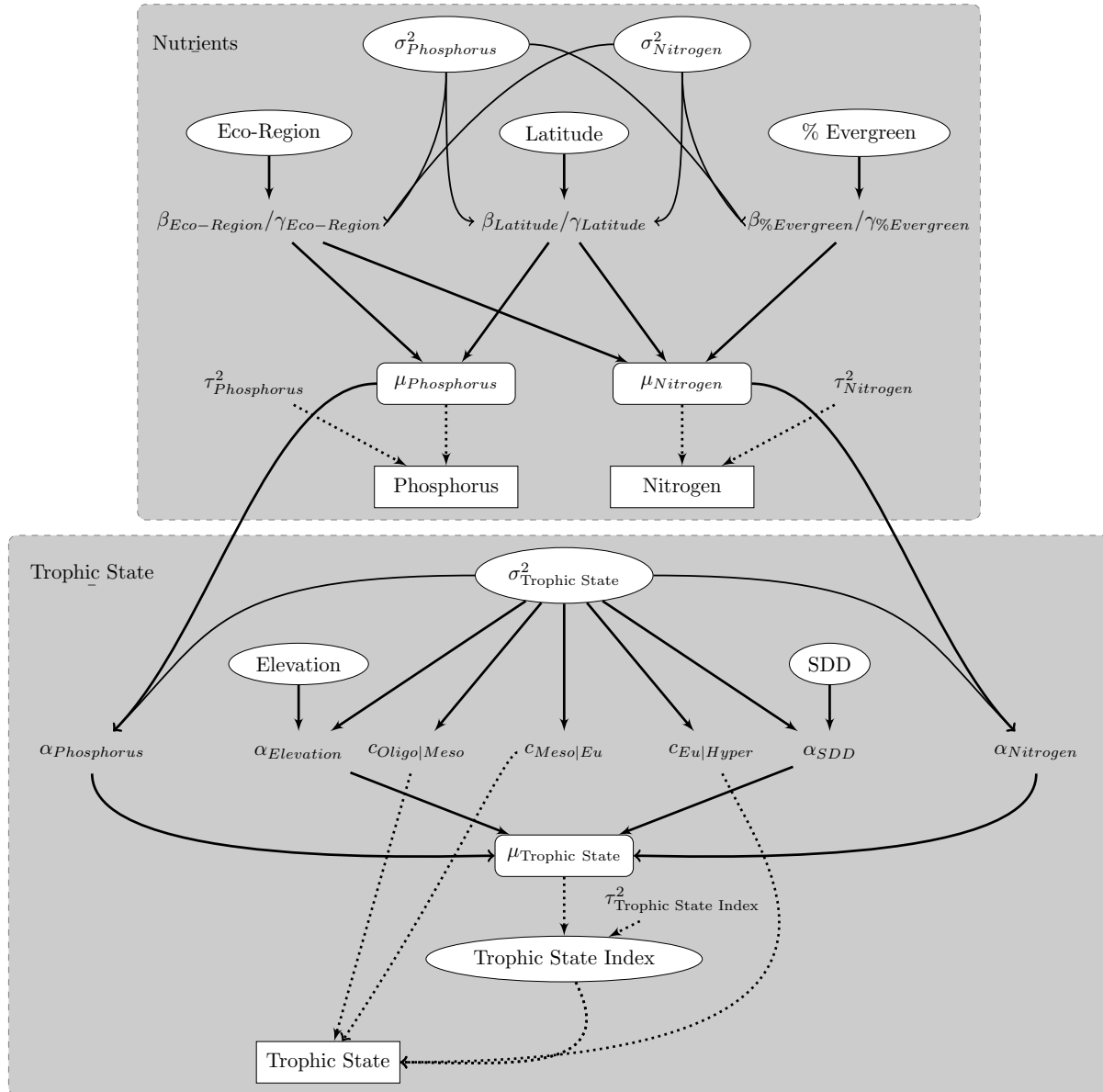


**Figure 2.** Random Forest model's output for phosphorus with GIS only variables as predictors. Shows model mean squared error as a function of the number of variables.
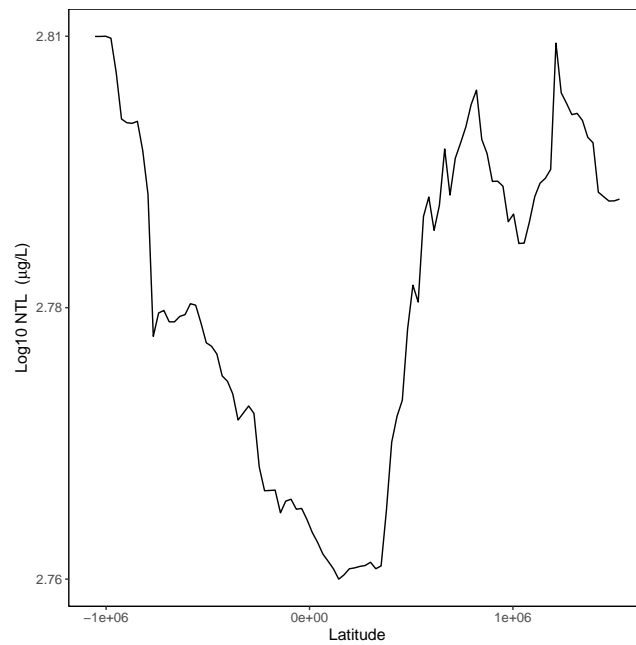
**Figure 3.** Random Forest model's output for nitrogen predictors. Importance plot for GIS variables. Shows percent increase in mean squared error. Higher values of percent increase in mean squared error indicates higher importance.
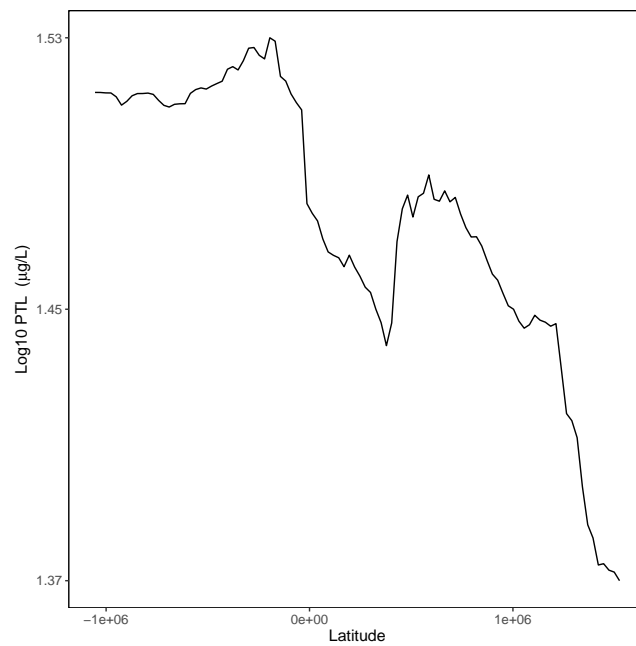
**Figure 4.** Random Forest model's output for nitrogen predictors. Importance plot for GIS variables. Shows percent increase in mean squared error. Higher values of percent increase in mean squared error indicates higher importance.

**Figure 5.** Directed Acyclic Graphical (DAG) model. The lower box depicts the POLR model with its four predictors of secchi disk depth (SDD), elevation, nitrogen, and phosphorus. The upper box is the extension to the POLR model to predict nitrogen and phosphorus using universally available GIS variables.

**Figure 6.** Partial dependency plot for predicted total nitrogen over the range latitude: the effect of latitude on predcited total nitrogen when the rest of the predictors are held constant.



**Figure 7.** Partial dependency plot for predicted total phosphorus over the range latitude: the effect of latitude on predicted total phosphorus when the rest of the predictors are held constant.

## REFERENCES

Hollister, Jeffrey W, W Bryan Milstead, and Betty J Kreakie. 2016. "Modeling Lake Trophic State: A Random Forest Approach." Ecosphere 7 (3).

Nojavan A., Farnaz, Betty J. Kreakie, Hollister Jeffrey W., and Song Qian. n.d. "Rethinking the Lake Trophic State Index." PeerJ.