

# The Implications of Simpson's Paradox for Cross-Scale Inference Among Lakes

Song S. Qian<sup>a,\*</sup>, Craig A. Stow<sup>b</sup>, Farnaz Nojavan A.<sup>c</sup>, Joseph Stachelek<sup>f</sup>,  
Yoonkyung Cha<sup>d</sup>, Ibrahim Alameddine<sup>e</sup>, Patricia Soranno<sup>f</sup>

<sup>a</sup>*Department of Environmental Sciences, University of Toledo, 2801 W. Bancroft Street, MS# 604, Toledo, OH, USA.*

<sup>b</sup>*Great Lakes Environmental Research Laboratory, National Oceanic and Atmospheric Administration, Ann Arbor, MI, USA*

<sup>c</sup>*Center for Industrial Ecology, Yale University, New Haven, CT, USA*

<sup>d</sup>*School of Environmental Engineering, University of Seoul, Seoul, South Korea*

<sup>e</sup>*Department of Civil and Environmental Engineering, American University of Beirut, Beirut, Lebanon*

<sup>f</sup>*Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI, USA*

---

## Abstract

Using cross-sectional data for making ecological inference started as a practical means of pooling data to enable meaningful empirical model development. For example, limnologists routinely use sample averages from numerous individual lakes to examine patterns across lakes. The basic assumption behind the use of cross-lake data is often that responses within and across lakes are identical. As data from multiple study units across a wide spatiotemporal scale are increasingly accessible for researchers, an assessment of this assumption is now feasible. In this study, we demonstrate that this assumption is usually unjustified, due largely to a statistical phenomenon known as the Simpson's paradox. Through

---

\*Corresponding author, Phone: 419.530.4230, Fax: 419.530.4421

Email address: [song.qian@utoledo.edu](mailto:song.qian@utoledo.edu) (Song S. Qian)

Preprint submitted to Water Research

June 17, 2019

comparisons of a commonly used empirical model of the effect of nutrients on algal growth developed using several data sets, we discuss the cognitive importance of distinguishing factors affecting lake eutrophication operating at different spatial and temporal scales. Our study proposes the use of the Bayesian hierarchical modeling approach to properly structure the data analysis when data from multiple lakes are employed.

16 *Keywords:* NLA, LAGOSSE, multilevel/hierarchical model, chlorophyll a

---

## 17 **1. Introduction**

18 Ecologists have a long history of using data from multiple lakes,  
19 summarized at various levels of spatial and temporal aggregation, to  
20 estimate empirical models (Vollenweider , 1968, 1975, Schindler, 1977,  
21 Wagner et al., 2011). Dillon and Rigler (1973) set an early precedent using  
22 reported sample averages from a combination of 46 North American lakes,  
23 lake years, and segments of lakes to estimate a simple linear regression  
24 model relating chlorophyll a (*chl*a) concentration to total phosphorus (TP)  
25 concentration. Numerous papers followed, applying regression approaches  
26 to estimate similar models using data from other lakes, sometimes  
27 comparing their estimated equations to the equation obtained by Dillon  
28 and Rigler (Jones and Bachmann, 1976, Canfield and Bachmann, 1981,  
29 Canfield, 1983, Prepas and Trew, 1983). The practice of estimating models  
30 using data from multiple lakes is common, fostered by increases in

31 computational capacity and corresponding advances in statistical software  
32 which now facilitates the estimation of nonlinear models, using large data  
33 sets (Filstrup et al., 2014).

34 These approaches are typically based on an implicit assumption that  
35 the *chla* and TP means from multiple lakes can be described by a  
36 dose-response equation (e.g., McCauley et al. (1989)) such as:

$$\log(\mu_{Chla}) = \beta_0 + \beta_1 \log(\mu_{TP}) + \varepsilon \quad (1)$$

37 where  $\mu_{Chla}$  is the mean of *chla* concentration for a specified time period  
38 (such as summer of a particular year) and lake (or lake segment),  $\mu_{TP}$  is the  
39 mean TP concentration for a corresponding, but not necessarily the same,  
40 time period (spring TP may be related to summer *chla*, for example),  $\beta_0$   
41 and  $\beta_1$  are the intercept and slope parameters, respectively, and  $\varepsilon$  is the  
42 model error term usually assumed to be normally distributed with a  
43 constant variance (Qian, 2016). Because the underlying “true” mean values  
44 are always unknown, sample averages are typically used as surrogates,  
45 although occasionally sample medians have been used (Reckhow 1988).  
46 This regression-based modeling approach has influenced lake management  
47 practices beyond the modeling of the *chla*-nutrient relationship. For  
48 example, Yuan and Pollard (2017) used data from the National Lake  
49 Assessment (NLA), a cross-lake data set including randomly selected lakes  
50 in all 48 contiguous states of the United States (Pollard et al., 2018), to

51 develop a dose-response model to describe the relationship between  
52 microcystin (MC) concentration and total nitrogen (TN) concentration.  
53 The resulting model was used to propose a national nitrogen criterion for  
54 controlling harmful algal blooms.

55 The implicit premise of this approach is that a relationship estimated  
56 using sample averages from many lakes can be applied to set criteria for  
57 individual lakes, because criteria compliance assessment is typically  
58 lake-specific. However, we see two potential problems with this supposition:

- 59 1. Using sample averages as surrogates for the “true,” unknown means,  
60 violates two assumptions of regression analysis: the variance of the  
61 response variable is constant and the predictor variables are observed  
62 without error. On the one hand, violating the equal variance  
63 assumption makes an estimated parameter and model error variances  
64 ambiguous; it is unclear what uncertainty bands calculated from these  
65 values, such as 95% confidence or prediction intervals, represent. On  
66 the other hands, the consequence of violating the observation error  
67 assumption has been well-studied; it is widely recognized that this  
68 “errors-in-variables” problem causes slope coefficient estimators to be  
69 biased toward zero (Fuller, 1987, Carroll et al., 2006).
- 70 2. Lake-specific factors may cause individual lakes to exhibit differing  
71 stressor-response relationships (Jones and Bachmann, 1976, Wagner  
72 et al., 2011, Malve and Qian, 2006). Using aggregated measures, such

73 as sample averages to estimate among-lake relationships can produce  
 74 results that poorly represent the individual lakes in the analysis. In  
 75 extreme cases, the sign of the estimated slope parameter can be  
 76 reversed (Figure 1), an example of Simpson’s Paradox (Simpson,  
 77 1951). Clearly, such a model should not be used to develop  
 78 lake-specific management strategies (Smith and Shapiro, 1981,  
 79 Reckhow, 1993, Liang et al., 2018).

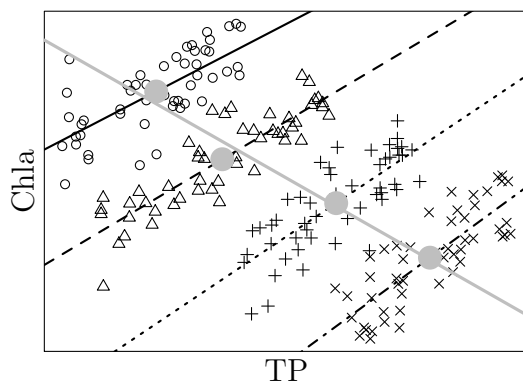


Figure 1: Hypothetical data from four lakes illustrate the worst case scenario for combining lake-means for developing empirical models. Within each lake, *chla* is positively correlated with *TP* (black lines). The correlation between lakes means of *chla* and *TP* is, however, negative (shaded dots and line). The best case scenario is realized when the four datasets overlap (four lakes are identical).

80 Simpson’s paradox is a well-discussed topic in social and political  
 81 sciences. An early case was the Berkeley graduate admission paradox  
 82 (Bickel et al., 1975), where the campus-wide aggregated graduate admission

83 rate showed a bias against female applicants, whereas disaggregated data  
84 showed neutral or favorable rates towards female applicants in most  
85 departments. More recently, the apparent switch of allegiance of the two  
86 major US political parties (blue states are more affluent than red states)  
87 was contradicted by data showing that wealthy people are more likely to  
88 vote for Republican candidates (Gelman, 2009). There are numerous  
89 statistical studies on the topic, with two that are particularly helpful in  
90 developing strategies to avoid the paradox. Lindley and Novick (1981)  
91 explained the paradox from a statistical inference perspective, that is,  
92 statistical inference is the application of a model developed based on data  
93 from the population to a new individual. They suggested that the cause of  
94 Simpson’s paradox is that the new individual is not “exchangeable” with  
95 individuals in the population. In Figure 1, we present two groups of  
96 models: models for individual lakes and the model of lake means. From a  
97 statistical inference perspective, both groups of models are valid. But the  
98 models are intended for two different populations: individual observations  
99 in a particular lake and lake means of *chla* and TP. The model developed  
100 using lake means may give the false impression that *chla* and TP are  
101 inversely correlated. Such inverse correlations can often be explained by  
102 factors not included in the model, as suggested by Pearl et al. (2016):  
103 Simpson’s paradox is a problem of confounding factors and thus can be  
104 easily resolved under a causal inference framework, where effects of these

105 confounders are explicitly accounted through the use of a causal diagram.  
106 This conclusion is supported by many cross-scale studies. For example, Li  
107 et al. (2019) show that parameters of a precipitation-stream flow model  
108 vary by region due to region-specific confounding factors.

109 In lake eutrophication studies, quantifying the effects of nutrients  
110 (nitrogen and phosphorous) on algal growth is almost always the primary  
111 concern, given that excessive nutrient input is a well-established cause of  
112 algal proliferation. If we can identify important confounding factors of this  
113 relationship, than adopting the causal inference approach is likely more  
114 suitable. When analyzing data from multiple lakes (as in Figure 1), each  
115 lake may have different confounding factors, statistical inference using a  
116 hierarchical modeling approach, such as the ones used in Cha and Stow  
117 (2014) may be more effective.

118 In this paper, we use two large data sets to illustrate the potential  
119 hazards of using data from multiple lakes without properly addressing the  
120 among-lake variation that is often defined as changes in regression model  
121 coefficients when the model is fit to data from different lakes. The  
122 among-lake variation can also be reflected in the changes in model  
123 coefficients when the same model is fit using two data sets collected using  
124 the same protocol, even when the number of lakes included in the data is  
125 large. We illustrate the effects of the among-lake variation on  
126 regression-based lake models by comparing models fit using lake sample

averages from several cross-sectional datasets. We then present a Bayesian hierarchical modeling (BHM) approach for the hierarchical data structure and an empirical Bayes interpretation of a BHM’s hyper-parameter distribution to facilitate the use of cross-lake data for lake-specific inference. As the BHM approach is consistent with the shrinkage estimator of Stein’s paradox (Qian et al., 2015), our paper provides a Stein’s paradox solution to a Simpson’s paradox problem.

## 2. Materials and Methods

### 2.1. Data

We used data from both the National Lakes Assessment (NLA) conducted by the US Environmental Protection Agency (EPA) (U.S. EPA, 2009, 2016) and the LAke multiscaled GeOSpatial and temporal database (LAGOSNE) (Soranno et al., 2017) to illustrate potential statistical issues that may arise when analyzing large data sets encompassing multiple lakes. The NLA consists of 1,152 lakes sampled in 2007 (NLA2007) and 1,099 lakes sampled in 2012 (NLA2012). Data were collected in each year using an identical sampling protocol. Lakes included in the NLA were selected using a probabilistic sampling design in an attempt to accurately represent the overall population of lakes in the United States. In contrast to the NLA, the LAGOSNE database contains information on lakes with monitoring data from federal, state, or citizen science monitoring programs



148 across 17 states in the northeast of the US. We used 27 lakes from  
 149 LAGOSNE that were also included in NLA2007 for detailed analysis. These  
 150 lakes have at least 10 observations in LAGOSNE (Figure 2). The selection  
 151 of these 27 lakes was for the purpose of methods comparison only. A  
 152 summary of the data is in Table 1.

Table 1: Summary of data used in the analysis

	NLA2007	NLA2012	LAGOSNE
No. of obs.	1328	1230	1340
No. of lakes	1152	1099	27
No. of obs per lake	1-2	1-2	17-192
No. of years	1	1	9-29

NLA2007: data from 2007 NLA; NLA2012: data from 2012 NLA;  
 LAGOSNE: data from 27 lakes in LAGOSNE with more than 10  
 observations that were also present in NLA2007.

153 These data sets were used to illustrate (1) the effects of among-lake  
 154 variation on regression-based lake modeling and (2) the Bayesian  
 155 hierarchical modeling approach to properly account for the among-lake  
 156 variation.

157 The two NLA data sets include a large number of lakes and were  
 158 collected to be representative of lakes in the US. Using these two data sets,  
 159 we illustrate how the among-lake variation may be reflected in regression

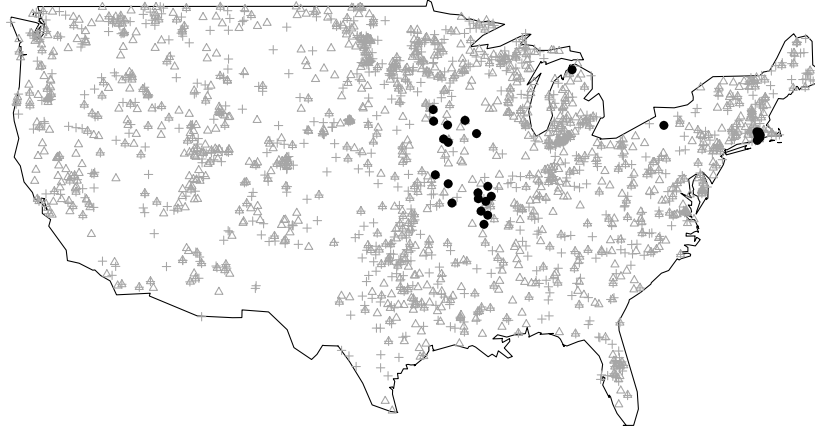


Figure 2: Locations of NLA2007 lakes (shaded pluses), NLA2012 lakes (shaded triangles), and the 27 lakes included in both NLA2007 and LAGOSNE (black dots)

models developed using the data sets separately, and fit to the combined  
data. To contrast the NLA, which includes only a small number of  
observations for each lake (such that lakes means are highly variable), we  
compare the three models fit using NLA data sets (models developed based  
on NLA2007, NLA2012, and NLA2007+NLA2012) to a model fit to a  
subset of LAGOSNE that includes 27 lakes that are represented in  
NLA2007 with at least 10 observations in each lake. For this comparison,  
we use lake mean concentrations of *chl<sub>a</sub>*, TP, and TN as the observations

168 for developing the regression model discussed in the next section.

169 Using data of the 27 lakes in LAGOSNE we show how Bayesian  
170 hierarchical modeling approach can be used to partially pool data from  
171 different lakes to avoid the potential problems of Simpson’s paradox (Figure  
172 1).

## 173 2.2. Statistical Modeling

### 174 2.2.1. Illustrating Among-Lake Variation in Model Coefficients

175 We first developed a regression model (equation (2)) to demonstrate the  
176 variability of model coefficients between data sets. The model used both  
177 TP, TN, and their interaction as predictor variables:

$$\log(chla_j) = \beta_0 + \beta_1 \log(TP_j) + \beta_2 \log(TN_j) + \beta_3 \log(TP_j) \log(TN_j) + \varepsilon_j \quad (2)$$

178 where  $chla_j$ ,  $TP_j$ , and  $TN_j$  are sample average concentrations for chl<sub>a</sub>, TP,  
179 and TN for the  $j$ th lake. Frequently, TP is used as the only predictor  
180 because phosphorus is usually assumed as the limiting nutrient; we did not  
181 make that *a priori* assumption for all the lakes in the data (Malve and  
182 Qian, 2006). Furthermore, TP and TN are often correlated, which can  
183 imply an interaction effect (Qian, 2016). For example, an oligotrophic lake  
184 may be limited by both phosphorus and nitrogen; thus increasing  
185 phosphorus may lead to an increased nitrogen demand, constituting a  
186 positive interaction. The most commonly used statistical modeling  
187 approach to account for the interaction effect is to include the product of

188 the two predictors (known as the interaction term in statistics (Qian,  
 189 2016)) in the regression model. For example, in an analysis of Finnish  
 190 lakes, Malve and Qian (2006) and Qian (2016) showed that including both  
 191 TP and TN, and their interaction term can lead to a more informative  
 192 model. Specifically, the magnitude of the coefficient  $\beta_3$  may be indicative of  
 193 a lake's trophic level (Qian, 2016). A lake is likely to be oligotrophic when  
 194  $\beta_3 > 0$  (both P and N are limiting), mesotrophic when  $\beta_3 \approx 0$  (P is likely  
 195 the limiting nutrient), and eutrophic when  $\beta_3 < 0$  (perhaps neither P nor N  
 196 is limiting). Because of the inclusion of the interaction term, the effects of  
 197 TP and TN on *chla* are no longer constants. The effect of TP depends on  
 198 the value of TN and vice versa. The meanings of software reported values  
 199 of  $\beta_1$  and  $\beta_2$  are the TP and TN effects for specific values of TN and TP,  
 200 respectively (Qian, 2016). Specifically, the reported  $\beta_1$  ( $\beta_2$ ) is the TP (TN)  
 201 effect when  $\log(TN) = 0$  ( $\log(TP) = 0$ ). In this paper, we centered both  
 202 predictors by subtracting the respective log means of TP and TN; such  
 203 that, the reported slopes (i.e.,  $\hat{\beta}_1$  and  $\hat{\beta}_2$ ) are the TP and TN effects when  
 204 the other predictor value is at the geometric mean of 27 LAGOSNE lakes.  
 205 Because the geometric means of 27 LAGOSNE lakes do not have the same  
 206 reference value for all lakes (e.g., the geometric mean of TP represents a  
 207 high phosphorus level for some lakes and a low level for other lakes), the  
 208 software reported  $\beta_1$  and  $\beta_2$  values are not comparable among lakes.  
 209 Consequently, we focus on the comparisons of  $\beta_0$  and  $\beta_3$ . See Qian (2016)

210 for more detailed explanations.

### 211 2.2.2. Using BHM to Account for Among-Lake Variation

212 Next, we developed a Bayesian hierarchical or multilevel model to  
 213 incorporate the hierarchical structure inherent in multi-lake data. We  
 214 constructed a two-tier multilevel model; at the lake level, we use a form of  
 215 equation (2):

$$\log(chla_{ij}) = \beta_{0j} + \beta_{1j} \log(TP_{ij}) + \beta_{2j} \log(TN_{ij}) + \beta_{3j} \log(TP_{ij}) \log(TN_{ij}) + \varepsilon_{ij} \quad (3)$$

216 where the subscript  $ij$  represents the  $i$ th observation from the  $j$ th lake.  
 217 Above the individual lake level, the BHM captures the variation of among  
 218 lake-specific model coefficients. As the regression model represents a basic  
 219 well-studied limnological relationship, we expect that the log-log linear  
 220 relationship to hold for all lakes, but model coefficients  $\beta_{0:3j}$  may differ by  
 221 lake. Statistically, these lakes are regarded as exchangeable with respect to  
 222 model coefficients because without additional information we would not  
 223 know how these coefficients might differ. Thus, the lake-specific model  
 224 coefficients are modeled as random variables from a common distribution:

$$\begin{pmatrix} \beta_{0j} \\ \beta_{1j} \\ \beta_{2j} \\ \beta_{3j} \end{pmatrix} \sim MVN \left[ \begin{pmatrix} \mu_{\beta_0} \\ \mu_{\beta_1} \\ \mu_{\beta_2} \\ \mu_{\beta_3} \end{pmatrix}, \Sigma \right] \quad (4)$$

225 where  $MVN$  represents a multivariate normal distribution. Equations (3)  
 226 and (4) combined form a two-tier hierarchical model. The multivariate  
 227 normal distribution on the right-hand-side of equation (4) is often known as  
 228 the hyper-parameter distribution. The rationale of using the BHM is  
 229 discussed by Qian et al. (2015) in the context of estimating mean  
 230 concentrations of water quality variables for multiple water bodies.  
 231 Compared to coefficients estimated using lake-specific data (one lake at a  
 232 time), BHM estimated model coefficients are more accurate overall. More  
 233 importantly, the hierarchical model specified in equations (3) and (4)  
 234 separates within-lake models (specified by  $\beta_{0:3j}$ ) from the among-lake model  
 235 ( $\mu_{\beta_{0:3j}}$ ). As a result, a lake-specific inference can be made more accurately  
 236 (Stow et al., 2009).

### 237 *2.3. Modeling Road Map*

238 Our analyses consist of two parts:

- 239 1. The model represented by equation (2) was fit to lake sample average  
 240 *chla*, TP, and TN concentrations from (1) NLA2007 data alone, (2)  
 241 NLA2012 alone, (3) combined NLA2007 and NLA2012 data, and (4)  
 242 LAGOSNE to illustrate the variability of the estimated model  
 243 coefficients as a function of the data set used.
- 244 2. The hierarchical model of equations (3) and (4) was fit using data  
 245 from the 27 lakes in LAGOSNE to demonstrate the use of a BHM to

246 properly account for the among-lake variation.

247 All models were fit with log TP and log TN centered at the respective  
248 means of log TP and TN concentrations of the 27 lakes in LAGOSNE. As a  
249 result, the intercept ( $\beta_0$ ) of these models represents the log mean *chla*  
250 concentrations when TP and TN are at the (log) mean levels of the 27 lakes  
251 (log TP mean of 3.112, or geometric mean of 22.5  $\mu\text{g/L}$ , and log TN mean  
252 of 6.296, or geometric mean of 542.7  $\mu\text{g/L}$ ).

253 All statistical models were implemented in R (R Core Team, 2018),  
254 using function `lm()` for linear regression models and the function `lmer` from  
255 package `lme4` (Bates and Maechler, 2010) for BHM in equations (3) and (4)  
256 (Gelman and Hill, 2007). Annotated R code can be found at GitHub  
257 (<https://github.com/songsqian/simpsons>).

### 258 **3. Results**

#### 259 *3.1. Variability in Model Coefficients*

260 The linear model fit to the 27 LAGOSNE lakes has a much smaller  $\hat{\beta}_3$ ,  
261 as compared to the same coefficient estimated for the three linear models fit  
262 to NLA2007, NLA2012, and NLA2007+NLA2012 (Figure 3, Table 2). In  
263 addition, the LAGOSNE model coefficients have much larger standard  
264 errors because the LAGOSNE model is based on 27 sets of lake sample  
265 average concentrations ( $n = 27$ ) whereas the three NLA models are based

266 on sample averages from over 1,000 lakes. The estimated model coefficients  
 267 based on NLA2007 and NLA2012 also differ, and the model coefficients  
 268 based on the combined NLA data are closer to coefficients of the model fit  
 269 to NLA2012. The interpretations of these model coefficients, especially the  
 270 slopes, are ambiguous.  $\beta_0$  is the expected log *chla* for lakes with TP and  
 271 TN concentrations near the respective geometric means of the 27  
 272 LAGOSNE lakes. However, the meanings of the three slopes of these  
 273 models are no longer clear. Mathematically,  $\beta_1$  is the expected change in  
 274  $\log(chla)$  for every unit change in  $\log(TP)$ , while TN is held unchanged. By  
 275 using a regression model, we assume that changes in  $\log(chla)$  due to  
 276 factors not included in the model will not affect the estimated slope and  
 277 can be lumped into the error term. This assumption, however, requires that  
 278 the within-lake and among-lake relationship between  $\log(chla)$  and  $\log(TP)$   
 279 be the same. As shown in the four hypothetical lakes in Figure 1, this  
 280 assumption is likely unrealistic.

281       The ambiguity of model coefficients, manifested in the differences  
 282 among the estimated coefficients of the four models, suggests that the  
 283 practice of using lake means for developing an empirical model is  
 284 potentially misleading. The difference in the estimated model coefficients  
 285 from the two data sets collected for the same purposes (NLA2007 and  
 286 NLA2012) suggests that the best case scenario discussed in the captions of  
 287 Figure 1 is highly unlikely.



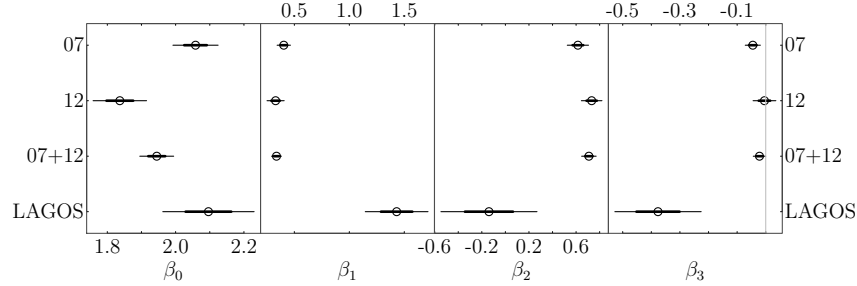


Figure 3: Model coefficients ( $\beta_{0:3}$ ) estimated using lake mean concentrations from NLA2007 (07), NLA2012 (12), NLA2007 and NLA2012 combined (07+12), and the 27 LAGOSNE lakes (LAGOS). Dots are the estimated means and thin and thick horizontal lines are the mean plus one and two standard errors, respectively. The shaded vertical line references  $\beta_3 = 0$ .

### 3.2. BHM for Among-Lake Variation

The hierarchical model fit to data from the 27 LAGOSNE lakes shows a large among-lake variation in model coefficients (Figure 4). The estimated intercepts ( $\hat{\beta}_0$ ) are the expected log *chl*a concentration for these 27 lakes when they all have the same TP and TN concentrations (the respective geometric means). As such, values of  $\beta_0$  in Figure 4 show the relative productivity of the 27 lakes (sorted based on their intercept values). The visible opposite trends between  $\beta_0$  and  $\beta_3$  are indicative of the value of  $\beta_3$  in understanding a lake's trophic level. Because the value of  $\beta_0$  is dependent on the baseline values of TP and TN, while the value of  $\beta_3$  is invariant, the interaction slope  $\beta_3$  is a more direct indicator of a lake's trophic status.

299 The wide range of  $\beta_3$  shows that these lakes have different trophic levels,  
300 indicating that nutrient effects on lake primary productivity vary by lake.

Table 2: Model Coefficients Estimated Using Different Methods

Models	07	12	07+12	LAGOS	BHM
$\beta_0$	2.058 (0.033)	1.837 (0.039)	1.9448 (0.025)	2.096 (0.067)	1.984 (0.098)
$\beta_1$	0.404 (0.030)	0.330 (0.039)	0.3376 (0.022)	1.430 (0.143)	0.850 (0.073)
$\beta_2$	0.616 (0.045)	0.732 (0.044)	0.7088 (0.031)	-0.139 (0.204)	0.390 (0.104)
$\beta_3$	-0.045 (0.013)	-0.004 (0.020)	-0.0218 (0.011)	-0.377 (0.075)	-0.014 (0.091)

Estimation standard errors are in parentheses. Models: “07” is the model fit to NLA2007 data, “12” is fit to NLA2012, “07+12” is fit to the combined NLA data, “LAGOS” is fit using the mean concentrations of the 27 lakes from LAGOSNE, BHM is the Bayesian hierarchical model (hyper-parameters,  $\mu_\beta$ ’s).

301 The difficulty in interpreting linear regression model slopes disappears  
302 when the coefficients are allowed to differ by lake. The hierarchical model  
303 estimated  $\beta_{0:3j}$  are lake-specific, while the hyper-parameters  $\mu_{\beta_{0:3}}$  are the  
304 means of the respective lake-specific coefficients. Consequently, the  
305 meaning of these estimated coefficients is unambiguous.

## 306 4. Discussion

307 Lakes in both NLA2007 and NLA2012 were selected based on a  
308 probabilistic sampling protocol such that analytical results can be

309 “(extrapolated) to national scales” (Pollard et al., 2018). It is tempting to  
310 interpret the difference in model coefficients between NLA2007 and  
311 NLA2012 (e.g., a decrease in  $\beta_0$ ) as a result of improved overall lake  
312 condition from 2007 to 2012. Because these coefficients were estimated  
313 using lake sample average concentrations of *chl<sub>a</sub>*, TP, and TN, we cannot  
314 directly interpret the differences in the models as a result of changes in lake  
315 conditions over time. A more reasonable explanation of these difference is  
316 the random sampling variability. Furthermore, the large variability in  
317 lake-specific model coefficients (Figure 4) suggests that an overall “average”  
318 model is unlikely to be informative, especially for developing management  
319 strategies that will be implemented to individual lakes.

320 Many early lake water quality models were based on simple mechanistic  
321 principles and model parameters were estimated using statistical methods  
322 (Reckhow and Chapra, 1983). These models relied on data from multiple  
323 lakes, with each lake or lake segment contributing one observation (Stow  
324 and Reckhow, 1996). As we accumulated a larger amount of data from  
325 multiple lakes, these simple modeling methods are increasingly being used  
326 as the basis for analyzing cross-sectional data. In the age of fast computers,  
327 the successful tools of the past can be easily applied to big data. Our study  
328 demonstrates the potential problems of treating “big” (multiple lakes) data  
329 using conventional methods. The hierarchical structure in the data (i.e.,  
330 from individual observations to lake-specific features to regional

331 characteristics shared by many lakes) should be properly reflected in our  
332 empirical models. The Bayesian hierarchical modeling approach provides a  
333 flexible tool for modeling the hierarchical structure inherent to most of our  
334 “big data.” When a dominant confounding factor can be identified, we can  
335 incorporate the confounding factor into the BHM (also known as the  
336 multilevel model) framework (Tang et al., 2019).

337       Without properly modeling the hierarchical structure, we risk  
338 misinterpreting the data (e.g., Figure 1), a situation that has long been  
339 recognized in statistics as the Simpson’s paradox (Simpson, 1951).  
340 Although the mathematics behind the Simpson’s paradox is  
341 straightforward, the implications of the paradox are still not widely  
342 recognized in our field. Frequently, we do not analyze data at different  
343 levels of aggregation, thereby we fail to notice the paradoxical phenomenon,  
344 which can lead to misinterpretation of the results. Lakes are naturally  
345 different (Figure 4); forcing a single model on all lakes is undesirable.

346       Developing “national” nutrient criteria using models based on lake  
347 average concentrations is likely counterproductive as nutrient  
348 concentrations are only one of many factors affecting a lake’s trophic status.  
349 A national standard would be inevitably too stringent for some lakes and  
350 too loose for others. When the among-lake variance is considered as in  
351 Yuan and Pollard (2017), the resulting criterion is most likely too stringent,  
352 and thereby unachievable, for most lakes. This result is not surprising as

353 the NLA program was designed to answer two questions (what is the  
354 current condition of lakes? and how is this condition changing over time?)  
355 that are not directly related to the quantification of the *chla*-nutrient  
356 relationship (Pollard et al., 2018).

357     The goals of the NLA monitoring program are similar to those of EPA’s  
358 Environmental Monitoring and Assessment Program (EMAP), which is  
359 optimized for estimating the mean and variance of individual  
360 environmental/ecological indicators over a national/regional scale, or of a  
361 stratified subpopulation (e.g., small lakes) (Overton and Stehman, 1996).  
362 These programs are purposefully designed to best support a limited number  
363 of objectives (Messer et al., 1991). As a result, when data from programs  
364 such as EMAP and NLA are used beyond their original design goals, we  
365 need to incorporate these data collection design parameters and plan our  
366 analysis accordingly.

367     When developing models for individual lakes, mathematical theories  
368 show that a Bayesian estimator with a proper (informative) prior is always  
369 better (compared to a non-Bayesian estimator) in terms of a model’s  
370 predictive accuracy (Efron and Morris, 1977, Efron, 1978). The difficulty in  
371 using a Bayesian method is in obtaining proper informative priors. The  
372 most important contribution of our paper is the recognition that such  
373 informative prior can be obtained by analyzing data from multiple lakes:  
374 the hyper-parameter distribution (right-hand-side of equation (4)) is

naturally such a proper prior. In other words, an important and valuable result of analyzing data from multiple lakes is the hyper-parameter distribution, which can be used as a proper informative prior for analyzing data from individual lakes that are not included in the data used to develop the hierarchical model. This conclusion is not limited to limnological modeling (Qian et al., 2015).

## 5. Conclusions

- Empirical models developed using lake average concentrations of *chla*, TP, and TN are unlikely coincide with models developed using data from individual lakes – a statistical phenomenon known as the Simpson’s paradox in statistics literature and “ecological fallacy” in social science literature.
- Regional differences in relevant natural (e.g., climate, weather, watershed soil) and cultural (e.g., land use) variables are attributed as the cause of the phenomenon. These relevant variables are known as confounding factors in causal analysis literature.
- When using cross-sectional data without detailed information about the confounding factors, a Bayesian hierarchical modeling approach is an appropriate analytic tool.

## 394 Acknowledgement

395 We thank Zutao Ouyang and colleagues at the Center for Global  
396 Change and Earth Observations at Michigan State University for feedback  
397 when the idea of the project was discussed. Constructive comments and  
398 recommendations from the associate editor and two anonymous reviewers  
399 are greatly appreciated. This work is partially supported by the Great  
400 Lakes Environmental Research Laboratory of the US National Oceanic and  
401 Atmospheric Administration (NOAA-GLERL contribution # 1918).

402 Douglas Bates and Martin Maechler. *lme4: Linear mixed-effects models*  
403 *using Eigen and Eigen++, 2010*. URL  
404 <http://CRAN.R-project.org/package=lme4>. R package version  
405 0.999375-33.

406 P.J. Bickel, E.A. Hammel, and J.W. O'Connell. Sex bias in graduate  
407 admissions: Data from Berkeley. *Science*, 187:398–404, 1975.

408 D.E. Canfield. Prediction of chlorophyll a concentrations in Florida lakes:  
409 The importance of phosphorus and nitrogen. *Journal of the American*  
410 *Water Resources Association*, 19(2):255–262, 1983.

411 D.E. Canfield and R.W. Bachmann. Prediction of total phosphorus  
412 concentrations, chlorophyll a, and secchi depths in natural and artificial

413 lakes. *Canadian Journal of Fisheries and Aquatic Sciences*, 38(4):  
414 414–423, 1981.

415 R.J. Carroll, D. Ruppert, L.A. Stefanski, and C.M. Crainiceanu.  
416 *Measurement Error in Nonlinear Models: A Modern Perspective, Second*  
417 *Edition*. Chapman & Hall/CRC Monographs on Statistics & Applied  
418 Probability. CRC Press, 2006.

419 Y. Cha and C.A. Stow. A bayesian network incorporating observation error  
420 to predict phosphorus and chlorophyll a in Saginaw Bay. *Environmental*  
421 *Modelling and Software*, 57:90–100, 2014.

422 P. J. Dillon and F. H. Rigler. Phosphorus-chlorophyll relationship in lakes.  
423 *Limnology and Oceanography*, 19:767–773, 1973.

424 B. Efron. Controversies in the foundations of statistics. *The American*  
425 *Mathematical Monthly*, 85(4):231–246, 1978.

426 B. Efron and C. Morris. Stein’s paradox in statistics. *Scientific American*,  
427 236:119–127, 1977.

428 C.T. Filstrup, T. Wagner, P.A. Soranno, E.H. Stanley, C.A. Stow, K.E.  
429 Webster, and J.A. Downing. Regional variability among nonlinear  
430 chlorophyll—phosphorus relationships in lakes. *Limnology and*  
431 *Oceanography*, 59(5):1691–1703, 2014.



- 432 W.A. Fuller. *Measurement Error Models*. Wiley Series in Probability and  
433 Statistics. Wiley, New York, 1987.
- 434 A. Gelman. *Red State, Blue State, Rich State, Poor State: Why Americans*  
435 *Vote the Way They Do - Expanded Edition*. Princeton University Press,  
436 2009. ISBN 9781400832118.
- 437 A. Gelman and J. Hill. *Data Analysis Using Regression and*  
438 *Multilevel/Hierarchical Models*. Cambridge University Press, New York,  
439 2007.
- 440 J.R. Jones and R.W. Bachmann. Prediction of phosphorus and chlorophyll  
441 levels in lakes. *Journal of Water Pollution Control Federation*, 48(9):  
442 2176–2182, 1976.
- 443 Yanzhong Li, Changming Liu, Wenjun Yu, Di Tian, and Peng Bai.  
444 Response of streamflow to environmental changes: A Budyko-type  
445 analysis based on 144 river basins over China. *Science of the Total*  
446 *Environment*, 664(10):824–833, 2019.
- 447 Z. Liang, H. Chen, S. Wu, X. Zhang, Y.H. Yu, and Y. Liu. Exploring  
448 dynamics of the chlorophyll a-total phosphorus relationship at the  
449 lake-specific scale: a bayesian hierarchical model. *Water, Air, & Soil*  
450 *Pollution*, 229(1):21, 2018.

- 451 D.V. Lindley and M.R. Novick. The role of exchangeability in inference.  
452 *The Annals of Statistics*, 9(1):45–58, 1981.
- 453 O. Malve and S.S. Qian. Estimating nutrients and chlorophyll a  
454 relationships in Finnish lakes. *Environmental Science and Technology*, 40  
455 (24):7848–7853, 2006.
- 456 E. McCauley, J.A. Downing, and S. Watson. Sigmoid relationships between  
457 nutrients and chlorophyll among lakes. *Canadian Journal of Fisheries  
458 and Aquatic Sciences*, 46:1171–1175, 1989.
- 459 J.J. Messer, R.A. Linthurst, and W.S. Overton. An EPA program for  
460 monitoring ecological status and trends. *Environmental Monitoring and  
461 Assessment*, 17(1):67–78, 1991.
- 462 W.S. Overton and S.V. Stehman. Desirable design characteristics for  
463 long-term monitoring of ecological variables. *Environmental and  
464 Ecological Statistics*, 3(4):349–361, 1996.
- 465 J Pearl, M. Glymour, and N.P. Jewell. *Causal Inference in Statistics*.  
466 Wiley, Chichester, UK, 2016.
- 467 A.I. Pollard, S.E. Hampton, and D.M. Leech. The promise and potential of  
468 continental-scale limnology using the U.S. Environmental Protection  
469 Agency’s National Lake Assessment. *Limnology and Oceanography  
470 Bulletin*, May:36–41, 2018.

471 E.E. Prepas and D.O. Trew. Evaluation of the phosphorus–chlorophyll  
 472 relationship for lakes off the precambrian shield in western canada.  
 473 *Canadian Journal of Fisheries and Aquatic Sciences*, 40(1):27–35, 1983.

474 S.S. Qian. *Environmental and Ecological Statistics with R*. Chapman and  
 475 Hall/CRC Press, 2nd edition, 2016.

476 S.S. Qian, C.A. Stow, and Y.K. Cha. Implications of Stein’s Paradox for  
 477 environmental standard compliance assessment. *Environmental Science*  
 478 *and Technology*, 49(10):5913–5920, 2015.

479 R Core Team. *R: A Language and Environment for Statistical Computing*.  
 480 R Foundation for Statistical Computing, Vienna, Austria, 2018. URL  
 481 <https://www.R-project.org/>.

482 K.H. Reckhow. A random coefficient model for chlorophyll-nutrient  
 483 relationships in lakes. *Ecological Modelling*, 70(1):35 – 50, 1993. ISSN  
 484 0304-3800.

485 K.H. Reckhow and S.C Chapra. *Engineering Approaches for Lake*  
 486 *Management: Data analysis and empirical modeling*, volume 1. Ann  
 487 Arbor Science, Butterworth Publishers, 1983.

488 D.E. Schindler. Evolution of phosphorus limitation in lakes. *Science*, 195:  
 489 260–262, 1977.

490 E.H. Simpson. The interpretation in contingency table. *Journal of Royal*  
491 *Statistical Society (B)*, 13:238–241, 1951.

492 Val H Smith and Joseph Shapiro. Chlorophyll-phosphorus relations in  
493 individual lakes. their importance to lake restoration strategies.  
494 *Environmental Science and Technology*, 15(4):444–451, 1981.

495 P.A. Soranno, L.C. Bacon, M. Beauchene, K.E. Bednar, E.G. Bissell, C.K.  
496 Boudreau, M.G. Boyer, M.T. Bremigan, S.R. Carpenter, J.W. Carr, K.S.  
497 Cheruvilil, S.T Christel, M. Claucherty, S.M. Collins, J.D. Conroy, J.A.  
498 Downing, J. Dukett, C.E. Fergus, C.T. Filstrup, C. Funk, M.J. Gonzalez,  
499 L.T. Green, C. Gries, J.D. Halfman, S.K. Hamilton, P.C. Hanson, E.N.  
500 Henry, E.M. Herron, C. Hockings, J.R. Jackson, K. Jacobson-Hedin, L.L.  
501 Janus, W.W. Jones, J.R. Jones, C.M. Keson, K.B.S. King, S.A.  
502 Kishbaugh, J.F. Lapierre, B. Lathrop, J.A. Latimore, Y. Lee, N.R.  
503 Lottig, J.A. Lynch, L.J. Matthews, W.H. McDowell, K.E.B. Moore, B.P  
504 Neff, S.J. Nelson, S.K. Oliver, M.L. Pace, D.C. Pierson, A.C. Poisson,  
505 A.I. Pollard, D.M. Post, P.O. Reyes, D.O. Rosenberry, K.M. Roy, L.G.  
506 Rudstam, O. Sarnelle, N.J. Schuldt, C.E. Scott, N.K. Skaff, N.J. Smith,  
507 N.R. Spinelli, J.J. Stachelek, E.H. Stanley, J.L. Stoddard, S.B. Stopyak,  
508 C.A. Stow, J.M. Tallant, P.N. Tan, A.P. Thorpe, M.J. Vanni, T. Wagner,  
509 G. Watkins, K.C. Weathers, K.E. Webster, J.D. White, M.K. Wilmes,  
510 and S. Yuan. LAGOS-NE: a multi-scaled geospatial and temporal

511 database of lake ecological context and water quality for thousands of us  
512 lakes. *GigaScience*, 6(12), 10 2017. ISSN 2047-217X.

513 Craig A. Stow and Kenneth H. Reckhow. Estimator bias in a lake  
514 phosphorus model with observation error. *Water Resources Research*, 32  
515 (1):165–170, 1996.

516 Craig A. Stow, E. Conrad Lamon, Song S. Qian, Patricia A. Soranno, and  
517 Kenneth H. Reckhow. *Bayesian Hierarchical/Multilevel Models for*  
518 *Inference and Prediction Using Cross-System Lake Data*, pages 111–136.  
519 Springer New York, New York, NY, 2009.

520 Tang, Q. and L. Peng and Y. Yang and S.S. Qian and B.P. Han. Total  
521 phosphorus-precipitation and Chlorophyll a-phosphorus relationships of  
522 lakes and reservoirs mediated by soil iron at regional scale. *Water*  
523 *Research*, 154: 136–143, 2019.

524 U.S. EPA. National Lakes Assessment: A Collaborative Survey of the  
525 Nation’s Lakes. Technical Report EPA 841-R-09-001, U.S. Environmental  
526 Protection Agency, Office of Water and Office of Research and  
527 Development, Washington D.C., 2009.

528 U.S. EPA. National Lakes Assessment 2012: A Collaborative Survey of  
529 Lakes in the United States. Technical Report EPA 841-R-16-113, U.S.

- 530 Environmental Protection Agency, Office of Water and Office of Research  
531 and Development, Washington D.C., December 2016.
- 532 R.A. Vollenweider. Scientific foundations of the eutrophication of lakes and  
533 flowing waters, with particular reference to nitrogen and phosphorus as  
534 factors in eutrophication. Organization for Economic Co-operation and  
535 Development, Technical Report DAS/CSI/68.27. 250p, 1968.
- 536 R.A. Vollenweider. Input-output models with special reference to  
537 phosphorus loading concept in limnology. *Schweizerische Zeitschrift für*  
538 *Hydrologie-Swiss*, 37:53–84, 1975.
- 539 T. Wagner, P.A. Soranno, K.E. Webster, and K.S. Cheruvilil. Landscape  
540 drivers of regional variation in the relationship between total phosphorus  
541 and chlorophyll in lakes. *Freshwater Biology*, 56:1811–1824, 2011.
- 542 Lester L. Yuan and Amina I. Pollard. Using national-scale data to develop  
543 nutrient–microcystin relationships that guide management decisions.  
544 *Environmental Science and Technology*, 51(12):6972–6980, 2017.

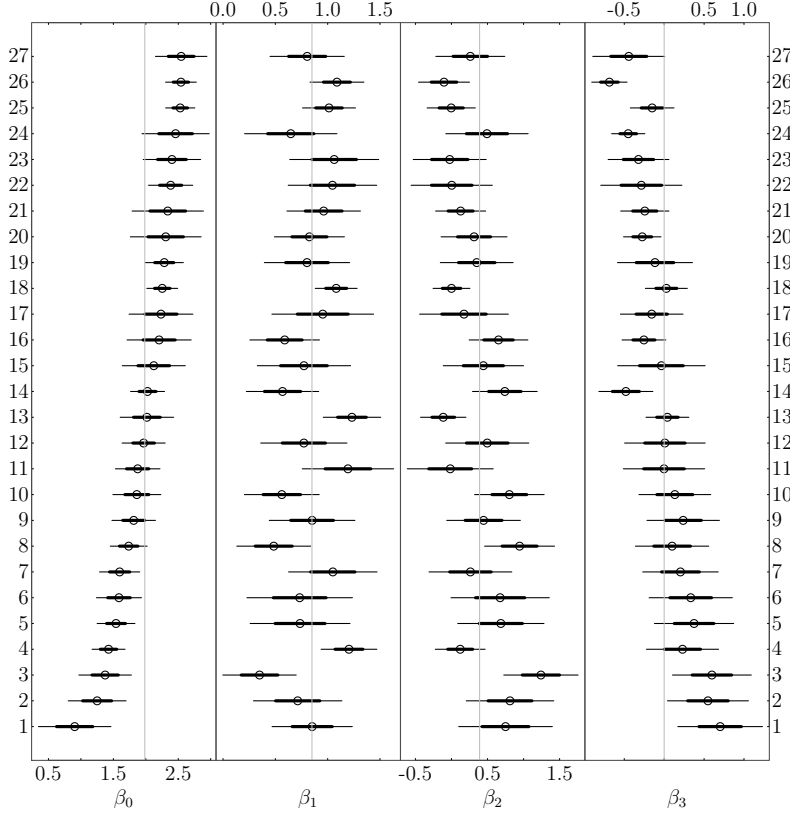


Figure 4: BHM estimated lake-specific model coefficients ( $\beta_{0j} - \beta_{3j}$ ) shown a strong negative correlation between  $\beta_{0j}$  and  $\beta_{3j}$ . Dots are the estimated means and thin and thick horizontal lines are the mean plus one and two standard errors, respectively. The shaded vertical lines for  $\beta_0, \beta_1$ , and  $\beta_2$  show the estimated respective hyper-parameters ( $\mu_{\beta_0}, \mu_{\beta_1}$ , and  $\mu_{\beta_2}$ ), the vertical line in the  $\beta_3$  panel references  $\beta_3 = 0$ .