

1

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

Song S. Qian,^{*,†} Craig A. Stow,[‡] Farnarz Nojavan A.,[¶] Joseph Stachelek,[§]
Yoonkyung Cha,^{||} Ibrahim Alameddine,[⊥] and Patricia Soranno[§]

[†]*Department of Environmental Sciences, The University of Toledo, Toledo, OH 43606*

2

[‡]*Great Lakes Environmental Research Laboratory, National Oceanic and Atmospheric
Administration, Ann Arbor, MI 48108*

[¶]*Center for Industrial Ecology, Yale University, New Haven, CT*

[§]*Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI*

^{||}*School of Environmental Engineering, University of Seoul, Seoul, South Korea*

[⊥]*Department of Civil and Environmental Engineering, American University of Beirut,
Beirut, Lebanon*

E-mail: song.qian@utoledo.edu

3

Abstract

4

Vollenweider’s revolutionary work in assessing the cause of lake eutrophication
not only implicated phosphorous as the main culprit of algal growth, but also
validated the approach of data synthesis by grouping data from multiple lakes.
Over the decades since Vollenweider’s report to OECD, limnologists routinely
use sample averages from numerous individual lakes to examine patterns across
lakes. The assumption behind the use of cross-lake data is often that responses
within and across lakes are identical. Using two large cross-lake datasets, we

10

11 demonstrate that this assumption is usually unjustified. Through comparisons
12 of an empirical model of the effect of nutrients on algal growth fit to several data
13 sets, we discuss the cognitive importance of distinguishing factors affecting lake
14 eutrophication operating at different spatial and temporal scales, and present an
15 analytic tool for properly structuring the data analysis when data from multiple
16 lakes are employed.

17 key words: NLA, LAGOS, multilevel/hierarchical model, chlorophyll a

18 **Introduction**

19 Limnologists have a long history of using data from multiple lakes, summarized at
20 various levels of spatial and temporal aggregation, to estimate empirical models. Dillon
21 and Rigler¹ set an early precedent using sample averages from a combination of 46
22 lakes, lake years, and segments of lakes to estimate a simple linear regression model
23 relating chlorophyll a (*chla*) concentration to total phosphorus (TP) concentration.
24 Numerous papers followed, applying regression approaches to estimate similar models
25 using data from other lakes, sometimes comparing their estimated equations to the
26 equation obtained by Dillon and Rigler²⁻⁵. The practice of estimating models using
27 data from multiple lakes is common, fostered by increases in computational capacity
28 and corresponding advances in statistical software which now facilitates the estimation
29 of nonlinear models, using large data sets⁶.

30 These approaches are typically based on an implicit assumption that the *chla* and
31 TP means from multiple lakes can be described by a dose-response equation such as:

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

32 where μ_{Chla} is the mean of *chla* concentration for a specified time period (such as
33 summer of a particular year) and lake (or lake segment), μ_{TP} is the mean TP con-

34 centration for a corresponding, but not necessarily the same, time period (spring TP
35 may be related to summer *chla*, for example), β_0 and β_1 are the intercept and slope
36 parameters, respectively, and ε is the model error term usually assumed to be normally
37 distributed with a constant variance. Because the underlying “true” mean values are
38 always unknown, sample averages are typically used as surrogates, although occasion-
39 ally sample medians have been used (Reckhow 1988). This regression-based modeling
40 approach has influenced lake management practices beyond the modeling of the *chla*-
41 nutrient relationship. For example, Yuan and Pollard⁷ used data from the National
42 Lake Assessment (NLA), a cross-lake data set including randomly selected lakes in all
43 48 contiguous states of the United States⁸, to develop a dose-response model to describe
44 the relationship between microcystin (MC) concentration and total nitrogen (TN) con-
45 centration. The resulting model was used to propose a national nitrogen criterion for
46 controlling harmful algal blooms.

47 The implicit premise of this approach is that a relationship estimated using sample
48 averages from many lakes can be applied to set criteria for individual lakes, because cri-
49 teria compliance assessment is typically lake-specific. However, there are two potential
50 problems with this supposition:

51 1. Using sample averages as surrogates for the “true,” unknown means, violates
52 two assumptions of regression analysis: the variance of the response variable is
53 constant and the predictor variables are observed without error. On the one
54 hand, violating the equal variance assumption makes estimated parameter and
55 model error variances ambiguous; it is unclear what uncertainty bands calculated
56 from these values, such as 95% credible or prediction intervals, represent. On the
57 other hands, violating the observation error assumption has been well-studied; it
58 is widely recognized that this “errors-in-variables” problem causes slope coefficient
59 estimators to be biased toward zero^{9,10}.

60 2. Lake-specific factors may cause individual lakes to exhibit differing stressor-response

61 relationships². Using aggregated measures, such as sample averages to estimate
 62 among-lake relationships can produce results that poorly represent the individual
 63 lakes in the analysis. In extreme cases, the sign of the estimated slope parameter
 64 can be reversed (Figure 1), a situation known as Simpson’s Paradox¹¹. Clearly,
 65 such a model should not be used to develop lake-specific management strate-
 66 gies¹²⁻¹⁴.

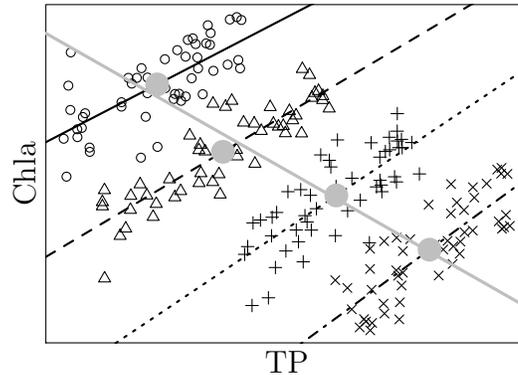


Figure 1: Hypothetical data from four lakes illustrate the worst case scenario for combining lake-means for developing empirical models. Within each lake, *chla* is positively correlated with *TP* (black lines). The correlation between lakes means of *chla* and *TP* is, however, negative (shaded dots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

67 Cha and Stow¹⁵ demonstrated a modeling approach that addresses problem 1; in
 68 this paper, we use two large data sets to illustrate the potential hazards of using data
 69 from multiple lakes without properly addressing the among-lake variation that is often
 70 defined as changes in regression model coefficients when the model is fit to data from
 71 different lakes. The among-lake variation can also be reflected in the changes in model
 72 coefficients when the same model is fit using two data sets collected using the same
 73 protocol, even when the number of lakes included in the data is large. We illustrate the
 74 effects of among-lake variation on regression-based lake models by comparing models
 75 fit using lake sample averages from several cross-sectional data sets. We then present a
 76 Bayesian hierarchical modeling (BHM) approach for the hierarchical data structure and

77 an empirical Bayes interpretation of a BHM’s hyper-parameter distribution to facilitate
78 the use of cross-lake data for lake-specific inference.

79 **Materials and Methods**

80 **Data**

81 We used data from both the National Lakes Assessment (NLA) conducted by the US
82 Environmental Protection Agency (EPA)^{16,17} and the LAke multiscaled GeOSpatial
83 and temporal database (LAGOSNE)¹⁸ to illustrate potential statistical issues that may
84 arise when analyzing large data sets encompassing multiple lakes. The NLA consists
85 of 1152 lakes sampled in 2007 (NLA2007) and 1099 lakes sampled in 2012 (NLA2012).
86 Data were collected in each year using an identical sampling protocol. Lakes included
87 in the NLA were selected using a probabilistic sampling design in an attempt to ac-
88 curately represent the overall population of lakes in the United States. In contrast to
89 the NLA, the LAGOSNE database contains information on lakes with monitoring data
90 from federal, state, or citizen science monitoring programs across 17 states in the north-
91 east of the US. We used 27 lakes from LAGOSNE that were also included in NLA2007
92 for detailed analysis. These lakes have at least 10 observations in LAGOSNE (Figure
93 2). The selection of these 27 lakes was for the purpose of methods comparison only. A
94 summary of the data is in Table 1.

Table 1: Summary of data used in the analysis

	NLA2007	NLA2012	LAGOSNE
# obs.	1328	1230	1340
# of lakes	1152	1099	27
# obs per lake	1-2	1-2	17-192
# of years	1	1	9-29

Data from LAGOSNE represent the 27 lakes with more than 10 observations that are also present in NLA2007.

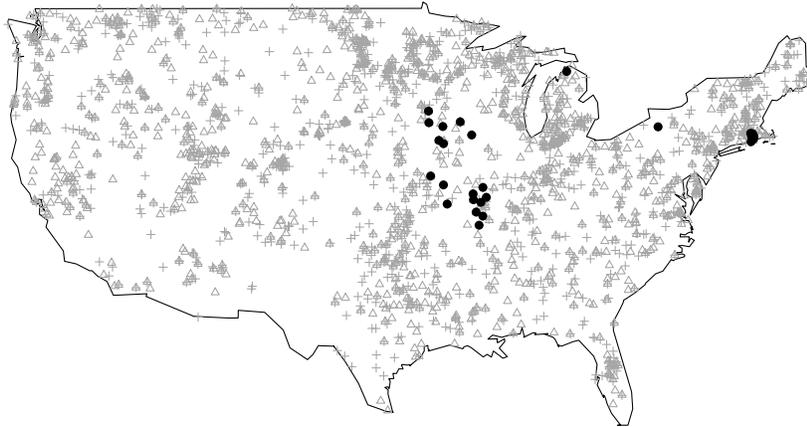


Figure 2: Locations of NLA2007 lakes (pluses), NLA2012 lakes (triangles), and the 27 lakes included in both NLA2007 and LARGOENE (black dots)

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

98 The two NLA data sets include a large number of lakes and were collected to be
 99 representative of lakes in the US. Using these two data sets, we illustrate how the among-
 100 lake variation may be reflected in regression models fit using the data sets separately,
 101 and fit to the combined data. To contrast the NLA which includes only a small number
 102 of observations for each lake (such that lakes means are highly variable), we compare
 103 the three models fit using NLA data sets to a model fit to a subset of LARGOENE
 104 that includes 27 lakes that are represented in NLA2007. For this comparison, we use
 105 lake mean concentrations of *chla*, TP, and TN as the observations for developing the
 106 regression model discussed in the next section.

107 Using data of the 27 lakes in LARGOENE we show how Bayesian hierarchical mod-

108 eling approach can be used to partially pool data from different lakes to avoid the
109 potential problems of Simpson’s paradox(Figure 1).

110 **Statistical Modeling**

111 **Illustrating Among-Lake Variation in Model Coefficients**

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

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

115 where $chla_j$, TP_j , and TN_j are sample average concentrations for chla, TP, and TN
116 for the j th lake. Frequently, TP is used as the only predictor because phosphorus is
117 usually assumed as the limiting nutrient; we did not make that *a priori* assumption for
118 all the lakes in the data¹⁹. Furthermore, TP and TN are often correlated, which can
119 imply an interaction effect²⁰. For example, an oligotrophic lake may be limited by both
120 phosphorus and nitrogen; thus increasing phosphorus may lead to an increased nitrogen
121 demand, constituting a positive interaction. In an analysis of Finnish lakes, Malve and
122 Qian¹⁹ and Qian²⁰ showed that including both TP and TN, and their interaction term
123 can lead to a more informative model. Specifically, the magnitude of the coefficient β_3
124 may be indicative of a lake’s trophic level²⁰. A lake is likely to be oligotrophic when
125 $\beta_3 > 0$ (both P and N are limiting), mesotrophic when $\beta_3 \approx 0$ (P is likely the limiting
126 nutrient), and eutrophic when $\beta_3 < 0$ (perhaps neither P nor N is limiting). Because of
127 the inclusion of the interaction term, the effects of TP and TN on $chla$ are no longer
128 constants. The effect of TP depends on the value of TN and vice versa. The meanings
129 of software reported values of β_1 and β_2 are the TP and TN effects for specific values
130 of TN and TP, respectively²⁰. Specifically, the reported β_1 (β_2) is the TP (TN) effect

131 when $\log(TN) = 0$ ($\log(TP) = 0$). In this paper, we centered both predictors by
 132 subtracting the respective log means of TP and TN ; such that, the reported slopes
 133 (i.e., $\hat{\beta}_1$ and $\hat{\beta}_2$) are the TP and TN effects when the other predictor value is at the
 134 geometric mean of 27 LAGOSNE lakes. Because the geometric means of 27 LAGOSNE
 135 lakes do not have the same reference value for all lakes (e.g., the geometric mean of
 136 TP represents a high phosphorus level for some lakes and a low level for other lakes),
 137 software reported β_1 and β_2 values are not comparable among lakes. Consequently, we
 138 focus on the comparisons of β_0 and β_3 .

139 Using BHM to Account for Among-Lake Variation

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

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

143 where the subscript ij represents the i th observation from the j th lake. Above the
 144 individual lake level, we capture the variation of among lake-specific model coefficients.
 145 As the regression model represents a basic well-studied limnological relationship, we
 146 expect that the log-log linear relationship to hold for all lakes, but that model coeffi-
 147 cients $\beta_{0:3j}$ may differ by lake. Statistically, these lakes are regarded as exchangeable
 148 because without additional information we would not know how these coefficients might
 149 differ. Thus, the lake-specific model coefficients are modeled as random variables from

150 a common distribution:

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

151 where MVN represents a multivariate normal distribution. Equations (3) and (4)
152 form a two-tier hierarchical model. The multivariate normal distribution on the right-
153 hand-side of equation (4) is often known as the hyper-parameter distribution. The
154 rationale of using the BHM is discussed by Qian et al.²¹ in the context of estimating
155 mean concentrations of water quality variables for multiple water bodies. Compared
156 to coefficients estimated using lake-specific data (one lake at a time), BHM estimated
157 model coefficients are more accurate overall. More importantly, the hierarchical model
158 specified in equations (3) and (4) separates within-lake models (specified by $\beta_{0:3j}$) from
159 the among-lake model ($\mu_{\beta_{0:3j}}$). As a result, a lake-specific inference can be made more
160 accurately²².

161 **Modeling Road Map**

162 Our analyses consist two parts:

- 163 1. The model of equation (2) was fit to lake sample average *chla*, TP, and TN
164 concentrations from (1) NLA2007 data alone, (2) NLA2012 alone, (3) combined
165 NLA2007 and NLA2012 data, and (4) LAGOSNE to illustrate the variability of
166 the estimated model coefficients.
- 167 2. The hierarchical model of equations (3) and (4) was fit using data from the 27
168 lakes in LAGOSNE to demonstrate the use of a BHM for properly account for
169 the among-lake variation.

170 All models were fit with log TP and log TN centered at the respective log means
171 of TP and TN of the 27 lakes in LAGOSNE. As a result, the intercept (β_0) of these
172 models represents the log mean *chla* concentrations when TP and TN are at the (log)
173 mean levels of the 27 lakes (log TP mean of 3.112, or geometric mean of 22.5 $\mu\text{g/L}$,
174 and log TN mean of 6.296, or geometric mean of 542.7 $\mu\text{g/L}$).

175 All statistical models were implemented in R²³, using function `lm()` for linear regres-
176 sion models and the function `lmer` from package `lme4`²⁴ for BHM in equations (3) and
177 (4). Annotated R code can be found at GitHub (<https://github.com/songsqian/simpsons>).

178 Results

179 Variability in Model Coefficients

180 The linear model fit to the 27 LAGOSNE lakes has a much smaller $\hat{\beta}_3$, as compared to
181 the same of the three linear models fit to NLA2007, NLA2012, and NLA2007+NLA2012
182 (Figure 3, Table 2). In addition, the LAGOSNE model coefficients have much larger
183 standard errors because the LAGOSNE model is based on 27 sets of lake sample av-
184 erage concentrations ($n = 27$) whereas the three NLA models are based on sample
185 averages from over 1000 lakes. The estimated model coefficients based on NLA2007
186 and NLA2012 also differ, and the model based on the combined NLA data is closer
187 to coefficients of the model fit to NLA2012. The interpretations of these model coeffi-
188 cients, especially the slopes, are ambiguous. β_0 is the expected log *chla* for lakes with
189 TP and TN concentrations near the respective geometric means of the 27 LAGOSNE
190 lakes. However, the meanings of the three slopes of these models are no longer clear.
191 Mathematically, β_1 is the expected change in $\log(\textit{chla})$ for every unit change in $\log(\textit{TP})$,
192 while TN is held unchanged. By using a regression model, we assume that changes in
193 $\log(\textit{chla})$ due to factors not included in the model will not affect the estimated slope
194 and can be lumped into the error term. This assumption, however, requires that the

195 within-lake and among-lake relationship between $\log(chla)$ and $\log(TP)$ be the same.
 196 As shown in the four hypothetical lakes in Figure 1, this assumption is likely unrealistic.
 197 The ambiguity of model coefficients manifested in the differences among the estimated
 198 coefficients of the four models, suggests that the practice of using lake means for de-
 199 veloping an empirical model is potentially misleading. The difference in the estimated
 200 model coefficients from the two data sets collected for the same purposes (NLA2007
 201 and NLA2012) suggests that the best case scenario (Figure 1) is highly unlikely.

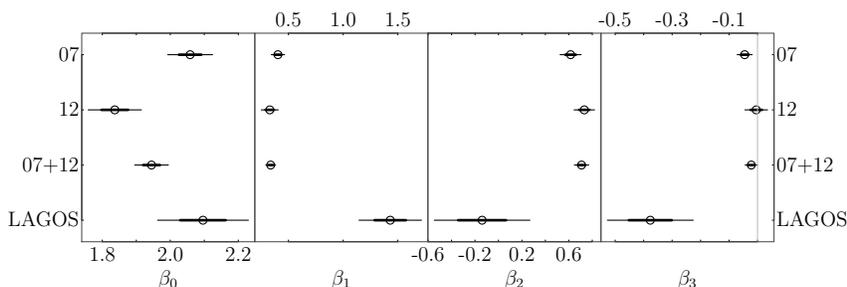


Figure 3: Model coefficients ($\beta_{0:3}$) estimated using lake mean concentrations from NLA2007 (07), NLA2012 (12), NLA2007 and NLA2012 combined (07+12), and the 27 LAGOSNE lakes (LAGOS). Dots are the estimated means and thin and thick horizontal lines are the mean plus one and two standard errors,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

202 BHM for Among-Lake Variation

203 The hierarchical model fit to data from the 27 LAGOSNE lakes shows a large among-
 204 lake variation in model coefficients (Figure 4). The estimated intercepts ($\hat{\beta}_0$) are the
 205 expected $\log chla$ concentration for these 27 lakes when they all have the same TP
 206 and TN concentrations (the respective geometric means). As such, values of β_0 in
 207 Figure 4 show the relative productivity of the 27 lakes (sorted based on their intercept
 208 values). The visible opposite trends between β_0 and β_3 are indicative of the value of
 209 β_3 in understanding a lake's trophic level. Because the value of β_0 is dependent on the
 210 baseline values of TP and TN, while the value of β_3 is invariant, the interaction slope

211 β_3 is a more direct indicator of a lake’s trophic status. The wide range of β_3 shows that
 212 these lakes have different trophic levels, indicating that nutrient effects on lake primary
 213 productivity vary by lake.

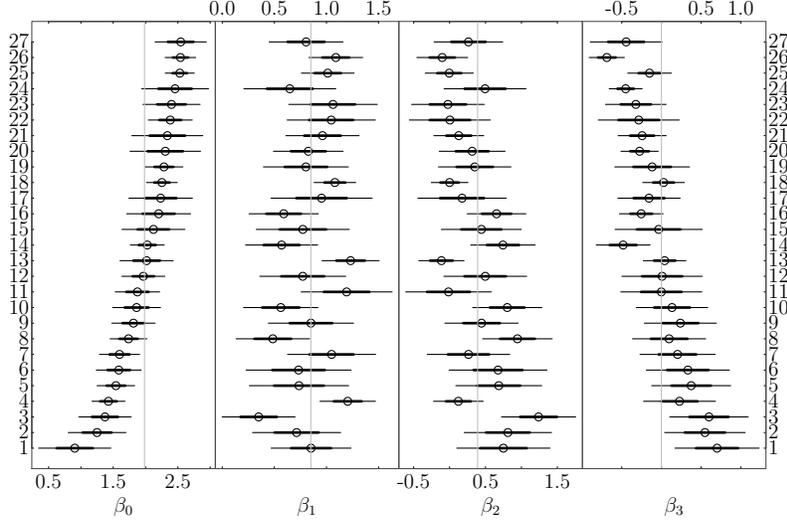


Figure 4: BHM estimated lake-specific model coefficients ($\beta_{0j} - \beta_{3j}$) shown a strong negative correlation between β_{0j} and β_{3j} . Dots are the estimated means and thin and thick horizontal lines are the mean plus one and two standard errors, respectively. The shaded vertical lines for β_0, β_1 , and β_2 show the estimated respective hyper-parameters ($\mu_{\beta_0}, \mu_{\beta_1}$, and μ_{β_2}), the vertical line in the β_3 panel references $\beta_3 = 0$.

Table 2: Model Coefficients Estimated Using Different Methods

Models	07	12	07+12	LAGOS	BHM
β_0	2.058 (0.033)	1.837 (0.039)	1.9448 (0.025)	2.096 (0.067)	1.984 (0.098)
β_1	0.404 (0.030)	0.330 (0.039)	0.3376 (0.022)	1.430 (0.143)	0.850 (0.073)
β_2	0.616 (0.045)	0.732 (0.044)	0.7088 (0.031)	-0.139 (0.204)	0.390 (0.104)
β_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, μ_{β} ’s).

214 The difficulty in interpreting linear regression model slopes disappears when the
 215 coefficients are allowed to differ by lake. The hierarchical model estimated $\beta_{0:3j}$ are lake-
 216 specific, while the hyper-parameters $\mu_{\beta_{0:3}}$ are the means of the respective lake-specific

217 coefficients. Consequently, the meaning of these estimated coefficients is unambiguous.

218 **Conclusions and Discussion**

219 We showed that empirical models based on mean concentrations of *chla*, TP, and
220 TN from NLA2007 and NLA2012 have visibly different coefficients. Lakes in both
221 NLA2007 and NLA2012 were selected based on a probabilistic sampling protocol such
222 that analytical results can be “(extrapolated) to national scales”⁸. It is tempting to
223 interpret the difference in model coefficients between NLA2007 and NLA2012 as a result
224 of improved overall lake condition from 2007 to 2012. Yet, Because these coefficients
225 were estimated using lake sample average concentrations of *chla*, TP, and TN, we cannot
226 directly interpret the differences in the models of NLA2007 and the model of NLA2012
227 as a direct result of changes in lake conditions over time. A more reasonable explanation
228 of these difference is the random sampling variability. Furthermore, the large variability
229 in lake-specific model coefficients as shown in Figure 4 suggests that an overall “average”
230 model is unlikely to be informative, especially for developing management strategies
231 that will be implemented to individual lakes.

232 Many early lake water quality models were based on simple mechanistic principles
233 and were parameterized using statistical methods²⁵. These models relied on data from
234 multiple lakes, with each lake or lake segment contributing one observation²⁶. As we ac-
235 cumulated a larger amount of data from multiple lakes, these simple modeling methods
236 are increasingly being used as the basis for analyzing cross-sectional data. In the age
237 of fast computers, the successful tools of the past can be easily applied to big data. In
238 this paper, we used a common regression model in the limnological literature to demon-
239 strate the potential problems of treating “big” (multiple lakes) data using conventional
240 methods. The hierarchical structure in the data (i.e., from individual observations to
241 lake-specific features to regional characteristics shared by many lakes) should be prop-

242 erly reflected in our empirical models. The Bayesian hierarchical modeling approach
243 provides a flexible tool for modeling the hierarchical structure inherent to most of our
244 “big data.”

245 Without properly modeling the hierarchical structure, we risk misinterpreting the
246 data (e.g., Figure 1), a situation has long been recognized in statistics as the Simp-
247 son’s paradox¹¹. Although the mathematics behind the Simpson’s paradox is straight-
248 forward, the implications of the paradox are still not widely recognized in our field.
249 Frequently, we do not analyze data at different levels of aggregation, thereby we fail
250 to notice the paradoxical phenomenon, which can lead to misinterpretation of the re-
251 sults. Lakes are naturally different (Figure 4); forcing a single model on all lakes is
252 undesirable.

253 When used to develop management strategies for eutrophication control, models
254 based on lake mean concentrations are likely to fail when used in compliance assess-
255 ment. Developing “national” nutrient criteria is likely counterproductive as nutrient
256 concentrations are only one of many factors affecting a lake’s trophic status. A na-
257 tional standard would be inevitably too stringent for some lakes and too loose for
258 others. When the among-lake variance is considered as in Yuan and Pollard⁷, the re-
259 sulting criterion is most likely too stringent, and thereby unachievable, for most lakes.
260 Consequently, a lake-specific approach is necessary.

261 When developing models for individual lakes, mathematical theories show that a
262 Bayesian estimator with a proper (informative) prior is always better (compared to a
263 non-Bayesian estimator) in terms of a model’s predictive accuracy^{27,28}. In fact, Bayes
264 himself showed that the Bayes estimator minimizes the squared error associated with
265 both observed means and the underlying true mean²⁹. In a regression problem, errors
266 associated with the observed means are the residuals. A regression model would min-
267 imize the residual sum of squares; whereas, a Bayesian regression model would also
268 minimize the error associated with the estimated model coefficients. The difficulty

269 in using a Bayesian method is in obtaining informative priors. The BHM approach
270 suggests that such informative prior can be obtained by analyzing data from multiple
271 lakes. The hyper-parameter distribution (right-hand-side of equation (4)) is naturally
272 such a proper prior. In other words, an important and valuable result of analyzing
273 data from multiple lakes is the hyper-parameter distribution, which can be used as a
274 proper informative prior for analyzing data from individual lakes that are not included
275 in the data used to develop the hierarchical model. This conclusion is not limited to
276 limnological modeling²¹.

277 Our analyses suggest that data such as NLA may be ill-suited for developing lake-
278 specific *chla*-nutrient models because of the limited lake-specific sample size. In fact,
279 with only 10% of the lakes were sampled twice⁸, fitting BHM is impossible. This out-
280 come is not surprising because the NLA program was designed to answer two questions
281 (what is the current condition of lakes? and how is this condition changing over time?)
282 that are not directly related to the quantification of the *chla* nutrient relationship⁸. The
283 goals of the NLA monitoring program are similar to those of EPA's Environmental Mon-
284 itoring and Assessment Program (EMAP), which is optimized for estimating the mean
285 and variance of individual environmental/ecological indicators over a national/regional
286 scale, or of a stratified subpopulation (e.g., small lakes)³⁰. These programs are purpose-
287 fully designed to best support a limited number of objectives³¹. As a result, when data
288 from programs such as EMAP and NLA are used beyond their original design goals,
289 we need to incorporate these data collection design parameters and plan our analysis
290 accordingly.

291 In this paper, our objectives were to (1) illustrate the potential problems of devel-
292 oping empirical models using cross-lake data and (2) demonstrate the use of BHM for
293 properly modeling the hierarchical structure of the data. Although the data we used
294 are ideal for both objectives, our BHM model from LAGOSNE may not be of any prac-
295 tical interest because the 27 lakes were selected to illustrate the potential issues and

296 for demonstrating methods. These lakes do not represent any particular subpopulation
297 of lakes. That is, the resulting models are of no particular practical purposes. For
298 the estimated hyper-parameter distribution to be practically meaningful, lakes used
299 for developing the hierarchical model should be selected to represent the subpopula-
300 tion of interest. As such, the values of large cross-lake data such as NLA lie in their
301 wide coverage that can be used to guide stratifying lakes into subpopulations, within
302 which lakes are “exchangeable,” to facilitate the proper data selection for lake-specific
303 inference. This process of careful data selection is necessitated by the recognition that
304 “correlation does not imply causation” (commonly attributed to the Irish philosopher
305 George Berkeley); statistical analysis of observational data must be done only after
306 properly balancing “confounding factors”³²⁻³⁴ and in the context of intended goals.

References

- 307
- 308 (1) Dillon, P. J.; Rigler, F. H. Phosphorus-Chlorophyll Relationship in Lakes. *Lim-*
309 *nology and Oceanography* **1973**, *19*, 767–773.
- 310 (2) Jones, J.; Bachmann, R. Prediction of Phosphorus and Chlorophyll Levels in
311 Lakes. *Journal of Water Pollution Control Federation* **1976**, *48*, 2176–2182.
- 312 (3) Canfield, D.; Bachmann, R. Prediction of Total Phosphorus Concentrations,
313 Chlorophyll a, and Secchi Depths in Natural and Artificial Lakes. *Canadian Jour-*
314 *nal of Fisheries and Aquatic Sciences* **1981**, *38*, 414–423.
- 315 (4) Canfield, D. Prediction of chlorophyll a concentrations in Florida lakes: The im-
316 portance of phosphorus and nitrogen. *Journal of the American Water Resources*
317 *Association* **1983**, *19*, 255–262.
- 318 (5) Prepas, E.; Trew, D. Evaluation of the Phosphorus–Chlorophyll Relationship for
319 Lakes Off the Precambrian Shield in Western Canada. *Canadian Journal of Fish-*
320 *eries and Aquatic Sciences* **1983**, *40*, 27–35.
- 321 (6) Filstrup, C.; Wagner, T.; Soranno, P.; Stanley, E.; Stow, C.; Webster, K.; Down-
322 ing, J. Regional variability among nonlinear chlorophyll—phosphorus relationships
323 in lakes. *Limnology and Oceanography* **2014**, *59*, 1691–1703.
- 324 (7) Yuan, L. L.; Pollard, A. I. Using National-Scale Data To Develop Nutri-
325 ent–Microcystin Relationships That Guide Management Decisions. *Environmental*
326 *Science and Technology* **2017**, *51*, 6972–6980.
- 327 (8) Pollard, A.; Hampton, S.; Leech, D. The promise and potential of continental-
328 scale limnology using the U.S. Environmental Protection Agency’s National Lake
329 Assessment. *Limnology and Oceanography Bulletin* **2018**, *May*, 36–41.

- 330 (9) Fuller, W. *Measurement Error Models*; Wiley Series in Probability and Statistics;
331 Wiley: New York, p 440.
- 332 (10) Carroll, R.; Ruppert, D.; Stefanski, L.; Crainiceanu, C. *Measurement Error in*
333 *Nonlinear Models: A Modern Perspective, Second Edition*; Chapman & Hall/CRC
334 Monographs on Statistics & Applied Probability; CRC Press, 2006; p 488.
- 335 (11) Simpson, E. The interpretation in contingency table. *Journal of Royal Statistical*
336 *Society (B)* **1951**, *13*, 238–241.
- 337 (12) Smith, V. H.; Shapiro, J. Chlorophyll-phosphorus relations in individual lakes.
338 Their importance to lake restoration strategies. *Environmental Science and Tech-*
339 *nology* **1981**, *15*, 444–451.
- 340 (13) Reckhow, K. A random coefficient model for chlorophyll-nutrient relationships in
341 lakes. *Ecological Modelling* **1993**, *70*, 35 – 50.
- 342 (14) Liang, Z.; Chen, H.; Wu, S.; Zhang, X.; Yu, Y.; Liu, Y. Exploring Dynamics of
343 the Chlorophyll a-Total Phosphorus Relationship at the Lake-Specific Scale: a
344 Bayesian Hierarchical Model. *Water, Air, & Soil Pollution* **2018**, *229*, 21.
- 345 (15) Cha, Y.; Stow, C. A Bayesian network incorporating observation error to pre-
346 dict phosphorus and chlorophyll a in Saginaw Bay. *Environmental Modelling and*
347 *Software* **2014**, *57*, 90–100.
- 348 (16) U.S. EPA, *National Lakes Assessment: A Collaborative Survey of the Nation's*
349 *Lakes*; 2009.
- 350 (17) U.S. EPA, *National Lakes Assessment 2012: A Collaborative Survey of Lakes in*
351 *the United States*; 2016.
- 352 (18) Soranno, P. A. et al. Building a multi-scaled geospatial temporal ecology database

- 353 from disparate data sources: fostering open science and data reuse. *GigaScience*
354 **2015**, *4*, 28.
- 355 (19) Malve, O.; Qian, S. Estimating Nutrients and Chlorophyll a Relationships in
356 Finnish Lakes. *Environmental Science and Technology* **2006**, *40*, 7848–7853.
- 357 (20) Qian, S. *Environmental and Ecological Statistics with R*, 2nd ed.; Chapman and
358 Hall/CRC Press, 2016.
- 359 (21) Qian, S.; Stow, C.; Cha, Y. Implications of Stein’s Paradox for Environmental
360 Standard Compliance Assessment. *Environmental Science and Technology* **2015**,
361 *49*, 5913–5920.
- 362 (22) Stow, C. A.; Lamon, E. C.; Qian, S. S.; Soranno, P. A.; Reckhow, K. H. In *Real*
363 *World Ecology: Large-Scale and Long-Term Case Studies and Methods*; Miao, S.,
364 Carstenn, S., Nungesser, M., Eds.; Springer New York: New York, NY, 2009; pp
365 111–136.
- 366 (23) R Core Team, R: A Language and Environment for Statistical Computing. R
367 Foundation for Statistical Computing: Vienna, Austria, 2018.
- 368 (24) Bates, D.; Maechler, M. lme4: Linear mixed-effects models using S4 classes. 2010;
369 R package version 0.999375-33.
- 370 (25) Reckhow, K.; Chapra, S. *Engineering Approaches for Lake Management: Data*
371 *analysis and empirical modeling*; Ann Arbor Science, Butterworth Publishers,
372 1983; Vol. 1.
- 373 (26) Stow, C. A.; Reckhow, K. H. Estimator Bias in a Lake Phosphorus Model with
374 Observation Error. *Water Resources Research* **1996**, *32*, 165–170.
- 375 (27) Efron, B.; Morris, C. Stein’s Paradox in Statistics. *Scientific American* **1977**, *236*,
376 119–127.

- 377 (28) Efron, B. Controversies in the foundations of statistics. *The American Mathemat-*
378 *ical Monthly* **1978**, *85*, 231–246.
- 379 (29) Berger, J. *Statistical Decision Theory and Bayesian Analysis*, 2nd ed.; Springer-
380 Verlag: New York, 1985.
- 381 (30) Overton, W.; Stehman, S. Desirable design characteristics for long-term monitor-
382 ing of ecological variables. *Environmental and Ecological Statistics* **1996**, *3(4)*,
383 349–361.
- 384 (31) Messer, J.; Linthurst, R.; Overton, W. An EPA program for monitoring ecological
385 status and trends. *Environmental Monitoring and Assessment* **1991**, *17(1)*, 67–78.
- 386 (32) Rubin, D. *Matched Sampling for Causal Effects*; Cambridge University Press,
387 Cambridge, UK, 2006.
- 388 (33) Qian, S.; Harmel, R. Applying statistical causal analyses to agricultural conser-
389 vation: a case study examining P loss impact. *Journal of the American Water*
390 *Resources Association* **2016**, *52*, 198–208.
- 391 (34) Nummer, S.; Qian, S.; Harmel, R. A meta-analysis on the effect of agricultural
392 conservation practices on nutrient loss. *Journal of Environmental Quality* **2018**,