ELSEVIER

#### Contents lists available at ScienceDirect

## Water Research

journal homepage: www.elsevier.com/locate/watres



# 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
- f Department of Fisheries and Wildlife, Michigan State University, East Lansing, MI, USA

#### ARTICLE INFO

#### Article history: Received 11 February 2019 Received in revised form 20 June 2019 Accepted 11 July 2019 Available online 13 July 2019

Keywords: NLA LAGOSSE Multilevel/hierarchical model Chlorophyll a

#### 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 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.

© 2019 Elsevier Ltd. All rights reserved.

## 1. Introduction

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

models using data from multiple lakes is common, fostered by increases in computational capacity and corresponding advances in statistical software which now facilitates the estimation of nonlinear models, using large data sets (Filstrup et al., 2014).

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

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

where  $\mu_{Chla}$  is the mean of chla concentration for a specified time period (such as summer of a particular year) and lake (or lake segment),  $\mu_{TP}$  is the mean TP concentration for a corresponding, but not necessarily the same, time period (spring TP may be related to summer *chla*, for example),  $\beta_0$  and  $\beta_1$  are the intercept and slope parameters, respectively, and  $\varepsilon$  is the model error term usually assumed to be normally distributed with a constant variance (Qian, 2016). Because the underlying "true" mean values are always unknown, sample averages are typically used as surrogates, although

<sup>\*</sup> Corresponding author.

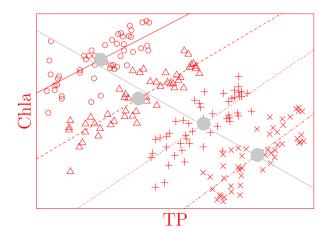
E-mail address: song.qian@utoledo.edu (S.S. Qian).

occasionally sample medians have been used (Reckhow, 1993). This regression-based modeling approach has influenced lake management practices beyond the modeling of the *chla*-nutrient relationship. For example, Yuan and Pollard (2017) used data from the National Lake Assessment (NLA), a cross-lake data set including randomly selected lakes in all 48 contiguous states of the United States (Pollard et al., 2018), to develop a dose-response model to describe the relationship between microcystin (MC) concentration and total nitrogen (TN) concentration. The resulting model was used to propose a national nitrogen criterion for controlling harmful algal blooms.

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

- 1. Using sample averages as surrogates for the "true," unknown means, violates two assumptions of regression analysis: the variance of the response variable is constant and the predictor variables are observed without error. On the one hand, violating the equal variance assumption makes an estimated parameter and model error variances ambiguous; it is unclear what uncertainty bands calculated from these values, such as 95% confidence or prediction intervals, represent. On the other hands, the consequence of violating the observation error assumption has been well-studied; it is widely recognized that this "errors-in-variables" problem causes slope coefficient estimators to be biased toward zero (Fuller, 1987; Carroll et al., 2006).
- 2. Lake-specific factors may cause individual lakes to exhibit differing stressor-response relationships (Jones and Bachmann, 1976; Wagner et al., 2011; Malve and Qian, 2006). Using aggregated measures, such as sample averages to estimate among-lake relationships can produce results that poorly represent the individual lakes in the analysis. In extreme cases, the sign of the estimated slope parameter can be reversed (Fig. 1), an example of Simpson's Paradox (Simpson, 1951). Clearly, such a model should not be used to develop lake-specific management strategies (Smith and Shapiro, 1981; Reckhow, 1993; Liang et al., 2018).

Simpson's paradox is a well-discussed topic in social and political sciences. An early case was the Berkeley graduate admission



**Fig. 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).

paradox (Bickel et al., 1975), where the campus-wide aggregated graduate admission rate showed a bias against female applicants, whereas disaggregated data showed neutral or favorable rates towards female applicants in most departments. More recently, the apparent switch of allegiance of the two major US political parties (blue states are more affluent than red states) was contradicted by data showing that wealthy people are more likely to vote for Republican candidates (Gelman, 2009). There are numerous statistical studies on the topic, with two that are particularly helpful in developing strategies to avoid the paradox. Lindley and Novick (1981) explained the paradox from a statistical inference perspective, that is, statistical inference is the application of a model developed based on data from the population to a new individual. They suggested that the cause of Simpson's paradox is that the new individual is not "exchangeable" with individuals in the population. In Fig. 1, we present two groups of models: models for individual lakes and the model of lake means. From a statistical inference perspective, both groups of models are valid. But the models are intended for two different populations: individual observations in a particular lake and lake means of chla and TP. The model developed using lake means may give the false impression that chla and TP are inversely correlated. Such inverse correlations can often be explained by factors not included in the model, as suggested by Pearl et al. (2016): Simpson's paradox is a problem of confounding factors and thus can be easily resolved under a causal inference framework, where effects of these confounders are explicitly accounted through the use of a causal diagram. This conclusion is supported by many cross-scale studies. For example, Li et al. (2019) show that parameters of a precipitation-stream flow model vary by region due to region-specific confounding factors.

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

In this paper, we use two large data sets to illustrate the potential hazards of using data from multiple lakes without properly addressing the among-lake variation that is often defined as changes in regression model coefficients when the model is fit to data from different lakes. The among-lake variation can also be reflected in the changes in model coefficients when the same model is fit using two data sets collected using the same protocol, even when the number of lakes included in the data is large. We illustrate the effects of the among-lake variation on regressionbased 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 across 17 states in the northeast of the US. We used 27 lakes from LAGOSNE that were also included in NLA2007 for detailed analysis. These lakes have at least 10 observations in LAGOSNE (Fig. 2). The selection of these 27 lakes was for the purpose of methods comparison only. A summary of the data is in Table 1.

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

The two NLA data sets include a large number of lakes and were collected to be representative of lakes in the US. Using these two data sets, we illustrate how the among-lake variation may be reflected in regression 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 lake 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 *chla*, TP, and TN as the observations for developing the regression model discussed in the next section.

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

#### 2.2. Statistical modeling

2.2.1. Illustrating among-lake variation in model coefficients

We first developed a regression model (equation (2)) to

**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 NI A2007

demonstrate the variability of model coefficients between data sets. The model used both TP, TN, and their interaction as predictor variables:

$$\log(chla_j) = \beta_0 + \beta_1 \log(TP_j) + \beta_2 \log(TN_j) + \beta_3 \log(TP_i)\log(TN_i) + \varepsilon_i$$
(2)

where  $chla_i$ ,  $TP_i$ , and  $TN_i$  are sample average concentrations for chla, TP, and TN for the *i*th lake. Frequently, TP is used as the only predictor because phosphorus is usually assumed as the limiting nutrient; we did not make that a priori assumption for all the lakes in the data (Malve and Qian, 2006). Furthermore, TP and TN are often correlated, which can imply an interaction effect (Qian, 2016). For example, an oligotrophic lake may be limited by both phosphorus and nitrogen; thus increasing phosphorus may lead to an increased nitrogen demand, constituting a positive interaction. The most commonly used statistical modeling approach to account for the interaction effect is to include the product of the two predictors (known as the interaction term in statistics (Oian, 2016)) in the regression model. For example, in an analysis of Finnish lakes, Malve and Qian (2006) and Qian (2016) showed that including both TP and TN, and their interaction term can lead to a more informative model. Specifically, the magnitude of the coefficient  $\beta_3$  may be indicative of a lake's trophic level (Oian, 2016). A lake is likely to be oligotrophic when  $\beta_3 > 0$  (both P and N are limiting), mesotrophic when  $\beta_3 \approx 0$  (P is likely the limiting nutrient), and eutrophic when  $\beta_3$  < 0 (perhaps neither P nor N is limiting). Because of the inclusion of the interaction term, the effects of TP and TN on chla are no longer constants. The effect of TP depends on the value of TN and vice versa. The meanings of software reported values of  $\beta_1$  and  $\beta_2$ are the TP and TN effects for specific values of TN and TP, respectively (Qian, 2016). Specifically, the reported  $\beta_1$  ( $\beta_2$ ) is the TP (TN)

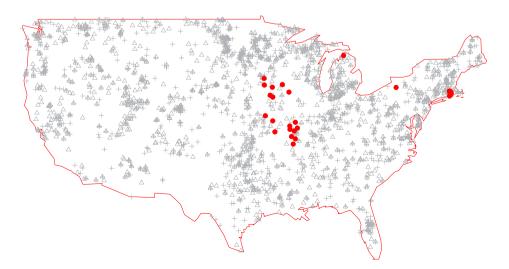


Fig. 2. Locations of NLA2007 lakes (shaded pluses), NLA2012 lakes (shaded triangles), and the 27 lakes included in both NLA2007 and LAGOSNE (black dots).

effect when  $\log(TN)=0$  ( $\log(TP)=0$ ). In this paper, we centered both predictors by subtracting the respective  $\log$  means of TP and TN; such that, the reported slopes (i.e.,  $\widehat{\beta}_1$  and  $\widehat{\beta}_2$ ) are the TP and TN effects when the other predictor value is at the geometric mean of 27 LAGOSNE lakes. Because the geometric means of 27 LAGOSNE lakes do not have the same reference value for all lakes (e.g., the geometric mean of TP represents a high phosphorus level for some lakes and a low level for other lakes), the software reported  $\beta_1$  and  $\beta_2$  values are not comparable among lakes. Consequently, we focus on the comparisons of  $\beta_0$  and  $\beta_3$ . See Qian (2016) for more detailed explanations.

## 2.2.2. Using BHM to account for among-lake variation

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

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

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

$$\begin{pmatrix} \beta_{0j} \\ \beta_{1j} \\ \beta_{2j} \\ \beta_{3j} \end{pmatrix} \sim MVN \begin{bmatrix} \begin{pmatrix} \mu_{\beta_0} \\ \mu_{\beta_1} \\ \mu_{\beta_2} \\ \mu_{\beta_3} \end{pmatrix}, \Sigma \end{bmatrix}$$

$$(4)$$

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

### 2.3. Modeling road map

Our analyses consist of two parts:

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

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

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

#### 3. Results

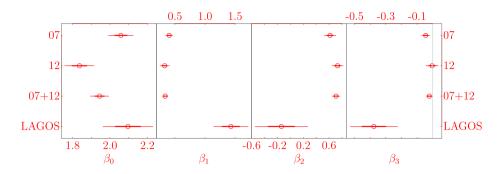
## 3.1. Variability in model coefficients

The linear model fit to the 27 LAGOSNE lakes has a much smaller  $\hat{\beta}_3$ , as compared to the same coefficient estimated for the three linear models fit to NLA2007, NLA2012, and NLA2007 + NLA2012 (Fig. 3, Table 2). In addition, the LAGOSNE model coefficients have much larger standard errors because the LAGOSNE model is based on 27 sets of lake sample average concentrations (n = 27) whereas the three NLA models are based on sample averages from over 1,000 lakes. The estimated model coefficients based on NLA2007 and NLA2012 also differ, and the model coefficients based on the combined NLA data are closer to coefficients of the model fit to NLA2012. The interpretations of these model coefficients, especially the slopes, are ambiguous.  $\beta_0$  is the expected log *chla* for lakes with TP and TN concentrations near the respective geometric means of the 27 LAGOSNE lakes. However, the meanings of the three slopes of these models are no longer clear. Mathematically,  $\beta_1$  is the expected change in log(chla) for every unit change in log(TP), while TN is held unchanged. By using a regression model, we assume that changes in log(chla) due to factors not included in the model will not affect the estimated slope and can be lumped into the error term. This assumption, however, requires that the within-lake and among-lake relationship between log(chla) and log(TP) be the same. As shown in the four hypothetical lakes in Fig. 1, this assumption is likely unrealistic.

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

## 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 (Fig. 4). The estimated intercepts  $(\hat{\beta}_0)$  are the expected log *chla* 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 Fig. 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. The wide range of  $\beta_3$  shows that these lakes have different trophic levels, indicating that nutrient effects on lake primary productivity vary by lake.



**Fig. 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  $\beta_3 = 0$ .

 Table 2

 Model coefficients estimated using different methods.

Models	07	12	07 + 12	LAGOS	ВНМ
$\beta_0$	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).

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

## 4. Discussion

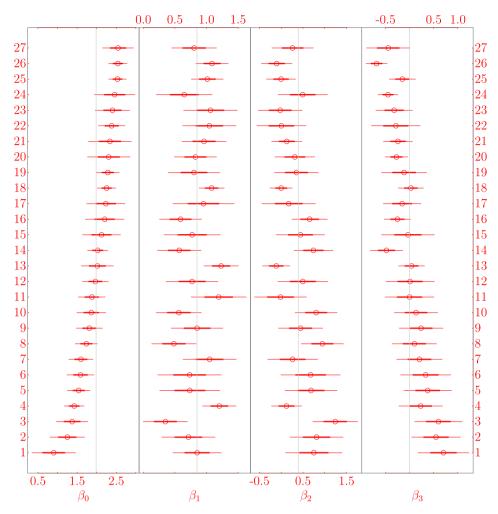
Lakes in both NLA2007 and NLA2012 were selected based on a probabilistic sampling protocol such that analytical results can be "(extrapolated) to national scales" (Pollard et al., 2018). It is tempting to interpret the difference in model coefficients between NLA2007 and NLA2012 (e.g., a decrease in  $\beta_0$ ) as a result of improved overall lake condition from 2007 to 2012. Because these coefficients were estimated using lake sample average concentrations of *chla*, TP and TN, we cannot directly interpret the differences in the models as a result of changes in lake conditions over time. A more reasonable explanation of these difference is the random sampling variability. Furthermore, the large variability in lake-specific model coefficients (Fig. 4) suggests that an overall "average" model is unlikely to be informative, especially for developing management strategies that will be implemented to individual lakes.

Many early lake water quality models were based on simple mechanistic principles and model parameters were estimated using statistical methods (Reckhow and Chapra, 1983). These models relied on data from multiple lakes, with each lake or lake segment contributing one observation (Stow and Reckhow, 1996). As we accumulated a larger amount of data from multiple lakes, these simple modeling methods are increasingly being used as the basis for analyzing cross-sectional data. In the age of fast computers, the successful tools of the past can be easily applied

to big data. Our study demonstrates the potential problems of analyzing "big" (multiple lakes) data using conventional methods. The hierarchical structure in the data (i.e., from individual observations to lake-specific features to regional characteristics shared by many lakes) should be properly reflected in our empirical models. The Bayesian hierarchical modeling approach provides a flexible tool for modeling the hierarchical structure inherent to most of our "big data." When a dominant confounding factor can be identified, we can incorporate the confounding factor into the BHM (also known as the multilevel model) framework (Tang et al., 2019).

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

Developing "national" nutrient criteria using models based on lake average concentrations is likely counterproductive as nutrient concentrations are only one of many factors affecting a lake's trophic status. A national standard would be inevitably too stringent for some lakes and too lenient for others. When the among-lake variance is considered as in Yuan and Pollard (2017), the resulting criterion is most likely too stringent, and thereby unachievable, for most lakes. This result is not surprising as the NLA program was designed to answer two questions ("what is the current condition of lakes?" and "how is this condition changing over time?") that are not directly related to the quantification of the *chla*-nutrient relationship (Pollard et al., 2018).



**Fig. 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  $(\mu_{\beta_0}, \mu_{\beta_1}, \mu_{\beta_2})$ , and  $\mu_{\beta_2}$ , the vertical line in the  $\beta_3$  panel references  $\beta_3 = 0$ .

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

When developing models for individual lakes, mathematical theories show that a Bayesian estimator with a proper (informative) prior is always better (compared to a non-Bayesian estimator) in terms of a model's predictive accuracy (Efron and Morris, 1977; Efron, 1978). The difficulty in using a Bayesian method is in obtaining proper informative priors. The most important contribution of our paper is the recognition that such informative prior can be obtained by analyzing data from multiple lakes: the hyperparameter 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.

## Acknowledgement

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

#### References

- Bates, D., Maechler, M., 2010. lme4: linear mixed-effects models using S4 classes. URL. http://CRAN.R-project.org/package=lme4. R package version 0.999375-33.
- Bickel, P.J., Hammel, E.A., O'Connell, J.W., 1975. Sex bias in graduate admissions: data from Berkeley. Science 187, 398–404.
- Canfield, D.E., 1983. Prediction of chlorophyll a concentrations in Florida lakes: the importance of phosphorus and nitrogen. J. Am. Water Resour. Assoc. 19 (2), 255–262.
- Canfield, D.E., Bachmann, R.W., 1981. Prediction of total phosphorus concentrations, chlorophyll a, and secchi depths in natural and artificial lakes. Can. J. Fish. Aquat. Sci. 38 (4), 414–423.
- Carroll, R.J., Ruppert, D., Stefanski, L.A., Crainiceanu, C.M., 2006. Measurement Error in Nonlinear Models: A Modern Perspective, second ed. Chapman & Hall/CRC Monographs on Statistics & Applied Probability. CRC Press.
- Cha, Y., Stow, C.A., 2014. A Bayesian network incorporating observation error to predict phosphorus and chlorophyll a in Saginaw Bay. Environ. Model. Softw 57, 90–100.
- Dillon, P.J., Rigler, F.H., 1973. Phosphorus-chlorophyll relationship in lakes. Limnol. Oceanogr. 19, 767–773.
- Efron, B., 1978. Controversies in the foundations of statistics. Am. Math. Mon. 85 (4), 231–246.
- Efron, B., Morris, C., 1977. Stein's paradox in statistics. Sci. Am. 236, 119–127.
- Filstrup, C.T., Wagner, T., Soranno, P.A., Stanley, E.H., Stow, C.A., Webster, K.E., Downing, J.A., 2014. Regional variability among nonlinear chlorophyll—phosphorus relationships in lakes. Limnol. Oceanogr. 59 (5), 1691–1703.
- Fuller, W.A., 1987. Measurement Error Models. Wiley Series in Probability and Statistics. Wiley, New York.
- Gelman, A., 2009. Red State, Blue State, Rich State, Poor State: Why Americans Vote the Way They Do, - Expanded Edition. Princeton University Press, ISBN 9781400832118.
- Gelman, A., Hill, J., 2007. Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press, New York.
- Jones, J.R., Bachmann, R.W., 1976. Prediction of phosphorus and chlorophyll levels in lakes. J. Water Pollut. Control Fed. 48 (9), 2176–2182.
- Li, Yanzhong, Liu, Changming, Yu, Wenjun, Tian, Di, Peng, Bai, 2019. Response of streamflow to environmental changes: a Budyko-type analysis based on 144 river basins over China. Sci. Total Environ. 664 (10), 824–833.
- Liang, Z., Chen, H., Wu, S., Zhang, X., Yu, Y.H., Liu, Y., 2018. Exploring dynamics of the chlorophyll a-total phosphorus relationship at the lake-specific scale: a bayesian hierarchical model. Water, Air, Soil Pollut. 229 (1), 21.
- Lindley, D.V., Novick, M.R., 1981. The role of exchangeability in inference. Ann. Stat. 9 (1), 45–58.
- Malve, O., Qian, S.S., 2006. Estimating nutrients and chlorophyll a relationships in Finnish lakes. Environ. Sci. Technol. 40 (24), 7848–7853.
- McCauley, E., Downing, J.A., Watson, S., 1989. Sigmoid relationships between nutrients and chlorophyll among lakes. Can. J. Fish. Aquat. Sci. 46, 1171–1175.
- Messer, J.J., Linthurst, R.A., Overton, W.S., 1991. An EPA program for monitoring ecological status and trends. Environ. Monit. Assess. 17 (1), 67–78.
- Overton, W.S., Stehman, S.V., 1996. Desirable design characteristics for long-term monitoring of ecological variables. Environ.d Ecol. Statatistics 3 (4), 349–361.
- Pearl, J., Glymour, M., Jewell, N.P., 2016. Causal Inference in Statistics. Wiley, Chichester, UK.
- Pollard, A.I., Hampton, S.E., Leech, D.M., 2018. The promise and potential of continental-scale limnology using the U.S. Environmental Protection Agency's National Lake Assessment. Limnol. Oceanogr. Bull. 36–41. May.
- Prepas, E.E., Trew, D.O., 1983. Evaluation of the phosphorus—chlorophyll relationship for lakes off the precambrian shield in western Canada. Can. J. Fish. Aquat.

- Sci. 40 (1), 27-35.
- Qian, S.S., 2016. Environmental and Ecological Statistics with R, second ed. Chapman and Hall/CRC Press.
- Qian, S.S., Stow, C.A., Cha, Y.K., 2015. Implications of Stein's Paradox for environmental standard compliance assessment. Environ. Sci. Technol. 49 (10), 5913—5920
- R Core Team, 2018. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. URL. https://www.R-project.org/.
- Reckhow, K.H., 1993. A random coefficient model for chlorophyll-nutrient relationships in lakes. Ecol. Model. 70 (1), 35–50. ISSN 0304-3800.
- Reckhow, K.H., Chapra, S.C., 1983. Engineering Approaches for Lake Management:
  Data Analysis and Empirical Modeling, vol. 1. Butterworth Publishers. Ann Arbor
  Science.
- Schindler, D.E., 1977. Evolution of phosphorus limitation in lakes. Science 195, 260–262.
- Simpson, E.H., 1951. The interpretation in contingency table. J. R. Stat. Soc. 13, 238–241.
- Smith, Val H., Shapiro, Joseph, 1981. Chlorophyll-phosphorus relations in individual lakes. their importance to lake restoration strategies. Environ. Sci. Technol. 15 (4), 444–451.
- Soranno, P.A., Bacon, L.C., Beauchene, M., Bednar, K.E., Bissell, E.G., Boudreau, C.K., Boyer, M.G., Bremigan, M.T., Carpenter, S.R., Carr, J.W., Cheruvelil, K.S., T. Christel, S., Claucherty, M., Collins, S.M., Conroy, J.D., Downing, J.A., Dukett, J., Fergus, C.E., Filstrup, C.T., Funk, C., Gonzalez, M.J., Green, L.T., Gries, C., Halfman, J.D., Hamilton, S.K., Hanson, P.C., Henry, E.N., Herron, E.M., Hockings, C., Jackson, J.R., Jacobson-Hedin, K., Janus, L.L., Jones, W.W., Jones, J.R., Keson, C.M., King, K.B.S., Kishbaugh, S.A., Lapierre, J.F., Lathrop, B., Latimore, J.A., Lee, Y., Lottig, N.R., Lynch, J.A., Matthews, L.J., McDowell, W.H., Moore, K.E.B., P. Neff, B., Nelson, S.J., Oliver, S.K., Pace, M.L., Pierson, D.C., Poisson, A.C., Pollard, A.I., Post, D.M., Reyes, P.O., Rosenberry, D.O., Roy, K.M., Rudstam, L.G., Sarnelle, O., Schuldt, N.J., Scott, C.E., Skaff, N.K., Smith, N.J., Spinelli, N.R., Stachelek, J.J., Stanley, E.H., Stoddard, J.L., Stopyak, S.B., Stow, C.A., Tallant, J.M., Tan, P.N., Thorpe, A.P., Vanni, M.J., Wagner, T., Watkins, G., Weathers, K.C., Webster, K.E., White, J.D., Wilmes, M.K., Yuan, S., LAGOS, N.E., 2017. A multiscaled geospatial and temporal database of lake ecological context and water quality for thousands of us lakes. GigaScience 6 (12), 10. ISSN 2047-217X.
- Stow, C.A., Reckhow, K.H., 1996. Estimator bias in a lake phosphorus model with observation error. Water Resour. Res. 32 (1), 165–170.
- Stow, C.A., Lamon, E.C., Qian, S.S., Soranno, P.A., Reckhow, K.H., 2009. Bayesian Hierarchical/Multilevel Models for Inference and Prediction Using Cross-System Lake Data. In: Miao, ShiLi, Carstenn, Susan, Nungesser, Martha (Eds.), Real World Ecology: Large-Scale and Long-Term Case Studies and Methods. Springer New York, New York, NY, pp. 111–136.
- Tang, Q., Peng, L., Yang, Y., Qian, S.S., Han, B.P., 2019. Total phosphorus-precipitation and Chlorophyll a-phosphorus relationships of lakes and reservoirs mediated by soil iron at regional scale. Water Res. 154, 136—143.
- U.S. EPA, 2009. National Lakes Assessment: A Collaborative Survey of the Nation's Lakes. Technical Report EPA 841-R-09-001. U.S. Environmental Protection Agency, Office of Water and Office of Research and Development, Washington D.C.
- U.S. EPA, December 2016. National Lakes Assessment 2012: A Collaborative Survey of Lakes in the United States. Technical Report EPA 841-R-16-113. U.S. Environmental Protection Agency, Office of Water and Office of Research and Development, Washington D.C.
- Vollenweider, R.A., 1968. Scientific Fundations of the Eutrophication of Lakes and Flowing Waters, with Particular Reference to Nitrogen and Phosphorus as Factors in Eutrophication. Organization for Economic Co-operation and Development. Technical Report DAS/CSI/68.27. 250pp.
- Vollenweider, R.A., 1975. Input-output models with special reference to phosphorus loading concept in limnology. Schweizerische Z. Hydrologie-Swiss 37, 53–84.
- Wagner, T., Soranno, P.A., Webster, K.E., Cheruvelil, K.S., 2011. Landscape drivers of regional variation in the relationship between total phosphorus and chlorophyll in lakes. Freshw. Biol. 56, 1811–1824.
- Yuan, L.L., Pollard, A.I., 2017. Using national-scale data to develop nutrient—microcystin relationships that guide management decisions. Environ. Sci. Technol. 51 (12), 6972–6980.