# Small Values in Big Data: The Continuing Need for Appropriate Metadata

Craig A. Stow[1], Katherine E. Webster[2], Tyler Wagner[3], Noah Lottig[4], Patricia A. Soranno[2], YoonKyung Cha[5]

[1]National Oceanic and Atmospheric Administration Great Lakes Environmental Research Laboratory, Ann Arbor, MI 48176 USA

[2]Michigan State University, Dept. Fisheries & Wildlife, East Lansing, MI 48824 USA

[3]U.S. Geological Survey, Pennsylvania Cooperative Fish and Wildlife Unit, The Pennsylvania State University, 402 Forest Resources Building, University Park, PA, 16802

[4]Univ Wisconsin, Center for Limnology, Boulder Jct, WI USA

[5]Univ Seoul, School Environmental Engineering, Seoul, South Korea

Disclaimer: This draft manuscript is distributed solely for purposes of scientific peer review. Its content is deliberative and predecisional, so it must not be disclosed or released by reviewers. Because the manuscript has not yet been approved for publication by the US Geological Survey (USGS), it does not represent any official finding or policy.

1  **Abstract**

2  Compiling data from disparate sources to address pressing ecological issues is increasingly

3  common. Many ecological datasets contain left-censored data – observations below an analytical

4  detection limit. Studies from single and typically small datasets show that common approaches

5  for handling censored data — e.g., deletion or substituting fixed values — result in systematic

6  biases. However, no studies have explored the degree to which the documentation and presence

7  of censored data influence outcomes from large, multi-sourced datasets. We describe left-

8  censored data in a lake water quality database assembled from 74 sources and illustrate the

9  challenges of dealing with small values in big data, including detection limits that are absent,

10  range widely, and show trends over time. We show that substitutions of censored data can also

11  bias analyses using 'big data' datasets, that censored data can be effectively handled with

12  modern quantitative approaches, but that such approaches rely on accurate metadata that describe

13  treatment of censored data from each source.

14

## Introduction

Data sharing is an increasing expectation in the sciences[1-3]. This outlook arises from the recognition that data are expensive and should be made widely available for maximum utility, as well as the view that information funded by taxpayers should be accessible. Although there have been concerns that users of such data are simply "datavores" or perhaps worse, "research parasites" [4], there are many scientific gains to be made from assembling data from diverse sources and harmonizing them into a consistent format for further research. The environmental sciences, in particular, stand to benefit as we investigate phenomena occurring across broad spatial and temporal scales[5-7].

Comprehensive metadata are essential to interpret large, integrated databases so that data provenance and context are retained[1,8], and to reduce the chance that patterns accidentally arise as artifacts of differing observational protocols. Complete metadata should accurately describe the "censored" observations, which result when measured samples have values that are either too high or low to be quantified (supplemental box). Samples that are below a lower detection limit are most common and are termed "left-censored". Examples include nutrient and chemical concentrations that fall below the detection limit of the analytical approach[9-10]. Though less common, "right-censoring" may also occur when, for example, concentrated aqueous samples are not adequately diluted before analysis or when Secchi depth, a measure of water clarity, exceeds the lake depth[11].

Analyzing data containing censored observations may be complicated by the fact that detection limits for the same characteristic can differ depending on the measurement protocols used, and may change over time. Ideally, metadata in a harmonized database would indicate which observations are censored and the detection limit for each censored observation. However,

38 even basic metadata can be lacking in data repositories containing data from many sources[8].

39 Thus, it is important to consider whether the censored observations are sufficiently well-

40 documented in ecological datasets to rigorously use them in analyses of compiled datasets.

41      Two common approaches for treating left-censored data include: 1) discarding the

42 censored observations or 2) substituting a value including: the detection limit, half the detection

43 limit, or zero. Under limited circumstances, these informal approaches may not strongly

44 influence the conclusions derived from the data analysis. For example, qualitative pattern

45 assessment may not be affected, particularly if the proportion of censored observations is low,

46 and their range is small relative to the overall data range. However, censored data contain

47 information, which will be improperly represented when observations are discarded or

48 substitution is used, possibly influencing inference, particularly when they comprise higher

49 proportions of the database. Additionally, even if the overall proportion of censored observations

50 is small, censoring may be disproportionately high in some groups within the data, causing

51 misleading comparisons.

52      Rigorous approaches to accommodate censored data have long been available[12-14]. Helsel

53 [15-18], Antweiler and Taylor[19], and Antweiler[20] stressed the challenges of analyzing censored data

54 and presented methods to analyze datasets containing censored observations. However, these

55 approaches still require accurate censoring metadata for all observations.

56      Our goal was to examine censored data properties in commonly-measured ecological

57 variables that have been harmonized into a large, integrated database to determine the effect of

58 censored data on ecological inference. Because such integrated databases are becoming

59 increasingly common, the potential biases due to censored data invites investigation. We used a

60 large, harmonized water quality database compiled from 76 sources[21-22]. Our objectives were to

61  quantify: a) the proportion of datasets and data values with sufficient metadata to confidently

62  identify censored observations; b) variation in reported detection limits across sources and

63  through time in the last several decades of water quality sampling; and c) the effect of three

64  strategies for dealing with censored observations on a simple water quality model and whether

65  the proportion of censored observations influences that effect.  Our results highlight the need for

66  accurate documentation and metadata.

67

68  **Methods**

69        We draw on our experience in developing LAGOS-NE (**LA**ke multi-scaled **GeOS**patial

70  & temporal database – Northeast and Midwest lakes), a lake water quality database with data

71  from 17 northeastern USA states[21].  LAGOS-NE version 1.087.1 includes contributions from 76

72  state, federal, tribal, university, citizen science, and non-profit monitoring programs with

73  chlorophyll a, total nitrogen, and total phosphorus (CHLa, TN, and TP, respectively)

74  measurements in lake surface waters. Data from two monitoring programs, consisting of 1 and 5

75  total observations, were omitted prior to our analysis.  The number of observations and programs

76  supplying data for each variable ranged, respectively, from 40,670 to 209,732 and from 33 to 66

77  (Table 1); most data were collected between 1970 and 2013.

78        During the creation of LAGOS-NE, codes that documented censor status and whether or

79  not the source program provided detection limits were assigned to each observation.  Data

80  providers indicated values were censored in multiple ways:  (a) explicit detection limits (DL)

81  were provided with each value; (b) DLs were assumed to be the reported value when tags such as

82  '<' were provided; and (c) DLs were provided in the metadata but not specified in the dataset.

83  Based on these codes, we summarized the number of programs and corresponding number of

84    observations that had DL information and the proportion of LAGOS-NE data that was comprised

85    of censored observations for each water quality variable. We used, respectively, statistical

86    summaries and cumulative frequency distributions compiled at decadal time steps to provide

87    insights into variation in DLs among programs and over time.

88        Prior to finalizing LAGOS-NE, we deleted a small number of non-censored that values

89    were reported as zero (351, 40 and 266 for CHLa, TN and TP, respectively).  We made the

90    decision to delete these, because it was unclear if these values were true zeroes, rounding

91    artifacts, or substituted values and because bivariate plots with related variables indicated, in

92    many cases, that these were outlier values.

93        To demonstrate the effect that data censoring can have on quantitative analyses we

94    simulated a large dataset with known censoring patterns. The simulated data represent a log-

95    linear relationship between TP and CHLa concentrations using parameter values previously

96    estimated from a subset of LAGOS-NE lakes[23]. We performed simulations where the proportion

97    of censoring was set to 5, 15, and 30% of the simulated data. For each of the three sets of

98    simulations, we generated 100 datasets consisting of 10,000 lakes each. The intercept, slope and

99    residual standard deviation used to generate the data were -0.24, 0.83, and 0.40, respectively. For

100   each simulated dataset, the response variable, CHLa, was left-censored at 5, 15, or 30%. We then

101   analyzed each dataset using linear regression where the censored values were estimated

102   iteratively and constrained to fall below the detection limit[24-25], and three naïve approaches

103   where: (1) censored values were omitted, (2) censored values were set to the detection limit, and

104   (3) censored values were set to half the detection limit. All models were fitted using Bayesian

105   estimation. Diffuse normal priors (N[0,1000]) were used for the intercept and slope parameters

106   and a diffuse uniform prior (Unif[0,10]) was used for the residual standard deviation using JAGS

107    in the R2jags package[26], run from within R version 3.3.0[27]. We ran three parallel Markov chains

108    beginning each chain with different values. From a total of 10,000 samples from the posterior

109    distribution the first 5,000 samples of each chain were discarded for a total of 15,000 samples

110    used to characterize the posterior distributions. We assessed convergence for all parameters both

111    visually (trace plots), as well as with the Brooks-Gelman-Rubin statistic. During each simulation

112    the estimated values of the intercept, slope, and residual standard deviation were compared to the

113    true values used in the data generating process to calculate the resultant biases.

114

115    **Results**

116        Depending on the water quality variable, 39.4 to 60.6 % of programs documented

117    censored observations either within the database or in accompanying metadata (Table 1a).

118    Despite substantial proportion of programs that did not provide DL information, their

119    contributions constituted less than 20 % of the observations in LAGOS-NE, suggesting that

120    larger lake monitoring programs typically had more information on censored data.  Further,

121    censored observations comprised a small percentage of the database, 2.4 % or less for all three

122    water quality variables (Table 1b).

123        The wide range of ways that censored data were identified in the original program

124    datasets complicated harmonization. For example, observations could be associated with specific

125    DLs, DLs could be documented program-wide, or DLs could be identified as tagged values or

126    even, in one case, inserted as negative numbers in the database. The percentage of observations

127    with specified DLs differed depending on the water quality variable.  For CHLa, TN, and TP,

128    respectively, 23, 66  and 28 % of observations had the DL specified for each observation; 19, 2,

129    and 42 % of observations had DLs assigned through metadata or as tags; and the remaining 38,

130 18, and 14 % of observations were from datasets with a mixture of censoring strategies. A few

131 of the latter programs provided databases with data collected over multiple decades and may

132 have changed specification of censored data within their database over time.

133       The extent to which individual programs substituted values when concentrations were

134 less than the DL cannot be fully evaluated. For censored observations that had associated DLs

135 specified, respectively, 7.5, 0, and 12.7 % of observations were equal to one-half the DL and

136 42.1, 1.6, and 16.0 % observations were equal to the DL for CHLa, TN, and TP. Some programs

137 reported non-censored observations with concentrations less than the reported DL, possibly

138 indicating that the reported DL was an overall method DL, not batch-specific. This disparity of

139 reporting approaches for censored observations was one of the most challenging aspects of data

140 harmonization.

141       Further complexity for data users of LAGOS-NE was the wide range of DLs (Table 1b).

142 Reported detection limits differed by over two orders of magnitude for CHLa and TP (Table 1b);

143 six DLs for TP were very high and exceeded 100 µg/L, with a maximum at 570. Despite large

144 ranges, however, median DLs were low, respectively, 1, 50 and 2 µg/L for CHLa, TN and TP.

145       Finally, we compared the overall distribution of DLs with those for data collected prior to

146 2000 and in the 2000 and 2010 decades (Figure 1). Temporal patterns in detection limits differed

147 among the three water chemistry variables. DLs for CHLa were most consistent over the three

148 time periods, with a only a small percentage having DLs exceeding 1 ug/L. In contrast, DLs for

149 TN and TP differed in cumulative frequency over time. For TN, DLs for samples collected prior

150 to 2000 included both lower and higher values compared to other time periods and overall

151 (Figure 1a). For TP, data collected prior to 2000 had lower DLs compared to later years with

152 70% of DL values less than 10 µg/L. The time period prior to 2000 did have a higher frequency

153     of DLs equal to and greater than 20 µg/L compared to later years, including half of the six DL's

154     over 100 and the two values exceeding 200. In subsequent decades, the DL for TP analyses

155     shifted towards a dominance of DL equal to 10 µg/L. These patterns suggest, at least for TP, that

156     while maximum detection limits have declined over time, the majority of earlier data was

157     analyzed under protocols with generally lower DLs. We speculate that this might be due to

158     increased automation in laboratories combined with a tradeoff of sacrificing lower sensitivity at

159     lower ends of the concentration range. The results provide cautions that systematic differences in

160     DL within the database have the potential to generate artifacts that interfere with trends and

161     patterns in the data, particularly influencing analyses based on low concentrations.

162       Our simulation study of the effects of different replacement strategies for censored data on

163     parameter estimation provide further evidence for careful consideration of how censored

164     observations are treated in large datasets. Regression lines generated from one of the 100

165     simulated data sets of 10,000 lakes help visualize the problem that occurs using various methods

166     to accommodate the censored observations (Figure 2a). In this specific result, the "true"

167     regression and censored model lines are essentially coincident, indicating that the censored

168     model closely replicates the truth. The lines generated by omitting the censored observations

169     and setting the censored observations to the detection limit are similar to one-another, both with

170     intercepts that are higher and slopes that are lower than those of the "true" model. In contrast, the

171     line that results from setting the censored observations to half the detection limit has an intercept

172     that is lower and a slope that is higher than the true model.

173       This specific result is indicative of the general pattern that becomes apparent from the 100

174     simulations (Figure 2b). Omitting censored observations or setting them to the detection limit

175     causes negatively biased slopes, positively biased intercepts, and negatively biased standard

176    deviations. However, when the censored observations are set to half the detection limit, the

177    slope, intercept, and standard deviation biases are reversed. For all three methods the size of the

178    bias increases with the proportion of censored observations. Concurrently, the censored model

179    remains unbiased, even when 30% of the observations were censored.

180

181    **Discussion**

182    We offer a cautionary tale regarding potential problems posed by censored data, for which

183    approaches to address them have been documented in the literature for many years.  However,

184    adding to the analytical issues raised in the past, the censored data in LAGOS-NE v1.087.1 are

185    likely characteristic of other large, harmonized, environmental databases and illustrate that

186    despite a history of documentation, problems persist, and new uncertainties introduced due to

187    differences in analytical procedures and data reporting among monitoring programs. While the

188    proportion of values clearly identifiable as below detection was small, there remained a

189    proportion of observations showing symptoms consistent with having been substituted, as well as

190    a small number that we labeled as "missing" because it was unclear if they were truly zero or if

191    their missingness was a detection limit artifact. This inability to clearly differentiate censored

192    observations puts users of compiled data in a difficult position; we discarded a small number of

193    observations for lack of a clearly superior alternative, given the limitations of the supporting

194    metadata.

195    Our results highlight the need for standard reporting of censored data for these common

196    water quality variables and identify complexities inherent in combining data from disparate

197    sources.  Additionally, our results support findings of Sprague et al.[8] regarding difficulties in

198    combining datasets.  In the case of LAGOS-NE, many of the limitations described in Sprague et

199    al.[8] were minimized because we solicited data directly from the program maintainers and

200    requested metadata information regarding aspects such as units, methods, chemical species and

201    detection limits and associated data tags[20]. In fact, if the dataset did not contain sufficient

202    metadata we did not consider it for inclusion in LAGOS-NE; however even with substantial

203    metadata, censored observation documentation was sometimes ambiguous.

204    Further, our simulation study showed how handling of censored data could influence

205    common analyses, such as regression modeling. The approach we have demonstrated is useful

206    for linear regression modeling; other approaches are available for different applications. For

207    example, the Bayesian hurdle model can use one set of predictor variables to predict which

208    response variable observations are below detection, and another set to estimate the value of the

209    response variable for those observations above the detection limit[28]. An important outcome of

210    our analysis shows that such biases do not diminish with sample size.  Thus, if quantified

211    estimates are needed, as they are for most statistical analyses of large datasets, then choosing

212    methods to appropriately incorporate the censored observations is necessary, and metadata

213    documentation of censoring is critical.

214    Harmonizing datasets from multiple sources offers great benefits, but also presents

215    challenges, many of which can be overcome with accurate metadata documenting the nuances of

216    the assembled data. The first major challenge that we documented is the wide range of strategies

217    for documenting DLs and censored observations among data sources. This challenge makes data

218    harmonization especially time-consuming. The second major challenge more for users of the

219    database is the changes in reported DL from the 1970's to present, the period when many

220    ecological datasets have been collected. These changes could bias trend detection in lower

221    concentrations of ecological variables such as nutrients. Although problems posed by improper

222  censored handling data are well-documented, and approaches to accommodate censored

223  observations are available when censored status is fully known, we find that the problem persists.

224  The temptation to treat left-censored values cavalierly may arise because, for many

225  environmental applications, low values indicate the absence of contamination, and thus are of

226  minimal concern. However, using substitution or discarding low values resulted in biased

227  estimation even when the proportion of censored values was small and the number of

228  observations was large. Our regression analysis example demonstrates that contemporary

229  computational approaches make rigorous treatment of censored observations straightforward, if

230  the metadata include adequate documentation. For censored data this documentation should

231  include a clear indication of which observations were censored and a specification of the

232  detection limit for each censored observation. Thorough compilation of detailed metadata in the

233  database harmonization process and attention to metadata during statistical analyses by the user

234  remain critical for successful research efforts relying on big data.

246 **References**

247 (1) Soranno, P. A.; Cheruvelil, K. S.; Elliott, K. C.; Montgomery, G. M., It's Good to Share:

248 Why Environmental Scientists' Ethics Are Out of Date. *Bioscience* **2015**, *65*, (1), 69-73.

249 (2) McNutt, M.; Lehnert, K.; Hanson, B.; Nosek, B. A.; Ellison, A. M.; King, J. L., Liberating

250 field science samples and data. *Science* **2016**, *351*, (6277), 1024-1026.

251 (3) Schimel, D., Open data. *Frontiers in Ecology and the Environment* **2017**, 15, (4), 175.

252 (4) McNutt, M., # IAmAResearchParasite. *Science* **2016,** *351*, (6277), 1005-1005.

253 (5) Heffernan, J. B.; Soranno, P. A.; Angilletta, M. J.; Buckley, L. B.; Gruner, D. S.; Keitt, T. H.;

254 Kellner, J. R.; Kominoski, J. S.; Rocha, A. V.; Xiao, J. F.; Harms, T. K.; Goring, S. J.;

255 Koenig, L. E.; McDowell, W. H.; Powell, H.; Richardson, A. D.; Stow, C. A.; Vargas, R.;

256 Weathers, K. C., Macrosystems ecology: understanding ecological patterns and processes at

257 continental scales. *Frontiers in Ecology and the Environment* **2014,** *12*, (1), 5-14.

258 (6) O'Reilly, C. M.; Sharma, S.; Gray, D. K.; Hampton, S. E.; Read, J. S.; Rowley, R. J.;

259 Schneider, P.; Lenters, J. D.; McIntyre, P. B.; Kraemer, B. M.; Weyhenmeyer, G. A.; Straile,

260 D.; Dong, B.; Adrian, R.; Allan, M. G.; Anneville, O.; Arvola, L.; Austin, J.; Bailey, J. L.;

261 Baron, J. S.; Brookes, J. D.; de Eyto, E.; Dokulil, M. T.; Hamilton, D. P.; Havens, K.;

262 Hetherington, A. L.; Higgins, S. N.; Hook, S.; Izmest'eva, L. R.; Joehnk, K. D.; Kangur, K.;

263 Kasprzak, P.; Kumagai, M.; Kuusisto, E.; Leshkevich, G.; Livingstone, D. M.; MacIntyre, S.;

264 May, L.; Melack, J. M.; Mueller-Navarra, D. C.; Naumenko, M.; Noges, P.; Noges, T.;

265 North, R. P.; Plisnier, P. D.; Rigosi, A.; Rimmer, A.; Rogora, M.; Rudstam, L. G.; Rusak, J.

266 A.; Salmaso, N.; Samal, N. R.; Schindler, D. E.; Schladow, S. G.; Schmid, M.; Schmidt, S.

267 R.; Silow, E.; Soylu, M. E.; Teubner, K.; Verburg, P.; Voutilainen, A.; Watkinson, A.;

268  Williamson, C. E.; Zhang, G. Q., Rapid and highly variable warming of lake surface waters

269  around the globe. *Geophys Res Lett* **2015,** *42*, (24), 10773-10781.

270  (7) LaDeau, S. L.; Han, B. A.; Rosi-Marshall, E. J.; Weathers, K. C., The Next Decade of Big

271  Data in Ecosystem Science. *Ecosystems* **2017,** *20*, (2), 274-283.

272  (8) Sprague, L. A.; Oelsner, G. P.; Argue, D. M., Challenges with secondary use of multi-source

273  water-quality data in the United States. *Water Res* **2017,** *110*, 252-261.

274  (9) Alexander, R. B.; Smith, R. A., Trends in the nutrient enrichment of US rivers during the late

275  20th century and their relation to changes in probable stream trophic conditions. *Limnol*

276  *Oceanogr* **2006,** *51*, (1), 639-654.

277  (10) Phillips, P. J.; Schubert, C.; Argue, D.; Fisher, I.; Furlong, E. T.; Foreman, W.; Gray, J.;

278  Chalmers, A., Concentrations of hormones, pharmaceuticals and other micropollutants in

279  groundwater affected by septic systems in New England and New York. *Sci Total Environ*

280  **2015,** *512*, 43-54.

281  (11) Carstensen, J., Censored data regression: Statistical methods for analyzing Secchi

282  transparency in shallow systems. *Limnol Oceanogr-Meth* **2010,** *8*, 376-385.

283  (12) Gilliom, R. J.; Helsel, D. R., Estimation of Distributional Parameters for Censored Trace

284  Level Water-Quality Data .1. Estimation Techniques. *Water Resour Res* **1986,** *22*, (2), 135-

285  146.

286  (13) Helsel, D. R.; Gilliom, R. J., Estimation of Distributional Parameters for Censored Trace

287  Level Water-Quality Data .2. Verification and Applications. *Water Resour Res* **1986,** *22*, (2),

288  147-155.

289  (14) Elshaarawi, A. H.; Dolan, D. M., Maximum-Likelihood Estimation of Water-Quality

290  Concentrations from Censored-Data. *Can J Fish Aquat Sci* **1989,** *46*, (6), 1033-1039.

291  (15) Helsel, D. R., More than obvious: Better methods for interpreting nondetect data. *Environ*

292      *Sci Technol* **2005,** *39*, (20), 419a-423a.

293  (16) Helsel, D. R., Fabricating data: How substituting values for nondetects can ruin results, and

294      what can be done about it. *Chemosphere* **2006,** *65*, (11), 2434-2439.

295  (17) Helsel, D., Much Ado About Next to Nothing: Incorporating Nondetects in Science. *Ann*

296      *Occup Hyg* **2010,** *54*, (3), 257-262.

297  (18) Helsel DR. 2012. Statistics for Censored Environmental Data using Minitab and R. John

298      Wiley & Sons, Inc. Hoboken, NJ.

299  (19) Antweiler, R. C.; Taylor, H. E., Evaluation of statistical treatments of left-censored

300      environmental data using coincident uncensored data sets: I. Summary statistics. *Environ Sci*

301      *Technol* **2008,** *42*, (10), 3732-3738.

302  (20) Antweiler, R. C., Evaluation of Statistical Treatments of Left-Censored Environmental Data

303      Using Coincident Uncensored Data Sets. II. Group Comparisons. *Environ Sci Technol* **2015,**

304      *49*, (22), 13439-13446.

305  (21) Soranno, P. A.; Bissell, E. G.; Cheruvelil, K. S.; Christel, S. T.; Collins, S. M.; Fergus, C.

306      E.; Filstrup, C. T.; Lapierre, J. F.; Lottig, N. R.; Oliver, S. K.; Scott, C. E.; Smith, N. J.;

307      Stopyak, S.; Yuan, S.; Bremigan, M. T.; Downing, J. A.; Gries, C.; Henry, E. N.; Skaff, N.

308      K.; Stanley, E. H.; Stow, C. A.; Tan, P. N.; Wagner, T.; Webster, K. E., Building a multi-

309      scaled geospatial temporal ecology database from disparate data sources: fostering open

310      science and data reuse. *Gigascience* **2015,** *4,* (1), 1-15.

311  (22) Soranno, P.A., L.C. Bacon, M. Beauchene, K.E. Bednar, E.G. Bissell, C.K. Boudreau, M.G.

312      Boyer, M.T. Bremigan, S.R. Carpenter, J.W. Carr and 70 co-authors. LAGOS-NE: A multi-

313    scaled geospatial temporal database of lake ecological context and water quality for

314    thousands of U.S. Lakes. *GigaScience*. **2017,** *12,* (12), 1-22.

315    (23) Wagner, T.; Soranno, P. A.; Webster, K. E.; Cheruvelil, K. S., Landscape drivers of regional

316    variation in the relationship between total phosphorus and chlorophyll in lakes. *Freshwater*

317    *Biol* **2011,** *56*, (9), 1811-1824.

318    (24) Gelman A, and Hill J. 2007. Data Analysis Using Regression and Multilvel/Hierarchical

319    Models. Cambridge University Press.

320    (25) Yun, J.; Qian, S. S., A Hierarchical Model for Estimating Long-Term Trend of Atrazine

321    Concentration in the Surface Water of the Contiguous US. *J Am Water Resour As* **2015,** *51*,

322    (4), 1128-1137.

323    (26) Su, Y.-S., and Yajima, M. 2015. R2jags: Using r to run 'jags'. R package version 0.5-7:

324    https://CRAN.R-project.org/package=R2jags.

325    (27) R Core Team. R: A Language and Environment for Statistical Computing. 2016. Vienna

326    Austria. https://www.R-project.org/

327    (28) Cha, Y.; Park, S. S.; Kim, K.; Byeon, M.; Stow, C. A., Probabilistic prediction of

328    cyanobacteria abundance in a Korean reservoir using a Bayesian Poisson model. *Water*

329    *Resour Res* **2014,** *50*, (3), 2518-2532.

330

331

332 **Table Captions**

333 **Table 1**. Overview of censored and non-censored data in the LAGOS-NE database for each

334 water quality variable. (a) The number and percentages of individual programs supplying

335 datasets with and without DL information and the corresponding number and percentage of

336 observations. (b) The number of censored observations within LAGOS-NE and summary

337 statistics of DL for censored values.

338

339 **Table 1**

| Measure | | Water quality variable | | |
|---|---|---|---|---|
| | | **CHLa** | **TN** | **TP** |
| **(a) Programs with and without DL information** | | | | |
| *Number of programs* | *n* | *58* | *33* | *66* |
| Percent with DL information | *%* | 43.1 | 39.4 | 60.6 |
| Percent with no DL information | *%* | 56.9 | 60.6 | 39.4 |
| | | | | |
| *Number of observations* | *n* | *209732* | *41670* | *158968* |
| Percent from programs with DL | *%* | 80.6 | 85.6 | 83.1 |
| Percent from programs with no DL | *%* | 19.4 | 14.4 | 16.9 |
| **(b) DL from censored observations** | | | | |
| *Number of censored observations* | *n* | *5088* | *192* | *3264* |
| | % of total | 2.43 | 0.46 | 2.05 |
| | | | | |
| *Concentration (µg/L)* | median | 1 | 84 | 10 |
| | mean | 0.99 | 145.3 | 9.0 |
| | min | 0.03 | 20 | 0.3 |
| | max | 10 | 280 | 570 |

340

341

342

343

344

345

346

**Figure Captions**

348

349 **Figure 1**

350 Cumulative frequency distribution plots of detection limits for censored observations in LAGOS-

351 NE. Distributions of all DLs and those within decadal time intervals are shown. The x-axis for

352 TP and CHLa plots, respectively, were truncated to 30 and 3 µg/L to better capture the majority

353 of observations, thus eliminating 84 and 18 observations. Summary statistics are in Table 1.

354

355 **Figure 2**

356 (a) One realization from a simulation representing the log-linear relationship between total

357 phosphorus (predictor variable) and chlorophyll *a* (response variable) in north temperate lakes.

358 Dots represent values from individual lakes (n = 10,000) and open dots represent censored

359 observations, where 30% of the observations are left-censored. Solid lines are posterior mean

360 regression lines from a censored regression model and three naïve regressions where censored

361 values were either substituted or omitted from the analysis. Note that the "Truth" fitted line is the

362 true underlying relationship and it is hardly visible because it is overlaid with the censored

363 regression model fit.

364

365 (b) The difference between the estimated and true values for the intercept, slope and residual

366 standard deviation used to simulate data for a simulation representing the log-linear relationship

367 between total phosphorus and chlorophyll *a* in north temperate lakes. There were five scenarios

368 evaluated, including a censored regression model and three naïve regressions where censored

369 values were either substituted or omitted from the analysis. Simulations were performed

20

370      assuming 5% (A), 15% (B), or 30% (C) of the observations being left-censored. The open

371      squares, triangles, and circles represent the mean difference across 100 iterations for the residual

372      standard deviation, slope, and intercept, respectively, and the horizontal bars represent the 2.5

373      and 97.5 percentiles across the 100 simulations.
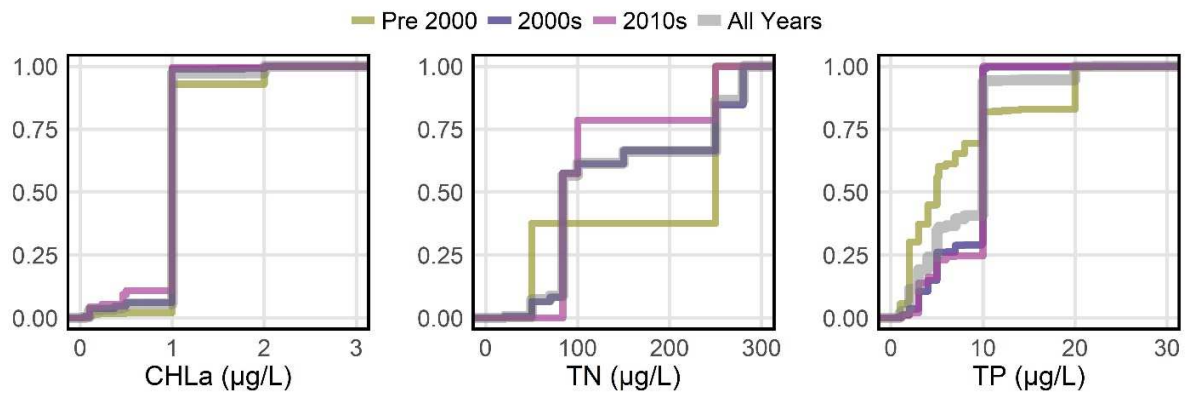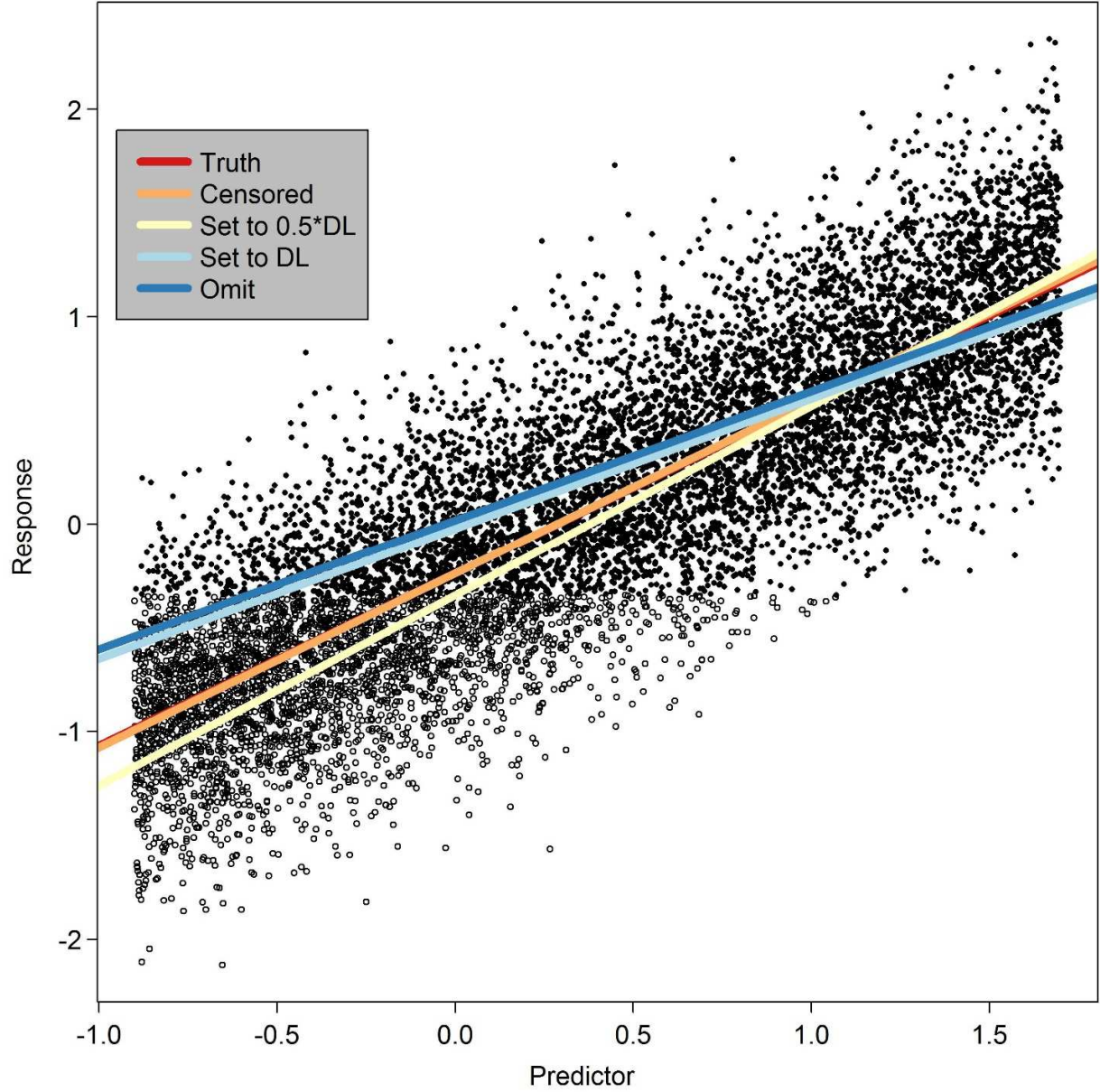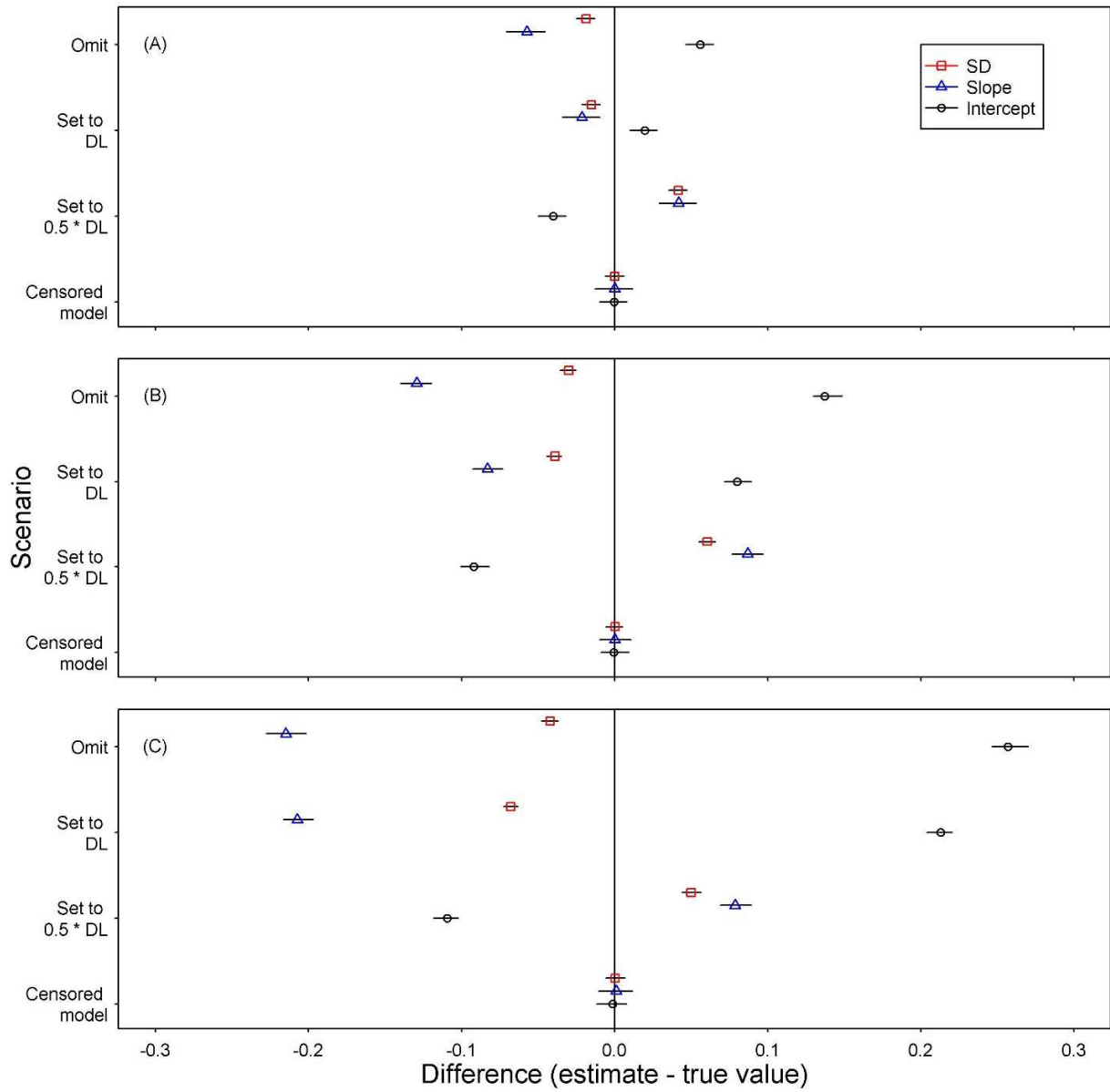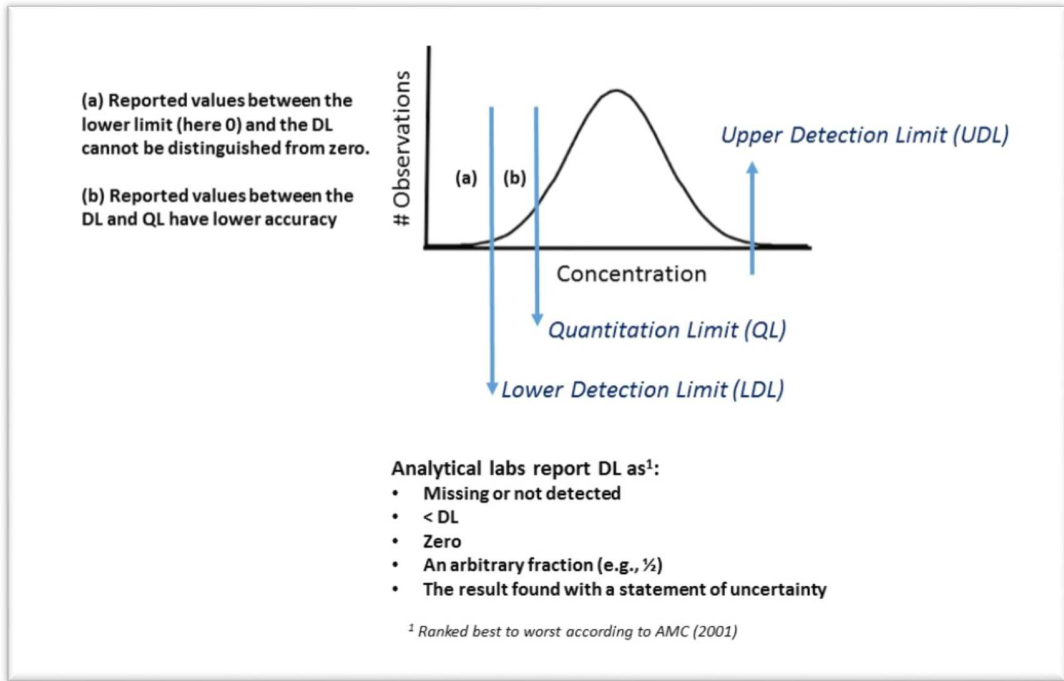
374    **Figure 1**



375

376

377

378

379

380

381

382

383　**Figure 2a**

384

385



386

387

**Figure 2b**

## Supplementary web panel:

393    Terminology used to define aspects of data quality.  Definitions from Helsel (2011).

| TERM | DEFINITION |
|------|------------|
| CENSORED DATA | Typically a low level concentration with a value between zero and the reporting limit; can also be a concentration above an upper threshold set by analytical constraints |
| REPORTING LIMIT | Concentration above which values are reported without qualification by either detection or quantitation limits |
| DETECTION LIMIT (DL) | Value below which a concentration cannot be distinguished from zero.  Related terms are LOD (limit of detection) and MDL (method detection limit) |
| QUANTITATION LIMIT (QL) | Value below which a reliable single number cannot be reported with precision.  Related term is  LOQ (limit of quantitation) |
| DATA SUBSTITUTION | Replacement of censored data in a dataset with, for example, zero, ½ the detection limit, or the detection limit. |
| TAG OR QUALIFIER | Field in a database that indicates whether a value is censored |
| MISSINGNESS | The manner in which data are missing from a sample of a population, which can cause artifacts in data analysis under certain conditions |



394
395    Analytical Methods Committee (AMC) (2001). What should be done with results below the detection
396    limit?  Mentioning the unmentionable. AMC Technical Brief (No 5): 2.