
QUANTIFYING UNCERTAINTY IN PARETO ESTIMATES OF GLOBAL LAKE AREA

A PREPRINT

 **Jemma Stachelek**

Earth System Observations
Los Alamos National Laboratory
Los Alamos, NM, USA, 87544
jsta@lanl.gov

ABSTRACT

Size is a critical factor determining the rate and occurrence of specific lake processes such as carbon sequestration and greenhouse gas emissions and emerging evidence suggests that small lakes in particular have particularly large CO_2 flux rates. Because we do not have a complete census of all lakes, upscaling estimates of such processes to small lakes at broad spatial scales requires the use of lake size-abundance distributions rather than empirical measurements of area. Existing lake census efforts are incomplete such that as lakes become smaller, they are more likely to be omitted either because they are too small to be resolved from remote sensing products or because of limited ground surveying effort (i.e. "censoring" of small lakes relative to large lakes). The present study explores one potential shortcoming of prior approaches estimating global lake area using lake size-abundance distributions. Namely, that these prior approaches rely on frequentist curve fitting techniques combined with an ad-hoc cutoff determination strategy (visual inspection to determine a likely censoring point). This yields an over-exact lake area estimate that is typically reported with no uncertainty bounds. I show how these shortcomings can be addressed with a Bayesian model that produces larger estimates of lake area uncertainty relative to the typical approach. When used as part of a sensitivity analysis, such an approach has the potential to enable more robust intercomparisons among studies of aquatic processes upscaling.

1 Introduction

Size is a critical factor determining the contribution of lakes to global biogeochemical cycling. As such, existing evidence suggests that lakes are key components of overall carbon burial, sequestration, and greenhouse gas emissions with small lakes having particularly high CO_2 flux rates (DelSontro et al., 2018; Keller et al., 2021). When upscaled to broad spatial extents, the high flux rates of these small lakes translate to a disproportionately large influence on overall CO_2 emissions (Pi et al., 2022). The calculation of such fluxes requires two terms 1) an areal flux rate and 2) the distribution of total lake area. As with all quantitative estimates at broad spatial scales, uncertainty estimates are of great interest both because they allow for intercomparisons among similar studies and because they allow for comparisons between ecosystems.

The techniques typically used to determine uncertainty in areal flux rates differ greatly from those used to determine uncertainty in total lake area distributions. In the former case, uncertainties can typically be calculated in straightforward manner from literature compilations (DelSontro et al., 2018; Keller et al., 2021). In the latter case, uncertainty in total lake area distributions is challenged by the fact that no existing database is a complete census of all lakes (Messenger et al., 2016) and the distribution of individual lake areas spans a range of approximately 7 orders of magnitude. The largest lakes are so big they could otherwise be classified as inland seas ($> 10^4 km^2$) while the smallest are barely larger than a regulation sized soccer (football) field ($> 10^{-3} km^2$).

A consequence of these challenges is that uncertainty varies within the distribution as a function of lake size. The area of the largest lakes is known with a high degree of certainty while the area of the smallest lakes is unknown below a

certain unknown threshold. The area of a small lake can be unknown either because it is too small to be resolved from remote sensing products or because of limited ground surveying effort. The omission or "censoring" of small lakes occurs because we know that small lakes exist but below a certain threshold, we have limited knowledge of their exact areas (Hamilton et al., 1992). Such censoring errors may have an outsized impact on upscaled estimates of aquatic processes such as CO_2 flux because small lakes have been shown to have particularly large CO_2 flux rates (DelSontro et al., 2018; Pi et al., 2022).

Estimating total lake area from a sample of lakes requires a conceptual model of how lakes are formed (i.e. the data generating process). Typically, lake areas are treated as arising from a fractal generating process due to the fact that landform topography, which determines the placement of lakes, can itself be treated as a fractal generating process. Indeed, many other geomorphological phenomena that are dependent on landform topography such as coastline length are often well-described by fractal generating processes (Newman, 2005). A challenge in modelling such data generating processes for lake areas is that large lakes likely follow a different data generating process than that of the smallest lakes. Whereas small lakes are constrained by landform topography, large lakes are essentially unconstrained by local landform topography and are instead constrained by the placement and arrangement of continents (Goodchild, 1988). As large lakes become even larger, they have a greater probability of intersecting a continent edge and becoming an estuary or embayment rather than a lake. Consequently, lake databases are said to be truncated on large lakes because we know that large lakes are essentially fixed in space and cannot occur in any given location (Hamilton et al., 1992).

Given that no existing database is a complete census of all lakes, yet we have near exact estimates of the area of large lakes, estimating total lake area requires a method of dealing both with the fact that lake databases are 1) truncated at large lakes and 2) censored at small lakes. Prior studies estimating global lake area have mostly not dealt with the first issue (but see Seekell and Pace (2011)). Note that in some instances the estimation process can be modified to limit its effect on the results (see Methods section). In contrast to the first issue, prior studies have typically addressed the second issue by specifying an ad-hoc cutoff value below which empirically determined lake areas are discarded and subsequently back-calculated. Hereafter, I refer to this strategy as the "cutoff" method.

In the present study, I explore two different methods for calculating uncertainty bounds around global lake area. First, I calculate these bounds using a frequentist approach for estimating total lake area (i.e. the cutoff method) via a simulation study. Second, I compare these frequentist uncertainty bounds with those calculated using a Bayesian framework. I carry out these demonstrations using a simulated dataset so that the ability of each method to recover the "true" parameter values can be evaluated. The use of a simulated dataset has the further benefit of avoiding potential confounding factors such as heterogeneity of survey effort or unknown data precision. The following analysis assumes that total global lake area has minimal temporal variability or at least minimal spatial trend. As a result, reported total lake areas can be thought of as "anticipated" or equilibrium long-run lake area rather than true totals representing any specific point in time.

2 Methods

2.1 Data description

Lake areas are typically treated as arising from a scale-invariant fractal generating process (Winslow et al., 2015; Downing et al., 2006; McDonald et al., 2012; Goodchild, 1988; Hamilton et al., 1992). This means that the number of lakes in one size class is proportional to the number of lakes in the preceding size class irrespective of their magnitudes. The numerical form describing such a process is a power-law function. One of the statistical tools often used to model data that follow a power-law function is the Pareto distribution which has a probability density function (pdf) of:

$$pdf(A) = \alpha x_{min}^{\alpha} A^{-(\alpha+1)} \quad (1)$$

where A is lake area, α controls the "shape" of the distribution and x_{min} controls the "scale" of the distribution (Shalizi, 2017). Lake area studies using the Pareto distribution do not typically use the pdf directly. Instead, they use the inverse (complementary) cumulative distribution function (ccdf) (i.e. quantile function):

$$ccdf(A) = \frac{x_{min}}{(1 - A)^{\frac{1}{\alpha}}} \quad (2)$$

The reason for using the ccdf is two-fold. First, it stabilizes model estimates in the lower tail of the distribution (Newman, 2005). This can be seen from the simulated data in Figure 1 where the Pareto pdf contains a lot of noise in the tail but the ccdf appears smoothed. The smoothing of the tail is also desirable because it functions as a way of dealing with the truncated nature of lake databases (i.e. the area of large lakes is known exactly). The second reason for

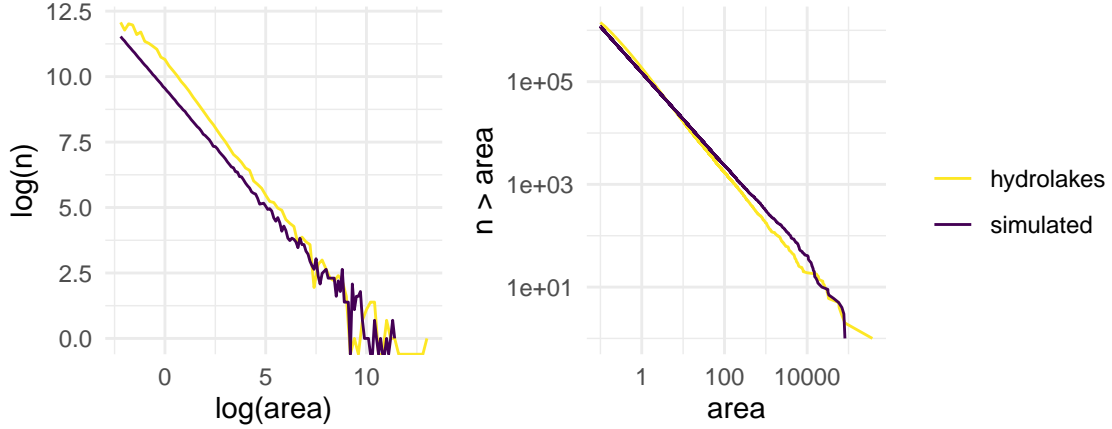


Figure 1: Realization of a Pareto (A) probability density function and (B) complementary cumulative distribution function.

using the ccdf is that it provides a computational shortcut for estimating the Pareto shape parameter a because it is numerically equivalent to the slope of the ccdf in log-log space (Downing et al., 2006).

For evaluation purposes, I generated a simulated dataset of 10,000 lake areas following the Pareto distribution using inverse transform sampling (Newman, 2005). Lakes in my simulated dataset have a minimum and maximum area of approximately 0.1 and 81000 km^2 respectively. This maximum was chosen to be approximately as large as Lake Superior but less than the Caspian Sea following Lehner and Döll (2004). The "true" total area of these lakes is approximately 230000 km^2 . I simulated a censored lake dataset by excluding lakes smaller than e^1 . This excludes (i.e. censors) approximately 60% of the total dataset. I approximated the "true" lake area total by constructing the empirical distribution function (edf) of the data which approximates the underlying Pareto cdf (Newman, 2005). Then I used this estimate of the cdf slope to generate cdf estimates for the censored lakes. I combined these cdf estimates with the edf values from the known lakes before calculating the sum of the inverted distribution (Figure ??).

I estimated the Pareto shape parameter a in a frequentist framework by calculating the the slope of the edf in log-log space using linear regression in R (Team et al., 2018). I evaluated uncertainty in both a and total lake area in a Bayesian framework using Stan (Stan Development Team, 2017). Instead of computing on the edf (as in the frequentist case), I computed directly on the pdf with the following Stan model:

```
data {
  int<lower=0> N;
  real x[N];
}
parameters {
  real<lower=0> alpha;
  real<lower=0> xmin;
}
model {
  real lpa[N];

  xmin ~ gamma(.001, .001);
  alpha ~ gamma(.001, .001);

  for (i in 1:N) {
    lpa[i] = pareto_lpdf(x[i] | xmin, alpha);
  }

  target += sum(lpa);
}
```

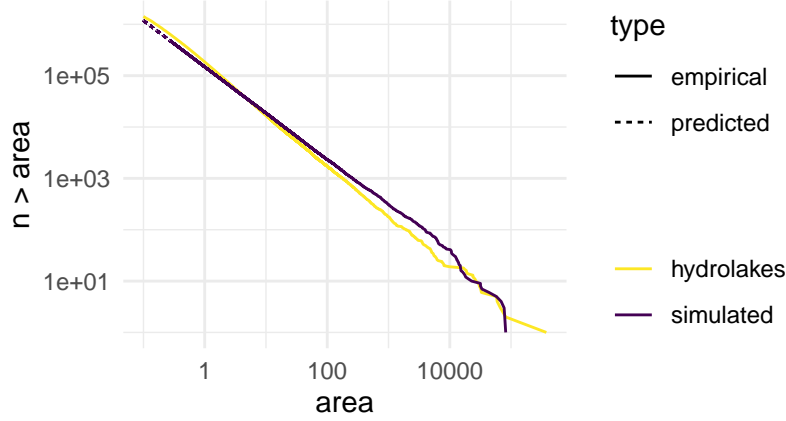


Figure 2: Censored lake area edf (solid line) and cdf estimate (dashed line).

Table 1: Frequentist uncertainty

Q5 (km^2)	Q50 (km^2)	Q95 (km^2)
2.271e+05	2.263e+05	123

I used uninformative gamma priors for both the x_{min} and a parameters following Scollnik (2007). I ran the model with four chains of 8,000 iterations and used the Stan defaults for burn-in and thinning which specify a burn-in of half the iterations and no thinning.

3 Results

Visual inspection of the frequentist method of computing on the edf appeared to produce a reasonable density estimate for small censored lakes (Figure ??). In addition, estimates of total lake area are somewhat close to the "true" value (Table 2). However, uncertainty around the frequentist estimates is unrealistically small (Figure 1).

Instead of the essentially fixed a and total lake area estimates produced by the frequentist approach, I found substantial variability in both a (95% CI: 0.86, 0.92) and total area using a Bayesian approach (Fig 3, 4). In particular, the Bayesian 95% credible intervals for both a and total lake area encapsulate the true values (Fig 3, 4). Despite better uncertainty estimates using a Bayesian approach, both the frequentist and Bayesian approaches underestimated the true value of a and total lake area.

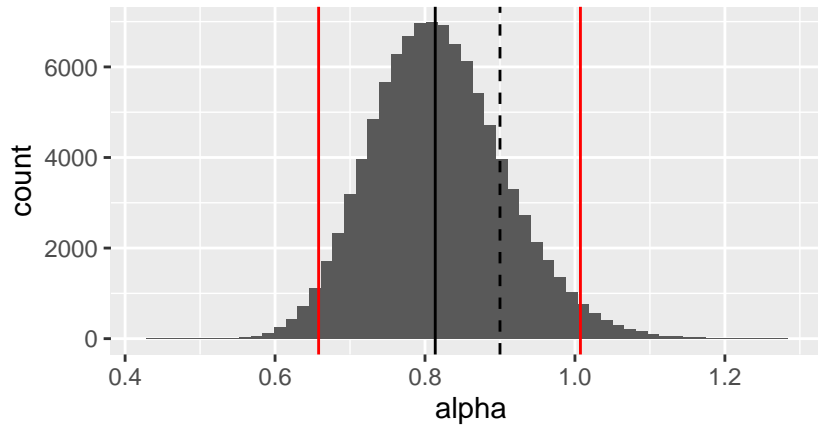


Figure 3: Median (black line) and central 95 percent interval estimates of alpha (red lines). Here the 'true' alpha is 0.9.

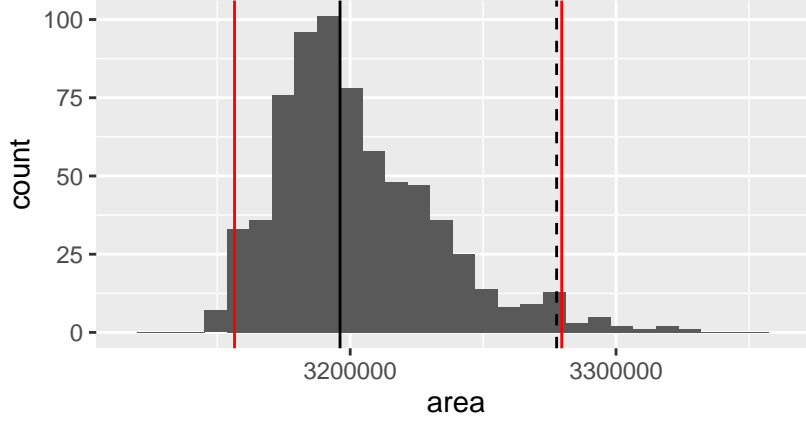


Figure 4: Median (black line) and central 95 percent interval estimates of total lake area (red lines). Here the true total lake area is marked with a dashed vertical line.

Table 2: ‘True’ total lake area from the uncensored edf and estimated lake area from a combination of the censored edf and the estimated cdf.

True Area (km^2)	Estimated Area (km^2)
2.271e+05	2.263e+05

4 Discussion

I have shown that the typical frequentist cutoff method produces reasonable estimates of the density of small censored lakes but that it does not capture uncertainty in total lake area (Table 2, Figure 1). Furthermore, I have shown that models fit using a Bayesian approach indicate substantial uncertainty in both total lake area and the underlying Pareto shape parameter a used to derive these estimates (Figure 3, 4).

Although the 95% credible interval of the Bayesian total lake area estimates encapsulate the true total lake area, the median value underestimates the true total lake area (Figure 4). It is likely that such underestimation, would increase with a greater proportion of censoring. This may explain the steady increase in estimates of global lake area through time from approximately 3 to 5 million km^2 as lake area databases have improved their accuracy (Lehner and Döll, 2004; Downing et al., 2006; Verpoorter et al., 2014). Future work on estimating global lake area should consider implementing a sensitivity analysis looking at the response of total area estimates to variation in the degree of censoring.

A caveat of the present study is that it is difficult to confirm whether or not any particular dataset follows a fractal generating or Pareto process without strong prior knowledge. Indeed data following many heavy-tailed distributions such as the lognormal or negative exponential can appear to be equivalent to power-law distributions such as the Pareto (Clauset et al., 2009). The results of the present study confirm that an apparent change in the shape of the lower tail of a lake area distribution does not necessarily indicate a change in the data generating process. Note that the simulation dataset analyzed herein *is* a truly fractal generated dataset with a homogeneous data generating process yet a cursory look would seem to indicate a change of data-generating process in the lower tail (Figure ??).

In addition to a sensitivity analysis of censoring, future work might consider more complex models that treat lake areas as a mixture of a Pareto distribution for small lakes and either a negative exponential or lognormal distribution for large lakes. Such an approach has been demonstrated by Bonabeau et al. (1999) and Scollnik (2007). Both studies show that the point at which the distribution mixtures converge can provide valuable inference. In the case of lakes, such a convergence point may indicate a change in the data generating process such as the point at which lake areas are controlled by continent placement rather than fractal landscape morphology (Goodchild, 1988; Hamilton et al., 1992).

References

Bonabeau, E., Dagorn, L., and Freon, P., 1999. Scaling in animal group-size distributions. *Proceedings of the National Academy of Sciences*, 96(8):4472–4477. doi:10.1073/pnas.96.8.4472.

- Clauset, A., Shalizi, C. R., and Newman, M. E., 2009. Power-law distributions in empirical data. *SIAM review*, 51(4): 661–703.
- DelSontro, T., Beaulieu, J. J., and Downing, J. A., 2018. Greenhouse gas emissions from lakes and impoundments: Upscaling in the face of global change: GHG emissions from lakes and impoundments. *Limnology and Oceanography Letters*. doi:10.1002/lol2.10073.
- Downing, J. A., Prairie, Y. T., Cole, J. J., Duarte, C. M., Tranvik, L. J., Striegl, R. G., McDowell, W. H., Kortelainen, P., Caraco, N. F., and Melack, J. M., 2006. The global abundance and size distribution of lakes, ponds, and impoundments. *Limnology and Oceanography*, 51(5):2388–2397.
- Goodchild, 1988. Lakes on fractal surfaces: A null hypothesis for lake-rich landscapes. *Mathematical Geology*.
- Hamilton, S. K., Melack, J. M., Goodchild, M. F., and Lewis, W., 1992. Estimation of the fractal dimension of terrain from lake size distributions. *Lowland floodplain rivers: Geomorphological perspectives*. Wiley, pages 145–163.
- Keller, P. S., Marcé, R., Obrador, B., and Koschorreck, M., 2021. Global carbon budget of reservoirs is overturned by the quantification of drawdown areas. *Nature Geoscience*. doi:10.1038/s41561-021-00734-z.
- Lehner, B. and Döll, P., 2004. Development and validation of a global database of lakes, reservoirs and wetlands. *Journal of Hydrology*, 296(1–4):1–22. doi:10.1016/j.jhydrol.2004.03.028.
- McDonald, C. P., Rover, J. A., Stets, E. G., and Striegl, R. G., 2012. The regional abundance and size distribution of lakes and reservoirs in the United States and implications for estimates of global lake extent. *Limnology and Oceanography*, 57(2):597–606. doi:10.4319/lo.2012.57.2.0597.
- Messenger, M. L., Lehner, B., Grill, G., Nedeva, I., and Schmitt, O., 2016. Estimating the volume and age of water stored in global lakes using a geo-statistical approach. *Nature Communications*, 7:13603. doi:10.1038/ncomms13603.
- Newman, M., 2005. Power laws, Pareto distributions and Zipf’s law. *Contemporary Physics*, 46(5):323–351. doi:10.1080/00107510500052444.
- Pi, X., Luo, Q., Feng, L., Xu, Y., Tang, J., Liang, X., Ma, E., Cheng, R., Fensholt, R., Brandt, M., Cai, X., Gibson, L., Liu, J., Zheng, C., Li, W., and Bryan, B. A., 2022. Mapping global lake dynamics reveals the emerging roles of small lakes. *Nature Communications*, 13(1):5777. doi:10.1038/s41467-022-33239-3.
- Scollnik, D. P. M., 2007. On composite lognormal-Pareto models. *Scandinavian Actuarial Journal*, 2007(1):20–33. doi:10.1080/03461230601110447.
- Seekell, D. A. and Pace, M. L., 2011. Does the pareto distribution adequately describe the size-distribution of lakes? *Limnology and Oceanography*, 56(1):350–356.
- Shalizi, C. R., 2017. *Advanced Data Analysis from an Elementary Point of View*.
- Stan Development Team, P., 2017. Stan modeling language users guide and reference manual, version 2.17. 0. *Technical report*.
- Team, R. C. et al. *R: A Language and Environment for Statistical Computing*, 2018.
- Verpoorter, C., Kutser, T., Seekell, D. A., and Tranvik, L. J., 2014. A global inventory of lakes based on high-resolution satellite imagery. *Geophysical Research Letters*, 41(18):6396–6402. doi:10.1002/2014GL060641.
- Winslow, L., Read, J., Hanson, P., and Stanley, E., 2015. Does lake size matter? Combining morphology and process modeling to examine the contribution of lake classes to population-scale processes. *Inland Waters*, 5(1):7–14. doi:10.5268/IW-5.1.740.