

Tools for identifying unexpectedly low microfossil count sums

Richard J. Telford^{*,a}

^a*Department of Biological Sciences, University of Bergen and Bjerknes Centre for Climate Research, Post Box 7803, N-5020 Bergen, Norway*

Abstract

Microfossil counts are a key data type in palaeoecology. Recent work has raised the possibility that some authors might misreport an important quality control parameter, the counts sums, occasionally dramatically so. This paper introduces methods that can flag assemblages with potentially misreported count sums and finds that some assemblage datasets fail these tests.

Keywords: Microfossils assemblages, Counts; Quantitative methods, Reproducible research, Meta-research

1. Introduction

Six percent of papers published in *Molecular and Cellular Biology* show evidence of inappropriate image duplication (Bik et al., 2018) either due to error or, more rarely, misconduct. It would be reckless to assume that the palaeoecological literature does not have equivalent problems: the low rate of retractions in the ecological and geological literature (Grieneisen and Zhang, 2012) may partially reflect our limited ability to detect errors and misconduct rather than their prevalence.

Several numerical tools have been developed to identify questionable data. Deviations from the distribution of digits expected from the Newcomb-Benford law has been used to detect issues with scientific and financial data (Barabesi et al., 2018). Carlisle (2017) identified papers where the baseline differences in means for different treatments were surprisingly high or low given the variance of the data. Brown et al. (2017) developed the granularity-related inconsistency of means (GRIM) test which assesses whether means of integer data are consistent with the reported sample size. However none of these tests are directly applicable to microfossil assemblage data, one of the most common types of palaeoecological data.

A common assertion in papers reporting microfossil assemblage data is that a minimum of N microfossils were counted, where N is often fifty for chironomids and several hundred for pollen and diatoms. This is important because larger count sums are associated with smaller uncertainties, both in relative abundance of taxa and derived statistics such as transfer function reconstructions (Heiri and Lotter, 2001; Payne and Mitchell, 2009). However,

^{*}Corresponding Author
Preprint submitted to *Journal of Quantitative Biology* (Richard J. Telford)

33 mistakes happen and metadata such as count sums can be forgotten once
34 percentages are calculated. In addition, given the time-consuming nature of
35 microfossil counting, especially when preservation is poor or concentrations are
36 low, there may be an incentive to misreport the minimum count sum. The risk
37 that count sums will be misreported is not just theoretical: Larocque-Tobler
38 et al (2015) reported that their chironomid count sums were at least fifty head
39 capsules; a subsequent corrigendum (Larocque-Tobler et al., 2016) acknowledged
40 that count sums were actually as low as nineteen. Telford (2019a) reports some
41 other cases where count sums may be much lower than reported.

42 If the genuine assemblage counts are archived for all taxa, it is trivial to
43 identify undercounts. Unfortunately, count data might be falsified so that it
44 appears to meet the reported count sum. A more common problem is that,
45 regrettably, many palaeoecologists archive percent data without an indication of
46 the count sum. This paper develops simple tests that can flag if the count data
47 might have been misrepresented, or if the count sum of percent data is perhaps
48 lower than reported.

49 The key insight that allows inference about the count sum is that assemblage
50 data are expected to follow the typical features of a rank abundance curve. In
51 particular, because there are many rare taxa in most communities (Darwin,
52 1859), most community or assemblage samples will include taxa represented
53 by a single individual, hereafter singletons, unless the sampling effort is high
54 relative to the taxonomic richness (Coddington et al., 2009). In the occasional
55 assemblages without singletons, the greatest common divisor should usually be
56 one, i.e. few assemblages should have all counts divisible by an integer larger
57 than one, especially if species richness is high.

58 In general, it is not possible to determine the sum from which percent were
59 calculated, but the properties of community and assemblage counts make it
60 possible to estimate it. Given that we expect the rarest taxon to be a singleton,
61 the count sum N can be estimated as $1/p_{min} \times 100$ where p_{min} is the percent
62 abundance of the rarest taxon. This method will fail for assemblages without
63 singletons. We also expect the percent p to be calculated from integer counts,
64 therefore there should be a count sum N such that $p_i/100 \times N$ is, within rounding
65 error, an integer for all i taxa. Possible values for N can be found by a direct
66 search algorithm over the range of plausible values of N . An infinite number of
67 possible count sums that are consistent with the percent exist, but the lowest will
68 give the correct value of N except in cases where the greatest common divisor is
69 greater than one. These tests are closely related to GRIM (Brown and Heathers,
70 2017), as all rely on the granularity of percentages calculated from integer data.

71 This paper aims to test the methods presented above and present some cases
72 with unexpected results. Some complications and caveats are discussed.

73 2. Methods

74 Publicly available datasets were downloaded from the Palaeodata Center,
75 Neotoma, Pangaea and other sources. A range of ecological and palaeoecological
76 data were sought to allow for differences in typical count sums and species

Table 1: Percent of bird counts with singletons and of greatest common divisor (GCD) of one, at different taxonomic levels.

Taxonomic level	Percent with singletons	Percent GCD = 1	Median richness
species	99.95	100.00	54
genus	99.85	100.00	46
family	95.69	100.00	25
order	82.62	99.42	9

richness. Fossil pollen assemblage in the Neotoma database (Williams et al., 2018) (downloaded 2019-05-03) with small count sums (< 50), that appeared to be percentages, or were documented as back-transformed from digitised data were excluded. Datasets with percent data where the minimum count sum was not reported in the associated paper were excluded. Datasets discussed by Telford (2019a) as possibly having under-reported count sums were also excluded to avoid double reporting. Datasets with possible misreporting are anonymised, but no attempt is made to diagnose whether errors or misconduct are responsible. This paper does not attempt to be an exhaustive survey of all the data available.

All analyses were done in R version 3.4.4 (R Core Team, 2018) and used the packages extraDistr version 1.8.10 (Wolodzko, 2018), numbers version 0.7.1 (Borchers, 2018), and countSum 0.0.3 (Telford, 2019b). Code to replicate all the analyses shown above is archived at <https://github.com/richardjtelford/count.check.ms>.

3. Results

3.1. Prevalence of singletons

The vast majority of the over 65,000 bird counts from the North American breeding bird survey (Pardieck et al., 2018) have singletons at the species level (Table 1). To explore the effect of taxonomic richness on the prevalence of singletons, I aggregate the birds counts to progressively lower taxonomic resolutions. As richness declines, the proportion of counts having singletons declines (Table 1), reaching a moderate proportion at the order level, where most of the counts lacking singletons have fewer than five taxa. The vast majority of counts have a greatest common divisor of one, except when taxonomic richness is low.

To test the sensitivity of the prevalence of counts without singletons to the count sum, I resample the counts with 400 or more observations to smaller count sums using a multivariate hypergeometric distribution. At the order level, the proportion of counts with singletons increases steeply until it reaches a maximum at about 200, thereafter it declines slowly (Fig. 1). The prevalence of counts with singletons is low for small counts because it is possible that only the common taxa are counted; with larger counts, the chance of counting a single individual of a rare taxon increases. With even larger counts the prevalence of counts with

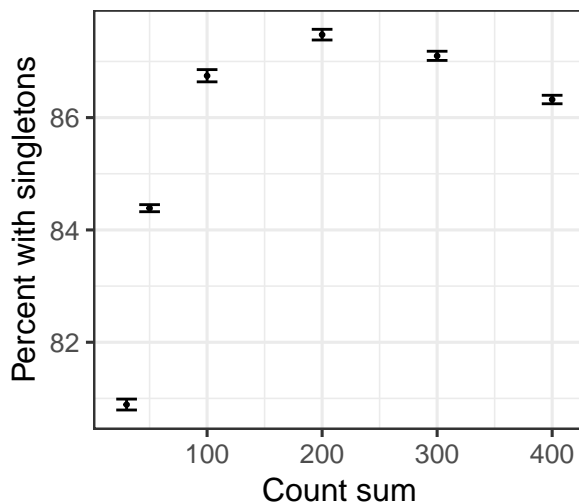


Figure 1: Effect of count sum on the proportion of bird counts without singleton at the order level. Results are the mean of ten trials, error bars are two standard errors.

singletons decreases due to saturation as singletons become doubletons and few new taxa are found.

The 862 diatom counts from Owen’s Lake span a wide range of count sums (1 – 701) because, if concentrations were low, the requirement that at least 300 valves were counted was replaced by a minimum area of slide to count (Bradbury, 1997). Discounting the 20 mono-specific assemblages (all with very low count sums), 97.39% of assemblages have singletons with a median number of singletons per assemblage of five. Assemblages with count sums above 50 have a higher chance of having singletons (98.18% vs 89.04%). 99.52% of assemblages have a greatest common divisor of one, the remainder have a greatest common divisor of two. Assemblages with a greatest common divisor above one all have low taxonomic richness.

The 1775 assemblages in the North American testate amoeba training set with available count data (Amesbury et al., 2018) have count sums between 52 and 456 tests (median = 151), and taxonomic richness ranges between 2 and 31 (median = 15) taxa. 94.87% of the assemblages have at least one singleton (median 4) and 99.72% of assemblages have a greatest common divisor of one. Of the five assemblages with a greatest common divisor above one, three have fewer than five taxa; four have a greatest common divisor of two, and the other one has a greatest common divisor of three.

The pollen data from the Neotoma database included over 182,000 assemblages in 4130 datasets. Nearly all assemblages have singletons (97.04%) and a greatest common divisor of one (99.64%). The assemblages with a greatest common divisor above one are not randomly distributed among the datasets: only 3.37% of the datasets have any such assemblages. Some datasets with a high proportion of assemblages with a greatest common divisor above one are

older datasets which may, despite the available metadata, have been digitised with the loss of rare taxa and precision. For example, one dataset includes 84 assemblages with a median richness of 16 taxa. Count sums vary between 69 and 3201. The greatest common divisor is three for all assemblages, that is all 1300 counts in the dataset are divisible by three. Excluding such datasets would further increase the proportion of assemblages with singletons. Some other Neotoma datasets are discussed in a subsequent section.

3.2. Estimating the percent sum

The Owen’s Lake dataset (Bradbury, 1997) includes percentages for each taxa, allowing the count sums estimated from the percent data to be verified. Excluding mono-specific assemblages, for which no meaningful estimate of the count sum is possible, the minimum percent method and the direct search method correctly estimate the count sum for, respectively, 97.39% and 99.52% of the assemblages. The few assemblages that fail the direct search method are species poor and have low count sums.

The 22 chironomid head capsule assemblage dataset from Last Chance Lake (Axford et al., 2017) also includes both the count sums and the percentage of each taxa. The percentages are given to two decimal places: to account for rounding errors, the countSum package uses the smallest and largest values consistent with the reported percentage. With the minimum percent method the estimated count sums are, within error, either identical to the reported count sum or twice as much (Fig. 2). With the direct search method, the estimated count sums are all exactly twice the reported count sum. The factor of two difference is because all the counts include half head capsules (but only sometimes is the rarest taxon represented by a half head capsule). Reporting half microfossils is common for several microfossil groups, including chironomids and pollen (especially for bisaccate conifers); occasionally other fractions are reported. Half counts make the estimated count too high by a factor of two, so do not risk incorrectly flagging count sums as being too small. The direct search method is more precise than the minimum percent method as the rounding error is relatively smaller on the larger percent values the method uses.

3.3. Unexpected count data

Some of the pollen datasets in the Neotoma database have an inexplicably high proportion of assemblages with a greatest common divisor greater than one. Here I focus on four datasets produced by the same research group. All the datasets, which are relatively recent, include information on the spike used to calculate pollen concentration, so cannot have been digitised. The four datasets have some assemblages with a greatest common divisor of two interspersed, sometimes regularly, amongst assemblages with the expected greatest common divisor of one. The assemblages with a greatest common divisor of one typically have several singletons (Table 2).

Assemblages with a greatest common divisor of one typically have a higher taxonomic richness than assemblages with a greatest common divisor of two

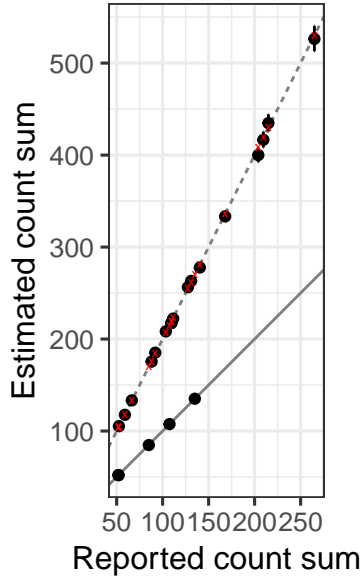


Figure 2: Estimated and reported chironomid count sums by the minimum percent method (solid symbols with error bars) and the direct search method (red crosses) from Last Chance Lake (Axford et al., 2017). Lines show the 1:1 (solid) and 2:1 (dashed) relationships.

Table 2: The number of assemblages, median number of taxa, percent of assemblages with a greatest common divisor (GCD) of one, and the median number of singletons in assemblages with a greatest common divisor of one, in four pollen datasets

Dataset	No. assemblages	Median no. taxa	% GCD = 1	Median no. singletons (GCD = 1)
1	53	36	49	16
2	83	36	61	17
3	27	19	70	7
4	41	19	93	10

(Fig. 3). In the datasets with a wide range of count sums, the richness of the assemblages with a greatest common divisor of two is comparable with the richness of the assemblages with a greatest common divisor of one that have a count sum that is half as large.

One assemblage with a greatest common divisor of two has 73 taxa. Assuming the distribution of counts is (at least locally) uniform, the probability of n counts being divisible by k is $1/k^n$: for a single assemblage with this many taxa, the probability is 1 in 10^{21} . The probability that these counts reflect the true nature of the pollen assemblages is exceedingly low; it is far more likely that the data have been mishandled at some stage in some way.

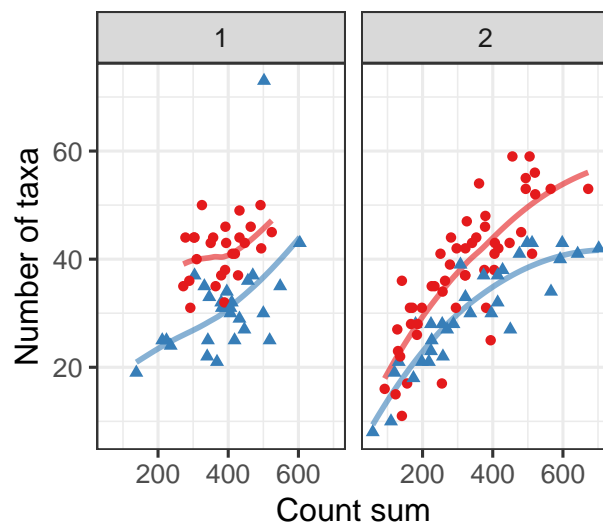


Figure 3: Number of taxa in against count sum for assemblages in two pollen datasets with a greatest common divisor of one (red circles) or two (blue triangles). Lines are loess with a span of 1. All four pollen datasets show the same pattern, but it is clearest in datasets 1 and 2 which have the highest proportion of assemblages with a greatest common divisor of two.

3.4. Unexpected percent data

Dataset diatom1 includes 35 diatom assemblages with a median richness of 19 taxa. The associated paper reports that at least 400 valves were counted from each assemblage, which means that singletons should have a relative abundance of 0.25% or less. However, the rarest taxon in two assemblages have a relative abundance of 3.23%, and the relative abundance of all other taxa in these assemblages are, within rounding error, integer multiples of this. If the count sums actually are at least 400, this would imply that each species count in these assemblages is an integer multiple of at least 13. This seems unlikely. Alternatively, the actual count sum for these assemblages could be as low as 31 valves. In total, six assemblages appear to have a count sum below 400 valves. The results section of the paper acknowledges that counts were below the standard 400 valves for part of the Late Glacial. All the assemblages identified with possible low count sums are in this part of the stratigraphy, demonstrating the utility of the method.

The archived data include 56 taxa that have a maximum abundance of at least 1.25%; this is about half of the 109 taxa reported in the paper. This pruning of rare taxa will increase the risk of counts without singletons, which would cause the minimum percent method to fail, but will have little impact on the direct search method.

Dataset diatom2 has 71 assemblages with a median richness of 57 taxa. The associated paper reports that the diatom counts included 400–500 valves per assemblage except for four diatom-poor assemblages with at least 100 valves.

212 The direct search method estimates that twelve assemblages have a count sum
213 of less than 400, and three have less than 100. The direct search and minimum
214 percent methods agree, within rounding error, on the estimated count sums,
215 implying that either all assemblages have singletons and count sums are low
216 in some assemblages, or that the greatest common divisor of the raw counts
217 is greater than one. This means that for the assemblage where the minimum
218 percent is 7.14, which has seven taxa, all counts would need to be multiples of 8
219 for the count sum to be at least 100.

220 Dataset *chironomid1* has 55 assemblages. Although the associated paper
221 reports that assemblages with count sums below 50 were discarded, the direct
222 search method estimates that three assemblages have count sums between 38
223 and 40. These assemblages are diverse (richness between 15 and 17 taxa), so it
224 is unlikely that the true counts are double that estimated here and the greatest
225 common divisor is two.

226 Palaeoceanographic dataset *marine1* has 160 assemblages with a median
227 taxonomic richness of 34.5 taxa. The associated paper reports that assemblages
228 with count sums below 100 were omitted. The direct search method estimates
229 that five assemblages have count sums below 100 (this excludes one mono-specific
230 assemblage), with count sums as low as 34. With this dataset, the direct search
231 and the minimum percent method give identical results for some assemblages but
232 highly divergent estimates for others, with the direct search method sometimes
233 giving estimates of several thousand, in one case over two orders of magnitude
234 higher than the minimum percent method. Some of the divergence appears to
235 be because some, mainly low diversity, assemblages lack singletons. The most
236 extreme divergences appear to be data entry errors with incorrect rounding,
237 perhaps during taxonomic revisions of the percent data rather than the raw
238 counts.

239 4. Discussion

240 Analysis of the breeding bird, Owen's Lake diatom, testate amoeba and
241 pollen datasets has shown that, as expected, the vast majority of community
242 and assemblage counts have singletons, and an even larger proportion have a
243 greatest common divisor of one. Therefore, these characteristics are potentially
244 useful for identifying data where the count sum is misreported.

245 My search through archived microfossil data found several datasets where
246 count sums are, or appear to be, smaller than that reported. Some of the count
247 sums appear to be an order of magnitude below what was reported. The non-
248 pollen datasets examined here are a convenience sample of available data that
249 met the inclusion criteria. As such, this analysis cannot be used to accurately
250 determine the prevalence of datasets with possible undercounts, but it appears
251 that a non-negligible fraction of the literature is affected. It is possible in
252 some cases that the assemblages with low counts were omitted or merged prior
253 to analysis, but the number of assemblages reported was not updated. Care
254 was taken to identify any such cases by, for example, examining stratigraphic
255 diagrams to see if the low count assemblages were included. In some cases, the

assemblages with low apparent counts can be seen in the published stratigraphic diagrams.

It is also possible that the apparently low count sums are a false positive. However, as taxonomically diverse assemblages with a greatest common divisor greater than one are rare and all the datasets with apparently low count sums contained several such assemblages, this is unlikely.

These tests require the assemblage data which are often not archived. However, if the count sums are very low it is possible to recognise this from a stratigraphic diagram. For example, a chironomid stratigraphy had an assemblage where all five taxa had a relative abundance of about 20%, suggesting a count sum of five, far below the reported fifty chironomids. Another assemblage had six taxa with a relative abundance of just over 10% and a seventh with twice as much, suggesting a count sum of eight chironomids. On request, the authors sent the data. None of the 15 assemblages in the dataset had more than fifty chironomids: the count sums varied between 5 and 40.

4.1. Challenges

Microfossil percent data are often given to two decimal places. This is sufficient precision for the direct search method to have utility with counts sums of several thousand. Some data are only given to the nearest percent, which means that neither method has utility with count sums above 100.

The direct search method will fail if the dataset includes taxa calculated with different count sums, for example pollen sums of trees, shrubs and upland herbs and a pollen and spore sum that also includes pteridiophytes. It should be possible to identify such cases from the meta-data and knowledge of usual practice with different proxies. It will also fail, as shown by dataset marine1, if percentages are incorrectly calculated or rounded.

Some taxonomic groups have microfossils that often come in groups of attached individuals. For example, each diatom cell has two valves which are the counting unit. These valves usually separate during processing, but paired valves are often found, and some taxa such as *Aulacoseira* produce long chains of strongly bound valves. This might make singletons slightly less likely. The results from Owen's Lake suggest that this is not an important problem, as most assemblages have many singletons.

4.2. Consequences

Small undercounts in a few assemblages will have minimal impact on the precision of any palaeoenvironmental reconstruction or other statistics derived from the assemblage. Substantial undercounts will potentially seriously effect the precision of the results. Such undercounts might constitute a data handling error, for example, if samples with low counts were supposed to be merged but that step was forgotten. Substantial and pervasive under-counting could be construed as scientific misconduct due to either negligence or falsification, and action to correct the literature is probably required.

Some papers do not report count sums. When this important quality metric is omitted, the reader should be able to assume that the standard minimum

count sum for the taxonomic group has been used (i.e. 50 for chironomids, several hundred for pollen and diatoms). If the actual count sums are materially below this, then this potentially constitutes falsification by omission (Fanelli, 2013).

5. Conclusions

A non-negligible fraction of the archived assemblage data includes counts that appear to have a count sum below that reported in the paper. Some of the apparent counts are small enough that the uncertainty of the count and derived statistics will be substantially larger than expected. These results highlight the importance of archiving both the raw data and the code required to process the data.

Acknowledgements

Some data were obtained from the Neotoma Paleocology Database (<http://www.neotomadb.org>), and the work of the data contributors and the Neotoma community is gratefully acknowledged. This work was partly supported by Norwegian Research Council project PalaeoDrivers (213607).

References

- Amesbury, M.J., Booth, R.K., Roland, T.P., et al., 2018. Towards a Holarctic synthesis of peatland testate amoeba ecology: Development of a new continental-scale palaeohydrological transfer function for North America and comparison to European data. *Quaternary Science Reviews* 201, 483–500. <https://doi.org/10.1016/j.quascirev.2018.10.034>
- Axford, Y., Levy, L.B., Kelly, M.A., et al., 2017. Timing and magnitude of early to middle Holocene warming in East Greenland inferred from chironomids. *Boreas* 46, 678–687. <https://doi.org/10.1111/bor.12247>
- Barabesi, L., Cerasa, A., Cerioli, A., et al., 2018. Goodness-of-fit testing for the Newcomb-Benford Law with application to the detection of customs fraud. *Journal of Business & Economic Statistics* 36, 346–358. <https://doi.org/10.1080/07350015.2016.1172014>
- Bik, E.M., Fang, F.C., Kullas, A.L., et al., 2018. Analysis and correction of inappropriate image duplication: The Molecular and Cellular Biology Experience. *Molecular and Cellular Biology* 38, e00309–18. <https://doi.org/10.1128/MCB.00309-18>
- Borchers, H.W., 2018. Numbers: Number-theoretic functions. R package version 0.7-1. <https://CRAN.R-project.org/package=numbers>
- Bradbury, J.P., 1997. A diatom-based paleohydrologic record of climate change for the past 800 k.y. from Owens Lake, California, in: An 800,000-year paleoclimatic record from core OL-92, Owens Lake, Southeast California. Geological Society of America. <https://doi.org/10.1130/0-8137-2317-5.99>
- Brown, N.J.L., Heathers, J.A.J., 2017. The GRIM test: A simple technique detects numerous anomalies in the reporting of results in psychology. *Social*

Psychological and Personality Science 8, 363–369. <https://doi.org/10.1177/1948550616673876>

Carlisle, J.B., 2017. Data fabrication and other reasons for non-random sampling in 5087 randomised, controlled trials in anaesthetic and general medical journals. *Anaesthesia* 72, 944–952. <https://doi.org/10.1111/anae.13938>

Coddington, J.A., Agnarsson, I., Miller, J.A., et al., 2009. Undersampling bias: The null hypothesis for singleton species in tropical arthropod surveys. *Journal of Animal Ecology* 78, 573–584. <https://doi.org/10.1111/j.1365-2656.2009.01525.x>

Darwin, C., 1859. On the origin of species by means of natural selection, or the preservation of favoured races in the struggle for life. John Murray, London.

Fanelli, D., 2013. Redefine misconduct as distorted reporting. *Nature* 494, 149. <https://doi.org/10.1038/494149a>

Grieneisen, M.L., Zhang, M., 2012. A comprehensive survey of retracted articles from the scholarly literature. *PLOS ONE* 7, 1–15. <https://doi.org/10.1371/journal.pone.0044118>

Heiri, O., Lotter, A.F., 2001. Effect of low count sums on quantitative environmental reconstructions: an example using subfossil chironomids. *Journal of Paleolimnology* 26, 343–350. <https://doi.org/10.1023/a:1017568913302>

Larocque-Tobler, I., Filipiak, J., Tylmann, W., et al., 2016. Corrigendum to “Comparison between chironomid-inferred mean-August temperature from varved Lake Żabińskie (Poland) and instrumental data since 1896 AD” [Quat. Sci. Rev. 111 (2015) 35–50]. *Quaternary Science Reviews* 140, 163–167. <https://doi.org/10.1016/j.quascirev.2016.01.020>

Larocque-Tobler, I., Filipiak, J., Tylmann, W., et al., 2015. Comparison between chironomid-inferred mean-August temperature from varved Lake Żabińskie (Poland) and instrumental data since 1896 AD. *Quaternary Science Reviews* 111, 35–50. <https://doi.org/10.1016/j.quascirev.2015.01.001>

Pardieck, K., Ziolkowski, D., Lutmerding, M., et al., 2018. North American breeding bird survey dataset 1966 - 2017, version 2017.0. [WWW Document]. <https://doi.org/10.5066/F76972V8>

Payne, R.J., Mitchell, E.A.D., 2009. How many is enough? Determining optimal count totals for ecological and palaeoecological studies of testate amoebae. *Journal of Paleolimnology* 42, 483–495. <https://doi.org/10.1007/s10933-008-9299-y>

R Core Team, 2018. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>

Telford, R.J., 2019a. Review and test of reproducibility of sub-decadal resolution palaeoenvironmental reconstructions from microfossil assemblages.

Telford, R.J., 2019b. CountSum: Check assemblage count sums and percent. R package version 0.0-3.

Williams, J.W., Grimm, E.C., Blois, J.L., et al., 2018. The neotoma paleoecology database, a multiproxy, international, community-curated data resource. *Quaternary Research* 89, 156–177. <https://doi.org/10.1017/qua.2017.105>

Wolodzko, T., 2018. ExtraDistr: Additional univariate and multivariate

386 distributions. R package version 1.8.10. <https://CRAN.R-project.org/package=>
387 extraDistr