

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 appear to fail these tests.

Keywords: Microfossil counts; quantitative methods

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 means test (GRIM) of 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). However, mistakes happen and metadata such as count sums

*Corresponding Author

Email address: `richard.telford@uib.no` (Richard J. Telford)

Preprint submitted to *Quaternary Science Reviews*

April 28, 2019

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

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

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

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

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

72 2. Methods

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

Table 1: Percent of bird counts without singletons and of greatest common divisor (GCD) greater than one of counts without singletons, at different taxonomic levels.

Taxonomic level	Percent without singletons	Percent GCD > 1	Mean richness
species	0.05	0.00	53
genus	0.15	0.00	45
family	4.31	0.07	24
order	17.38	3.37	9

Datasets where the minimum count sum was not reported in the associated paper were excluded. Datasets discussed by Telford (2019) as possibly having under-reported count sums were also excluded to avoid double reporting. Datasets with possible mis-reporting 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 countChecker 0.0.2 (Telford 2018). 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

A tiny minority of the over 65,000 bird counts from the North American breeding bird survey (Pardieck et al. 2018) lack 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 lacking singletons increases (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 without singletons 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 re-sample the counts with 400 or more observations to smaller count sums using a multivariate hypergeometric distribution. At the order level, the proportion of counts without singletons declines steeply until it reaches a minimum at about 200, thereafter it rises slowly (Fig. 1). The prevalence of counts without singletons is high 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 without singletons increases due to saturation as singletons become doubletons and few new taxa are found.

The 862 diatoms counts from Owen’s Lake span a wide range of count sums (1–701) as, if concentrations were low, the requirement that at least 300 valves

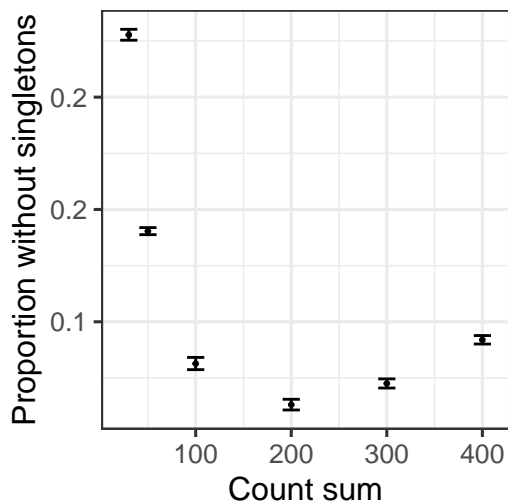


Figure 1: Effect of count sum on the proportion of bird counts without singleton orders. Results are the mean of ten trials, error bars are two standard errors.

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), the median number of singletons per count is five and only 2.6% of assemblages lack singletons. Assemblages with count sums below 50 have a higher chance of lacking singletons (11% vs 1.8%). Of the assemblages without singletons, 81.8% have a greatest common divisor of one. The maximum greatest common divisor is 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. 94.9% of the assemblages have at least one singleton (median 4). In the assemblages that lack singletons, 94.5% 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.

3.2. Estimating the percent sum

The Owen's Lake dataset (Bradbury 1997) also includes percentages for each taxa, allowing the count sums estimated from the percent data to be verified. Excluding monospecific 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, 97.4% and 99.5% of the assemblages. The few assemblages that fail the direct search method are species poor and have low count sums.

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

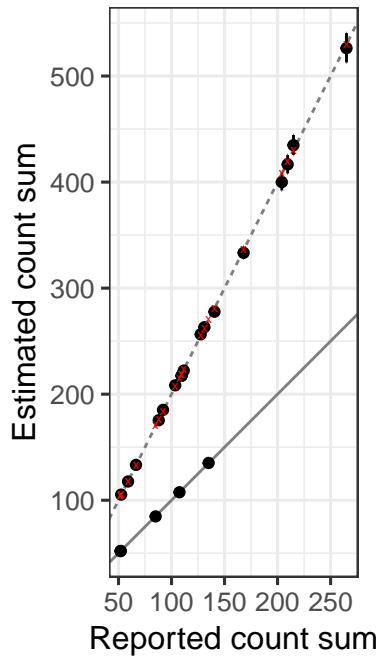


Figure 2: Estimated and reported chironomid count sums by the mininum percent method (solid symbols with error bars) and the direct search method (red crosses) from Last Chance Lake (Axford et al. 2017). Lines show the 1:1 (solid) and 2:1 (dashed) relationships.

150 3.3. Unexpected count data

151 Dataset pollen1 is a 54 count pollen stratigraphy. The published stratigraphy
 152 and the archived dataset generally resemble each other, but there are several

discrepancies. The associated paper reports that count sums are 200 (excluding fern and fungal spores). The archived data have counts sums between 194 and 206, with the exception of one count of 240 (which may be so high because of a possible data entry error). Count sums that are smaller than 200 may be because some rare taxa have been omitted from the archived data. This should normally only cause small discrepancies. More importantly, these count sums include trilete spores, probably from a fern or fern ally, with an abundance as high as 94 spores. Unexpectedly, only one of the counts has a singleton, a *Pinus* grain. The other counts have minimum abundances of two or four pollen grains. With the exception of the single singleton, all 562 species counts are divisible by two. 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 in some way.

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 3.2%, 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 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. The direct search method estimates that twelve assemblages have a count sum of less than 400, and three have less than 100. The direct search and minimum percent methods agree, within rounding error, on the estimated count sums, implying that either all assemblages have singletons and count sums are low in some assemblages, or that the greatest common divisor of the raw counts is greater than one. This means that for the assemblage with the minimum percent is 7.1, which has seven taxa, all counts would need to be multiples of 8 for the count sum to be at least 100.

Dataset chironomid1 has 55 assemblages. Although the associated paper reports that assemblages with count sums below 50 were discarded, the direct search method estimates that three assemblages had count sums between 38 and 40. These assemblages are diverse (richness between 15 and 17 taxa), so it is

197 unlikely that the true counts are double that estimated here and the greatest
198 common divisor is two.

199 Palaeoceanographic dataset marine1 has 160 assemblages with a median
200 taxonomic richness of 34.5 taxa. The associated paper reports that assemblages
201 with count sums below 100 were omitted, and that there were 20 assemblages with
202 count sums below 200. The direct search method estimates that five assemblages
203 have count sums below 100 (this excludes one mono-specific assemblage), with
204 count sums as low as 34. With this dataset, the direct search and the minimum
205 percent method give identical results for some assemblages but highly divergent
206 estimates for others, with the direct search method sometimes giving estimates
207 of several thousand, in one case over two orders of magnitude higher than the
208 minimum percent method. Some of the divergence appears to be because some,
209 mainly low diversity, assemblages lack singletons. The most extreme divergences
210 appear to be data entry errors with incorrect rounding, perhaps during taxonomic
211 revisions of the percent data rather than the raw counts.

212 4. Discussion

213 My search through archived microfossil count data found several datasets
214 where count sums are, or appear to be, smaller than that reported. Some of the
215 count sums appear to be an order of magnitude below what was reported. The
216 datasets examined are a convenience sample of available count data that met the
217 inclusion criteria. As such, this analysis cannot be used to accurately determine
218 the prevalence of datasets with possible undercounts, but it appears that a
219 non-negligible fraction of the literature is affected. It is possible in some cases
220 that the assemblages with low counts were omitted or merged prior to analysis,
221 but the number of assemblages reported was not updated. Care was taken to
222 identify any such cases by, for example, examining stratigraphic diagrams to see
223 if the low count assemblages were included. In some cases, the assemblages with
224 low apparent counts can be seen in the published stratigraphic diagrams.

225 It is also possible that the apparently low count sums are a false positive.
226 This can only happen in assemblages that have a greatest common divisor greater
227 than one. The analysis of the breeding bird, testate amoeba and Owen's Lake
228 datasets shows that assemblages without singletons are rare, and assemblages
229 with a greatest common divisor greater than one are very rare, especially for
230 taxonomically diverse assemblages. If a dataset contains several taxonomically
231 diverse assemblages that have a greatest common divisor of two (or especially if
232 greater than two) doubts should be raised about whether the counts sums are
233 accurately described.

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

238 The direct search method will fail if the dataset includes taxa calculated
239 with different count sums, for example pollen sums of trees, shrubs and upland
240 herbs and a pollen and spore sum that also includes pteridiophytes. It should

241 be possible to identify such cases from the meta-data and knowledge of usual
242 practice with different proxies. It will also fail, as show by dataset marine1, if
243 percentages are incorrectly calculated or rounded.

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

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

259 Some papers do not report count sums. When this important quality metric
260 is omitted, the reader should be able to assume that the standard minimum
261 count sum for the taxonomic group has been used (i.e. 50 for chironomids, several
262 hundred for pollen and diatoms). If the actual count sums are materially below
263 this, then this potentially constitutes falsification by omission (Fanelli 2013).

264 5. Conclusions

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

271 Acknowledgements

272 Some data were obtained from the Neotoma Paleoecology Database (<http://www.neotomadb.org>),
273 and the work of the data contributors and the Neotoma community is gratefully
274 acknowledged.

275 References

276 Amesbury MJ, Booth RK, Roland TP et al (2018) Towards a Holarctic synthe-
277 sis of peatland testate amoeba ecology: Development of a new continental-scale

278 palaeohydrological transfer function for North America and comparison to Euro-
279 pean data. *Quaternary Science Reviews* 201:483–500. doi: 10.1016/j.quascirev.2018.10.034

280 Axford Y, Levy LB, Kelly MA et al (2017) Timing and magnitude of early to
281 middle Holocene warming in East Greenland inferred from chironomids. *Boreas*
282 46:678–687. doi: 10.1111/bor.12247

283 Barabesi L, Cerasa A, Cerioli A, Perrotta D (2018) Goodness-of-fit test-
284 ing for the Newcomb-Benford Law with application to the detection of cus-
285 toms fraud. *Journal of Business & Economic Statistics* 36:346–358. doi:
286 10.1080/07350015.2016.1172014

287 Bik EM, Fang FC, Kullas AL et al (2018) Analysis and correction of inap-
288 propriate image duplication: The Molecular and Cellular Biology Experience.
289 *Molecular and Cellular Biology* 38:e00309–18. doi: 10.1128/MCB.00309-18

290 Borchers HW (2018) Numbers: Number-theoretic functions

291 Bradbury JP (1997) A diatom-based paleohydrologic record of climate change
292 for the past 800 k.y. from Owens Lake, California. In: An 800,000-year paleocli-
293 matic record from core OL-92, Owens Lake, Southeast California. Geological
294 Society of America

295 Brown NJL, Heathers JAJ (2017) The GRIM test: A simple technique detects
296 numerous anomalies in the reporting of results in psychology. *Social Psychological*
297 *and Personality Science* 8:363–369. doi: 10.1177/1948550616673876

298 Carlisle JB (2017) Data fabrication and other reasons for non-random sam-
299 pling in 5087 randomised, controlled trials in anaesthetic and general medical
300 journals. *Anaesthesia* 72:944–952. doi: 10.1111/anae.13938

301 Coddington JA, Agnarsson I, Miller JA et al (2009) Undersampling bias:
302 The null hypothesis for singleton species in tropical arthropod surveys. *Journal*
303 *of Animal Ecology* 78:573–584. doi: 10.1111/j.1365-2656.2009.01525.x

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

306 Fanelli D (2013) Redefine misconduct as distorted reporting. *Nature* 494:

307 Grieneisen ML, Zhang M (2012) A comprehensive survey of retracted ar-
308 ticles from the scholarly literature. *PLOS ONE* 7:1–15. doi: 10.1371/jour-
309 nal.pone.0044118

310 Heiri O, Lotter AF (2001) Effect of low count sums on quantitative envi-
311 ronmental reconstructions: an example using subfossil chironomids. *Journal of*
312 *Paleolimnology* 26:343–350. doi: 10.1023/a:1017568913302

313 Larocque-Tobler I, Filipiak J, Tylmann W et al (2015) Comparison be-
314 tween chironomid-inferred mean-August temperature from varved Lake Żabińskie
315 (Poland) and instrumental data since 1896 AD. *Quaternary Science Reviews*
316 111:35–50. doi: 10.1016/j.quascirev.2015.01.001

317 Larocque-Tobler I, Filipiak J, Tylmann W et al (2016) Corrigendum to “Com-
318 parison between chironomid-inferred mean-August temperature from varved
319 Lake Żabińskie (Poland) and instrumental data since 1896 AD” [*Quat. Sci.*
320 *Rev.* 111 (2015) 35–50]. *Quaternary Science Reviews* 140:163–167. doi:
321 10.1016/j.quascirev.2016.01.020

322 Pardieck K, Ziolkowski D, Lutmerding M, Hudson M-A (2018) North Ameri-
323 can breeding bird survey dataset 1966 - 2017, version 2017.0. <https://doi.org/10.>

324 5066/F76972V8
325 R Core Team (2018) R: A language and environment for statistical computing.
326 R Foundation for Statistical Computing, Vienna, Austria
327 Telford RJ (2019) Review and test of reproducibility of sub-decadal resolution
328 palaeoenvironmental reconstructions from microfossil assemblages
329 Telford RJ (2018) CountChecker: Check assemblage count sums and percent
330 Wolodzko T (2018) ExtraDistr: Additional univariate and multivariate dis-
331 tributions