

1 Issues in the statistical detection of data
2 fabrication and data errors in the scientific
3 literature: simulation study and reanalysis of
4 Carlisle, 2017

5 *Scott W. Piraino*

6 *20 August, 2017*

7 **Abstract**

8 **Background:** The detection of fabrication or error within the scientific literature is an
9 important and underappreciated problem. Retraction of scientific articles is rare, but
10 retraction may also be conservative, leaving open the possibility that many fabricated or
11 erroneous findings remain in the literature as a result of lack of scrutiny. A recently statistical
12 analysis of randomized controlled trials [1] has suggested that the reported statistics from
13 these trials deviate substantially from expectation under truly random assignment, raising
14 the possibility of fraud or error. It has also been proposed that the method used could be
15 implemented to prospectively screen research, for example by applying the method prior to
16 publication.

17 **Methods and Findings:** To assess the properties of the method proposed in [1], I carry
18 out both theoretical and empirical evaluations of the method. Simulations suggest that the
19 method is sensitive to assumptions that could reasonably be violated in real randomized

20 controlled trials. This suggests that deviation for expectation under this method can not
21 be used to measure the extent of fraud or error within the literature, and raises questions
22 about the utility of the method for prospective screening. Empirical analysis of the results
23 of the method on a large set of randomized trials suggests that important assumptions may
24 plausibly be violated within this sample. Using retraction as a proxy for fraud or serious
25 error, I show that the method faces serious challenges in terms of precision and sensitivity for
26 the purposes of screening, and that the performance of the method as a screening tool may
27 vary across journals and classes of retractions.

28 **Conclusions:** The results in [1] should not be interpreted as indicating large amount of fraud
29 or error within the literature. The use of this method for screening of the literature should
30 be undertaken with great caution, and should recognize critical challenges in interpreting the
31 results of this method.

32 Introduction

33 Meta-research, a scientific endeavor aimed at studying and improving the process of science
34 itself, has gained increasing interests among scientists. This interest has partially been driven
35 by theoretical [2,3] and empirical work [4,5] that suggests concerns about the validity of
36 the published scientific literature. One area of the scientific process that is amenable to
37 meta-research is the detection of data validity/data integrity issues within the literature.
38 Methods such as statcheck [6] and granularity testing [7] and its variants [8] have been
39 developed to identify possible data validity issues by checking summary statistics reported in
40 published research for consistency. In some cases, it has been proposed that the method be
41 applied in an automated manner at various stages of the scientific process, for instance, prior
42 to publication [6,9].

43 One class of methods for the detection of data validity issues is based on detecting whether
44 data or summary statistics are consistent with their expected statistical distribution [10–15].

45 Under this framework, large deviations from the expected distribution of reported data are
46 interpreted as indication of possible data integrity issues. In several cases within the literature,
47 this method has been used to flag publications that were later determined to be based on
48 fabricated data [11,12]. One variation on this, developed by Carlisle [11,13], uses reported
49 summary statistics on baseline variables from randomized clinical trials to score published
50 trials in terms of statistical deviation from that expected if subjects were truly assigned at
51 random to various experimental groups. Large deviations potentially suggest issues with the
52 validity of the reported summary statistics.

53 If methods for the detection of data validity issues are to play an increasing role in the
54 scientific process, it is critical that scientists have a good understanding of the appropriate
55 interpretation of these kinds of procedures. Of particular concern is that scientists may
56 interpret the fact that a study or numerical result is flagged by these methods as substantial
57 evidence of some type of flaw even when the method can sometimes flag an analysis for other
58 reasons [16,17]. Especially if such methods are used to systematically screen research, it will
59 be essential for scientists to have a grasp on the limitations that these methods may face.
60 In order to understand these limitations, it is useful to distinguish between multiple types
61 of numerical results that may be identified by these methods. In what follows, I make a
62 distinction between two different threats to data validity: data fabrication and data errors.
63 Data fabrication may be said to occur when authors of published research intentionally alter
64 data that they report in a way that is not consistent with how the data was collected or by
65 reporting fictitious results about data that was never actually collected. Data errors may be
66 said to occur when authors of published research unintentionally report data in a way that is
67 not consistent with what was actually observed, for example by unintentional typographical
68 errors or accidental errors in numerical calculations. Often, methods aimed at detecting
69 data validity issues can be expected to flag both data fabrication and data errors, without
70 distinguishing between the two. In principle, this fact does not preclude the use of these
71 methods for screening scientific research, because both fabrication and honest errors should

72 be detected and corrected. However, the fact that these methods can not distinguish between
73 errors and fabrication presents important interpretational challenges, since parties involved in
74 the process are likely to respond differently if they interpret a flag by one of these methods as
75 evidence of fabrication vs evidence of error. As a result, it is critical to manage expectations
76 about what these methods show and how they should be applied.

77 Potentially of more concern for the application of these methods is the possibility that some
78 of the numerical results flagged by these methods are not erroneous, or in other words, are
79 false positives. This may happen when there are aspects of data generation or reporting
80 which are entirely legitimate but which are not accounted for by the method used to detect
81 data validity issues. For example, it has been suggested that statcheck, which focuses on
82 p-values, may result in false positives in cases where p-values are corrected (e.g. for multiple
83 comparisons) but this correction is not taken into account [16,18]. Numerical results that are
84 flagged as a result of these types of benign issues are problematic from a screening perspective
85 because they can result in a waste of resources if all flagged analysis are investigated for
86 potential data validity issues, as well as bringing unfair suspicion upon honest scientists.
87 Understanding the relative frequencies of these different categories: fabrication, honest errors,
88 and false positives, among flagged results is essential for the proper interpretation of these
89 methods.

90 Although these issues are generally applicable to methods aimed at identifying data validity
91 issues, they are particularly timely in light of a recent analysis by Carlisle [1], which applied
92 a data validity detection method to a large sample of randomized controlled trials. This
93 analysis has already generated significant attention both within the scientific literature [9] as
94 well as the in the press [19,20]. The importance of [1] can be seen as relating to two related
95 issues:

96 First, deviation from the expected statistical distribution of results across many clinical
97 trials has implications for the global rate of fabrication or errors within the literature. This

98 interpretation is apparent in the coverage of [1] (see [9], speculating that the results of
99 [1] possibly indicate a “tsunami” of previously unrecognized fabrication in the literature,
100 or [19], addressing the possible frequency of fabrication in literature based on figures from
101 [1]). Fabrication and errors may be difficult to detect and may persist un-noticed in the
102 literature [21], leaving open the possibility that these occurrences are not as rare as scientists
103 might hope or desire [17,22]. This fact, combined with the observation that automated
104 methods for error detection sometimes flag large proportions of the literature compared to
105 what might be expected is potentially alarming. Certainly it appears that some observers
106 have considered this interpretation [9,19,20]. The possibility that a large proportion of of
107 the scientific literature contains errors or fabrication would raise important question for
108 the scientific community, and the suspicion that this is the case likely underlies part of the
109 recent interest in meta-research. As a result, understanding the implications of methods and
110 results such as those presented by Carlisle [1] for the overall rate of data validity issues in the
111 literature is highly relevant.

112 Second, the method utilized by [1] is already being used to screen papers submitted to
113 Anaesthesia [9], the journal in which [1] was published. The appropriate interpretation of
114 methods for the detection of data validity issues in terms of screening the published literature,
115 either retrospective or particularly prospectively (e.g. as a condition of publication) is an
116 essential questions that remains to be addressed in the meta-research literature. Additionally,
117 the editors of Anaesthesia have decided [9] to contact the journals associated with trials
118 flagged by the analysis in [1], suggesting the need for serious investigations into these papers
119 on the basis of the analysis in [1]. This suggests that understanding of the appropriate
120 interpretation of [1] is urgently needed. In this article, I analyze the method utilized by
121 Carlisle [1] to give insight into how it should be interpreted and what conclusions can be
122 drawn from about these two critical questions.

123 Results

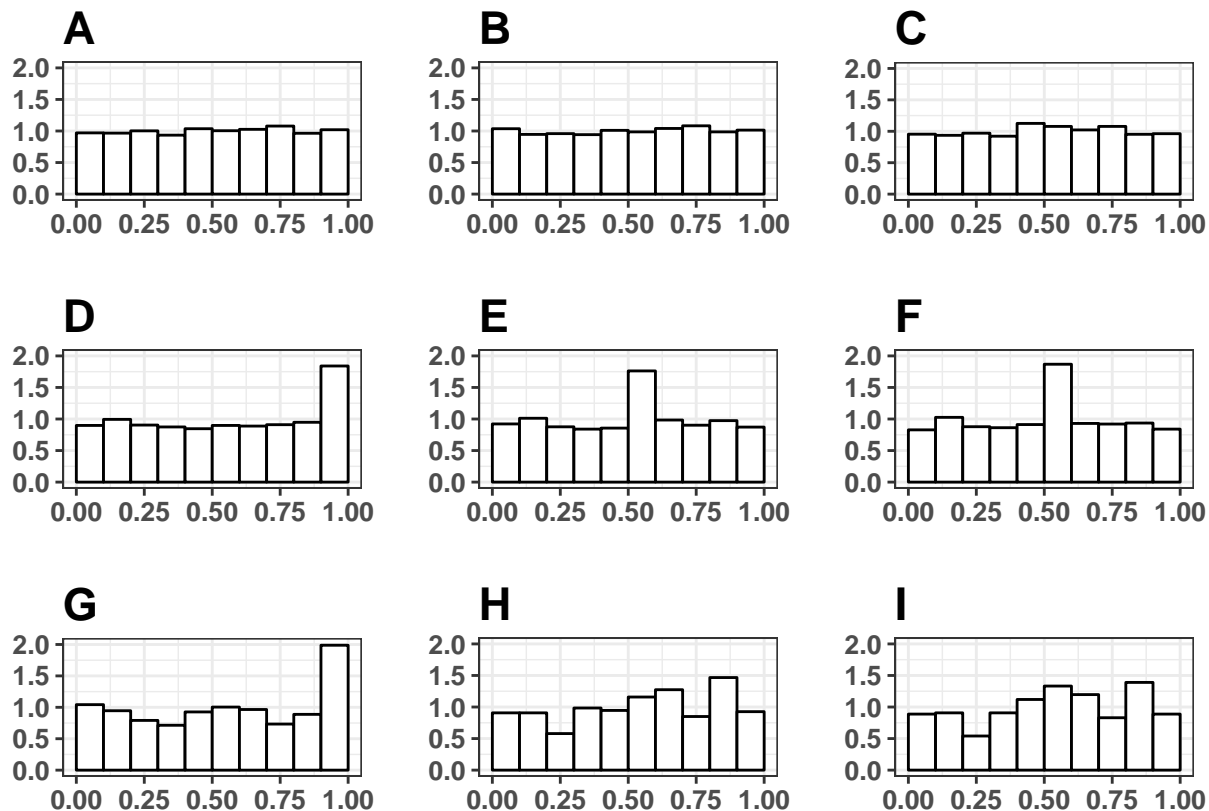
124 To facilitate understanding of the theoretical and empirical results I present, I briefly review
125 the method utilized by Carlisle (which I refer to as the CM) [1]. For a single randomized
126 controlled trial, the CM first involves manually extracting summary statistics on baseline
127 (pre-treatment) variables from all groups which are randomized. For each variable, a p-value is
128 calculated which tests the null hypothesis that the population means of the variable are equal
129 across the groups. If the groups were truly assigned at random, then the null hypothesis
130 is expected to be true for all of the variables. To combine the tests for all variables, the
131 CM as applied in [1] utilized several methods for combining p-values that test a common
132 null hypothesis, but [1] focuses on Stouffer's method [23], which transforms the p-values
133 to z-scores and calculates their sum. Under the assumption that the p-values included are
134 independent, this sum is then compared to its own null distribution to derive a global p-value.
135 Below, I highlight several stages at which this process may go wrong, along with re-analyses
136 of the data used in [1] showing that these issues plausibly effected the analysis.

137 Calculation of variable-level p-values from summary statistics

138 The CM as implemented in [1] involves calculating p-values for for the differences in means of
139 individual baseline variables within each trial using summary statistics, and then aggregating
140 the p-values across each trial. Issues in the calculation of the p-value for each variable may
141 impact the validity of the downstream analysis. In order to test the ability of the method
142 used by Carlisle [1] to recalculate p-values from summary statistics, I simulate data from
143 two identically distributed groups and apply two of the p-value calculation methods used by
144 Carlisle, a Monte Carlo method and ANOVA. The null hypothesis is true in these simulations,
145 so the distribution of p-values should be uniform. Deviations from uniformity could indicate
146 problems with the recalculated p-values, and could explain the deviations from uniformity
147 that Carlisle [1] observed. To assess the robustness of these methods to assumption violations,

148 I include two potentially problematic issues as part of my simulations. First, I simulate
149 data from log-normal distributions instead of normal, as assumed in Carlisle's analysis [1,13].
150 Second, I include rounding of reported summary statistics.

151 Fig 1 presents the distribution of simulated p-values. P-values generated from data with an
152 underlying log-normal distribution and rounding to 2 digits (Fig 1A-C) display a roughly
153 uniform distribution. Following Carlisle [1], I consider the closest p-value to 0.5 from multiple
154 methods. When the underlying distribution is log-normal with 2 digit rounding this p-value
155 has a slight excess near the center of the distribution compared to uniform. When some
156 of the p-values are subject to extreme rounding (Fig 1D-F) the distributions display large
157 excesses of p-values compared to uniform either near 1 (for ANOVA (Fig 1D)) or near 0.5
158 (for Carlisle's Monte Carlo method (Fig 1E) or the closest to 0.5 of ANOVA and Monte Carlo
159 (Fig 1E)). The observed p-values (Fig 1G-I) from [1] display some of these properties, with
160 some large spikes of p-values near certain values, as well as possibly some excess in the center
161 of the distributions.



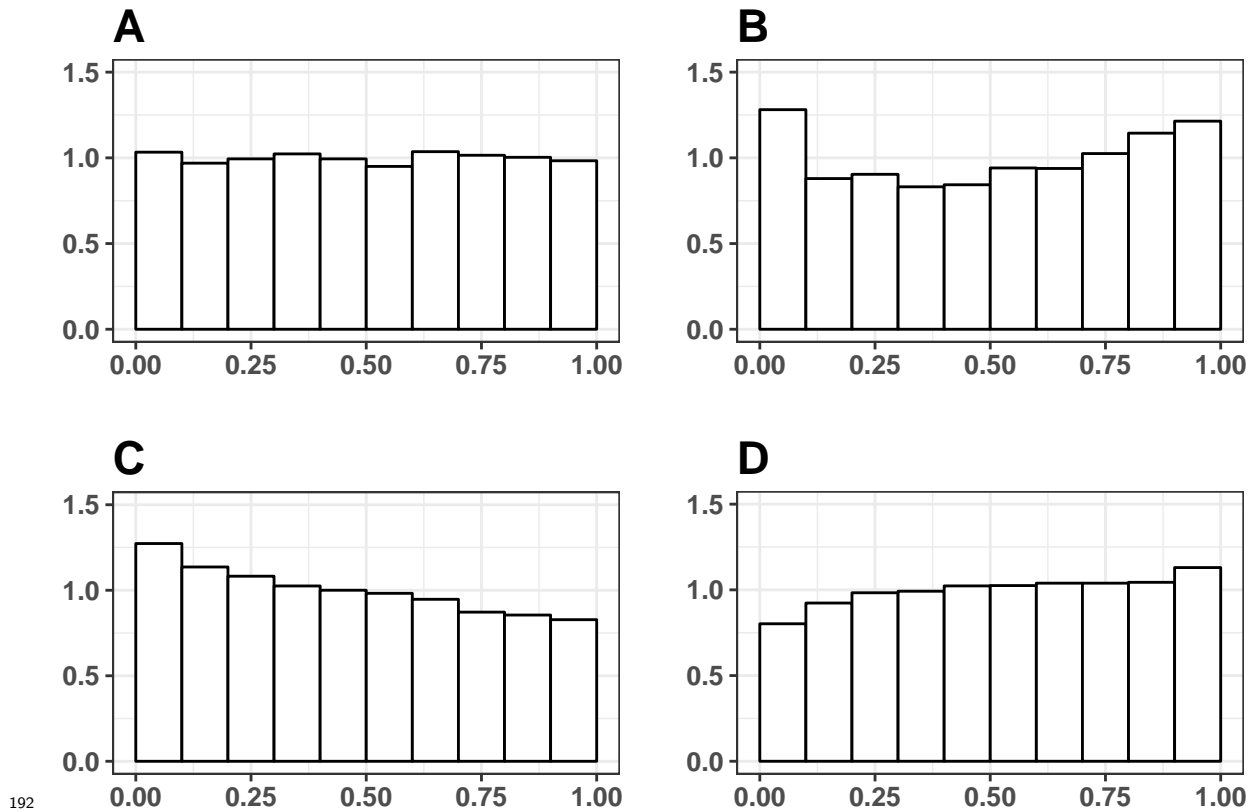
162

163 **Fig 1:** Distribution of simulated (A-F) and observed (G-I) p-values. The first row shows
164 simulated p-values from summary statistics generated from log-normal distributions with
165 moderate rounding for ANOVA (A) Carlisle's Monte Carlo method (B) and the closed of
166 those two to 0.5 (C). The second row shows simulated p-values from the model with summary
167 statistics generated from a normal distribution where 90% of statistics have moderate rounding
168 and 10% have extreme rounding for ANOVA (D) Monte Carlo (E) and the closest of the two
169 to 0.5 (F). The third row shows observed p-values from the data collected by Carlisle from
170 the Journal of the American Medical Association (JAMA) for ANOVA (G) Monte Carlo (H)
171 and the closest of the two to 0.5 (I). For all rows, the first column shows ANOVA p-values,
172 the second Monte Carlo p-values, and the third the closest of ANOVA and Monte Carlo to
173 0.5.

174 **Factors effecting trial-level p-values**

175 Even if the variable-level p-values are validly calculated, issues may arise when multiple
176 variable p-values are aggregated at the level of each trial. In Fig 2 I present simulations
177 showing deviation from the expected null distribution of aggregated p-values (Fig 2A) under
178 three conditions unrelated to data validity (Fig 2B-D). Fig 2B shows the distribution of
179 trial p-values when the baseline variables that are aggregated are correlated with each other.
180 The p-values show a pattern of excess p-values near 0 and 1, just as [1] observed. Fig 2C
181 shows the distribution of p-values when there is imperfect randomization, resulting in residual
182 confounding influence of the baseline variables. In this case, the p-values are right-skewed.
183 Fig 2D shows the distribution of p-values when treatment assignment is randomized within
184 strata that are associated with the baseline variables, resulting in left-skewed p-values. In
185 all three cases, the p-value distributions have an excess of extreme p-values relative the the
186 expected uniform distribution. If the CM is applied to a study for the purposes of screening,
187 and one of these factors is application to the study, the CM may produce produce an extreme

188 p-value for that study as a result of one of these factors rather than as a result of data
189 validity issues. Likewise, for the global assessment of the prevalence of data validity issues
190 in randomized control trials, deviations like those observed in [1] may be the result of a
191 combination of these factors rather than a high prevalence of data validity issues.

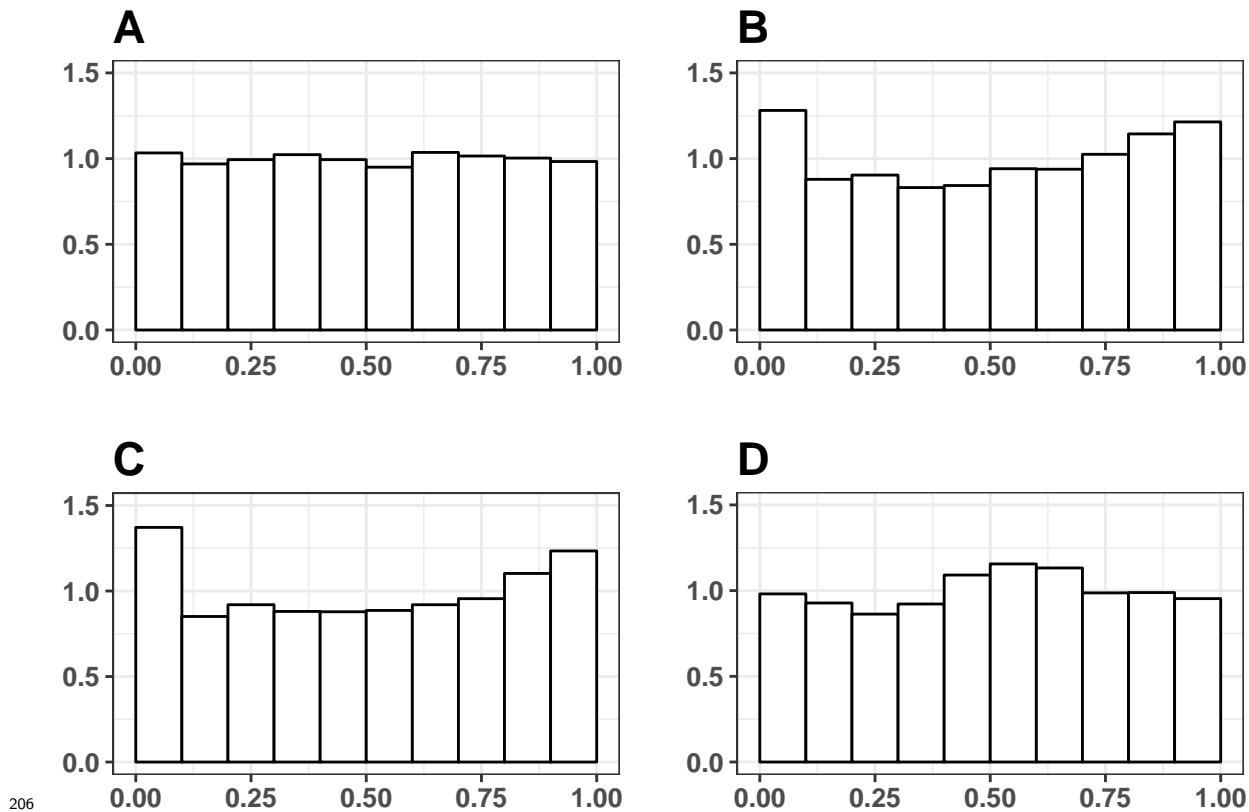


192
193 **Fig 2:** Histograms of p-values simulated from four different models. Null, with an expected
194 uniform distribution (A) correlated baseline variables (B) imperfect randomization (C) and
195 stratification (D).

196 **The Carlisle analysis is plausibly impacted by these issues**

197 To determine if these issues plausibly played a role in the analysis conducted by Carlisle in [1],
198 I reanalyzed the data from the supplement of [1]. Fig 3 compared theoretical distributions
199 derived from simulations (Fig 3A and B, top row) with p-values from [1] (Fig 3C and D,
200 bottom row). Fig 3A and B give the simulation distributions for null p-values and correlated

201 baseline variable p-values, respectively. Fig 3C shows the distribution of p-values for all trials
202 in [1], aggregated by Stouffer's method. As Carlisle [1] notes, this distribution has an excess of
203 p-values near 1 and 0 relative to the null (Fig 3A). However, this distribution is remarkably
204 similar to the simulated distribution with correlated baseline variables (Fig 3B), suggesting
205 that correlated variables could plausibly explain the deviations for uniformity.



206
207 **Fig 3:** Distribution of trial-level p-values for simulated null distribution (A), simulated
208 p-values with correlated variables (B), observed trial-level p-values (C), and observed p-values
209 for the first variable in each trial (D).

210 One objection to this is that the similarity between Fig 3B and C is sensitive to the simulation
211 parameters. Indeed, I chose the parameters for the simulation intentionally to make the
212 point that correlation can result in a similar p-value distribution. Other parameter settings
213 can produce distributions which are less similar. In general, it is difficult to assess how
214 realistic the simulation parameters are. For example, it could be argued that the correlation
215 I used in this simulation (0.33) is higher than generally expected. On the other hand, this

216 does not preclude correlation as an explanation for the results observed by Carlisle [1]. For
217 instance, I assume all trials report 5 baseline variables. Trials that report more variables
218 can have extreme deviations from uniform with lower correlations. Likewise, even if most
219 correlations are lower, there may be some trials with extremely high correlation, or it could
220 be that multiple factors (correlated variables, stratification, confounding) combine to form
221 the observed distribution.

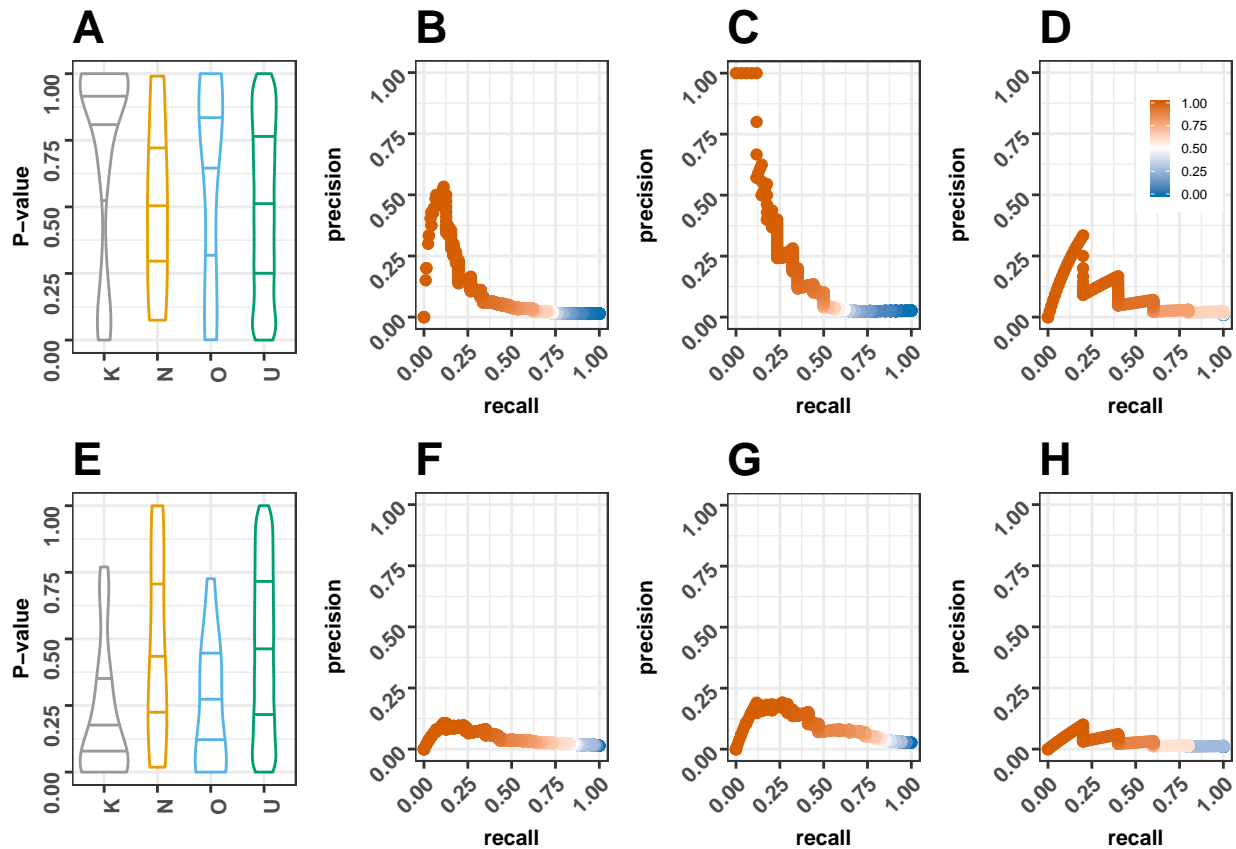
222 In general, it is difficult to definitively identify the cause of the deviation from uniformity
223 using simulations alone. To overcome this, I examined the distribution of the first p-value for
224 each trial (Fig 3D). Some causes of deviation from uniformity, such as fabrication or error, are
225 expected to manifest on the level of individual variable-level p-values. Other causes, such as
226 correlated variables, are expected to manifest when the p-values are aggregated. Comparison
227 of the first variable p-values (Fig 3D) with the aggregated p-values (Fig 3C) can suggest
228 what effects these different sets of explanations may have. The first variable p-values have a
229 qualitatively different appearance compared to the aggregated p-values, lacking the excess
230 of p-values near 0 and 1. This raises the possibility that the excess extreme p-values are
231 due to some issue with the aggregation process, rather than with the individual p-values
232 themselves. The first variable p-values also display an excess of p-values in the center of
233 the distribution relative to uniform. This may result from issues with the calculation of the
234 individual p-values, as discussed above.

235 **Evaluation of the ability of the CM to identify retracted trials**

236 The above analysis suggests that extreme trial-level p-values derived from the CM don not
237 necessarily indicate data validity errors. However, this does not necessarily preclude the
238 usefulness of the CM for the detection of data fabrication and data errors. If the CM can
239 identify known cases of fabrication or error in practice, then that empirical usefulness could
240 form the basis for interpretation of the CM. Indeed, Carlisle analyzed retracted trials and

241 showed that the CM p-values for retractions are more extreme compared to unretracted trials
242 [1]. I extend this analysis by evaluation the distributions of the trial-level p-values from [1]
243 across several retraction categories (Fig 4A and E).

244 Using information contained in the supplemental materials of [1], I place each trial in one
245 of four categories, based on its retraction status. I first divide the trials into those that
246 have been retracted vs those that have not. I further divide the retracted trials into three
247 categories, starting by dividing them based on mention of fabrication in text descriptions of
248 the retractions extracted by Carlisle. For those trials where fabrication is mentioned (and for
249 which it is likely the reason for the retraction), I categorize the trials based on the presence
250 of certain author names in the retraction descriptions. The CM has previously been used by
251 Carlisle [11,13,15] to identify studies by several authors as potentially fraudulent. Several
252 of these sets of studies are highlighted in the text of [1]. I separately classify putatively
253 fraud-based retractions based on the presence of these authors or other known authors of
254 prominent anesthesia-related fraud cases to assess the possibility that the association of
255 trial-level p-values and retraction status differs between these groups.



256

257 **Fig 4:** Assessment of the association between both one-sided (A-D, top row) and two-sided
258 (E-H, bottom row) p-values with retraction status. Violin plots of one-sided (A) and two-sided
259 (E) p-values for each of four trial categories. Horizontal lines represent 25th, 50th, and 75th
260 percentiles. The categories represent unretracted trials (“U”, green), trials retracted without
261 indication of fabrication (“O”, blue), trials likely retracted for fabrication that were prominent
262 examples known in the anesthesia community based on author names (“K”, grey), and trials
263 likely retracted for fabrication that were not prominent examples known in the anesthesia
264 community (“N”, orange). Panels B-D plot precision-recall curves for one-sided trial p-values
265 for all trials (B), for trials in the journal Anesthesia and Analgesia (C), and for trials in the
266 Journal of the American Medical Association (D). Panels E-G plot precision-recall curves
267 for the inverse of the two-sided trial p-values ($1 - p$) for all trials (F), for trials in the
268 journal Anesthesia and Analgesia (G), and for trials in the Journal of the American Medical
269 Association (H). Color of the points in the precision recall curves indicate the threshold value

270 used.

271 Fig 4 shows the distribution of trial-level p-values across these four groups. Fig 4A shows the
272 distribution of one-sided p-values, while fig 4E shows two-sided p-values. Consistent with
273 the observations by Carlisle [1], trials by previously suspected anesthesiology-related authors
274 (based on author name) (“K”, grey), and to a lesser extent trials retracted for putatively
275 non-retraction related reasons both have abnormal p-value distributions, with the one-side
276 p-values displaying an excess of p-values near 1 and a smaller excess near 0 (Fig 4A), and
277 the two-sided p-values shifted toward 0 (Fig 4E). Both unretracted trials (“U”, green), and
278 trials putatively retracted for fabrication that were not prominently known in the anesthesia
279 community (“N”, orange) have p-value distributions much closer to uniform.

280 I also evaluate the ability of the trial p-values to identify retracted trials in terms of the
281 precision (also called positive predictive value) and recall (also called sensitivity) of the
282 p-values at various thresholds. I plot the results in precision-recall curves (Fig 4B-D and
283 F-H), which displays precision-recall pairs when a p-value threshold is used to classify trials
284 as retracted vs unretracted, for many different thresholds. The two-sided p-values (Fig 4 F-H,
285 bottom row) show generally poor performance that is inferior to the one-sided p-values (Fig
286 4B-D, top row), so I focus further discussion on the one-sided p-values. As Loadman and
287 McCulloch [9] note in their commentary on [1], high recall (sensitivity) is not achieved without
288 sacrificing precision. For all trials (Fig 4B, second column), precision is moderate, with the
289 maximum slightly in excess of 0.5. Precision is also low at the highest p-values, suggesting
290 that, while the retracted p-values are shifted towards 1, there are still unretracted trials that
291 have high p-values as well. I also identified variability in the performance of the trial p-values
292 across journals. For example, trials published in the journal Anesthesia and Analgesia (Fig
293 4C) had precision near 1 for the highest p-values thresholds, while trials published in the
294 Journal of the American Medical Association (Fig 4D) had poorer performance compared
295 to the aggregate of all trials. The vast majority of trials are not retracted, which means
296 that related to classification such as precision and recall can sometimes be based on small

297 numbers of trials within particular categories. As a result, the results I present here should
298 be considered with caution. Never the less, I beleive that these analyses raise important
299 issues with regard to the applicability of the CM. The variability in classification properties
300 in different journals, along with the observed differences between fabrication previously
301 identified by Carlisle vs new instances of fabrication, raises important questions about the
302 generalizability of the CM, an issue which I dicuss more in depth below.

303 **Discussion**

304 **Implications for global error rates**

305 I first address the implications of the results I present here for the issue of global error rates.
306 Readers of [1] may be concerned by the results presented there if they interpret the analysis
307 to suggest that fraud or error are rampant in the literature, a possiblity that has already
308 been aluded to by some observers [9,19,20]. The analysis that I have presented here indicates
309 that the analysis by Carlisle [1] is not informative of the rate of data validity issues (either
310 fabrication or error) within the literature. The pattern observed by Carlisle [1] in the global
311 distribution of trial-level p-values can plausibly arise for benign reasons. When considering
312 only a single p-values per trial, which avoids some of the problematic assumptions made in
313 [1], the p-value distribution does not display the pattern that Carlisle identifies as potentially
314 indicative of error, suggesting that this critique is not simply speculation.

315 **Implications for use of the CM for screening**

316 The theoretical arguments I give also have implications for the use of the CM in screening.
317 In particular, my theoretical results suggest that screening should not rely on probability
318 statements based on the CM. For example, in [1], Carlisle sometimes thresholds the trial-level

319 p-values (e.g. $p < 1/10000$). If the p-values produced by the CM were valid p-values, then it
320 would be tempting to make statements like “ $p < 1/10000$ would only happen once in 10,000
321 trials, if the trials were truly randomized”. My results suggest that these types of statements
322 are not valid. For instance, using simulated p-values from correlated variables, 0.0024% of
323 p-values are less than $1e-04$, a 24 fold increase over the nominal rate.

324 If certain trials are particularly susceptible to these types of issues (e.g. a subset of studies
325 with highly correlated variables, extreme rounding, strong stratification) this inflation could
326 be exacerbated, without necessarily being obvious to the user of the CM. Likewise, if multiple
327 CM p-values are used together, as Carlisle [1] suggests could be done using multiple trials
328 from the same author, the inflation of error rates compared to their nominal values could
329 be further increased. For example, using correlated p-values, a single p-value of 0.01 has
330 a p-value under correlation of 0.0291, a 2.91 fold inflation, while for two p-values of 0.01
331 combined by multiplication, the inflation is 8.47 fold. This suggests that the p-values produced
332 by the CM have a problem in terms of calibration. If users of the CM target a particular
333 confidence level, in the presence of assumption violations the p-values produced by the CM
334 may not necessarily meet their nominal rates. In addition, extreme assumption violation may
335 produce extreme p-values, so using conservative thresholds does not necessarily alleviate this
336 problem.

337 In addition to problems with calibration, my analysis raises issues with the empirical
338 performance of the CM in terms of its ability to classify known instances of error or misconduct.
339 A global analysis, aggregating across types of retractions and across journals, indicates that
340 when applying the CM there is a strong trade-off between precision and recall (sensitivity)
341 as others [9] have speculated. This suggests that acceptable precision will result in low
342 recall, suggesting that screening initiatives that utilize the CM may not result in significant
343 proportions of errors being identified. If the CM generally identifies few errors, its benefits as
344 a screening tool may be modest. In addition, even the optimal precision achieved by the CM
345 in the full sample of trials is moderate. Retraction are rare, so a moderate precision does

346 not imply that the CM is uninformative. However, there are important implications for the
347 use of CM for screening. First, moderate precision warrants caution in the interpretation of
348 results from the CM. Users should be aware that even at “conservative” thresholds, many
349 of the flagged trials may not be erroneous or fraudulent. Second, parties that may consider
350 using the CM for screening may consider false positives to be associated with increased costs
351 of using the method, such as increased effort need to evaluate flagged trials or the potential
352 of delaying publication of valid research over a false positive.

353 My analysis also reveals heterogeneity in the performance of the CM across categories of
354 retractions and across journals. I discuss three possible explanations for this, all plausible.
355 First, It may be that this heterogeneity arises from heterogeneity in the behavior of researchers
356 who submit to different journals. If the processes by which error or fabrication occur tend
357 to be different across the different journals and retraction categories, this could explain the
358 observed heterogeneity.

359 Second, heterogeneity in precision-recall curves may arise due to issues in the detection
360 of errors. Not all erroneous publications are retracted, and it may be that retractions are
361 generally a conservative marker of error, such that there are many potentially erroneous
362 trials within the literature that could go un-retracted. If this is, then the precision and
363 recall rates calculated based on retraction could give a pessimistic picture of the CM. This is
364 particularly the case if several of the un-retracted trials that have extreme CM p-values are
365 erroneous but undetected, which might be expected if the CM is an effective measure of error.
366 Assuming this is the case, precision and recall using retraction as a metric may underestimate
367 the values that would be obtained using the unseen labels of true error. Under this model,
368 heterogeneity in precision and recall performance is really due to heterogeneity in the extent
369 to which retraction detects error. This suggests the possibility that the journals where the
370 CM performs well are more indicative of the true performance of the CM, while the journals
371 where it performs poorly simply underestimate performance because the trials with extreme
372 p-values that remain un-retracted truly are erroneous, but simply have not been detected as

373 such. Deeper looks at trials that produce extreme CM p-values are warranted to assess this
374 possibility.

375 Finally, it is possible the performance of the CM is overestimated in in the journals where
376 it performs best. The journals where the CM performs well tend to be anesthesia journals
377 that also contain retractions that may have been known to Carlisle during the development
378 of the CM, and in some cases the journals contain retractions that occurred directly as a
379 result of being identified by the CM. Retractions in this category show more extreme CM
380 p-values compared to other fabrication-related retractions (Fig 4A and E). This suggests
381 the possibility that the CM may be “overfit” to these particular trials. If the CM was used
382 to identify some retracte trials, it may be that erroroneous trials that have extreme CM
383 p-values were more likely to be identified, while erroroneous trials that have less extreme CM
384 p-values received less scrutiny within this sample, and therefore remain un-retracted. This
385 may result in inflated recall values, due to the existence of un-retracted trials with moderate
386 CM p-values.

387 This third possibility may work synergistically with the first. For instance, if by chance the
388 anesthesiology field happened to have several prominent examples of fabrication that display
389 the property targeted in the CM, then it is possible that methods similar to the CM were
390 more likely to emerge within this field. As a result, tests of these methods might be more
391 likely to include anesthesia trials from this period, which happened to have an excess of
392 trials displaying these properties, thus resulting in overestimation of the performance of these
393 methods.

394 Taken together, this analysis suggests that caution is warranted if the CM is to be used for
395 screening. Notably, some of the assumptions made my Carlisle are nessesitated by the fact
396 that the analysis in [1] by nessesity was based on summary statistics. This significantly
397 complicates the application for the CM. For example, addressing correlation among baseline
398 variables would be very difficult using only summary statistics, but could be facilitated

399 by analysis of the raw data. If journals choose to implement screening procedures prior
400 to publication, authors of papers would be able to respond to the results of the CM by
401 reanalyzing the raw data, thus potentially giving more definitive answers as to the validity of
402 certain assumptions made by the CM. Critically, in order for this strategy to function well,
403 authors, reviewers, and editors need to have a solid understanding of how various assumptions
404 can impact the results of the CM. This paper can serve as a starting point for members of the
405 scientific community who need to interpret results in this context. Likewise, when institutions
406 such as funders or journals consider whether or how the CM could play a role in decision
407 making, the results presented here can give insight into the possible costs and benefits of
408 various implementations.

409 **Methods**

410 **P-value calculations**

411 I recalculated p-values from summary statistics using the “anovaSummarized” command
412 in the package “CarletonStats” [24], as well as a custom Monte-Carlo method used by
413 Carlisle in [1]. To replicate the method used by Carlisle in my own simulations, I modified
414 code provided by Carlisle in the comments at ([http://steamtraen.blogspot.com/2017/06/
415 exploring-john-carlises-bombshell.html](http://steamtraen.blogspot.com/2017/06/exploring-john-carlises-bombshell.html)). I used the “metap” [25] package for combining
416 p-values by Stouffer’s method [23] using the “sumz” function.

417 **Carlisle data**

418 I obtained Data from the supplemental tables [1], and loaded the data into R using the
419 “readxl” [26] package.

420 **Precision/recall analysis**

421 I computed precision/recall curves using the package “PRROC” [27], with the CM p-value
422 was the metric and a binary indicator of retraction (1 = retracted, 0 otherwise) as the target.

423 **Categorization of retractions**

424 Table S1 of [1] contains notes by Carlisle with information about individual trials, including
425 details of retractions. I categorize the trials by detecting the presence of certain terms or
426 word stems within these notes. I categorize a trial as having been retracted by the presence
427 of “RETRACTED” within these notes, and un-retracted otherwise. I categorize a retraction
428 as coming from a prominent anesthesia related author that is known for fabrication based
429 on the presence of one of four names in the notes. I use the names “Sato”, “Boldt”, “Fuji”
430 and “Reuben”. I categorize a retracted trial that lacks one of these names as having been
431 caused by fabrication based on the presence of “fabricat” in the notes. All other retracted
432 trial I categorize as having occurred for reasons other than fabrication. For one trial, manual
433 review suggested that the note mentions “fabricat” but without definitively attributing the
434 trial to fabrication. As a result, I label this trial as having occurred for reasons other than
435 fabrication.

436 **Computational analysis**

437 I conducted all analyses in R [28] version 3.3.2. I used the package “ggplot2” [29] for
438 visualization.

439 **Reproducibility and computational details**

440 Code used to generate the analyses and figures included in this article are available at
441 https://github.com/ScottWPiraino/carlisle_reanalysis.

442 References

- 443 1. Carlisle JB. Data fabrication and other reasons for non-random sampling in 5087 ran-
444 domised, controlled trials in anaesthetic and general medical journals. *Anaesthesia*. 2017;
445 doi:10.1111/anae.13938
- 446 2. Ioannidis JPA. Why Most Published Research Findings Are False. *PLoS Medicine*. 2005;2:
447 e124. doi:10.1371/journal.pmed.0020124
- 448 3. Simmons JP, Nelson LD, Simonsohn U. False-positive psychology: undisclosed flexibility in
449 data collection and analysis allows presenting anything as significant. *Psychological science*.
450 2011;22: 1359–66. doi:10.1177/0956797611417632
- 451 4. Open Science Collaboration. Estimating the reproducibility of psychological science.
452 *Science*. 2015;349: aac4716–aac4716. doi:10.1126/science.aac4716
- 453 5. Chang AC, Li P, Hanson TJ, Larsson E, Mai KT, Marozzi A, et al. Is Economics Research
454 Replicable? Sixty Published Papers from Thirteen Journals Say "Usually Not". *Finance*
455 *and Economics Discussion Series 2015-083* Washington: Board of Governors of the Federal
456 Reserve System. 2015; doi:10.17016/FEDS.2015.083
- 457 6. Nuijten MB, Hartgerink CHJ, Assen MALM van, Epskamp S, Wicherts JM. The prevalence
458 of statistical reporting errors in psychology (1985-2013). *Behavior Research Methods*. 2016;48:
459 1205–1226. doi:10.3758/s13428-015-0664-2
- 460 7. Brown NJL, Heathers JAJ. *The GRIM Test*. *Social Psychological and Personal-*
461 *ity Science*. SAGE PublicationsSage CA: Los Angeles, CA; 2016; 194855061667387.
462 doi:10.1177/1948550616673876
- 463 8. Anaya J. The GRIMMER test: A method for testing the validity of reported measures of
464 variability. *PeerJ Inc*. 2016; doi:10.7287/peerj.preprints.2400v1
- 465 9. Loadsman JA, McCulloch TJ. Widening the search for suspect data - is the flood of

- 466 retractions about to become a tsunami? *Anaesthesia*. 2017; doi:10.1111/anae.13962
- 467 10. Al-Marzouki S, Evans S, Marshall T, Roberts I. Are these data real? Statistical methods
468 for the detection of data fabrication in clinical trials. *BMJ (Clinical research ed)*. 2005;331:
469 267–70. doi:10.1136/bmj.331.7511.267
- 470 11. Carlisle JB. A meta-analysis of prevention of postoperative nausea and vomiting: ran-
471 domised controlled trials by Fujii et al. compared with other authors. *Anaesthesia*. 2012;67:
472 1076–1090. doi:10.1111/j.1365-2044.2012.07232.x
- 473 12. Simonsohn U. Just post it: the lesson from two cases of fabricated data detected by
474 statistics alone. *Psychological science*. 2013;24: 1875–88. doi:10.1177/0956797613480366
- 475 13. Carlisle JB, Dexter F, Pandit JJ, Shafer SL, Yentis SM. Calculating the probability of
476 random sampling for continuous variables in submitted or published randomised controlled
477 trials. *Anaesthesia*. 2015;70: 848–858. doi:10.1111/anae.13126
- 478 14. Bolland MJ, Avenell A, Gamble GD, Grey A. Systematic review and statistical anal-
479 ysis of the integrity of 33 randomized controlled trials. *Neurology*. 2016;87: 2391–2402.
480 doi:10.1212/WNL.0000000000003387
- 481 15. Carlisle JB, Loadman JA. Evidence for non-random sampling in randomised, controlled
482 trials by Yuhji Saitoh. *Anaesthesia*. 2017;72: 17–27. doi:10.1111/anae.13650
- 483 16. Baker M. Stat-checking software stirs up psychology. *Nature*. 2016;540: 151–152.
484 doi:10.1038/540151a
- 485 17. Buranyi S. The hi-tech war on science fraud [Internet]. 2017. Available: [https://www.](https://www.theguardian.com/science/2017/feb/01/high-tech-war-on-science)
486 [theguardian.com/science/2017/feb/01/high-tech-war-on-science](https://www.theguardian.com/science/2017/feb/01/high-tech-war-on-science)
- 487 18. Schmidt T. Sources of false positives and false negatives in the STATCHECK algorithm:
488 Reply to Nuijten et al. (2016). 2016; Available: <http://arxiv.org/abs/1610.01010>
- 489 19. Oransky I. Tracking retractions as a window into the scientific process Two in 100 clinical
490 trials in eight major journals likely contain inaccurate data: Study [Internet]. 2017. Available:

- 491 <http://retractionwatch.com/2017/06/05/two-100-clinical-trials-eight-major-journals-likely-contain-inaccura>
- 492 20. Buranyi S, Devlin H. Dozens of recent clinical trials may contain wrong or falsified data,
493 claims study. 2017.
- 494 21. Casadevall A, Steen RG, Fang FC. Sources of error in the retracted scientific literature.
495 The FASEB Journal. 2014;28: 3847–3855. doi:10.1096/fj.14-256735
- 496 22. Fanelli D, Kerridge I, Hill S, McNeill P, Doran E. How Many Scientists Fabricate
497 and Falsify Research? A Systematic Review and Meta-Analysis of Survey Data. Tregenza
498 T, editor. PLoS ONE. Southern Illinois University. Doctoral dissertation; 2009;4: e5738.
499 doi:10.1371/journal.pone.0005738
- 500 23. Darlington RB, Hayes AF, Darlington R. Combining Independent p Values: Extensions
501 of the Stouffer and Binomial Methods. Psychological Methods. 2000;5: 496–515. Available:
502 <https://pdfs.semanticscholar.org/5f0b/8e577027d1e1a411c4da9e2fd5a9892895e1.pdf>
- 503 24. Chihara L. CarletonStats: Functions for Statistics Classes at Carleton College. 2016.
- 504 25. Dewey M. Metap: Meta-analysis of significance values. 2017.
- 505 26. Wickham H, Bryan J. readxl: Read Excel Files. 2017.
- 506 27. Keilwagen J, Grosse I, Grau J, Haldemann B, Posch S. Area under Precision-Recall
507 Curves for Weighted and Unweighted Data. Chen Z, editor. PLoS ONE. Public Library of
508 Science; 2014;9: e92209. doi:10.1371/journal.pone.0092209
- 509 28. R Core Team. R: A Language and Environment for Statistical Computing. 2016.
- 510 29. Wickham H. ggplot2: Elegant Graphics for Data Analysis. 2009.