

1 **Title:** fragilityindex: An R Package for Statistical Fragility Estimates in Biomedicine

2

3 **Authors:** Kipp W. Johnson^{1,2}, Eli Rappaport^{1,2}, Khader Shameer^{1,2}, Benjamin S.

4 Glicksberg^{1,2}, Joel T. Dudley^{1,2*}

5

6 **Affiliations:**

7 1. Institute for Next Generation Healthcare, Mount Sinai Health System, New York,

8 NY

9 2. Department of Genetics and Genomic Sciences, Icahn Institute for Genomics and

10 Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, NY

11

12 *Corresponding author: joel.dudley@mssm.edu

13

14

15

16

17

18

19

20

21

22

23

24 **Abstract**

25 Poor reproducibility is a growing crisis in biomedical research. The fragility index
26 was introduced as a convenient measure to estimate how fragile statistical results in
27 clinical trials are to small perturbations in event outcome counts. There is currently
28 no freely available R package to produce this calculation. Furthermore, the original
29 definition of the method is applicable only to 2x2 contingency tables.

30 As such, we developed an R package to calculate fragility index. We have also
31 extended the concept of a statistical fragility index to two of the most commonly
32 used methods in clinical research, survival analysis via weighted log-rank tests and
33 logistic regression, and implemented these technique sin this R package. We
34 describe example applications of these methods to existing publically available
35 datasets. This R package is freely available under the AGPL license on CRAN
36 (<https://cran.r-project.org/web/packages/fragilityindex/index.html>). The most
37 recent versions may be downloaded and installed via Github
38 (<https://github.com/kippjohnson/fragilityindex>).

39

40

41

42

43

44

45

46

47

48 **Introduction**

49 Awareness of the problem of poor reproducibility in biomedical science has grown
50 in the past several years. There have been numerous attempts to explain this
51 phenomenon, such as an over-reliance on null-hypothesis significance testing (1, 2),
52 financial conflicts of interest, or simply prevailing biases in hypotheses tested (3). In
53 2014, Walsh et al. showed that the statistical significance of randomized clinical
54 trials in medical research are often fragile, or sensitive, to small perturbations in
55 patient outcomes with a simple statistic called the fragility index (4). The fragility
56 index was originally defined for use in 2x2 contingency tables with dichotomous
57 outcomes. In the case of clinical trials, fragility index is defined as the minimum
58 number of patients it would take to make a statistically significant result non-
59 significant if the patients' outcomes were reversed. Despite being quite simple, this
60 is a relevant statistic that can facilitate better interpretation of results from clinical
61 trials. For example, Ridgeon et al. found that 40% of trials in critical care medicine
62 had a fragility index of 1 or less (5), which is of large concern in the medical field.

63 This statistic has since been embraced by physicians in a number of clinical
64 fields as an easy-to-interpret measure of the robustness of clinical trials (5, 6, 7).
65 However, the adoption of this metric in routine statistical analyses has been limited
66 due to the unavailability of a user-friendly statistical package. Here, we present an R
67 package (8) "fragilityindex" which can perform fragility index estimations and has
68 tools which extend the concept of the fragility index to two other statistical

69 techniques often employed in medical research: survival analysis with weighted log-
70 rank tests and logistic regression.

71

72 **Materials and Methods**

73 *Fragility Index for contingency tables*

74 Fragility index calculation is most often associated with contingency tables and
75 there are many studies that portray its utility. For example, consider the following
76 scenario: in a clinical trial, there are two groups of patients. In group one, 15/40
77 patients have an adverse event. Within group two, 6/40 patients have the adverse
78 event. The corresponding exact test P-value for this situation is 0.041. However, if a
79 single additional patient in group two experienced the adverse event, we obtain a
80 new p-value of 0.078. This result is no longer reportable as “significant,” and the
81 fragility index for this trial is 1. We consider this to be a fragile clinical trial result,
82 since a single swapped outcome would change the conclusion of the study. In our
83 package, we can conduct this test as follows:

84

```
> fragility.index(15, 6, 40, 40, verbose=TRUE)
```

```
fragility.index p.value
```

```
85 1      0 0.041  
   2      1 0.078
```

86

87 To further illustrate the utility of our tool, we applied this analysis to the results
88 from all clinical studies publically available on ClinicalTrials.gov (<https://www.clinicaltrials.org/aact-database>, March 27th, 2016 accession). Before the data can

90 be read into R, “line feed” characters within variables need to be removed from the
91 datasets, which can be accomplished with the following line of bash code (taken
92 from the README file included within the data):

93

```
94 > tr -d '\012' < clinical_study_noclob.txt > clinical_study_noclob_nolf.txt
```

95

96 As many of these studies are not in a useable format for the current analysis, we
97 filtered them based on criteria described in Figure 1. Briefly, we only considered
98 clinical trials where there were significant results (defined as $p < 0.05$) from Fisher
99 Exact tests performed on 2x2 tables. There are several issues with formatting in the
100 dataset that must be adjusted before the resulting values can be used in our tests.
101 Furthermore, some entries have impossible results (i.e. non-count values) where
102 perhaps the statistical test was mis-specified on ClinicalTrials.gov. Due to limitations
103 in the dataset, aggregate results should thus be considered with these limitations in
104 mind. In the manuscript, we highlight the application of this function to a particular
105 clinical trial whose values we manually verified. We provide the complete R code to
106 generate the final fragility indices from the clinical data in the Supplemental
107 Materials.

108

109 **Figure 1:** Workflow of filtering steps for available clinical trial data used in fragility
110 index analysis for contingency tables. Data was acquired from ClinicalTrials.gov,
111 specifically the publically available Aggregate Analysis of ClinicalTrials.gov (AACT)
112 database.

113

114

115 *Fragility Index for survival analysis*

116 To increase power, many clinical trials compare the outcomes of two groups studied
117 longitudinally with survival analysis. We have developed a new technique to enable
118 fragility index calculation in survival analysis utilizing the *survdiff()* function from
119 the R package “*survival*” (9). A survival time series consists of paired outcomes (0/1)
120 and exposure times for patients. To compute this fragility index analogously to the
121 fragility index for 2x2 tables, we add a new individual to the input dataset with an
122 exposure time equal to the average exposure time of all previous individuals. We
123 indicate this individual died at the average exposure time. The individual’s covariate
124 values, if applicable for the *survdiff()* function, are randomly sampled from another
125 individual in the dataset. We then repeat the process for multiple iterations and
126 average the outcome of each run into an aggregate, overall fragility index estimate
127 for the dataset. This fragility index may thus be interpreted as the mean number of
128 individuals it would take who died at the average exposure time of all individuals in
129 the study to take a result from significant to non-significant.

130 To show how fragility index calculation can be extended to survival analyses,
131 we performed our method on open-source TCGA, using a Colorectal
132 Adenocarcinoma dataset with 276 patients (10). This dataset is comprised of
133 patient survival rates along with their genetic profiles, particularly which genes
134 have copy number alterations (CNA; amplifications or deletions). We then selected
135 two genes, BCL2L1 and POFUT1, in which there were the most patients (35/276)
136 with amplifications in both. The downloading and pre-processing steps are
137 described fully in the supplementary material. Briefly, the TCGA dataset for

138 colorectal adenocarcinoma (“coadread”) was downloaded from
139 http://www.cbioportal.org/study?id=coadread_tcga_pub#summary using the GUI
140 option to “download data” on November 7th, 2016. From this file, we use the table
141 entitled “data_clinical.txt” for survival information. For the CNA patient survival
142 data, we used the GUI on the same website to check off genes BCL2L1 and POFUT1
143 and then used the GUI option to download this data, saving the file as
144 “CNA_genes.txt”. We use patient identifiers from this file in the resulting analysis.
145 Calculation of the survival analysis fragility index from these datasets can be
146 performed with the R code located in the Supplementary Materials, using the R
147 package “*survival*” for weighted log-rank tests.

148

149 *Fragility Index for Logistic Regression β -coefficients*

150

151 We present a method to calculate logistic regression coefficient fragility, or how
152 many events it would take to change a significant logistic regression coefficient to
153 non-significant at the given confidence level. To do this, we replace responses
154 (which should be binary events, 0/1) with the opposite event until the coefficient is
155 non-significant. If the initial regression coefficient (β) is positive, we change a 1
156 event to a 0. If the regression coefficient is negative, we change a 0 event to a 1. We
157 then count the number of times a replacement must occur to produce a non-
158 significant result and use this amount as the fragility index. To account for
159 variability in the response replacement, we then repeat this process a number of
160 times and take the mean of all of the computed fragility indices as the representative

161 fragility index. As with all multiple regression models, the issue of potential
162 (multi)colinearity, or high correlation between variables, must be addressed before
163 running the analysis. This is especially important as multicollinear variables may
164 have a significant impact upon the “fragility” of the regression coefficients by
165 making them unstable. Using existing packages, we must include this step as a
166 check before running the fragility index for regression coefficients in order to most
167 accurately interpret results.

168

169 **Results**

170 *Fragility Index for contingency tables*

171 We calculated the fragility index for a curated list of all available clinical trial studies
172 from [clinicaltrials.gov](https://www.clinicaltrials.gov) (<https://www.ctti-clinicaltrials.org/aact-database>, March
173 27th, 2016 accession). Of 3,083 Fisher’s exact test results, we excluded results that
174 were not in a 2x2 format (n=19), were not statistically significant ($p \geq 0.05$, n= 2004),
175 or had improper values (n=37), leaving 1,023 usable tests from 136 unique clinical
176 trials (Figure 2). We randomly selected one test from each trial to calculate fragility
177 index. We found five tests from five separate trails that have a fragility index of 1.
178 For example, in a clinical trial (NCT00058019) evaluating the use of Ixabepilone (a
179 chemotherapeutic agent) for non-Hodgkin’s lymphoma, one tested outcome was the
180 response rate between chemosensitive and chemoresistant individuals. They found
181 that Ixabepilone administered to chemosensitive individuals (14/39) had a better
182 response rate ($p=0.022$) than chemoresistant individuals (0/12). We can apply the
183 fragility index assessment on these outcomes:

184

```
> fragility.index(15, 6, 40, 40, verbose=TRUE)
```

```
fragility.index p.value
```

```
185 1      0 0.041
```

```
2      1 0.078
```

186

187 **Figure 2:** The distribution of fragility indices for the 136 tests used as
188 representative examples from 136 distinct clinical trials.

189

190

191 Accordingly, a single addition in response outcome (i.e. a fragility index of 1)
192 between the two groups makes the association not reportable as statistically
193 significant ($p=0.083$). Reproducing or replicating the study to determine whether
194 this effect was due simply to chance may be recommended. We provide full data
195 acquisition steps and the fragility index distribution for all usable clinical trials in
196 the Supplemental Materials.

197

198 *Fragility Index for Survival Analysis*

199 We have applied this analysis to a Colorectal Adenocarcinoma dataset from TCGA
200 (10) which assesses the effect of gene Copy Number Alteration on survival.

201 Specifically, we calculated the fragility index for the significant association of
202 amplifications of genes *BCL2L1* and *POFUT1* and survival ($p=0.03$) as follows
203 (detailed steps to replicate analysis in Supplemental Materials):

204

```
205 survivalfragility(Surv(months, status) ~ cna.status, data=infile,  
206                   niter=100, progress.bar=TRUE)  
207 |+++++| 100% elapsed = 06s  
208 [1] 1.36
```

209

210 Although the survival analysis is significant for amplification of these 2 genes, we
211 find that it is relatively fragile as it is non significant by the addition of only 1-2
212 individuals who die at the mean exposure time.

213

214 *Fragility Index for Logistic Regression β -coefficients*

215 To illustrate the utility of calculating a fragility index for logistic regression, we
216 utilized data from UCI Machine Learning Repository (11), specifically a multivariate
217 heart disease dataset (Cleveland site). This dataset is comprised of an outcome
218 variable (heart disease status) and 13 predictor variables (full description can be
219 found on their website). We provide the full R code to reproduce this example
220 analysis in the Supplementary Materials.

221

```
222  
223 # before running logisticfragility function, first test for (multi)colinearity  
224 (see Supplementary Materials)  
225  
226 > logisticfragility(num ~ age + sex + cp + trestbps + chol + fbs + restecg +  
227 thalach + exang + oldpeak + slope + ca + thal, data = mydata, covariate="all",  
228 niter=100, progress.bar=FALSE)  
229  
230 # only significant associations displayed  
231 coefficient fragility.index  
232 1 (Intercept)          17.21  
233 3 sex                  29.81  
234 4 cp                   35.73  
235 5 trestbps              8.06  
236 9 thalach               4.99  
237 10 exang                9.37
```

238 13 ca 69.69
239 14 thal 61.83
240

241

242 We find that the “ca” (# of major vessels colored by flourosopy) and “thal” (defect
243 type) coefficients are the least fragile as heart disease predictors, while “thalach”
244 (maximum heart rate achieved), “trestbps” (resting blood pressure), and “exang”
245 (presence of exercised induced angina) are the most fragile. As previously indicated
246 (6, 12), we do find a positive correlation between degree of significance and fragility
247 (i.e. higher significance generally leads to more stable fragility scores) and, as such,
248 both metrics should be considered when interpreting outcome associations.

249

250 **Discussion**

251 We have developed the first R statistical package for the calculation of fragility index
252 in multiple common settings, having extended and incorporated this technique to
253 other commonly used statistical tests in the medical literature. Our extensions of the
254 fragility index to survival analysis and logistic regression may prove to be helpful
255 tools for researchers attempting to analyze clinical results. We hope that increased
256 awareness of the fragility of some results will lead re-searchers to be more
257 cognizant of the fragility of their claims. Ultimately, we hope this process will
258 contribute to enhancing the reproducibility of the biomedical literature.

259

260 **Acknowledgements**

261 We would like to thank Institute of Next Generation Healthcare, Mount Sinai Health
262 System, New York, NY for infrastructure and support. We would also like to thank
263 the authors of the open-source datasets used as examples for this study.

264

265 **Funding**

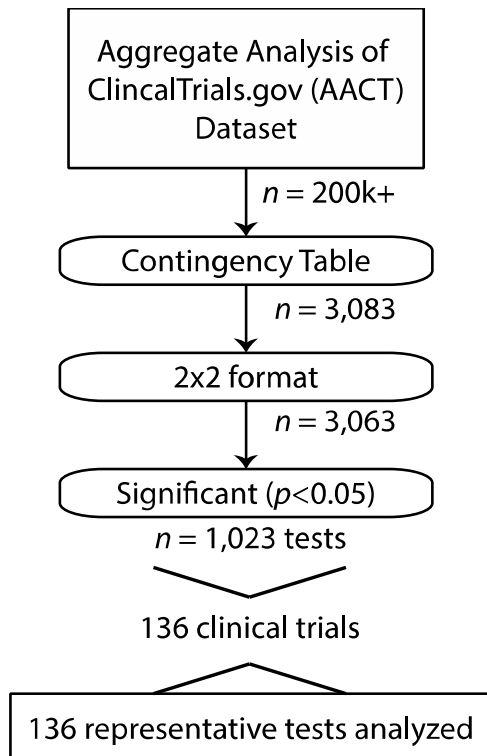
266 This scholar work is supported by the following grants from National Institutes of
267 Health: National Institute of Diabetes and Digestive and Kidney Diseases
268 (R01DK098242); National Cancer Institute (U54CA189201); Illuminating the
269 Druggable Genome; Knowledge Management Center sponsored by National Insti-
270 tutes of Health Common Fund; National Cancer Institute (U54-CA189201-02);
271 National Center for Advancing Translational Sciences and Clinical and Translational
272 Science Award (UL1TR000067)¹¹

273

274 Conflict of Interest: KWJ, BSG, and KS: none declared. JTD: Dudley has received
275 consulting fees or honoraria from Janssen Pharmaceuticals, GlaxoSmithKline,
276 AstraZeneca, and Hoffman-La Roche; is a scientific advisor to LAM Therapeu-tics;
277 and holds equity in NuMedii Inc., Ayasdi Inc., and Ontomics, Inc.

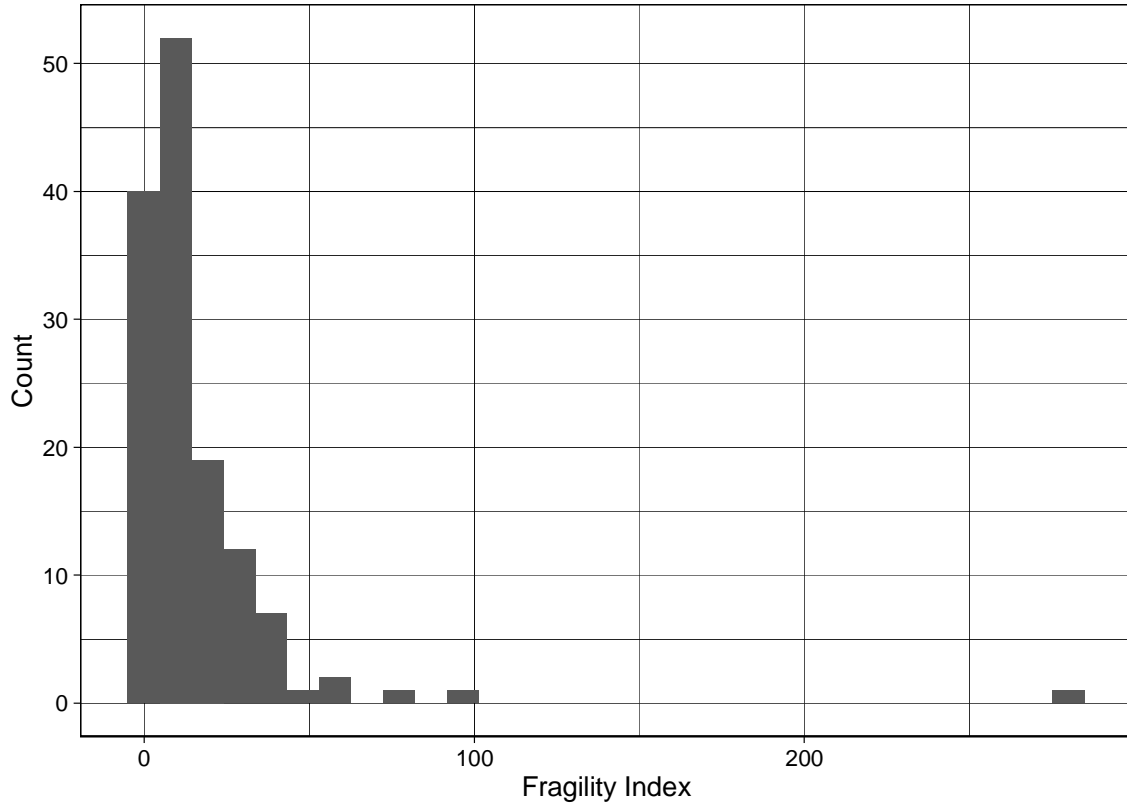
278 **Figure 1**

279



280

281 **Figure 2**
282



283

284 References

- 1 Leggett NC, Thomas NA, Loetscher T, Nicholls ME. The life of p: "just significant" results are on the rise. *Q J Exp Psychol (Hove)* 2013;66: 2303-2309.
- 2 Glaser DN. The controversy of significance testing: misconceptions and alternatives. *Am J Crit Care* 1999;8: 291-296.
- 3 Ioannidis JP. Why most published research findings are false. *PLoS Med* 2005;2:e124.
- 4 Walsh M, Srinathan SK, McAuley DF, Mrkobrada M, Levine O, Ribic C, et al. The statistical significance of randomized controlled trial results is frequently fragile: a case for a Fragility Index. *J Clin Epidemiol* 2014;67: 622-628.
- 5 Ridgeon EE, Young PJ, Bellomo R, Mucchetti M, Lembo R, Landoni G. The Fragility Index in Multicenter Randomized Controlled Critical Care Trials. *Crit Care Med* 2016;44: 1278-1284.
- 6 Docherty KF, Campbell RT, Jhund PS, Petrie MC, McMurray JJ. How robust are clinical trials in heart failure? *Eur Heart J* 2017;38: 338-345.
- 7 Ahmed W, Fowler RA, McCredie VA. Does Sample Size Matter When Interpreting the Fragility Index? *Crit Care Med* 2016;44: e1142-e1143.
- 8 R Development Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2016.
- 9 Therneau TM and Grambsch PM. *Modeling survival data: Extending the Cox Model*. New York: Springer, 2000. ISBN: 0-387-98784-3.
- 10 Cancer Genome Atlas N. Comprehensive molecular characterization of human colon and rectal cancer. *Nature* 2012;487: 330-337.
- 11 Lichman M. *UCI Machine Learning Repository*. Irvine, CA, 2013.
- 12 Carter RE, McKie PM, Storlie CB. The Fragility Index: a P value in sheep's clothing? *European Heart Journal* 2017; 38: 346-348.

Aggregate Analysis of
ClinicalTrials.gov (AACT)
Dataset

$n = 200k+$

Contingency Table

$n = 3,083$

2x2 format

$n = 3,063$

Significant ($p < 0.05$)

$n = 1,023$ tests

136 clinical trials

136 representative tests analyzed

