

Canons and Sparrows II*: The Enhanced Bernoulli Exact Method for Determining Statistical Significance and Effect Size in the Meta-Analysis of k 2×2 Tables

Lawrence Marc Paul (✉ Impaul@sciencesupportconsulting.com)

Bell Laboratories

Methodology

Keywords: Meta-analysis, Categorical Analysis, Dichotomous Analysis, Sparse Tables, Mantel-Haenszel, DerSimonian, Exact Solution, Inverse Variance, Convolution, Heterogeneity, Rare Events

Posted Date: January 7th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-139437/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Emerging Themes in Epidemiology on August 3rd, 2021. See the published version at <https://doi.org/10.1186/s12982-021-00101-8>.

1 **Canons and Sparrows II^{*} : The Enhanced Bernoulli Exact Method for**
2 **Determining Statistical Significance and Effect Size**
3 **in the Meta-Analysis of k 2 x 2 Tables**

4

5

6 **Lawrence M. Paul**

7

8

9

10

11

12

13

14

15

16

17 Bell Laboratories, Retired

18

19 E-mail: lpaul@sciencesupportconsulting.com

20

21

22

23 **Running head:** Non-parametric exact tests for meta-analysis of categorical data

24 * *“Little experience is sufficient to show that the traditional machinery of statistical*
25 *processes is wholly unsuited to the needs of practical research. Not only does it take a*
26 *cannon to shoot a sparrow, but it misses the sparrow. The elaborate mechanism built on*
27 *the theory of infinitely large samples is not accurate enough for simple laboratory data.”*
28 *(R. A. Fisher, 1925)*

29

30 **ABSTRACT**

31

32 **Background:** The use of meta-analysis to aggregate the results of multiple studies has
33 increased dramatically over the last 40 years. For homogeneous meta-analysis, the
34 Mantel-Haenszel technique has typically been utilized. In such meta-analyses, the effect
35 size across the contributing studies of the meta-analysis differ only by statistical error. If
36 homogeneity cannot be assumed or established, the most popular technique developed to
37 date is the inverse-variance DerSimonian & Laird (DL) technique [1]. However, both of
38 these techniques are based on large sample, asymptotic assumptions. At best, they are
39 approximations especially when the number of cases observed in any cell of the
40 corresponding contingency tables is small.

41 **Results:** This research develops an exact, non-parametric test for evaluating statistical
42 significance and a related method for estimating effect size in the meta-analysis of $k \times 2 \times 2$
43 tables for any level of heterogeneity as an alternative to the asymptotic techniques.
44 Monte Carlo simulations show that even for large values of heterogeneity, the Enhanced

45 Bernoulli Technique (EBT) is far superior at maintaining the pre-specified level of Type I
46 Error than the DL technique. A fully tested implementation in the R statistical language
47 is freely available from the author. In addition, a second related exact test for estimating
48 the Effect Size was developed and is also freely available.

49 **Conclusions:** This research has developed two exact tests for the meta-analysis of
50 dichotomous, categorical data. The EBT technique was strongly superior to the DL
51 technique in maintaining a pre-specified level of Type I Error even at extremely high
52 levels of heterogeneity. As shown, the DL technique demonstrated many large violations
53 of this level. Given the various biases towards finding statistical significance prevalent in
54 epidemiology today, a strong focus on maintaining a pre-specified level of Type I Error
55 would seem critical.

56 **Keywords:** Meta-analysis; Categorical Analysis; Dichotomous Analysis; Sparse Tables;
57 Mantel-Haenszel; DerSimonian; Exact Solution; Inverse Variance, Convolution,
58 Heterogeneity, Rare Events

59

60 **BACKGROUND**

61

62 The use of meta-analysis in epidemiological research has been increasing at a very rapid
63 rate. A review of the National Library of Medicine’s online database (“Pub Med”) shows
64 that in 1977 there was only a single research article with the term “meta-analysis” in its
65 title. This number had increased to 138 in 1991, 966 in 2005 and to 17,205 in 2019 (see
66 Figure 1).

67

68

Insert About Here

69 **Fig. 1 Number of articles containing "meta-analysis" in the title by year of**
70 **publication.**

71

72 Part of this growth may be due to the widespread availability of powerful personal
73 computer software making meta-analysis techniques more feasible to implement. More
74 importantly, the need to draw meaningful conclusions from an aggregation of small
75 studies may help explain this exponential growth.

76 The use of meta-analytic techniques is controversial when the contributing studies are not
77 randomized control trials (RCT). Many researchers feel that it is highly misleading to
78 attempt to combine a series of disparate studies [2] while others maintain that, with
79 proper safeguards, meta-analysis allows an extremely useful pooling of smaller studies
80 [3], [4]. A discussion of the appropriateness of meta-analysis is beyond the scope of this
81 paper. Rather, the focus here will be on minimizing unnecessary error in testing the
82 overall statistical significance of a meta-analysis and in estimating the Effect Size.

83 **Overview of 2 x 2 x k Categorical Meta-Analysis**

84 The "2 x 2 x k" categorical meta-analysis paradigm is probably the most frequently
85 encountered situation in meta-analysis. It consists of a series of k contributing studies
86 each described by a 2 x 2 contingency table. Every cell of each 2 x 2 table contains the
87 number of occurrences of an event (e.g., disease cases) for the particular combination of
88 row and column variables. For the sake of illustration, we can associate the two columns

A non-parametric exact test for meta-analysis of categorical data

89 of each table with Disease Manifestation vs. No Disease Manifestation and the two rows
90 with Exposure vs. No Exposure. Table 1 represents the results of one of these k studies
91

Table 1 Typical contributing study (one of k) in a dichotomous meta-analysis			
	Disease Status		
Exposure	Disease Manifestation	No Disease Manifestation	Total
Exposure	4	96	100
No Exposure	2	98	100
Total	6	194	200

92

93 In most meta-analyses, there are typically two distinct components: 1) A statistical test
94 of the overall difference between the Exposure and No Exposure groups across the k
95 contributing studies and 2) A method to pool the observed differences between groups
96 across the k studies in order to estimate the true difference (the Effect Size).

97 Surprisingly, in recent years, many epidemiologists employing meta-analytic techniques
98 have greatly deemphasized the first component. Borenstein *et al.* [3] conclude:

99 “... However, meta-analysis also allows us to move beyond the question of
100 statistical significance, and address questions that are more interesting and also
101 more relevant.” (pp. 11-12).

102 Similarly, Higgins *et al.* [4] rather dismissively state:

103 “... If review authors decide to present a p value with the results of a meta-
104 analysis, they should report a precise p value, together with the 95% confidence
105 interval” (pp. 371-372).

106 A method is developed that maintains the Type I error (“false alarm rate”) at the desired
107 level, but which has good power to detect true differences across a large range of event
108 probability, number of contributing studies, sample size and level of heterogeneity. An
109 argument can be made that maintaining the Type I error at a pre-specified level is more
110 important than the power ($1 - \text{Type II error rate}$) to detect true differences between
111 conditions. The framers of modern statistical testing called such errors “Errors of the
112 First Kind” and placed a special emphasis on them. Neyman & Pearson in 1933 stated:

113 “A new basis has been introduced for choosing among criteria available for
114 testing any given statistical hypothesis, H_0 , with regard to an alternative H_t . If Θ_1
115 and Θ_2 are two such possible criteria and if in using them there is the same
116 chance, ϵ , of rejecting H_0 when it is in fact true, we should choose that one of the
117 two which assures the minimum chance of accepting H_0 when the true hypothesis
118 is H_t .” [5][p. 336]

119 Thus, while Neyman & Pearson supported the effort to choose criteria that yield the
120 greatest power to detect true differences, this effort is secondary to maintaining a pre-
121 specified level of Type I error. A second exact method is developed to estimate the effect
122 size of any statistically significant finding.

123 “Rare” Events and Meta-Analysis

124 The probability of occurrence of a disease is often categorized as “rare” although no
125 specific definition exists. As an example, Higgins *et al.* state that “There is no single risk
126 at which events are classified as ‘rare’ ” but gives as examples 1 in a 100 or 1 in a 1,000
127 (see [6], p. 520). An obvious related issue is observing zero cases in one or more cells of

A non-parametric exact test for meta-analysis of categorical data

- 128** a contingency table. Table 2 shows the expected cell sizes from various realistic
- 129** combinations of disease probability and contributing study sample size.

Table 2 Expected number of disease cases as a function of disease probability and individual study sample size (each arm)				
Disease/Condition	Approximate Disease Prob.	Individual Study Sample Size (each arm)		
		100	500	1000
Myocardial Infarction Incidence Rate for Age ≥ 60 years	.011 [7]	1.1	5.5	11
Parkinson’s Disease Incidence Rate (60 -69 Age Group)	.00058 [8]	.06	.29	.58
Alzheimer’s Disease Incidence Rate (60 -74 Age Group)	.002 [9]	.2	1	2
Lung Cancer Incidence Rate for White Males	.00051 [10]	.05	.26	.51

130

131 Table 2 supports the notion that “rare” events are a focus of many epidemiological

132 studies.

133 For homogeneous meta-analysis (i.e. where the effect across studies may be assumed to

134 be the same within statistical variation), the two techniques typically used for categorical

135 data are the Mantel-Haenszel and Peto techniques. Both of these techniques rely on the

136 Mantel-Haenszel Chi Square to test for the overall statistical significance. For

137 heterogeneous meta-analyses, the asymptotic DerSimonian-Laird (DL) inverse variance

138 technique is typically used [1].

139 The problem in applying large sample asymptotic techniques to meta-analyses involving

140 small numbers of cases will be illustrated in the older and much more developed domain

141 of homogeneous meta-analyses. Mantel developed what is probably the most widely

142 used technique for homogeneous meta-analyses [11]. In applying his technique, he

143 showed that a minimum of approximately five cases was required in each of the 4 cells of
144 each of the 2 x 2 tables for each of the k studies comprising the meta-analysis [12]. This
145 is the same heuristic requirement typically used without any particular justification for
146 the simple chi-square test. Mantel & Fleiss reviewed the options when a reasonable
147 number of cases was not present in all cells:

148 “The investigators could have obtained data from very many more tables to make
149 things more asymptotic for use of M-H [note: this is the Mantel-Haenszel
150 technique], or they could readily have applied a more exact procedure for the data
151 at hand” [p. 134].

152 R. A. Fisher made essentially the same plea in 1925 in the preface to the first edition of
153 his well-known *Statistical Methods for Research Workers* [13]:

154 “Little experience is sufficient to show that the traditional machinery of statistical
155 processes is wholly unsuited to the needs of practical research. Not only does it
156 take a cannon to shoot a sparrow, but it misses the sparrow. The
157 elaborate mechanism built on the theory of infinitely large samples is not accurate
158 enough for simple laboratory data. Only by systematically tackling small sample
159 problems on their merits does it seem possible to apply accurate tests to practical
160 data.”

161 Both criticisms suggest the use of exact methods to handle the sparseness of the
162 underlying contingency tables at least for the disease examples contained in Table 2. All
163 but two of the combinations of individual study sample size and disease probability
164 shown in Table 2 would yield fewer than five cases per cell leading to violations of the
165 minimum cell size in the Mantel-Haenszel (MH) Chi Square test, and thus the test would
166 be potentially flawed. In addition, these two cases were for sample size equal to 500 and
167 1,000 which may not represent many realistic studies. While this limitation of the MH
168 Chi Square test was known to Mantel and others (e.g. [12]), it seems to generally have
169 been forgotten for meta-analysis of $2 \times 2 \times k$ categorical data. The continued use of an
170 asymptotic test in situations not suited for its use is unacceptable given the computer
171 power that is now available to all researchers.

172 **Heterogeneity vs. Homogeneity in Meta-Analyses**

173 The term “heterogeneity” refers to the fact that studies done at different times and by
174 different researchers might be expected to yield different results. The expectation is that
175 a variable of interest may be dependent, at least in part, to one or more other variables.
176 The meta-analysis researcher, J. P. T. Higgins stated “As Heterogeneity is to be expected
177 in a meta-analysis: it would be surprising if multiple studies, performed by different
178 teams in different places with different methods, all ended up estimating the same
179 underlying parameter.” [[14], p. 158]. While researchers may agree that heterogeneity is
180 to be expected, there is very little agreement on how to quantify this variability. The
181 obvious candidate is τ^2 , the estimated variability between studies. However, τ^2 is not
182 invariant across study designs and its interpretation may not be intuitive. Alternatives

183 include I^2 , the ratio of the inter-study variability to the total variability and the Q statistic,
184 which is mathematically related to I^2 (see, e.g., [15]).

185 In the technique described in this paper, heterogeneity will be mathematically
186 manipulated through τ^2 using the logit distribution as developed by Bhaumik *et al.* [16].

187 Namely:

188
$$x_{Ci} \sim B(p_{Ci}, n_{Ci}), \quad x_{Ei} \sim B(p_{Ei}, n_{Ei}), \quad [1]$$

189
$$\text{logit}(p_{Ci}) = \mu + \varepsilon_{1i}, \quad \text{logit}(p_{Ei}) = \mu + \theta + \varepsilon_{1i} + \varepsilon_{2i} \quad [2]$$

190
$$\varepsilon_{1i} \sim N(0, \gamma^2) \quad [3]$$

191
$$\varepsilon_{2i} \sim N(0, \tau^2) \quad [4]$$

192 Where:

193 B is the Binomial Distribution

194 N is the Normal Distribution

195 X_{Ci} , X_{Ei} are the observed number of cases in the control and exposure groups

196 respectively of the i th study

197 p_{Ci} , p_{Ei} are the event probabilities in the control and exposure groups respectively of

198 the i th study

199 n_{Ci} , n_{Ei} are the sample sizes in the two groups of the i th study

200 μ corresponds to the background event (disease) probability in the exposure and

201 control groups

202 θ corresponds to the overall ratio of the event probabilities for the exposure group

203 relative to the control group

A non-parametric exact test for meta-analysis of categorical data

204 γ^2 is a variance corresponding to the uncertainty of the observed disease probability in
205 both the exposure and control groups of the k contributing studies.
206 τ^2 is a variance corresponding to the heterogeneity or the “heterogeneity parameter” only
207 in the exposure group.
208 ε_{1i} is a Normal distribution deviation in background event (disease) probability for both
209 the exposure and control groups of each of the contributing studies.
210 ε_{2i} is a Normal distribution deviation in background disease probability due to
211 heterogeneity in the exposure group of each of the contributing studies

212

213 **The Basic Principles of the Dersimonian-Laird (DL) Method**

214

215 As stated above, this research specifically contrasts an exact method for conducting meta-
216 analyses in k 2 x 2 tables with heterogeneity with the most popular approach which was
217 developed by DerSimonian and Laird (DL) [1].

218 For each contributing study, the DL technique calculates the logarithm of the sample
219 odds ratio and a corresponding estimate of the variance of this measure based on the
220 asymptotic distribution of these logarithms. Adjustments are made for entries in the
221 individual 2 x 2 tables that contain a zero-cell count. Equations 5-8 below capture the
222 core DL approach. In Equation 5, an estimate of the interstudy variability, τ^2 , is first
223 derived from Cochran’s Q statistic and the weights assigned to each of the k contributing
224 studies, ω_i . These weights are equal to the inverse of the square of the standard error of
225 the estimate of the odds ratio, $\hat{\theta}_i$, in each of the k contributing studies.

226

227
$$\hat{\tau}^2 = \frac{Q - (k-1)}{\sum \omega_i - \left(\frac{\sum \omega_i^2}{\sum \omega_i} \right)} \quad (5)$$

228

229

230 As shown in Equation 6, a new set of weights, ω'_i , are then calculated based on the
231 estimated value of $\hat{\tau}^2$ from Equation 5 and the standard errors of the contributing studies.

232

233

234
$$\omega'_i = \frac{1}{SE(\hat{\theta}_i)^2 + \tau^2} \quad [6]$$

235

236 These new weights are then used to calculate estimates of both the overall log odds ratio,
237 $\hat{\theta}_{DL}$ and its standard error as shown in Equation 7 and 8.

238

239

240
$$\hat{\theta}_{DL} = \frac{\sum \omega'_i \hat{\theta}_i}{\sum \omega'_i} \quad [7]$$

241

242
$$SE(\hat{\theta}_{DL}) = \frac{1}{\sqrt{\sum \omega'_i}} \quad [8]$$

243

244 A test of statistical significance is then based on a large sample normal distribution. The

245 DL technique requires asymptotic assumptions regarding both the Q statistic used to

246 estimate the interstudy variability, τ^2 , and the normal distribution required to test for
247 statistical significance. A more subtle issue is the possibility of distorting correlations
248 between the individual estimates of the effect size for each contributing study, θ_i , and
249 the individual weights used for each of these contributing effect sizes.

250 **RESULTS**

251 **A Non-Parametric Exact Test of Overall Statistical Significance for Dichotomous** 252 **Categorical Meta-Analysis**

253 Jakob Bernoulli's notion of what is now called a Bernoulli Trial offers the basis for a
254 non-parametric approach to aggregating multiple epidemiological studies based on
255 dichotomous categorical data. The enhancements to the Bernoulli method developed in
256 this paper offer a practical exact method for assessing the overall statistical significance
257 and effect size of a dichotomous meta-analysis.

258 One of the many important contributions of this outstanding 17th century mathematician
259 was the idea of the fixed probability of an event over a sequence of independent trials
260 which led to what is now called Bernoulli Trials and to the related Binomial Distribution.
261 In brief, Bernoulli viewed a set of statistical events as a series of independent coin flips
262 with each flip having a probability p of obtaining a head and $q = 1 - p$ of obtaining a tail.
263 This hypothetical coin is often treated as a fair coin where both p and q equal $.5$. The
264 simplest Bernoulli Trials approach encompasses a series of n flips and answers questions
265 of the type: what is the probability of observing x heads in n such flips? (See for example
266 Rosner [17]). In epidemiology, one could consider each of the k contributing studies of a
267 meta-analysis as a single Bernoulli Trial with $p = .5$. Then the combination of the k

268 studies could be analyzed as a binomial distribution. This is the standard Sign Test (see,
269 for example, [18]).

270 For example, for a meta-analysis of 20 studies, if 15 out of 20 studies had more cases in
271 the exposure group than in the control group, we could ask: What is the probability that
272 15 or more of the 20 studies could have shown a larger effect in the exposure group
273 strictly by chance alone? If this cumulative probability is less than a pre-specified level
274 of Type I error (e.g. .05), one would reject the null hypothesis and conclude there
275 probably exists a statistically reliable relationship between exposure and the end point
276 used.

277 The principal reason that this approach has seen little use in practical epidemiology is that
278 it suffers from two critical deficits. First, the dichotomous Bernoulli heads vs. tails
279 approach doesn't deal with the third possibility of a tie. The author of this study believes
280 that no truly useful method to date has been offered to deal with those situations when
281 there are an identical number of events in each of the exposure and the control arms of a
282 study other than to discard the study. Second, a truly exact EBT method requires a
283 complete convolution of the frequency distributions of the contributing studies in order to
284 derive the combined frequency distribution. Even for equal sample size, each of the k
285 contributing studies could have a different Bernoulli probability, p , requiring a full
286 convolution to determine the null distribution of the total number of times there were
287 more cases in the exposure group relative to the control group across the k contributing
288 studies. Before dealing with the ties problem, the determination of the combined
289 distribution will be outlined.

290 Combining the Individual Studies Contributing to the Meta-Analysis. A critical
291 problem is finding a method for combining the individual study binomial distributions of
292 the k contributing studies each with a possibly different p value into an overall frequency
293 distribution.
294 Prior to the widespread availability of computing power, the convolution of a large
295 number of individual binomial distributions was typically handled by approximate
296 methods given the unwieldy nature of the calculations. Even with the advent of available
297 computer power, convolution is still often impractical. As an example, for a meta-
298 analysis involving 24 studies each with a unique binomial distribution p , there are over 2
299 million unique combinations of the studies that need to be considered just to calculate the
300 single discrete probability that exactly 12 of the 24 studies have more cases in the
301 exposure group than in the control group.¹ However, an exact algorithm was laid out in a
302 readily implementable fashion by Butler and Stephens in a 1993 technical report [19]
303 which can easily be implemented even on a personal computer. The algorithm yields the
304 exact probability distribution of the convolution of individual binomial distributions
305 corresponding to the specific studies contributing to the meta-analysis. The method
306 makes use of a recurrence relationship inherent in the binomial distribution which allows
307 the semi-automatic calculation of its probabilities without resort to the simple but
308 overwhelmingly inefficient enumeration of all of the possible combinations of studies.
309 This easily established relationship can be stated as:
310

¹ This number of combinations is simply $C(24,12) = 2,704,156$.

A non-parametric exact test for meta-analysis of categorical data

311
$$P(X = 0) = (1 - p)^n \quad \text{if } j = 0$$

312

313
$$P(X = j) = \left\{ \frac{(n - j + 1)}{j} \right\} * \left\{ \frac{p}{(1 - p)} \right\} * P(X = j - 1) \quad \text{if } j \geq 1$$

314

315 Figure 2 compares the estimated number of computer executable steps required in the
316 Butler & Stephens method compared to a traditional convolution.

317

318 **Insert About Here**

319 **Fig. 2 Estimated computer executable steps per Butler & Stephens vs. traditional**
320 **convolution.**

321

322 As can be seen, a traditional convolution is only tractable when the number of
323 contributing studies is less than or equal to approximately 20.

324 **The Ties Problem.** The next problem in adapting the standard Bernoulli Trials
325 technique to practical meta-analysis is a procedure to deal with the situation where there
326 are an identical number of cases in both the exposure and control arms of a study
327 contributing to the meta-analysis. In studies with small sample sizes and/or low disease
328 probabilities, the highest probability tie is typically the “0/0” tie in which no cases are
329 observed in either the exposure or the control arms.

330 A first step in dealing with ties is to more clearly define the criteria for a “success”. The
331 present EBT approach defines a success as there being a **strictly** greater number of cases
332 in the exposure group relative to the control group. Under this definition, the same

333 number of cases in both arms of the study or more cases in the control arm of the study is
334 considered a “failure”. In essence, this is a trinomial situation. There are successes,
335 failures and ties. We are simply combining the failures where there are more cases in the
336 control group relative to the exposure group and tie situations and calling the
337 combination “failures.”

338 Equation 9 below forms the basis of the EBT method. The Greek capital letter “ Π ” has
339 been chosen to specify the probabilities of there being more cases in one arm of the study
340 relative to the other to differentiate these parameters from the underlying disease
341 probabilities:

$$342 \quad \Pi_{E_i} + \Pi_{C_i} + \mathit{prob}(tie)_i = 1 \quad [9]$$

343 where:

344 Π_{E_i} = probability of there being strictly more cases in the exposure group relative to the
345 control group in Study i

346 Π_{C_i} = probability of there being strictly more cases in the control group relative to the
347 exposure group in Study i

348 $\mathit{prob}(tie)_i$ = probability of finding exactly same number of cases in both groups of Study
349 i .

350 Assuming that Π_{E_i} and Π_{C_i} would be equal under the null hypothesis of no difference
351 between exposure and control groups and rearranging terms, we have:

352

$$2\Pi_{E_i} + \mathit{prob}(tie)_i = 1 \quad [10]$$

353 Solving for Π_{E_i} we have:

354
$$\Pi_{E_i} = \frac{1 - \mathit{prob}(\mathit{tie})_i}{2} \quad [11]$$

355 Thus, the only requirement for calculating the Π_{E_i} parameter for each contributing study
 356 is to first determine the probability of all tie situations for the study.

357 This is a very straightforward procedure. To determine $\mathit{prob}(\mathit{tie})_i$ for each of the
 358 contributing studies, all of the tie situations need to be enumerated and then summed
 359 together.

360 As a simple example, assume that Study i has 100 participants in each of its exposure and
 361 control arms and that the underlying event (disease) probability p is .01.

362 The probability that there are no cases among these 100 participants in the exposure arm
 363 would then be:

364 $\mathit{Prob}(0 \text{ cases}) = .01^0 * (1 - .01)^{100} = .99^{100} = .37$

365 Similarly, the probability of there being no cases in the control arm would also be .37.

366 Thus, the probability of a “0,0” tie would be $.37^2 = .13$ which is surprisingly large.

367 Table 3 lists the probabilities for the first five tie situations and sums these probabilities
 368 to determine $\mathit{prob}(\mathit{tie})_i$.²

Table 3 Probability of observing exactly the same number of cases in both the exposure and control groups for background event probability equal to .01 and sample size equal to 100 as a function of the number of observed cases	
Number of Cases in Each Group	Probability

² The probabilities for six or more ties decrease to extremely small values. However, in actuality, the EBT method calculates all possible ties in calculating $\mathit{prob}(\mathit{tie})_i$

A non-parametric exact test for meta-analysis of categorical data

0	.13
1	.137
2	.034
3	.004
4	.0002
5	.00001
Total	.309

369

370 As shown in Table 3, there is over a 30% probability of obtaining a tie for 0 cases
371 through 5 cases in both the exposure and control groups. Applying Equation [11] to this
372 hypothetical study, we see that, under the null hypothesis of equal probabilities, Π_{E_i} and
373 Π_{C_i} are both equal to .35. Thus, due to ties, the nominal .50 value for Π_{E_i} and Π_{C_i} has
374 been greatly reduced.

375 The EBT technique is indeed a “vote counting” method and such methods have been
376 greatly disparaged by Rothman [20] among others. However, unlike a simple Sign Test,
377 the EBT method is based on a reasonable approach to the ties problem and combines the
378 individual P_{E_i} values by doing the equivalent of a formal convolution of the frequency
379 distributions of the individual contributing studies.

380

381 **A Non-Parametric Exact Method for the Estimation of Effect Size for Dichotomous**
382 **Categorical Meta-Analysis**
383 **Basic Estimation Technique**

A non-parametric exact test for meta-analysis of categorical data

384 A second exact technique was developed to estimate the effect size for dichotomous
385 categorical meta-analysis. As a starting point, one might simply form the ratio of the
386 average observed event probabilities, p_{E_i} and p_{C_i} , in the exposure and control groups
387 respectively of each study and average these ratios across the k contributing studies.
388 This simple approach, however, is highly biased. As shown in the underlying model that
389 is described in Equations 1 – 4, the number of observed “successes” in the exposure and
390 control arms of the k contributing studies each depend on an identical source of variation
391 captured by ϵ_{1i} in the model. The exposure group, however, contains an additional
392 source of variation, captured by ϵ_{2i} in the model. Figure 3 illustrates the problem of
393 estimating the effect size by simply forming the ratio of p_E to p_C .

394 **Insert About Here**

395 **Fig. 3 Demonstration of inappropriateness of comparing the p_E and p_C**
396 **distributions to estimate Effect Size.**

397

398 Even for the relative risk of 1.0 depicted in the figure, the exposure distribution will have
399 positive excursions that are not compensated for by equally robust negative excursions at
400 least for small (rare) values of event probability.

401 The differential skew of the p_{E_i} distribution relative to the p_{C_i} distribution was used to
402 address this issue. The additional skew in the exposed group due to the source of ϵ_{2i} was
403 estimated by taking the difference between the total exposure group skew and the
404 expected skew from a pure binomial with the same observed event probability. The

A non-parametric exact test for meta-analysis of categorical data

405 observed average p_E across the k contributing studies was then reduced by a factor
406 proportional to this difference in skew levels.

407 **Monte Carlo Simulation of The EBT And DL Techniques for Statistical Significance** 408 **and Effect Size Estimation**

409 A series of Monte Carlo simulations was conducted to evaluate the EBT statistical
410 significance test and the effect size estimation techniques and to compare them to the
411 typically used DerSimonian-Laird Inverse Variance technique. The simulation was
412 written and executed in *R: A Programming Environment for Data Analysis and Graphics*.
413 [21]. The DerSimonian results were calculated using the “meta” package in R.
414 Five levels of relative risk (ratio of exposure group to control group event probability) of
415 1.0, 1.25, 1.5, 1.75, and 2.0 were crossed with three levels of disease background event
416 probability (.005, .01, and .05), and three levels of sample size (50,100 and 200). Finally,
417 the number of studies entering into each meta-analysis was chosen to be 5, 10, 20, or 40
418 studies.

419 In addition, the heterogeneity between the contributing studies, τ^2 [Equation 4], was
420 evaluated at 0 (homogeneity), .4, and .8. This last value of .8 represents a very large
421 variance among the studies and was partially chosen to be able to compare the results
422 with previous work ([22] and [16]). As an example, at $\tau^2 = .8$ a nominal exposure group
423 event probability p_E of .05 would vary from of .009 to .24 which is almost a 30:1 ratio.
424 Finally, the common variability in both the exposure and control groups represented by
425 γ^2 in Equation 1 was chosen to be .5 to allow direct comparison with earlier work ([22]
426 and [16]).

427 The statistical significance and effect size were evaluated using both the EBT and
 428 DerSimonian techniques for each replication. All simulation runs were conducted with
 429 10,000 replications. A value of .05 was used as the pre-specified level of Type I Error.
 430 The “Mid-P” technique advocated by Agresti [23] and others was used to determine the p
 431 values in a less conservative manner leading to more realistic power levels.

432

433 **Results from the Monte Carlo Simulations: Testing Statistical Significance**

434 Figures 4 and 5 show the results using the EBT and DL methods respectively.
 435 To simplify presentation, only scenarios in which the expected number of cases was
 436 greater than or equal to two were utilized. Table 4 shows the included scenarios.

Table 4 Scenarios included in the analysis of statistical significance		
Background Event Probability	Sample Size	Expected Number of Observed Cases
.01	200	2
.05	50	2.5
.05	100	5
.05	200	10

437

438

439

Insert About Here

440 **Fig. 4 Power for number of studies, relative risk, and heterogeneity for EBT. A-D**
 441 **correspond to studies equal to 5, 10, 20 and 40.**

442

443

Insert About Here

444 **Fig. 5 Power for number of studies, relative risk, and heterogeneity for DL. A-D**
445 **correspond to studies equal to 5, 10, 20 and 40.**

446

447 The basic finding was that the EBT method maintained the prespecified level of Type I
448 error for both the homogeneous and heterogeneous scenarios while the DL method had
449 many violations of this level for heterogeneous scenarios. This limitation of the DL
450 method is highlighted in the inserts in Figure 5 which show the power when the Relative
451 Risk is one (i.e. the false alarm rate). For the homogeneous scenario where $\tau^2 = 0$, both
452 the EBT and the DL methods respect the prespecified Type I error level. However, for τ^2
453 $= .4$ and for $\tau^2 = .8$, the DL method exhibits large violations of this level. Not
454 unexpectedly, as the number of contributing studies increases, the power for Relative
455 Risk greater than one increases for both the EBT and DL methods for values of Relative
456 Risk of greater than 1.0. A separate analysis showed that the standard deviation of the
457 power estimates in Figures 4 and 5 was less than or equal to .42% (i.e. .0042).
458 In actuality, comparing the power between the EBT & DL techniques for event ratios
459 greater than 1.0 is not possible due to the large number of violations of the pre-specified
460 Type 1 Error violations for the DL technique.

461 Figure 6 is a comparison of Type I Error (false alarm rate) for the current EBT technique
462 and the DL technique as a function of heterogeneity (τ^2).

463

Insert About Here

464 **Fig. 6 Type I error for EBT and DL methods as a function of heterogeneity.**

465 As can be clearly seen, the current EBT technique is relatively resistant to the effects of
466 increasing heterogeneity over a very large range. The DL technique, however exhibits a

467 monotonically increasing sensitivity to heterogeneity. A related aspect of any meta-
468 analysis technique's ability to perform well in the face of heterogeneity is its resistance to
469 "contamination" from one or a small number of "rogue studies." Since the EBT method
470 does not directly allow such rogue studies to directly affect the test statistic, it should be
471 much more resistant to these distortions.

472 The large costs of discreteness have been studied by Agresti [24] and others.

473 A first cost of discreteness results when the number of contributing studies is small. The
474 general issue of overcoverage is highlighted in Figure 7.

475 **Insert About Here**

476 **Fig. 7 Interval overcoverage as a function of the number of contributing studies.**

477

478 The overcoverage is greatest for the smallest number of k contributing studies, and
479 generally decreases as the number of contributing studies increases. As Figure 7
480 demonstrates, even an unrealistic level of 500 contributing studies is still associated with
481 a relatively large level of overcoverage. While such discreteness clearly reduces power,
482 it could be argued that a statistically significant finding based on extremely sparse tables
483 and a handful of studies requires stronger evidence. Unfortunately, the majority of meta-
484 analyses consist of fewer than two or three studies as Kontopantelis *et al.* have shown in
485 their extensive analysis of all meta-analyses in the Cochrane Library [25].

486 Additional Monte Carlo testing was done for unbalanced designs (unequal sample sizes
487 in the exposure and control arms of the contributing studies) and meta-analyses with
488 unequal sample sizes across contributing studies. Table 5 shows the sample sizes for the
489 two groups for a typical unbalanced design in which the control group sample size is

490 twice the exposure group sample size. The sum of the two sample sizes across both arms
 491 of the study was chosen to be 200 yielding an average sample size of 100 to allow
 492 comparison with the balanced designs of Figures 4 and 5.
 493

Table 5 Sample sizes for simulation of unbalanced designs										
Number of studies = 10										
	Study #									
Group	1	2	3	4	5	6	7	8	9	10
Exposure	66	66	66	66	66	66	66	66	66	66
Control	134	134	134	134	134	134	134	134	134	134

494
 495 Table 6 below shows the results of the simulation for heterogeneity values $\tau^2 = 0$ and $\tau^2 =$
 496 .8, Event (“disease”) Probability of .05, Number of Studies = 10, and Sample Size (avg.)
 497 = 100 at the same five levels of Relative Risk used above. The simulation run consisted
 498 of 10,000 replications as in Figures 4 and 5.
 499

500

Table 6 Power (%) for the unbalanced design of table 5 τ^2 (heterogeneity) equal to 0 and .8; event probability = .05; Number of studies equal to 10; Sample size (per study arm) equal to 100										
Technique	Heterogeneity									
	0					.8				
	Risk Ratio									
	1.0	1.25	1.5	1.75	2.0	1.0	1.25	1.5	1.75	2.0
EBT	2.1	14.4	43.3	71.0	88.4	4.2	11.0	21.6	35.0	46.9
DerSimonian & Laird	2.2	16.5	57.2	87.5	97.7	11.5	24.0	41.6	59.2	72.8

501

502 As the results in Table 6 indicate, the Type I Error (Relative Risk = 1.0) remained below
 503 the specified value of five percent for the EBT technique and was far above this point for
 504 the DerSimonian technique when the heterogeneity was equal to .8.

505 Table 7 below shows the sample sizes for the exposure and control groups for each of the
 506 contributing studies for a design with unequal sample size across the contributing studies.

507 This particular design was chosen as a relatively extreme case. As can be seen, the
 508 average sample size for both the groups was maintained at 100 to allow comparison of
 509 the simulation results with the equal sample size scenarios of Figures 4 and 5.

510

Table 7 Sample sizes for simulation of unequal sample size designs										
Number of studies equal to 10										
	Study #									
Group	1	2	3	4	5	6	7	8	9	10
Exposure	175	25	175	25	175	25	175	25	175	25
Control	175	25	175	25	175	25	175	25	175	25

511

512

513 Table 8 below shows the results of the simulation for a heterogeneity values of $\tau^2 = 0$ and

514 $\tau^2 = .8$, Event (“disease”) Probability of .05, and Sample Size (individual study arm

515 average) = 100, at the same five levels of Relative Risk as used above. The simulation

516 run consisted of 10,000 replications as in Figures 4 and 5.

517

518

519

<p align="center">Table 8 Power (%) for the unbalanced design of table 7</p> <p align="center">τ^2 (heterogeneity) equal to 0 and to .8; Event probability = .05;</p> <p align="center">Number of studies equal to 10; Sample size (avg. per individual study arm) equal to 100</p>										
	Heterogeneity									
	0					.8				
	Risk Ratio									
Technique	1.0	1.25	1.5	1.75	2.0	1.0	1.25	1.5	1.75	2.0
EBT	2.0	14.9	43.0	70.0	87.7	4.3	11.4	22.4	34.5	47.9
DerSimonian & Laird	2.4	16.6	56.6	87.2	97.8	11.4	25.1	41.6	58.2	73.5

520

521 Most importantly, the EBT Technique was superior at protecting the pre-specified level

522 of Type I Error relative to the DL technique at a heterogeneity level of .8.

523 A clear finding of the Monte Carlo simulations common to both meta-analysis techniques

524 studies is the apparent fruitlessness of searching for small effect sizes. Both the EBT and

525 DL techniques are very poor at reliably finding statistically significant results until the

526 relative risk approaches 2.0. While this finding does not directly bear on the issues

527 studied in this report, it does serve as a cautionary tale to those who continue to try to

528 tease out very small effects especially from sparse data.

529 **Results from the Monte Carlo Simulations: Effect Size Estimation**

A non-parametric exact test for meta-analysis of categorical data

530 Figures 8 and 9 capture the basic findings for estimating the Effect Size.

531

532

Insert About Here

533 **Fig. 8 Effect Size as function of Relative Risk and heterogeneity. A & B correspond**
534 **to the EBT and DL methods.**

535

536 Again, only simulation scenarios in which the expected number of observed cases was
537 greater than or equal to two were utilized. Since the effect of the number of studies
538 contributing to the meta-analysis was small for this effect size estimation, results were
539 averaged across this variable. Figure 8 shows the results for the EBT method and Figure
540 9 for the DL method. As shown in the figures, both methods were reasonably successful
541 at estimating the levels of relative risk. However, the EBT method generally
542 underestimated the relative risk for $\tau^2 = 0$ and overestimated it for $\tau^2 = .8$ while the DL
543 method tended to underestimate the relative risk. Finally, the interquartile range for the
544 DL method was considerably smaller than for the EBT method as shown in Figures 10
545 and 11.

546

547

Insert About Here

548 **Fig. 9 Semi-Interquartile Range as function of Relative Risk and heterogeneity. A**
549 **& B correspond to the EBT and DL methods.**

550

551

552 **CONCLUSIONS AND SUGGESTIONS FOR THE FUTURE**

A non-parametric exact test for meta-analysis of categorical data

553 This research has developed an exact test for the meta-analysis of dichotomous,
554 categorical data and a related method to estimate the size of the effect.

555 **The Enhanced Binomial Technique (EBT) to Assess Statistical Significance**

556 The EBT technique was greatly superior to the DerSimonian technique in maintaining a
557 pre-specified level of Type I Error. As shown, the DerSimonian technique demonstrated
558 many large violations of this level when heterogeneity was present. Given the various
559 biases towards finding statistical significance prevalent in epidemiology today, a strong
560 focus on maintaining a pre-specified level of Type I Error would seem critical (see, e.g.,
561 [26]). The EBT approach is greatly superior at maintaining this pre-specified value of
562 Type I Error in the face of even extreme heterogeneity.

563

564 **The Enhanced Binomial Technique (EBT) to Estimate Effect Size**

565 A related but separate method was developed to estimate the effect size. This new
566 technique was comparable to the often-used DL method although both methods
567 demonstrated some accuracy issues. The DL method exhibited a somewhat smaller
568 Semi-IQR variability. The fact that the EBT method was clearly superior in assessing
569 statistical significance while the DL method demonstrated a smaller variability in
570 estimating effect Size supports the possible utility of separating these two procedures as
571 outlined at the beginning of this article. One possibility is to use the EBT and DL
572 methods for statistical significance assessment and effect size estimation respectively.

573

574 While statistical programs providing exact solutions already exist such as Cytel's
575 StatXact, they are beyond the means of most practicing statisticians and epidemiologists.

A non-parametric exact test for meta-analysis of categorical data

576 For example, Cytel Inc. currently lists a price of over \$900 USD for their current version,

577 StatXact 11 [27]

578 The techniques developed here are written in the almost universal statistical language of

579 R and are freely available from the author. As such, it is hoped that other researchers

580 would be able to extend and improve these initial versions.

581 As outlined in this report, the use of meta-analysis in epidemiology is increasing very

582 rapidly and appears to be meeting an important need. Fortunately, inexpensive and

583 readily available computer power has also vastly increased in the past forty years. For

584 example, task speed as measured in Million Instructions per Second (“MIPS”) has

585 increased from .64 for the IBM370 mainframe computer in 1972 to 238,000 for an Intel

586 Pentium processor personal computer in 2014. [28]. By using the techniques developed

587 here and the computer power available to all researchers today, the determination of

588 statistical significance and the estimation of effect size can be readily accomplished

589 without unnecessary error.

590 **Declarations**

591 **Ethics approval and consent to participate**

592 No experimental participants

593

594 **Consent for publication**

595 No consent required

596

597 **Availability of data and material**

598 Both software programs are freely available from the author

599 Competing interests

600 The author has no competing interests

601

602 Funding

603 No funding was obtained for this work

604

605 Authors' contributions

606 Work was fully done by L. Paul

607

608 Acknowledgements

609 No acknowledgements

610

611 REFERENCES

612

1. **DerSimonian R, Laird N. Meta-analysis in clinical trials. *Controlled Clinical Trials*. 1986; 7.3: p. 177-188.**
2. **Shapiro S. Meta-analysis/Shmeta-analysis. *Am J Epidemiol*. 1994 Nov 1; 140(9): p. 771-8.**
3. **Borenstein M, et al.. *Introduction to meta-analysis.*: John Wiley & Sons; 2009.**
4. **Higgins JP, editor. *Cochrane handbook for systematic reviews of interventions*. Chichester: Wiley-Blackwell; 2008.**
5. **Neyman J, Pearson ES. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*. 1933: p. 289-337.**

6. Higgins J, Deeks JJ, Altman DG. Special topics in statistics. In Higgins J, editor. *Cochrane Handbook for Systematic Reviews of Interventions: Cochrane Book Series*. Chichester: John Wiley; 2008. p. 481-529.
7. Mons U, Müezzinler A, Gellert C, Schöttker B, Abnet C, Bobak M, et al. Impact of smoking and smoking cessation on cardiovascular events and mortality among older adults: meta-analysis of individual participant data from prospective cohort studies of the CHANCES consortium. *bmj*. 2015 Apr 20; 350(h1551).
8. Hirsch L, Jette N, Frolkis A, Steeves T, Pringsheim T. The incidence of Parkinson's disease: a systematic review and meta-analysis. *Neuroepidemiology*. 2016; 46(4): p. 292-300.
9. Alzheimer's Association 2015 Alzheimer's disease facts and figures. *Alzheimer's & dementia. Journal of the Alzheimer's Association*. 2015 Mar; 11(3): p. 332.
10. Torre L, Siegel R, Ward E, Jemal A. Global cancer incidence and mortality rates and trends—an update. *Cancer Epidemiology and Prevention Biomarkers*. 2016 Jan; 25(1): p. 16-27.
11. Mantel N, Haenszel W. Statistical aspects of the analysis of data from retrospective studies of disease. *J Natl Cancer Inst*. 1959 Apr; 22(4): p. 719-48.
12. Mantel N, Fleiss J. Minimum expected cell size requirements for the Mantel-Haenszel one-degree-of-freedom chi-square test and a related rapid procedure. *Am J Epidemiol*. 1980; 112(1): p. 129-34.
13. Fisher RA. *Statistical methods for research workers*. Edinburgh: Oliver and Boyd; 1925.
14. Higgins JPT. Commentary: Heterogeneity in meta-analysis should be expected and appropriately quantified. *International Journal of Epidemiology*. 2008; 37: p. 1158-1160.
15. Higgins JPT, Thompson SG, Deeks JJ, Altman DG. Measuring inconsistency in meta-analyses. *British Medical Journal*. 2003; 327(7414): p. 557.
16. Bhaumik DK, Amaty A, Normand SLT, Greenhouse J, Kaizar E, Neelon B, et al. Meta-analysis of rare binary adverse event data. *Journal of the American Statistical Association*. 2012; 107(498): p. 555-567.
17. Rosner B. *Fundamentals of Biostatistics*. 5th ed.: Duxbury Press; 1999.

18. Siegel S, Castellan NJ. Nonparametric statistics for the behavioral sciences. 2nd ed.; 1988.
19. Butler K, Stephens M. Distribution of a Sum of Binomial Random Variables. Technical Report No. 467 prepared under contract N00014-92-5-1264 (NR-042-267) for The Office of Naval Research. Palo Alto: Stanford University; 1993.
20. Rothman K, Greenland S. Meta-analysis. Some methods to avoid: Qualitative tally (vote counting) and Quality scoring. New York: Lippincott Williams Wilkins; 1998.
21. The R Project for Statistical Computing.. Available from: <https://www.r-project.org/>.
22. Paul LM. Cannons and sparrows: an exact maximum likelihood non-parametric test for meta-analysis of $k \times 2 \times 2$ tables. Emerging themes in epidemiology. 2018 December; 15(1).
23. Agresti A. A survey of exact inference for contingency tables. Statistical science. 1992; 7(1).
24. Agresti A. Dealing with discreteness: making exact confidence intervals for proportions, differences of proportions, and odds ratios more exact. Statistical Methods in Medical Research. 2003 Feb; 12(1).
25. Kontopantelis E, Springate D, Reeves D. A re-analysis of the Cochrane Library data: the dangers of unobserved heterogeneity in meta-analyses. PloS one. 2013 Jul 26; 8(7).
26. Ioannidis JPA. Why most published research findings are false. PLoS Med. 2005: p. e124.
27. Cytel Inc. Cytel. [Online].; 2020 [cited 2020 November 17]. Available from: <https://store.cytel.com/products/statxact?hsCtaTracking=be2ed66d-9346-4239-ad8b-c19193bfcda0%7C40b5c432-854f-4116-a2d2-079223b15428>.
28. Wikipedia. [Document: "Instructions per Second"].; 2016 [cited 2016 June 3]. Available from: https://en.wikipedia.org/wiki/Instructions_per_second.

Figures

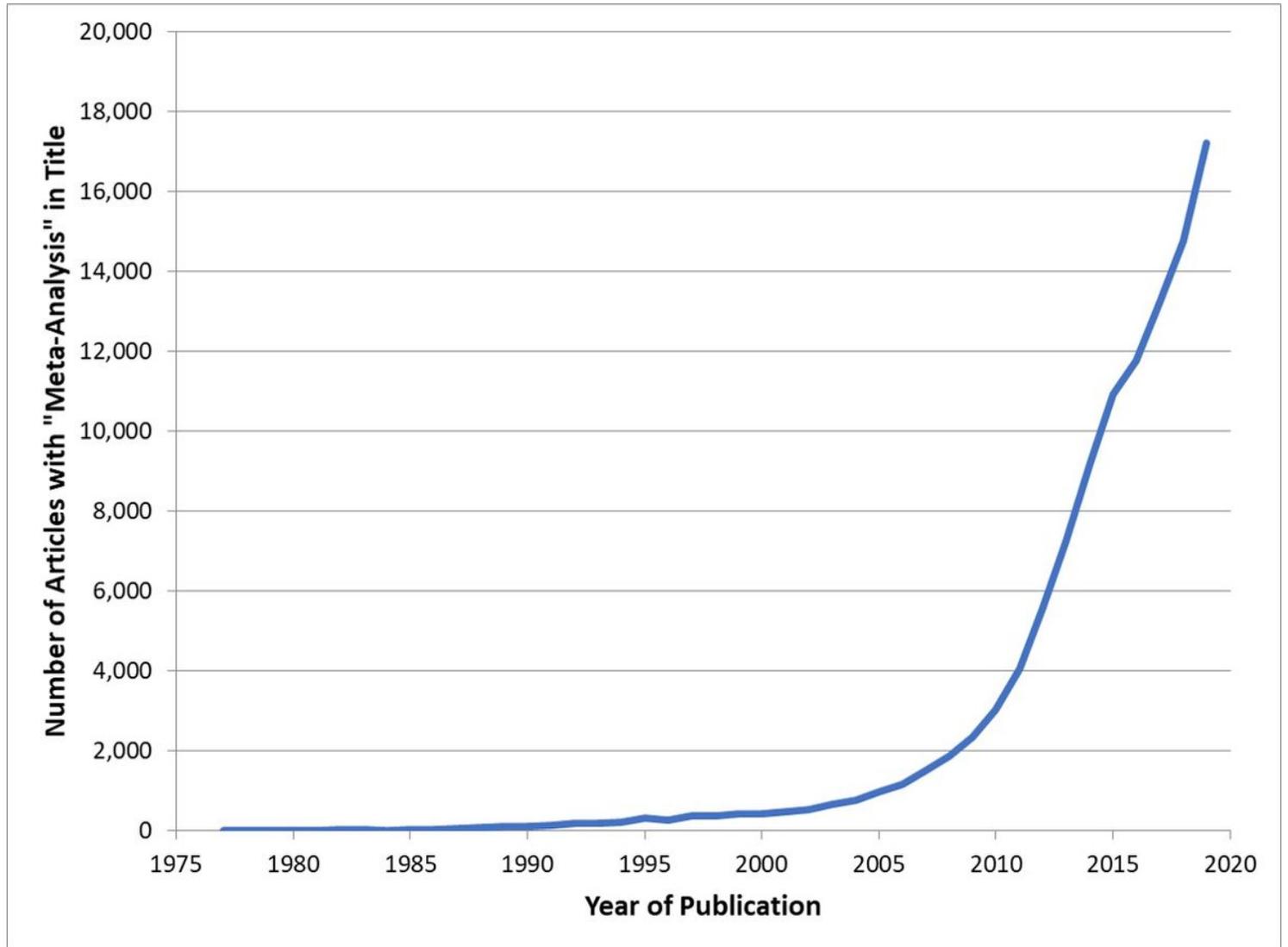


Figure 1

Number of articles containing "meta-analysis" in the title by year of publication.

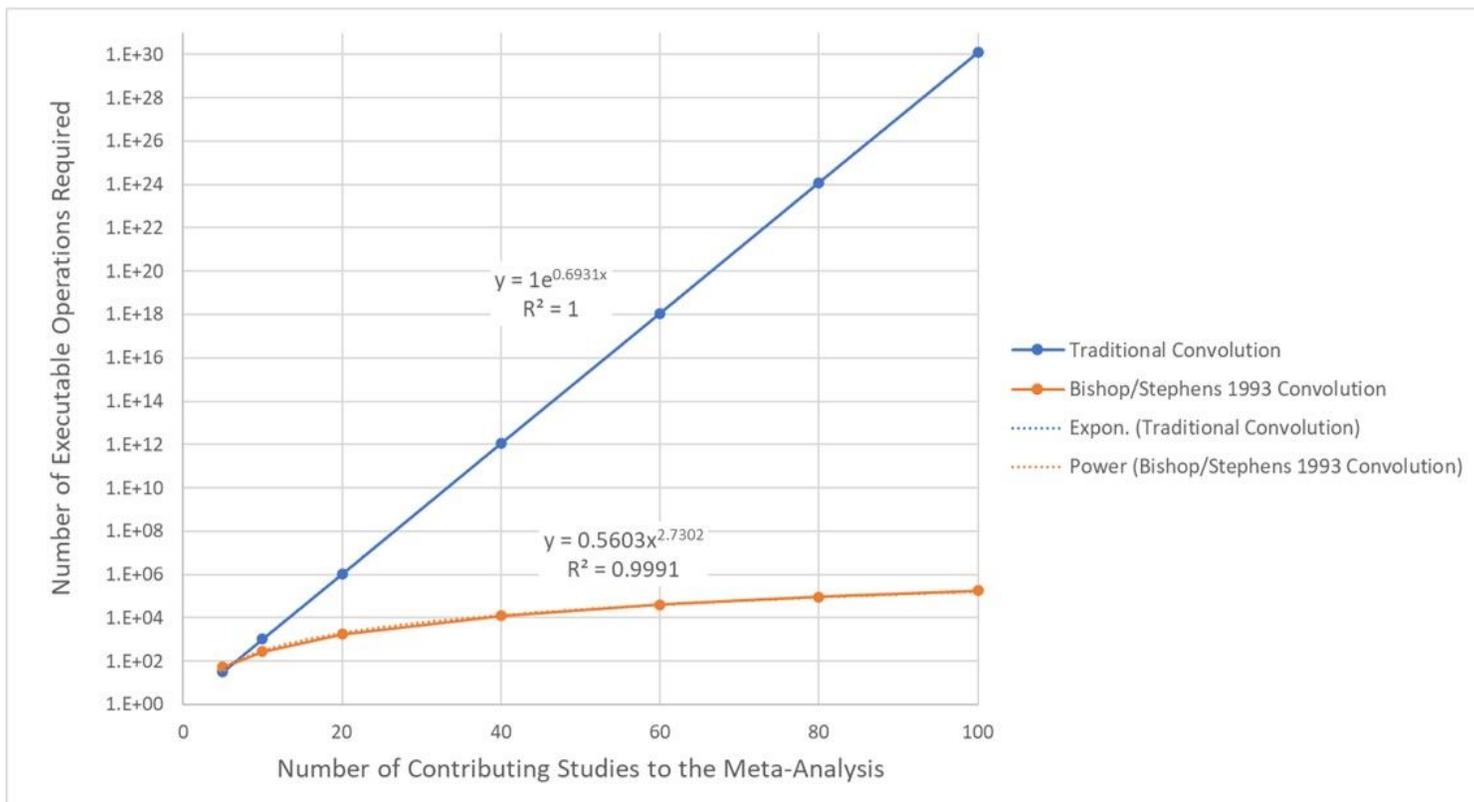


Figure 2

Estimated computer executable steps per Butler & Stephens vs. traditional convolution.

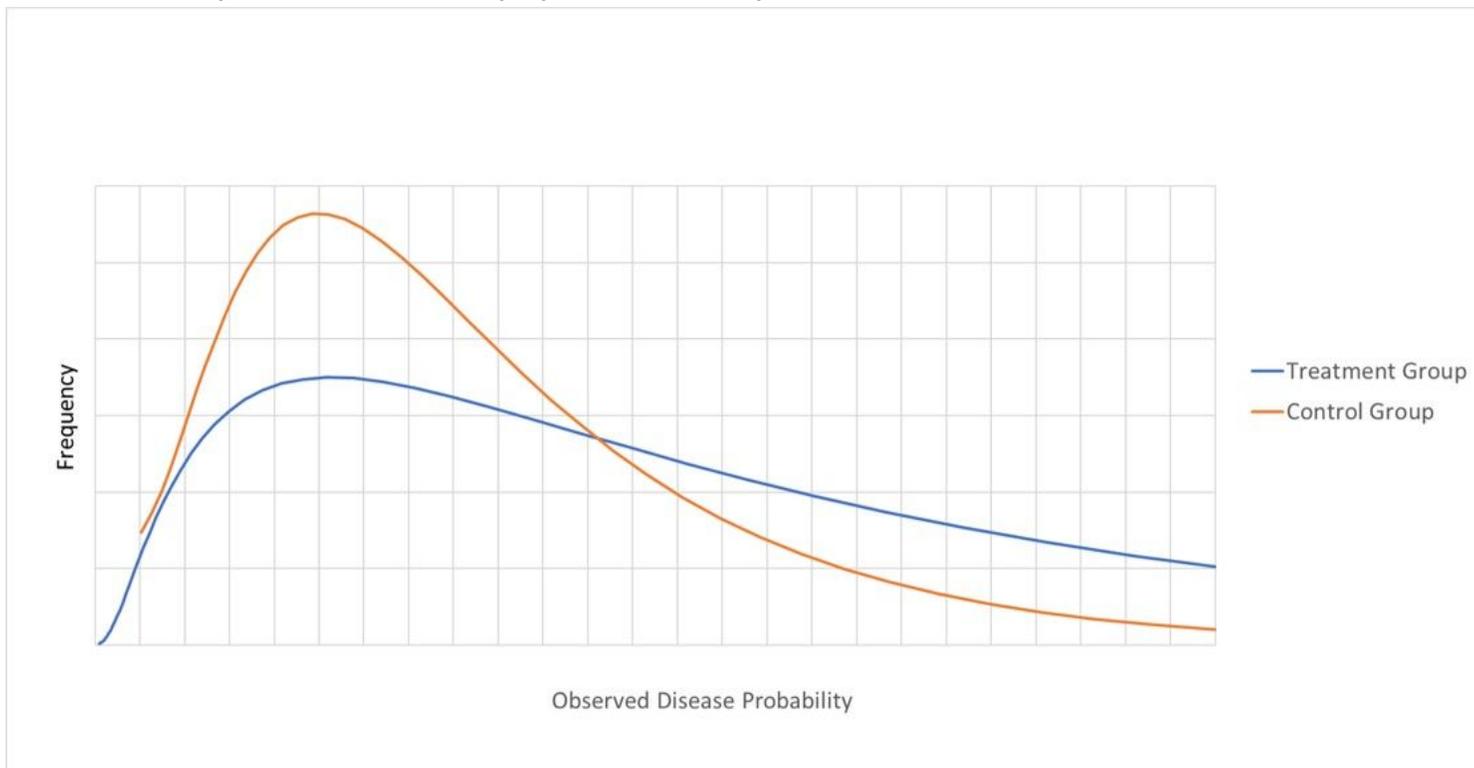


Figure 3

Demonstration of inappropriateness of comparing the PE and PC distributions to estimate Effect Size.

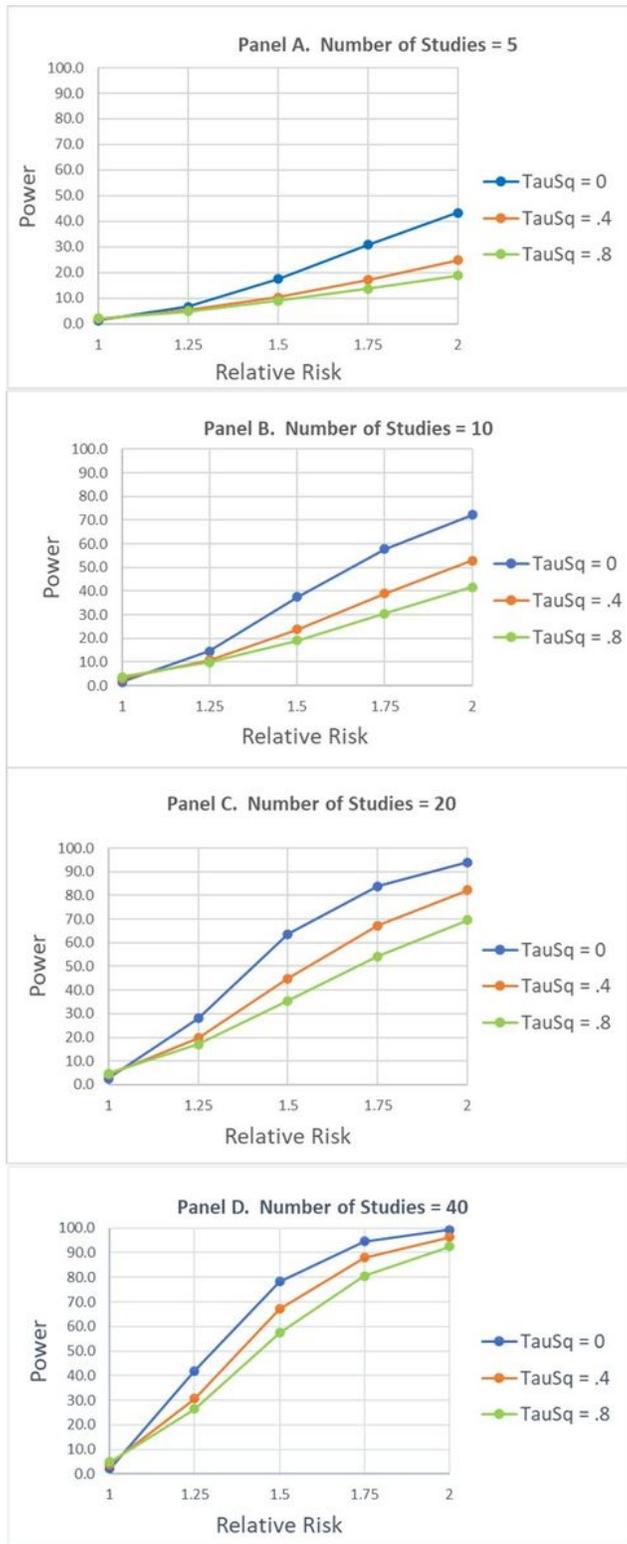


Figure 4

Power for number of studies, relative risk, and heterogeneity for EBT. A-D correspond to studies equal to 5, 10, 20 and 40.

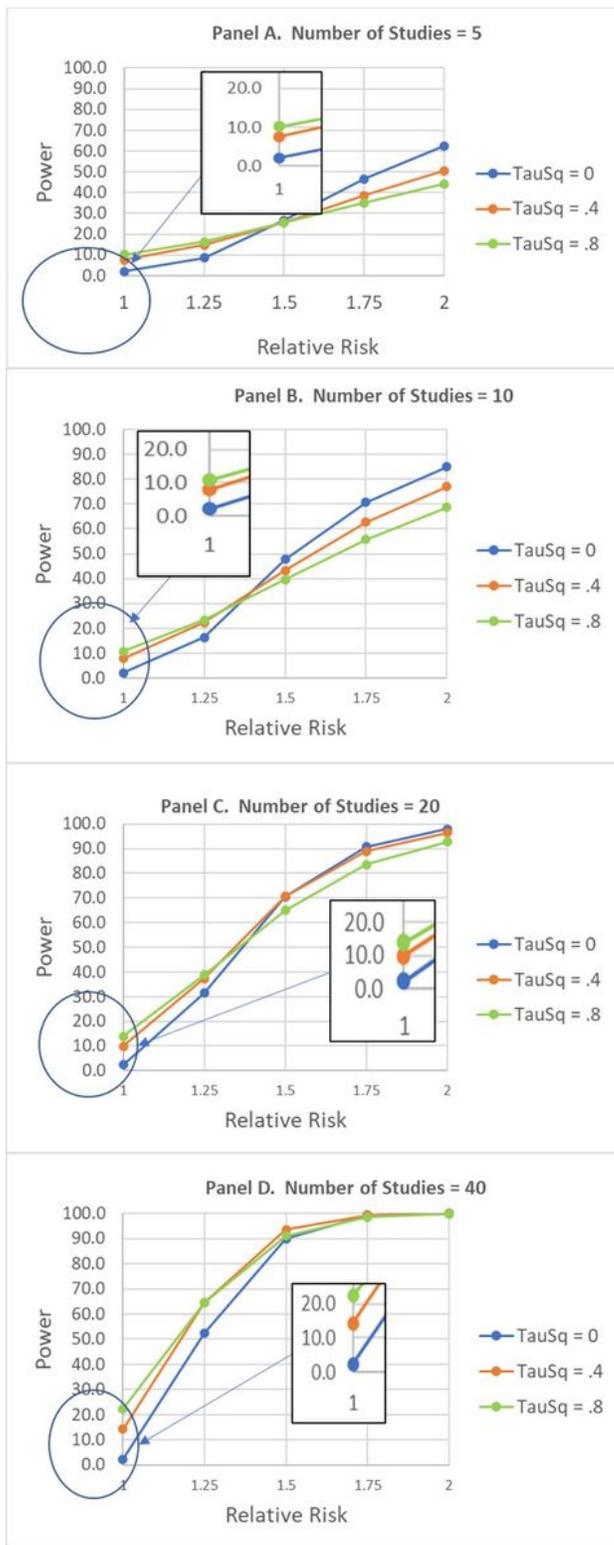


Figure 5

Power for number of studies, relative risk, and heterogeneity for DL. A-D correspond to studies equal to 5, 10, 20 and 40.

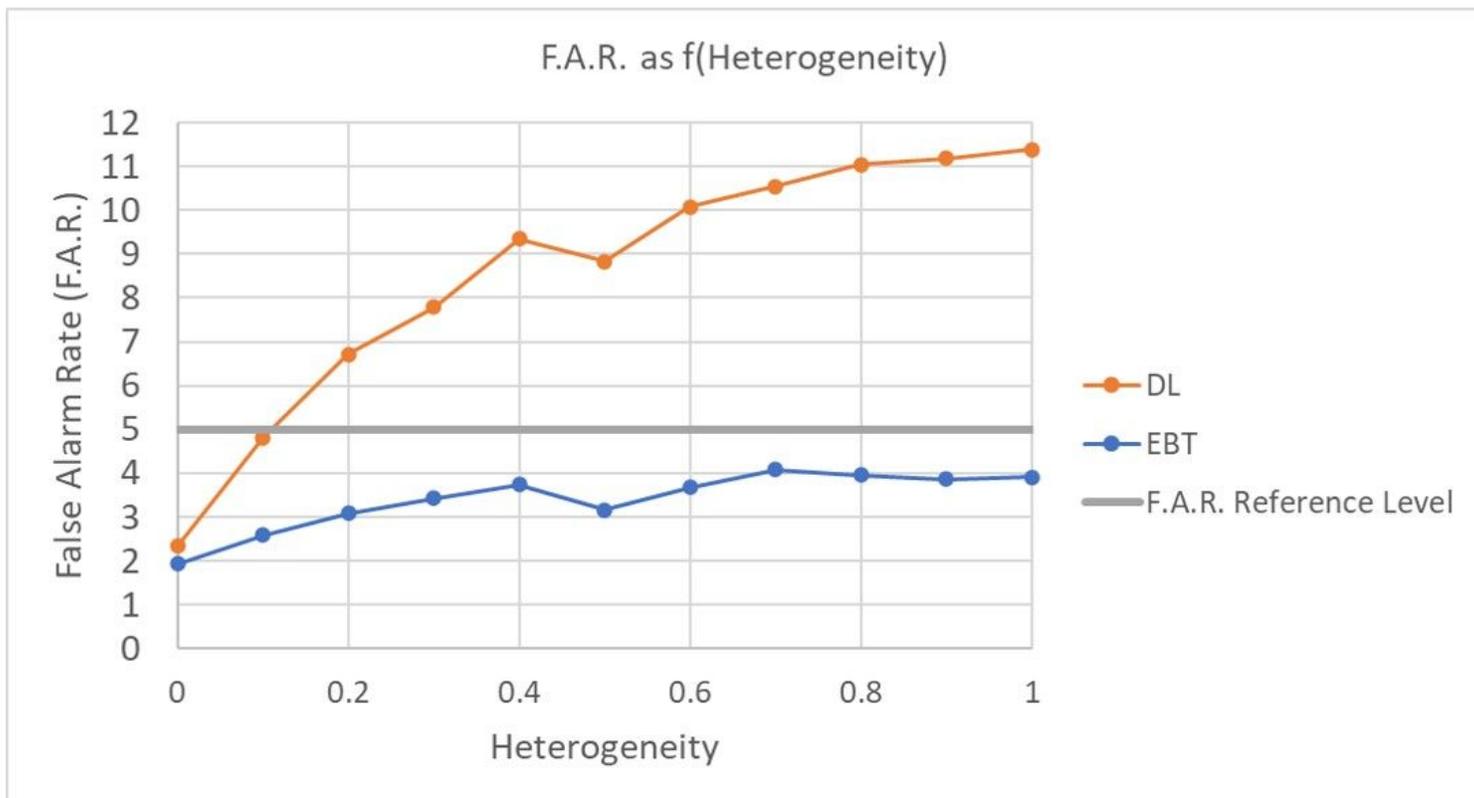


Figure 6

Type I error for EBT and DL methods as a function of heterogeneity.

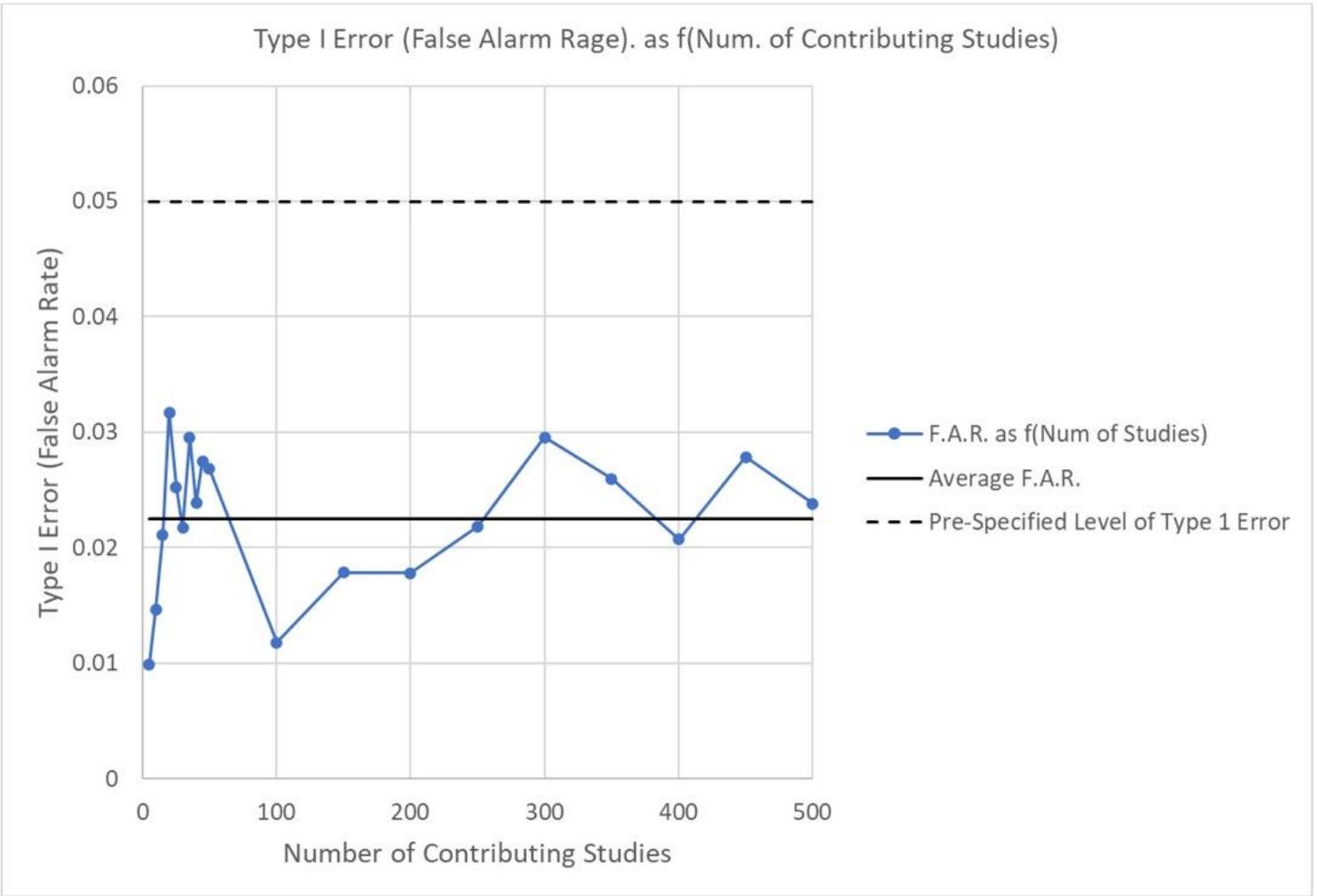


Figure 7

Interval overcoverage as a function of the number of contributing studies.

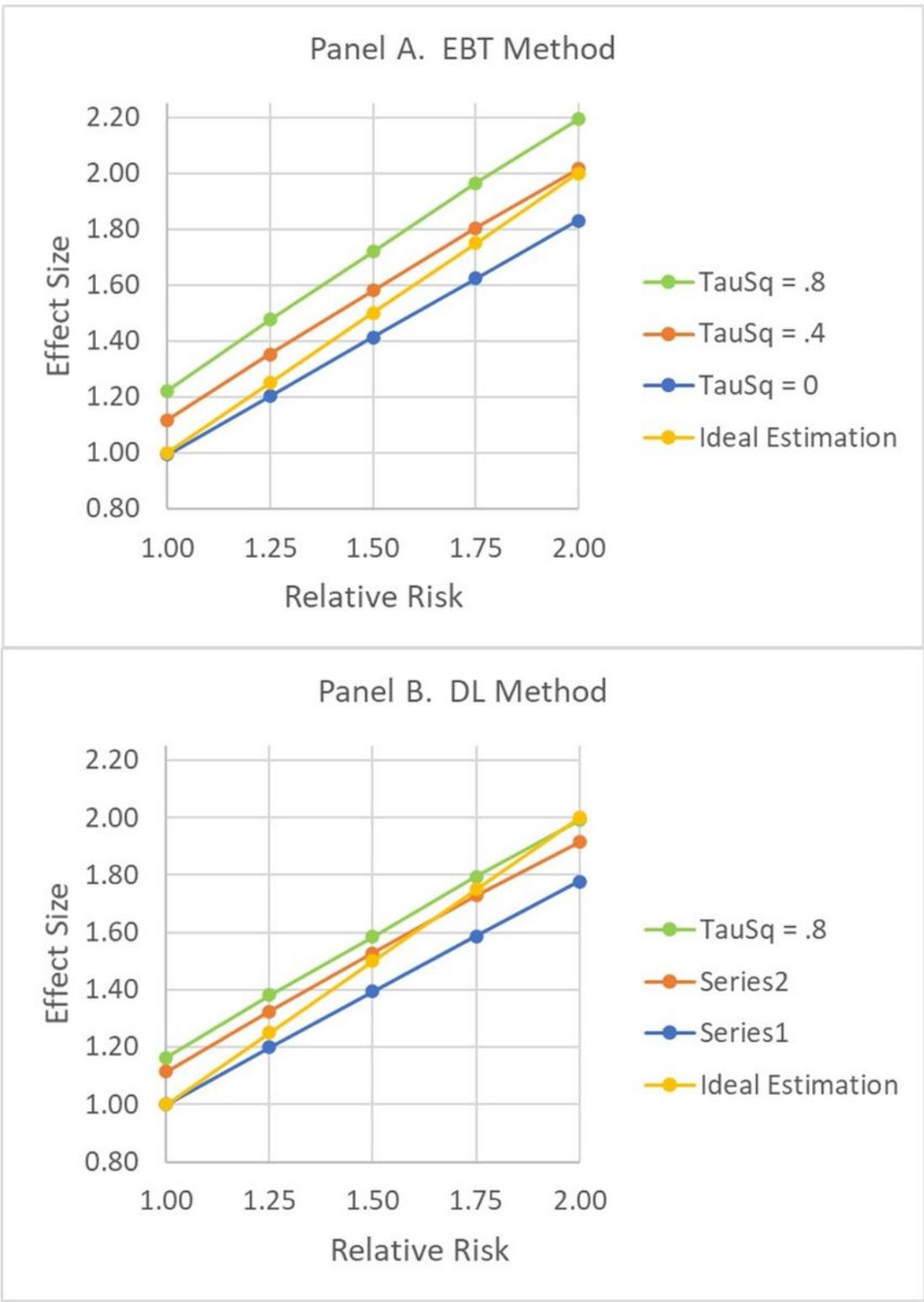


Figure 8

Effect Size as function of Relative Risk and heterogeneity. A & B correspond to the EBT and DL methods.

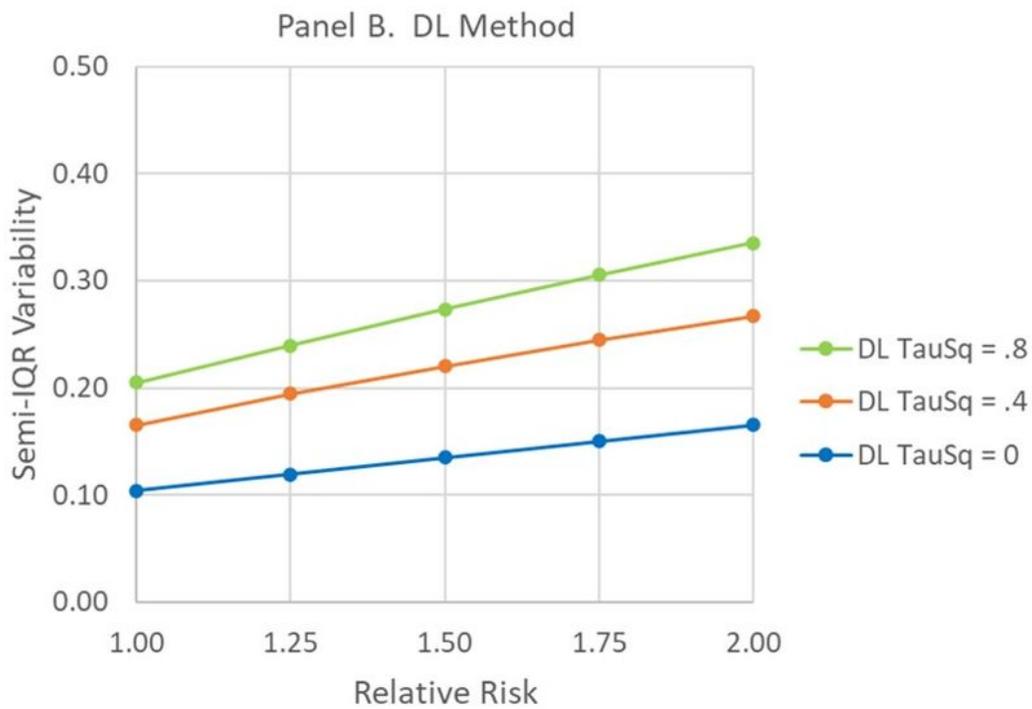
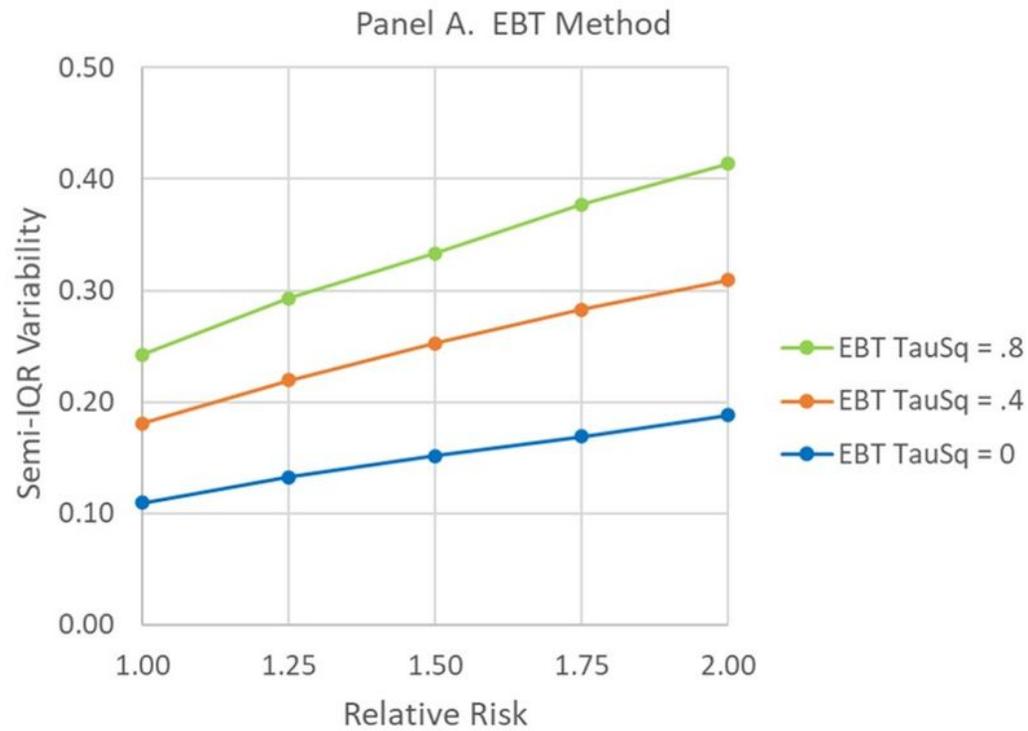


Figure 9

Semi-Interquartile Range as function of Relative Risk and heterogeneity. A & B correspond to the EBT and DL methods.