

# Contemporary frequentist views of the $2 \times 2$ binomial trial

Enrico Ripamonti, Chris Lloyd, Piero Quatto

University of Milan-Bicocca and University of Melbourne

*Abstract.* The  $2 \times 2$  table is the simplest of data structures yet it is of immense practical importance. It is also just complex enough to provide a theoretical testing ground for general frequentist methods. Yet after 70 years of debate, its correct analysis is still not settled. Rather than recount the entire history, our review is motivated by contemporary developments in likelihood and testing theory as well as computational advances. We will look at both conditional and unconditional tests. Within the conditional framework, we explain the relationship of Fisher's test with variants such as mid- $p$  and Liebermeister's test, as well as modern developments in likelihood theory, such as  $p^*$  and approximate conditioning. Within an unconditional framework, we consider four modern methods of correcting approximate tests to properly control size by accounting for the unknown value of the nuisance parameter: maximisation (M), partial maximisation (B), estimation (E), and estimation followed by maximisation (E+M). Under the conditional model, we recommend Fisher's test. For the unconditional model, amongst standard approximate methods Liebermeister's tests comes closest to controlling size. However, our best recommendation is the E procedure applied to the signed root likelihood statistic, as this performs very well in terms of size and power and is easily computed. We support our assertions with a numerical study.

*Key words and phrases:* approximate conditioning, binomial trial, conditional test, exact tests, Fisher test, Liebermeister test, mid- $p$  test, parametric bootstrap, unconditional test.

## 1. INTRODUCTION

Testing for a treatment effect in the  $2 \times 2$  binomial trial is a seminal topic in Statistics, founded on the original contributions of renowned statisticians Karl Pearson, Jerzy Neyman and Ronald Fisher. Not only is the design of great practical importance, notably in clinical trials, but the admittedly simple model is still rich enough to expose some of the tensions and limitations of frequentist statistics. Consequently, there are now literally dozens of test procedures that have been proposed. Reviews have appeared at regular intervals, see [Gart \(1969\)](#), [Yates \(1984\)](#), [Martin Andres \(1991\)](#), [Agresti \(1992\)](#), [Agresti \(2001\)](#) and [Lydersen et al \(2009\)](#). Why is it worth discussing tests for  $2 \times 2$  tables at all and why again now?

Firstly, the  $2 \times 2$  table is the basic data structure in clinical trials with binary outcome. Modern adaptive designs allow changes to treatments, sample sizes and

---

(e-mail: [enrico.ripamonti@unimib.it](mailto:enrico.ripamonti@unimib.it)). 1

40 even hypotheses, but their analysis relies on combining the evidence from different  
 41 arms, stages and hypotheses. The  $p$ -values from each  $2 \times 2$  table are fed into a more  
 42 complex  $p$ -value (using combination functions, multiple comparison adjustments  
 43 and the closed testing principle) whose statistical properties are inherited from  
 44 the component  $p$ -values. So while the  $2 \times 2$  table might appear a toy example it  
 45 is the building block of various modern methods/designs.

46 Secondly, over the past ten years, there has been considerable progress in the  
 47 foundational theory of exact or almost exact frequentist inference as well as meth-  
 48 ods of implementation. Some of these methods require complex computations.  
 49 Amongst unconditional methods there are several ways of correcting an approx-  
 50 imate test to be exact, see [Lloyd \(2008b\)](#). For conditional methods, the seminal  
 51 test of [Fisher \(1935\)](#) is limited by discreteness which can be mitigated by the  
 52 well known mid- $p$  correction ([Lancaster, 1961](#)). Further developments include so-  
 53 called approximative conditioning ([Pierce and Peters, 1992, 1999](#)) as well as the  
 54 famous  $p^*$  formula ([Barndorff-Nielsen, 1983](#)).

55 The objective of this paper is to present a detailed discussion of tests based  
 56 on the  $2 \times 2$  binomial trial with an emphasis on more contemporary theories  
 57 and proposals. We include unconditional and conditional perspectives without  
 58 taking a definite position on which is better. We do not study Bayesian methods,  
 59 even though one of the methods we include ([Liebermeister, 1877](#)) was originally  
 60 motivated from a Bayesian approach. Our overall aim is to place the different  
 61 methods within a coherent framework, to assess and compare their properties,  
 62 both theoretically and numerically and to arrive at clear recommendations.

63 In assessing the tests, we focus on four main criteria. First, we require the test  
 64 to be based on a  $p$ -value which measures, possibly approximately, the probability  
 65 of some observed event. Second, does the test exaggerate the evidence against the  
 66 null? This is based on comparing the nominal size with the actual *size profile* but  
 67 also the quoted  $p$ -value with its true *profile* (see [Lloyd, 2008a](#)). Third, we look  
 68 at the extent to which the test under-estimates the evidence against the null,  
 69 commonly called conservatism. Conservative tests tend to have lower power and  
 70 we confirm this with a numerical study. Lastly, we impose certain natural mono-  
 71 tonicity constraints ([Barnard, 1947](#); [Röhmel and Mansmann, 1999](#); [Skipka et al., 2004](#))  
 72 on the test statistics. These constraints also have favourable computational  
 73 implications.

74 The plan of the article is as follows. In the next section we establish the ba-  
 75 sic model, the notation and the theoretical framework for assessing the different  
 76 methods. In section 3 we present the conditional approach and in section 4 the  
 77 unconditional approach, in both cases emphasising modern perspectives and de-  
 78 velopments. In section 5, we assess conditional and unconditional tests within  
 79 their own frameworks and in section 6 we report the results of a numerical study  
 80 on the size and power of 28 different unconditional tests, where some clear con-  
 81 clusions do emerge. Our final recommendations are articulated in section 7.

## 2. THEORY RELEVANT TO EXACT TESTS

### 82 2.1 Notation

83 We suppose that  $n_0$  patients are given a comparison treatment,  $y_0$  of whom  
 84 respond positively with probability  $p_0$ , and  $n_1$  are given a new treatment,  $y_1$  of

85 whom respond positively with probability  $p_1$ . We henceforth call a positive re-  
 86 sponse a success. Provided patients respond independently, we have the standard  
 87 binomial model:

$$(2.1) \quad Y_0 \sim Bi(n_0, p_0) \quad \text{and} \quad Y_1 \sim Bi(n_1, p_1) \quad \text{with} \quad Y_0 \perp\!\!\!\perp Y_1$$

88 We mainly focus on one-sided hypotheses

$$(2.2) \quad H_0 : p_1 \leq p_0 \quad \text{vs.} \quad H_1 : p_1 > p_0$$

89 though the theory in this section applies to two-sided tests without modification.  
 90 We denote the total number of successes by  $S = Y_0 + Y_1$  and the proportion of  
 91 successes under treatment and control as  $\hat{p}_1 = y_1/n_1$  and  $\hat{p}_0 = y_0/n_0$ .

92 For the sake of giving general definitions and results, we will refer to the data  
 93  $(Y_0, Y_1)$  as  $Y$ , taking values in a sample space  $\mathcal{Y}$  and the parameter as  $\omega = (\theta, \varphi)$ ,  
 94 where  $\theta$  is the interest parameter and  $\varphi$  a nuisance parameter vector. We wish to  
 95 test the null hypothesis that  $\theta \in \Theta_0$  without specifying the value of the nuisance  
 96 parameter  $\varphi$ . For the binomial trial,  $\theta$  can be taken as any contrast of  $p_1$  and  $p_0$   
 97 (such as the difference, or the log-odds ratio), the nuisance parameter is  $\varphi = p_0$   
 98 and for the hypotheses in (2.2) the null parameter space is  $\Theta_0 = \{\theta : \theta \leq 0\}$ .

## 99 2.2 Size and power

100 All tests can be expressed in the form *reject the null if  $P(Y)$  is less than or*  
 101 *equal to  $\alpha$* , where  $\alpha$  is the nominal size of the test and  $P(Y)$  is called a  $p$ -value.  
 102 The probability of rejecting the null hypothesis is

$$(2.3) \quad \beta(\theta, \varphi) := \Pr[P(Y) \leq \alpha | \theta, \varphi]$$

103 The size of the test is typically defined as  $a(\varphi) = \sup_{\theta \in \Theta_0} \beta(\theta, \varphi)$  and the test is  
 104 *valid* if  $a(\varphi)$  is less than  $\alpha$  for all  $\varphi$  (Lehmann, 1959). The power of the test is  
 105 the probability of rejecting the null when  $\theta \notin \Theta_0$  and is desired to be as large as  
 106 possible, subject to validity.

## 107 2.3 What is an exact test?

108 Ideally we would want the size to equal  $\alpha$  but for discrete models  $a(\varphi)$  is a  
 109 polynomial and can never equal a constant. If we further maximise with respect  
 110 to  $\varphi$  then a bound can be given but again, because of discreteness, this bound  
 111 almost never equals  $\alpha$  exactly. In summary, if we define an exact test to have  
 112 “exact size  $\alpha$ ” then such tests almost never exist for discrete models.

113 For this reason, Lloyd (2008a) instead looks at the  $p$ -value, specifically at the  
 114 so-called *profile* of a  $p$ -value which is defined as

$$(2.4) \quad \pi(y, \varphi) = \sup_{\theta \in \Theta_0} \Pr[P(Y) \leq P(y); \theta, \varphi]$$

115 The putative property of a  $p$ -value is that an observed value of say 0.042 means  
 116 that something unusual has happened and the probability of it happening under  
 117 the null is 0.042. Thus we would like  $\pi(y, \varphi)$  to equal  $P(y)$  for all  $\varphi$ . Again,  
 118 because of discreteness this is impossible, so again it appears as if no exact  $p$ -  
 119 value can exist. However, if  $\sup_{\varphi} \pi(y, \varphi) \leq P(y)$  for all  $y$  in  $\mathcal{Y}$  we call the  $p$ -value  
 120 *guaranteed*. This is identical to  $P(Y)$  being stochastically no smaller than uniform

121 for all  $\theta \in \Theta_0$ . It is the analog of a test being valid and a guaranteed  $p$ -value does  
 122 imply a valid test. However, the advantage of basing the theory on  $p$ -values rather  
 123 than test size is that there always exists a  $p$ -value for which

$$(2.5) \quad \sup_{\varphi} \pi(y; \varphi) = P(y) \quad \forall y \in \mathcal{Y}$$

124 as proven by [Röhmel and Mansmann \(1999\)](#). Such a  $p$ -value is called *exact*. It is  
 125 further shown that amongst  $p$ -values that impose the same ordering on the sample  
 126 space there always exists a smallest  $p$ -value and that this  $p$ -value is exact. The  
 127 construction of this  $p$ -value is simple and will be given in section 4.1. The theory  
 128 is completely general and applies to the conditional or unconditional model, as  
 129 well as to one-sided or two-sided tests.

## 130 2.4 Most powerful tests

131 [Lehmann \(1959\)](#) established the existence of both exact and optimal tests,  
 132 which is relevant to our purposes. The main class of models where a most powerful  
 133 test exists is the natural exponential family, where the joint density or probability  
 134 function of the data  $Y$  can be written as:

$$(2.6) \quad f_{\theta, \varphi}(y) = \exp\{\theta T(y) + S'(y)\varphi + \varsigma(\theta, \varphi)\}$$

135 where  $T(y)$  is a scalar and  $S(y)$  is the sufficient statistic for  $\varphi$ . For model 2.6,  
 136 uniformly most powerful unbiased (UMPU) tests exist for both one and two-sided  
 137 alternatives. These procedures are based on tail probabilities of the conditional  
 138 distribution of  $T$  given  $S$  but their UMPU properties are also unconditional.

139 For the binomial trial, the model is of exponential form with  $\theta = \text{logit}(p_1) -$   
 140  $\text{logit}(p_0)$ , the statistic  $T(Y) = Y_1$  and its distribution given  $S = Y_0 + Y_1 = s$  is

$$(2.7) \quad \Pr(Y_1 = y_1; s, \theta) = e^{\theta y_1} \binom{n_1}{y_1} \binom{n_0}{s - y_1} / \kappa(\theta, s)$$

141 where  $\max\{0; s - n_0\} \leq y_1 \leq \min\{s; n_1\}$  and  $\kappa(\theta, s)$  is a normalising constant.  
 142 When  $\theta = 0$ , the distribution is hypergeometric. However, this model is discrete.

143 For discrete exponential models, Lehmann's UMPU tests involve randomisation.  
 144 This also arises in certain non-exponential continuous models, such as uni-  
 145 form when the support depends on the parameter value. In any case, randomi-  
 146 sation is never used in practice. The lack of an optimal test explains the many  
 147 alternative tests of the  $2 \times 2$  table that have been proposed in the literature. Nev-  
 148 ertheless, a key insight that comes from the theory is that optimal tests should  
 149 be based on the conditional distribution of  $T$  given  $S$  and that for fixed  $S$  there  
 150 is more evidence against the null hypothesis when  $T$  is larger.

## 151 2.5 Monotonicity properties

152 For some models, there are basic logical properties that we expect any pro-  
 153 cedure to have. For instance, any statistical procedure for assessing reliability  
 154 should produce a less favorable assessment if you add any errors to the dataset.  
 155 Such conditions can be expressed mathematically (see [Harris and Som, 1991](#) and  
 156 [Kabaila, 2005](#)) and come down to test statistics having certain monotonicity  
 157 properties. For the  $2 \times 2$  table, the evidence for  $p_1 > p_0$  is stronger if  $Y_1$  is larger

158 for fixed  $Y_0$  and if  $Y_0$  is smaller for fixed  $Y_1$ . Equivalently, the  $p$ -value should be  
 159 non-increasing in  $Y_1$  for fixed  $S = Y_0 + Y_1$  and non-decreasing in  $S$  for fixed  $Y_1$ ,  
 160 as noted by Berger and Sidik (2003). While the condition may appear obvious,  
 161 standard approaches, such as the standard Z-test can violate it. Likelihood ratio  
 162 tests typically satisfy any required monotonicity conditions.

163 These monotonicity properties have two important consequences. First, the  
 164 maximum probability  $a(\varphi) = \sup_{\theta \in \Theta_0} \beta(\theta, \varphi)$  of rejecting the null is achieved  
 165 at the boundary point  $\theta = \theta_0$  (Röhmel, 2005). This not only ensures that the  
 166 test is unbiased, but facilitates computations. Firstly, there is no need to search  
 167 over  $\theta$ . Secondly, the tail set  $\{P(T, S) \leq P(t, s)\}$  can be simply determined using  
 168 the fact that  $P(T, S)$  is a non-increasing function of  $T$  for fixed  $S$ , as noted by  
 169 Finner and Strassburger (2002).

170 The above conditions mention non-decreasing rather than strictly increasing.  
 171 What about ties? For discrete data, it is never advantageous to have ties. It was  
 172 shown by Röhmel and Mansmann (1999) that if a guaranteed  $p$ -value has any  
 173 ties then breaking the ties appropriately can often make the  $p$ -value smaller while  
 174 still being guaranteed. Similar results for confidence limits were demonstrated in  
 175 Kabaila and Lloyd (2006).

## 176 2.6 Criteria for comparison

177 There are several different criteria that can be used to assess the effectiveness  
 178 of a test. If prior information summarised as a distribution is available on the  
 179 unknown parameters then an exact Bayesian solution immediately follows. The  
 180 frequentist properties of the Bayes tests are rarely poor, but neither are they  
 181 exact. In decision theoretic approaches, various loss functions can be defined and  
 182 minimised within a specified space of decision functions.

183 Even within the pure frequentist paradigm that we assume here, there is no  
 184 non-randomised test with maximum power and controlled size for discrete models.  
 185 It seems unsatisfactory that this paradigm does not support an optimal analysis  
 186 for a simple data structure like the  $2 \times 2$  table. However, based on the four criteria  
 187 to be listed below, we will find that there is indeed a practically optimal approach.

188 We now state four criteria that we will use to assess the different tests. The  
 189 first two relate to their statistical accuracy i.e. to the test size and power. These  
 190 two descriptors are central to frequentist theory and also to all trial regulation  
 191 authorities. Tests should firstly be *valid*, or equivalently the  $p$ -value should be  
 192 guaranteed. Gross violations of the size restriction is a serious defect of any test  
 193 in our review. Ideally, the  $p$ -value should also be exact. This means that the test  
 194 does not under-estimate the evidence against the null and will tend to lead to  
 195 higher power. Restricting attention to valid tests means that the power achieved  
 196 by different tests can be compared, without the complicating possibility that any  
 197 extra power is purchased by size violations.

198 The other two criteria are more foundational. Tests should be based on a  $p$ -  
 199 value that measures the probability of an observed event. This not only leads  
 200 to a transparent test decision but provides quantitative information about how  
 201 unusual the data is under the null hypothesis. Finally, where the model supports  
 202 it on logical grounds, tests should satisfy certain monotonicity conditions. For  
 203 the  $2 \times 2$  table, these conditions were listed in the previous section.

### 3. MODERN PERSPECTIVES ON CONDITIONAL TESTS

204 It was argued by Fisher (1935) that the number of successes  $S$  should be  
 205 treated as fixed, see Choi et al. (2015) for an overview and historical perspective.  
 206 We evaluate the merits of this key modelling decision in section 5.

#### 3.1 Fisher's exact test

208 If we treat  $S = s$  as fixed then the model is given by (2.7). The distribution is  
 209 stochastically increasing in  $\theta$  and so we reject  $H_0 : \theta \leq 0$  for larger values of  $y_1$ ,  
 210 and the  $p$ -value is  $\Pr[Y_1 \geq y_1 | S = s]$  calculated from (2.7) maximised over  $\theta \leq 0$ .  
 211 Because of stochastic monotonicity, the maximum occurs when  $\theta = 0$ . Fisher's  
 212  $p$ -value  $P_F(y_1; n_1, n_0, s)$  is this tail sum of hypergeometric probabilities:

$$(3.1) \quad P_F(y_1; n_1, n_0, s) = \sum_{y \geq y_1} \binom{n_1}{y} \binom{n_0}{s-y} / \binom{n_0 + n_1}{s}.$$

213 The test is exact, in the sense that no approximation or estimation of unknown  
 214 parameters is involved. Fisher's  $p$ -value is also exact in the technical sense of  
 215 equation 2.5, assuming the model for  $Y_1$  given  $s$ . The test generated by this  
 216  $p$ -value is therefore valid within this same conditional model.

217 The size of the test can be calculated exactly, since the hypergeometric distri-  
 218 bution has no unknown parameters. For given values of  $n_0, n_1, s$  and target size  
 219  $\alpha$ , let  $c_s$  be the smallest integer value  $c$  such that  $P_F(c; n_1, n_0, s) \leq \alpha$ . So the test  
 220 rejects the null exactly when  $y_1 \geq c_s$ . It follows that the size of Fisher's test is

$$(3.2) \quad \alpha_s = \sup_{\theta \leq 0} \Pr(Y_1 \geq c_s | s) = P_F(c_s; n_1, n_0, s)$$

221 In words, the true size  $\alpha_s$  is equal to the largest observable  $p$ -value less than  $\alpha$ . So  
 222 the test is not exact in the sense of having exactly the correct size. The smaller  
 223 the support of the distribution the less likely it is that  $\alpha_s$  will be close to the  
 224 chosen  $\alpha$ . The support is smaller when the observed value of  $s$  is more extreme.  
 225 In the extreme cases where  $s = 0$  or  $s = n$ , the conditional test never rejects the  
 226 null and the true size is  $\alpha_s = 0$ .

#### 3.2 Randomised version of Fisher's exact test

227 Based on the earlier mentioned theory of Lehmann (1959), there is a ran-  
 228 domised version of Fisher's exact test which is UMPU for the one-sided test (see  
 229 also Tocher, 1950) and for the two-sided alternative. For the one-sided alternative,  
 230 this comes down to using the randomised  $p$ -value  
 231

$$(3.3) \quad P_R(y_1, U; n_0, n_1, s) = P_F(y_1; n_1, n_0, s) - U \Pr(Y_1 = y_1; n_0, n_1, s)$$

232 where  $U$  is a uniformly distributed random number in the interval  $(0, 1)$  (e.g.,  
 233 Cox and Hinkley, 1974, p.101). At the expense of introducing the random number  
 234  $U$  into the inference, we obtain a  $p$ -value with exact uniform distribution and a  
 235 test with exact size  $\alpha$ . Apparently, this  $p$ -value is always smaller than  $P_F$  and  
 236 so is less conservative. In fact, the test is UMP amongst unbiased tests that are  
 237 functions of  $(T, S, U)$  (Lehmann, 1959) which suggests that the shortcomings of  
 238 Fisher's test are all due to discreteness. Of course, randomisation is almost never  
 239 used in practice because we feel that conclusions should not depend on the random



240 number  $u$ . If one takes the data to be  $(T, S, U)$  then the sufficiency principle  
 241 states that inference should not depend on  $U$ . The conditionality principle would  
 242 also recommend conditioning on the value of  $U$  which, since  $U$  is independent  
 243 of  $(T, S)$ , again means just using  $(T, S)$ . There are alternative decision theoretic  
 244 perspectives where the inference can be a distribution and randomisation is used  
 245 to generate from this distribution. In this approach,  $U$  is not considered part of  
 246 the data. However, this paper takes a frequentist approach.

### 247 3.3 Lancaster's and Liebermeister's $p$ -value

248 [Lancaster \(1961\)](#) proposed an alternative solution to the problem of conser-  
 249 vatism of any discrete test, which has seen a fair degree of application. Lancaster's  
 250 mid- $p$ -value only counts half of the observed null probability of the observed sam-  
 251 ple point in the tail probability. Equivalently, it is obtained by subtracting half the  
 252 observed probability from the usual tail probability. Referring to (3.3), the mid  
 253  $p$ -value  $P_{\text{mid}}(Y_1; n_0, n_1, s)$  is given explicitly by replacing  $U$  by its mean value of  
 254 0.5. While not uniformly distributed like the randomised  $p$ -value, it has the exact  
 255 mean (0.5) and variance of a uniform distribution ([Agresti, 2002](#)). Stronger the-  
 256 oretical justification for the one-sided mid- $p$  are provided by [Hwang et al. \(1991\)](#)  
 257 and recently by [Wells \(2010\)](#).

258 Another test closely related to Fisher's was proposed by [Liebermeister \(1877\)](#).  
 259 It is based on a Bayesian argument and turns out to equal Fisher's  $p$ -value but  
 260 with a fictitious success added to the treatment group and a failure to the control  
 261 group so it can be expressed as  $P_F(y_1 + 1; n_0 + 1, n_1 + 1, s + 1)$ . It was shown  
 262 by [Seneta and Phipps \(2001\)](#) that it is always between  $P_F(y_1 + 1; n_0, n_1, s)$  and  
 263  $P_F(y_1; n_0, n_1, s)$ , though not necessarily half way between. Like Lancaster's  $p$ -  
 264 value, tests based on Liebermeister's  $p$ -value are less conservative than those  
 265 based on Fisher's  $p$ -value.

### 266 3.4 Modern approximations

267 During the 1980s, new developments in likelihood theory led to the proposal of  
 268 the  $p^*$  formula by [Barndorff-Nielsen \(1983\)](#). The theory is complex but is based  
 269 on a saddlepoint approximation to the density of the maximum likelihood esti-  
 270 mator, conditional on a very generally formulated approximate ancillary statistic.  
 271 Suppose we want to test a null hypothesis that the parameter  $\delta = p_1 - p_0$  is less  
 272 than or equal to  $\delta_0$ . Until this point, the null value  $\delta_0$  has been zero. A general  
 273 form for the  $p^*$  test statistic is

$$(3.4) \quad r^*(\delta_0) = r(\delta_0) + r(\delta_0)^{-1} \log \left( \frac{q(\delta_0)}{r(\delta_0)} \right)$$

274 where  $r(\delta_0)$  is the signed root likelihood ratio statistic for testing  $\delta \leq \delta_0$  and  $q(\delta_0)$   
 275 is very complex in its general formulation but for the  $2 \times 2$  table reduces to

$$(3.5) \quad q(\delta_0) = \frac{\{\tilde{w}_0(\text{logit}(\tilde{p}_1) - \text{logit}(\hat{p}_1)) - \tilde{w}_1(\text{logit}(\tilde{p}_0) - \text{logit}(\hat{p}_0))\}}{\sqrt{\tilde{w}_1/n_1 + \tilde{w}_0/n_0}}$$

276 where  $\tilde{w}_j = \tilde{p}_j(1 - \tilde{p}_j)$  and  $\tilde{p}_j$  is the ML estimate of  $p_j$  under the null, as shown  
 277 in [Lloyd \(2010b\)](#). The corresponding  $p$ -value is denoted  $p^*(\delta_0) = 1 - \Phi(r^*(\delta_0))$ .  
 278 An advantage of this approach is that it is available in closed form. The normal  
 279 approximation is held to be accurate to  $O(n^{-1})$  in the medium deviation range  
 280 ([Davison et al, 2006](#)).

281 The appeal of  $p^*$  is that it depends continuously on the null value  $\delta_0$ . Amongst  
 282 other consequences, this means we can invert the test to get a confidence interval  
 283 for  $\delta$ . In contrast, Fisher’s method only works for testing  $\delta_0 = 0$ , since no condi-  
 284 tional distribution free of unknown nuisance parameters exists as  $\delta_0$  moves away  
 285 from 0. The  $p^*$  method gives an answer close to the exact conditional solution  
 286 when one exists but generalises, albeit approximately, to models and hypotheses  
 287 where no exact conditional inference is possible. The  $p$ -value based on  $r^*$  is an  
 288 approximation to the probability of a well-defined event, unlike the Lancaster or  
 289 Liebermeister  $p$ -values.

290 The approach does present several problems however. First, the formula for  $r^*$   
 291 breaks down when either  $r = 0$  or  $q = 0$ . So to properly investigate its exact fre-  
 292 quentist properties, it must be redefined. We define  $r^* = r$  whenever the absolute  
 293 value of  $r$  is less than 0.1 or when  $q = 0$ . These problems are completely ignored  
 294 in the literature.

295 Another lesser problem is that, even with these modifications,  $p^*$  is not nec-  
 296 essarily guaranteed (which is a fundamental criterion in our review) and further  
 297 numerical work is required to evaluate its degree of liberalism. Second, for some  
 298 quite natural models such as logistic regression with interest on the intercept, the  
 299 conditional  $p$ -value becomes degenerate, even though  $p^*$  does not. In this case,  
 300 what is the relation between the conditional degenerate  $p$ -value and  $p^*$  which is  
 301 supposed to approximate it? According to [Pierce and Peters \(1999\)](#), in such cases  
 302  $p^*$  is an “approximately conditional”  $p$ -value.

### 303 3.5 Approximately conditional $p$ -values

304 One novel proposal to mitigate conservatism is to use a less discrete conditional  
 305 distribution by conditioning on a range of values for the conditioning statistic  
 306 rather than the exact value. This leads to a distribution with finer support. On  
 307 the other hand, the nuisance parameter is no longer eliminated. Consider the  
 308  $p$ -value  $P(t, s) = \Pr[T \geq t | S = s]$  calculated under the null. In the current  
 309 context this would equal Fisher’s  $p$ -value. Define a neighbourhood  $N_r(s)$  around  
 310 the observed value of  $S = s$ , for example  $\{s - r, \dots, s + r\}$ . Then an approximately  
 311 conditional  $p$ -value is defined as

$$(3.6) \quad \Pr(P(T, S) \leq P(t, s) | S \in N_r(s_{obs})) = \sum_{s \in N_r(s_{obs})} \Pr[P(T, s) \leq p_{obs} | S = s] \Pr[S = s | s \in N_r(s_{obs}); \varphi]$$

312 When the size  $r$  of the neighbourhood equals zero, this gives the conditional  
 313  $p$ -value  $P(t, s)$  since there is only one term in the sum. When  $r > 0$  it is ap-  
 314 proximately conditional. Residual dependence on the nuisance parameter could  
 315 in principle be handled by any of the methods that we will explain in section 4.1  
 316 below.

317 The main problem is a lack of recommendation for the size of the neighbour-  
 318 hood  $N_r(s)$  as well as its shape when  $S$  is higher dimensional. Certainly a different  
 319 choice of the neighborhood leads to a different  $p$ -value. A second problem is that  
 320 the  $p$ -value still depends on the nuisance parameter  $\varphi$ . A third logical problem is  
 321 a phenomenon known as spurious deflation, see [Lloyd \(2010a\)](#). It is too early to  
 322 dismiss approximately conditional  $p$ -values though theoretical problems remain.



323 They are at least based on the probability of a well-defined event and can be  
 324 guaranteed by maximising with respect to the nuisance parameter.

### 325 3.6 Unconditional assessment of conditional tests

326 The tests just described are based on the distribution of  $T(Y)$  given  $S(Y)$ .  
 327 When  $S$  really is fixed by design, it seems pertinent to assess the size and power  
 328 treating it as fixed. When  $S$  is not fixed by design, it is still sometimes argued  
 329 that conditional assessment is appropriate. Certainly though, in future hypothet-  
 330 ical repetitions the value of  $S$  will vary and to allow for this we have to use  
 331 the unconditional model. So both conditional and unconditional assessment have  
 332 plausible arguments in their favour. But how are the two approaches related?

333 The conditional probability of rejection we will denote by

$$(3.7) \quad \beta(\theta|s) := \Pr[P(T, s) \leq \alpha | \theta, S(Y) = s]$$

334 where we have used the fact that the distribution of  $T$  given  $S(Y) = s$  does  
 335 not depend on  $\varphi$ . With slight abuse of notation, the unconditional probability of  
 336 rejection is the mean value

$$(3.8) \quad \beta(\theta, \varphi) = \sum_s \beta(\theta|s) \Pr(S = s; \theta, \varphi)$$

337 of  $\beta(\theta|S)$  with respect to the distribution of  $S$  which depends again on  $(\theta, \varphi)$ . For  
 338 the  $2 \times 2$  table,  $S = Y_0 + Y_1$  has the distribution of a sum of two binomials and  
 339 the summation is from  $s = 0, \dots, n_0 + n_1$ .

340 When  $\theta = \theta_0$ ,  $\beta(\theta_0|s)$  is the conditional size which we earlier denoted  $\alpha_s$ . The  
 341 unconditional size is the mean value of  $\alpha_s$  with respect to  $S$ . For the Fisher test,  
 342  $\alpha_s$  is almost always strictly less than  $\alpha$  for all values of  $s$  and so unconditional  
 343 size is also less than  $\alpha$ . So Fisher's test is unconditionally conservative by design.  
 344 For the Lancaster or Liebermeister test, their conditional size  $\alpha_s$  is typically less  
 345 than  $\alpha$  for some values of  $s$  and larger than  $\alpha$  for others. The unconditional size  
 346 is the mean value which is typically quite close to  $\alpha$  because of the averaging,  
 347 though it can exceed  $\alpha$ .

## 4. MODERN PERSPECTIVES ON UNCONDITIONAL TESTS

348 For testing the null hypothesis  $\theta \leq 0$  against  $\theta > 0$  there are several commonly  
 349 used test statistics based on the unconditional joint binomial model.

350 From a historical point of view, the best known is the chi-squared statistic  
 351 (Pearson, 1900). Alternatively, tests can be based on the difference  $\hat{p}_1 - \hat{p}_0$  di-  
 352 vided by a standard error. When this standard error is estimated under the null  
 353 hypothesis (i.e. assuming  $p_1 = p_0$ ) it gives rise to the so-called pooled statistic.  
 354 This can be shown to be a particular case of the Rao's score statistic and is iden-  
 355 tical to a one-sided version of the chi-square statistic. When the standard error is  
 356 estimated without any restrictions on  $p_1$  and  $p_0$  it is called the unpooled statistic,  
 357 which is the Wald statistic based on the interest parameter  $\theta = p_1 - p_0$ . There  
 358 are other Wald-type statistics that can also be used, for instance based on the  
 359 difference between the logarithm or the logit of the estimated success rates  $\hat{p}_1$  and  
 360  $\hat{p}_0$  but they are not in common use. The main alternative to these statistics is the  
 361 likelihood ratio test, or its one-sided version known as the signed root likelihood  
 362 ratio (SRLR) test. Formulas for these well known statistics are in Appendix 1.

363 The problem is that none of these tests are exact. Suppose we start with an  
 364 approximate  $p$ -value  $P(Y)$  based on test statistic  $Z(Y)$ . We remind the reader  
 365 of the definition of the profile  $\pi(y, \varphi)$  of a  $p$ -value  $P(Y)$  given in equation 2.4.  
 366 In words, it is the null probability of the  $p$ -value being equal or smaller than its  
 367 observed value  $P(y)$ , the *true significance* if you will. Ideally, it should equal  $P(y)$ .  
 368 It is worth noting that the tail set  $\{P(Y) \leq P(y)\}$  in the definition of  $\pi(y, \varphi)$  can  
 369 equally be expressed as  $\{Z(Y) \geq Z(y)\}$  and it is partly a matter of taste how the  
 370 formulas below are presented.

371 There are several methods of using the profile function to define either an exact  
 372 or almost exact version of the original  $p$ -value  $P(Y)$ . These ideas are mostly quite  
 373 recent and can be implemented with modern computational resources.

#### 374 4.1 The maximisation procedure

375 The “worst case”  $p$ -value is  $P^*(y) = \sup_{\varphi} \pi(y, \varphi)$ . While this is a completely  
 376 general method, with general optimality properties stated below, the seminal  
 377 paper recommending maximising out nuisance parameters was [Barnard \(1945\)](#).  
 378 For the  $2 \times 2$  table, we take the nuisance parameter  $\varphi$  to be the common value  
 379 of  $p_1 = p_0$  under the null, and this becomes

$$(4.1) \quad P^*(y_1, y_0) = \sup_{0 \leq p \leq 1} \Pr[Z(Y_1, Y_0) \geq Z(y_1, y_0); p_0 = p_1 = p]$$

380 computed by enumerating all pairs  $(y_1, y_0)$  in the tail set  $\{Z(Y_1, Y_0) \geq Z(y_1, y_0)\}$ ,  
 381 summing their null probabilities based on the independent binomial distribution  
 382 and then maximising with respect to  $p$ . We will call this adjustment the M-step.  
 383 The maximised  $p$ -value has the following incredibly strong optimality property:  
 384 amongst all statistics that are non-increasing functions of  $Z(Y)$ ,  $P^*(Y)$  is the  
 385 smallest function that is a guaranteed  $p$ -value. It is also exact in the sense of  
 386 (2.5). So, if a test is not expressible as an M  $p$ -value based on some test statistic,  
 387 then it can be improved by the M-step.

388 The test statistic  $Z(Y_1, Y_0)$  may be any of the three mentioned in the previ-  
 389 ous section. Since  $P^*(y_1, y_0)$  only depends on the way  $Z(y_1, y_0)$  ranks the sample  
 390 space, dependence on the choice of  $Z$  is modest, so long as it is chosen to be one  
 391 of the standard test statistics. Moreover, while slightly different answers can be  
 392 obtained, each  $p$ -value is exact in the sense of equation 2.5. The maximization  
 393 procedure can even be applied to any of the conditional  $p$ -values from the previ-  
 394 ous section, thus converting a conditional test into an exact unconditional test  
 395 ([Boschloo, 1970](#); [McDonald et al., 1977](#); [Mehrotra et al., 2003](#)).

396 We remind the reader that when  $Z(y_1, y_0)$  does not satisfy the monotonicity  
 397 properties, the tail probability in (4.1) should in principle be maximised over  
 398  $\{p_1 \leq p_0\}$ , as first pointed out by [Röhmel \(2005\)](#). Fisher’s  $p$ -value, as well as  
 399 Lancaster and Lieberman are monotonic. Among the three standard statistics,  
 400 only the SRLR statistic is necessarily monotonic.

#### 401 4.2 The restricted maximization procedure

Maximising the profile function over the entire nuisance parameter space seems  
 extreme when many nuisance parameter values will be very unlikely in light of  
 the data. This might lead to unnecessarily conservative inference. When this  
 occurs, it can be traced to the existence of spikes in the profile, often at values  
 of  $\varphi$  far from its estimated value. Such problems may be avoided by using the

procedure proposed by Berger and Boos (1994), which narrows the set of values in the domain of the parameter  $\varphi$  to a confidence set before taking the maximum:

$$P_{BB}(Y) = \sup_{\varphi \in C_\gamma} \Pr[Z(Y) \geq Z(y); \theta_0, \varphi] + \gamma$$

402 where  $C_\gamma$  is a  $100(1 - \gamma)$  percent confidence interval for  $\varphi$ . In the present case,  
 403 probability on the right hand side is calculated under  $p_1 = p_0 = p$  and a confidence  
 404 interval for  $\varphi = p$  is the well known Clopper-Pearson interval. Normally,  $\gamma$  is taken  
 405 to be very small, e.g., 0.001. Again, the tail set could be expressed in terms of the  
 406  $p$ -value instead of the test statistic. This restricted maximization, which we will  
 407 call the  $B$ -step, produces a guaranteed  $p$ -value and typically a smaller  $p$ -value  
 408 than using the  $M$ -step.

### 409 4.3 The estimation procedure

410 A cruder alternative to accounting for the nuisance parameter by maximization  
 411 is to replace it with an estimate (Storer and Kim, 1990). In its most general form  
 412 this gives what we will call the E  $p$ -value

$$(4.2) \quad P_E(y) = \pi(y, \hat{\varphi}_0)$$

413 where  $\hat{\varphi}_0$  is an estimator of  $\varphi$  under the null hypothesis. This is a parametric  
 414 bootstrap  $p$ -value if the bootstrap is viewed as a general recommendation to use  
 415 the data to estimate the null distribution of the test statistic. For simple models  
 416 like the  $2 \times 2$  table no simulation is required. The value of the E  $p$ -value is obtained  
 417 from equation (4.1) but, rather than maximise with respect to  $p$ , it is replaced by  
 418 the estimate  $\hat{p}$ . The main problem of this approach is that the resulting  $p$ -value  
 419 is not necessarily guaranteed (Berger and Boos, 1994).

420 The E-step can be performed more than once by iterating the construction of  
 421 the significance profile in 2.4. The three methods, M-step, B-step and E-step can  
 422 also be combined. Lloyd (2008a) proposed applying the M-step to  $P_E(y)$ , known  
 423 as the the E+M  $p$ -value. More explicit formulas for all these methods are given  
 424 in Appendix 2.

### 425 4.4 Numerical illustration

426 To clarify exactly how these three adjustments work, we illustrate their appli-  
 427 cation when  $Z(Y_1, Y_0)$  is the pooled z-test. The fictitious data is  $y_1 = 13$  responses  
 428 out of  $n_1 = 100$  and  $y_0 = 2$  responses out of  $n_0 = 50$ . The observed value of the  
 429 test statistic is  $Z(y_0, y_1) = 1.732$ . Practitioners would typically quote the ob-  
 430 served  $p$ -value  $1 - \Phi(1.732) = 0.0416$ . How accurate is this? Figure 1 displays the  
 431 true significance, as measured by the profile  $\pi(y = (2, 13); p)$ , with the quoted  
 432 value 0.0416 as a horizontal line. It deviates from the quoted value, mainly for  
 433 larger values of  $p$  but also for  $p = 0.5$ . The nuisance parameter  $\varphi$  here is again  
 434 the assumed common value  $p$  of  $p_1 = p_0$  under the null.

435 The maximum of the profile is  $P^* = 0.0677$  and occurs at  $p = 0.968$ . On the  
 436 basis of this M-step, we quote the  $p$ -value 0.0677 instead of the original 0.0416.  
 437 This value is much larger because of the presence of a spike in the profile but  
 438 considering that  $\hat{p} = 15/150 = 0.1$ , one might wonder about taking account of the  
 439 possibility that  $p = 0.968$ . This motivates the alternative B-step. With  $\gamma = 0.01$ ,  
 440 an exact 99% confidence interval for  $p$  is (0.0471, 0.1796) marked on the plot as

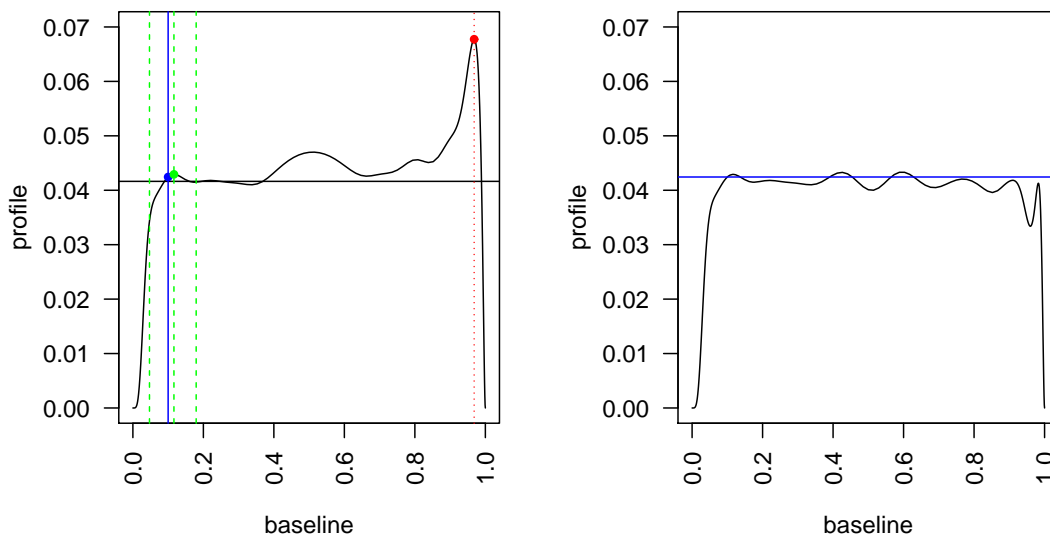


FIG 1. *Left. Profile for pooled  $Z$   $p$ -value for data ( $y_1 = 13, n_1 = 100; y_0 = 2, n_0 = 50$ ), illustrating the  $M$  (in red),  $B$  (in green) and  $E$ -step (in blue)  $p$ -values that this profile generates. Right. Profile of the  $E$   $p$ -value for the same data set.*

441 dashed green vertical lines. The maximum over this restricted range is 0.0429.  
 442 After adding the penalty  $\gamma = 0.01$  we quote  $P_{BB} = 0.0529$ . The  $E$ -step involves  
 443 estimating  $p$  by  $\hat{p} = 0.1$  and the value of the profile at this point, marked by a  
 444 vertical blue line, is  $P_E = 0.0424$ .

445 It was noted previously that estimated  $p$ -values may not be guaranteed. In the  
 446 right panel of Figure 1 we have calculated the profile of  $P_E(Y)$  (which required  
 447 calculating all possible values of  $P_E(y)$ ) as well as the quoted value as a horizontal  
 448 line. In such cases, the quoted  $p$ -value is extremely close to the true significance  
 449 profile. This behaviour is typical for  $E$   $p$ -values in this context (Lloyd, 2008b).  
 450 The  $E+M$   $P$ -value is the maximum of this profile and equals 0.0427, achieved at  
 451  $p = 0.42$ . The latter  $M$ -step removes the practically tiny amount of conservatism  
 452 or liberality that may be present, and the resulting  $p$ -value is exact.

#### 453 4.5 Two-sided and multi-dimensional tests

454 One-sided tests (such as superiority or non-inferiority trials) are very common  
 455 in biomedical contexts, which is why the theory presented to this point is oriented  
 456 towards one-sided tests. However, the tail set  $\{Z(Y) \geq Z(y)\}$  could be based on a  
 457 two-sided test statistic  $Z$  if desired. This is perhaps even clearer when the theory is  
 458 expressed in terms of the equivalent  $p$ -value, where the tail set is  $\{P(Y) \leq P(y)\}$ .

459 Also suppressed in the theory is the dimension of the nuisance parameter  $\varphi$ ,  
 460 which is unspecified. In principle then, the theoretical framework is completely  
 461 general. The  $M$ -step,  $B$ -step and  $E$ -steps are applied in exactly the same way for  
 462 one or two-sided tests and for any number of nuisance parameters. The  $M$ -step  
 463 retains the same optimality properties stated in section 4.1 and the  $B$ -step always  
 464 produces a guaranteed  $p$ -value. But both these methods become computationally

465 infeasible for many nuisance parameters. Only the E-step maintains the same  
 466 computational burden as the dimension of  $\varphi$  increases. In the context of  $2 \times 2$   
 467 tables all three methods are computable for realistic sample sizes.

## 5. STRUCTURED ASSESSMENT OF COMPETING TESTS

468 In this section we review the main arguments for and against the conditional  
 469 and unconditional model, without taking a position on which is better. We then  
 470 compare proposed tests within the conditional framework and come to a clear  
 471 recommendation. Within the unconditional framework, there are literally dozens  
 472 of plausible tests. We assess these by their theoretical properties as detailed in  
 473 section 2.6, as well as their computational burden. Moreover, we support our final  
 474 recommendation with a numerical study.

### 5.1 Conditional or unconditional?

475  
 476 All statistical models involve some conditioning; those things we consider in-  
 477 cidental to the data, for instance the sample size, do not have their distributions  
 478 modeled but, rather, are considered fixed. The dispute between the use of con-  
 479 ditional or unconditional tests has a long history and many of the battles have  
 480 been fought around the  $2 \times 2$  table (see [Agresti, 1992, 2001](#) for a review).

481 In Fisher's famous tea-tasting experiment, the total number of positive re-  
 482 sponses  $s$  was fixed by design. However, Fisher later argued that it should be  
 483 considered fixed regardless, arguing that  $S$  has much in common with the sample  
 484 size. Conditioning on the total successes was later proposed not only for compar-  
 485 ative trials (e.g., [Gail and Gart, 1973; Gart, 1969](#)) but also in matched case-control  
 486 studies (e.g., [Hirji et al., 1988](#)) and for tables of higher dimension (see [Hirji, 2006](#)).  
 487 The theory extends naturally to generalised linear models with canonical link. For  
 488 non-canonical link, conditioning on approximate ancillary statistics has led to the  
 489  $p^*$  formula discussed earlier. The argument for conditioning in  $2 \times 2$  tables cannot  
 490 be understood without reference to this wider context.

491 In the narrower context of  $2 \times 2$  tables, the conditional model has a single  
 492 free variable  $y_1$  and a single parameter  $\theta$  and the theory is very simple. The  
 493 unconditional model has two free variables  $(y_1, y_0)$  and an additional nuisance  
 494 parameter. Even though the model is still very simple, it is rich enough to expose  
 495 all of the difficulties and limitations of frequentist inference. There are plausible  
 496 arguments for either model, which we now elucidate.

### 5.2 Arguments for conditional inference

498 The first argument for conditional inference is based on Lehmann's theory.  
 499 He showed that in full rank exponential families the use of conditioning, with  
 500 randomisation, leads to an unconditional test. This suggests that the conditional  
 501 likelihood contains all the relevant information with respect to the parameter of  
 502 interest. While seldom used in practice, randomisation reveals this basic structure,  
 503 just as embedding real numbers within the complex number system brings insights  
 504 into solving polynomials.

505 A second argument is that the total number of successes  $S$  has the same ger-  
 506 mane properties as the sample size. However,  $S$  differs from the sample size in  
 507 that its distribution depends on the parameters. More refined arguments were  
 508 based on the idea that  $S$ , by itself, is uninformative about the interest parameter

509  $\theta$ . Formalising this notion leads to various definitions of approximate ancillarity  
 510 and sufficiency (Barndorff-Nielsen, 1973; Cox, 1980; Godambe, 1980). However,  
 511 all these definitions have unsatisfactory implications for some statistical models  
 512 and it is fair to say that no consensus emerged. There is also debate about the  
 513 ancillarity argument itself (Berkson, 1978).

A third argument, due to the second author, points to an epistemologically  
 undesirable property of unconditional inference. Consider the most extreme out-  
 come, when there are all successes for treatment and all failures for control; the  
 $p$ -value equals the probability of this single most extreme outcome  $y_E = (n_1, 0)$ .  
 The unconditional  $p$ -value can then be decomposed as

$$Pr(Y = y_E) = Pr(Y_1 = n_1 | S = n_1) \times Pr(S = n_1; p).$$

514 The first factor is Fisher's  $p$ -value; the second factor is a pure probability about  
 515  $S$  which depends on  $p$  but whose maximum value is small. Why should the event  
 516  $S = n_1$  be counted against the null hypothesis? In words, why should a total  
 517 of  $n_1$  successes out of  $n_0 + n_1$  trials be counted as evidence that the treatment  
 518 works? Unconditional inference commits us to this inference.

519 The fourth and last argument for conditioning is simplicity and convenience:  
 520 conditioning on  $S$  eliminates  $\varphi$  from the model and provides an incredibly simple  
 521 model 2.7 with a single variable  $y_1$  depending on the parameter of interest  $\theta$ .  
 522 All good statistical modelling involves treating incidental aspects of the data  
 523 generating mechanism as fixed so that we can focus on the issue at hand. While  
 524 eliminating  $\varphi$  from the model is attractive, there are other methods that do not  
 525 involve conditioning, as detailed in section 4 (see also Basu, 1977, for an earlier  
 526 inventory of methods). So conditioning, while one option, is not necessary to  
 527 account for the nuisance parameter.

### 528 5.3 Comparison of tests under the conditional model

529 In this section, let us accept the conditional model 2.7 as the model for the  
 530 number of treatment successes  $y_1$  given the total successes  $s$ .

531 Fisher's  $p$ -value is the probability of an observed event. It answers the ques-  
 532 tion: out of  $y_1 + y_0 = s$  successes, how often would at least  $y_1$  of them be in  
 533 the treatment group if the treatment has no effect? It also decreases in  $y_1$  for  
 534 fixed  $s$  and increases in  $s$  for fixed  $y_1$ , which are the key monotonicity properties  
 535 in section 2.5. So the key logical hurdles are passed. On the other hand, for a  
 536 fixed target nominal size  $\alpha$ , the test size  $a_s$  given in 3.2 is less than nominal,  
 537 sometimes much less. This is the source of the common claim that Fisher's test  
 538 is conservative. The claim is spurious.

539 Within the conditional framework, some size conservatism is an inevitable  
 540 consequence of discreteness but an exact  $p$ -value still exists, as explained in sec-  
 541 tion 2.3. Fisher's  $p$ -value is exact in this strong technical sense. It is the smallest  
 542 possible valid  $p$ -value that is monotone increasing in  $y_1$ . So within the conditional  
 543 framework, Fisher's test is not unnecessarily conservative. Indeed, any other test  
 544 that is valid will be even more conservative and any test that is less conservative  
 545 will be invalid.

546 The mid- $p$  and Liebermeister proposals are both attractive, but their condi-  
 547 tional size can exceed nominal (Hirji et al., 1991; Seneta and Phipps, 2001) and  
 548 the  $p$ -values are never guaranteed. Seneta and Phipps (2001) compared the size



549 attained by Fisher's, Lieberman's and Lancaster's test. These authors showed  
 550 that Lieberman's test is the closest to the nominal level (even though it is not  
 551 valid, exceeding the nominal level) followed by Lancaster's and Fisher's test. So  
 552 if closeness of attained size, rather than validity, were our key criterion we might  
 553 be moved towards Lieberman's test. However, we consider validity of the test  
 554 and the guaranteed property of a  $p$ -value a key criterion.

555 In addition, both mid- $p$  and Lieberman suffer from the drawback that they  
 556 are not the probability of any observed event. While we might consider approxi-  
 557 mations to a guaranteed test, neither is an approximation to the Fisher  $p$ -value.  
 558 Certainly, neither can be justified within the conditional model. Evaluated un-  
 559 conditionally, their performance may be acceptable and we will present some  
 560 numerical results in section 6. However, there are competing tests within the  
 561 unconditional framework that we will ultimately prefer.

562 Finally, the randomised version of Fisher's test is UMPU. So Fisher's test may  
 563 be thought of as the closest valid discrete approximation to the UMPU test.  
 564 Thus, any conservatism of Fisher  $p$ -value is an inevitable artifact of discreteness.  
 565 In summary, within the conditional framework there appears to be no alternative  
 566 to Fisher's test. As we shall see later though, the criticisms of Fisher's test are  
 567 mainly made from an unconditional perspective.

568 An area for future research is to clarify the properties of  $p^*$  and approximately  
 569 conditional  $p$ -values. The former are not necessarily valid and their degree of  
 570 liberalism should be better assessed. For the latter, it remains unresolved how to  
 571 determine a general neighborhood for conditioning.

#### 572 **5.4 Arguments for unconditional inference**

573 The most persuasive argument for the unconditional model is that in future  
 574 repetitions of the experiment the value of  $S$  will vary. If we want practical assess-  
 575 ment of the future performance of the test - which is the key aim of frequentist  
 576 inference - then we should allow  $S$  to vary. At the very least, this suggests aug-  
 577 menting any conditional test with a statement of its unconditional properties.

578 There are two specific arguments against the conditional model. The first is  
 579 that conditional inference does not easily generalise to non-canonical parameters.  
 580 In the context of  $2 \times 2$  tables, we can perform conditional tests of the log-odds ratio  
 581 but not of the risk difference, as pointed out in section 3.4. Moreover, even with  
 582 canonical parameters the relevant conditional distribution can be degenerate,  
 583 leading to a test with zero size and power. While  $p^*$  methods were developed to  
 584 address these problems, the fact that it gives a non-degenerate answer in this  
 585 latter case is problematic.

586 The second argument against the conditional model is conservatism. Basing  
 587 tests on the unconditional model allows greater unconditional power. This is  
 588 partly because the distributions involved are much less discrete but also because  
 589 the conditional size  $a_s$  of a test need not be less than  $\alpha$  for all  $s$ , so long as its  
 590 mean value is less than  $\alpha$ . It is worth noting though that there are cases where  
 591 conditional tests are more powerful, see [Mehrotra et al. \(2003\)](#). Extending the  
 592 investigation to the case of three binomials (which arises in a three-arm clinical  
 593 trials), the conditional and unconditional approach seem to achieve similar power  
 594 ([Mehta and Hilton, 1993](#)).

595 There does not exist a conclusive argument for or against conditioning, either

596 in general or for  $2 \times 2$  tables. Many might argue that this dilemma reveals a  
 597 fundamental weakness in frequentist inference. For  $2 \times 2$  tables, if the conditioning  
 598 argument is accepted, then Fisher's exact  $p$ -value is exact in the sense of equation  
 599 2.5. If the conditioning argument is not accepted, then there is a much wider field  
 600 of candidate tests which have to be compared. This includes ostensibly conditional  
 601 tests that are made unconditional by the M, B or E steps.

## 602 5.5 Comparison of tests under the unconditional model

603 There are many tests in current use: Fisher, mid- $p$ , Lieberman, pooled-  
 604 Z, unpooled-Z, the likelihood ratio and various Wald tests. All of these can be  
 605 assessed within the unconditional model. None of them are exact. All can be  
 606 adjusted using the  $M$ -step,  $B$ -step or  $E$ -step. A numerical study below will il-  
 607 luminate the properties of the basic and adjusted tests. However, we can say  
 608 quite a lot about the three adjustments based on theoretical considerations. The  
 609 example and figure in section 4.4 serves as an excellent heuristic.

610 Firstly, all  $M$ -step  $p$ -values are exact in the sense of 2.5 and subject to the  
 611 ordering of the sample space induced by the initial test cannot be improved. If  
 612 there is a spike in the profile then the maximised  $p$ -value will tend to be larger and  
 613 power will be degraded. If there is no spike then maximisation will just recalibrate  
 614 the test to remove its conservatism or liberality.

615 Partially maximised  $p$ -values tend to be smaller when there is a spike and pay  
 616 an insurance premium  $\gamma$  to achieve this, even if there is not a spike. From their  
 617 definitions, it can be asserted that  $P_B(y) < P_M(y) + \gamma$  but when there is a spike  
 618  $P_B(y)$  will be much smaller than  $P_M(y)$ . The  $B$ -step  $p$ -value is guaranteed but is  
 619 not exact: only  $M$   $p$ -values can be exact, and applying the  $M$ -step to the  $B$   $p$ -value  
 620 will reduce it slightly (but by no more than  $\gamma$ ). For more complex models where  
 621  $\varphi$  is a vector, construction of the confidence region  $C_\gamma$  is left unspecified and so  
 622 partial maximisation is not a well defined procedure. Indeed, for many models no  
 623 exact confidence region for  $\varphi$  exists and the method cannot be formally applied.

624 The estimated  $p$ -value is the smallest of the three. It can be easily shown  
 625 that  $P_E(y) < P_B(y) - \gamma < P_M(y)$ . The cost is that  $P_E(y)$  is not guaranteed  
 626 and tests based on it can be invalid. However, empirically it is found that the  
 627 profile of  $P_E(y)$  is very flat, much closer than any asymptotic argument might  
 628 suggest (Lloyd, 2008b). Consequently,  $P_{E+M}(y) \approx P_E(y)$ . This supports the use  
 629 of  $P_{E+M}(y)$  in principle and  $P_E(y)$  in practice. Of course, when the original profile  
 630 is quite flat, all the  $p$ -values will be close. However, for all of the standard tests  
 631 the profile can be far from flat.

632 A final issue worth mentioning is the choice of initial test to generate the  
 633 profile. Maximised  $p$ -values depend quite a lot on this choice, partially maximised  
 634  $p$ -values much less and estimated  $p$ -values hardly at all. This is a very attractive  
 635 feature of  $P_E(y)$ , as it effectively removes any consequences of the user's choice  
 636 of initial test. All these assertions will be verified in the numerical study below.

637 We now turn to computational issues. All three adjustments require the set  
 638  $\{P(Y) \leq P(y)\}$  to be enumerated. Potentially, this requires evaluation of the  
 639 generating  $p$ -value for all possible data sets. So a simpler generating  $p$ -value has  
 640 great computational advantages. Amongst the simplest are Fisher's exact  $p$ -value  
 641 and the maximised version was recommended by Boschloo (1970). The  $M$  and  $B$

642 steps require similar computation, while the  $E$ -step is faster, since the nuisance  
 643 parameter  $\varphi$  is estimated rather than maximised. For  $2 \times 2$  tables, all three can  
 644 be calculated in a few seconds for sample sizes up to 1000.

645 The theoretically attractive  $E + M$   $p$ -value requires computing all possible  
 646 values of  $P_E(y)$ . This is currently limited to modest sample sizes of a few 100.  
 647 Nevertheless, if computation were not an issue we would recommend the  $E+M$   
 648  $p$ -value, based on the LR test because of its monotonicity properties and consis-  
 649 tently high power of the resulting  $E + M$   $p$ -value.

650 For more complex models where  $\varphi$  is a vector, the  $M$  and  $B$ -steps are not  
 651 practical to compute. The  $E$ -step is not adversely affected by the dimension of  $\varphi$   
 652 and can be implemented for generalised linear models using importance sampling,  
 653 see [Lloyd \(2012\)](#).

## 6. A NUMERICAL STUDY

654 To illustrate, verify and compare the unconditional performance of the tests  
 655 reviewed in this article, we conducted a numerical study. Full details are provided  
 656 in the online appendix but it is pertinent to give representative results here. We  
 657 considered eight test statistics: pooled, unpooled, log Wald, SRLR,  $p^*$ , Lieber-  
 658 meister, mid-p, and Fisher. Only the last of these is guaranteed and the others  
 659 can all be liberal for some parameter values. We calculated the unconditional  
 660 size and power of the tests using five different versions of the basic statistics: raw,  
 661 M, B (with  $\gamma = .001$ ), E, and E+M. We fixed the nominal size  $\alpha = 0.05$ , the  
 662 control sample size  $n_0 = 40$ , and the treatment sample size  $n_1 = 60$ . This choice  
 663 is broadly representative of the patterns we have observed across all unbalanced  
 664 designs.

665 In Table 1 we report the exact size of the 40 tests using two measures: maximum  
 666 size with respect to  $p_0$  (upper section) and mean size with respect to  $p_0$  (lower  
 667 section). In the max part of the table, violations over 0.051 are highlighted in  
 668 red (over 0.06 is bold). Amongst the raw tests, mid-p is very close to exact but  
 669 this is not the always case for other sample sizes. M, B and E+M tests are all  
 670 theoretically valid (which is confirmed numerically in the table), whereas the E  
 671 test does occasionally violate size by a non-trivial amount, but only when the  
 672 original statistic is the unpooled or Fisher.

673 In the lower part of the table, we use colour coding to highlight the largest  
 674 possible mean size subject to validity. This would imply a flatter profile and  
 675 would typically lead to higher power. The B procedure is never worse than the M  
 676 procedure, but is occasionally only slight advantageous. By contrast, the E and  
 677 E+M procedures are very stable across test statistics, and the advantage is more  
 678 pronounced.

679 In Table 2 we show the power results for three selected values of  $p_0$ , and  
 680 corresponding values of  $p_1$  chosen so that the power is in a practically interesting  
 681 range. As previously reported in the literature, it emerges that the B procedure  
 682 leads to more powerful tests than the M procedure. E and E+M tests are always  
 683 best or the equal best tests; the added value of these procedures is that they seem  
 684 to work well across all circumstances. The E procedure should be recommended in  
 685 applications, with the caution that it can occasionally lead to very slight violation  
 686 of size when sample sizes are unbalanced. The E+M procedure is guaranteed, at  
 687 the cost of a higher computational burden.

TYPE	pooled	unpooled	log.wald	lr	p*	lieberm.	midp	fisher
raw	<b>0.066</b>	<b>0.088</b>	0.040	<b>0.088</b>	<b>0.088</b>	<b>0.064</b>	0.051	0.033
M	0.049	0.048	0.047	0.048	0.046	0.050	0.049	0.049
B	0.047	0.048	0.047	0.048	0.050	0.047	0.047	0.049
E	0.050	<b>0.054</b>	<b>0.054</b>	0.050	0.050	<b>0.057</b>	<b>0.057</b>	<b>0.057</b>
E+M	0.050	0.050	0.050	0.050	0.046	0.050	0.050	0.050
raw	0.051	0.052	0.033	0.053	0.053	0.047	0.042	0.028
M	<b>0.037</b>	<b>0.030</b>	<b>0.037</b>	<b>0.029</b>	<b>0.027</b>	0.042	0.041	0.041
B	<b>0.041</b>	<b>0.036</b>	<b>0.039</b>	<b>0.038</b>	<b>0.039</b>	<b>0.042</b>	0.042	0.041
E	<b>0.045</b>	0.045	0.044	0.043	0.043	0.045	0.045	0.046
E+M	<b>0.045</b>	<b>0.045</b>	<b>0.044</b>	<b>0.043</b>	<b>0.043</b>	<b>0.044</b>	<b>0.043</b>	<b>0.043</b>

TABLE 1

Maximum size (above) and mean size (below) calculated for 8 test statistics with 5 methods;  
 $n_0 = 40$ ,  $n_1 = 60$ ,  $\alpha = 0.05$ .

p0	p1	TYPE	pooled	unpooled	log.wald	lr	p*	lieberm.	midp	fisher	
0.10	0.25	M	0.583	0.583	0.628	0.561	<b>0.508</b>	0.595	0.628	0.613	
		B	0.628	0.583	0.628	0.598	0.583	0.628	0.628	<b>0.636</b>	
		E	0.628	0.633	0.629	0.633	0.629	0.629	0.629	0.629	0.652
		E+M	0.628	<b>0.633</b>	0.629	0.633	0.629	0.629	0.628	0.628	<b>0.636</b>
0.50	0.70	M	0.610	<b>0.555</b>	0.610	<b>0.540</b>	<b>0.540</b>	0.634	0.623	0.634	
		B	0.623	0.603	0.623	0.623	0.632	0.634	0.634	0.634	
		E	0.634	0.634	0.634	0.634	0.634	0.634	0.634	0.634	
		E+M	0.634	0.634	0.634	0.634	0.634	0.634	0.634	0.634	0.634
0.75	0.90	M	0.609	<b>0.536</b>	0.589	0.586	0.586	<b>0.653</b>	0.608	0.628	
		B	0.608	0.608	0.608	0.608	0.628	0.628	0.628	0.628	
		E	<b>0.653</b>	0.673	0.672	0.651	0.651	0.653	0.653	0.672	
		E+M	<b>0.653</b>	<b>0.653</b>	0.651	0.651	0.651	0.653	0.651	0.651	0.651

TABLE 2

Power for three different combinations of  $p_0$  and  $p_1$ , for 8 test statistics and 5 methods;  
 $n_0 = 40$ ,  $n_1 = 60$ ,  $\alpha = 0.05$ .

## 7. CONCLUSIONS

688 In this paper, we have reviewed both the conditional and unconditional ap-  
689 proach to the  $2 \times 2$  table. Both approaches lead to valid frequentist inference  
690 within their own different model frameworks.

691 Fisher originally suggested that the statistical model should depend on the  
692 study design. Indeed, when the total sum of successes in a binomial trial is fixed  
693 by design, it seems natural to consider the conditional approach and to evaluate  
694 power conditionally. When the sum is not naturally fixed by design, as is more  
695 often the case, it seems at least pertinent to adopt an unconditional perspective.  
696 Treating the sum of successes as fixed when it is not is defensible, but is based  
697 on general conditionality arguments whose application to the  $2 \times 2$  table is not  
698 completely clear.

699 Thus, rather than conclusively support one approach against the other, we  
700 have reviewed and assessed alternative procedures within the conditional and  
701 unconditional models separately.

702 From the conditional perspective, there does exist an optimal test, namely the  
703 randomised version of Fisher's test. This test would be the gold standard for bino-  
704 mial endpoints, but cannot be recommended in practice, because randomisation  
705 introduces extra variation into the analysis. It is worth observing that within a de-  
706 cision theory framework randomised tests do not violate the sufficiency principle  
707 (e.g., [Lehmann and Romano \(2005\)](#), p. 58) since the randomisation distribution  
708 is the same for all data that give the same value of a sufficient statistic. At the  
709 point where a random number is drawn to complete the test and reject or accept  
710 the null hypothesis, the classical Fisherian sufficiency principle is contradicted.  
711 Regardless of these theoretical arguments however, randomisation is almost never  
712 used in practice.

713 The unrandomised version, which is Fisher's exact test, has long been criticised  
714 for its conservatism, both conditionally and unconditionally. However, Fisher's  
715 test is valid and satisfies our definition of exactness. In addition, if one accepts  
716 the conditional model, then conservatism is an inevitable consequence of the dis-  
717 creteness. Fisher's p-value is exact and is the smallest possible valid p-value that  
718 is monotone increasing in  $y_1$ . By contrast, both Lancaster's and Lieberman's  
719 tests mitigate the conservatism of Fisher's exact test, but both are necessarily  
720 liberal within the conditional framework. Moreover, neither is the probability of  
721 any observed event. Hence our recommendation is clear: within the conditional  
722 framework, Fisher's is really the only test to recommend.

723 The newest methods that are motivated from a conditional approach are  $p^*$   
724  $p$ -values and approximate conditional p-values, which follow recent developments  
725 in likelihood theory. These approaches have the virtue of extending conditional  
726 methods to models where exact conditioning is not possible, and provide close  
727 approximation to conditional procedures when it is possible. However, for the  $2 \times 2$   
728 table an optimal exact approach is available and there is no need to approximate  
729 it. More generally, the validity of such p-values is not guaranteed and they are  
730 not easy to write down explicitly or compute.

731 Amongst unconditional methods, we have explicated four methods for con-  
732 structing p-values with good size control. From the least to the most compu-  
733 tationally intensive these methods are E, M, B and E+M. The last three are  
734 guaranteed to be valid, while the first is very close to valid in practice.

735 M  $p$ -values account for the worst possible scenario and represent the most  
736 classical approach to handling the nuisance parameter, which is the key difficulty  
737 in the unconditional approach. Maximization is a straightforward procedure and  
738 relatively easy to compute, at least for models with a single nuisance parameter;  
739 it guarantees exactness and validity, but allowing for the worst case often leads  
740 to unnecessarily low power.

741 B  $p$ -values are a simple method for overcoming the power loss of accounting  
742 for an unlikely worst case, especially for highly unbalanced designs. Though not  
743 exact, they are necessarily valid, while typically being smaller than M  $p$ -values.  
744 The difficulty of extending to more general and higher dimensional models and  
745 the lack of specification for the level of confidence is a theoretical weakness.  
746 However, for  $2 \times 2$  tables, there is no practical impediment to their use and we  
747 would recommend B  $p$ -values over M  $p$ -values (with  $\gamma = 0.001$ ). B  $p$ -values are  
748 today available in some statistical softwares.

749 The easiest  $p$ -value to compute (even for sample sizes of several 1,000) is the  
750 E  $p$ -value, which leads to a flat profile. It is very close to exact and could be  
751 recommended in practice, provided that tiny violations of the size constraint are  
752 acceptable. The method extends to general models in a straightforward manner.

753 E+M tests combine the use of a flatter significance profile (E step) with guar-  
754 anteeing validity (M step). They are typically slightly more powerful than B  
755  $p$ -values and they do not require user choice of a level of confidence. Another  
756 attractive feature is that E+M (and E) methods lead to almost the same fi-  
757 nal inference, regardless of the choice of the test statistic. The only weakness of  
758 the E+M method is its computational burden. R-code for all these methods is  
759 available from the authors.

760 This paper has reviewed contemporary approaches to the  $2 \times 2$  table from a  
761 frequentist point of view. One reason of this choice is practical, as the clinical  
762 trials regulators mainly employ frequentist protocols. However, for the sake of  
763 comparison, we briefly mention Bayesian methods.

764 The Bayesian paradigm introduces *a priori* information in the inferential frame-  
765 work, and the plausibility of a hypothesis based on the data is assessed in terms  
766 of *a posteriori* probabilities. The key idea is treating parameters as random vari-  
767 ables, so that, unlike the frequentist approach, one can integrate out any nuisance  
768 parameters. By imposing a continuous distribution on the parameters, the prob-  
769 lem of discreteness is naturally solved, and, in case of the  $2 \times 2$  table, computations  
770 are very simple.

771 Since conclusions are dependent on the prior specification, the frequentist prop-  
772 erties of Bayesian methods cannot be stated in general. Nevertheless, there are  
773 links between Bayesian and frequentist inference for  $2 \times 2$  tables. For instance,  
774 Lieberman's test can be generated from a Bayesian argument and has quite  
775 good unconditional properties. The one-sided  $p$ -value from the pooled  $z$ -test can  
776 also be derived as the posterior probability of the null hypothesis based on in-  
777 dependent Jeffries priors for  $(p_0, p_1)$ , see [Howard \(1998\)](#). Confidence intervals,  
778 which is not a topic we have touched on, can be replaced by highest posterior  
779 density intervals, see [Brown et al. \(2001\)](#) and [Brown et al. \(2002\)](#).



## REFERENCES

- 780 AGRESTI, A.(1992). A survey of exact inference for contingency tables. *Statistical Science* **7**  
781 131-153.
- 782 AGRESTI, A.(2001). Exact inference for categorical data: recent advances and continuing contro-  
783 versies. *Statistics in Medicine* **20** 2709–2722.
- 784 AGRESTI, A.(2002). *Categorical data analysis*. John Wiley and Sons, Hoboken(NJ).
- 785 BARNARD, G.A.(1945). A new test for  $2 \times 2$  tables. *Nature* **156** 388.
- 786 BARNARD, G.A.(1947). Significance tests for  $2 \times 2$  tables. *Biometrika* **34** 123–138.
- 787 BARNDORFF-NIELSEN, O.E.(1973). On M-ancillarity. *Biometrika* **60** 447–455.
- 788 BARNDORFF-NIELSEN, O.E.(1983). On a formula for the distribution of the maximum likelihood  
789 estimator. *Biometrika* **70** 343–356.
- 790 BASU, D.(1977). On the elimination of nuisance parameters. *Journal of the American Statistical*  
791 *Association* **72** 355–366.
- 792 BERGER, R.L. AND BOOS, D.D.(1994). P values maximized over a confidence set for the nuisance  
793 parameter. *Journal of the American Statistical Association* **89** 1012–1016.
- 794 BERKSON, J.(1978). In dispraise of the exact test: Do the marginal totals of the  $2 \times 2$  table  
795 contain relevant information respecting the table proportions? *Journal of Statistical Planning*  
796 *and Inference* **2** 27–42.
- 797 BERGER, R.L. AND SIDIK, K.(2003). Exact unconditional tests for a  $2 \times 2$  matched-pairs design.  
798 *Statistical Methods in Medical Research* **12** 91–108.
- 799 BOSCHLOO, R.D.(1970). Raised conditional level of significance for the  $2 \times 2$  table when testing  
800 the equality of two probabilities. *Statistica Neerlandica* **24** 1–9.
- 801 BROWN, L.D., CAI, T.T., AND DASGUPTA, A.(2002). Confidence intervals for a binomial pro-  
802 portion and asymptotic expansions. *Annals of Statistics* **30** 160-201.
- 803 BROWN, L.D., CAI, T.T., AND DASGUPTA, A.(2001). Interval estimation for a binomial pro-  
804 portion. *Statistical Science* **16** 160-201.
- 805 CHOI, L., BLUME, J.D., AND DUPONT, W.D.(2015). Elucidating the foundations of statistical  
806 inference with  $2 \times 2$  tables. *PloSone* **10** e0121263.
- 807 COX, D.R.(1980). Local ancillarity. *Biometrika* **62** 292–282.
- 808 COX, D. R., AND HINKLEY, D. V. (1974). *Theoretical Statistics*. Chapman and Hall, London.
- 809 DAVISON, A. C., FRASER, D. A. S., AND REID, N.(2006). Improved likelihood inference for  
810 discrete data. *Journal of the Royal Statistical Society* **B68** 495-508.
- 811 FINNER, H. AND STRASSBURGER, K.(2002). Structural properties of UMPU-tests for  $2 \times 2$  tables  
812 and some applications. *Journal of Statistical Planning and Inference* **104** 103–120.
- 813 FISHER, R.A.(1935). *The design of experiments*. 1st ed., Oliver and Boyd, London.
- 814 GAIL, M.H. AND GART. J.J.(1973). The determination of sample sizes for use with the exact  
815 conditional test in  $2 \times 2$  comparative trials. *Biometrics* **29** 441–448.
- 816 GART. J.J.(1969). An exact test for comparing matched proportions in crossover designs.  
817 *Biometrika* **56** 75–80.
- 818 GODAMBE, V.P.(1980). On sufficiency and ancillarity in the presence of nuisance parameters.  
819 *Biometrika* **67** 155–162.
- 820 HARRIS, B. AND SOM, A.P.(1991). Theory and counterexamples for confidence limits on system  
821 reliability. *Statistics and Probability Letters* **11** 411–417.
- 822 HIRJI, K.F.(2006). *Exact analysis of discrete data*. Chapman and Hall, CRC Press, Boca Raton  
823 (FL).
- 824 HIRJI, K.F, MEHTA, C.R., AND PATEL, N.R.(1988). Exact inference for matched-case control  
825 studies. *Biometrics* **44** 803–814.
- 826 HIRJI, K.F, TAN, S.J., AND ELASHOFF, R.M.(1991). A quasi-exact test for comparing two  
827 binomial proportions. *Statistics in Medicine* **10** 1137–1153.
- 828 HOWARD, J.V.(1998). The  $2 \times 2$  table: A discussion from a Bayesian viewpoint. *Statistical Science*  
829 **13** 351–367.
- 830 HWANG, J.G., AND YANG, M.C.(2001). An optimality theory for mid p-values in  $2 \times 2$  contin-  
831 gency tables. *Statistica Sinica* **11** 807–826.
- 832 KABAILA, P.V.(2005). Computation of exact confidence limits from discrete data. *Computa-*  
833 *tional Statistics* **20** 401–414.
- 834 KABAILA, P.V. AND LLOYD, C.J.(2006). Improved Buehler limits based on refined designated  
835 statistics. *Journal of Statistical Planning and Inference* **136** 3145–3155.
- 836 LANCASTER, H.O.(1961). Significance tests in discrete distributions. *Journal of the American*  
837 *Statistical Association* **56** 223–224.

- 838 LEHMANN, E.(1959). *Testing Statistical Hypotheses*. Wiley, New York.
- 839 LEHMANN, E. AND ROMANO, J.P.(2005). *Testing Statistical Hypotheses*. Springer, New York.
- 840 LIEBERMEISTER, C.(1877). *Über Wahrscheinlichkeitsrechnung in Anwendung auf therapeutische*  
841 *Statistik*. Breitkiof and Härtel.
- 842 LLOYD, C.J.(2008a). A new and more powerful unconditional test of no treatment effect from  
843 binary matched pairs. *Biometrics* **64** 716–723.
- 844 LLOYD, C.J.(2008b). Exact  $p$ -values for discrete models obtained by estimation and maximiza-  
845 tion. *Australian and New Zealand Journal of Statistics* **50** 329–345.
- 846 LLOYD, C.J.(2010a).  $p$ -values based on approximate conditioning and  $p^*$ . *Journal of Statistical*  
847 *Planning and Inference* **140** 1073–1081.
- 848 LLOYD, C.J.(2010b). Bootstrap and second order tests of risk difference. *Biometrics* **66** 975–982.
- 849 LLOYD, C.J.(2012). Computing highly accurate or exact  $p$ -values using importance sampling.  
850 *Computational Statistics and Data Analysis* **56** 1784–1794.
- 851 LYDERSEN, S., FAGERLAND, M. W., AND LAAKE, P.(2009). Recommended tests for association  
852 in  $2 \times 2$  tables. *Statistics in Medicine* **28** 1159–1175.
- 853 MARTIN ANDRES, A.(1991). A review of classic non-asymptotic methods for comparing two  
854 proportions by means of independent samples. *Communications in Statistics-Simulation and*  
855 *Computation* **20** 551–583.
- 856 McDONALD, L.L., DAVIS, B.M., AND MILLIKEN, G.A.(1977). A nonrandomized unconditional  
857 test for comparing two proportions in  $2 \times 2$  contingency tables. *Technometrics* **19** 145–158.
- 858 MEHTA, C.R., AND HILTON, J.F.(1993). Exact power of conditional and unconditional tests:  
859 going beyond the  $2 \times 2$  contingency table. *American Statistician* **47** 91–98.
- 860 MEHROTRA, D.V., CHAN, I.S.F., AND BERGER, R.L.(2003). A cautionary note on exact uncon-  
861 ditional inference for a difference between two independent binomial proportions. *Biometrics*  
862 **59** 441–450.
- 863 PEARSON, K.(1900). On the criterion that a given system of deviations from the probable in the  
864 case of a correlated system of variables is such that it can be reasonably supposed to have  
865 arisen from random sampling. *Philosophical Magazine* **5** 157–175.
- 866 PIERCE, D.A. AND PETERS, D.(1992). Practical use of higher order asymptotics for multipa-  
867 rameter exponential families. *Journal of the Royal Statistical Society, Series B* **54** 701–737.
- 868 PIERCE, D.A. AND PETERS, D.(1999). Improving on exact tests by approximate conditioning.  
869 *Biometrika* **86** 265–277.
- 870 RÖHMEL, J. AND MANSMANN, U.(1999). Unconditional non-asymptotic one-sided tests for in-  
871 dependent binomial proportions when the interest lies in showing non-inferiority and/or su-  
872 periority. *Biometrical Journal* **2** 149–170.
- 873 RÖHMEL, J.(2005). Problems with existing procedures to calculate exact unconditional  $p$ -values  
874 for non-inferiority/superiority and confidence intervals for two binomials and how to resolve  
875 them. *Biometrical Journal* **47** 37–47.
- 876 SENETA, E., AND PHIPPS, M.C. (2001). On the comparison of two observed frequencies. *Bio-*  
877 *metrical Journal* **43** 23–43.
- 878 SKIPKA, G., MUNK, A. AND FREITAG, G.(2004). Unconditional exact tests for the difference of  
879 binomial probabilities contrasted and compared. *Computational Statistics and Data Analysis*  
880 **47** 757–773.
- 881 STORER, B.E. AND KIM, C.(1990). Exact properties of some exact test statistics for comparing  
882 two binomial proportions. *Journal of the American Statistical Association* **85** 146–155.
- 883 TOCHER, K. D.(1950). Extension of the Neyman-Pearson theory of tests to discontinuous vari-  
884 ates. *Biometrika* **37** 130-144.
- 885 WELLS, M.T. (2010). Optimality results for mid  $p$ -values. In *Borrowing strength: Theory pow-*  
886 *ering applications - A Festschrift for Lawrence D. Brown (pp. 184-198)*, 2nd ed. Institute of  
887 Mathematical Statistics.
- 888 YATES, F. (1984). Tests of significance for  $2 \times 2$  contingency tables. *Journal of the Royal Sta-*  
889 *tistical Society, Series A* **147** 426-463.