

Intervals for Post-Test Probabilities: A Comparison of Five Methods

DOUGLAS MOSSMAN, M.D.

Professor and Director, Division of Forensic Psychiatry,
Wright State University School of Medicine;
Adjunct Professor, University of Dayton School of Law

JAMES O. BERGER, PH.D.

Arts and Sciences Professor of Statistics
Institute of Statistics and Decision Sciences, Duke University

*Submitted March 2001 to Medical Decision Making.
Please do not cite or quote without the express permission of Dr. Mossman.*

Address for correspondence:

W. S. U. Department of Psychiatry, P. O. Box 927, Dayton, Ohio 45401-0927
Telephone: (937) 276-8325 — Fax: (937) 275-2817 — e-mail: dmossman@pol.net

Abstract

Background. Several medical articles discuss methods of constructing confidence intervals for single proportions and the likelihood ratio, but scant attention has been given to the systematic study of intervals for the posterior odds or the “positive predictive value” of a test. **Methods.** We describe five methods of constructing confidence intervals for post-test probabilities when estimates of sensitivity, specificity, and the pre-test probability of a disorder are derived from empirical data. We then evaluate each method to see how well the intervals’ coverage properties correspond to their nominal value. **Results.** When the estimates of pre-test probabilities, sensitivity, and specificity are derived from more than 80 subjects and are not close to zero or 1, all methods generate intervals with appropriate coverage properties. When these conditions are not met, however, the best-performing method is an “objective Bayesian” approach implemented by a simple simulation using a spreadsheet. **Conclusions.** Physicians and investigators can generate accurate confidence intervals for post-test probabilities in small-sample situations using the objective Bayesian approach.

* * * * *

Key words: positive predictive value, Bayes’s Theorem, posterior probability, posterior odds, post-test probability, confidence interval, objective Bayesian

Introduction

Mr. Smith is worried that he has disorder D , and consults Dr. Jones. Dr. Jones orders a test T that has only two outcomes, “positive” and “negative,” and Mr. Smith’s test result comes back “positive.” Dr. Jones tells his patient that he has a 69% chance of having the disorder, given this test result.

Mr. Smith asks his doctor how he came to this conclusion. Dr. Jones — an unusual physician who uses Bayes’s Theorem and discusses it with his patients — answers, “Twenty percent of men who are your age have D . Test T has a 90% sensitivity and a 90% specificity. Applying Bayes’s Theorem, your positive test result means your chance of having D is 69%.”

Mr. Smith likes his doctor’s answer, because he is a *very* unusual patient who wants his care-givers to formulate clinical decisions with explicit mathematical rigor. He expects his physicians to use Bayes’s Theorem, but is troubled by one aspect of Dr. Jones’s calculation.

“I understand how you reached your conclusion,” Mr. Smith responds, “but your information doesn’t allow you to say that my chance of having D is *exactly* 69%. Your figure for the 20% prevalence came from a published medical study that evaluated 50 men in my age group. Ten of these men had D , but this finding only lets doctors *estimate* the prevalence of D as being 20% in men my age.” Dr. Jones agrees with his patient’s summary of the literature on the prevalence of D .

“Likewise,” Mr. Smith continues, “your figures for sensitivity and specificity came from another published medical study of the diagnostic test T . That article reported that 36 of 40 people who had D also had positive tests, but only four of 40 people who did not have D had positive tests. These results just let the researchers *estimate* the test’s sensitivity and specificity as 90%.”

“The published estimates of prevalence, sensitivity, and specificity are subject to random sampling error,” continues Mr. Smith. “So what I want to know is this: What is the 95 percent confidence interval for my probability of having D , given my positive test result and the imprecision in the estimates? Knowing whether the interval is narrow or broad might affect my decisions about getting other tests or choosing treatment.”

Dr. Jones is not sure how to compute the answer to his patient’s question, and the existing medical literature provides little help. Several writers have discussed methods of constructing confidence intervals for single proportions,¹ for the likelihood ratio,²⁻⁴ and the odds ratio.⁵ To the best of our knowledge, however, the medical literature does not prescribe or discuss methods for constructing intervals when proportions are combined in Bayes’s Theorem.

This article attempts to fill this gap. The following section describes five methods of constructing intervals, and applies these methods to post-test probabilities calculated from

random samples and data-based estimates of prevalence (or pre-test probability of a disorder), sensitivity, and specificity. Subsequent sections describe calculations and simulation studies designed to evaluate how well the coverage of each method's intervals corresponds to its nominal value.

Description of Five Methods for Constructing Intervals

Symbols and Assumptions

Let p_0 represent a population's prevalence (or pre-test probability) of the disorder D . Assume the diagnostic test has just two outcomes, "positive" and "negative." Let p_1 represent the probability of a positive test when a patient from the population has D , or the *true positive rate* or *sensitivity* of test T . Let p_2 represent the probability of a positive test when a patient from the population does not have D , or the *false positive rate* of test T ; the test's specificity will equal $1-p_2$. Bayes's Theorem then says that the probability ϕ of having D given a positive test result (sometimes called the "positive predictive value") is

$$\begin{aligned}\phi &= \frac{\text{prevalence} \cdot \text{sensitivity}}{\text{prevalence} \cdot \text{sensitivity} + [1 - \text{prevalence}] \cdot [1 - \text{specificity}]} \\ &= \frac{p_0 p_1}{p_0 p_1 + (1 - p_0) p_2}\end{aligned}\tag{1}$$

Clinicians and researchers cannot know the true values of p_0, p_1, p_2 , but they often can estimate these values — $\hat{p}_0, \hat{p}_1, \hat{p}_2$ — from empirical data on disorder prevalence and test accuracy. The prevalence is estimated by $\hat{p}_0 = x_0/n_0$, where n_0 is the number of subjects evaluated in a "random" sample from the relevant population, and x_0 is the number of subjects who had D ; test sensitivity is estimated by $\hat{p}_1 = x_1/n_1$, where n_1 equals the number of randomly chosen subjects known to have D who underwent testing with T , and x_1 is the number of these subjects who had positive test results; and the false positive rate is estimated by $\hat{p}_2 = x_2/n_2$, where n_2 is the number of randomly chosen subjects known not to have D who were tested, and x_2 is the number of these subjects who had positive test results. Dr. Jones estimated $\hat{\phi}$, Mr. Smith's chance of having disorder D , as follows:

$$\hat{\phi} = \frac{\hat{p}_0 \hat{p}_1}{\hat{p}_0 \hat{p}_1 + (1 - \hat{p}_0) \hat{p}_2} = \frac{\frac{x_0}{n_0} \frac{x_1}{n_1}}{\frac{x_0}{n_0} \frac{x_1}{n_1} + \left(1 - \frac{x_0}{n_0}\right) \frac{x_2}{n_2}} = \frac{\frac{10}{50} \frac{36}{40}}{\frac{10}{50} \frac{36}{40} + \left(1 - \frac{10}{50}\right) \frac{4}{40}} = 0.692\tag{2}$$

1 The following subsections describe five ways to produce intervals for ϕ using x_0/n_0 ,
 2 x_1/n_1 , and x_2/n_2 . The first of these procedures — the objective Bayesian method — is the
 3 one that we ultimately recommend.

4 *Objective Bayesian Method*

5 In general, the Bayesian approach assigns prior distributions to p_0 , p_1 , and p_2 , from
 6 which one determines a posterior distribution and interval for ϕ . The structure of the prior
 7 distributions might reflect physicians' or investigators' knowledge of (or "subjective" beliefs
 8 about) the nature of the population or the test data, in which case the analysis would be
 9 termed a "subjective Bayesian" analysis. Often, however, intervals for ϕ will be derived
 10 solely from the data, in which case the appropriate Bayesian approach is "objective
 11 Bayesian" analysis. Objective Bayesians choose "noninformative" or default prior
 12 distributions for the parameters, prior distributions that depend only on the data model
 13 under consideration, and not subjective beliefs.

14 The intervals for ϕ that result from the Bayesian approach are usually called
 15 "credible intervals," the distinction in terminology reflecting a distinction in interpretation:
 16 for a 95% credible interval computed in a specific application, one can assert that ϕ has a
 17 95% chance of being in the interval. In contrast, 95% frequentist coverage means only that,
 18 in repeated use, a confidence procedure will result in intervals that contain ϕ 95% of the
 19 time; one cannot assert confidence for specific instances. The difference can be quite
 20 dramatic and have important practical consequences.⁶ However, to avoid engaging in a
 21 debate concerning the Bayesian versus frequentist approaches, we will not further discuss
 22 this here. Instead, we will simply view the credible intervals developed from the objective
 23 Bayesian approach as frequentist confidence intervals to be considered for their frequentist
 24 coverage properties. A substantial literature⁷ suggests that the objective Bayesian
 25 approach is the best way to develop intervals with good frequentist coverage in complex
 26 situations.

27 The beta distribution is useful^{1,8-10} for describing the *a priori* distribution of the
 28 parameter p , arising as the probability of a success on an individual trial, of the binomial
 29 distribution. The density of a $Beta(a,b)$ distribution is given by

$$\beta_{(a,b)}(p) = C p^{a-1} (1-p)^{b-1} \quad (3)$$

30 for $0 < p < 1$, and zero otherwise.¹⁰ The constant C is chosen so the total area under the
 31 distribution equals 1.

32 Objective Bayesians have primarily used two possible "noninformative" prior
 33 distributions for p . Laplace¹¹ utilized the "uniform" prior distribution, $Beta(1,1)$, while
 34 Jeffreys¹² proposed using the $Beta(1/2,1/2)$ prior distribution — now called the "Jeffreys prior"
 35 — for a proportion p . The Jeffreys prior is known to give confidence sets with very good
 36 coverage properties for a single proportion p .^{13,14} In this article, we study its use for
 37 multiple proportions; in particular, we use it as the prior distribution for each of the

1 proportions p_i found in Equation 1.

2 After observing (x_i, n_i) , the posterior distribution of p_i is the $Beta(x_i + 1/2, n_i - x_i + 1/2)$
 3 distribution. Kleiter has described a complicated formula that uses a weighted-sum of beta
 4 distributions to produce the posterior distribution for the Bayes parameter ϕ .⁹ However, it
 5 is much easier to obtain the desired confidence interval for ϕ by working directly with the
 6 $Beta(x_i + 1/2, n_i - x_i + 1/2)$ distributions via simulation to produce what we will henceforth call
 7 “ ϕ_β intervals.” Indeed, the following simple Monte Carlo procedure will yield an estimate of
 8 the “equal-tailed” $100(1-\alpha)\%$ confidence set for ϕ . One uses a computer to:

- 9 (1) draw, at random, values from each of the $Beta(x_i + 1/2, n_i - x_i + 1/2)$ distributions,
 10 $i=0,1,2$;
- 11 (2) combine these three values in Equation 1 to obtain a value of ϕ ;
- 12 (3) repeat this process N times (where N is a large number, say, 10,000) to generate N
 13 values of ϕ , and sort these values from lowest to highest;
- 14 (4) find the integers nearest to $N\alpha/2$ and $N(1-\alpha/2)$, and choose the corresponding
 15 sorted values of ϕ as the lower and upper limits of the confidence set for ϕ .

16 Commonly available spreadsheet programs contain random number generators and
 17 inverse beta distribution commands, and with a little practice, it takes just a few seconds to
 18 create a spreadsheet that performs the just-described procedure. For the Smith–Jones
 19 example from the introduction, the 95% ϕ_β interval is [0.431, 0.887].

21 Two Log-Odds Methods

22 If one divides Equation 1 by $1-\phi$ and rearranges terms, the resulting version of
 23 Bayes’s Theorem expresses O_{post} — the post-test odds of having D , given a positive test
 24 result — as the product of the pre-test odds, O_{ante} , and the likelihood ratio, LR :

$$O_{post} = \frac{\phi}{1-\phi} = \frac{P_0}{1-p_0} \times \frac{P_1}{P_2} = O_{ante} \times LR \quad (4)$$

25 The likelihood ratio, O_{ante} , and O_{post} can range from zero to infinity. Symmetric
 26 confidence intervals for LR constructed from \hat{p}_1 and \hat{p}_2 are unlikely to be accurate and may
 27 include values less than zero.^{2,3,15} One commonly-invoked remedy^{3,16} is to construct
 28 confidence intervals for LR using its natural logarithm, $\ln LR$. Extending this reasoning to
 29 O_{post} , one would obtain upper and lower bounds of the $100(1-\alpha)\%$ interval, \hat{O}_{post}^U and \hat{O}_{post}^L ,
 30 from

$$\hat{O}_{post}^{U,L} \approx e^{\ln\left(\frac{\hat{\phi}}{1-\hat{\phi}}\right) \pm z_{\alpha/2} \sqrt{\text{var}(\ln \hat{O}_{post})}} \quad (5)$$

1 where $z_{\alpha/2}$ is the normal deviate of the upper $\alpha/2$ percentile of the normal distribution, and

$$\text{var}(\ln \hat{O}_{post}) \approx \frac{n_0}{x_0(n_0 - x_0)} + \frac{1}{x_1} - \frac{1}{n_1} + \frac{1}{x_2} - \frac{1}{n_2} \quad (6)$$

2 Equation 6 arises from application of the standard “delta method” together with some
3 algebra.¹⁷ One uses the relationship $\phi = \frac{\hat{O}_{post}}{1 + \hat{O}_{post}}$ to convert $[\hat{O}_{post}^U, \hat{O}_{post}^L]$ to $[\hat{\phi}^U, \hat{\phi}^L]$,
4 the interval for ϕ .

5 Equations 5 and 6 are undefined if $x_0, x_1,$ or x_2 equal zero, or if $x_0 = n_0$; also, if $x_1 = n_1$
6 and $x_2 = n_2$, the likelihood ratio makes no contribution to $\text{var}(\ln \hat{O}_{post})$. Here are two simple
7 ways to address this.

- 8 (1) When calculating $\hat{\phi}$ in Equation 5 and \hat{O}_{post}^U and \hat{O}_{post}^L in Equation 6, substitute $1/2$
9 for $x_0, x_1,$ or x_2 whenever these terms equal zero, and substitute $(n_i - 1/2)$ for x_i
10 whenever $x_i = n_i$. Hereinafter, these “substituted” results will be designated “ ϕ_S
11 intervals.” For the Smith–Jones example, the 95% ϕ_S interval is [0.413, 0.878].
- 12 (2) Adapt the *LR* variance-approximation procedure of Walter¹⁸ and by Pettigrew and
13 colleagues¹⁹ to the present problem; that is, estimate O_{post} from

$$\ln(\hat{O}_{post}) \approx \ln\left(\frac{x_0 + 1/2}{n_0 - x_0 + 1/2}\right) + \ln\left(\frac{x_1 + 1/2}{n_1 + 1/2}\right) - \ln\left(\frac{x_2 + 1/2}{n_2 + 1/2}\right) \quad (7)$$

14 calculate $\text{var}(\ln \hat{O}_{post})$ from

$$\text{var}(\ln \hat{O}_{post}) \approx \frac{n_0 + 1/2}{(x_0 + 1/2)(n_0 - x_0 + 1/2)} + \frac{1}{x_1 + 1/2} - \frac{1}{n_1 + 1/2} + \frac{1}{x_2 + 1/2} - \frac{1}{n_2 + 1/2} \quad (8)$$

17

18

19

20

and use these results in Equation 5. Intervals obtained using this approach are
henceforth referred to as “ $\phi_{1/2}$ intervals.” For the Smith–Jones example, the 95% $\phi_{1/2}$
interval is [0.410, 0.864].

1 Notice that one can quickly compute ϕ_S intervals and $\phi_{1/2}$ intervals using an inexpensive
2 pocket calculator.

3 *A Method Using the Gart-Nam Approximation for the Likelihood Ratio*

4 Gart and Nam² describe an iterative procedure for calculating the $100(1-\alpha)\%$
5 confidence interval for LR using likelihood-based estimates supplemented by skewness
6 correction derived from Bartlett's score method.^{20,21} Centor has discussed the advantages
7 of the Gart-Nam procedure and provides computer code that performs it,²² and the present
8 authors have developed a spreadsheet that accomplishes the same task.

9 One can include Gart and Nam's interval-estimation procedure in a Monte Carlo
10 simulation similar to that described for the objective Bayesian method to obtain what we
11 shall term " ϕ_{GN} intervals" for ϕ . Instead of sampling from the beta distributions for p_1 , and
12 p_2 , one uses the formulation of Bayes's Theorem in Equation 4, and samples from the
13 distribution for LR implied by the Gart-Nam algorithm. One obtains a set of N values of
14 \hat{O}_{post} that are readily converted to N values of ϕ , the central $N(1-\alpha)$ of which will represent
15 the $100(1-\alpha)\%$ confidence interval for ϕ . For the Smith-Jones example, the 95% ϕ_{GN}
16 interval is [0.439, 0.891].

17 *A Delta Method Derived from Kleiter*

18 In addition to providing an analytic expression for the posterior distribution of ϕ ,
19 Kleiter suggests approximating this posterior using a secondary beta distribution, with
20 parameters a and b estimated from the data utilizing the delta method. Equal-tailed
21 $100(1-\alpha)\%$ confidence intervals can then be found using the inverse beta distribution
22 function available on many spreadsheets. The Appendix describes the derivation and
23 details of this interval estimation method. Hereinafter, the results obtained with this delta
24 method will be referred to as " ϕ_δ values." For the Smith-Jones example, the 95% ϕ_δ
25 interval is [0.419, 0.888].

27 **Comparing Intervals**

28 The first issue of interest to practitioners is the typical size of the confidence intervals
29 and the magnitude of the differences that can be expected from the five just-described
30 procedures. Table 1 gives the 95% confidence intervals for these procedures under nine
31 possible data sets. These data sets utilize values for \hat{p}_0 , \hat{p}_1 , and \hat{p}_2 and sample sizes that
32 physicians typically encounter in clinical medicine. The most important fact to note from
33 Table 1 is simply that the 95% confidence intervals for ϕ are quite large, even for the
34 moderately large sample sizes ($n_i=80$). The impact of the uncertainty in ϕ can thus be
35 quite significant.

36 *[place Table 1 about here]*

37 The top portion of Table 1 shows that, with moderately large samples and when \hat{p}_0 ,
38 \hat{p}_1 , and \hat{p}_2 are not close to zero or 1, the methods produce very similar intervals; the

1 differences usually would be too small to influence a practical decision-making task. When
 2 \hat{p}_0 , \hat{p}_1 , and/or \hat{p}_2 equal $1/10$ or $9/10$, however, the intervals differ; their dissimilarities also
 3 become more salient when the data samples are small (middle portion of Table 1, $n_i=20$).
 4 The bottom portion of Table 1 shows how the five methods behave when the sample size is
 5 small and either \hat{p}_0 or \hat{p}_2 equals 0; here, the differences in calculated intervals become
 6 pronounced. In particular, the objective Bayesian intervals are considerably shorter, a
 7 desirable property that might have significant clinical consequences.

8 The other key feature of the confidence intervals that needs to be ascertained is the
 9 extent to which they meet their nominal coverage goals. If one asserts, for example, that
 10 $[\hat{\phi}^L, \hat{\phi}^U]$ is the central or “equal-tailed” 95% confidence interval for ϕ , one is saying that
 11 there is a 2.5% chance that ϕ falls below the interval and a 2.5% chance that ϕ falls above
 12 the interval.¹² So, one can evaluate a particular interval-generating procedure by
 13 computing the actual percentage of the time that, in repeated use, the procedure produces
 14 a confidence interval that lies above or below the actual ϕ .

15 To evaluate the coverage of the $\hat{\phi}_\beta$ and $\hat{\phi}_{GN}$ interval estimation procedures, we
 16 choose triads of p_0 , p_1 , and p_2 which yield a known value of ϕ , and choose values of n_0 , n_1 ,
 17 and n_2 to represent sample sizes from the population and subpopulations relevant to our
 18 inferences about p_0 , p_1 , and p_2 . We then “simulate” research studies with a computer,
 19 making a large number J of random drawings of sizes n_0 , n_1 , and n_2 from the populations.
 20 Each of the J random drawings yields data from which we can construct confidence
 21 intervals for ϕ , and we then see how often the J intervals obtained with each method
 22 undershoot or exceed the true value of ϕ . A “good” interval construction method will
 23 produce $100(1-\alpha)\%$ intervals that correspond to their nominal value, so that roughly $J\alpha/2$
 24 intervals will undershoot ϕ and roughly $J\alpha/2$ intervals will exceed ϕ .²³

25 For the $\hat{\phi}_{1/2}$, $\hat{\phi}_S$, and $\hat{\phi}_\delta$ methods, which use explicit formulae to calculate intervals,
 26 one can use an analytic evaluative approach based on the binomial probabilities of the data
 27 that correspond to the given values of p_0 , p_1 , and p_2 . The probability of each possible triad
 28 of data is simply the product of the three respective binomial probabilities. Thus, for each
 29 possible data triad, one can simply compute the confidence interval and note whether the
 30 interval exceeds or lies below the true value of ϕ (computed by Equation 1). Summing the
 31 data triad probabilities over all those data for which the interval exceeds or lies below ϕ
 32 yields the desired information about interval coverages properties.

33 Table 2 provides coverage results for three “scenarios” involving hypothetical,
 34 known populations for which the values of the triplet (p_0, p_1, p_2) were $(1/4, 3/4, 1/4)$, $(1/10, 9/10, 1/10)$,
 35 or $(1/2, 9/10, 1/10)$, the same values used in the upper and middle parts of Table 1. For the
 36 simulation studies of the $\hat{\phi}_\beta$ and $\hat{\phi}_{GN}$ methods, $J=200,000$ random draws of 20, 40, 80, or
 37 160 were made from each of the three “populations” whose characteristics are represented
 38 by p_0 , p_1 , and p_2 , so that $n_0=n_1=n_2$ always equaled 20, 40, 80, or 160, but x_{0j} , x_{1j} , and x_{2j}
 39 varied randomly with each draw. The values chosen for n_i allowed us to explore the small-
 40 and larger-sample performance behavior of each interval construction method. For the $\hat{\phi}_{1/2}$,

ϕ_S , and ϕ_δ methods, we assumed the same “population” characteristics, but made exact computations according to the above-discussed analytic calculation procedure.

[place Table 2 about here]

The results from evaluating the five methods appear in Table 2, the interpretation of which is most easily explained through an example. The six numbers in the first row of the first column describe the $\phi_{1/2}$ method’s performance when 20 “subjects” are randomly sampled from a population for which $p_0=1/4$, $p_1=3/4$, and $p_2=1/4$. If $n_0=n_1=n_2=20$, there are $(20+1)^3 = 9,261$ possible combinations of samples, and for each possible combination, the 90%, 95%, and 99% confidence intervals for ϕ (i.e., intervals with $\alpha=0.10$, $\alpha=0.05$, and $\alpha=0.01$) were calculated, along with the probability of obtaining each of those intervals.*

Because $\phi = \frac{(1/4)(3/4)}{(1/4)(3/4) + (1-1/4)(1/4)} = 1/2$, we hope to find that the lower limit of the 99%

intervals is *greater* than $1/2$ about 0.5% of the time, and that the upper limit of the 99% intervals is *less* than $1/2$ about 0.5% of the time. The results in Table 2 show, however, that only 0.12% of the $\phi_{1/2}$ method’s 99% intervals would be expected to lie above $1/2$ and that only 0.03% of the intervals would lie below $1/2$. That is, when $n_0=n_1=n_2=20$, the $\phi_{1/2}$ method’s 99% intervals are too wide, as are the 95% and 90% intervals.

The next three rows show that as the n_i increase, the $\phi_{1/2}$ intervals get closer to the nominal values. In the fifth row, however, which contains results involving reference populations where $p_0=1/10$, $p_1=9/10$, and $p_2=1/10$, the $\phi_{1/2}$ method did poorly again. Especially when the populations are small, the $\phi_{1/2}$ method produced confidence intervals that were too broad or badly skewed.

Why would these problems occur? The answer comes from using Table 2 to compare the results for the $\phi_{1/2}$ method with those for the ϕ_β method, and then examining the lower portions of Table 1:

- Table 2 shows that intervals using the ϕ_β method come fairly close to what we would desire given our earlier interpretation of the assertion “ (ϕ^L, ϕ^U) is the central $100(1-\alpha)\%$ confidence interval for ϕ ,” even when sample sizes are small and the reference populations’ values of p_0 , p_1 , or p_2 are $1/10$ or $9/10$. When, for example, $p_0=1/10$ and $n_0=20$, approximately 12.2% of the 20-member draws from the population would contain *no* members with the disorder, so that $\hat{p}_0=x_0/n_0=0$. When random 20-member drawings are made simultaneously from populations in

* To illustrate: in random independent samplings from a population in which $p_0=1/4$, $p_1=3/4$, and $p_2=1/4$, the probability of obtaining $\hat{p}_0=x_0/n_0=4/20$ is 0.190, the probability of obtaining $\hat{p}_1=x_1/n_1=17/20$ is 0.134, and the probability of obtaining $\hat{p}_2=x_2/n_2=4/20$ is 0.202; the joint probability of these three results is 0.00514. Using the $\phi_{1/2}$ method, these results yield $\hat{\phi}=0.465$, with a 95% interval of [0.1965, 0.7548].

1 which $p_0=1/10$, $p_1=9/10$, and $p_2=1/10$, the probability of getting a draw in which either
 2 $\hat{p}_0=x_0/n_0=0$, $\hat{p}_1=x_1/n_1=1$, or $\hat{p}_2=x_2/n_2=0$, or some combination thereof occurs, is
 3 $1-[1-(0.9)^{20}]^3=0.322$. This means that, for small samples, the coverage properties
 4 are very strongly influenced by the difficult-to-handle cases where $\hat{p}_0=0$, $\hat{p}_1=1$,
 5 and/or $\hat{p}_2=0$.[†]

- 6 • The lower portion of Table 1 shows that when $x_0/n_0=0/20$ or $x_2/n_2=0/20$, the ϕ_β
 7 intervals are quite different from — and, as was mentioned earlier, narrower than —
 8 the $\phi_{1/2}$ intervals. Thus, the difference in how the two methods handle these types of
 9 cases probably accounts for the major differences in the two procedures' coverage
 10 properties.

11 Taken as a whole, the results in Table 2 show that the objective Bayesian method is
 12 the best performer. Even with the smaller samples, the ϕ_β intervals are consistently near
 13 the desired frequency of under- and over-shooting ϕ , and there is little systematic tendency
 14 for ϕ_β intervals to be biased downward or upward. The ϕ_{GN} method is the “runner up”:
 15 although it yields intervals that usually are as close to nominal values as the ϕ_β intervals, it
 16 is much more difficult to implement and also typically gives larger intervals. Of the three
 17 methods that make use of explicit formulae, the ϕ_δ method does best, though it sometimes
 18 does poorly with small samples. Both log-odds methods do well in the larger-sample
 19 evaluations when reference probabilities do not lie close to 0 or 1. When these conditions
 20 do not obtain, however, the ϕ_S and $\phi_{1/2}$ intervals usually are too broad and tend to be
 21 biased downward.

22 How might the ϕ_β method work under “extreme” circumstances, *i.e.*, in situations
 23 where reference probabilities lie very close to 0 or 1? Of course, there are an infinite
 24 number of combinations of such circumstances, but the lowest portion of Table 2 gives us
 25 a flavor of the method's limitations. This portion summarizes findings from 200,000
 26 simulations when $p_0=1/100$, $p_1=99/100$, and $p_2=1/100$. Unless the sample size is nearly 200, the
 27 coverage of this procedure fails to attain the nominal level. Interpretation of this finding is
 28 far from clear, however. The problem could be the rather simplistic use, in this multivariate
 29 situation, of the *Beta*($1/2, 1/2$) priors, which were designed for binomial parameters alone.
 30 The objective Bayesian approach can, in principle, be used to determine that prior which
 31 would provide intervals having optimal frequentist coverage, but determination of this prior
 32 would be a formidable undertaking and would result in a much more complicated
 33 procedure. Furthermore, in extreme cases where the parameters are close to the boundary
 34 of the parameter space, the difference between the Bayesian and the frequentist

[†] When $p_0=1/10$, $p_1=9/10$, $p_2=1/10$, and the sample size is 40, the chance that random drawing will lead to $\hat{p}_0=0$, $\hat{p}_1=1$, and/or $\hat{p}_2=0$ is $1-[1-(0.9)^{40}]^3=0.0437$. This suggests that the intervals constructed under these circumstances figure importantly in the results concerning 40-member reference populations. When the sample size is 80, however, the chance is only $1-[1-(0.9)^{80}]^3=6.55 \times 10^{-4}$.

1 interpretations of confidence intervals becomes highly significant and cannot be ignored in
2 practice.

4 A Medicolegal Example

5 Malingering is a commonly-encountered phenomenon in medical evaluations that
6 occur in legal contexts, where being (or appearing to be) ill may lead to financial
7 compensation or permit an individual to avoid legal responsibility.²⁴ For example, studies
8 and surveys suggest that roughly one-sixth of persons evaluated concerning their
9 competence to stand trial (CST) feign incapacitating mental disorders.²⁵⁻²⁷ Each year, an
10 estimated 30,000–60,000 CST evaluations are performed each year in the United
11 States,^{28,29} which implies that U.S. mental health professionals evaluate 5,000–10,000
12 defendants each year who are feigning incompetence to stand trial.

13 Mossman and Hart have suggested that by using Bayes's Theorem, mental health
14 professionals could improve their interpretations of data on malingering and make better
15 presentations of evidence on malingering in court.³⁰ Rogers and Salekin²⁷ argue against
16 doing this: they assert that imprecision in the estimated "base rate" of malingering (in this
17 article's symbols, \hat{p}_0) and the estimated accuracy of malingering instruments (*i.e.*, \hat{p}_1 and
18 \hat{p}_2) would result in confidence intervals for the probability of malingering (*i.e.*, ϕ) that were
19 so broad as to make Bayesian estimates useless, misleading, and unfair to evaluatees.

20 The preceding sections give us the tools to evaluate Rogers and Salekin's claim.
21 One of the malingering measures discussed by Mossman and Hart is a simple "Symptom
22 Validity Test" (SVT) for detecting feigned memory deficits. In one study describing the
23 SVT, 17 of 20 malingering simulators made four or more errors, but only one of 40
24 genuine patients did this poorly;³¹ based on this small study then, $\hat{p}_1 = x_1/n_1 = 17/20$ and
25 $\hat{p}_2 = x_2/n_2 = 1/40$. As an example of how one might estimate p_0 , we use results reported by
26 Gothard and colleagues,²⁵ who examined findings in 55 CST evaluatees and determined that
27 25 were competent, 23 were incompetent, and seven were malingering. Because the
28 purpose of testing for malingered incompetence involves distinguishing defendants who are
29 feigning illness from defendants who are truly incompetent, the relevant portion of the
30 sample does not include the 25 defendants who were deemed competent; thus
31 $\hat{p}_0 = x_0/n_0 = 7/30$.

32 We now can estimate ϕ and construct a confidence interval. Suppose a defendant-
33 evaluatee has a "positive" result on the SVT, *i.e.*, has made four or more errors on the test.
34 (Using Equation 2, we find that $\hat{\phi} = 0.912$, and the ϕ_β method gives a 95% confidence
35 interval for ϕ of [0.633, 0.990]. This interval is too broad to let us say that a positive SVT is
36 "proof beyond a reasonable doubt" of malingering. In U.S. jurisdictions, however,
37 adjudicatory competence typically is assumed unless a preponderance of the evidence
38 shows that a defendant's mental condition makes him unable to proceed.^{32,33} Thus, even
39 with the small sample sizes used in this example, the SVT's confidence interval is narrow
40 enough to support testimony that the defendant probably was malingering.

Discussion

Some version of Equation 1 is a commonplace in texts on medical decision-making³⁴ and Bayesian statistics.^{6,8} Medical articles commonly report tests' positive and negative predictive values, but medical publications rarely (if ever) report confidence intervals for these values. Although the probabilities used to calculate positive and negative predictive values are themselves imprecise, little systematic attention has been devoted to the problem of "propagating" these imprecisions⁹ to generate confidence intervals for post-test probabilities.

This article suggests that investigators and clinicians can construct such intervals in several ways. Intervals based on the natural logarithm of the posterior odds can be calculated quickly using inexpensive hand-held calculators; these intervals are reasonably accurate when the estimates of pre-test probabilities, sensitivity, and specificity are based on large numbers of subjects and do not lie too close to either 0 or 1. When these conditions are not met, or when better accuracy is desired, investigators can produce useful confidence intervals quickly if they have access to one of the many available spreadsheets with a command that computes the inverse of the cumulative beta distribution function. The ϕ_β method should, under many (though not all) conditions, generate intervals whose coverage properties are close to the nominal value.

Acknowledgments

The authors thank the anonymous reviewers for their many helpful comments on an earlier version of this article, and German Molina for his computer assistance in performing the simulation studies. Professor Berger's research was partially supported by the National Science Foundation, Grants DMS-9802261 and DMS-0103265.

References

1. Newcombe RG. Two-sided confidence intervals for the single proportion: comparison of seven methods. *Statistics in Medicine* 1998;17:857-872.
2. Gart JJ, Nam J. Approximate interval estimation of the ratio of binomial parameters: A review and corrections for skewness. *Biometrics* 1988;44:323-338.
3. Simel DL, Samsa GP, Matchar DB. Likelihood ratios with confidence: sample size estimation for diagnostic test studies. *Journal of Clinical Epidemiology* 1991;44:763-770.
4. Agresti A. Introduction to categorical data analysis. New York: Wiley, 1996.

- 1 5. Fleiss JL, Davies M. Jackknifing functions of multinomial frequencies, with an
2 application to a measure of concordance. *American Journal of Epidemiology*
3 1982;115:841-845.
- 4 6. Berger JO. *Statistical decision theory and Bayesian analysis*, 2nd edition. New York:
5 Springer-Verlag, 1985.
- 6 7. Yang R, Berger JO: A catalog of noninformative priors. Viewable at
7 www.stat.duke.edu/~berger/papers/catalog.html.
- 8 8. Iverson GR. *Bayesian Statistical Inference*. Newbury Park, CA: Sage Publications,
9 1984.
- 10 9. Kleiter GD. Propagating imprecise probabilities in Bayesian networks. *Artificial*
11 *Intelligence* 1996;88:143-161.
- 12 10. Freund JE, Williams FJ. *Dictionary/Outline of Basic Statistics*. Mineola, NY: Dover,
13 1991.
- 14 11. Laplace PS. *Théorie Analytique des Probabilités*. Paris: V. Courcier, 1812.
- 15 12. Jeffreys H. *Theory of probability* (3rd edition). Oxford: Clarendon Press, 1961.
- 16 13. Brown L, Cai T, DasGupta A. Confidence intervals for a binomial proportion and
17 edgeworth expansions. *Annals of Statistics* (in press, 2001).
- 18 14. Brown L, Cai T, DasGupta A. Interval estimation for a binomial proportion.
19 *Statistical Science* (in press, 2001)
- 20 15. Peirce JC, Cornell RG. Integrating stratum-specific likelihood ratios with the analysis
21 of ROC curves. *Med Decis Making* 1993;13:141-51.
- 22 16. Katz D, Baptista J, Azen SP, Pike MC. Obtaining confidence intervals for the risk
23 ratio in cohort studies. *Biometrics* 1978;34:469-474.
- 24 17. Lehmann E, Casella G. *Theory of Point Estimation*. New York: Springer, 1998.
- 25 18. Walter SD. The distribution of Levin's measure of attributable risk. *Biometrika*
26 1975;62:371-375.
- 27 19. Pettigrew HM, Gart JJ, Thomas DG. The bias and higher cumulants of the
28 logarithm of a binomial variate. *Biometrika* 1986;73:425-435.
- 29 20. Bartlett MS. Approximate confidence intervals, II. More than one unknown
30 parameter. *Biometrika* 1953;40:306-317.
- 31 21. Bartlett MS. Approximate confidence intervals, III. A bias correction. *Biometrika*
32 1955;43:993-998.
- 33 22. Centor RM. Estimating confidence intervals of likelihood ratios. *Med Decis Making*
34 1992;12:229-233.

- 1 23. Sims CA, Zha T. Error Bands for Impulse Responses. Federal Reserve Bank of
2 Atlanta Working Paper 95-6, September 1995.
- 3 24. American Psychiatric Association. Diagnostic and Statistical Manual of Mental
4 Disorders, Fourth Edition. Washington, DC: American Psychiatric Association,
5 1994.
- 6 25. Gothard S, Rogers R, Sewell KW. Feigning incompetency to stand trial: An
7 investigation of the Georgia Court Competency Test. *Law and Human Behavior*
8 1995;19:363-373
- 9 26. Rogers R, Sewell KW, Goldstein A. Explanatory models of malingering: A
10 prototypical analysis. *Law and Human Behavior* 1994;18:543-552.
- 11 27. Rogers R, Salekin RT. Beguiled by Bayes: A reanalysis of Mossman and Hart's
12 estimates of malingering. *Behavioral Sciences and the Law*. 1998;16:147-153.
- 13 28. Nicholson RA, Kugler KE. Competent and incompetent criminal defendants: a
14 quantitative review of comparative research. *Psychological Bulletin* 1991;109:355-
15 370.
- 16 29. Poythress NG, Nicholson R, Otto RK, Edens JF, Bonnie RJ, Monahan J, Hoge SK.
17 The MacArthur Competence Assessment Tool—Criminal Adjudication. Odessa, FL:
18 Psychological Assessment Resources, 1999.
- 19 30. Mossman D, Hart KJ. Presenting evidence of malingering to courts: insights from
20 decision theory. *Behavioral Sciences and the Law*. 1996;14:271-291.
- 21 31. Guilmette TJ, Hart KJ, Guiliano AJ, Leininger BE. A comparison of the Fifteen Item
22 Test and the Hiscock Forced-Choice Procedure in detecting simulated memory
23 impairment. *Clinical Neuropsychologist*. 1994;8:283-294.
- 24 32. *Cooper v. Oklahoma*, 517 U.S. 348 (1996).
- 25 33. Ohio Revised Code § 2945.37(G).
- 26 34. Sox HC, Blatt MA, Higgins MC, Marton KI. *Medical decision making*. Boston:
27 Butterworths, 1988.

Appendix: Calculation of a and b for Kleiter's Delta Method

For a beta distribution with parameters a and b , the mean μ is $\frac{a}{a+b}$ and the variance σ^2 is $\frac{ab}{(a+b)^2(a+b+1)}$.¹² Some algebra gives these solutions for a and b :

$$a = \mu \left[\frac{\mu(1-\mu)}{\sigma^2} - 1 \right]; \quad b = (1-\mu) \left[\frac{\mu(1-\mu)}{\sigma^2} - 1 \right] \quad (\text{A1})$$

To obtain estimates of these values from the data, one can use $\hat{\phi}$ from Equation 2 as the estimate of the mean μ , and the variance σ^2 can be approximated using the delta method, namely:

$$\text{var}(\hat{\phi}) \approx \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_0} \right)^2 \text{var}(\hat{p}_0) + \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_1} \right)^2 \text{var}(\hat{p}_1) + \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_2} \right)^2 \text{var}(\hat{p}_2) \quad (\text{A1})$$

Differentiating Equation 1 with respect to p_0 , p_1 , and p_2 yields the partial derivatives

$$\begin{aligned} \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_0} \right) &= \frac{\hat{p}_1 \hat{p}_2}{[\hat{p}_0 \hat{p}_1 + (1 - \hat{p}_0) \hat{p}_2]^2} \\ \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_1} \right) &= \frac{(1 - \hat{p}_0) \hat{p}_0 \hat{p}_2}{[\hat{p}_0 \hat{p}_1 + (1 - \hat{p}_0) \hat{p}_2]^2} \\ \left(\frac{\partial \hat{\phi}}{\partial \hat{p}_2} \right) &= \frac{-(1 - \hat{p}_0) \hat{p}_0 \hat{p}_1}{[\hat{p}_0 \hat{p}_1 + (1 - \hat{p}_0) \hat{p}_2]^2} \end{aligned} \quad (\text{A3})$$

Three adjustments help one avoid undefined terms and improve this method's results:

- when $x_i=0$, change x_i to $1/2$;
- when $x_i=n_i$, change x_i to $n_i - 1/2$;
- add $1/2$ to the values of a and b calculated from Equations A1–A3, so that one calculates the $100(1-\alpha)\%$ interval from the inverse of the $Beta(a+1/2, b+1/2)$ distribution.

Table 1. — Results from five methods of constructing 95 percent intervals for ϕ . For explanation of symbols, see text.

Method	Data and sample sizes		
	$x_0/n_0, x_1/n_1, x_2/n_2$ 8/80, 72/80, 8/80	$x_0/n_0, x_1/n_1, x_2/n_2$ 20/80, 60/80, 20/80	$x_0/n_0, x_1/n_1, x_2/n_2$ 40/80, 72/80, 8/80
$\phi_{1/2}$ (log odds, add $1/2$)	[0.264, 0.714]	[0.341, 0.651]	[0.797, 0.949]
ϕ_β (objective Bayesian)	[0.264, 0.726]	[0.346, 0.658]	[0.808, 0.952]
ϕ_S (log odds, substituted)	[0.272, 0.728]	[0.344, 0.656]	[0.803, 0.952]
ϕ_{GN} (Gart-Nam)	[0.284, 0.740]	[0.349, 0.662]	[0.807, 0.955]
ϕ_δ (delta method)	[0.266, 0.734]	[0.343, 0.657]	[0.812, 0.955]
	$x_0/n_0, x_1/n_1, x_2/n_2$ 2/20, 18/20, 2/20	$x_0/n_0, x_1/n_1, x_2/n_2$ 5/20, 15/20, 5/20	$x_0/n_0, x_1/n_1, x_2/n_2$ 10/20, 18/20, 2/20
$\phi_{1/2}$ (log odds, add $1/2$)	[0.112, 0.842]	[0.210, 0.769]	[0.632, 0.970]
ϕ_β (objective Bayesian)	[0.107, 0.872]	[0.216, 0.780]	[0.679, 0.980]
ϕ_S (log odds, substituted)	[0.122, 0.878]	[0.216, 0.784]	[0.648, 0.978]
ϕ_{GN} (Gart-Nam)	[0.149, 0.902]	[0.241, 0.795]	[0.699, 0.982]
ϕ_δ (delta method)	[0.099, 0.901]	[0.208, 0.792]	[0.695, 0.982]
	$x_0/n_0, x_1/n_1, x_2/n_2$ 0/20, 18/20, 2/20	$x_0/n_0, x_1/n_1, x_2/n_2$ 0/20, 15/20, 5/20	$x_0/n_0, x_1/n_1, x_2/n_2$ 10/20, 18/20, 0/20
$\phi_{1/2}$ (log odds, add $1/2$)	[0.009, 0.799]	[0.004, 0.570]	[0.675, 0.998]
ϕ_β (objective Bayesian)	[0.000, 0.649]	[0.000, 0.326]	[0.857, 1.000]
ϕ_S (log odds, substituted)	[0.010, 0.837]	[0.004, 0.588]	[0.570, 1.000]
ϕ_{GN} (Gart-Nam)	[0.032, 0.824]	[0.014, 0.507]	[0.701, 0.984]
ϕ_δ (delta method)	[0.006, 0.833]	[0.003, 0.438]	[0.695, 0.982]