

LIKELIHOOD RATIO INTERVALS WITH BAYESIAN TREATMENT OF UNCERTAINTIES: COVERAGE, POWER AND COMBINED EXPERIMENTS

J. CONRAD

CERN, PH-EP Dept., CH-1211 Geneva 23, Switzerland
(now at: Royal Institute of Technology (KTH), Particle and Astroparticle Physics,
AlbaNova University Center, SE-10691 Stockholm, Sweden)
E-mail: Jan.Conrad@cern.ch

F. TEGENFELDT

Iowa State University
Ames, IA 5011-3160, USA
E-mail: Fredrik.Tegenfeldt@cern.ch

In this note we present studies of coverage and power for confidence intervals for a Poisson process with known background calculated using the Likelihood ratio (aka Feldman & Cousins) ordering with Bayesian treatment of uncertainties in nuisance parameters. We consider the variant where the Bayesian integration is performed in both the numerator and the denominator, and also the modification where the integration is done only in the numerator whereas in the denominator the likelihood is taken at the maximum likelihood estimate of the parameters. Furthermore we discuss how measurements can be combined in this framework and give an illustration with limits on the branching ratio of a rare B-meson decay recently presented by CDF/D0. A set of C++ classes has been developed which can be used to calculate confidence intervals for single or combining multiple experiments using the above algorithms and considering a variety of parameterizations to describe the uncertainties.

1. Introduction

A popular technique to calculate confidence intervals in recent years is the one suggested by Feldman & Cousins¹. The method consists of constructing an acceptance region for each possible hypothesis (in the way proposed by Neyman²) and fixing the limits of the region by including experimental outcomes according to rank which is given by the likelihood ratio^a:

$$R(s, n)_{\mathcal{L}} = \frac{\mathcal{L}(n|s + b)}{\mathcal{L}(n|s_{best} + b)} \quad (1)$$

where s is the hypothesis, n the experimental outcome, b the expected background, s_{best} is the hypothesis most compatible with n and \mathcal{L} the Likelihood function. The expected background b is an example of a so-called *nuisance parameter*, i.e. a parameter which is not of primary interest but which still affects the calculated confidence interval. Another example of such a nuisance parameter could be the signal efficiency. In the originally proposed method by Feldman & Cousins, only the presence of background was considered and it was assumed to be

exactly known. The question on how to treat uncertainties in nuisance parameters in confidence interval calculation, in particular in context of the frequentist construction, has drawn considerable attention in the recent years. In 1992 Cousins & Highland³ proposed a method which is based on a Bayesian treatment of the nuisance parameters. The main idea is to use a probability density function (pdf) in which the average is taken over the nuisance parameter:

$$P(n|s, \epsilon) \longrightarrow \int P(n|s, \epsilon')P(\epsilon'|\epsilon)d\epsilon' := q(n|s, \epsilon) \quad (2)$$

where ϵ' is the true value of the nuisance parameter, ϵ denotes its estimate and s and n symbolize the signal hypothesis and the experimental outcome respectively.

Cousins & Highland only treated the case of Gaussian uncertainties in the signal efficiency. The method has since been generalized by Conrad et al.⁴ to operate with the Feldman & Cousins ordering scheme and taking into account both efficiency and background uncertainties as well as correlations. This generalized method has already been used in a number of particle and astroparticle physics experiments (see references in Tegenfeldt & Conrad⁵). FHC² denotes this generalized method in the remain-

^aThroughout this note we consider Poisson distributions with experimental outcome n , hypothesis parameter s and (possibly not exactly) known background b .

der of this note. If there are significantly less events than expected background, FHC² tends to result in confidence intervals which become smaller with increasing uncertainties. Hill⁶ therefore proposed a modification where the ordering of the likelihood ratio is defined as:

$$R(s, n)_{\mathcal{L}} = \frac{q(n|s + b)}{\mathcal{L}(\max(0, n_{obs} - \hat{b}) + \hat{b})} \quad (3)$$

here \hat{b} is the maximum likelihood estimate of b given the subsidiary observation of b . MBT (“Modified Bayesian Treatment”) denotes this modification in the remainder of this note.

In this contribution, we discuss coverage and power of FHC² and MBT as well as the combination of different experiments with and without correlations. We start by introducing the C++ library which has been developed to be able to do the necessary calculations.

2. POLE++

For the coverage studies presented in this paper, a reasonably fast and efficient code is required. Hence, a user-friendly and flexible C++ library of classes was developed based on the FORTRAN routine presented by Conrad⁷. The library is independent of external libraries and consists of two main classes, *Pole* and *Coverage*. The first class takes as input the number of observed events, the efficiency and background with uncertainties and calculates the limits using the method described in this paper. The integrals are solved analytically. *Coverage* generates user-defined pseudo-experiments and calculates the coverage using *Pole*. Presently the library supports Gauss, log-Normal and flat pdf for description of the nuisance parameters. Several experiments with correlated or uncorrelated uncertainties in the nuisance parameters can be combined. The pole++ library can be obtained from <http://cern.ch/tegen/statistics.html>

3. Coverage and Power

The most crucial property of methods for confidence interval construction is the coverage, which states that a fraction $(1-\alpha)$ of infinitely many repeated experiments should yield confidence intervals that include the true hypothesis irrespective of what the true hypothesis is. For confidence interval construction (according to Neyman) without uncertainties

in nuisance parameters this property is fulfilled by construction. In the present case however, we have to test the coverage employing Monte Carlo experiments.

Power on the other hand is a concept which is defined in the context of hypothesis testing: the power of a hypothesis testing method is the probability that it will reject the null hypothesis, s_0 , given that the alternative hypothesis s_{true} is true. This concept is rather difficult to generalize to confidence intervals since the alternative hypothesis is not uniquely defined. We use the following definition for power:

$$\Pi(s_{true})_{s_0} = \sum_{n \notin Acc(s_0)} q(n|s_{true}, \epsilon) \quad (4)$$

and view power as a function of s_{true} . $Acc(s_0)$ here denotes the acceptance region of s_0 . This seems an intuitively appealing measure: given the choice between different methods, the method which has minimally overlapping acceptance regions should be taken.

Typical examples of the coverage as a function of signal hypothesis are shown in Figure 1. It can be seen that the introduction of a continuous variable leads to a considerable smoothing of the coverage plot. A modest amount of over-coverage is introduced, similarly for the MBT method and the FHC² method. For high Gaussian uncertainties in efficiency ($\sim 40\%$) the over-coverage of MBT is less pronounced than that for FHC². More detailed coverage studies of the FHC² method have been presented by Tegenfeldt & Conrad⁵. The power of the FHC² and MBT methods is compared in Figure 1 for 40% uncertainties in the efficiency. FHC² has higher power for hypotheses rather far away from the null hypotheses. This is true only for large signals and comparably large uncertainties (and for not too large differences between s_0 and s_{true}), otherwise differences are negligible.

4. Combining Different Experiments

Combination of experiments can be divided into two cases. The simpler case is the one of completely uncorrelated experiments: in this case the pdf used in the construction is given by a multiplication of the pdfs of the single experiments:

$$q(\vec{n}|s) = \prod_{i=1}^{n_{exp}} q(n_i|s, \epsilon_i) \quad (5)$$

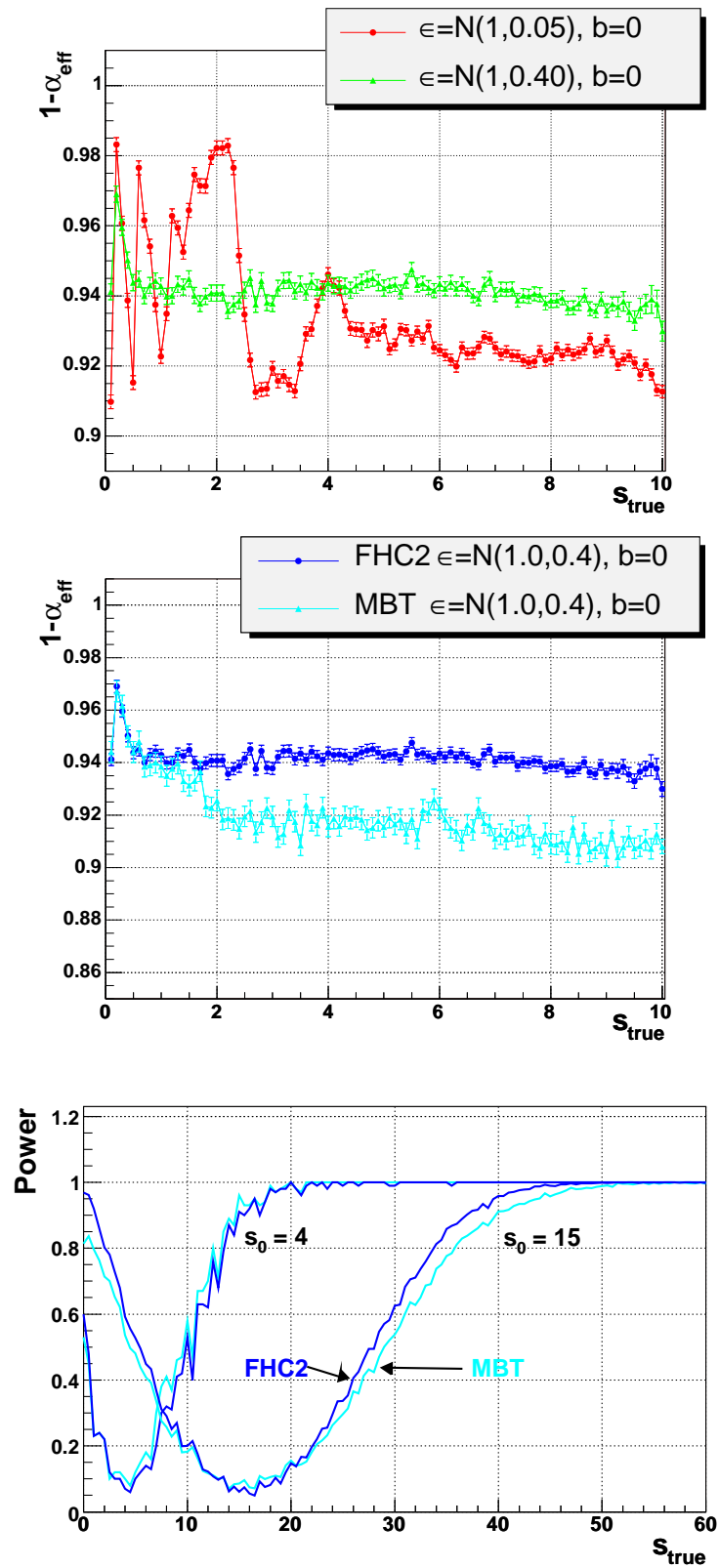


Fig. 1. Examples for the coverage and power of the discussed methods. Uppermost figure: coverage of the FHC² method assuming 5% and 40% Gaussian uncertainties in efficiency. Middle figure: the coverage for the FHC² method compared to the MBT method for 40% Gaussian efficiency uncertainties. Lowest figure: the power of the two methods compared for 40% Gaussian uncertainties in efficiency.

If correlations between uncertainties in nuisance parameters have to be considered, multivariate pdfs have to be employed:

$$q(\vec{n}|s, \vec{\epsilon}) = \int_0^\infty \dots \int_0^\infty \prod_{i=1}^{n_{exp}} P(n|s, \epsilon'_i) P(\vec{\epsilon}'|\vec{\epsilon}) \prod_{i=1}^{n_{exp}} d\epsilon'_i \quad (6)$$

We illustrate the effect of combining different experiments with the example of the CDF limit on the branching ratio for $B_s^0 \rightarrow \mu^+ \mu^-$, see Table 1. In this case, two CDF data sets are combined with an uncorrelated uncertainty in the background expectation and an uncertainty in the efficiency which can be factorized into a correlated and uncorrelated part⁸. Bernhard et al.⁸ presented a fully Bayesian combination, which is included in the table for comparison. The limit obtained using the FHC² method is slightly smaller than the fully Bayesian upper limit.

Table 1. The CDF single and combined limits on $B_s^0 \rightarrow \mu^+ \mu^-$ calculated by FHC². CDF1 and CDF2 denote the two different data sets used for single limits. The quoted uncertainties are for the single experiments, the efficiency uncertainties change to 13.1 and 11.1% for the uncorrelated part if experiments are combined. The number in the parentheses is the result of the purely Bayesian calculation⁷.

| | CDF 1 | CDF 2 |
|----------------------------|-----------|-------|
| background uncertainty [%] | 14.8 | 19.7 |
| eff. uncertainty [%] | 10.4 | 11.3 |
| corr. eff. uncertainty [%] | 15.5 | |
| 95% CL [10^{-7}] | 2.5 | 4.3 |
| 95% combined [10^{-7}] | 1.7 (2.0) | |

5. Discussion & Conclusion

There are two main caveats when interpreting the presented results: first of all, the methods (more or less implicitly) assume a flat prior probability for the true nuisance parameter. Thus, conclusions on the coverage and power are true only for that prior. This assumption seems particularly harmful in the case of combined experiments, a case for which we did not calculate the coverage. Results presented at this conference by Heinrich⁹ indicate that the assumption of a flat prior for nuisance parameters in each

channel leads to significant under-coverage for fully Bayesian confidence intervals. Heinrich also shows that this behavior can be remedied with an appropriate choice of prior (in his particular example: $1/\epsilon$). For the methods presented here this might imply that there is under-coverage in the case of several combined experiments. A second caveat is that we test the coverage only for 90% confidence level. At this conference Cranmer¹⁰ presented results that indicate under-coverage for very high confidence levels ($> 5 \sigma$) if uncertainties in the background are treated in the Bayesian way. Tests of coverage for high confidence levels and combined experiments are currently under way.

With these caveats in mind, we conclude that the Bayesian treatment of nuisance parameters introduces a moderate amount of over-coverage. The MBT method has less over-coverage for the case with large Gaussian uncertainties in the signal efficiencies. We also compared the power of the two suggested methods. For large uncertainties and large true signals, the FHC² method has higher power for hypotheses relatively far away from the null hypothesis.

Acknowledgments

We would like to thank the conference organizers, in particular Louis Lyons for organizing this useful and very enjoyable conference.

References

1. G. J. Feldman and R. D. Cousins, Phys. Rev. **D57**, 3873 (1998).
2. J. Neyman, Phil. Trans. Royal Soc. London **A**, 333 (1937).
3. R. D. Cousins and V. L. Highland, Nucl. Instrum. Meth. A **320**, 331 (1992).
4. J. Conrad, O. Botner, A. Hallgren and C. P. de los Heros, Phys. Rev. **D67**, 012002 (2003).
5. F. Tegenfeldt and J. Conrad, Nucl. Instrum. Meth. A **539**, 407 (2005).
6. G. C. Hill, Phys. Rev. **D67**, 118101 (2003).
7. J. Conrad, Comp. Phys. Comm. **158** 117 (2004).
8. R. Bernhard *et al.* [CDF Collaboration], arXiv:hep-ex/0508058
9. J. Heinrich, these proceedings.
10. K. Cranmer, these proceedings.