Searches and Discoveries using Likelihoods

https://indico.desy.de/indico/event/25594/



Terascale Statistics School DESY / Zoom 10 July 2020



Glen Cowan Physics Department Royal Holloway, University of London g.cowan@rhul.ac.uk www.pp.rhul.ac.uk/~cowan

Outline

0) Brief review of statistical tests and setting limits.

- 1) Tests based on the profile likelihood ratio for discovery and limits.
- 2) Experimental sensitivity

(Frequentist) statistical tests

Consider test of a parameter μ , e.g., proportional to cross section. Result of measurement is a set of numbers *x*.

To define test of μ , specify *critical region* w_{μ} , such that probability to find $x \in w_{\mu}$ is not greater than α (the *size* or *significance level*):

 $P(\mathbf{x} \in w_{\mu}|\mu) \leq \alpha$ data space Ω such that the critical region corresponds to $p_{\mu} < \alpha$. Often use, e.g., $\alpha = 0.05$. If observe $\mathbf{x} \in w_{\mu}$, reject μ . data space Ω **use** data space Ω

Test statistics and *p*-values

Often define the test with a statistic $q_{\mu}(\mathbf{x})$ such that the boundary of the critical region is $q_{\mu}(\mathbf{x}) = c_{\alpha}$ for some constant c_{α} .

For examples of statistics based on the profile likelihood ratio, see, e.g., CCGV, EPJC 71 (2011) 1554; arXiv:1007.1727.

Usually define q_{μ} such that higher values represent increasing incompatibility between the data and the hypothesized μ , so that the *p*-value of μ is

$$p_{\mu} = \int_{q_{\mu,obs}}^{\infty} f(q_{\mu}|\mu) dq_{\mu}$$

observed value of q_{μ} pdf of q_{μ} assuming μ

Equivalent formulation of test: reject μ if $p_{\mu} \leq \alpha$.

Confidence interval from inversion of a test

Carry out a test of size α for all values of μ .

The values that are not rejected $(p_{\mu} > \alpha)$ constitute a *confidence interval* for μ at confidence level CL = $1 - \alpha$.

The confidence interval will by construction contain the true value of μ with probability of at least $1 - \alpha$.

The interval depends on the choice of the critical region of the test.

Put critical region where data are likely to be under assumption of the relevant alternative to the μ that's being tested.

Test $\mu = 0$, alternative is $\mu > 0$: test for discovery.

G. Cowan

Test $\mu = \mu_0$, alternative is $\mu = 0$: testing all μ_0 gives upper limit. To find boundary of confidence region, set $p_{\mu} = \alpha$ and solve for μ .

p-value for discovery

Large q_0 means increasing incompatibility between the data and hypothesis, therefore *p*-value for an observed $q_{0,obs}$ is

$$p_0 = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0) \, dq_0$$

will get formula for this later



From *p*-value get equivalent significance,

$$Z = \Phi^{-1}(1-p)$$

Significance from *p*-value

Often define significance Z as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same p-value.



$$p = \int_{Z}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^{2}/2} \, dx = 1 - \Phi(Z) \qquad \text{1 - TMath::Freq}$$

 $Z = \Phi^{-1}(1-p)$ TMath::NormQuantile

G. Cowan

Prototype search analysis

Search for signal in a region of phase space; result is histogram of some variable *x* giving numbers:

$$\mathbf{n}=(n_1,\ldots,n_N)$$

Assume the n_i are Poisson distributed with expectation values

$$E[n_i] = \mu s_i + b_i$$
strength parameter

where

G. Cowan

Prototype analysis (II)

Often also have a subsidiary measurement that constrains some of the background and/or shape parameters:

$$\mathbf{m} = (m_1, \ldots, m_M)$$

Assume the m_i are Poisson distributed with expectation values

$$E[m_i] = u_i(\boldsymbol{\theta})$$

nuisance parameters ($\boldsymbol{\theta}_{s}, \boldsymbol{\theta}_{b}, b_{tot}$)

Likelihood function is

$$L(\mu, \theta) = \prod_{j=1}^{N} \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)} \quad \prod_{k=1}^{M} \frac{u_k^{m_k}}{m_k!} e^{-u_k}$$

G. Cowan

The profile likelihood ratio

Base significance test on the profile likelihood ratio:



Define critical region of test of μ by the region of data space that gives the lowest values of $\lambda(\mu)$.

Important advantage of profile LR is that its distribution becomes independent of nuisance parameters in "large sample" limit:

Wilks' theorem: $f(-2\ln\lambda(\mu)|\mu)$ tends to chi-square, number of d.o.f. = number of parameters of interest (here 1).

Test statistic for discovery

Try to reject background-only ($\mu = 0$) hypothesis using

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \ge 0\\ 0 & \hat{\mu} < 0 \end{cases}$$

i.e. here only regard upward fluctuation of data as evidence against the background-only hypothesis.

Note that even though here physically $\mu \ge 0$, we allow $\hat{\mu}$ to be negative. In large sample limit its distribution becomes Gaussian, and this will allow us to write down simple expressions for distributions of our test statistics.

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Distribution of q_0 in large-sample limit

Assuming approximations valid in the large sample (asymptotic) limit, we can write down the full distribution of q_0 as

$$f(q_0|\mu') = \left(1 - \Phi\left(\frac{\mu'}{\sigma}\right)\right)\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}\exp\left[-\frac{1}{2}\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)^2\right]$$

The special case $\mu' = 0$ is a "half chi-square" distribution:

$$f(q_0|0) = \frac{1}{2}\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}e^{-q_0/2}$$

In large sample limit, $f(q_0|0)$ independent of nuisance parameters; $f(q_0|\mu')$ depends on nuisance parameters through σ .

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Cumulative distribution of q_0 , significance

From the pdf, the cumulative distribution of q_0 is found to be

$$F(q_0|\mu') = \Phi\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)$$

The special case $\mu' = 0$ is

$$F(q_0|0) = \Phi\left(\sqrt{q_0}\right)$$

The *p*-value of the $\mu = 0$ hypothesis is

$$p_0 = 1 - F(q_0|0)$$

Therefore the discovery significance Z is simply

$$Z = \Phi^{-1}(1 - p_0) = \sqrt{q_0}$$

G. Cowan

Example of a *p*-value ATLAS, Phys. Lett. B 716 (2012) 1-29



Test statistic for upper limits

cf. Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554. For purposes of setting an upper limit on μ use

$$q_{\mu} = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \end{cases} \quad \text{where} \quad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})}$$

I.e. when setting an upper limit, an upwards fluctuation of the data is not taken to mean incompatibility with the hypothesized μ :

From observed
$$q_{\mu}$$
 find *p*-value: $p_{\mu} = \int_{q_{\mu,\text{obs}}}^{\infty} f(q_{\mu}|\mu) dq_{\mu}$

Large sample approximation: $p_{\mu} = 1 - \Phi(\sqrt{q_{\mu}})$ 95% CL upper limit on μ is highest value for which *p*-val

95% CL upper limit on μ is highest value for which *p*-value is not less than 0.05.

G. Cowan 2020 DESY Terascale Statistics School / Searches and Discoveries using Likelihoods

 \sim

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Monte Carlo test of asymptotic formulae

Consider again $n \sim \text{Poisson}(\mu s + b), m \sim \text{Poisson}(\tau b)$ Use q_{μ} to find *p*-value of hypothesized μ values.

E.g. $f(q_1|1)$ for *p*-value of $\mu = 1$. Typically interested in 95% CL, i.e., *p*-value threshold = 0.05, i.e., $q_1 = 2.69$ or $Z_1 = \sqrt{q_1} = 1.64$. Median[$q_1 | 0$] gives "exclusion sensitivity".

for s = 6, b = 9.

G. Cowan



How to read the green and yellow limit plots For every value of $m_{\rm H}$, find the upper limit on μ .

Also for each $m_{\rm H}$, determine the distribution of upper limits $\mu_{\rm up}$ one would obtain under the hypothesis of $\mu = 0$.

The dashed curve is the median μ_{up} , and the green (yellow) bands give the $\pm 1\sigma(2\sigma)$ regions of this distribution.



Expected (or median) significance / sensitivity

When planning the experiment, we want to quantify how sensitive we are to a potential discovery, e.g., by given median significance assuming some nonzero strength parameter μ' .



So for *p*-value, need $f(q_0|0)$, for sensitivity, will need $f(q_0|\mu')$,

G. Cowan

Expected discovery significance for counting experiment with background uncertainty Discovery sensitivity for counting experiment with *b* known:

I. Discovery sensitivity for counting experiment with *b* known:

(a)
$$\frac{s}{\sqrt{b}}$$

(b) Profile likelihood ratio test & Asimov:

$$\sqrt{2\left((s+b)\ln\left(1+\frac{s}{b}\right)-s\right)}$$

II. Discovery sensitivity with uncertainty in b, σ_b :

(a)
$$\frac{s}{\sqrt{b+\sigma_b^2}}$$

(b) Profile likelihood ratio test & Asimov:

$$\left[2\left((s+b)\ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2}\ln\left[1 + \frac{\sigma_b^2s}{b(b+\sigma_b^2)}\right]\right)\right]^{1/2}$$

G. Cowan

Counting experiment with known background Count a number of events $n \sim Poisson(s+b)$, where s = expected number of events from signal,

b = expected number of background events.

To test for discovery of signal compute *p*-value of s = 0 hypothesis,

$$p = P(n \ge n_{\text{obs}}|b) = \sum_{n=n_{\text{obs}}}^{\infty} \frac{b^n}{n!} e^{-b} = 1 - F_{\chi^2}(2b; 2n_{\text{obs}})$$

Usually convert to equivalent significance: $Z = \Phi^{-1}(1-p)$ where Φ is the standard Gaussian cumulative distribution, e.g., Z > 5 (a 5 sigma effect) means $p < 2.9 \times 10^{-7}$.

To characterize sensitivity to discovery, give expected (mean or median) Z under assumption of a given s.

 s/\sqrt{b} for expected discovery significance For large s + b, $n \to x \sim \text{Gaussian}(\mu, \sigma)$, $\mu = s + b$, $\sigma = \sqrt{(s + b)}$. For observed value x_{obs} , *p*-value of s = 0 is $\text{Prob}(x > x_{\text{obs}} | s = 0)$,:

$$p_0 = 1 - \Phi\left(\frac{x_{\rm obs} - b}{\sqrt{b}}\right)$$

Significance for rejecting s = 0 is therefore

$$Z_0 = \Phi^{-1}(1 - p_0) = \frac{x_{\text{obs}} - b}{\sqrt{b}}$$

Expected (median) significance assuming signal rate *s* is

$$\mathrm{median}[Z_0|s+b] = \frac{s}{\sqrt{b}}$$

G. Cowan

Better approximation for significance Poisson likelihood for parameter *s* is

> $L(s) = \frac{(s+b)^n}{n!} e^{-(s+b)}$ For now no nuisance

To test for discovery use profile likelihood ratio:

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{s} \ge 0 \ , \\ 0 & \hat{s} < 0 \ . \end{cases} \qquad \lambda(s) = \frac{L(s, \hat{\hat{\theta}}(s))}{L(\hat{s}, \hat{\theta})}$$

So the likelihood ratio statistic for testing s = 0 is

$$q_0 = -2\ln\frac{L(0)}{L(\hat{s})} = 2\left(n\ln\frac{n}{b} + b - n\right) \quad \text{for } n > b, \ 0 \text{ otherwise}$$

G. Cowan

2020 DESY Terascale Statistics School / Searches and Discoveries using Likelihoods

params.

Approximate Poisson significance (continued)

For sufficiently large s + b, (use Wilks' theorem),

$$Z = \sqrt{2\left(n\ln\frac{n}{b} + b - n\right)} \quad \text{for } n > b \text{ and } Z = 0 \text{ otherwise.}$$

To find median[*Z*|*s*], let $n \rightarrow s + b$ (i.e., the Asimov data set):

$$Z_{\rm A} = \sqrt{2\left((s+b)\ln\left(1+\frac{s}{b}\right) - s\right)}$$

This reduces to s/\sqrt{b} for s << b.

G. Cowan

 $n \sim \text{Poisson}(s+b)$, median significance, assuming *s*, of the hypothesis s = 0

CCGV, EPJC 71 (2011) 1554, arXiv:1007.1727



"Exact" values from MC, jumps due to discrete data.

Asimov $\sqrt{q_{0,A}}$ good approx. for broad range of *s*, *b*.

 s/\sqrt{b} only good for $s \ll b$.

Extending s/\sqrt{b} to case where b uncertain

The intuitive explanation of s/\sqrt{b} is that it compares the signal, *s*, to the standard deviation of *n* assuming no signal, \sqrt{b} .

Now suppose the value of *b* is uncertain, characterized by a standard deviation σ_b .

A reasonable guess is to replace \sqrt{b} by the quadratic sum of \sqrt{b} and σ_b , i.e.,

$$\operatorname{med}[Z|s] = \frac{s}{\sqrt{b + \sigma_b^2}}$$

This has been used to optimize some analyses e.g. where σ_b cannot be neglected.

Profile likelihood with b uncertain

This is the well studied "on/off" problem: Cranmer 2005; Cousins, Linnemann, and Tucker 2008; Li and Ma 1983,...

Measure two Poisson distributed values:

 $n \sim \text{Poisson}(s+b)$ (primary or "search" measurement) $m \sim \text{Poisson}(\tau b)$ (control measurement, τ known)

The likelihood function is

$$L(s,b) = \frac{(s+b)^n}{n!} e^{-(s+b)} \frac{(\tau b)^m}{m!} e^{-\tau b}$$

Use this to construct profile likelihood ratio (*b* is nuisance parmeter): $\lambda(0) = \frac{L(0, \hat{b}(0))}{L(0, \hat{b}(0))}$

$$\lambda(0) = \frac{L(0, b(0))}{L(\hat{s}, \hat{b})}$$

G. Cowan

Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

$$\begin{aligned} \hat{s} &= n - m/\tau ,\\ \hat{b} &= m/\tau ,\\ \hat{b}(s) &= \frac{n + m - (1 + \tau)s + \sqrt{(n + m - (1 + \tau)s)^2 + 4(1 + \tau)sm}}{2(1 + \tau)} \end{aligned}$$

and in particular to test for discovery (s = 0),

$$\hat{\hat{b}}(0) = \frac{n+m}{1+\tau}$$

G. Cowan

Asymptotic significance

Use profile likelihood ratio for q_0 , and then from this get discovery significance using asymptotic approximation (Wilks' theorem):

$$Z = \sqrt{q_0}$$
$$= \left[-2\left(n\ln\left[\frac{n+m}{(1+\tau)n}\right] + m\ln\left[\frac{\tau(n+m)}{(1+\tau)m}\right]\right) \right]^{1/2}$$

for $n > \hat{b}$ and Z = 0 otherwise.

Essentially same as in:

Robert D. Cousins, James T. Linnemann and Jordan Tucker, NIM A 595 (2008) 480– 501; arXiv:physics/0702156.

Tipei Li and Yuqian Ma, Astrophysical Journal 272 (1983) 317–324.

Asimov approximation for median significance

To get median discovery significance, replace *n*, *m* by their expectation values assuming background-plus-signal model:

$$n \to s + b$$

$$m \to \tau b$$

$$Z_{\rm A} = \left[-2\left((s+b) \ln\left[\frac{s+(1+\tau)b}{(1+\tau)(s+b)}\right] + \tau b \ln\left[1+\frac{s}{(1+\tau)b}\right] \right) \right]^{1/2}$$
Or use the variance of $\hat{b} = m/\tau$, $V[\hat{b}] \equiv \sigma_b^2 = \frac{b}{\tau}$, to eliminate τ :

$$Z_{\rm A} = \left[2\left((s+b) \ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2} \ln\left[1+\frac{\sigma_b^2 s}{b(b+\sigma_b^2)}\right] \right) \right]^{1/2}$$

G. Cowan

Limiting cases

Expanding the Asimov formula in powers of *s/b* and σ_b^2/b (= 1/ τ) gives

$$Z_{\rm A} = \frac{s}{\sqrt{b + \sigma_b^2}} \left(1 + \mathcal{O}(s/b) + \mathcal{O}(\sigma_b^2/b) \right)$$

So the "intuitive" formula can be justified as a limiting case of the significance from the profile likelihood ratio test evaluated with the Asimov data set. Testing the formulae: s = 5



G. Cowan

Using sensitivity to optimize a cut



Figure 1: (a) The expected significance as a function of the cut value x_{cut} ; (b) the distributions of signal and background with the optimal cut value indicated.

Summary

Here only saw frequentist approach:

Test hypothesis *H* by computing *p*-value, reject if $p_H < \alpha$.

Reject *H* only if $p_H \leq \alpha$ for all values of nuisance parameters.

To define test (or *p*-value), need to ask what are the relevant alternatives to *H*.

Discovery means *p*-value of no-signal $H < \alpha \sim 3 \times 10^{-7}$.

Confidence interval = set of hypothesized parameter values μ for which $p_{\mu} > \alpha \sim 0.05$ (edge of conf. int. = "limit").

Characterize sensitivity with mean expected significance from test under assumption of relevant alternative.

Bayesian equivalents (beyond scope of this talk): limits based on posterior pdf; discovery using Bayes factor.

Extra slides

Why 5 sigma?

Common practice in HEP has been to claim a discovery if the *p*-value of the no-signal hypothesis is below 2.9 × 10⁻⁷, corresponding to a significance $Z = \Phi^{-1} (1 - p) = 5$ (a 5 σ effect).

There a number of reasons why one may want to require such a high threshold for discovery:

The "cost" of announcing a false discovery is high.

Unsure about systematics.

Unsure about look-elsewhere effect.

The implied signal may be a priori highly improbable (e.g., violation of Lorentz invariance).

Why 5 sigma (cont.)?

But the primary role of the *p*-value is to quantify the probability that the background-only model gives a statistical fluctuation as big as the one seen or bigger.

It is not intended as a means to protect against hidden systematics or the high standard required for a claim of an important discovery.

In the processes of establishing a discovery there comes a point where it is clear that the observation is not simply a fluctuation, but an "effect", and the focus shifts to whether this is new physics or a systematic.

Providing LEE is dealt with, that threshold is probably closer to 3σ than 5σ .

Choice of test for limits

In some cases $\mu = 0$ is no longer a relevant alternative and we want to try to exclude μ on the grounds that some other measure of incompatibility between it and the data exceeds some threshold.

If the measure of incompatibility is taken to be the likelihood ratio with respect to a two-sided alternative, then the critical region can contain both high and low data values.

→ unified intervals, G. Feldman, R. Cousins, Phys. Rev. D 57, 3873–3889 (1998)

The Big Debate is whether to use one-sided or unified intervals in cases where small (or zero) values of the parameter are relevant alternatives. Professional statisticians have voiced support on both sides of the debate.

Unified (Feldman-Cousins) intervals

We can use directly

$$t_{\mu} = -2\ln\lambda(\mu)$$
 where



38

as a test statistic for a hypothesized μ .

Large discrepancy between data and hypothesis can correspond either to the estimate for μ being observed high or low relative to μ .

This is essentially the statistic used for Feldman-Cousins intervals (here also treats nuisance parameters).

G. Feldman and R.D. Cousins, Phys. Rev. D 57 (1998) 3873.

Lower edge of interval can be at $\mu = 0$, depending on data.

G. Cowan

Distribution of t_{μ}

Using Wald approximation, $f(t_{\mu}|\mu')$ is noncentral chi-square for one degree of freedom:

$$f(t_{\mu}|\mu') = \frac{1}{2\sqrt{t_{\mu}}} \frac{1}{\sqrt{2\pi}} \left[\exp\left(-\frac{1}{2}\left(\sqrt{t_{\mu}} + \frac{\mu - \mu'}{\sigma}\right)^2\right) + \exp\left(-\frac{1}{2}\left(\sqrt{t_{\mu}} - \frac{\mu - \mu'}{\sigma}\right)^2\right) \right]$$

Special case of $\mu = \mu'$ is chi-square for one d.o.f. (Wilks).

The *p*-value for an observed value of t_{μ} is

$$p_{\mu} = 1 - F(t_{\mu}|\mu) = 2\left(1 - \Phi\left(\sqrt{t_{\mu}}\right)\right)$$

and the corresponding significance is

$$Z_{\mu} = \Phi^{-1}(1 - p_{\mu}) = \Phi^{-1} \left(2\Phi \left(\sqrt{t_{\mu}} \right) - 1 \right)$$

G. Cowan

Upper/lower edges of F-C interval for μ versus *b* for $n \sim \text{Poisson}(\mu+b)$



Lower edge may be at zero, depending on data. For n = 0, upper edge has (weak) dependence on *b*.

G. Cowan

40

Feldman-Cousins discussion

The initial motivation for Feldman-Cousins (unified) confidence intervals was to eliminate null intervals.

The F-C limits are based on a likelihood ratio for a test of μ with respect to the alternative consisting of all other allowed values of μ (not just, say, lower values).

The interval's upper edge is higher than the limit from the onesided test, and lower values of μ may be excluded as well. A substantial downward fluctuation in the data gives a low (but nonzero) limit.

This means that when a value of μ is excluded, it is because there is a probability α for the data to fluctuate either high or low in a manner corresponding to less compatibility as measured by the likelihood ratio.