Classical interval estimation, limits systematics and beyond - Part 2



IN2P3 School of Statistics Zoom / 19 January 2021

https://indico.in2p3.fr/event/20220/



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

Outline

- Nuisance parameters, systematic uncertainties
- Prototype analysis with profile likelihood ratio
- Expected discovery significance (with systematics)

Systematic uncertainties and nuisance parameters In general, our model of the data is not perfect:



Can improve model by including additional adjustable parameters.

$$P(x|\mu) \to P(x|\mu, \theta)$$

Nuisance parameter ↔ systematic uncertainty. Some point in the parameter space of the enlarged model should be "true".

Presence of nuisance parameter decreases sensitivity of analysis to the parameter of interest (e.g., increases variance of estimate).

Profile Likelihood

Suppose we have a likelihood $L(\mu, \theta) = P(x|\mu, \theta)$ with Nparameters of interest $\mu = (\mu_1, ..., \mu_N)$ and M nuisance parameters $\theta = (\theta_1, ..., \theta_M)$. The "profiled" (or "constrained") values of θ are:

$$\hat{\hat{\boldsymbol{\theta}}}(\boldsymbol{\mu}) = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\mu}, \boldsymbol{\theta})$$

and the profile likelihood is: $L_{\rm p}(\boldsymbol{\mu}) = L(\boldsymbol{\mu}, \hat{\boldsymbol{\theta}})$

The profile likelihood depends only on the parameters of interest; the nuisance parameters are replaced by their profiled values.

The profile likelihood can be used to obtain confidence intervals/regions for the parameters of interest in the same way as one would for all of the parameters from the full likelihood.

Profile Likelihood Ratio – Wilks theorem

Goal is to test/reject regions of μ space (param. of interest).

Rejecting a point μ should mean $p_{\mu} \leq \alpha$ for all possible values of the nuisance parameters θ .

Test $\boldsymbol{\mu}$ using the "profile likelihood ratio": $\lambda(\boldsymbol{\mu}) = \frac{L(\boldsymbol{\mu}, \hat{\boldsymbol{\theta}})}{L(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\theta}})}$

Let $t_{\mu} = -2 \ln \lambda(\mu)$. Wilks' theorem says in large-sample limit: $t_{\mu} \sim \text{chi-square}(N)$

where the number of degrees of freedom is the number of parameters of interest (components of μ). So *p*-value for μ is

$$p_{\boldsymbol{\mu}} = \int_{t_{\boldsymbol{\mu},\text{obs}}}^{\infty} f(t_{\boldsymbol{\mu}} | \boldsymbol{\mu}, \boldsymbol{\theta}) \, dt_{\boldsymbol{\mu}} = 1 - F_{\chi_N^2}(t_{\boldsymbol{\mu},\text{obs}})$$

G. Cowan / RHUL Physics

Profile Likelihood Ratio – Wilks theorem (2)

If we have a large enough data sample to justify use of the asymptotic chi-square pdf, then if μ is rejected, it is rejected for any values of the nuisance parameters.

The recipe to get confidence regions/intervals for the parameters of interest at $CL = 1 - \alpha$ is thus the same as before, simply use the profile likelihood:

$$\ln L_{\rm p}(\boldsymbol{\mu}) = \ln L_{\rm max} - \frac{1}{2} F_{\chi_N^2}^{-1} (1 - \alpha)$$

where the number of degrees of freedom N for the chi-square quantile is equal to the number of parameters of interest.

If the large-sample limit is not justified, then use e.g. Monte Carlo to get distribution of t_{μ} .

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 / RHUL Physics

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 / RHUL Physics

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.

G. Cowan / RHUL Physics

Test statistic for discovery

Suppose relevant alternative to background-only ($\mu = 0$) is $\mu \ge 0$. So take critical region for test of $\mu = 0$ corresponding to high q_0

and $\hat{\mu} > 0$ (data characteristic for $\mu \ge 0$).

That is, to test background-only hypothesis define statistic

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

i.e. here only large (positive) observed signal strength is 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.

G. Cowan / RHUL Physics

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 .

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

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

use e.g. asymptotic formula



From *p*-value get equivalent significance,

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

G. Cowan / RHUL Physics

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 / RHUL Physics

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

Monte Carlo test of asymptotic formula

- $n \sim \text{Poisson}(\mu s + b)$
- $m \sim \text{Poisson}(\tau b)$
- $\mu =$ param. of interest
- *b* = nuisance parameter

Here take *s* known, $\tau = 1$.

Asymptotic formula is good approximation to 5σ level ($q_0 = 25$) already for $b \sim 20$.



How to read the p_0 plot

The "local" p_0 means the *p*-value of the background-only hypothesis obtained from the test of $\mu = 0$ at each individual $m_{\rm H}$, without any correct for the Look-Elsewhere Effect.

The "Expected" (dashed) curve gives the median p_0 under assumption of the SM Higgs (μ = 1) at each $m_{\rm H}$.



The blue band gives the width of the distribution $(\pm 1\sigma)$ of significances under assumption of the SM Higgs.

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

Test statistic for upper limits

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

Large sample approximation:

$$p_{\mu} = 1 - \Phi\left(\sqrt{q_{\mu}}\right)$$

To find upper limit at CL = $1-\alpha$, set $p_{\mu} = \alpha$ and solve for μ .

G. Cowan / RHUL Physics

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".

Here asymptotic formulae good for s = 6, b = 9.



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 σ) regions of this distribution.



ATLAS, Phys. Lett. B 716 (2012) 1-29

Expected discovery significance for counting experiment with background uncertainty

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 / RHUL Physics

Counting experiment with known background

Count a number of events $n \sim \text{Poisson}(s+b)$, 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.

G. Cowan / RHUL Physics

 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 / RHUL Physics

Better approximation for significance

Poisson likelihood for parameter s is

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

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

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 / RHUL Physics

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):

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

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

 $n \sim \text{Poisson}(s+b)$, 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 $s \ll b$.

G. Cowan / RHUL Physics

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 parameter): $L(0, \hat{\hat{b}}(0))$

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

G. Cowan / RHUL Physics

Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

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

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

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

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.

G. Cowan / RHUL Physics

Asimov approximation for median significance

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

$$\begin{split} n &\to s + b \\ m &\to \tau b \end{split}$$

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

Limiting cases

Expanding the Asimov formula in powers of s/b and σ_b^2/b (= $1/\tau$) 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 / RHUL Physics

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.

G. Cowan / RHUL Physics

Summary on discovery sensitivity

Simple formula for expected discovery significance based on profile likelihood ratio test and Asimov approximation:

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

For large *b*, all formulae OK.

For small b, s/\sqrt{b} and $s/\sqrt{(b+\sigma_b^2)}$ overestimate the significance.

Could be important in optimization of searches with low background.

Formula maybe also OK if model is not simple on/off experiment, e.g., several background control measurements (check this).

Finally

Two lectures only enough for a brief introduction to:

Limits (confidence intervals/regions)

Systematics (nuisance parameters)

A bit beyond... (sensitivity)

Final thought: once the basic formalism is fixed, most of the work focuses on writing down the likelihood, e.g., $P(x|\theta)$, and including in it enough parameters to adequately describe the data (true for both Bayesian and frequentist approaches).



p-values in cases with nuisance parameters

Suppose we have a statistic q_{θ} that we use to test a hypothesized value of a parameter θ , such that the *p*-value of θ is

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

But what values of v to use for $f(q_{\theta} | \theta, v)$?

Fundamentally we want to reject θ only if $p_{\theta} < \alpha$ for all v.

$$\rightarrow$$
 "exact" confidence interval

But in general for finite data samples this is not true; one may be unable to reject some θ values if all values of v must be considered (resulting interval for θ "overcovers").

Profile construction ("hybrid resampling")

K. Cranmer, PHYSTAT-LHC Workshop on Statistical Issues for LHC Physics, 2008. oai:cds.cern.ch:1021125, cdsweb.cern.ch/record/1099969.

Approximate procedure is to reject θ if $p_{\theta} \leq \alpha$ where the *p*-value is computed assuming the value of the nuisance parameter that best fits the data for the specified θ :

$\hat{\hat{ u}}(heta)$	"double hat" notation means profiled
	value, i.e., parameter that maximizes
	likelihood for the given θ .

The resulting confidence interval will have the correct coverage for the points $(\theta, \hat{\hat{v}}(\theta))$.

Elsewhere it may under- or overcover, but this is usually as good as we can do (check with MC if crucial or small sample problem).

G. Cowan / RHUL Physics

Low sensitivity to μ

It can be that the effect of a given hypothesized μ is very small relative to the background-only (μ = 0) prediction.

This means that the distributions $f(q_{\mu}|\mu)$ and $f(q_{\mu}|0)$ will be almost the same:



G. Cowan / RHUL Physics

Having sufficient sensitivity

In contrast, having sensitivity to μ means that the distributions $f(q_{\mu}|\mu)$ and $f(q_{\mu}|0)$ are more separated:



That is, the power (probability to reject μ if $\mu = 0$) is substantially higher than α . Use this power as a measure of the sensitivity.

G. Cowan / RHUL Physics

Spurious exclusion

Consider again the case of low sensitivity. By construction the probability to reject μ if μ is true is α (e.g., 5%).

And the probability to reject μ if $\mu = 0$ (the power) is only slightly greater than α .



This means that with probability of around $\alpha = 5\%$ (slightly higher), one excludes hypotheses to which one has essentially no sensitivity (e.g., $m_{\rm H} = 1000$ TeV).

"Spurious exclusion"

Ways of addressing spurious exclusion

The problem of excluding parameter values to which one has no sensitivity known for a long time; see e.g.,

Virgil L. Highland, *Estimation of Upper Limits from Experimental Data*, July 1986, Revised February 1987, Temple University Report C00-3539-38.

In the 1990s this was re-examined for the LEP Higgs search by Alex Read and others

T. Junk, Nucl. Instrum. Methods Phys. Res., Sec. A 434, 435 (1999); A.L. Read, J. Phys. G 28, 2693 (2002).

and led to the "CL_s" procedure for upper limits.

Unified intervals also effectively reduce spurious exclusion by the particular choice of critical region.

The CL_s procedure

In the usual formulation of CL_s , one tests both the $\mu = 0$ (*b*) and $\mu > 0$ ($\mu s+b$) hypotheses with the same statistic $Q = -2\ln L_{s+b}/L_b$:



G. Cowan / RHUL Physics

The CL_s procedure (2)

As before, "low sensitivity" means the distributions of Q under b and s+b are very close:



G. Cowan / RHUL Physics

The CL_s procedure (3)

The CL_s solution (A. Read et al.) is to base the test not on the usual *p*-value (CL_{s+b}), but rather to divide this by CL_b (~ one minus the *p*-value of the *b*-only hypothesis), i.e.,



G. Cowan / RHUL Physics

Choice of test for limits (2)

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 λ

$$(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\hat{\theta}})}$$

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 μ = 0, depending on data.

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



Lower edge may be at zero, depending on data.

For n = 0, upper edge has (weak) dependence on b.

Example: fitting a straight line

Data:
$$(x_i, y_i, \sigma_i), i = 1, ..., n$$
.

Model: y_i independent and all follow $y_i \sim \text{Gauss}(\mu(x_i), \sigma_i)$

$$\mu(x;\theta_0,\theta_1)=\theta_0+\theta_1x,$$

assume x_i and σ_i known.

Goal: estimate θ_0

Here suppose we don't care about θ_1 (example of a "nuisance parameter")



Maximum likelihood fit with Gaussian data

In this example, the y_i are assumed independent, so the likelihood function is a product of Gaussians:

$$L(\theta_0, \theta_1) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] ,$$

Maximizing the likelihood is here equivalent to minimizing

$$\chi^{2}(\theta_{0},\theta_{1}) = -2 \ln L(\theta_{0},\theta_{1}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i};\theta_{0},\theta_{1}))^{2}}{\sigma_{i}^{2}}.$$

i.e., for Gaussian data, ML same as Method of Least Squares (LS)

θ_1 known a priori

$$L(\theta_0) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] .$$

 $\chi^{2}(\theta_{0}) = -2 \ln L(\theta_{0}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i}; \theta_{0}, \theta_{1}))^{2}}{\sigma_{i}^{2}}.$

For Gaussian y_i , ML same as LS

 $\begin{array}{l} \text{Minimize } \chi^2 \to \text{estimator } \widehat{\theta}_0 \ . \\ \text{Come up one unit from } \chi^2_{\min} \\ \text{to find } \sigma_{\widehat{\theta}_0} \ . \end{array}$



ML (or LS) fit of θ_0 and θ_1

$$\chi^{2}(\theta_{0},\theta_{1}) = -2 \ln L(\theta_{0},\theta_{1}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i};\theta_{0},\theta_{1}))^{2}}{\sigma_{i}^{2}}.$$

Standard deviations from tangent lines to contour

 $\chi^2 = \chi^2_{\rm min} + 1 \; .$

Correlation between $\hat{\theta}_0, \ \hat{\theta}_1$ causes errors to increase.



If we have a measurement $t_1 \sim \text{Gauss}(\theta_1, \sigma_{t_1})$

$$L(\theta_0, \theta_1) = \frac{1}{\sqrt{2\pi\sigma_t}} e^{-(t_1 - \theta_1)^2 / 2\sigma_{t_1}^2} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right]$$

$$\chi^2(\theta_0, \theta_1) = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} + \frac{(t_1 - \theta_1)^2}{\sigma_{t_1}^2}$$

The information on θ_1 improves accuracy of $\hat{\theta}_0$.



Profiling

The $\ln L = \ln L_{max} - \frac{1}{2}$ contour in the (θ_0 , θ_1) plane is a confidence region at CL = 39.3%.

Furthermore if one wants to know only about, say, θ_0 , then the interval in θ_0 corresponding to $\ln L = \ln L_{\max} - \frac{1}{2}$ is a confidence interval at CL = 68.3% (i.e., ±1 std. dev.).

I.e., form the interval for θ_0 using

$$\ln L(\theta_0, \hat{\hat{\theta}}_1(\theta_0)) = \ln L_{\max} - 1/2$$

where θ_1 is replaced by its "profiled" value

$$\hat{\hat{\theta}}_1(\theta_0) = \operatorname*{argmax}_{\theta_1} L(\theta_0, \theta_1)$$



G. Cowan / RHUL Physics

Reminder of Bayesian approach

In Bayesian statistics we can associate a probability with a hypothesis, e.g., a parameter value θ .

Interpret probability of θ as 'degree of belief' (subjective). Need to start with 'prior pdf' $\pi(\theta)$, this reflects degree of belief about θ before doing the experiment.

Our experiment has data x, \rightarrow likelihood $L(x|\theta)$.

Bayes' theorem tells how our beliefs should be updated in light of the data *x*:

$$p(\theta|x) = \frac{L(x|\theta)\pi(\theta)}{\int L(x|\theta')\pi(\theta') d\theta'} \propto L(x|\theta)\pi(\theta)$$

Posterior pdf $p(\theta|x)$ contains all our knowledge about θ .

Bayesian approach: $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$ We need to associate prior probabilities with θ_0 and θ_1 , e.g.,

 $\pi(\theta_0, \theta_1) = \pi_0(\theta_0)\pi_1(\theta_1) \quad \leftarrow \text{suppose knowledge of } \theta_0 \text{ has}$ no influence on knowledge of θ_1

$$\pi_0(\theta_0) = \text{const.} \qquad \leftarrow \text{`non-informative', in any} \\ \text{case much broader than } L(\theta_0)$$

$$\pi_{1}(\theta_{1}) = p(\theta_{1}|t_{1}) \propto p(t_{1}|\theta_{1})\pi_{\mathrm{Ur}}(\theta_{1}) = \frac{1}{\sqrt{2\pi}\sigma_{t}}e^{-(t_{1}-\theta_{1})^{2}/2\sigma_{t}^{2}} \times \mathrm{const.}$$
prior after t_{1} , Ur = "primordial" Likelihood for control before y prior measurement t_{1}

Bayesian example: $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$

Putting the ingredients into Bayes' theorem gives:



Note here the likelihood only reflects the measurements *y*.

The information from the control measurement t_1 has been put into the prior for θ_1 .

We would get the same result using the likelihood $P(y,t|\theta_0,\theta_1)$ and the constant "Ur-prior" for θ_1 .

Marginalizing the posterior pdf

We then integrate (marginalize) $p(\theta_0, \theta_1 | \mathbf{y})$ to find $p(\theta_0 | \mathbf{y})$:

$$p(\theta_0|\mathbf{y}) = \int p(\theta_0, \theta_1|\mathbf{y}) \, d\theta_1$$

In this example we can do the integral (rare). We find

$$p(\theta_0|\mathbf{y}) = \frac{1}{\sqrt{2\pi\sigma_{\theta_0}}} e^{-(\theta_0 - \hat{\theta}_0)^2/2\sigma_{\theta_0^2}}$$

 $\hat{\theta}_0 = \text{same as MLE}$

 $\sigma_{\theta_0} = \sigma_{\hat{\theta}_0}$ (same as for MLE)

For this example, numbers come out same as in frequentist approach, but interpretation different.

G. Cowan / RHUL Physics

Marginalization with MCMC

Bayesian computations involve integrals like

$$p(\theta_0|x) = \int p(\theta_0, \theta_1|x) d\theta_1$$
.

often high dimensionality and impossible in closed form, also impossible with 'normal' acceptance-rejection Monte Carlo.

Markov Chain Monte Carlo (MCMC) has revolutionized Bayesian computation.

MCMC (e.g., Metropolis-Hastings algorithm) generates correlated sequence of random numbers:

cannot use for many applications, e.g., detector MC; effective stat. error greater than if all values independent .

Basic idea: sample multidimensional θ but look only at distribution of parameters of interest.

MCMC basics: Metropolis-Hastings algorithm Goal: given an *n*-dimensional pdf $p(\theta)$, generate a sequence of points $\theta_1, \theta_2, \theta_3, \dots$

- 1) Start at some point $\vec{\theta}_0$
- 2) Generate $\vec{\theta} \sim q(\vec{\theta}; \vec{\theta}_0)$

Proposal density $q(\theta; \theta_0)$ e.g. Gaussian centred about θ_0

3) Form Hastings test ratio $\alpha = \min \left| 1, \frac{\pi}{n} \right|$

$$1, \frac{p(\vec{\theta})q(\vec{\theta}_{0};\vec{\theta})}{p(\vec{\theta}_{0})q(\vec{\theta};\vec{\theta}_{0})} \bigg]$$

- 4) Generate $u \sim \text{Uniform}[0, 1]$
- 5) If $u \le \alpha$, $\vec{\theta_1} = \vec{\theta}$, \leftarrow move to proposed point else $\vec{\theta_1} = \vec{\theta_0} \leftarrow$ old point repeated 6) Iterate

Metropolis-Hastings (continued)

This rule produces a *correlated* sequence of points (note how each new point depends on the previous one).

- Still works if $p(\theta)$ is known only as a proportionality, which is usually what we have from Bayes' theorem: $p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)\pi(\theta)$.
- The proposal density can be (almost) anything, but choose so as to minimize autocorrelation. Often take proposal density symmetric: $q(\theta; \theta_0) = q(\theta_0; \theta)$

Test ratio is (*Metropolis*-Hastings): $\alpha = \min \left[1, \frac{p(\vec{\theta})}{p(\vec{\theta}_0)}\right]$

I.e. if the proposed step is to a point of higher $p(\theta)$, take it; if not, only take the step with probability $p(\theta)/p(\theta_0)$. If proposed step rejected, repeat the current point.

Example: posterior pdf from MCMC

Sample the posterior pdf from previous example with MCMC:



Bayesian method with alternative priors

Suppose we don't have a previous measurement of θ_1 but rather, an "expert" says it should be positive and not too much greater than 0.1 or so, i.e., something like

$$\pi_1(\theta_1) = \frac{1}{\tau} e^{-\theta_1/\tau} , \quad \theta_1 \ge 0 , \quad \tau = 0.1 .$$

From this we obtain (numerically) the posterior pdf for θ_0 :

