

Probability & Statistics

Lecture 3



Data Science in Fundamental Physics
Santiago de Compostela
3-5 June 2024

<https://igfae.usc.es/datascience2024/school/>



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

Outline

Lecture 1: Probability, Bayes vs. Frequentist
Frequentist parameter estimation
Hypothesis tests

Lecture 2: p -values
Confidence limits
Systematic uncertainties
Bayesian parameter estimation

→ Lecture 3: Significance, sensitivity
Bayes factors
Models for anomalies

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:

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

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,

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

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 $s + b$, $n \rightarrow x \sim \text{Gaussian}(\mu, \sigma)$, $\mu = s + b$, $\sigma = \sqrt{s + b}$.

For observed value x_{obs} , p -value of $s = 0$ is $\text{Prob}(x > x_{\text{obs}} | s = 0)$,:

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

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

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


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.

To test for discovery use profile likelihood ratio:

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


So the likelihood ratio statistic for testing $s = 0$ is

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

Approximate Poisson significance (continued)

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

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

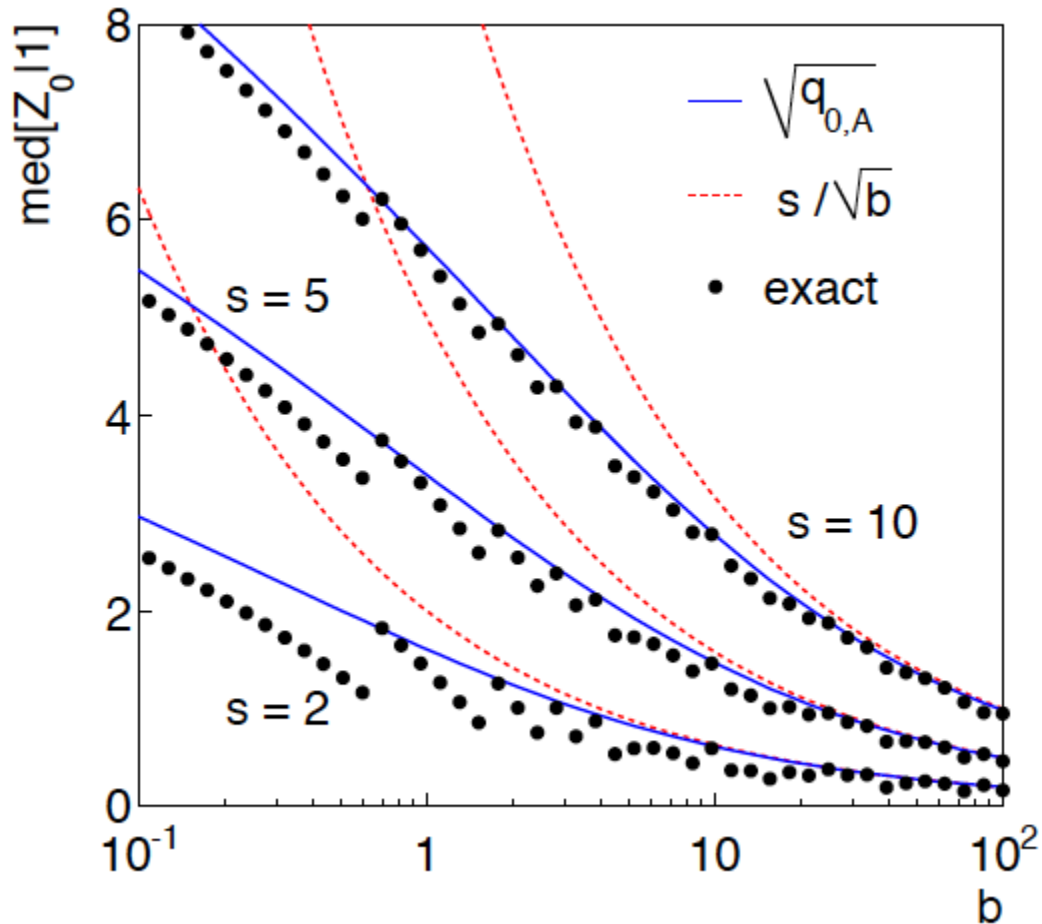
To find median[$Z|s$], let $n \rightarrow s + b$ (i.e., the Asimov data set):

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

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 $s \ll b$.

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.,

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

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):

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

Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

$$\hat{s} = n - m/\tau ,$$

$$\hat{b} = m/\tau ,$$

$$\hat{b}(s) = \frac{n + m - (1 + \tau)s + \sqrt{(n + m - (1 + \tau)s)^2 + 4(1 + \tau)sm}}{2(1 + \tau)} .$$

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

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

Asymptotic significance

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

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

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:

$$n \rightarrow s + b$$

$$m \rightarrow \tau b$$

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

Or use the variance of $\hat{b} = m/\tau$, $V[\hat{b}] \equiv \sigma_b^2 = \frac{b}{\tau}$, to eliminate τ :

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

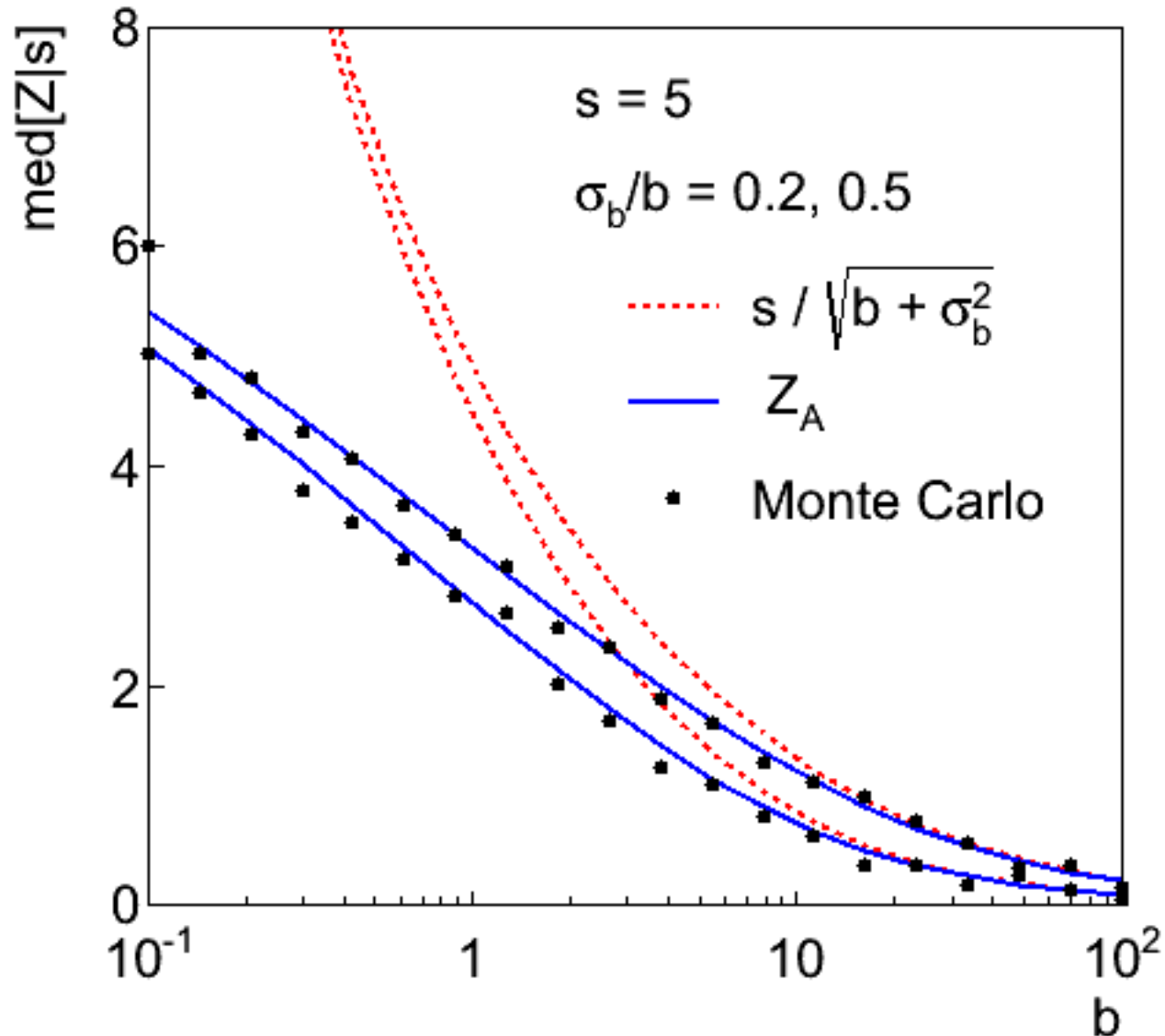
Limiting cases

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

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

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: $s = 5$



Using sensitivity to optimize a cut

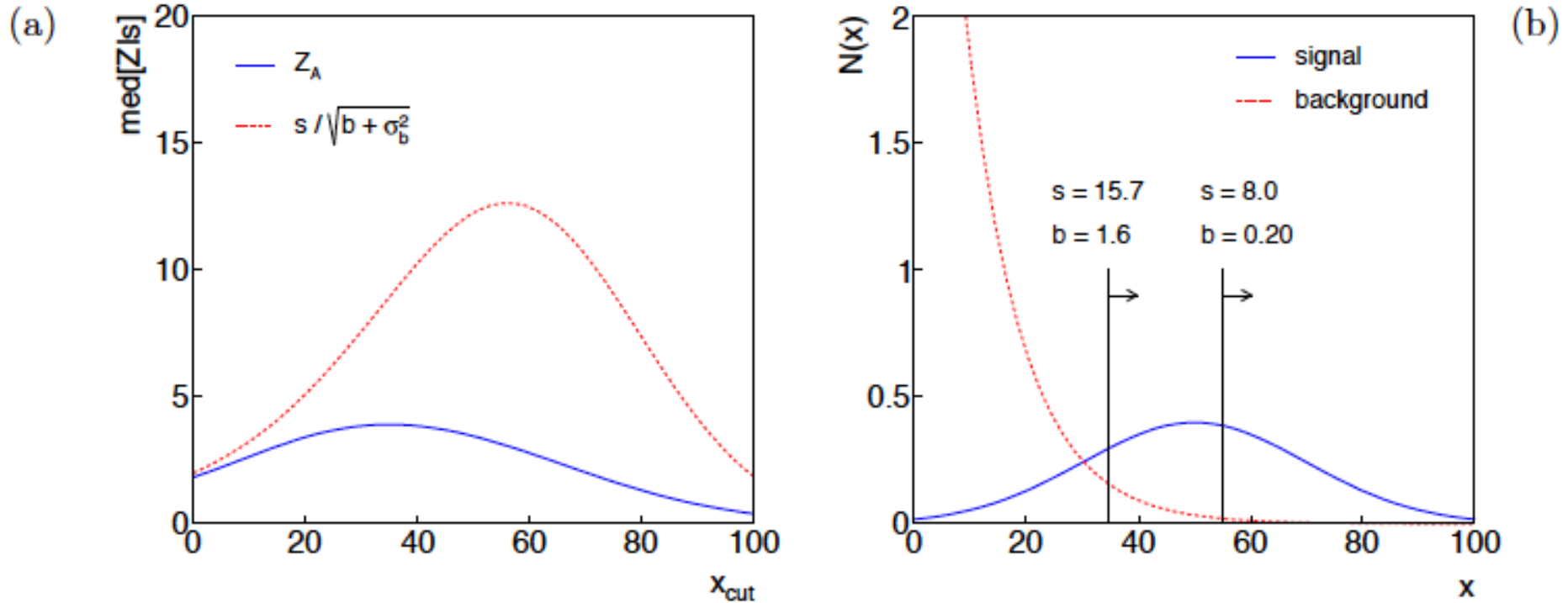


Figure 1: (a) The expected significance as a function of the cut value x_{cut} ; (b) the distributions of signal and background with the optimal cut value indicated.

Bayesian model selection

Fundamentally the probability of a hypothesis H_i in the Bayesian approach is given by its posterior probability given the data: $P(H_i|\mathbf{x})$.

Finding this requires assignment of prior probabilities to all hypotheses that are considered.

We can give the posterior *odds* (ratio of probabilities) for any pair of hypotheses H_i and H_j (use Bayes' theorem; factors of $P(\mathbf{x})$ cancel):

$$\frac{P(H_i|\mathbf{x})}{P(H_j|\mathbf{x})} = \frac{P(\mathbf{x}|H_i)}{P(\mathbf{x}|H_j)} \frac{\pi(H_i)}{\pi(H_j)}$$

posterior odds

Bayes factor

prior odds

See: Kass and Raftery, *Bayes Factors*, J. Am Stat. Assoc 90 (1995) 773.

The Bayes factor

The Bayes factor B_{ij} is the likelihood ratio of the two hypotheses:

$$B_{ij} = \frac{P(\mathbf{x}|H_i)}{P(\mathbf{x}|H_j)} \quad = \text{posterior odds if one takes prior odds equal to one.}$$

The Bayes factor is regarded as measuring the weight of evidence of the data in support of H_i over H_j . and can be used much like a p -value (or Z value).

The Jeffreys scale, analogous to the 5σ rule in Particle Physics:

B_{10}	Evidence against H_0

1 to 3	Not worth more than a bare mention
3 to 20	Positive
20 to 150	Strong
> 150	Very strong

Marginal likelihood (evidence)

If the model H_i contains internal parameters θ_i , then these must be characterized by a prior pdf $\pi_i(\theta_i|H_i)$ and marginalized:

$$P(\mathbf{x}|H_i) = \int P(\mathbf{x}, \theta_i|H_i) d\theta_i = \int P(\mathbf{x}|H_i, \theta_i)\pi_i(\theta_i|H_i) d\theta_i$$

This is called the “marginal likelihood” or “evidence” of H_i .

It is independent of the overall prior probability of H_i

$$\pi(H_i) = \int \pi(H_i, \theta_i) d\theta_i$$

but it depends on the prior pdf for the model’s internal parameters θ_i :

$$\pi_i(\theta_i|H_i) = \frac{\pi(H_i, \theta_i)}{\pi(H_i)}$$

Bayes factor for models with internal parameters

The Bayes factor is thus the ratio of marginal likelihoods for the two models:

$$B_{ij} = \frac{P(\mathbf{x}|H_i)}{P(\mathbf{x}|H_j)} = \frac{\int P(\mathbf{x}|H_i, \boldsymbol{\theta}_i) \pi_i(\boldsymbol{\theta}_i|H_i) d\boldsymbol{\theta}_i}{\int P(\mathbf{x}|H_j, \boldsymbol{\theta}_j) \pi_j(\boldsymbol{\theta}_j|H_j) d\boldsymbol{\theta}_j}$$

Simplifying the notation, the numerator and denominator are both of the form

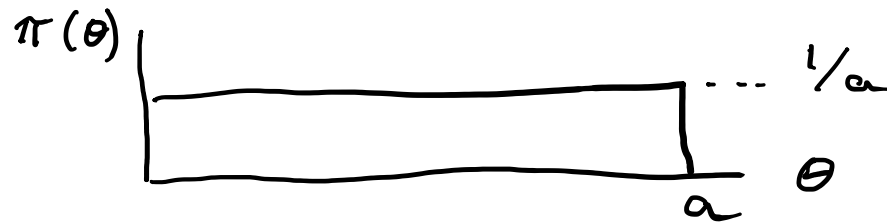
$$m = \int P(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

For high-dimensional $\boldsymbol{\theta}$ these integrals can be very difficult to compute (more on this later).

Priors for Bayes factors

Prior pdfs for the marginal likelihoods used in Bayes factors cannot be improper, i.e., they cannot be defined only up to an arbitrary normalization constant, in which case B_{ij} would not be well defined.

Suppose we try to take a “non-informative” prior to be constant out to some large cut-off, in the hope that the Bayes factor will decouple from it:



In such cases we find that the Bayes factor remains sensitive to the cut-off even for $a \rightarrow \infty$.

So all priors used for Bayes factors must reflect a meaningful degrees of uncertainty about the parameters.

Bayes factor for Poisson counting experiment

Suppose $n \sim \text{Poisson}(s + b)$ with b known. We want to compare

$$H_0 : s = 0 ,$$

$$H_1 : s > 0 .$$

The likelihoods of H_0 and H_1 are

$$L(n|H_0) = \frac{b^n}{n!} e^{-b}$$

$$L(n|s, H_1) = \frac{(s + b)^n}{n!} e^{-(s+b)}$$

Bayes factor for Poisson counting experiment (2)

Suppose the prior pdf for the parameter s in H_1 is:

$$\pi(s|H_1) = \frac{1}{s_{\max}} \quad (0 \leq s \leq s_{\max})$$

The posterior probability for s given n is, assuming H_1 ,

$$\begin{aligned} p(s|n, H_1) &= \frac{L(n|s, H_1)\pi(s|H_1)}{\int L(n|s, H_1)\pi(s|H_1) ds} \\ &= \frac{(s+b)^n e^{-(s+b)}}{\int_0^{s_{\max}} (s+b)^n e^{-(s+b)} ds} \quad (0 \leq s \leq s_{\max}) \\ &= \frac{(s+b)^n e^{-(s+b)}}{\gamma(n+1, s_{\max}+b) - \gamma(n+1, b)} \end{aligned}$$

γ = lower incomplete gamma function

Bayes factor for Poisson counting experiment (3)

In the limit $s_{\max} \rightarrow \infty$ this goes to

$$p(s|n, H_1) = \frac{(s+b)^n e^{-(s+b)}}{\Gamma(n+1) - \gamma(n+1, b)}$$

where $\gamma(a, x) = \int_0^x t^{a-1} e^{-t} dt$

is the lower incomplete gamma function.

Thus the posterior pdf for s given n under assumption of H_1 decouples from s_{\max} in the limit $s_{\max} \rightarrow \infty$, and hence we can use this limiting case e.g. for finding an upper limit (credibility interval) for s .

Bayes factor for Poisson counting experiment (4)

The hypothesis H_0 has no internal parameters so its marginal likelihood is simply $m_0 = L(n|H_0)$.

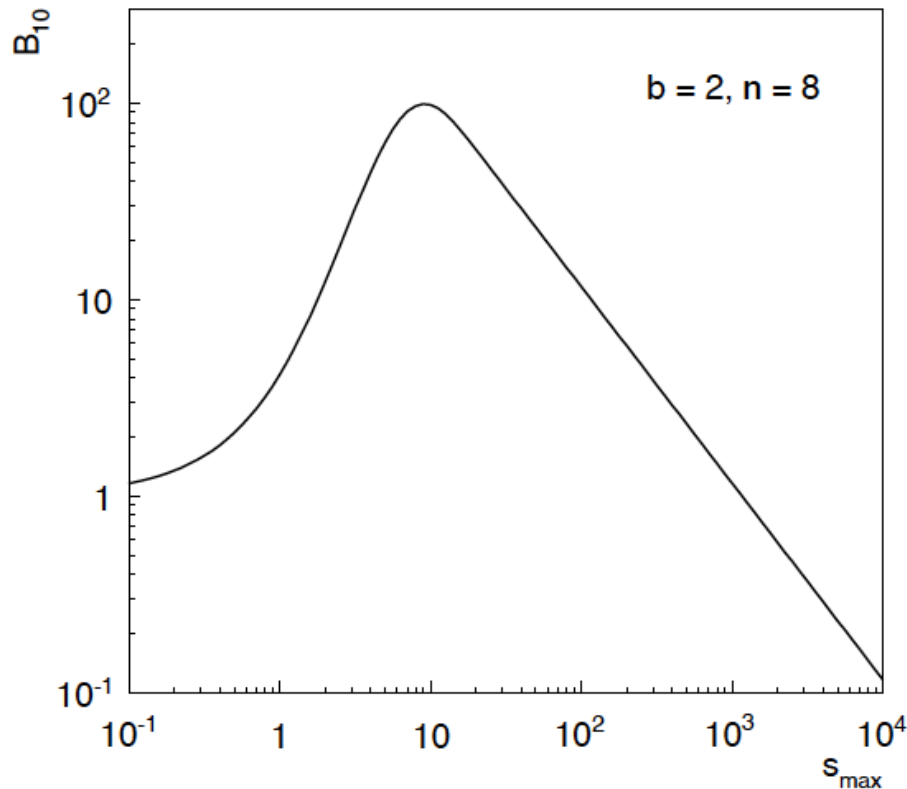
The marginal likelihood of H_1 is

$$\begin{aligned} m_1 &= \int L(n|s, H_1) \pi(s|H_1) ds \\ &= \frac{1}{n! s_{\max}} \int_0^{s_{\max}} (s+b)^n e^{-(s+b)} ds \\ &= \frac{1}{n! s_{\max}} (\gamma(n+1, s_{\max}+b) - \gamma(n+1, b)) \end{aligned}$$

Bayes factor for Poisson counting experiment (5)

So the Bayes factor is

$$B_{10} = \frac{m_1}{m_0} = \frac{1}{s_{\max}} \frac{\gamma(n+1, s_{\max} + b) - \gamma(n+1, b)}{b^n e^{-b}}$$



Example: $b = 2, n = 8$

As s_{\max} increases the data start to favour H_1 .

As s_{\max} increases further, the larger volume of H_1 's parameter space penalizes it (Ockham's razor).

Numerical determination of Bayes factors

Both numerator and denominator of B_{ij} are of the form

$$m = \int L(\vec{x}|\vec{\theta})\pi(\vec{\theta}) d\vec{\theta} \quad \leftarrow \text{‘marginal likelihood’}$$

Various ways to compute these, e.g., using sampling of the posterior pdf (which we can do with MCMC).

Harmonic Mean (and improvements)

Importance sampling

Parallel tempering (\sim thermodynamic integration)

Nested Sampling (MultiNest), ...

Kass and Raftery, *Bayes Factors*, J. Am. Stat. Assoc. 90 (1995) 773-795.

Cong Han and Bradley Carlin, *Markov Chain Monte Carlo Methods for Computing Bayes Factors: A Comparative Review*, J. Am. Stat. Assoc. 96 (2001) 1122-1132.

Phil Gregory, *Bayesian Logical Data Analysis for the Physical Sciences*, Cambridge University Press, 2005.

“Errors on errors”

The uncertainties on estimated systematic errors (“errors on errors”) can in general play an important role in many analyses, see:

G. Cowan, *Statistical Models with Uncertain Error Parameters*, Eur. Phys. J. C (2019) 79:133, arXiv:1809.05778

E. Canonero, A. Brazzale and G. Cowan, *Higher-order asymptotic corrections and their application to the Gamma Variance Model*, Eur. Phys. J. C (2023) 83:1100, arXiv:2304.10574

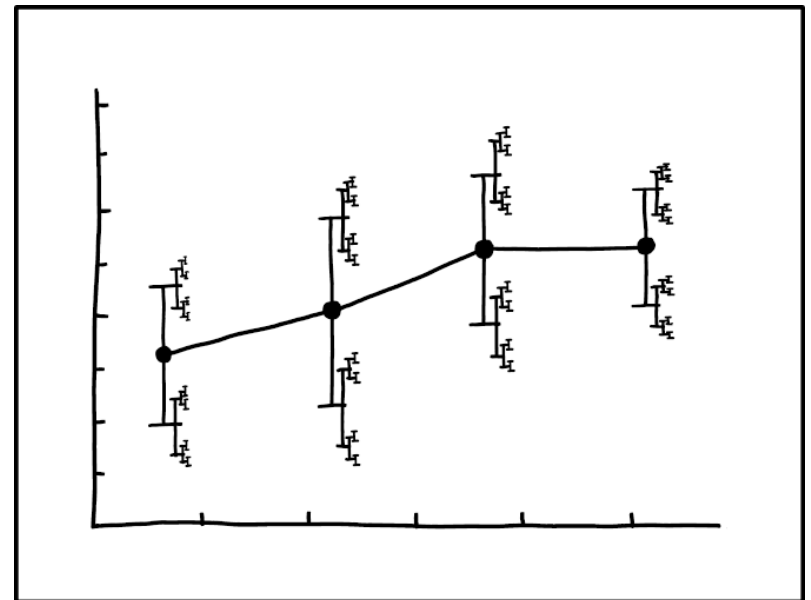
It turns out that models that use errors on errors have qualitatively new, interesting, desirable features:

Sensitivity to outliers reduced.

Confidence intervals sensitive to goodness of fit.

Effect on goodness of fit, p -values, significance.

<https://xkcd.com/2110/>



I DON'T KNOW HOW TO PROPAGATE
ERROR CORRECTLY, SO I JUST PUT
ERROR BARS ON ALL MY ERROR BARS.

Prototype example: curve fitting, averages

Suppose independent
 $y_i \sim \text{Gauss}, i = 1, \dots, N$, with

$$E[y_i] = \varphi(x_i; \boldsymbol{\mu})$$

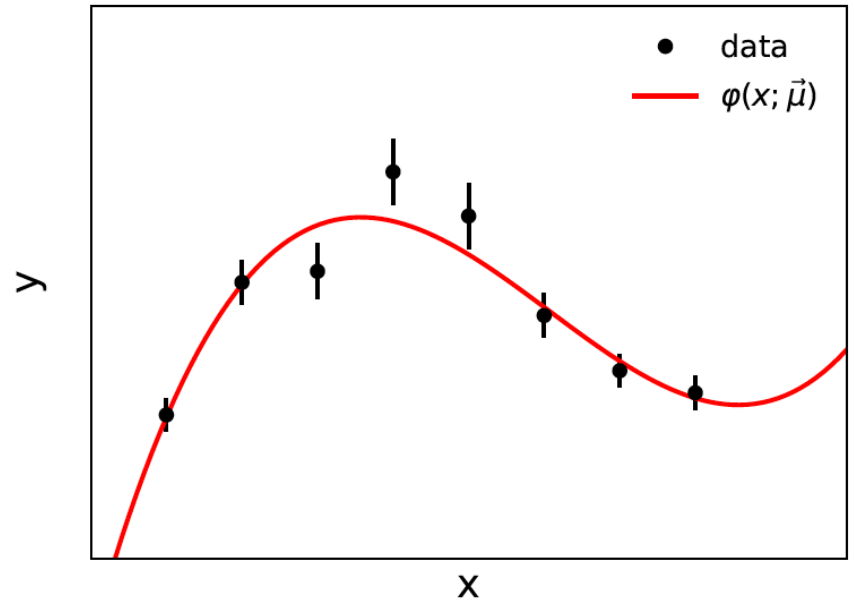
$$V[y_i] = \sigma_i^2$$

$\boldsymbol{\mu} = (\mu_1, \dots, \mu_M)$ are M parameters in the fit function $\varphi(x; \boldsymbol{\mu})$.

If we take the σ_i as known, we have the usual log-likelihood

$$\ln L(\boldsymbol{\mu}) = -\frac{1}{2} \sum_{i=1}^N \frac{(y_i - \varphi(x_i; \boldsymbol{\mu}))^2}{\sigma_i^2}$$

which leads to the Least Squares estimators for $\boldsymbol{\mu}$.



Goodness of fit for Least Squares

In the least-squares approach, the statistic

$$q = -2 \ln \frac{L(\hat{\boldsymbol{\mu}})}{L(\hat{\boldsymbol{\varphi}})} = \sum_{i=1}^N \frac{(y_i - \varphi(x_i; \hat{\boldsymbol{\mu}}))^2}{\sigma_i^2}$$

 Likelihood of saturated model $L(\varphi_1, \dots, \varphi_N)$

provides a measure of goodness of fit. The p -value of the composite hypothesis $\varphi(x_i; \boldsymbol{\mu})$ is

$$p = \int_{q_{\text{obs}}}^{\infty} f(q) dq$$

If the $y_i \sim \text{Gauss}(\varphi(x_i; \boldsymbol{\mu}), \sigma_i)$ then $f(q)$ is chi-squared for $N-M$ degrees of freedom, independent of $\boldsymbol{\mu}$ (Wilks).

Special case: $\varphi(x_i; \boldsymbol{\mu}) = \mu$, i.e., test if the y_i have a common mean μ

$$\rightarrow q \sim \text{chi-square}(N-1) \rightarrow p = 1 - F_{\chi_{N-1}^2}(q)$$

What if the σ_i are not known?

The LS approach assumes that the standard deviations σ_i of the measurements are known.

σ_i = statistical error, usually well estimated from sample size.

σ_i = systematic error:

- related to stat. error of control measurement – well estimated

- related to size of MC event sample – well estimated

- systematic uncertainty from modelling of experiment – could be poorly estimated

- reflects uncertainty resulting from some mathematical approximation (theory error) – could be poorly estimated

In general, we should allow that the σ_i may not be exactly known.

Gamma Variance Model

G. Cowan, EPJC (2019) 79:133

If the σ_i^2 are uncertain, we can take them as adjustable parameters.

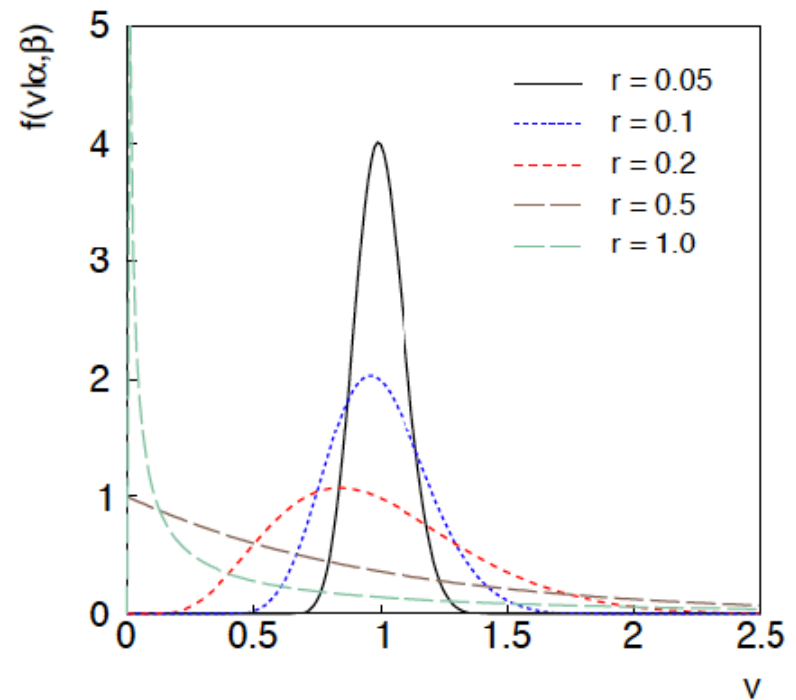
The estimated variances $v_i = s_i^2$ are modeled as gamma distributed.

The likelihood becomes

$$L(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-(y_i - \varphi(x_i; \boldsymbol{\mu}))^2 / 2\sigma_i^2} \frac{\beta_i^{\alpha_i}}{\Gamma(\alpha_i)} v_i^{\alpha_i - 1} e^{-\beta_i v_i}$$

Want $E[v_i] = \sigma_i^2$ $r_i = \frac{\sigma_{v_i}}{2E[v_i]} \approx \frac{\sigma_{s_i}}{E[s_i]}$

→ $\alpha_i = \frac{1}{4r_i^2}$ $\beta_i = \frac{\alpha_i}{\sigma_i^2}$



Profile log-likelihood

One can profile over the σ_i^2 in close form.

The log-profile-likelihood is

$$\ln L'(\boldsymbol{\mu}) = \ln L(\boldsymbol{\mu}, \widehat{\boldsymbol{\sigma}^2}) = -\frac{1}{2} \sum_{i=1}^N \left(1 + \frac{1}{2r_i^2} \right) \ln \left[1 + 2r_i^2 \frac{(y_i - \varphi(x_i; \boldsymbol{\mu}))^2}{v_i} \right]$$

Quadratic terms replace by sum of logs.

Equivalent to replacing Gauss pdf for y_i by Student's t , $\nu_{\text{dof}} = 1/2r_i^2$

Simple program for Student's t average: `stave.py`

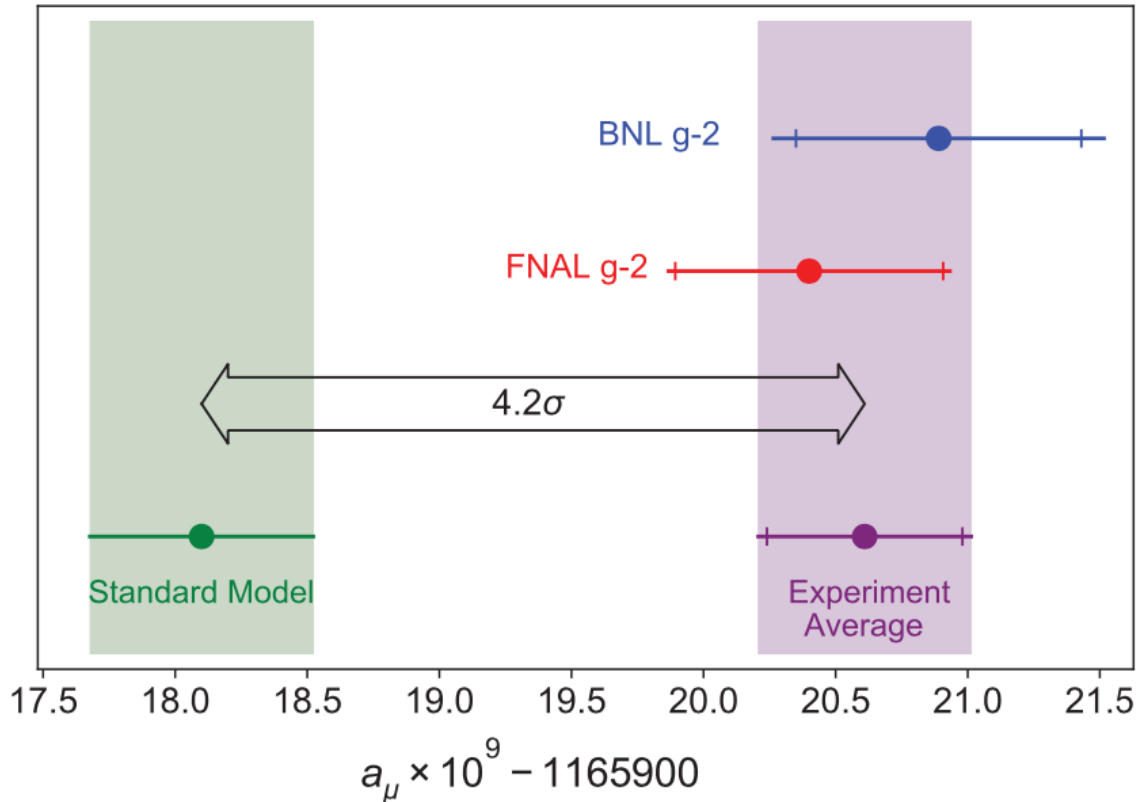
<http://www.pp.rhul.ac.uk/~cowan/stat/stave/>

Application to the muon $g - 2$ anomaly

G. Cowan, Effect of Systematic Uncertainty Estimation on the Muon $g - 2$ Anomaly, EPJ Web of Conferences 258, 09002 (2022), arXiv:2107.02652

The recently measured muon $g - 2$ (ave. of 2006, 2021) disagrees with the Standard Model prediction with a significance of 4.2σ .

Muon $g-2$ Collab., PRL 126, 141801 (2021)



Discrepancy significantly reduced by 2021 lattice-based prediction of Borsanyi et al. (BMW).

Current goal is to investigate sensitivity of significance to error assumptions, so for now focus on the 4.2σ problem.

Muon $g - 2$ ingredients

Using $a_\mu = (g - 2)/2$ $y = a_\mu \times 10^9 - 1165900$

the ingredients of the 4.2σ effect are:

$$y_{\text{exp}} = 20.61 \pm 0.41$$

(ave. of BNL 2006 and FNAL 2021)

$$0.37 \text{ (stat.)} \pm 0.17 \text{ (sys.)}$$

B. Abi et al. (Muon $g-2$ Collaboration), *Measurement of the Positive Muon Anomalous Magnetic Moment to 0.46 ppm*, Phys. Rev. Lett. 126, 141801 (2021).

G. W. Bennett et al. (Muon $g - 2$ Collaboration), *Final report of the E821 muon anomalous magnetic moment measurement at BNL*, Phys. Rev. D 73, 072003 (2006).

$$y_{\text{SM}} = 18.10 \pm 0.43$$

(SM pred. by Muon $g-2$ theory initiative)

$$0.40 \text{ (Had. Vac. Pol.)} \pm 0.18 \text{ (Had. Light-by-Light)}$$

T. Aoyama, N. Asmussen, M. Benayoun, J. Bijnens, and T. Blum et al., *The anomalous magnetic moment of the muon in the standard model*, Phys. Rep. 887, 1 (2020).

Suppose σ_{SM} uncertain

Suppose measurement errors well known, but that the SM theory error σ_{SM} (estimated 0.43) could be uncertain.

This is the largest systematic and probably hardest to estimate.

Treat estimate $v_{\text{SM}} = (0.43)^2$ of variance σ_{SM}^2 as gamma distributed, width from relative uncertainty parameter r_{SM} .

Maximum-likelihood for mean from minimum of

$$Q(\mu) = -2 \ln \frac{L(\mu)}{L_{\text{sat}}}$$
$$= \frac{(y_{\text{exp}} - \mu)^2}{\sigma_{\text{exp}}^2} + \left(1 + \frac{1}{2r_{\text{SM}}^2}\right) \ln \left[1 + 2r_{\text{SM}}^2 \frac{(y_{\text{SM}} - \mu)^2}{v_{\text{SM}}}\right]$$

p -value/significance of common-mean hypothesis

Significance (goodness of fit) from $q = Q(\hat{\mu})$.

Because of non-quadratic term in $Q(\mu)$, distribution of q departs from chi-square(1) for increasing r_{SM} .

Best to get distribution of q from Monte Carlo (and speed up with Bartlett correction – see EPJC (2019) 79:133).

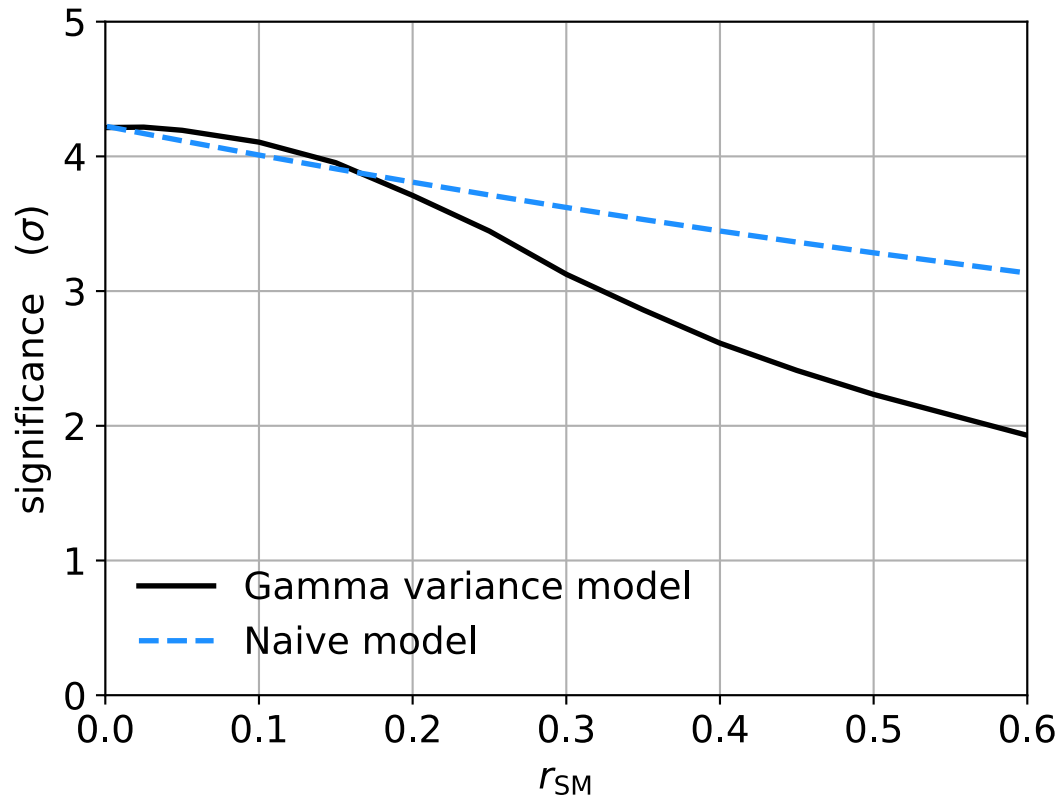
For $r_{\text{SM}} > 0$ distribution of q depends on σ_{SM}^2 . For MC use Maximum-Likelihood estimate (“profile construction”):

$$\widehat{\sigma}_{\text{SM}}^2 = \frac{v_{\text{SM}} + 2r_{\text{SM}}^2(y_{\text{SM}} - \hat{\mu})^2}{1 + 2r_{\text{SM}}^2}$$

$$\text{MC} \rightarrow f(q) \rightarrow p = \int_{q, \text{obs}}^{\infty} f(q) dq \rightarrow \text{significance } Z = \Phi^{-1}(1 - p/2)$$

↖ # of sigmas

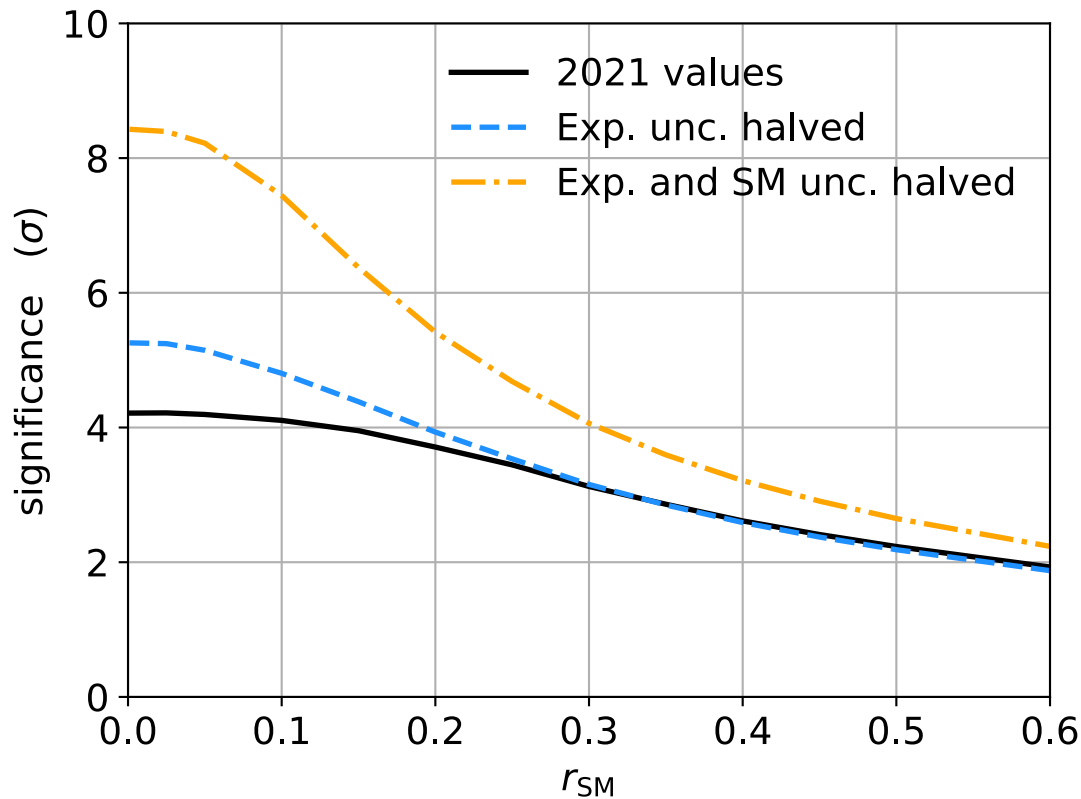
Significance of discrepancy versus r_{SM}



Naive model: use least squares but let $\sigma_{\text{SM}} \rightarrow (1 + r_{\text{SM}})\sigma_{\text{SM}}$

Gamma variance model gives greater decrease in significance for $r_{\text{SM}} \gtrsim 0.2$, e.g., 3.1σ for $r_{\text{SM}} = 0.3$, 2.0σ for $r_{\text{SM}} = 0.6$.

Significance of discrepancy versus r_{SM}



Establishing 4σ effect requires $r_{\text{SM}} \lesssim 0.3$ even if nominal exp. and SM uncertainties become half of present values.

Discussion on muon $g-2$

Including uncertainties on estimates of uncertainties can have large effect on hypothesis test, esp. for high significance.

To establish e.g. a 5σ effect it is crucial to have both:

- small uncertainties

- accurate estimates of those uncertainties ($\sim 20\%$ level)

This is ultimately because the tails of the Gaussian fall off so quickly.

Gamma Variance Model \sim Student's t likelihood with $\nu = 1/2r^2$ degrees of freedom \rightarrow longer tails than Gaussian.

Ongoing discussion with Bogdan Malaescu of Muon $g-2$ Theory Initiative on the HVP uncertainty, see, e.g.,

B. Malaescu et al., https://indico.him.uni-mainz.de/event/11/contributions/80/attachments/50/51/amuWorkshop_Correlations_Malaescu.pdf

M. Davier et al., Eur. Phys. J. C 80 (2020) 241 , arXiv:1908.00921

Discussion on Gamma Variance Model

Other features of Gamma Variance Model (see EPJC (2019) 79:133 and the extra slides)

- averages/fits become less sensitive to outliers;

- confidence intervals linked to goodness of fit;

- straightforward to include multiple correlated error sources.

But... is part of the reason for requiring 5σ for discovery not to account for uncertainties in assigned errors? Is there a trade-off between “errors on errors” and the requirement for discovery?

Best to have most realistic model. If the estimated errors are indeed uncertain, this should be reflected in the model.

Bottom line – it is very difficult to establish convincing evidence for a new physics if relevant uncertainties are estimated in an ad hoc way. We need robust procedures for their assignment.

Finally

Three lectures only enough for a brief discussion of:

Parameter estimation

Hypothesis tests (\rightarrow path to Machine Learning)

Limits (confidence intervals/regions)

Systematics (nuisance parameters)

Bayesian methods, MCMC

A bit beyond... (“errors on errors”)

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

Extra Slides

Harmonic mean estimator

E.g., consider only one model and write Bayes theorem as:

$$\frac{\pi(\boldsymbol{\theta})}{m} = \frac{p(\boldsymbol{\theta}|\mathbf{x})}{L(\mathbf{x}|\boldsymbol{\theta})}$$

$\pi(\boldsymbol{\theta})$ is normalized to unity so integrate both sides,

$$m^{-1} = \int \frac{1}{L(\mathbf{x}|\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} = E_p[1/L]$$

posterior
expectation



Therefore sample $\boldsymbol{\theta}$ from the posterior via MCMC and estimate m with one over the average of $1/L$ (the harmonic mean of L).

M.A. Newton and A.E. Raftery, *Approximate Bayesian Inference by the Weighted Likelihood Bootstrap*, Journal of the Royal Statistical Society B 56 (1994) 3-48.

Called the “worst Monte Carlo method ever”

<https://radfordneal.wordpress.com/2008/08/17/the-harmonic-mean-of-the-likelihood-worst-monte-carlo-method-ever/>

Improvements to harmonic mean estimator

The harmonic mean estimator is numerically very unstable; formally infinite variance (!). A variant (cf. Gelfand and Dey):

Rearrange Bayes thm; multiply both sides by arbitrary pdf $f(\boldsymbol{\theta})$:
$$\frac{f(\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathbf{x})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} = \frac{f(\boldsymbol{\theta})}{m}$$

Integrate over $\boldsymbol{\theta}$:
$$m^{-1} = \int \frac{f(\boldsymbol{\theta})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} = E_p \left[\frac{f(\boldsymbol{\theta})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} \right]$$

Improved convergence if tails of $f(\boldsymbol{\theta})$ fall off faster than $L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$

Note harmonic mean estimator is special case $f(\boldsymbol{\theta}) = \pi(\boldsymbol{\theta})$.

A.E. Gelfand and D.K. Dey, *Bayesian model choice: asymptotics and exact calculations*, Journal of the Royal Statistical Society B 56 (1994) 501-514.

Adaptive Harmonic Mean Integration

A. Caldwell et al., International Journal of Modern Physics A Vol. 35, No. 24 (2020) 2050142



Want to compute $I \equiv \int_{\Omega} f(\lambda) d\lambda$ ($\Omega = \text{support of } f$)

E.g. $f(\lambda) = L(\lambda) \pi(\lambda) = \text{unnormalized target density}$; we can sample from this with MCMC.

Define integral over subvolume Δ of Ω with volume V_{Δ}

$$I_{\Delta} \equiv \int_{\Delta} f(\lambda) d\lambda \quad r \equiv \frac{I_{\Delta}}{I}$$

Adaptive Harmonic Mean Integration (2)

If $f(\lambda)$ not small in Δ , then we can find I_Δ from harmonic mean:

$$E \left[\frac{1}{f(\lambda)} \right]_{\lambda \in \Delta} = \int_{\Delta} \frac{1}{f(\lambda)} \frac{f(\lambda)}{I_\Delta} d\lambda = \frac{1}{I_\Delta} \int_{\Delta} d\lambda = \frac{V_\Delta}{I_\Delta} \approx \frac{1}{N_\Delta} \sum_{\lambda_i \in \Delta} \frac{1}{f(\lambda_i)}$$

Sample λ from $f(\lambda)$ using MCMC, estimate $r = I_\Delta/I$ with fraction of points found in Δ :

$$\hat{r} = \frac{N_\Delta}{N_\Omega}$$

Use these to estimate I :

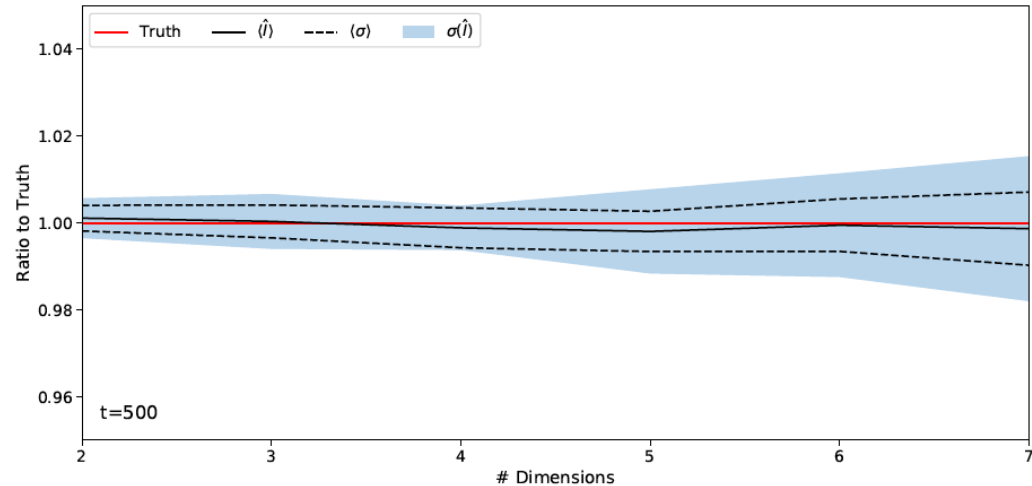
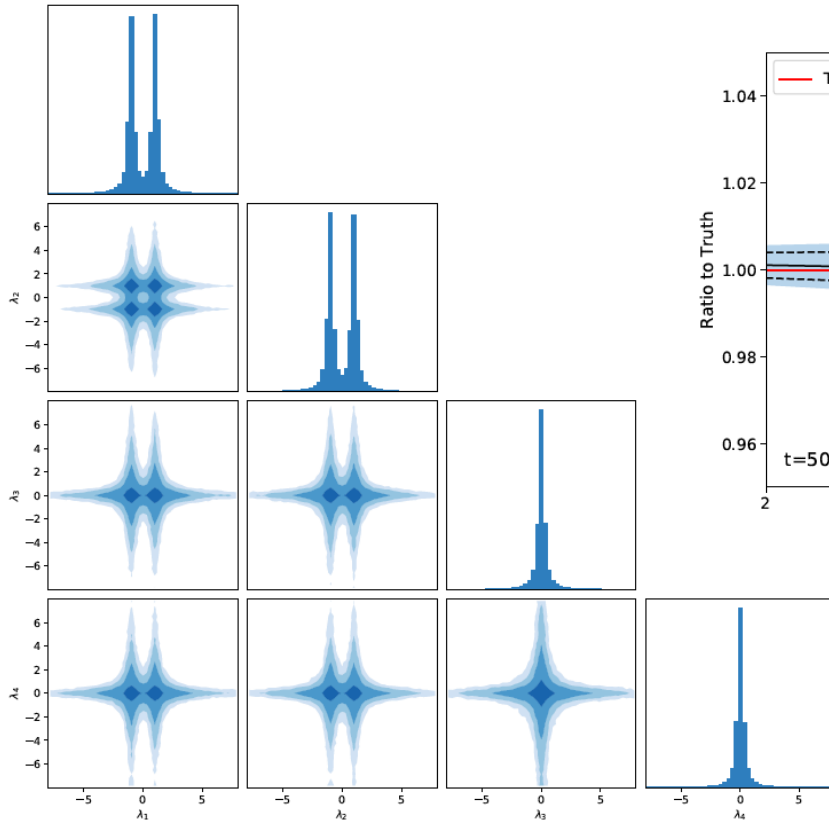
$$\hat{I} = \frac{\hat{I}_\Delta}{\hat{r}} = \frac{N_\Omega V_\Delta}{\sum_{\lambda_i \in \Delta} \frac{1}{f(\lambda_i)}}$$

“The task of estimating our integral, therefore reduces to choosing one or several subspaces Δ — typically small regions around local modes of $f(\lambda)$. The full space Ω over which the integration ought to be performed can be large or even infinite, while this does not affect the outcome of our integral estimate.”

A. Caldwell et al., IJMP A Vol. 35, No. 24 (2020) 2050142

Adaptive Harmonic Mean Integration (3)

Testing AHMI with multimodal multidimensional Cauchy pdf



Challenging pdf because of long tails.

Good results for up to 7 dimensions for MCMC sample size of 10^6 .

Software: Bayesian Analysis Toolkit

<https://github.com/bat/BAT.jl>



Importance sampling

Need pdf $f(\boldsymbol{\theta})$ which we can evaluate at arbitrary $\boldsymbol{\theta}$ and also sample with MC.

The marginal likelihood can be written

$$m = \int \frac{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{f(\boldsymbol{\theta})} f(\boldsymbol{\theta}) d\boldsymbol{\theta} = E_f \left[\frac{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{f(\boldsymbol{\theta})} \right]$$

Sample $\boldsymbol{\theta} \sim f(\boldsymbol{\theta})$, compute average of $L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})/f(\boldsymbol{\theta})$.

Best convergence when $f(\boldsymbol{\theta})$ approximates shape of $L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$.

Use for $f(\boldsymbol{\theta})$ e.g. multivariate Gaussian with mean and covariance estimated from posterior.

Nested sampling

J. Skilling, Bayesian Analysis, No. 4, pp. 833-860 (2006)

We want to compute $Z = \text{evidence} = \int L dX$ $L = L(\theta)$
 $dX = \pi(\theta)d\theta$

Can add up portions of X (equivalently, θ) space in any order. Use

$$X(\lambda) = \int_{L(\theta) > \lambda} \pi(\theta) d\theta$$

X near 1 means low λ , all of θ space included.

Write inverse function as $L(X(\lambda)) \equiv \lambda$ so that the desired result is

$$Z = \int_0^1 L(X) dX$$

Elements of θ space are sorted by decreasing likelihood.

Nested sampling (2)

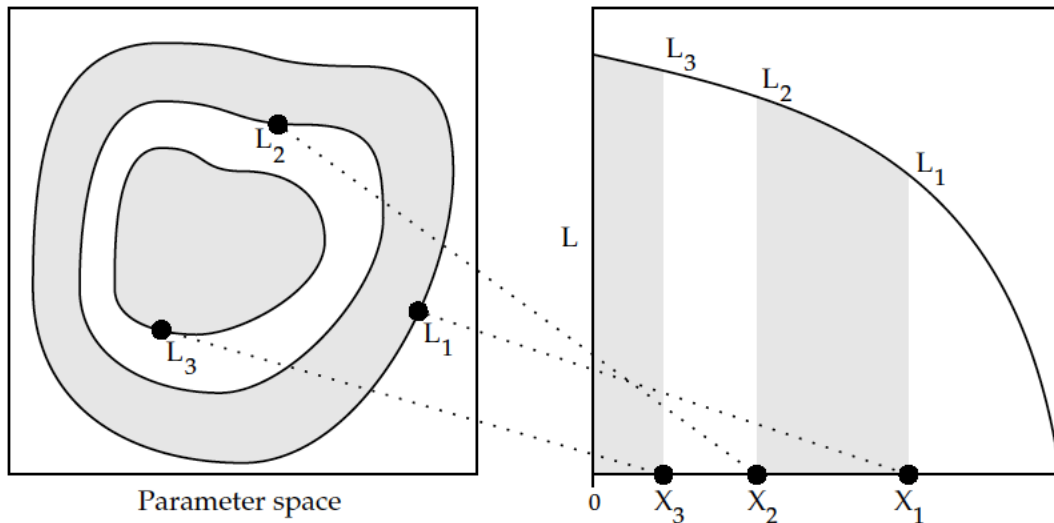


Figure 3: Nested likelihood contours are sorted to enclosed prior mass X .

The evidence Z
is the area under
the curve of $L(X)$.

Computational challenge is to sample θ space from prior subject to constraint $L(\theta) > \lambda$. Software: MultiNest

Farhan Feroz, Mike Hobson, Mon. Not. Roy. Astron. Soc., 384, 2, 449-463 (2008);
arXiv:0704.3704,

F. Feroz, M.P. Hobson, M. Bridges, Mon. Not. Roy. Astron. Soc. 398: 1601-1614, 2009;
arXiv:0809.3437

F. Feroz, M.P. Hobson, E. Cameron, A.N. Pettitt, arXiv:1306.2144

Full likelihood for gamma variance model

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}, \boldsymbol{\sigma}_u^2) = P(\mathbf{y} | \boldsymbol{\mu}, \boldsymbol{\theta}) \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_{u_i}^2}} e^{-(u_i - \theta_i)^2 / 2\sigma_{u_i}^2}$$

$$\times \frac{\beta_i^{\alpha_i}}{\Gamma(\alpha_i)} v_i^{\alpha_i - 1} e^{-\beta_i v_i}, \quad \begin{aligned} \alpha_i &= 1/4r_i^2 \\ \beta_i &= \alpha_i / \sigma_{u_i}^2 \end{aligned}$$

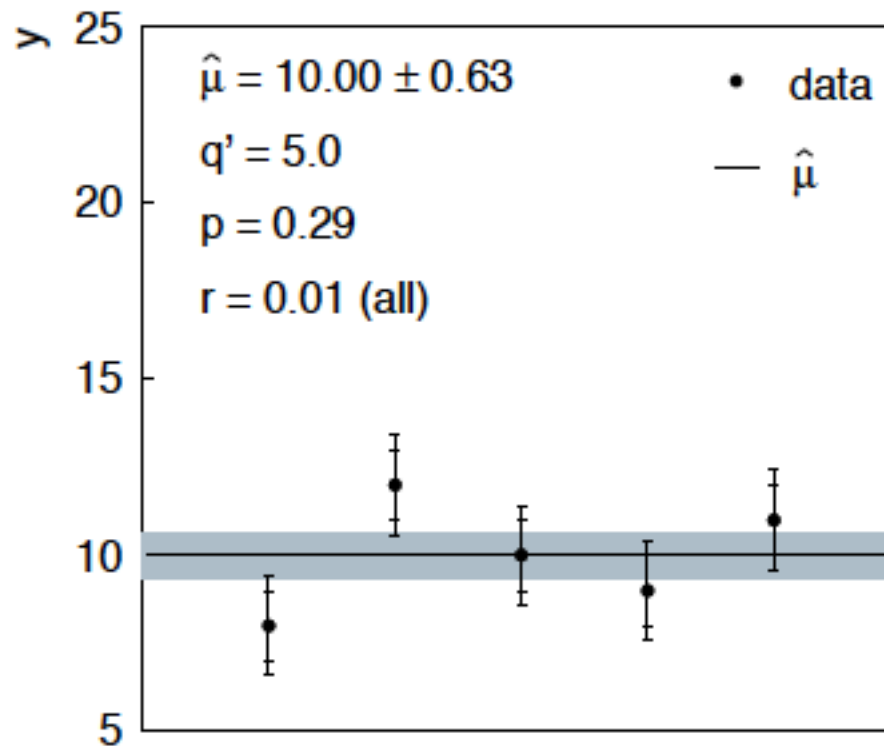
Treated like data: y_1, \dots, y_L (the primary measurements)
 u_1, \dots, u_N (estimates of nuisance par.)
 v_1, \dots, v_N (estimates of variances of estimates of NP)

Adjustable parameters: μ_1, \dots, μ_M (parameters of interest)
 $\theta_1, \dots, \theta_N$ (nuisance parameters)
 $\sigma_{u,1}, \dots, \sigma_{u,N}$ (sys. errors = std. dev. of of NP estimates)

Fixed parameters: r_1, \dots, r_N (rel. err. in estimate of $\sigma_{u,i}$)

Sensitivity of average to outliers

Suppose we average 5 values, $y = 8, 9, 10, 11, 12$, all with stat. and sys. errors of 1.0, and suppose negligible error on error (here take $r = 0.01$ for all).



inner error bars

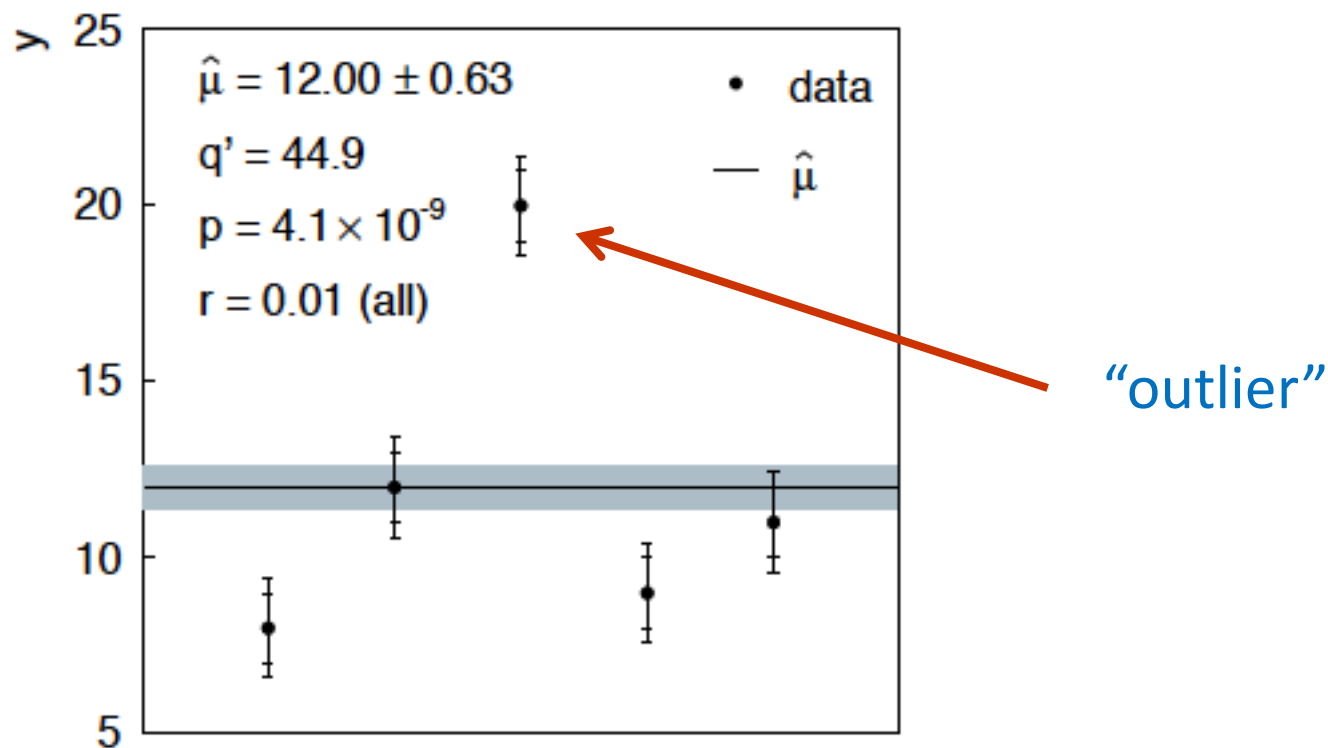
$$= \sigma_{y,i}$$

outer error bars

$$= (\sigma_{y,i}^2 + \sigma_{u,i}^2)^{1/2}$$

Sensitivity of average to outliers (2)

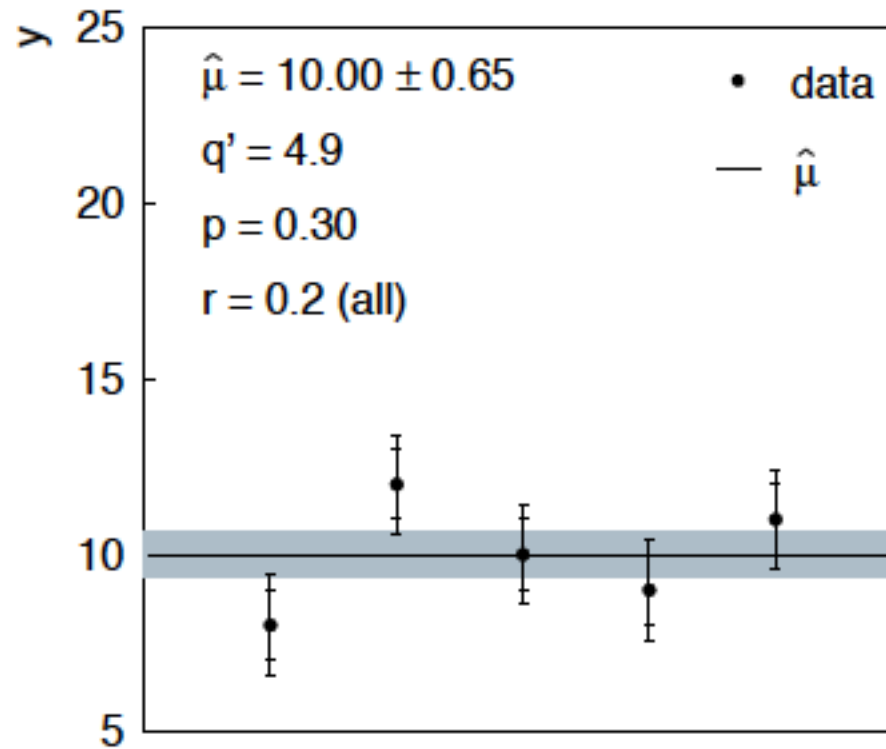
Now suppose the measurement at 10 had come out at 20:



Estimate pulled up to 12.0, size of confidence interval \sim unchanged (would be exactly unchanged with $r \rightarrow 0$).

Average with all $r = 0.2$

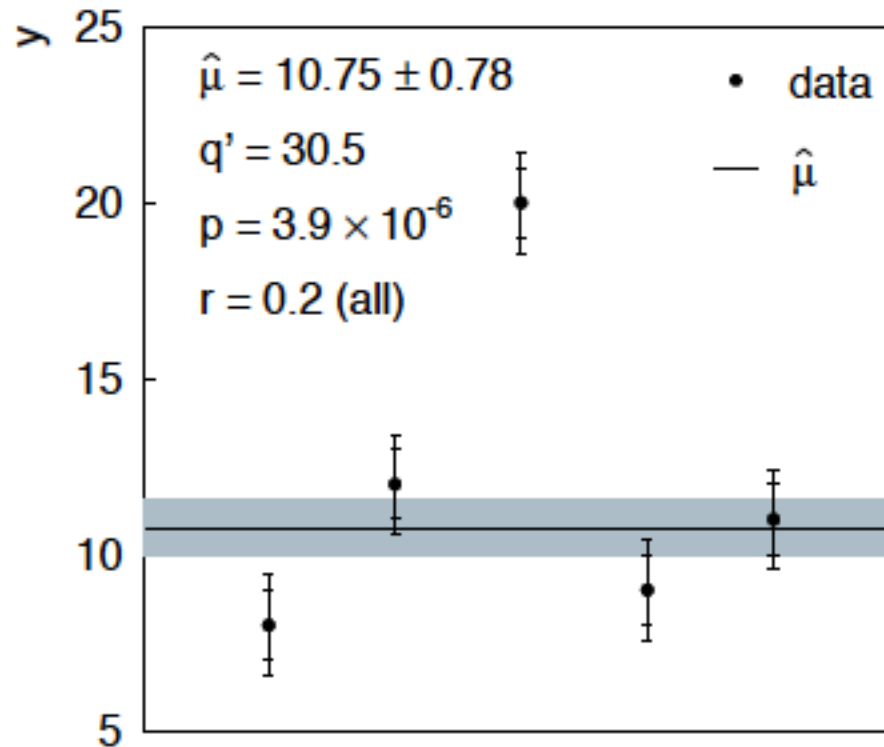
If we assign to each measurement $r = 0.2$,



Estimate still at 10.00, size of interval moves $0.63 \rightarrow 0.65$

Average with all $r = 0.2$ with outlier

Same now with the outlier (middle measurement $10 \rightarrow 20$)



Estimate $\rightarrow 10.75$ (outlier pulls much less).

Half-size of interval $\rightarrow 0.78$ (inflated because of bad g.o.f.).

Naive approach to errors on errors

Naively one might think that the error on the error in the previous example could be taken into account conservatively by inflating the systematic errors, i.e.,

$$\sigma_{u_i} \rightarrow \sigma_{u_i} (1 + r_i)$$

But this gives

$$\hat{\mu} = 10.00 \pm 0.70 \quad \text{without outlier (middle meas. 10)}$$

$$\hat{\mu} = 12.00 \pm 0.70 \quad \text{with outlier (middle meas. 20)}$$

So the sensitivity to the outlier is not reduced and the size of the confidence interval is still independent of goodness of fit.

Correlated uncertainties

The phrase “correlated uncertainties” usually means that a single nuisance parameter affects the distribution (e.g., the mean) of more than one measurement.

For example, consider measurements y , parameters of interest μ , nuisance parameters θ with

$$E[y_i] = \varphi_i(\mu, \theta) \approx \varphi_i(\mu) + \sum_{j=1}^N R_{ij} \theta_j$$

That is, the θ_i are defined here as contributing to a bias and the (known) factors R_{ij} determine how much θ_j affects y_i .

As before suppose one has independent control measurements $u_i \sim \text{Gauss}(\theta_i, \sigma_{ui})$.

Correlated uncertainties (2)

The total bias of y_i can be defined as
$$b_i = \sum_{j=1}^N R_{ij} \theta_j$$

which can be estimated with
$$\hat{b}_i = \sum_{j=1}^N R_{ij} u_j$$

These estimators are correlated having covariance

$$U_{ij} = \text{cov}[\hat{b}_i, \hat{b}_j] = \sum_{k=1}^N R_{ik} R_{jk} V[u_k]$$

In this sense the present method treats “correlated uncertainties”, i.e., the control measurements u_i are independent, but nuisance parameters affect multiple measurements, and thus bias estimates are correlated.