

# Reflections on 20+ years of F-C: Hypothesis testing of a point null vs a continuous alternative

**Bob Cousins**

**Univ. of California, Los Angeles**

**PhyStat-DM at Stockholm**

**August 2, 2019**

Based on (a small part of) my writeup,

“Lectures on Statistics in Theory: Prelude to Statistics in Practice”

<https://arxiv.org/abs/1807.05996> and references therein.

*[Section numbers in today's slides refer to this arxiv post.](https://arxiv.org/abs/1807.05996)*

# Notation

**x denotes observable(s), can be multi-D**

**More generally, x is any convenient or useful function of the observable(s), and is called a “statistic” or “test statistic”**

**$\mu$  denotes parameter(s) (sometimes we use  $\theta$ )**

**$p(x|\mu)$  is probability/pdf characterizing everything that determines the probabilities (densities) of the observations, from laws of physics to experiment setup and protocol**

**$p(x|\mu)$  is called the “statistical model” or simply “the model” by statisticians.**

# Basic notions of confidence intervals (Sec. 6.2)

In two sentences:

Given the model  $p(x|\mu)$  and the observed value  $x_{\text{obs}}$ , for what values of  $\mu$  is  $x_{\text{obs}}$  an “extreme” value of  $x$ ?

Include in the confidence interval  $[\mu_1, \mu_2]$  those values of  $\mu$  for which  $x_{\text{obs}}$  is *not* “extreme”.

# Basic notions of confidence intervals (Sec. 6.2)

In two sentences:

Given the model  $p(x|\mu)$  and the observed value  $x_{\text{obs}}$ , for what values of  $\mu$  is  $x_{\text{obs}}$  an “extreme” value of  $x$ ?

Include in the confidence interval  $[\mu_1, \mu_2]$  those values of  $\mu$  for which  $x_{\text{obs}}$  is *not* “extreme”.

To be well-defined, the first point needs to be supplemented:

1) In order to define “extreme”, one needs to choose an *ordering principle* for  $x$  applicable to each  $\mu$ : *high rank means not extreme*.

2) One also needs to specify what *fraction* of values of  $x$  are *not* considered extreme. Called the *confidence level C.L.*;  $\alpha = 1 - \text{C.L.}$

# Basic notions of confidence intervals (cont.)

## Three common ordering choices in 1D

(when  $p(x|\mu)$  is such that higher  $\mu$  implies higher average  $x$ ):

1. Order  $x$  from largest to smallest.  
Leads to confidence intervals known as *upper limits* on  $\mu$ .
2. Order  $x$  from smallest to largest. Leads to *lower limits* on  $\mu$ .
3. Order  $x$  using *central* quantiles of  $p(x|\mu)$ .  
Gives *central* confidence intervals for  $\mu$ .

These orderings apply only when  $x$  is 1D

# Basic notions of confidence intervals (cont.)

So, one-sentence definition of confidence interval:

The *confidence interval*  $[\mu_1, \mu_2]$  contains those values of  $\mu$  for which  $x_{\text{obs}}$  is *not* “extreme” at the chosen C.L., *given the ordering*.

See Section 6.8 (and F-C paper) for graphical equivalent that we call “Neyman’s construction”, and “confidence belts”

# Confidence Intervals and Coverage (Sec. 6.11)

Let  $\mu_t$  be the unknown true value of  $\mu$ . In repeated experiments, confidence intervals will have different endpoints  $[\mu_1, \mu_2]$ , since the endpoints are functions of the randomly sampled  $x$ .

A little thought will convince you that a fraction C.L. =  $1 - \alpha$  of confidence intervals so obtained will contain (“cover”) the fixed but unknown  $\mu_t$ . I.e.,

$$P(\mu_t \in [\mu_1, \mu_2]) = \text{C.L.} = 1 - \alpha. \quad (\text{Definition of coverage})$$

# Confidence Intervals and Coverage (Sec. 6.11)

Let  $\mu_t$  be the unknown true value of  $\mu$ . In repeated experiments, confidence intervals will have different endpoints  $[\mu_1, \mu_2]$ , since the endpoints are functions of the randomly sampled  $x$ .

A little thought will convince you that a fraction C.L. =  $1 - \alpha$  of confidence intervals so obtained will contain (“cover”) the fixed but unknown  $\mu_t$ . I.e.,

$$P(\mu_t \in [\mu_1, \mu_2]) = \text{C.L.} = 1 - \alpha. \quad (\text{Definition of coverage})$$

In this (frequentist) equation,  $\mu_t$  is *fixed* and unknown. The endpoints  $\mu_1, \mu_2$  are the random variables (!).

Coverage is a property of the *set* of confidence intervals, not of any one interval.

See backup re Neyman’s point that expts need not be the same.

Discrete observations and/or nuisance parameters typically make exact coverage unobtainable – see writeup.



# Fourth ordering: Likelihood ratios (Sec. 6.7)

4. Order  $x$  using *likelihood ratio*  $\mathcal{L}(x|\mu) / \mathcal{L}(x|\mu_{\text{best fit}})$ , advocated by F-C.

Unified approach to the classical statistical analysis of small signals

Gary J. Feldman<sup>\*</sup>

*Department of Physics, Harvard University, Cambridge, Massachusetts 02138*

Robert D. Cousins<sup>†</sup>

*Department of Physics and Astronomy, University of California, Los Angeles, California 90095*

Phys. Rev. D57 3873 (1998):

***Ordering applies (in principle) for arbitrary dimensions of  $x$ ,  $\mu$ .***

**We looked “everywhere” in literature on confidence intervals, did not see this ordering used for intervals. *Was it really new?***

**Instructive twist as our paper was in proof!**

**For that we must first turn to...**

# Hypothesis testing

**Many special cases, including:**

- a) A given functional form (“model”) vs another functional form. Also known as “model selection”.**
- b) Within the same functional form, a single value of a parameter (say 0 or 1) vs all other values. The model with the single value is *nested* within the model with all other values.**
- c) Goodness of Fit: A given functional form against all other (unspecified) functional forms (aka “model checking”)**

**(Section 2.3)**

# Hypothesis testing

Many special cases, including:

- a) A given functional form (“model”) vs another functional form. Also known as “model selection”.
- b) Within the same functional form, a single value of a parameter (say 0 or 1) vs all other values. The model with the single value is *nested* within the model with all other values.
- c) Goodness of Fit: A given functional form against all other (unspecified) functional forms (aka “model checking”)

(Section 2.3)

# Nested Hypothesis Testing is common in HEP

There is an undetermined parameter  $\mu$  in  $H_1$ , and  $H_0$  corresponds to a particular parameter value  $\mu_0$  (e.g., zero, SM prediction, or  $\infty$ ).

$H_0: \mu = \mu_0$  (the “point null”, or “sharp hypothesis”)

$H_1: \mu \neq \mu_0$  (the “continuous alternative”).

Common examples:

1) Signal strength  $\mu$  of new physics: null  $\mu_0 = 0$ , alternative  $\mu > 0$

2) Higgs boson  $\rightarrow \gamma\gamma$  before observation, signal strength  $\mu$ :  
null  $\mu_0 = 0$ , alternative  $\mu > 0$

3) Higgs boson  $\rightarrow \gamma\gamma$  after observation:  
null  $\mu_0 = \text{SM prediction}$ , alternative is any other  $\mu \neq \mu_0$

(Section 7.3)

# Nested Hypothesis Testing in Neutrino Physics

1a)  $\nu$  mixing angle  $\theta_{23}$  *before 1998*: null  $\theta_{23} = 0$ , alternative  $\theta_{23} \neq 0$

1b)  $\nu$  mixing angle  $\theta_{23}$  *after 1998*: null  $\theta_{23} = 45^\circ$ , alternative  $\theta_{23} \neq 45^\circ$

2a) CP violation phase  $\delta$  before it is observed:

Two-point-null: “ $\delta = 0$  or  $\delta = \pi$ ” vs alternative: all other  $\delta$

2b) After two-point null for  $\delta$  is rejected: maybe a theorist has a “predicted” value of  $\delta$  to test

# Classical Frequentist Hypothesis Testing

For null hypothesis  $H_0$ , order possible observations  $x$  from least extreme to most extreme, using an ordering principle (which can depend on  $H_1$  as well). Choose a cutoff  $\alpha$  (smallish number).

Then “reject”  $H_0$  if observed  $x_{\text{obs}}$  is in the *most* extreme fraction  $\alpha$  of observations  $x$  (generated under  $H_0$ ). By construction,

$\alpha$  = probability (with  $x$  generated according to  $H_0$ ) of rejecting  $H_0$  when it is true, i.e., false discovery claim (Type I error)

[See elsewhere for Type II error prob  $\beta$  ]

# Nested Hypothesis Testing: Duality with Intervals

In classical/frequentist formalism (but *not* Bayesian formalism), the theory of these hypo tests maps to that of confidence intervals:

Having observed data  $x_{\text{obs}}$ , suppose the 95% C.L. confidence interval for  $\mu$  is  $[\mu_1, \mu_2]$ .

This contains all values of  $\mu$  for which observed  $x_{\text{obs}}$  is ranked in the *least* extreme 95% of possible outcomes  $x$  according to  $p(x|\mu)$  and the ordering principle in use.

# Nested Hypothesis Testing: Duality with Intervals

In classical/frequentist formalism (but *not* Bayesian formalism), the theory of these hypo tests maps to that of confidence intervals:

Having observed data  $x_{\text{obs}}$ , suppose the 95% C.L. confidence interval for  $\mu$  is  $[\mu_1, \mu_2]$ .

This contains all values of  $\mu$  for which observed  $x_{\text{obs}}$  is ranked in the *least* extreme 95% of possible outcomes  $x$  according to  $p(x|\mu)$  and the ordering principle in use.

Now suppose we wish to test  $H_0$  vs  $H_1$  at Type I error prob  $\alpha = 5\%$ . We reject  $H_0$  if  $x_{\text{obs}}$  is ranked in the *most* extreme 5% of  $x$  according to  $p(x|\mu)$  and the ordering principle in use.



# Nested Hypothesis Testing: Duality with Intervals

In classical/frequentist formalism (but *not* Bayesian formalism), the theory of these hypo tests maps to that of confidence intervals:

Having observed data  $x_{\text{obs}}$ , suppose the 95% C.L. confidence interval for  $\mu$  is  $[\mu_1, \mu_2]$ .

This contains all values of  $\mu$  for which observed  $x_{\text{obs}}$  is ranked in the *least* extreme 95% of possible outcomes  $x$  according to  $p(x|\mu)$  and the ordering principle in use.

Now suppose we wish to test  $H_0$  vs  $H_1$  at Type I error prob  $\alpha = 5\%$ . We reject  $H_0$  if  $x_{\text{obs}}$  is ranked in the *most* extreme 5% of  $x$  according to  $p(x|\mu)$  and the ordering principle in use.

Comparing the two procedures, we see:

**Reject  $H_0$  at  $\alpha=5\%$  iff  $\mu_0$  is *not* in 95% C.L. conf. interval  $[\mu_1, \mu_2]$ .**

Use of the duality is referred to as “**inverting a test**” to obtain confidence intervals, and vice versa. (Section 7.4)

# Duality in Nested Hypothesis Testing

While F-C was “in proof”, Gary realized that “our” intervals were simply those obtained by “inverting” the classic “exact” LR hypothesis test (which specifies LR ordering) in Kendall and Stuart.

It was all on 1¼ pages, plus profiling nuisance parameters!

See Gary’s Fermilab talk, “Journeys of an Accidental Statistician”,  
<http://users.physics.harvard.edu/~feldman/Journeys.pdf>

This was of course good!  
It led to rapid inclusion in PDG RPP.

CHAPTER 22

## LIKELIHOOD RATIO TESTS AND TEST EFFICIENCY

### The LR statistic

**22.1** The ML method discussed in Chapter 18 is a constructive method of obtaining estimators which, under certain conditions, have desirable properties. A method of test construction closely allied to it is the likelihood ratio (LR) method, proposed by Neyman and Pearson (1928). It has played a role in the theory of tests analogous to that of the ML method in the theory of estimation.

As before, we have the LF

$$L(x|\theta) = \prod_{i=1}^n f(x_i|\theta),$$

where  $\theta = (\theta_r, \theta_s)$  is a vector of  $r + s = k$  parameters ( $r \geq 1, s \geq 0$ ) and  $x$  may also be a vector. We wish to test the hypothesis

$$H_0 : \theta_r = \theta_{r0}, \quad (22.1)$$

which is composite unless  $s = 0$ , against

$$H_1 : \theta_r \neq \theta_{r0}.$$

We know that there is generally no UMP test in this situation, but that there may be a UMPU test – cf. **21.31**.

The LR method first requires us to find the ML estimators of  $(\theta_r, \theta_s)$ , giving the unconditional maximum of the LF

$$L(x|\hat{\theta}_r, \hat{\theta}_s), \quad (22.2)$$

and also to find the ML estimators of  $\theta_s$ , when  $H_0$  holds,<sup>1</sup> giving the conditional maximum of the LF

$$L(x|\theta_r, \hat{\theta}_s). \quad (22.3)$$

$\hat{\theta}_s$  in (22.3) has been given a double circumflex to emphasize that it does not in general coincide with  $\hat{\theta}_s$  in (22.2). Now consider the likelihood ratio<sup>2</sup>

$$l = \frac{L(x|\theta_r, \hat{\theta}_s)}{L(x|\hat{\theta}_r, \hat{\theta}_s)}. \quad (22.4)$$

Since (22.4) is the ratio of a conditional maximum of the LF to its unconditional maximum, we clearly have

$$0 \leq l \leq 1. \quad (22.5)$$

Intuitively,  $l$  is a reasonable test statistic for  $H_0$ : it is the maximum likelihood under  $H_0$  as a fraction of its largest possible value, and large values of  $l$  signify that  $H_0$  is reasonably acceptable. The critical region for the test statistic is therefore

$$l \leq c_\alpha, \quad (22.6)$$

where  $c_\alpha$  is determined from the distribution  $g(l)$  of  $l$  to give a size- $\alpha$  test, that is,

$$\int_0^{c_\alpha} g(l) dl = \alpha. \quad (22.7)$$

Neither maximum value of the LF is affected by a change of parameter from  $\theta$  to  $\tau(\theta)$ , the ML estimator of  $\tau(\theta)$  being  $\tau(\hat{\theta})$  – cf. **18.3**. Thus the LR statistic is invariant under reparametrization.

# Famous confusion re Gaussian $p(x|\mu)$ where $\mu \geq 0$

It is *crucial* to distinguish between the data  $x$ , which *can* be negative (no problem), and a parameter  $\mu$  such as mass or signal strength, for which negative values *do not exist in the model*.

I.e., for mass  $\mu < 0$ ,  $p(x|\mu)$  does not exist: You would not know how to simulate the physics of detector response for *mass*  $< 0$ .

Constraint  $\mu \geq 0$  has *nothing* to do with a Bayesian prior for  $\mu$  !!! It's in the *model* (and hence in  $\mathcal{L}(\mu)$ ).

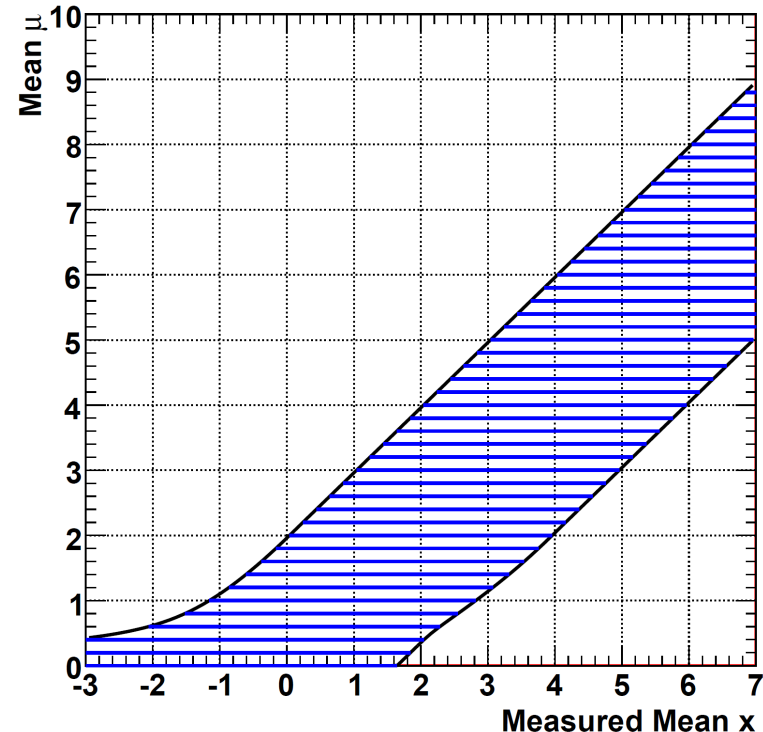
# Famous confusion re Gaussian $p(x|\mu)$ where $\mu \geq 0$

It is *crucial* to distinguish between the data  $x$ , which *can* be negative (no problem), and a parameter  $\mu$  such as mass or signal strength, for which negative values *do not exist in the model*.

I.e., for mass  $\mu < 0$ ,  $p(x|\mu)$  does not exist: You would not know how to simulate the physics of detector response for *mass*  $< 0$ . Constraint  $\mu \geq 0$  has *nothing* to do with a Bayesian prior for  $\mu$  !!! It's in the *model* (and hence in  $\mathcal{L}(\mu)$ ).

The confusion is encouraged since we often refer to  $x$  as the “measured value of  $\mu$ ”, and say that  $x < 0$  is “unphysical” – bad habits!

A proper confidence belt has  $x$  of both signs, only non-negative  $\mu \geq 0$ . Example: Construction on right is LR ordering advocated by F-C (Sections 6.9, 14)



# Rollout of F-C

Posted to arxiv Nov. 1997.

Published in Phys. Rev. D on April 1, 1998

Gary goes to Takayama Japan for Neutrino '98.

**Official Super-Kamiokande Press Release from Japan MEDIA  
ADVISORY for afternoon June 5, 1998, Takayama, Japan:**

**EVIDENCE FOR MASSIVE NEUTRINOS”**

**[Atmospheric neutrino oscillations]**

**Long email to me from Gary on June 5, 1998, detailing widespread interest in F-C, noting:**

“

**Most people seem to have heard about our paper, or, if not, are starting to ask about it.**

**Long email to me from Gary on June 5, 1998, detailing widespread interest in F-C, noting:**

“

**Most people seem to have heard about our paper, or, if not, are starting to ask about it.**

***The most disconcerting thing is that I keep getting introduced as ‘Feldman, of Feldman and Cousins.’***

”

# 20 years of experience with F-C

Lots of experience in HEP, many find it useful, especially when:

- ★ A model parameter is bounded (mass, cross section, sin/cosine of an angle, etc.); and/or
- ★ When log-likelihood is non-Gaussian (so Wilks's Theorem is inaccurate); multiply connected confidence regions; and/or
- ★ The interesting parameter space is  $>1D$ , where LR ordering a la F-C and K&S is particularly useful, and other orderings are poorly defined (metric dependent)



# 20 years of experience with F-C

Lots of experience in HEP, many find it useful, especially when:

- ★ A model parameter is bounded (mass, cross section, sin/cosine of an angle, etc.); and/or
- ★ Log-likelihood is non-Gaussian (so Wilks's Theorem is inaccurate); multiply connected confidence regions; and/or
- ★ The interesting parameter space is  $>1D$ , where LR ordering a la F-C and K&S is particularly useful, and other orderings are poorly defined (metric dependent)

DM experiments have one or more, so various usage.

**BTW, for data with a “5-sigma discovery”, the F-C “unified approach” reproduces same answer as usual one-tailed test.**

# 20 years of experience with F-C (cont.)

Main foundational (philosophical) issue, already discussed in the F-C paper, is illustrated by Poisson case with non-zero expected background, zero events observed.

See Section 9.1 of arxiv post (violation of Likelihood Principle, common in frequentist statistics).

**Main practical issues:**

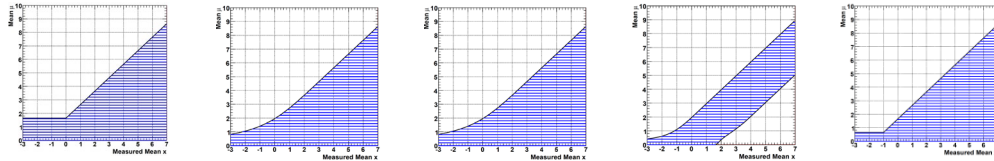
- 1) **Computational time, especially in presence of nuisance parameters.**
- 2) **In common with other frequentist methods, there is no automatic way to “eliminate” nuisance parameters that is always satisfactory. (Section 12)**

**Comparison to other “contenders” in a prototype problem:**

[http://www.physics.ucla.edu/~cousins/stats/cousins\\_bounded\\_gaussian\\_virtual\\_talk\\_12sep2011.pdf](http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf)



# Bayes, Fisher, Neyman, Neutrino Masses, and the LHC



**Bob Cousins**  
**Univ. of California, Los Angeles**

**Virtual Talk**

**12 September 2011**

[http://www.physics.ucla.edu/~cousins/stats/cousins\\_bounded\\_gaussian\\_virtual\\_talk\\_12sep2011.pdf](http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf)

**Back to main theme of this talk:**

**Hypothesis testing of a point null vs a continuous alternative**

# ***Bayesian* Hypothesis Testing (Model Selection)**

***Forget the duality with intervals (!).***

**Typically follows Chapter 5 of book by Harold Jeffreys:  
Bayes's Theorem is applied to the models themselves after  
integrating out *all* parameters, including parameter of interest!**

**Presented too often as “logical” and therefore simple to use,  
with great benefits such as automatic “Ockham’s razor”, etc.**

**In fact, it is *full of subtleties*. E.g., Jeffreys and followers use  
*different priors* for integrating out parameter in model selection  
than for *same* parameter in parameter estimation.**

**(Sections 5, 10, Appendix A)**

# ***Bayesian* Hypothesis Testing (Model Selection)**

Here I will mainly just say: Beware! There are posted/published applications HEP that are silly (*by Bayesian standards*).

A pentaquark example in PRL provoked me to write a Comment: <https://arxiv.org/abs/0807.1330> .

**For testing point null vs continuous alternative, in asymptotic limit of large sample size, your answer (e.g. probability  $H_0$  is true, or an odds ratio called the Bayes Factor) *remains proportional to the prior pdf of parameter of interest.***

This is *totally different* behavior compared to interval estimation, where the effect of prior  $p(\mu)$  typically becomes negligible as sample size increases without bound.

# *Bayesian* hypothesis testing for nested case

$$H_0: \theta = \theta_0 \text{ vs } H_1: \theta \neq \theta_0$$

Let  $\pi_0$  be prior prob for  $H_0$ . Then  $\pi_1 = 1 - \pi_0$  is prior prob for  $H_1$ .

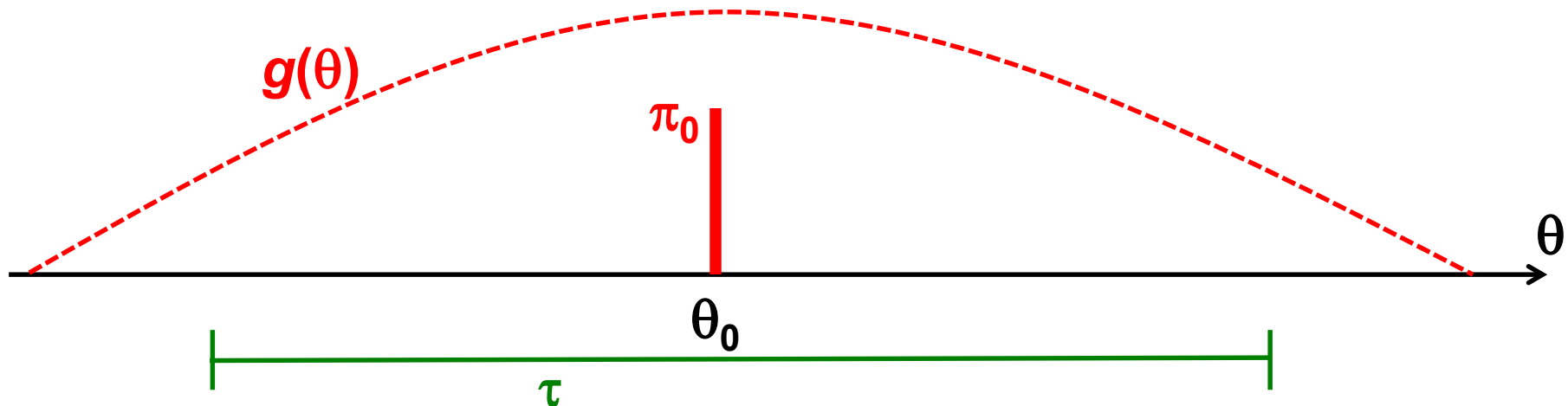
Conditional on  $H_1$  true: prior pdf for  $\theta$ ,  $g(\theta)$ .

$\pi_0$  is like bit of Dirac  $\delta$ -ftn (“probability mass”) at  $\theta = \theta_0$ .

In practice can have a little width:

$\varepsilon_0$  = scale of width of null value(s) of  $\theta$

scale  $\tau$ : extent of prior plausible values in  $g(\theta)$



Gaussian model  $p(x|\theta)$  with rms  $\sigma_{\text{tot}}$ , sampled value  $x_{\text{obs}}$ .

ML Estimate for  $\theta$  is  $\hat{\theta} = x_{\text{obs}}$ .

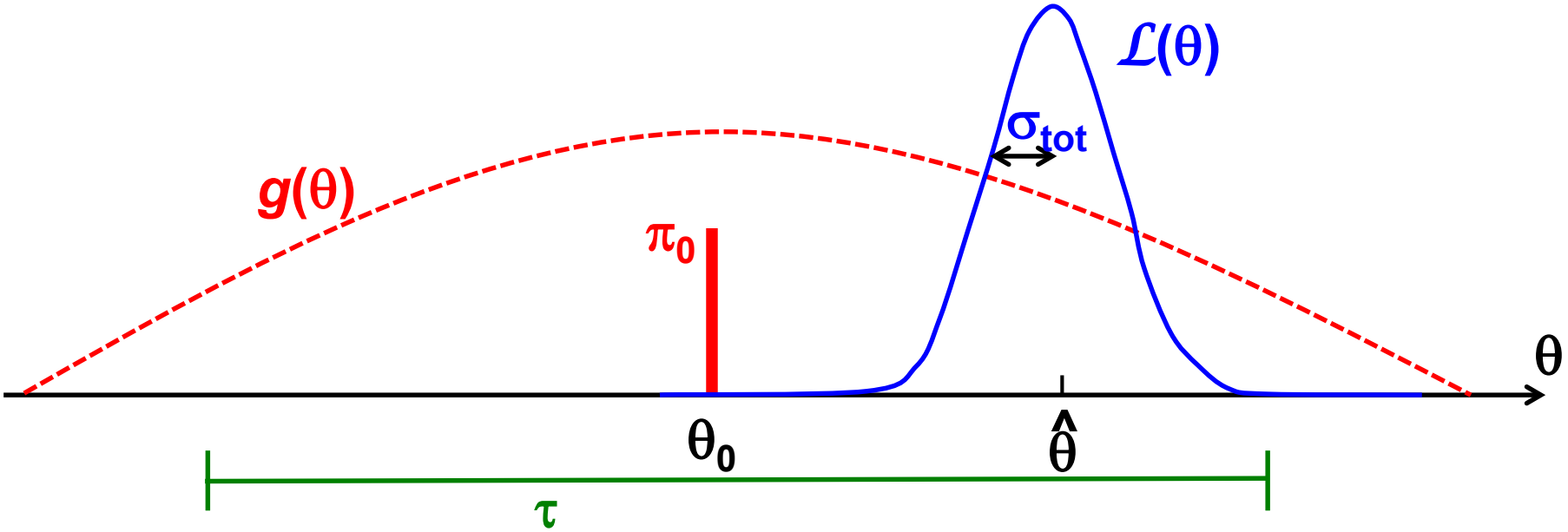
Departure from null in sigma:  $Z = (\hat{\theta} - \theta_0)/\sigma_{\text{tot}}$

Sketch has  $Z \approx 5$ .

Three independent scales: gets interesting when, as shown,

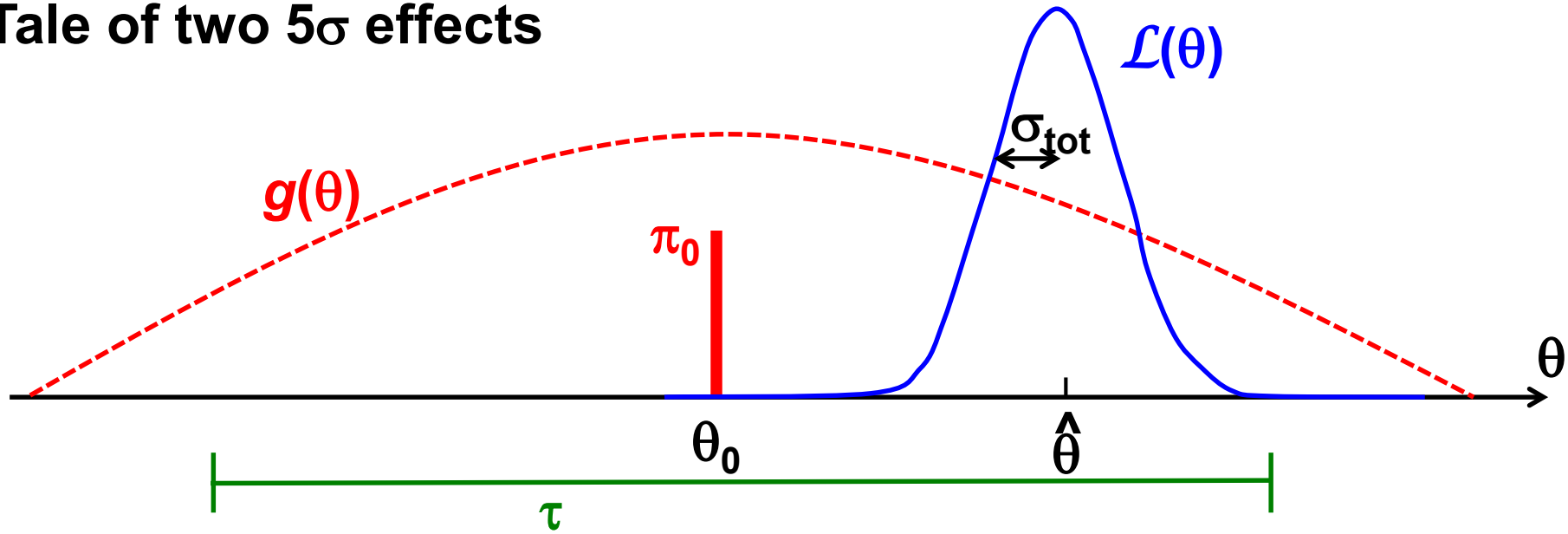
$$\epsilon_0 \ll \sigma_{\text{tot}} \ll \tau.$$

Bayesian posterior prob for  $H_0$ , and Bayes Factor are prop to  $\tau / \sigma_{\text{tot}}$  (1/Ockham factor), independent of  $Z$ !

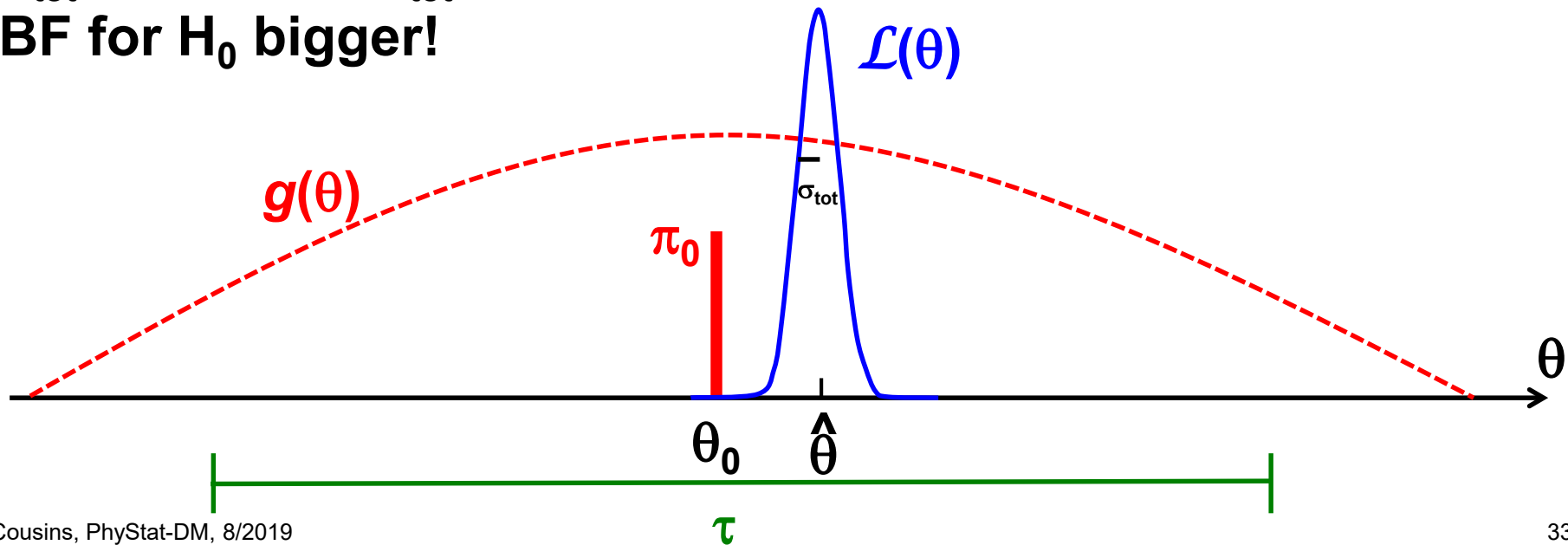




# Tale of two $5\sigma$ effects



$\sigma_{\text{tot}}$  smaller,  $\tau / \sigma_{\text{tot}}$  larger:  
BF for  $H_0$  bigger!



# Jeffreys-Lindley paradox

The fact that BF varies as  $\tau/\sigma_{\text{tot}}$  while Z is fixed is called (at least in extreme cases) the Jeffreys-Lindley paradox.

Also implies that, for experiments obtaining the *same* Z, the Bayesian answers depend on sample size (since  $\sigma_{\text{tot}}$  typically goes as  $1/\sqrt{\text{sample size}}$ ). Very different behavior!

For a review and comparison to p-values in discovery of Higgs boson, see my paper:

“The Jeffreys-Lindley Paradox and Discovery Criteria in High Energy Physics”

(Published in Synthese – long story)

<https://arxiv.org/abs/1310.3791> .

# Priors in Bayesian Hypothesis Testing

For testing  $H_0: \theta = \theta_0$  vs  $H_1: \theta \neq \theta_0$ , improper priors  $g(\theta)$  for  $\theta$  that work fine for estimation (such as Jeffreys priors) become a disaster.

E.g. scale  $\tau$  of  $g(\theta)$  diverges so  $H_0$  always preferred.

Adding cut-off to make prior  $g(\theta)$  proper just gives direct dependence on (arbitrary?)  $\tau$ .

(Contrast with Bayesian point/interval estimation!)

Silly things like prior flat in log of mass as a way to represent “ignorance” are *strongly* informative!

(See any serious Bayesian literature.)

# A side note on priors for “Scales”

As my writeup mentions (following Eadie et al. nearly 50 years earlier) various words including “estimation” have different meanings to statisticians than to physicists. Beware!

Since then, I have realized that a disastrous mistake seems to be made by some physicists regarding the word “scale”.

Recall (or learn about): So-called “objective priors” or “default priors” (often called by the misnomer “noninformative priors”) in Bayesian estimation are based on the *measurement model*, i.e., the measuring apparatus and the protocol (stopping rule, etc.). (Jeffreys’s Rule, Bernardo-Berger Reference Priors, etc.)

E.g., if the measuring apparatus has Gaussian resolution for some parameter (say mass-squared), then the default prior for that parameter (for estimation) is uniform, with all that implies.

# A side note on priors for “Scales” (cont.)

To a statistician, whether or not a parameter is a *scale parameter* again depends on the *measurement model* (!).

Parameter  $\theta$  is a scale parameter if the model  $p(x|\theta)$  has the form:  
 $p(x|\theta) = (1/\theta) f(x/\theta)$ .

From this one can partly derive and partly argue\* that invariance of prior form under change of scale parameter implies the “noninformative prior”  $p(\theta) = 1/\theta$ , i.e., a prior uniform in  $\log(\theta)$ .

\*See pp. 85-87 of Jim Berger’s book on decision theory for subtleties of derivation. See also Jeffreys pp. 120-123, which he abandons later in the book.

# A side note on priors for “Scales” (cont.)

To a statistician, whether or not a parameter is a *scale parameter* again depends on the *measurement model* (!).

Parameter  $\theta$  is a scale parameter if the model  $p(x|\theta)$  has the form:  
 $p(x|\theta) = (1/\theta) f(x/\theta)$ .

From this one can partly derive and partly argue\* that invariance of prior form under change of scale parameter implies the “noninformative prior”  $p(\theta) = 1/\theta$ , i.e., a prior uniform in  $\log(\theta)$ .

To a physicist, a “scale” is a quantity that sets the size of physical quantities like mass, length.

E.g., “What is the DM mass scale?”

So it seems that some physicists make the mistake of saying, “Since mass is a *scale*, I use the prior uniform in  $\log(\text{mass})$ .”  
OOPS! This “scale” is not a statistician’s “scale parameter”!

**See Comment** <https://arxiv.org/abs/1703.04585>

\*See pp. 85-87 of Jim Berger’s book on decision theory for subtleties of derivation. See also Jeffreys pp. 120-123, which he abandons later in the book.

# My advocacy for $>10$ years (Section 16):

**Have in place tools to allow computation of results using a variety of recipes, for problems up to intermediate complexity:**

- Bayesian with analysis of sensitivity to prior
- Profile likelihood ratio (Minuit MINOS)
- Frequentist construction with approximate treatment of nuisance parameters
- Other “favorites” such as LEP’s  $CL_s$  (an HEP invention)

# My advocacy for $>10$ years (Section 16):

**Have in place tools to allow computation of results using a variety of recipes, for problems up to intermediate complexity:**

- Bayesian with analysis of sensitivity to prior
- Profile likelihood ratio (Minuit MINOS)
- Frequentist construction with approximate treatment of nuisance parameters
- Other “favorites” such as LEP’s  $CL_s$  (an HEP invention)

**The community can (and should) then demand that a result shown with one’s preferred method also be shown with the other methods, and with sampling properties studied.**

**When the methods all agree, we are in asymptotic nirvana.**

**When the methods disagree, we are reminded that the results are answers to different questions, and we learn something! E.g.:**

- Bayesian methods can have poor frequentist properties
- Frequentist methods can badly violate likelihood principle



# There is a large literature on frequentist properties of Bayesian (inspired) procedures

**Google on:**

**probability matching priors**

**Welch and Peers 1963**

**calibrated Bayes**

**Bayes non-Bayes compromise**

**prior predictive p-value**

**posterior predictive p-value**

**etc.**

**A nice introductory review is by M.J. Bayarri and J.O. Berger, “The Interplay of Bayesian and Frequentist Analysis”, Statist. Sci. 19 58-80 (2004), doi:10.1214/088342304000000116**

**We should be doing more of this in HEP, in my opinion.**

# Coverage of Bayesian estimation procedures

**Pre-data, Bayesians have the model  $p(x|\mu)$ .**

**Thus, quite apart from imagined repeated experiments (to which they may object) or frequentist definition of probability (to which they may object), a Bayesian can calculate:**

**As a function of  $\mu$ , what is the coverage probability of the credible interval  $[\mu_1, \mu_2]$  that they will report: what is the probability, given the model  $p(x|\mu)$  (with whatever definition of  $p$  they use), that their procedure will lead to an interval in which  $\mu \in [\mu_1, \mu_2]$ .**

***This is a crucial diagnostic to report to the consumer, especially if default priors are used! (Jim B. says reference priors will work.)***

**(Of course, one can also average this coverage over  $\mu$ , weighted by either the prior or the posterior.)**

# Evaluation of properties of Bayesian hypothesis testing procedures

Similarly, quite apart from imagined repeated experiments or frequentist definition of  $p$ , a Bayesian can calculate:

As a function of assumed  $H_0$  and  $H_1$  and any parameters, what is the distribution of the Bayes Factors that they will report: what is the probability, given each model  $p(x|H_i, \mu)$  (with whatever definition of  $p$  they use), that their procedure will obtain various values of the Bayes Factor (or posterior probabilities).

*This is also a crucial diagnostic to report to the consumer, especially if attempts at “noninformative” priors are used!*

*(enlightening for seeing relationship between Bayes Factors and  $p$ -values)*

**Thanks to all (see note), including my  
“sponsor”, U.S. DOE Office of Science**

# BACKUP

# Coverage: The experiments in the ensemble do not have to be the same.

**Neyman pointed this out in his 1937 paper (in which his  $\alpha$  is the modern  $1 - \alpha$ ):**

It is important to notice that for this conclusion to be true, it is not necessary that the problem of estimation should be the same in all the cases. For instance, during a period of time the statistician may deal with a thousand problems of estimation and in each the parameter  $\theta_1$  to be estimated and the probability law of the  $\mathbf{X}$ 's may be different. As far as in each case the functions  $\underline{\theta}(\mathbf{E})$  and  $\bar{\theta}(\mathbf{E})$  are properly calculated and correspond to the same value of  $\alpha$ , his steps (a), (b), and (c), though different in details of sampling and arithmetic, will have this in common—the probability of their resulting in a correct statement will be the same,  $\alpha$ . Hence the frequency of actually correct statements will approach  $\alpha$ .

# Above is all “pre-data” characterization of the test

## How to characterize *post-data*?

### p-values and Z-values

In N-P theory,  $\alpha$  is *specified in advance*.

Suppose after obtaining data, you notice that with  $\alpha=0.05$  previously specified, you reject  $H_0$ , but with  $\alpha=0.01$  previously specified, you accept  $H_0$ .

In fact, you determine that with the data set in hand,  $H_0$  would be rejected for  $\alpha \geq 0.023$ . This interesting value has a name:

***After data are obtained, the p-value is the smallest value of  $\alpha$  for which  $H_0$  would be rejected, had it been specified in advance.***

**This is numerically (if not philosophically) the same as definition used e.g. by Fisher and often taught: “p-value is probability under  $H_0$  of obtaining  $x$  as extreme or more extreme than observed  $x_0$ .”**

# Interpreting p-values and Z-values

**It is crucial to realize that that value of  $\alpha$  (0.023 in the example) was typically *not* specified in advance, so p-values do *not* correspond to Type I error probs of experiments reporting them.**

**In HEP, p-value typically converted to Z-value (unfortunately commonly called “the significance S”), equivalent number of Gaussian sigma.\***

**E.g., for one-tailed test,  $p = 2.87\text{E-}7$  is  $Z = 5$ .**



# Interpreting p-values and Z-values (cont.)

Interpretation of p-values (and hence Z-values) is a long, contentious story – beware!

Widely bashed. I give some reasons why in <https://arxiv.org/abs/1807.05996> .

**I defend their use in HEP. See <https://arxiv.org/abs/1310.3791>.)**

Whatever they are, p-values are *not* the probability that  $H_0$  is true!

- They are calculated *assuming that  $H_0$  is true*, so they can hardly tell you the probability that  $H_0$  is true!
- Calculation of “probability that  $H_0$  is true” requires prior(s)!

**Please help educate press officers and journalists!  
(and physicists) !**

# Whatever you call non-subjective priors, they do *not* represent ignorance!

**Dennis V. Lindley** *Stat. Sci* 5 85 (1990), “the mistake is to think of them [Jeffreys priors or Bernardo/Berger’s reference priors] as representing ignorance”

This Lindley quote is emphasized by Christian Robert, *The Bayesian Choice*, (2007) p. 29.

**Jose Bernardo**: “[With non-subjective priors,] The contribution of the data in constructing the posterior of interest should be “dominant”. Note that this does not mean that a non-subjective prior is a mathematical description of “ignorance”. Any prior reflects some form of knowledge.”

Nonetheless, Berger (1985, p. 90) argues that Bayesian analysis with noninformative priors (older name for objective priors) such as Jeffreys and Bernardo/Berger “is the single most powerful method of statistical analysis, in the sense of being the *ad hoc* method most likely to yield a sensible answer for a given investment of effort”.

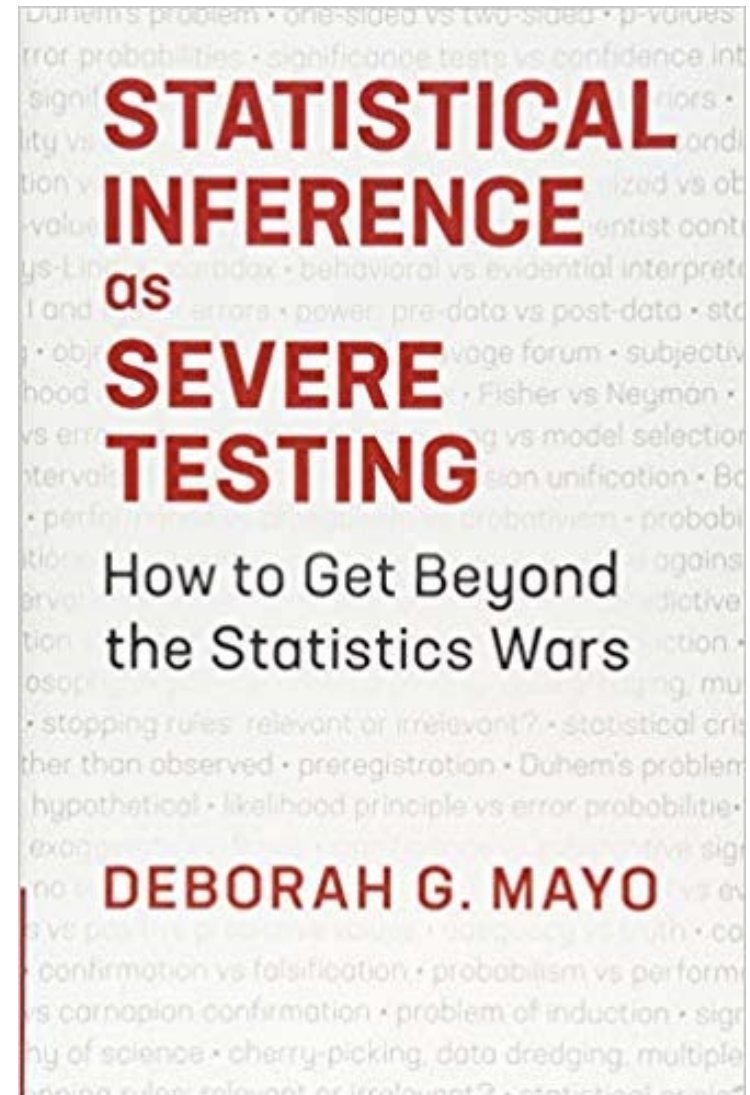
# Recent book exploring Bayesian-frequentist divide

Much interesting history and up to-date discussion of both theory and practice, including, for example, internal debates among Bayesians.

Very well-referenced.

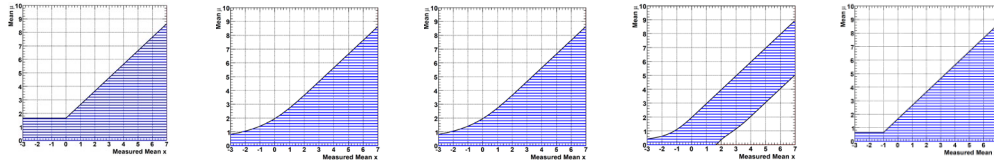
Mayo has long advocated “error statistics”, and in particular the concept of how *severely* a hypothesis has been tested in a test that “passes”.

(I plan to think more about how this maps on to what we do in HEP.)





# Bayes, Fisher, Neyman, Neutrino Masses, and the LHC



**Bob Cousins**  
**Univ. of California, Los Angeles**

**Virtual Talk**

**12 September 2011**

[http://www.physics.ucla.edu/~cousins/stats/cousins\\_bounded\\_gaussian\\_virtual\\_talk\\_12sep2011.pdf](http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf)

# Negatively Biased Relevant Subsets Induced by the Most-Powerful One-Sided Upper Confidence Limits for a Bounded Physical Parameter

Robert D. Cousins\*

Department of Physics and Astronomy  
University of California, Los Angeles, California 90095, USA

September 9, 2011

## Abstract

Suppose an observable  $x$  is the measured value (negative or non-negative) of a “true mean”  $\mu$  (physically *non-negative*) in an experiment with a Gaussian resolution function with known fixed rms deviation  $\sigma$ . The most powerful one-sided upper confidence limit at 95% confidence level (C.L.) is  $\mu_{\text{UL}} = x + 1.64\sigma$ , which I refer to as the “original diagonal line”. Perceived problems in HEP with small or non-physical upper limits for  $x < 0$  historically led, for example, to substitution of  $\max(0, x)$  for  $x$ , and eventually to abandonment in the Particle Data Group’s Review of Particle Physics of this diagonal line relationship between  $\mu_{\text{UL}}$  and  $x$ . Recently Cowan, Cranmer, Gross, and Vitells (CCGV) have advocated a concept of “power constraint” that when applied to this problem yields variants of diagonal line, including  $\mu_{\text{UL}} = \max(-1, x) + 1.64\sigma$ . Thus it is timely to consider again what is problematic about the original diagonal line, and whether or not modifications cure these defects. In a 2002 Comment, statistician Leon Jay Gleiser pointed to the literature on *recognizable* and *relevant subsets*. For upper limits given by the original diagonal line, the sample space for  $x$  has recognizable relevant subsets in which the quoted 95% C.L. is *known* to be negatively biased (anti-conservative) by a finite amount for *all* values of  $\mu$ . This issue is at the heart of a dispute between Jerzy Neyman and Sir Ronald Fisher over fifty years ago, the crux of which is the relevance of pre-data coverage probabilities when making post-data inferences. The literature describes illuminating connections to Bayesian statistics as well. Methods such as that advocated by CCGV have 100% unconditional coverage for certain values of  $\mu$  and hence formally evade the traditional criteria for negatively biased relevant subsets; I argue that concerns remain. Comparison with frequentist intervals advocated by Feldman and Cousins also sheds light on the issues.

arXiv:1109.2023v1 [physics.data-an] 9 Sep 2011