

Reflections on the 20th year of F-C: Hypothesis testing of a point null vs a continuous alternative

Bob Cousins

Univ. of California, Los Angeles

PhyStat-nu at CERN

January 25, 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.

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”

μ denotes parameter(s) (sometimes we use θ)

$p(x|\mu)$ 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_{obs} , for what values of μ is x_{obs} an “extreme” value of x ?

Include in the confidence interval $[\mu_1, \mu_2]$ those values of μ for which x_{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_{obs} , for what values of μ is x_{obs} an “extreme” value of x ?

Include in the confidence interval $[\mu_1, \mu_2]$ those values of μ for which x_{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 μ : *high rank means not extreme*.

2) Need also 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 μ implies higher average x):

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

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 μ for which x_{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 μ_t be the unknown true value of μ . 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 μ_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 μ_t be the unknown true value of μ . 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 μ_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, μ_t is *fixed* and unknown. The endpoints μ_1, μ_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* $L(x|\mu) / L(x|\mu_{\text{best fit}})$, advocated by F-C.

Unified approach to the classical statistical analysis of small signals

Gary J. Feldman^{*}

Department of Physics, Harvard University, Cambridge, Massachusetts 02138

Robert D. Cousins[†]

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

Phys. Rev. D57 3873 (1998):

Ordering applies (in principle) for arbitrary dimension of x , μ .

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 μ in H_1 , and H_0 corresponds to a particular parameter value μ_0 (e.g., zero, SM prediction, or ∞).

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

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

Common examples:

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

2) Higgs boson $\rightarrow \gamma\gamma$ before observation, signal strength μ :
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) ν mixing angle θ_{23} *before 1998*: null $\theta_{23} = 0$, alternative $\theta_{23} \neq 0$

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

2a) CP violation phase δ before it is observed:

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

2b) After two-point null for δ is rejected: maybe a theorist has a “predicted” value of δ 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 α (smallish number).

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

α = 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 β]

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_{obs} , suppose the 90% C.L. confidence interval for μ is $[\mu_1, \mu_2]$.

This contains all values of μ for which observed x_{obs} is ranked in the *least extreme 90%* 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_{obs} , suppose the 90% C.L. confidence interval for μ is $[\mu_1, \mu_2]$.

This contains all values of μ for which observed x_{obs} is ranked in the *least extreme* 90% 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 $\alpha = 10\%$. We reject H_0 if x_{obs} is ranked in the *most extreme* 10% 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_{obs} , suppose the 90% C.L. confidence interval for μ is $[\mu_1, \mu_2]$.

This contains all values of μ for which observed x_{obs} is ranked in the *least extreme* 90% 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 $\alpha = 10\%$. We reject H_0 if x_{obs} is ranked in the *most extreme* 10% of x according to $p(x|\mu)$ and the ordering principle in use.

Comparing the two procedures, we see:

Reject H_0 at $\alpha=10\%$ iff μ_0 is *not* in 90% 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 (θ_r, θ_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 θ_s , when H_0 holds,¹ giving the conditional maximum of the LF

$$L(x|\theta_{r0}, \hat{\theta}_s). \quad (22.3)$$

$\hat{\theta}_s$ 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 ratio²

$$l = \frac{L(x|\theta_{r0}, \hat{\theta}_s)}{L(x|\hat{\theta}_r, \hat{\theta}_s)}. \quad (22.4)$$

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_α is determined from the distribution $g(l)$ of l to give a size- α 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 θ to $\tau(\theta)$, the ML estimator of $\tau(\theta)$ being $\tau(\hat{\theta})$ – cf. 18.3. Thus the LR statistic is invariant under reparametrization.

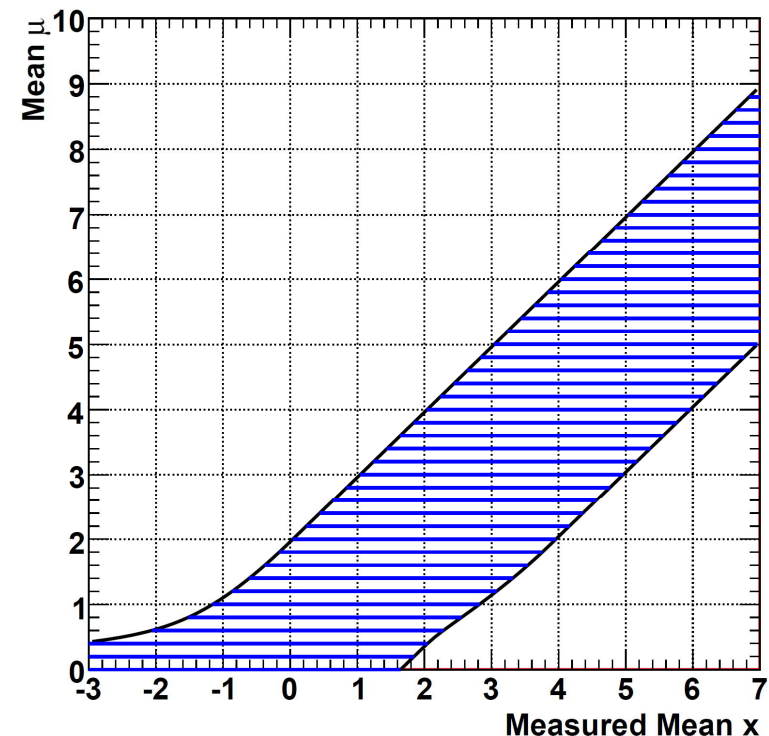
Famous confusion re Gaussian $p(x|\mu)$ where $\mu \geq 0$

It is *crucial* to distinguish between the data x , which *can* be negative (no problem), and a parameter μ such as mass, 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 μ !!! It's in the *model* (and hence in $L(\mu)$).

The confusion is encouraged since we often refer to x as the “measured value of μ ”, 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”**

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
- ★ 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)

Neutrino community gets three gold stars, so major user!

(And in fact F-C were working on the NOMAD neutrino experiment at CERN in 1998.)

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

***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 lots of data, 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 typically becomes negligible.

Bayesian hypothesis testing for nested case

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

Let π_0 be prior prob for H_0 . Then $\pi_1 = 1 - \pi_0$ is prior prob for H_1 .

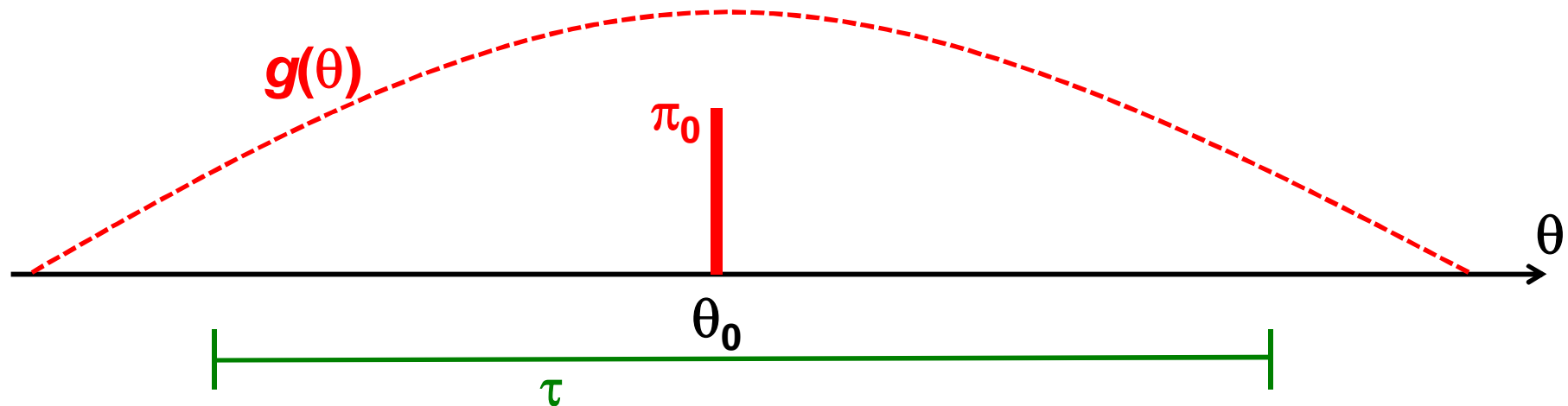
Conditional on H_1 true: prior pdf for θ , $g(\theta)$.

π_0 is like bit of Dirac δ -ftn (“probability mass”) at $\theta = \theta_0$.

In practice can have a little width:

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

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



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

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

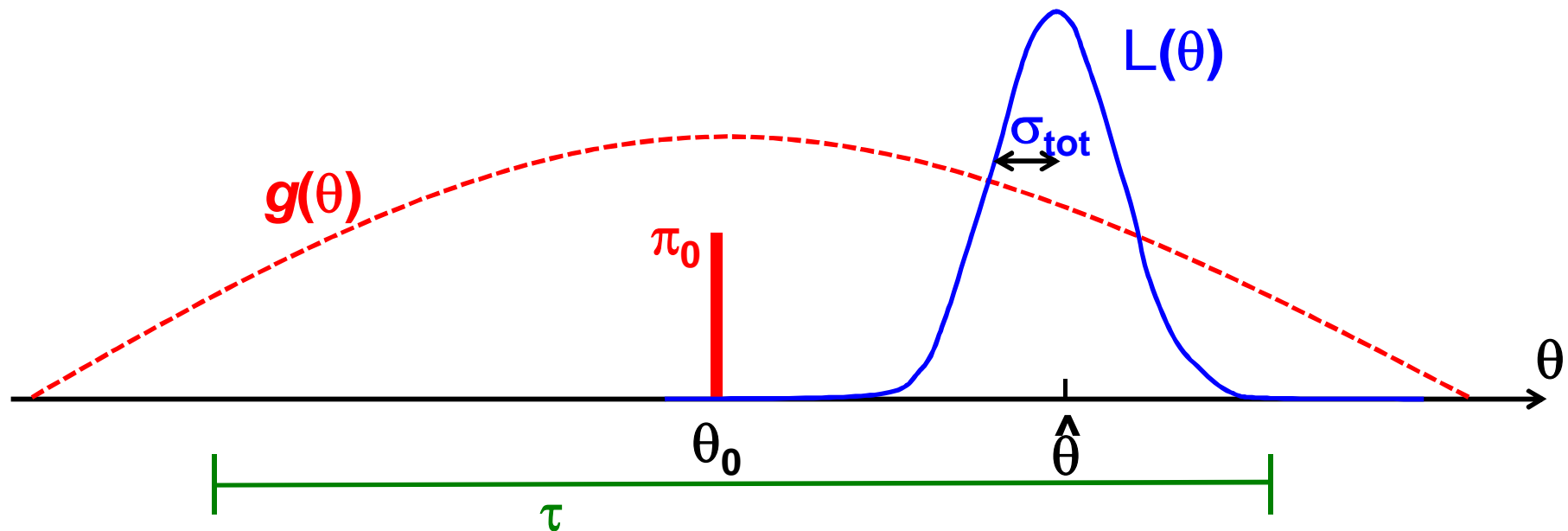
Departure from null in sigma: $Z = (\hat{\theta} - \theta_0)/\sigma_{\text{tot}}$

Sketch has $Z \approx 5$.

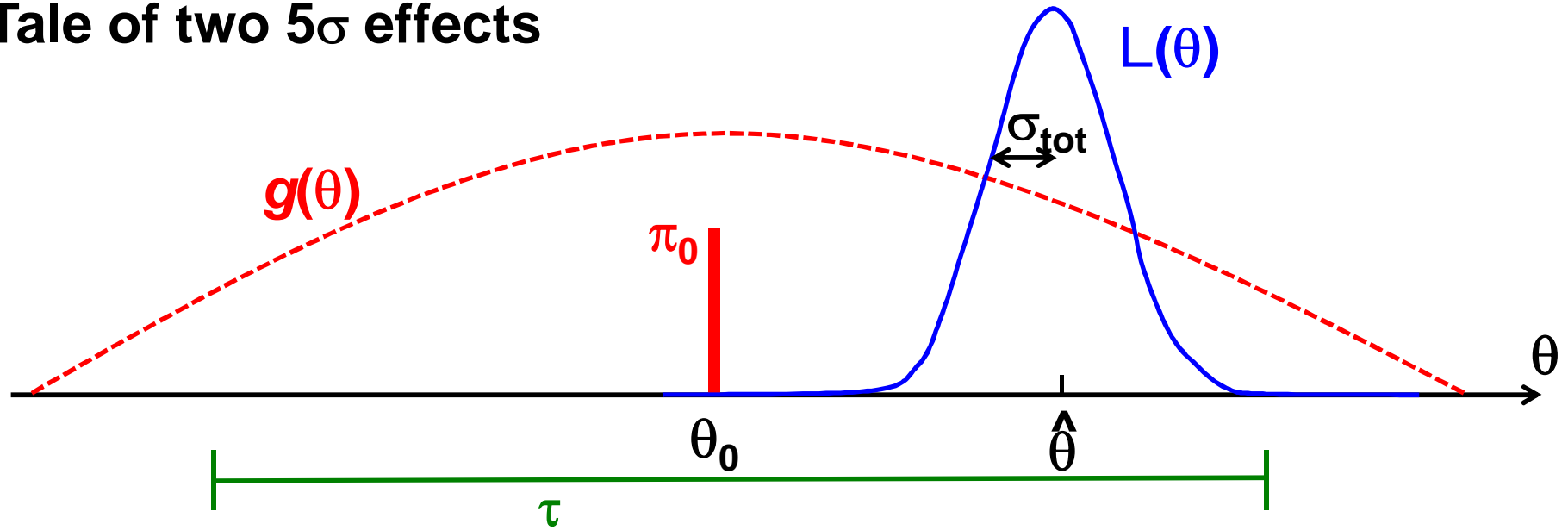
Three independent scales: gets interesting when, as shown,

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

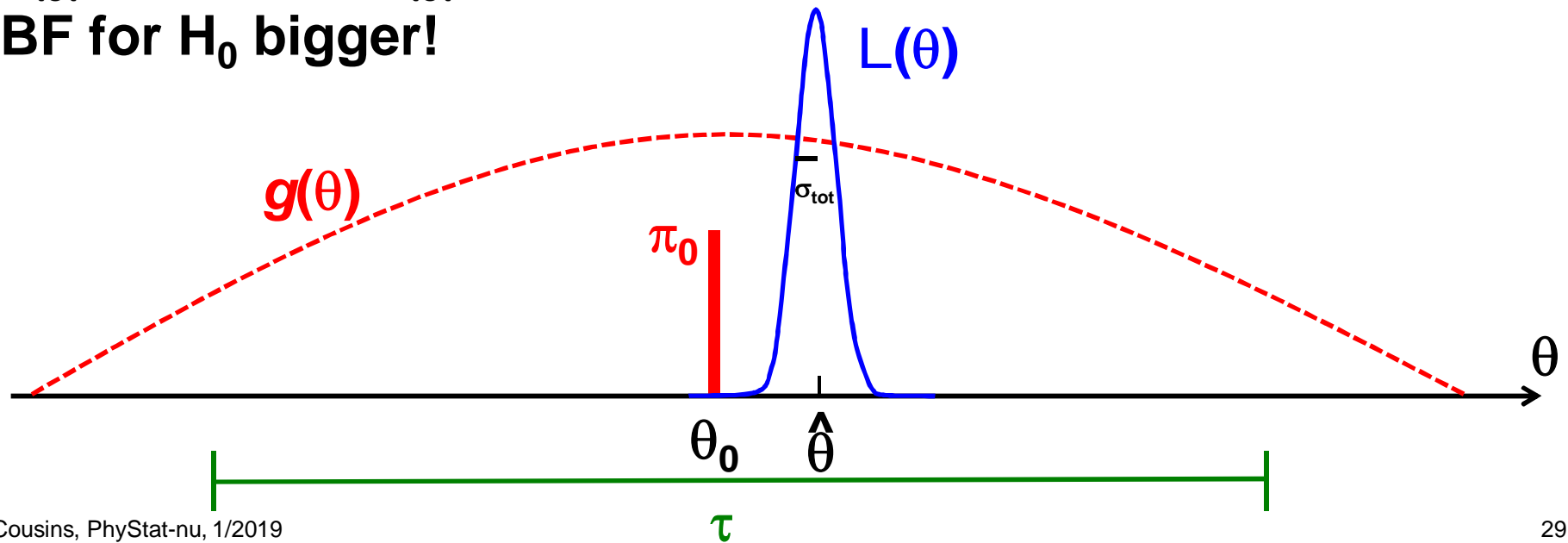
Bayesian posterior prob for H_0 , and Bayes Factor are prop to τ/σ_{tot} (1/Ockham factor), independent of Z !



Tale of two 5σ effects



σ_{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 σ_{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 for θ that work fine for estimation become a disaster.

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

Adding cut-off to make prior $g(\theta)$ proper just gives direct dependence on (arbitrary?) τ . (Contrast with estimation!)

Silly things like prior flat in log of mass as a way to represent “ignorance” are *strongly* informative!
(See any serious Bayesian literature.)

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 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

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 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 or frequentist definition of p (to which they may object), they can calculate:

As a function of μ , 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 μ , 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 (to which they may object), one 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 probs).

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)

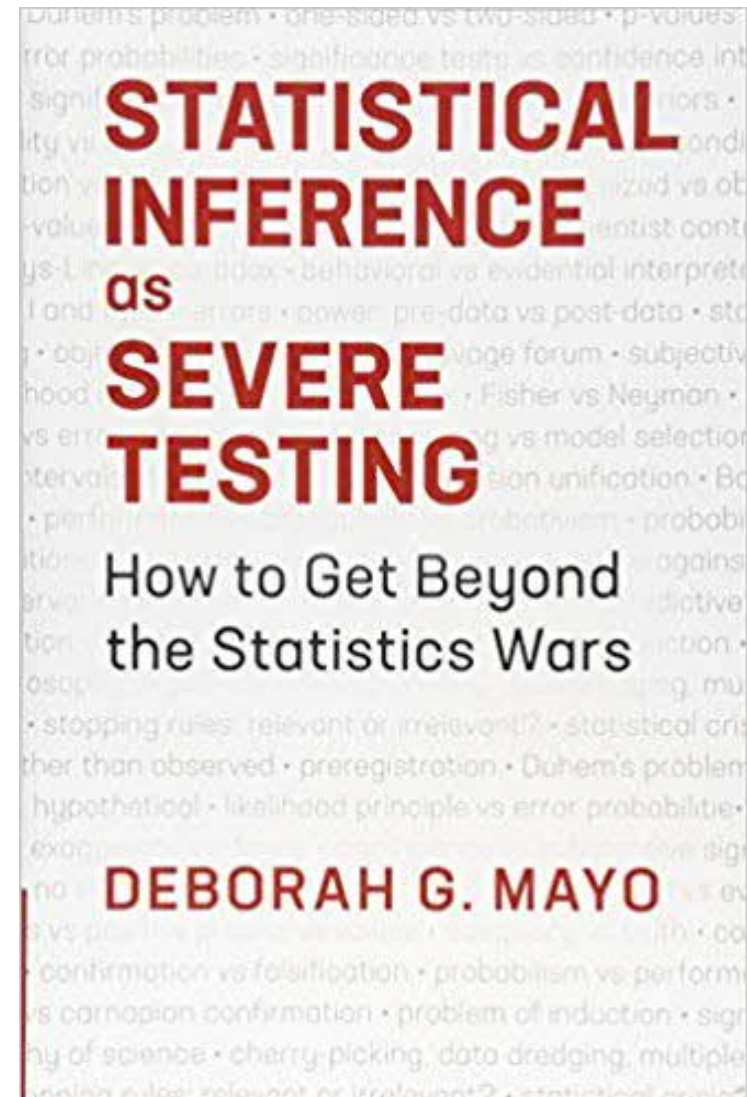
Recent book exploring Bayesian-frequentist divide

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.)



**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 α 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 θ_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 α , 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, α . Hence the frequency of actually correct statements will approach α .

Above is all “pre-data” characterization of the test
How to characterize *post-data*?
p-values and Z-values

In N-P theory, α 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 α 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 α (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 later.

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”

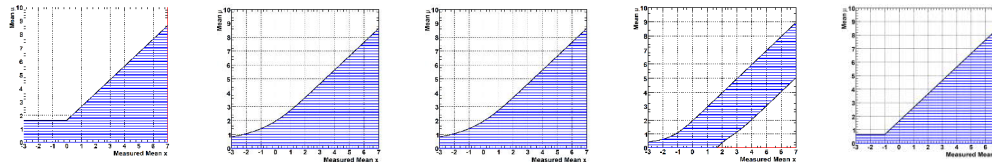
This Lindley quote is emphasized by Christian Robert, *The Bayesian Choice*, (2007) p. 29.

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 Barnardo/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”.



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