

“Statistics in Theory”: Prelude to Statistics in Practice

Robert (Bob) Cousins

Univ. of California, Los Angeles

<https://www.physics.ucla.edu/~cousins/homepage/bio.html>

<http://www.physics.ucla.edu/~cousins/homepage/research.html>

**8th African School of
Fundamental Physics and Applications
Marrakesh, July 8-11, 2024**

**For more complete writeup with references, see arXiv post
<https://arxiv.org/abs/1807.05996>**

Preface

Many of us teach advanced “data analysis” courses that last an entire academic term. How to condense?

Here I concentrate on the “theoretical” underpinnings:

What you must know in order to choose appropriate methods.

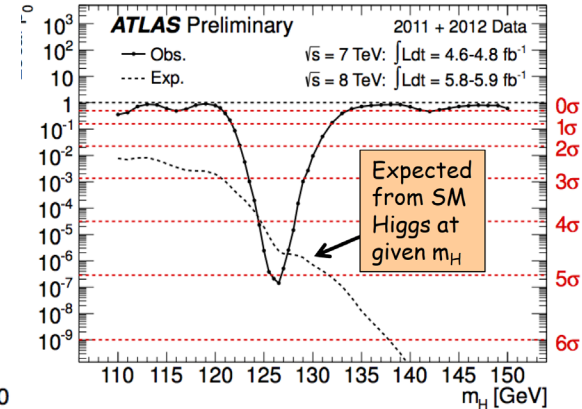
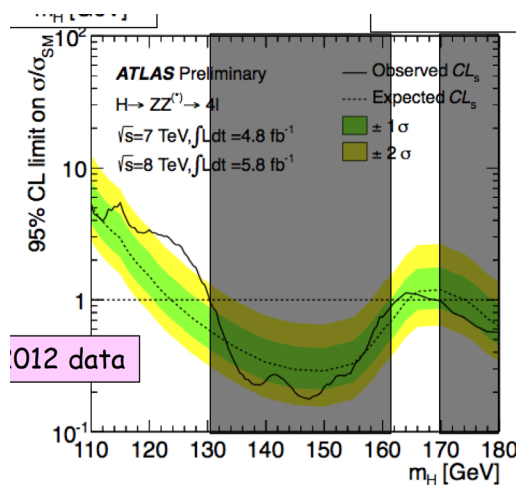
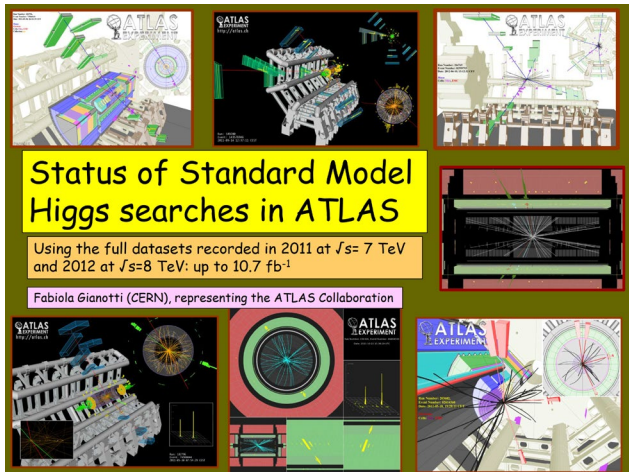
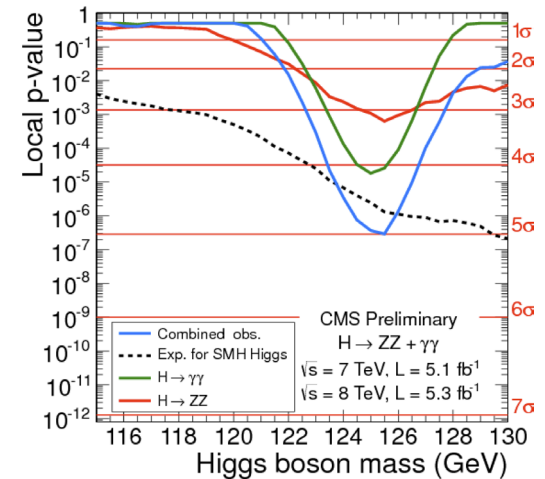
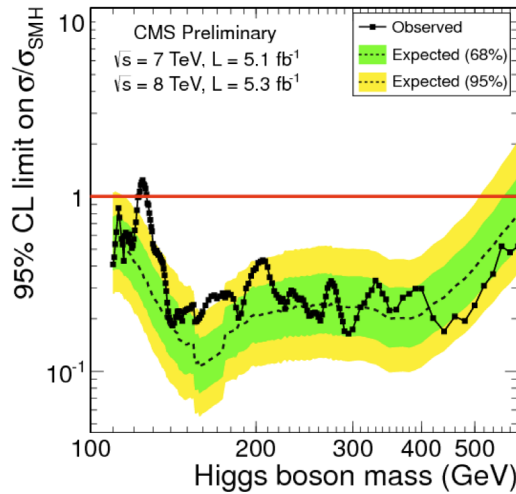
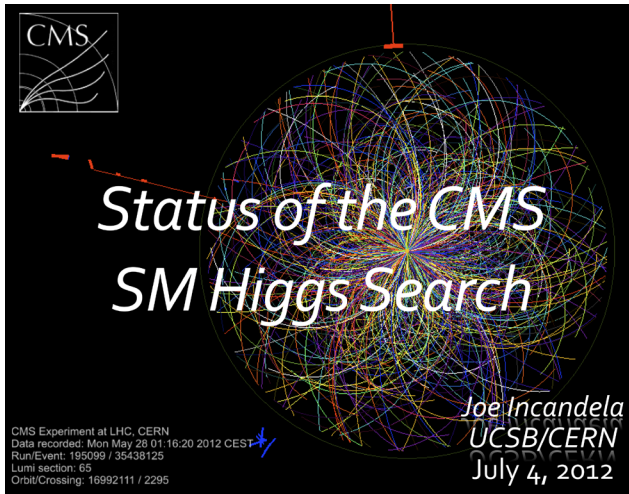
My hope is that by studying these slides you will learn to avoid common pitfalls (and even silly statements) that can trip up professionals in the field.

This is a dense talk – you will not pick it all up in real time. It should however be extremely useful to you to study this talk, referring to the references.

I initially focus on definitions and the Bayesian approach. That helps to understand what the frequentist approach (common in HEP) is *not* ! Then, compare and contrast.

I will end on Thursday with an example of statistics in practice:

Statistics in practice on July 4, 2012: Imagine being in the audience for the talks on the discovery of the Higgs boson



A goal this week is to help you understand plots like these ...not just what is plotted, but also deeper issues.

Why Foundations Matter

In the “final analysis”, we often make approximations, take a pragmatic approach, or follow a convention.

To inform such actions, it is important to understand some foundational aspects of statistical inference.

In Quantum Mechanics, we are used to the fact that for all of our practical work, one’s philosophical interpretation (e.g., of collapse of the wave function) does not matter.

Why Foundations Matter (cont.)

In statistical inference, however, *foundational differences result in different answers*: one cannot ignore them!

The professional statistics community went through the topics of many of our discussions starting in the 1920's, and revisited them in the resurgence of Bayesian methods in recent decades.

I will attempt to summarize some of the things we should understand from that debate.

Most importantly: understand both approaches!

Machine learning inherits these issues

The increasingly common uses of machine learning (boosted decision trees, deep neural nets, etc) in physics do not make the foundational issues go away. If anything, ML adds more issues while being amazingly useful.

E.g., the output of a neural net is a “statistic” (function of the data) that is conceptually on a similar level (and can even correspond to) more traditional statistics such as likelihood ratios, which we will discuss.

The language in the machine learning community is sometimes different from that in the statistical inference community, while concepts can be fundamentally the same.

Mastering the traditional foundations of statistical inference, and mapping the language of ML it, can thus lead to useful inquiries regarding what is the underlying “philosophy” of the “machine”, etc.

Definitions are Important

As in physics, much confusion can be avoided by being precise about definitions, and much confusion can be generated by being imprecise, or by assuming every-day definitions in a technical context.

You have learned in physics to see confusion in the statement,

“I did a lot of *work* today by carrying this big stone around the building and then putting it back in its original place.”

You should see just as much confusion in these two statements:

- 1) “The *confidence* level tells you how much confidence one has that the true value is in the confidence interval.”
- 2) “A *noninformative* prior probability density does not insert any information.”

Example adapted from Eadie et al. (James06, p. 2)

Physicists say... when Statisticians say:

Determine

Estimate

Gaussian

**Breit-Wigner,
Lorentzian**

Estimate

(Informed) Guess

Normal

Cauchy

Key tasks: Important to distinguish!

- *Point estimation*: **what single “measured” value of a parameter do you report?**

Key tasks: Important to distinguish!

- *Point estimation*: what single “measured” value of a parameter do you report?
- *Interval estimation*: what interval (giving a measure of uncertainty of the parameter inference) do you report?

Key tasks: Important to distinguish!

- **Point estimation:** what single “measured” value of a parameter do you report?
- **Interval estimation:** what interval (giving a measure of uncertainty of the parameter inference) do you report?
- **Hypothesis testing:** Many special cases:
 - a) A given functional form (“model”) vs another functional form. Also known as “model selection”.
 - b) A single value of a parameter (say 0 or 1) vs all other values
 - c) Goodness of Fit: A given functional form against all other (unspecified) functional forms (aka “model checking”)

Key tasks: Important to distinguish!

- **Point estimation:** what single “measured” value of a parameter do you report?
- **Interval estimation:** what interval (giving a measure of uncertainty of the parameter inference) do you report?
- **Hypothesis testing:** Many special cases:
 - a) A given functional form (“model”) vs another functional form. Also known as “model selection”.
 - b) A single value of a parameter (say 0 or 1) vs all other values
 - c) Goodness of Fit: A given functional form against all other (unspecified) functional forms (aka “model checking”)
- **Decision making:** What action should I take (tell no one, issue press release, propose new experiment, ...) based on the observed data? Rarely done formally in HEP, but important to understand outline of formal theory, to avoid confusion with inference and to inform informal application.

Key tasks: Important to distinguish! (cont.)

In frequentist statistics, the above hypothesis testing case,

(b) A single value of a parameter (say 0 or 1) vs all other values,

maps identically onto interval estimation.

This is called the duality of “inversion of a hypothesis test to get confidence interval”, and vice versa. I just mention it now but discuss it in more detail later.

In contrast, in Bayesian statistics, testing case (b) is an especially controversial form of case (a) model selection.

The model with fixed value of parameter is lower-dimensional in parameter space than the model with parameter not fixed.

Again, I just mention this now to foreshadow a very deep issue, where frequentist and Bayesian methods do not converge in the limit of large data sets.

Comments on key tasks

- ***Point estimation:*** Long history. In the end it is not clear what the criteria are for “best” estimator. Decision Theory can be used to specify criteria and choose among point estimators. IMO, point estimation is not a key issue; typically the *maximum likelihood* (ML) point estimator serves our needs in HEP.
- ***Interval estimation:*** In HEP, it is fairly mandatory that there is a *confidence level* that gives *frequentist coverage* probability of a method, even if it’s a Bayesian-inspired recipe. For many problems in HEP, there is reasonable hope of approximate reconciliation between Bayesian and frequentist methods,.

Point estimation and interval estimation can be approached consistently by insisting that the interval estimate contain the point estimate; in that case one constructs the point estimate by taking the limit of interval estimates as intervals get smaller (limit of confidence level going to zero).

Comments on key tasks (cont.)

- ***Hypothesis testing:***
 - Bayesian methods attempt to calculate probability that a hypothesis is true.
 - Frequentist methods use p-values (often bashed).

Can be dramatic differences between frequentist and Bayesian hypothesis testing methods, even asymptotically. Beware!

See my paper on Jeffreys-Lindley paradox,

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

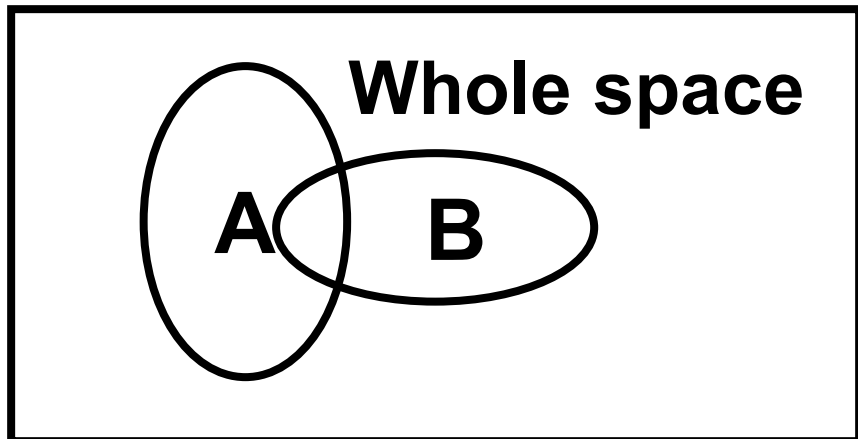
Much more in subsequent slides.

“Probability”

- **Abstract mathematical probability P can be defined in terms of sets and axioms that P obeys. Conditional probabilities are related by (next slide) **Bayes' Theorem (or “Bayes' Rule”)**,**

$$P(B|A) = P(A|B) P(B) / P(A).$$

P, Conditional P, and “Derivation” of Bayes’ Theorem



$$P(A) = \frac{\text{Area of } A}{\text{Area of Whole space}}$$

$$P(B) = \frac{\text{Area of } B}{\text{Area of Whole space}}$$

$$P(A|B) = \frac{\text{Area of } A \cap B}{\text{Area of } B}$$

$$P(B|A) = \frac{\text{Area of } A \cap B}{\text{Area of } A}$$

$$P(A \cap B) = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}}$$

$$P(A) \times P(B|A) = \frac{\text{Area of } A}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } B} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

$$P(B) \times P(A|B) = \frac{\text{Area of } B}{\text{Area of Whole space}} \times \frac{\text{Area of } A \cap B}{\text{Area of } A} = \frac{\text{Area of } A \cap B}{\text{Area of Whole space}} = P(A \cap B)$$

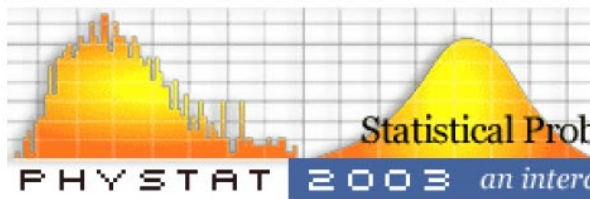
$$\Rightarrow P(B|A) = P(A|B) \times P(B) / P(A)$$

“Probability” (cont.)

Two established* incarnations of mathematical P are:

- 1) **Frequentist P** : limiting frequency in ensemble of imagined repeated samples (as usually taught in Q.M.).
 P (constant of nature) and P (SUSY is true) do not exist (in a useful way) for this definition of P (at least in 1 universe).
 - 2) **(Subjective) Bayesian P** : subjective (personalistic) degree of belief. (de Finetti, Savage)
 **P (constant of nature) and P (SUSY is true) exist for You.
Shown to be basis for coherent personal decision-making.**
- ***It is important to be able to work with either definition of P , and to know which one you are using!***

*Of course they are still argued about, but to less practical effect, I think.



Bayesians, Frequentists, and Physicists

Bradley Efron

*Department of Statistics and Department of Health Research and Policy,
Stanford University, Stanford, CA 94305, USA*

“Bayes' rule is satisfying, convincing, and fun to use. But using Bayes' rule does not make one a Bayesian; *always* using it does, and that's where difficulties begin.”

I'll give a simple example for each definition of P. One of the sillier things one sometimes sees in HEP is the use of a frequentist example of Bayes' Theorem as a foundational argument for “Bayesian” statistics.

Aside: What is the “Whole Space”?

For probabilities to be well-defined, the “whole space” needs to be defined. Can be hard for both frequentists and Bayesians!

Thus the “whole space” itself is more properly thought of as a conditional space, conditional on the assumptions going into the model (Poisson process, whether or not total number of events was fixed, etc.).

Furthermore, it is widely accepted that restricting the “whole space” to a relevant (“conditional”) subspace can sometimes improve the quality of statistical inference. The important topic of such “conditioning” in frequentist inference will be discussed in detail later.

I will not clutter the notation with explicit mention of the assumptions defining the “whole space”, but some prefer to do so – in any case, it is important to keep them in mind.

Example of Bayes' Theorem Using Frequentist P

In high-energy collisions, dedicated algorithms are used to detect the presence of clusters (“jets”) of particles containing bottom quarks, i.e., to “tag b jets”.

A b-tagging algorithm is developed and one measures:

$P(\text{btag} \mid \text{b-jet})$, i.e., efficiency for tagging b jets

$P(\text{btag} \mid \text{not a b-jet})$, i.e., efficiency for background

$P(\text{no btag} \mid \text{b-jet}) = 1 - P(\text{btag} \mid \text{b-jet})$,

$P(\text{no btag} \mid \text{not a b-jet}) = 1 - P(\text{btag} \mid \text{not a b-jet})$

Question: Given a selection of jets tagged as b-jets, what fraction of them is truly b-jets? I.e., what is $P(\text{b-jet} \mid \text{btag})$?

Example of Bayes' Theorem Using Frequentist P

In high-energy collisions, dedicated algorithms are used to detect the presence of clusters (“jets”) of particles containing bottom quarks, i.e., to “tag b jets”.

A b-tagging algorithm is developed and one measures:

$P(\text{btag} \mid \text{b-jet})$, i.e., efficiency for tagging b jets

$P(\text{btag} \mid \text{not a b-jet})$, i.e., efficiency for background

$P(\text{no btag} \mid \text{b-jet}) = 1 - P(\text{btag} \mid \text{b-jet})$,

$P(\text{no btag} \mid \text{not a b-jet}) = 1 - P(\text{btag} \mid \text{not a b-jet})$

Question: Given a selection of jets tagged as b-jets, what fraction of them is truly b-jets? I.e., what is $P(\text{b-jet} \mid \text{btag})$?

Answer: *Cannot be determined from the given information!*

Need in addition: $P(\text{b-jet})$, the true fraction of *all* jets that are b-jets. Then Bayes' Thm inverts the conditionality:

$$P(\text{b-jet} \mid \text{btag}) \propto P(\text{btag} \mid \text{b-jet}) P(\text{b-jet})$$

Example of Bayes' Theorem Using Frequentist P (cont.)

In HEP, as noted,

$P(\text{btag} \mid \text{b-jet})$ is called the *efficiency* for tagging b's.

Meanwhile

$P(\text{b-jet} \mid \text{btag})$ is often called the *purity* of a sample of b-tagged jets.

As this is a pretty “easy” distinction, it is helpful to keep it in mind when one encounters cases where it is perhaps tempting to make the logical error of equating $P(A|B)$ and $P(B|A)$.

Note: Looking ahead, when we talk about frequentist hypothesis testing later in the lectures, we will mention names for analogous probabilities in other fields of science.

E.g., in medicine, $P(\text{Covid test is positive} \mid \text{patient has Covid})$ is called the *sensitivity* of the Covid test, or (unfortunately IMO), the *true positive rate*.

Example of Bayes' Theorem Using Bayesian P

In a *background-free* experiment, a theorist uses a “model” to predict a signal with Poisson mean of 3 events. From Poisson formula we know

$$P(0 \text{ events} \mid \text{model true}) = 3^0 e^{-3} / 0! = 0.05$$

$$P(0 \text{ events} \mid \text{model false}) = 1.0$$

$$P(>0 \text{ events} \mid \text{model true}) = 0.95$$

$$P(>0 \text{ events} \mid \text{model false}) = 0.0$$

The experiment is performed and *zero events are observed*.

Question: Given the result of the expt, what is the probability that the model is true? I.e., What is $P(\text{model true} \mid 0 \text{ events})$?

Example of Bayes' Theorem Using Bayesian P

In a *background-free* experiment, a theorist uses a “model” to predict a signal with Poisson mean of 3 events. From Poisson formula we know

$$P(0 \text{ events} \mid \text{model true}) = 3^0 e^{-3} / 0! = 0.05$$

$$P(0 \text{ events} \mid \text{model false}) = 1.0$$

$$P(>0 \text{ events} \mid \text{model true}) = 0.95$$

$$P(>0 \text{ events} \mid \text{model false}) = 0.0$$

The experiment is performed and *zero events are observed*.

Question: Given the result of the expt, what is the probability that the model is true? I.e., What is $P(\text{model true} \mid 0 \text{ events})$?

Answer: *Cannot be determined from the given information!*

Need in addition: $P(\text{model true})$, the *degree of belief* in the model *prior* to the experiment. Then Bayes' Thm inverts the conditionality:

$$P(\text{model true} \mid 0 \text{ events}) \propto P(0 \text{ events} \mid \text{model true}) P(\text{model true})$$

Example of Bayes' Theorem Using Bayesian P (cont.)

$$P(0 \text{ events} \mid \text{model true}) = 0.05$$

$$P(0 \text{ events} \mid \text{model false}) = 1.0$$

0 events observed

Apply Bayes' Thm in a little more detail, with normalization:

Let "A" \leftrightarrow "0 events"; let "B" \leftrightarrow "model true". Recall:

$$P(B|A) = P(A|B) \times P(B) / P(A) .$$

Similarly, with $P(\text{not } B) = 1 - P(B)$

$$P(\text{not } B|A) = P(A| \text{not } B) \times P(\text{not } B) / P(A) .$$

$P(B|A) + P(\text{not } B|A) = 1$, so $P(A)$ is normalization. Can easily show

$$P(B|A) = 0.05 P(B) / (1 - 0.95 P(B)), \text{ i.e.,}$$

$$P(\text{model true} \mid 0 \text{ events})$$

$$= 0.05 P(\text{model true}) / (1 - 0.95 P(\text{model true})).$$

Example of Bayes' Theorem Using Bayesian P (cont.)

$P(\text{model true} \mid 0 \text{ events})$

$$= 0.05 P(\text{model true}) / (1 - 0.95 P(\text{model true})).$$

Limiting cases of very high and very low prior belief on model:

1) Let “model” be Standard Model, prior $P(\text{model true}) = 1 - \varepsilon_1$
for $\varepsilon_1 \ll 1 \Rightarrow P(\text{model true} \mid 0 \text{ events}) \approx 1 - 20\varepsilon_1$

Still very high degree of belief! Even if someone says,
“ $P(0 \text{ events} \mid \text{model true}) = 5\%$, and 0 events observed
means there is 5% chance the S.M. is true.” (UGH!)

2) Let “model” be large extra dimensions,
prior $P(\text{model true}) = \varepsilon_2 \ll 1$,
 $\Rightarrow P(\text{model true} \mid 0 \text{ events}) \approx 0.05 \varepsilon_2$

Low prior belief becomes even lower.

N.B. More realistic examples are of course more complex.

A Note re *Decisions*

Suppose that as a result of the previous experiment, your degree of belief in the model is $P(\text{model true} \mid 0 \text{ events}) = 1\%$.

And you need to *decide* on an action, e.g., announcing in a press release that the model is false, or making no announcement while taking more data.

Question: What should you *decide*?

A Note re *Decisions*

Suppose that as a result of the previous experiment, your degree of belief in the model is $P(\text{model true} \mid 0 \text{ events}) = 1\%$.

And you need to *decide* on an action, e.g., announcing in a press release that the model is false, or making no announcement while taking more data.

Question: What should you *decide*?

Answer: *Cannot be determined from the given information!*

You need in addition:

The *utility* function (or its negative, the *loss* function), which quantifies the relative costs (to You) of

- **Type I error:** announcing that the model is false, when it is true (thus eventually harming your reputation);
- **Type II error:** not announcing that the model is false when it *is* false, thus potentially allowing another experiment to make the announcement first.

A Note re *Decisions* (cont.)

Thus, Your *decision*, requires two subjective inputs: Your prior probabilities, and the relative costs to You of outcomes.

Statisticians often focus on decision-making.

In HEP, the tradition thus far is to communicate experimental results (well) short of formal decision calculations.

It should become clear later in lectures:

Frequentist (classical) “hypothesis testing” (especially with conventions like 95% C.L. or 5σ) is not a complete theory of decision-making!

It is important to keep this in mind, since the “accept/reject” language of classical hypothesis testing (later in lectures) is too simplistic for “deciding” in important situations.

Probability, Probability Density, Likelihood

These are key building blocks in both frequentist and Bayesian statistics, and it is crucial distinguish among them.

In the following, we let x be an observed or measured quantity; sometimes we use n if the observation is integer-valued and we want to emphasize that.

A “(statistical) model” is an expression specifying probabilities or probability densities for observing x .

We use μ for parameters (sometimes vector-valued) in the model. (Statistical literature prefers θ .)

Then the most common examples in HEP (discussed in the prerequisite reading in Leo’s book) are:

- **Binomial probability of n_{on} successes in n_{tot} trials, each with binomial parameter ρ :**

$$\text{Bi}(\text{non} \mid n_{\text{tot}}, \rho) = \frac{n_{\text{tot}}!}{n_{\text{on}}! (n_{\text{tot}} - n_{\text{on}})!} \rho^{n_{\text{on}}} (1 - \rho)^{(n_{\text{tot}} - n_{\text{on}})}$$

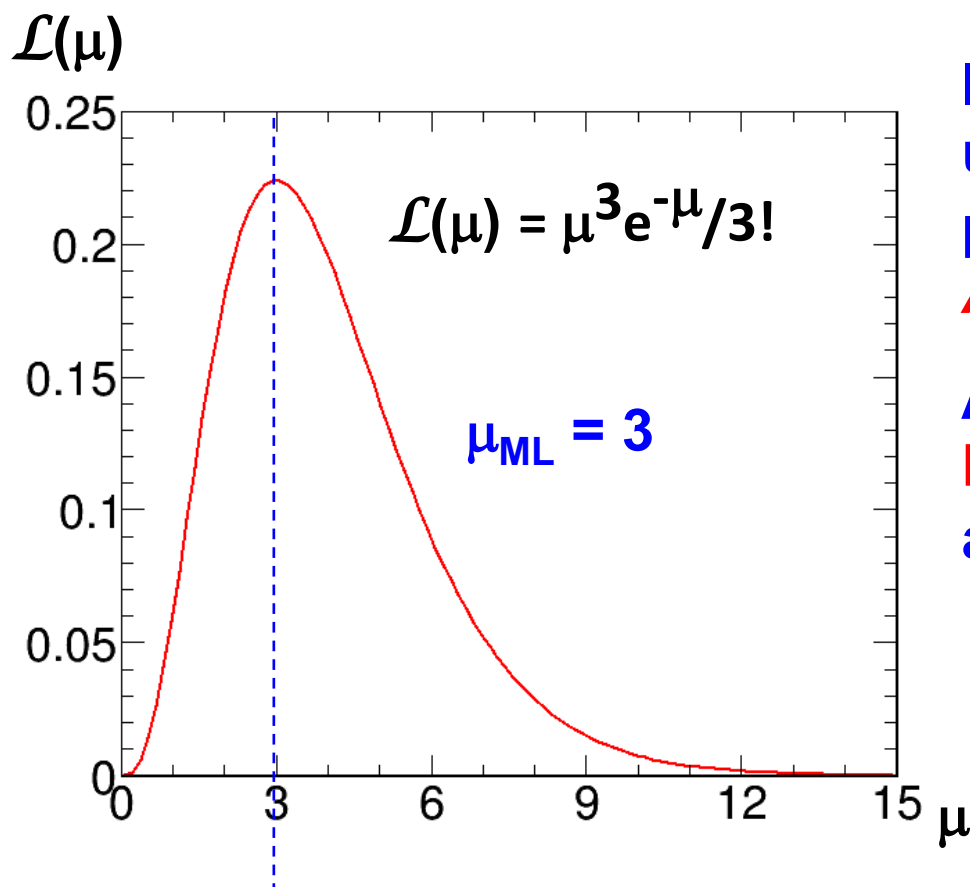
- **Poisson *probability* $P(n|\mu) = \mu^n \exp(-\mu)/n!$**
- **Gaussian *probability density function* (pdf) $p(x|\mu, \sigma)$:
 $p(x|\mu, \sigma)dx$ is differential of probability dP .**

**In Poisson case, suppose $n=3$ is observed.
Substituting *observed value* $n=3$ into $P(n|\mu)$ yields the
*likelihood function***

$$\mathcal{L}(\mu) = \mu^3 \exp(-\mu)/3!$$

Example likelihood function $\mathcal{L}(\mu) = \mu^3 \exp(-\mu)/3!$

Its maximum is at $\mu_{\text{ML}} = 3$.



**It is tempting to consider area under \mathcal{L} , but $\mathcal{L}(\mu)$ is *not* a probability density in μ :
*Area under \mathcal{L} is meaningless.***

**As we shall see,
Likelihood Ratios $\mathcal{L}(\mu_1) / \mathcal{L}(\mu_2)$
are useful and frequently used.**

Notation reminder

x denotes observable(s)

More generally, x is any convenient or useful function of the observed data, and is called a “statistic” or “test statistic”

μ denotes parameter(s)

$p(x|\mu)$ 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.

Change of observable variable (“metric”) x in pdf $p(x|\mu)$

For pdf $p(x|\mu)$ and 1-to-1 change of variable (metric) from x to $y(x)$, volume element modified by Jacobian. In 1D, $p(y) |dy| = p(x) |dx|$.

$$p(y(x)|\mu) = p(x|\mu) / |dy/dx|.$$

Jacobian modifies probability *density*, guaranties that

$$P(y(x_1) < y < y(x_2)) = P(x_1 < x < x_2),$$

(or equivalent with decreasing $y(x)$). I.e., guarantees that

Probabilities are invariant under change of variable x .

E.g., for $x \leftrightarrow \tau$ and $y(x) \leftrightarrow \Gamma = 1/\tau$, must have

$$P(\tau \in [\tau_1, \tau_2]) = P(\Gamma \in [1/\tau_2, 1/\tau_1])$$

Mode of probability *density* is *not* invariant (so, e.g., criterion of maximum probability density is ill-defined).

Likelihood *ratio* $\mathcal{L}(\mu_1) / \mathcal{L}(\mu_2)$ is invariant under change of variable x to $y(x)$. (Jacobian in denominator cancels that in numerator).

Probability Integral Transform

“...seems likely to be one of the most fruitful conceptions introduced into statistical theory in the last few years”

– Egon Pearson (1938) commenting on his father’s work.

Given continuous $x \in (a,b)$, and its pdf $p(x)$, define

$$y(x) = \int_a^x p(x') dx' .$$

Then $y \in (0,1)$ and easy to show that $p(y) = 1$ (uniform) for all y . (!)

So there always exists a metric y in which the pdf is uniform.

Many issues become more clear (or trivial) after this transformation. (If x is discrete, some complications.)*

Probability Integral Transform

“...seems likely to be one of the most fruitful conceptions introduced into statistical theory in the last few years”

– Egon Pearson (1938) commenting on his father’s work.

Given continuous $x \in (a,b)$, and its pdf $p(x)$, let

$$y(x) = \int_a^x p(x') dx' .$$

Then $y \in (0,1)$ and easy to show that $p(y) = 1$ (uniform) for all y . (!)

So there always exists a metric y in which the pdf is uniform.

Many issues become more clear (or trivial) after this transformation. (If x is discrete, some complications.)*

A look ahead, just mentioned here: The specification of a Bayesian prior pdf $p(\mu)$ for parameter μ is thus equivalent to the choice of the metric $g(\mu)$ in which the pdf is uniform. This is a *deep issue*, not always recognized by users of uniform prior pdf’s in HEP!

***And the inverse transformation provides for efficient M.C. generation of $p(x)$ starting from $\text{RAN}()$.**

Change of parameter μ in pdf $p(x|\mu)$

The pdf for x given parameter $\mu=3$ is the *same* as the pdf for x given $1/\mu=1/3$, or given $\mu^2=9$, or given any specified function of μ .

They all imply the same μ , and hence the same pdf for x .

In slightly confusing notation, that is what we mean by changing parameter from μ to $f(\mu)$, and saying that

$$p(x|f(\mu)) = p(x|\mu).$$

Thus the likelihood function $\mathcal{L}(\mu)$ is *invariant (!)* under reparametrization from parameter μ to $f(\mu)$:

$$\mathcal{L}(f(\mu)) = \mathcal{L}(\mu).$$

This reinforces the fact that $\mathcal{L}(\mu)$ is *not* a pdf in μ .

Bayes' Theorem Generalized to Probability Densities

Recall $P(B|A) \propto P(A|B) P(B)$.

For Bayesian P, continuous parameters such as μ are *random variables* with pdf's.

**Let pdf $p(\mu|x)$ be the conditional pdf for parameter μ , given data x .
As usual $p(x|\mu)$ is the conditional pdf for data x , given parameter μ . **Then Bayes' Thm becomes****

$p(\mu|x) \propto p(x|\mu) p(\mu)$.

Substituting in a particular set of observed data, x_0 :

$p(\mu|x_0) \propto p(x_0|\mu) p(\mu)$.

$$p(\mu|x_0) \propto p(x_0|\mu) p(\mu).$$

Recognizing the likelihood (variously written as $\mathcal{L}(x_0|\mu)$, $\mathcal{L}(\mu)$, or unfortunately even $\mathcal{L}(\mu|x_0)$), then

$$p(\mu|x_0) \propto \mathcal{L}(x_0|\mu) p(\mu), \text{ where:}$$

$p(\mu|x_0)$ = *posterior* pdf for μ , given the results of this expt

$\mathcal{L}(x_0|\mu)$ = *likelihood* function of μ from the experiment

$p(\mu)$ = *prior* pdf for μ , before applying the results of this expt

Note! There is one (and only one) probability density in μ on each side of the eqn, consistent with $\mathcal{L}(x_0|\mu)$ *not* being a density in μ .

Quick intro to “Bayesian” analysis

All equations up until now are true for *any* definition of probability P that obeys the axioms, including frequentist P , as long as the probabilities exist (for example if μ is sampled from an ensemble with known “prior” pdf).

The word “Bayesian” refers *not* to these equations, but to the choice of definition of P as *personal subjective degree of belief*.

Bayesian P applies to hypotheses and constants of nature (frequentist P does not), so many Bayesian-only applications.

Since Bayesian analysis *requires* a prior pdf, big issues in Bayesian analysis include:

- What prior pdf to use, and how sensitive is the result?
- How to interpret posterior probability if the prior pdf is not Your personal subjective belief?

Frequentist tools can be highly relevant to both questions!

Use of Bayesian posterior pdf $p(\mu|x_0)$

***Point estimation:* Some Bayesians use the posterior mode (aka maximum posterior density) as the point estimate of μ (though metric-dependent), others say point estimation is misguided.**

Since the Jacobian moves the mode around under change of parameter (say from lifetime τ to decay rate $\Gamma=1/\tau$), care must be used to interpret it. (Posterior median can be used in 1D.)

Use of Bayesian posterior pdf $p(\mu|x_0)$

Point estimation: Some Bayesians use the posterior mode (aka maximum posterior density) as the point estimate of μ (though metric-dependent), others say point estimation is misguided.

Since the Jacobian moves the mode around under change of parameter (say from lifetime τ to decay rate $\Gamma=1/\tau$), care must be used to interpret it. (Posterior median can be used in 1D.)

Interval estimation: Credibility of μ being in any interval $[\mu_1, \mu_2]$ can be calculated by integrating $p(\mu|x_0)$ over the interval.

Use of Bayesian posterior pdf $p(\mu|x_0)$

Point estimation: Some Bayesians use the posterior mode (aka maximum posterior density) as the point estimate of μ (though metric-dependent), others say point estimation is misguided.

Since the Jacobian moves the mode around under change of parameter (say from lifetime τ to decay rate $\Gamma=1/\tau$), care must be used to interpret it. (Posterior median can be used in 1D.)

Interval estimation: Credibility of μ being in any interval $[\mu_1, \mu_2]$ can be calculated by integrating $p(\mu|x_0)$ over the interval.

Hypothesis testing: Unlike frequentist statistics, testing credibility of whether or not μ equals a particular value μ_0 is *not* performed by examining intervals.*

One starts over with Bayesian model selection (later topic).

*assuming regular pdf p . Dirac δ -fns in p correspond to model selection, with its issues.

Use of Bayesian posterior pdf $p(\mu|x_0)$

Point estimation: Some Bayesians use the posterior mode (aka maximum posterior density) as the point estimate of μ (though metric-dependent), others say point estimation is misguided.

Since the Jacobian moves the mode around under change of parameter (say from lifetime τ to decay rate $\Gamma=1/\tau$), care must be used to interpret it. (Posterior median can be used in 1D.)

Interval estimation: Credibility of μ being in any interval $[\mu_1, \mu_2]$ can be calculated by integrating $p(\mu|x_0)$ over the interval.

Hypothesis testing: Unlike frequentist statistics, testing credibility of whether or not μ equals a particular value μ_0 is *not* performed by examining intervals.*

One starts over with Bayesian model selection (later topic).

Decision making: All Decisions about μ require not only $p(\mu|x_0)$ but also further input: the utility function.

*assuming regular pdf p . Dirac δ -fns in p correspond to model selection, with its issues.

Can “subjective” be taken out of “degree of belief”?

There are compelling arguments (Savage, De Finetti et al.) that Bayesian reasoning with *personal subjective P* is the uniquely “coherent” way (with technical definition of coherent) of updating *personal* beliefs upon obtaining new data.

The huge question is: can the Bayesian formalism be used by scientists to report the results of their experiments in an “objective” way (however one defines “objective”), and does any of the glow of coherence remain when subjective P is replaced by something else?

An idea vigorously pursued by physicist Harold Jeffreys in mid-20th century:

Can one define a prior $p(\mu)$ that contains as little information as possible?

“Uniform Prior” Requires a Choice of Metric

The really *really* thoughtless idea*, recognized by Jeffreys as such, but dismayingly common historically in HEP:

Just choose prior $p(\mu)$ uniform in whatever metric you happen to be using! (UGH!)

Recall that the probability integral transform *always* allows one to find a metric in which p is uniform (for continuous μ).

Thus the question “What is the prior pdf $p(\mu)$?” is equivalent to the question, “For what function $y(\mu)$ is $p(y)$ uniform?”

The choice $y(u) = u$ needs to be justified.

(It does not represent ignorance!)

*despite having a fancy name, Laplace’s Principle of Insufficient Reason

Jeffreys's Choice of Metric for Uniform Prior

For estimation, Harold Jeffreys answered the question using a prior uniform in a metric related to the Fisher information, calculated from curvature of the log-likelihood function averaged over sample space. Jeffreys priors:

Poisson signal mean μ , no background: $p(\mu) = 1/\sqrt{\mu}$

Poisson signal mean μ , mean background b : $p(\mu) = 1/\sqrt{\mu+b}$

Unbounded or bounded mean μ of Gaussian: $p(\mu) = 1$

RMS deviation of a Gaussian when mean fixed: $p(\sigma) = 1/\sigma$

Binomial parameter ρ , $0 \leq \rho \leq 1$: $p(\rho) = \rho^{-1/2}(1 - \rho)^{-1/2} = \text{Beta}(1/2, 1/2)$

Note: Jeffreys priors are commonly improper: cannot be normalized to 1. Considered by proponents not to be a disaster for estimation, as long as posterior pdf is proper, as is typical.

(Still a disaster for model selection.)

Jeffreys's Choice of Metric for Uniform Prior (cont.)

If parameter μ is changed to $f(\mu)$, the recipe for obtaining Jeffreys prior for $f(\mu)$ yields a different-looking prior that corresponds to the *same choice of uniform metric*.

So $p(\mu)$ is replaced by $p(f(\mu))$ that is correctly related by Jacobian, and *probabilities* (integrals of pdfs over equivalent endpoints) using Jeffreys prior are invariant under choices of different parameterizations.

What to call such Non-Subjective Priors?

- “Noninformative priors”? “Uninformative priors”? Traditional among statisticians, even though *they know it is misnomer*. (You should too!)
- “Vague priors”? “Ignorance priors”? “Default priors”?
- “Reference priors”? (Unfortunately also refers to a specific recipe of Bernardo)
- “Objective priors”? Despite the highly questionable use of the word, Jeffreys prior and its generalization by Bernardo and Berger are now widely referred to as “objective priors”.
- Kass and Wasserman J. Amer. Stat. Assn. 91 1343 (1996) give the best (neutral) name in my opinion:
Priors selected by “formal rules”.
 - *Required reading for anyone using Bayesian methods!*

Whatever the name, prior in one metric determines it in all other metrics: be careful in choice of metric in which it is uniform!

Jeffreys Prior (cont.)

For one-parameter models, the “Jeffreys prior” is the most common choice among statisticians for a “default” prior -- so common that statisticians are referring to the Jeffreys prior when they say “flat prior”.

The obvious generalization to multi-parameter models turns out to be problematic, so alternatives have been developed, notably by Bernardo (with Berger). In 1D, they provide a different rationale for Jeffreys’s prior, namely the prior that leads to a posterior pdf that is most dominated by the likelihood.

There are many subtleties. Beware!

Jeffreys's Prior (cont.)

A key point: priors such as the Jeffreys prior are based on the likelihood function and **thus inherently derived from the *measurement apparatus and procedure***, not from thinking about the parameter!

This may seem strange, but does give advantages, particularly for frequentist (!) coverage, as discussed briefly later.

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”.

Priors in high dimensions

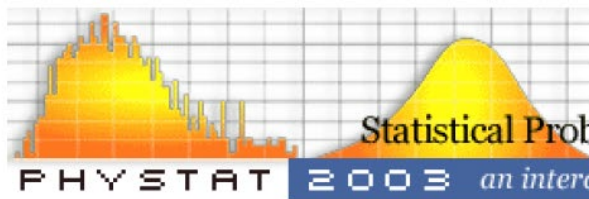
Is there a sort of informational phase space that can lead us to a sort of probability Dalitz plot? I.e., the desire is that structure in the posterior pdf represents information in the data, *not* the effect of Jacobians. *Notoriously hard problem!*

Be careful: Uniform priors push the probability away from the origin to the boundary! (Volume element in 3D goes as $r^2 dr$.)

State of the art for “objective” priors may be “reference priors” of Bernardo and Berger, but multi-D tools have been lacking.

Subjective priors also very difficult to construct in high dimensions: human intuition is poor.

- **Subjective Bayesian Michael Goldstein:** “meaningful prior specification of beliefs in probabilistic form over very large possibility spaces is very difficult and may lead to a lot of arbitrariness in the specification”.



Bayesians, Frequentists, and Physicists

Bradley Efron

*Department of Statistics and Department of Health Research and Policy,
Stanford University, Stanford, CA 94305, USA*

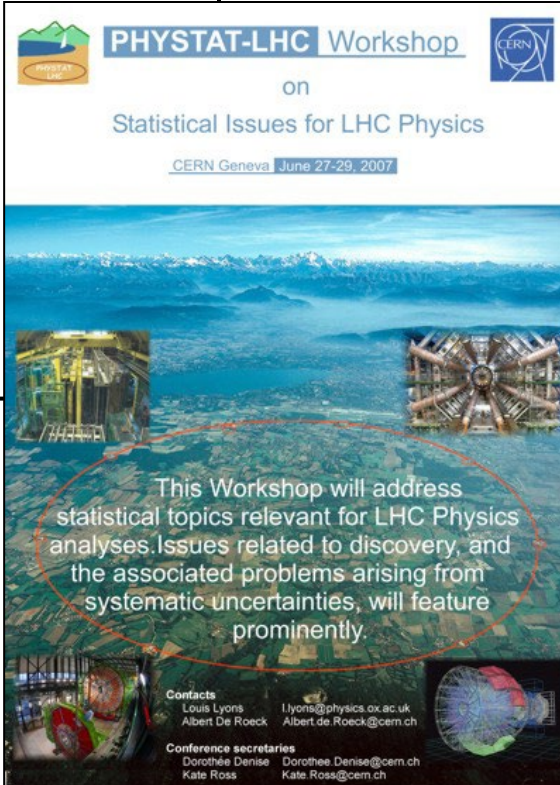
“Perhaps the most important general lesson is that the facile use of what appear to be uninformative priors is a dangerous practice in high dimensions.”

Sir David Cox at PhyStat-LHC 2007


Five faces of Bayesian statistics

- empirical Bayes; number of similar parameters with a frequency distribution
- neutral (reference) priors: Laplace, Jeffreys, Jaynes, Berger and Bernardo
- information-inserting priors (evidence-based)
- personalistic priors
- technical device for generating frequentist inference

**Currently in HEP, the main application is #5.
In particular, for *upper limits* we use uniform prior for
Poisson mean for frequentist reasons (See my AJP
paper, <http://aapt.scitation.org/doi/10.1119/1.17901> .)
Unfortunately some in HEP have also added a 6th:**



The poster for the PHYSTAT-LHC Workshop features a background image of a mountain range under a blue sky. In the foreground, there are two circular diagrams: one on the left showing a particle detector cross-section and one on the right showing a particle detector end view. A red circle highlights the central text. The CERN logo is in the top right corner.

PHYSTAT-LHC Workshop 

on
Statistical Issues for LHC Physics

CERN Geneva June 27-29, 2007

This Workshop will address statistical topics relevant for LHC Physics analyses. Issues related to discovery, and the associated problems arising from systematic uncertainties, will feature prominently.

Contacts
Louis Lyons llyons@physics.ox.ac.uk
Albert De Roeck Albert.deRoeck@cern.ch


Conference secretaries
Dorothee Denise Dorothee.Denise@cern.ch
Kate Ross Kate.Ross@cern.ch

Further information and registration at <http://cern.ch/physstat-lhc>

Cox's list, as I have seen it augmented in HEP

Six faces of Bayesian statistics

- empirical Bayes; number of similar parameters with a frequency distribution
- neutral (reference) priors: Laplace, Jeffreys, Jaynes, Berger and Bernardo
- information-inserting priors (evidence-based)
- personalistic priors
- technical device for generating frequentist inference
- Priors uniform in arbitrary variables, or in “the parameter of interest” (UGH!). This has no justification in modern subjective or objective Bayesian theory.

PHYSTAT-LHC Workshop 

on
Statistical Issues for LHC Physics

CERN Geneva June 27-29, 2007

This Workshop will address statistical topics relevant for LHC Physics analyses. Issues related to discovery, and the associated problems arising from systematic uncertainties, will feature prominently.

Contacts
Louis Lyons llyons@physics.ox.ac.uk
Albert De Roeck Albert.deRoeck@cern.ch

Conference secretaries
Dorothee Denise Dorothee.Denise@cern.ch
Kate Ross Kate.Ross@cern.ch

Further information and registration at <http://cern.ch/physstat-lhc>

Sensitivity Analysis

Since a Bayesian result depends on the prior probabilities, which are either personalistic or with elements of arbitrariness, it is widely recommended by Bayesian statisticians to study the *sensitivity* of the result to varying the prior.

I think that historically, too little emphasis was given to this by Bayesian advocates in HEP.

Sensitivity Analysis

An “objective Bayesian’s” point of view:

“Non-subjective Bayesian analysis is just a part -- an important part, I believe – of a healthy *sensitivity* analysis to the prior choice...”

– J.M. Bernardo, J. Roy. Stat. Soc., Ser. B 41 113 (1979)

Sensitivity analysis: A subjective Bayesian's point of view:

WHY BE A BAYESIAN?

Michael Goldstein

Dept. of Mathematical Sciences, University of Durham, England

From the Proceedings: “...Again, different individuals may react differently, and the sensitivity analysis for the effect of the prior on the posterior is the analysis of the scientific community...”

From his transparencies:

“Sensitivity Analysis is at the heart of scientific Bayesianism.”

The Institute for Particle Physics Phenomenology
will host a Conference on

ADVANCED STATISTICAL TECHNIQUES IN PARTICLE PHYSICS
at
The University of Durham, UK, March 18 - 22, 2002

Topics to be covered include:

Setting Limits Signal Significance Systematics
Combining Results Unfolding Convolution Simulation Issues
Multivariate Event Classification Techniques for Blind Analysis
Statistical Issues to do with Parton Distributions

Organising Committee
Roger Barlow (Manchester)
Bob Cousins (UCLA)
Glen Cowen (RHUL)
Fred James (CERN)
Dean Karlen (Carleton)

Jim Linnemann (Michigan State)
Louis Lyons (Oxford)
Bill Murray (RAL)
Harrison Prosper (Florida State)
Pekka Sinervo (Toronto)

Local Organising Committee
James Stirling
Mike Whalley
Linda Wilkinson

Further information and registration procedures can be obtained via
WWW at <http://www.ippp.dur.ac.uk/statistics/>

Bayesian Must-Read for HEP/Astro/Cosmo (incl discussion!)

Robert E. Kass and Larry Wasserman, “The Selection of Prior Distributions by Formal Rules,” J. Amer. Stat. Assoc. 91 1343 (1996).

Telba Z. Irony and Nozer D. Singpurwalla, “Non-informative priors do not exist: A dialogue with Jose M. Bernardo,” J. Statistical Planning and Inference 65 159 (1997).

James Berger, “The Case for Objective Bayesian Analysis,” Bayesian Analysis 1 385 (2006)

Michael Goldstein, “Subjective Bayesian Analysis: Principles and Practice,” Bayesian Analysis 1 403 (2006)

J.O. Berger and L.R. Pericchi, “Objective Bayesian Methods for Model Selection: Introduction and Comparison,” in Model Selection, Inst. of Mathematical Statistics Lecture Notes-Monograph Series, ed. P. Lahiri, vol 38 (2001) pp .135-207

Memorable Quotes Therein from Jim Berger

“The Case for Objective Bayesian Analysis,” Bayesian Analysis 1. See pp. 397, 459.

I call such analyses *pseudo-Bayes* because, while they utilize Bayesian machinery, they do not carry with them any of the guarantees of good performance that come with either true subjective analysis (with a very extensive elicitation effort) or (well-studied) objective Bayesian analysis. I will briefly discuss the problem with each of these pseudo-Bayes procedures.

I shall resist the temptation of saying more, because model selection is a can of worms for both objectivists and subjectivists.

Memorable Quotes Therein from Jim Berger

“The Case for Objective Bayesian Analysis,” Bayesian Analysis 1. See pp. 397, 459.

I call such analyses *pseudo-Bayes* because, while they utilize Bayesian machinery, they do not carry with them any of the guarantees of good performance that come with either true subjective analysis (with a very extensive elicitation effort) or (well-studied) objective Bayesian analysis. I will briefly discuss the problem with each of these pseudo-Bayes procedures.

I shall resist the temptation of saying more, because model selection is a can of worms for both objectivists and subjectivists.

Pseudo-Bayes analyses pop up from time to time in HEP. Flat priors, etc. (The worst are in Model Selection.) See my Comment at <https://arxiv.org/abs/0807.1330> and my “pseudo-Bayes detection” slides 62-68 at Tokyo PhyStat-nu, <http://indico.ipmu.jp/indico/event/82/session/9/contribution/16/material/slides/0.pdf> .

An excellent discussion by Harrison Prosper is in Ch. 12 of Data Analysis in High Energy Physics, Ed. By O. Behnke et al.

What can be computed without using a prior, with only the frequentist definition of P?

Not $P(\text{constant of nature is in some } \textit{specific} \text{ interval} \mid \text{data})$

Not $P(\text{Supersymmetry is true} \mid \text{data})$

Not $P(\text{Standard Model is false} \mid \text{data})$

Rather:

1) **Confidence Intervals** for constants of nature, parameter values, as defined in the 1930's by Jerzy Neyman.

Statements are made about probabilities in *ensembles* of intervals (fraction containing unknown true value)

2) **Likelihood ratios**, the basis for a large set of techniques for point estimation, interval estimation, and hypothesis testing.

Both can be constructed using the frequentist definition of P.

Confidence Intervals

“Confidence intervals”, and this phrase to describe them, were invented by Jerzy Neyman in 1934-37. Statisticians mean Neyman’s intervals (or an approximation) when they say “confidence interval”. In HEP the language is a little loose.

I highly recommend using “confidence interval” (and “confidence regions” when multi-D) only to describe intervals and regions corresponding to Neyman’s construction, described below, or by recipes of any origin that yield good approximations thereof.

Basic notions of confidence intervals

Conceptual idea in two sentences:

Given the model $p(x|\mu)$ and the observed value x_0 , for what values of μ is x_0 an “extreme” value of x ?

Include in the confidence interval $[\mu_1, \mu_2]$ those values of μ for which x_0 is *not* “extreme”.

Basic notions of confidence intervals

Conceptual idea in two sentences:

Given the model $p(x|\mu)$ and the observed value x_0 , for what values of μ is x_0 an “extreme” value of x ?

Include in the confidence interval $[\mu_1, \mu_2]$ those values of μ for which x_0 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.

Basic notions of confidence intervals (cont.)

Some 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.
So smallest values of x are most extreme.
Given x_0 , the confidence interval containing μ for which x_0 is not extreme will typically not contain largest values of μ .
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 $p(x|\mu)$, with the quantiles shorter in x (least integrated probability of x) containing higher-ranked x , with lower-ranked x added as the central quantile gets longer. Gives *central* confidence intervals for μ .

N.B. These three apply only when x is 1D.

(4th ordering, likelihood ratio used by F-C, still to come.)

Basic notions of confidence intervals (cont.)

Given model $p(x|\mu)$ and ordering of x , one chooses a fraction of highest-ranked values of x that are *not* considered as “extreme”.

This fraction is called the *confidence level (C.L.)*, say 68% or 95%.

We also define $\alpha = 1 - \text{C.L.}$, the lower-ranked, “extreme” fraction.

The *confidence interval* $[\mu_1, \mu_2]$ contains those values of μ for which x_0 is *not* “extreme” at the chosen C.L., *given the ordering*.

E.g., at 68% C.L., $[\mu_1, \mu_2]$ contains those μ for which x_0 is in the highest-ranked (least extreme) 68% values of x .*

**In this talk, 68% is more precisely 68.27%; 84% is 84.13%; etc.*

Basic notions of confidence intervals (cont.)

The endpoints of *central* confidence intervals at C.L. are the same as upper/lower limits with $1 - (1 - \text{C.L.})/2$.

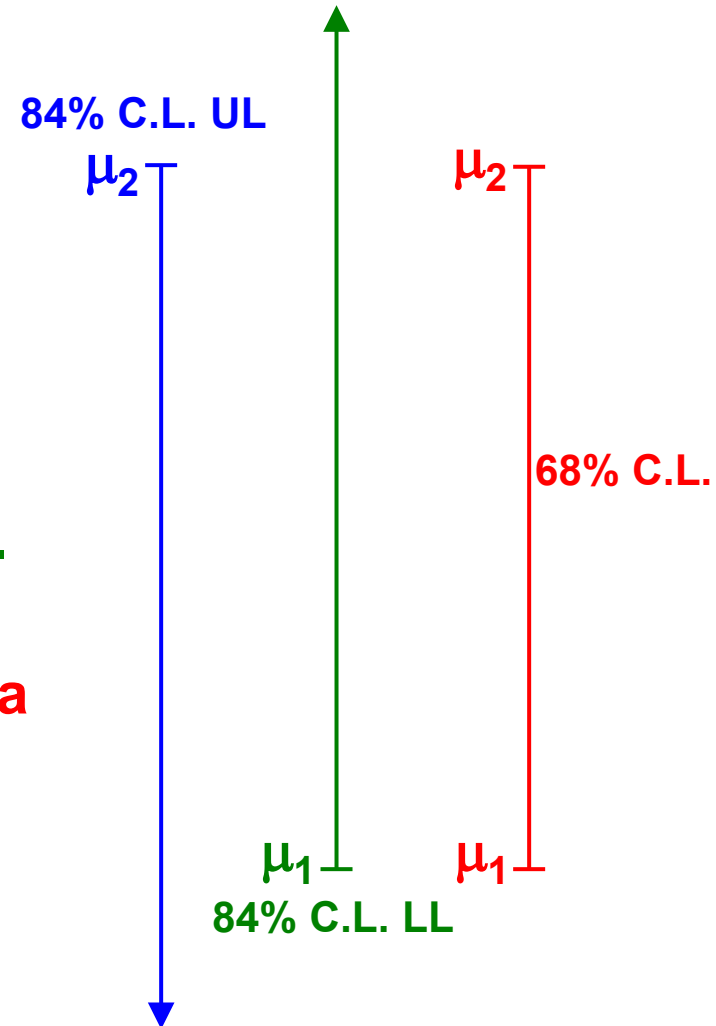
E.g.:

84% C.L. *upper* limit μ_2 excludes μ for which x_0 is in the lowest 16% values of x .

84% C.L. *lower* limit μ_1 excludes μ for which x_0 is in the highest 16% values of x .

Then $[\mu_1, \mu_2]$ includes μ for which x_0 is in the central 68% quantile of x values. It is a 68% C.L. *central* confidence interval (!)

Examples follow, first with continuous x , then with discrete x .

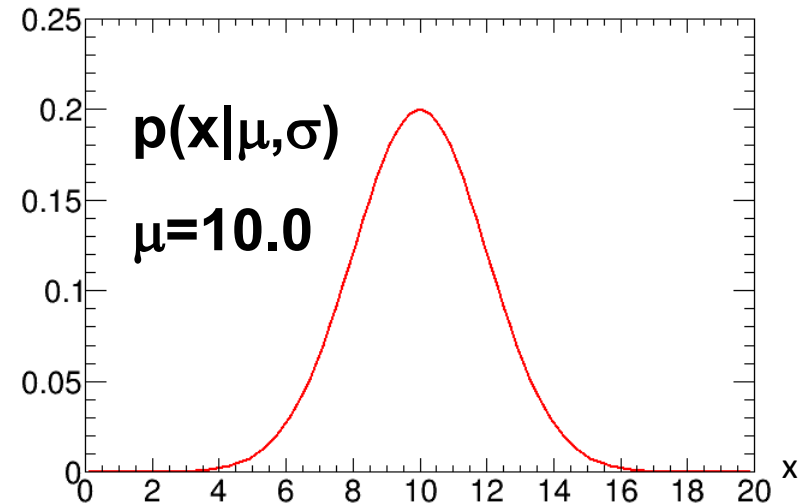


Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$

$$p(x|\mu,\sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\sigma(\mu) = (0.2) \mu$$

Plot of $p(x|\mu,\sigma)$ with $\mu=10.0$, $\sigma = 2.0$:

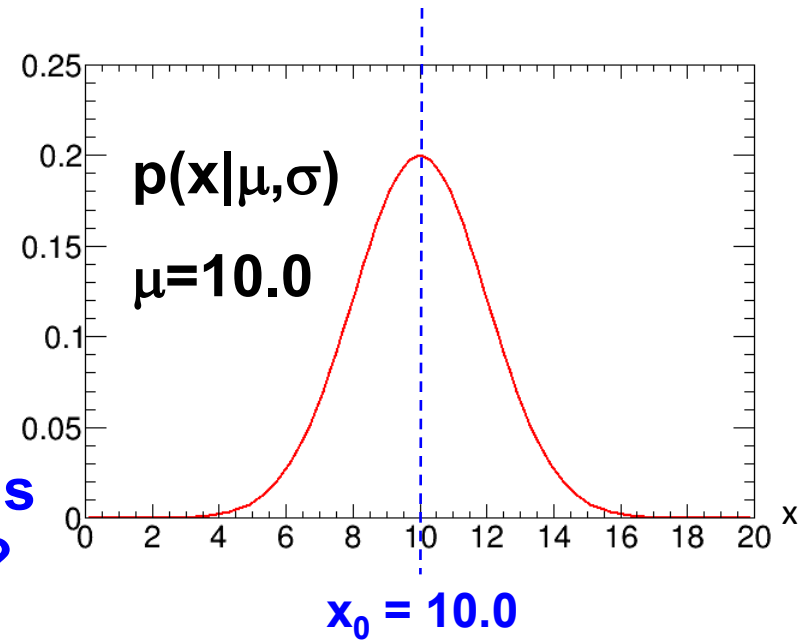


Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$

$$p(x|\mu,\sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\sigma(\mu) = (0.2) \mu$$

Plot of $p(x|\mu,\sigma)$ with $\mu=10.0$, $\sigma = 2.0$:



Suppose μ is unknown, and $x_0 = 10.0$ is observed. What can one say about μ ?

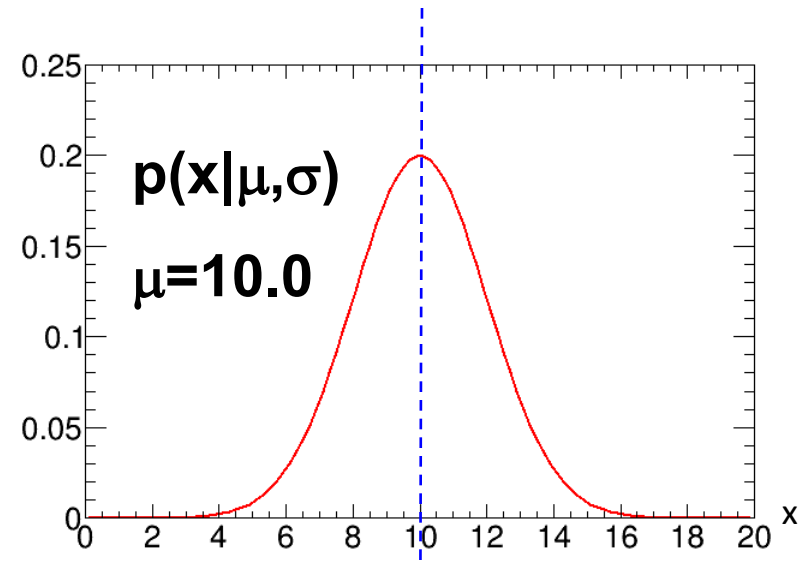
μ

Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$

$$p(x|\mu,\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

$$\sigma(\mu) = (0.2) \mu$$

Plot of $p(x|\mu,\sigma)$ with $\mu=10.0$, $\sigma = 2.0$:



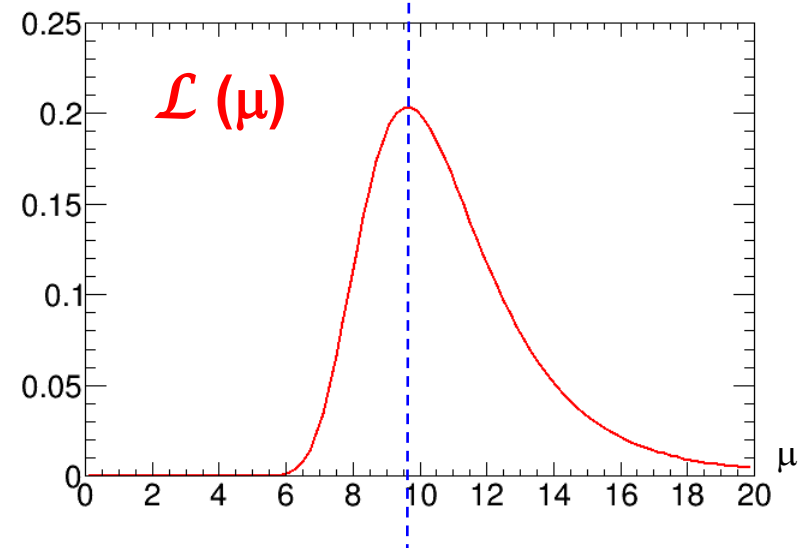
Suppose $x_0 = 10.0$ is observed.

$$\mathcal{L}(\mu) = \frac{1}{\sqrt{2\pi(0.2\mu)^2}} e^{-(x_0-\mu)^2/2(0.2\mu)^2}$$

$x_0 = 10.0$

Plot of $\mathcal{L}(\mu)$ for observed $x_0 = 10.0$:

$$\mu_{ML} = 9.63$$



What is *confidence interval* for μ ?

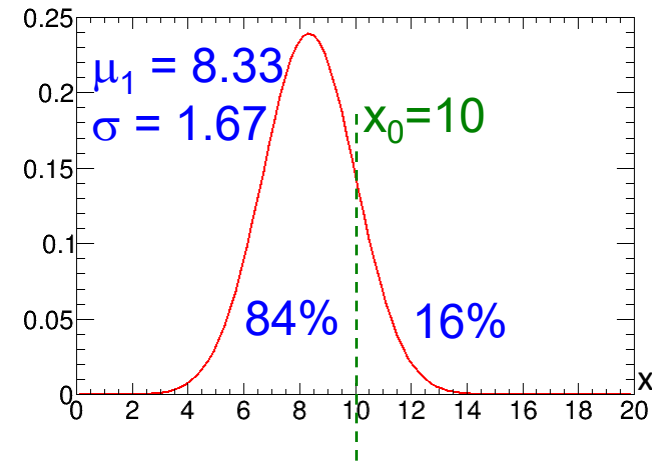
**Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$
Observed $x_0 = 10.0$.**

Find μ_1 such that 84% of $p(x|\mu_1,\sigma=0.2\mu_1)$ is below $x_0 = 10.0$; 16% of prob is above.

Solve: $\mu_1 = 8.33$.

$[\mu_1, \infty]$ is 84% C.L. confidence interval

μ_1 is 84% C.L. lower limit for μ .



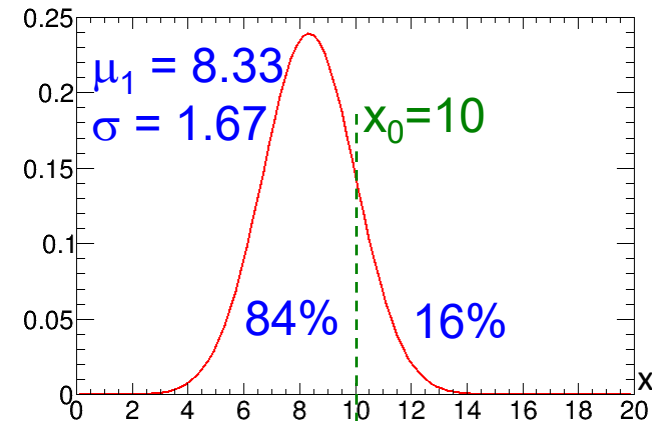
Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$ Observed $x_0 = 10.0$.

Find μ_1 such that 84% of $p(x|\mu_1,\sigma=0.2\mu_1)$ is below $x_0 = 10.0$; 16% of prob is above.

Solve: $\mu_1 = 8.33$.

$[\mu_1, \infty]$ is 84% C.L. confidence interval

μ_1 is 84% C.L. *lower limit* for μ .

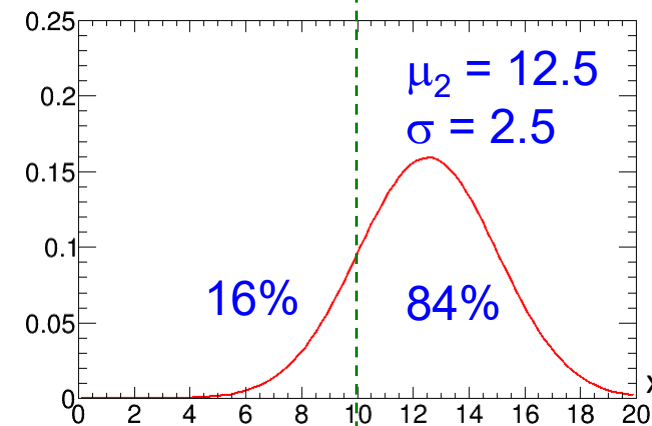


Find μ_2 such that 84% of $p(x|\mu_2,\sigma=0.2\mu_2)$ is above $x_0 = 10.0$; 16% of prob is below.

Solve: $\mu_2 = 12.5$.

$[-\infty, \mu_2]$ is 84% C.L. confidence interval

μ_2 is 84% C.L. *upper limit* for μ .



Then 68% C.L. *central* confidence interval is $[\mu_1, \mu_2] = [8.33, 12.5]$.

**Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$
Observed $x_0 = 10.0$.**

So the 68% C.L. *central* confidence interval is [8.33,12.52].

This is “exact”. Follows reasoning of E.B. Wilson, JASA 1927!

Note difference from (“Wald-like”) reasoning that proceeds as:

- 1) For $x_0 = 10.0$, minimum- χ^2 point estimate of μ is $\hat{\mu} = 10.0$.**
- 2) Then estimate $\hat{\sigma} = 0.2 \times \hat{\mu} = 2.0$.**
- 3) Then $\hat{\mu} \pm \hat{\sigma}$ yields interval [8.0,12.0].**

For (“exact”) confidence intervals, the reasoning must always involve probabilities for x calculated *considering particular possible true values of parameters*, as on previous slide!

Clearly the validity of the Wald-like approximate reasoning depends on how much $\sigma(\mu)$ changes for μ relevant to problem at hand. Beware!

Confidence intervals for binomial parameter ρ : Directly relevant to efficiency calculation in HEP

Recall $\text{Bi}(n_{\text{on}} | n_{\text{tot}}, \rho)$ for binomial probability of n_{on} successes in n_{tot} trials, each with **binomial parameter ρ** :

$$\text{Bi}(n_{\text{on}} | n_{\text{tot}}, \rho) = \frac{n_{\text{tot}}!}{n_{\text{on}}! (n_{\text{tot}} - n_{\text{on}})!} \rho^{n_{\text{on}}} (1 - \rho)^{(n_{\text{tot}} - n_{\text{on}})}$$

In repeated trials, n_{on} has **mean $n_{\text{tot}} \rho$** and

rms deviation $\sqrt{n_{\text{tot}} \rho (1 - \rho)}$

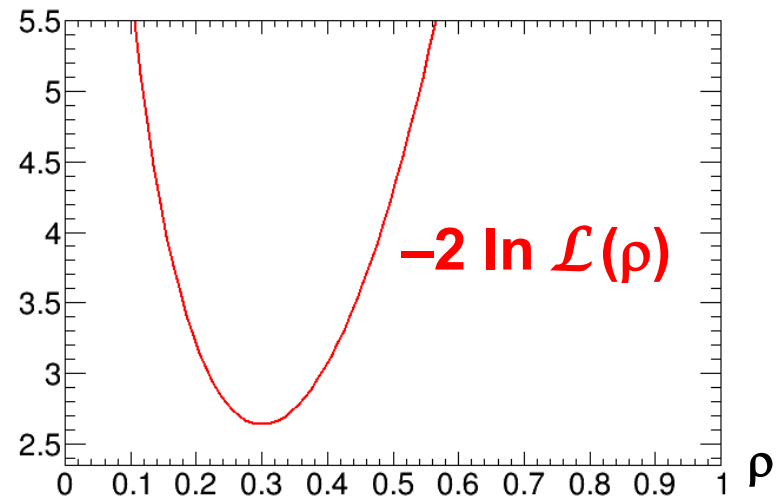
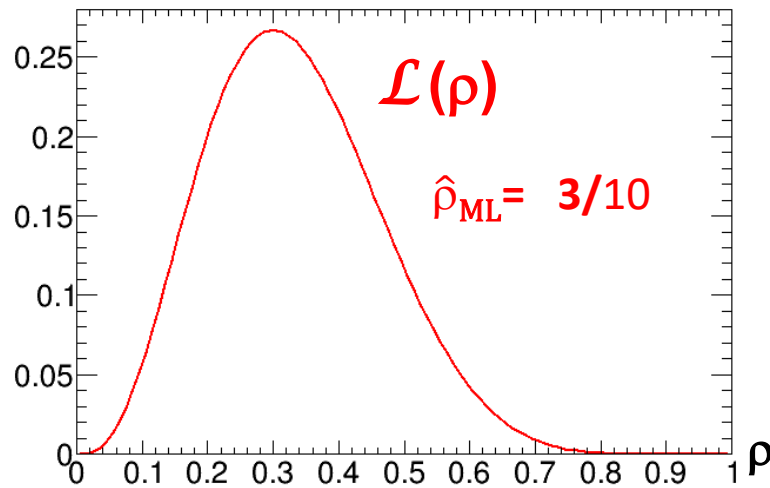
With observed successes n_{on} , **the M.L. point estimate $\hat{\rho}$ of ρ is (next slide)**

$$\hat{\rho} = n_{\text{on}} / n_{\text{tot}} .$$

What confidence interval $[\rho_1, \rho_2]$ should we report for ρ ?

Confidence intervals for binomial ρ (cont.)

Suppose $n_{\text{on}}=3$ successes in $n_{\text{tot}}=10$ trials.



Let's find "exact" 68% C.L.* *central* confidence interval $[\rho_1, \rho_2]$.

Recall shortcut above for central intervals:

Find *lower limit* ρ_1 with C.L. = $1 - (1 - 68\%)/2 = 84\%$

Find *upper limit* ρ_2 with C.L. = 84%

Then $[\rho_1, \rho_2]$ is 68% C.L. central confidence interval

**Recall in this talk, 68% is more precisely 68.27; 84% is 84.13%; etc.*

$n_{on} = 3$, $n_{tot} = 10$.

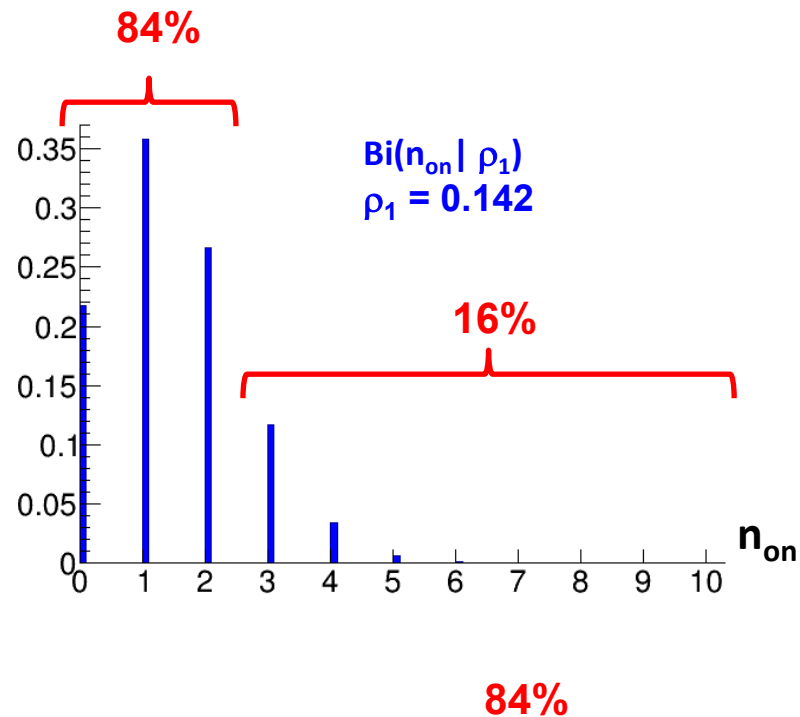
Find ρ_1 such that

$Bi(n_{on} < 3 \mid \rho_1) = 84\%$

$Bi(n_{on} \geq 3 \mid \rho_1) = 16\%$

(lower limit at 84% C.L.)

Solve: $\rho_1 = 0.142$



84%

$n_{on} = 3$, $n_{tot} = 10$.

Find ρ_1 such that

$Bi(n_{on} < 3 \mid \rho_1) = 84\%$

$Bi(n_{on} \geq 3 \mid \rho_1) = 16\%$

(lower limit at 84% C.L.)

Solve: $\rho_1 = 0.142$

And find ρ_2 such that

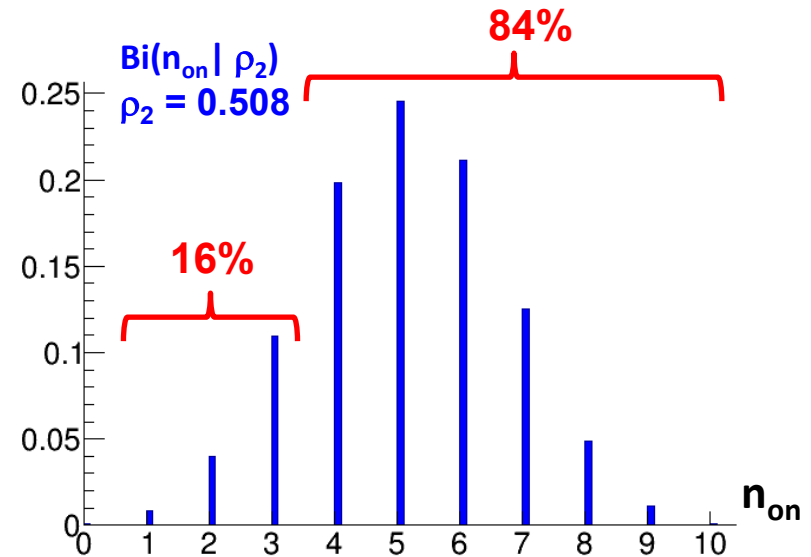
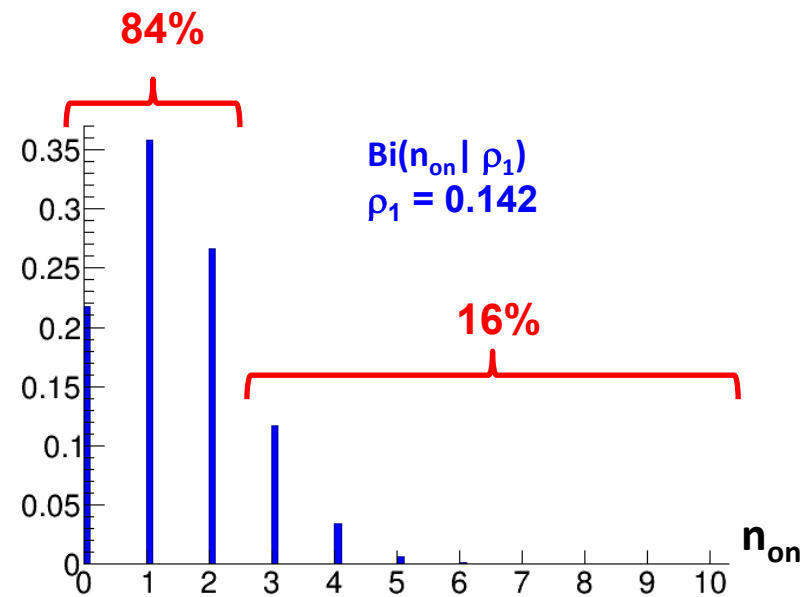
$Bi(n_{on} > 3 \mid \rho_2) = 84\%$

$Bi(n_{on} \leq 3 \mid \rho_2) = 16\%$

(upper limit at 84% C.L.)

Solve: $\rho_2 = 0.508$

Then $[\rho_1, \rho_2] = (0.142, 0.508)$
is *central* confidence interval
with 68% C.L. *Same as*
Clopper and Pearson (1934)



For Poisson example, see Fig. 3a,b; R. Cousins, Am. J. Phys. 63 398 (1995) DOI: 10.1119/1.17901

In HEP, such Clopper-Pearson intervals are the standard for a binomial parameter

In Particle Data Group's Review of Particle Physics since 2002.

Many tables and online calculators for C-P exist, e.g., <http://statpages.org/confint.html> .

But discreteness of x leads to an issue: C-P is criticized by some as “wastefully conservative” – see CHT paper below.

For a comprehensive review of both central and non-central confidence intervals for a binomial parameter and for the ratio of Poisson means, see Cousins, Hyme, and Tucker, <http://arxiv.org/abs/0905.3831> . Many are implemented in <https://root.cern.ch/doc/master/classTEfficiency.html> .

For related construction of upper/lower limits and central interval for Poisson mean, see R. Cousins, Am. J. Phys. 63 398 (1995)

Gaussian approximation for binomial conf. int.

As above, n_{on} has mean $n_{tot} \rho$ and rms deviation $\sqrt{n_{tot} \rho (1 - \rho)}$. One can approximate binomial by Gaussian with mean and rms

$$\mu(\rho) = n_{tot} \rho$$

$$\sigma(\rho) = \sqrt{n_{tot} \rho (1 - \rho)}$$

Idea is *not* to substitute $\hat{\rho}$ for ρ (big mistake), but rather to follow E.B. Wilson (1927): use above recipe for upper and lower limits:

1) Find ρ_1 such that Gauss($x \geq 3$ | mean ρ_1 , $\sigma(\rho_1)$) = 0.16

2) Find ρ_2 such that Gauss($x \leq 3$ | mean ρ_2 , $\sigma(\rho_2)$) = 0.16

This consistently uses the σ associated with each ρ . Leads to a quadratic equation with solution $[\rho_1, \rho_2] = [0.18, 0.46]$ which is the approximate 68% C.L. confidence interval known as the *Wilson score interval*. (See CHT paper.)

Avoid the Wald interval – no reason to use it

This “Wilson score interval” needs only the quadratic formula but is for some reason relatively unknown. It is tempting instead to substitute $\hat{\rho} = n_{\text{on}}/n_{\text{tot}}$ for ρ in the expression for σ :

$\hat{\sigma} = \sqrt{n_{\text{tot}} \hat{\rho} (1 - \hat{\rho})}$, obtaining the potentially disastrous “Wald interval”: $[\rho_1, \rho_2] = \hat{\rho} \pm \hat{\sigma}$.

The Wald interval does not use the correct logic for frequentist confidence! In fact when $n_{\text{on}} = 0$ (or $n_{\text{on}} = n_{\text{tot}}$), this gives $\hat{\sigma} = 0$.

Incredibly, failure of the Wald interval when $n_{\text{on}} = 0$ (or $n_{\text{on}} = n_{\text{tot}}$) has been used as a *foundational argument* in favor of Bayesian intervals in at least four public HEP postings (one retracted) and one published astro paper! (Typically the authors did not understand Bayesian statistics either, and used flat prior...)

HEP applications of conf. intervals for binomial param

1. **As mentioned, directly relevant to efficiency calculations.**
2. **Using a famous math identity, directly applicable to confidence intervals for *ratio of Poisson means*. Cousins, Hyme, and Tucker, <http://arxiv.org/abs/0905.3831> .**
3. **Then, applicable to significance (Z_{B_i}) of excess in a signal bin when sideband is used to estimate background. Cousins, Linnemann, and Tucker, <http://arxiv.org/abs/physics/0702156> .**
4. **Can even stretch #3 (using “rough correspondence”) to problem of signal bin when Gaussian estimate of mean background exists.**

Binomial example in Leo's book

The prerequisite reading for my lectures was Chapter 4 of William R. Leo, on *Techniques for Nuclear and Particle Physics experiments*, **except for Example 4.5 on page 100.**

One wants a 95% C.L. lower limit for the binomial parameter (his efficiency ε_0) if there are N successes in N trials, with N=100.

I just described how to do this in frequentist statistics:

Find *lower limit* ρ_1 with C.L. 95%: Find ρ_1 such that

$\text{Bi}(n_{\text{on}} \geq 100 \mid \rho_1) = 5\% \Rightarrow \text{Bi}(n_{\text{on}} = 100 \mid \rho_1) = 0.05 \Rightarrow (\rho_1)^{1/100} = 0.05$
 $\Rightarrow \rho_1 = 0.9705$. *Almost identical to Leo's solution, but (!):*

If you now read Leo's solution carefully, you should be able to see that "With some reflection..." is actually changing paradigms without notice and calculating a *Bayesian* credible interval that uses a uniform prior pdf, which is however not the "noninformative" prior that most statisticians would use.

Reaction rate example in Leo's book

In fact, I should have also excluded the top of page 99.

Bottom of page 98:

Let us assume therefore that the process has some mean reaction rate λ . Then the probability for observing no counts in a time period T is $P(0 | \lambda) = \exp(-\lambda T)$.

This is fine, a frequentist probability. But then he says:

This, now, can also be interpreted as the probability distribution for λ when no counts are observed in a period T .

Whoa, probability distribution for λ ? Changing paradigms?
Do you see the mistake? It is the classic mistake of inverting conditional probability, equating $P(\lambda | 0)$ with $P(0 | \lambda)$. The math is the same as changing to Bayesian paradigm and assuming uniform prior for λ . What is the justification???

Issues for upper-lower limits and *central* confidence intervals

For decades, issues with upper limits and central confidence intervals have been discussed in prototype problems in HEP:

- 1. Gaussian measurement resolution near a physical boundary (e.g. neutrino mass-squared is non-negative)**
- 2. Poisson signal mean measurement when observed number of events is less than mean expected background (so naïve “background-subtracted” cross section is negative)**

Many ideas put forward, PDG settled on three. Some history:
http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf

Today, I mostly stick to frequentist confidence intervals in this situation.

Beyond upper/lower limits and *central* confidence intervals

More general choices for ordering x in $p(x|\mu)$:

- For each μ , order x_0 using *likelihood ratio* $\mathcal{L}(x_0|\mu) / \mathcal{L}(x_0|\mu_{\text{best fit}})$.
Advocated in HEP by Feldman and Cousins in 1998
(and in Kendall and Stuart long before and since).
Applicable in both 1D and multi-D for x .

N.B. Recall that likelihood *ratios* as in F-C are independent of metric in x since Jacobian cancels.

Beyond upper/lower limits and *central* confidence intervals

More general choices for ordering x in $p(x|\mu)$:

- For each μ , order x_0 using *likelihood ratio* $\mathcal{L}(x_0|\mu) / \mathcal{L}(x_0|\mu_{\text{best fit}})$.
Advocated in HEP by Feldman and Cousins in 1998
(and in Kendall and Stuart long before and since).
Applicable in both 1D and multi-D for x .

N.B. Recall that likelihood *ratios* as in F-C are independent of metric in x since Jacobian cancels.

Note:

Ordering x by the probability *density* $p(x|\mu)$ is *not* recommended!
Recall that change of metric from x to $y(x)$ leads to Jacobian $|dy/dx|$ in $p(y|\mu) = p(x|\mu) / |dy/dx|$.

So ordering by $p(y|\mu)$ is different than ordering by $p(x|\mu)$ and all that follows depends on arbitrary choice of metric.

Neyman's Construction of Confidence Intervals

The general method for constructing “confidence intervals”, and the name, were invented by Jerzy Neyman in 1934-37.



The next few slides give basic outline.

It takes a bit of time to sink in – given how often confidence intervals are misinterpreted, the argument is perhaps a bit too ingenious.

In particular, you should understand that the confidence level does *not* tell you “how confident you are that the unknown true value is in the specific interval you report” – only a *subjective* Bayesian credible interval has that property!

Neyman's Construction of Confidence Intervals

Given $p(x|\mu)$ from a model:
For each value of μ , one
draws a horizontal *acceptance interval* $[x_1, x_2]$ such that
 $p(x \in [x_1, x_2] | \mu) = \text{C.L.} = 1 - \alpha$.
("Ordering principle" for x is
used to well-define.)

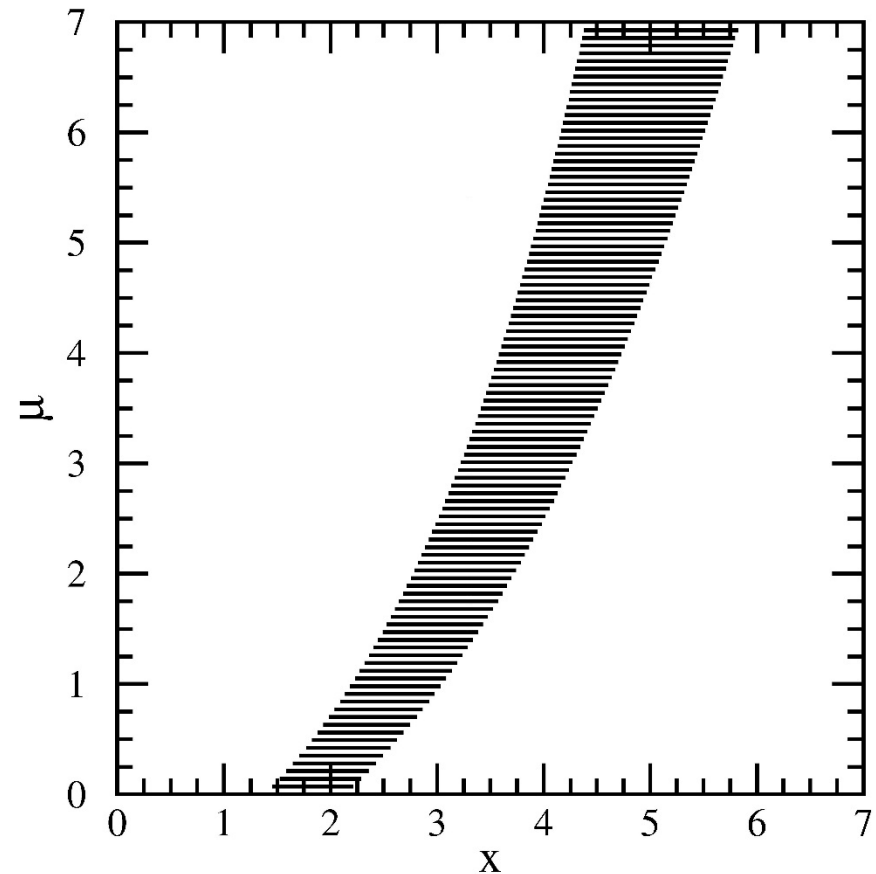


Figure from G. Feldman, R Cousins, Phys Rev D57 3873 (1998)

Neyman's Construction of Confidence Intervals

Given $p(x|\mu)$ from a model:
For each value of μ , one draws a horizontal *acceptance interval* $[x_1, x_2]$ such that $p(x \in [x_1, x_2] | \mu) = \text{C.L.} = 1 - \alpha$. (“Ordering principle” for x is used to well-define.)

Upon observing x , obtaining the value x_0 , one draws the vertical line through x_0 .

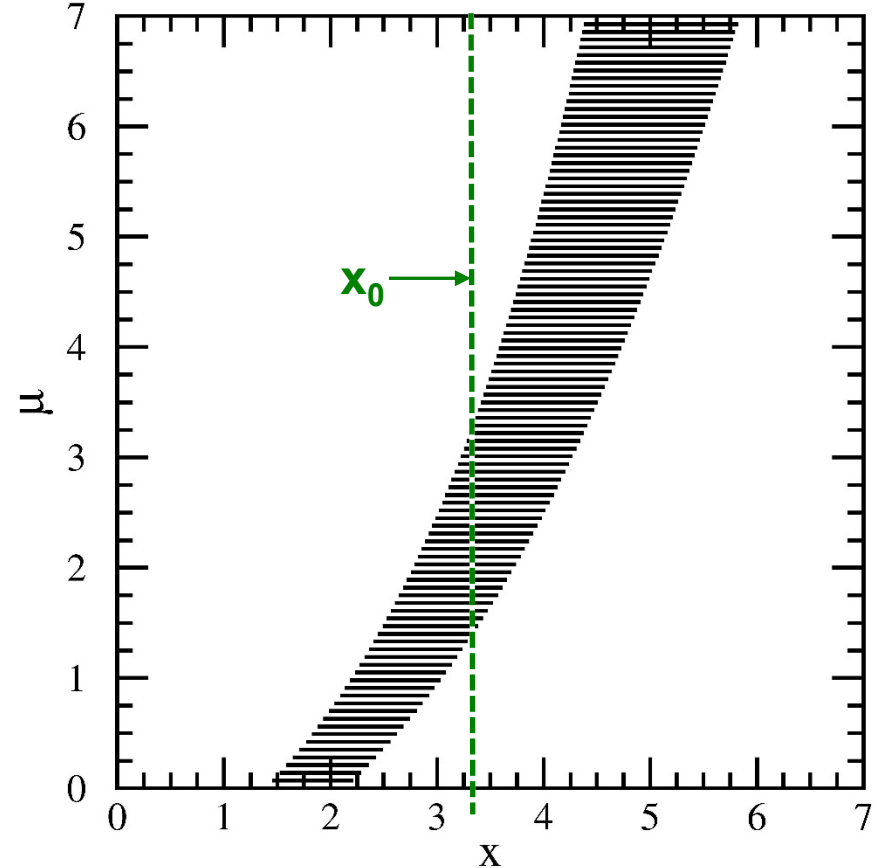
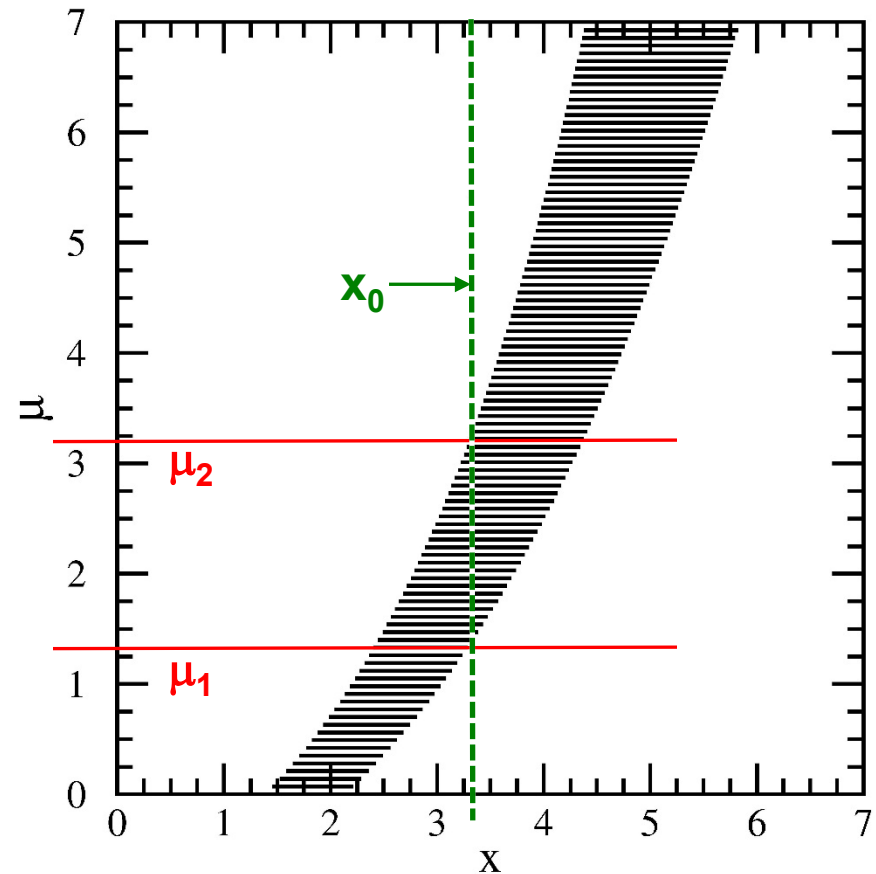


Figure from G. Feldman, R Cousins, Phys Rev D57 3873 (1998)

Neyman's Construction of Confidence Intervals

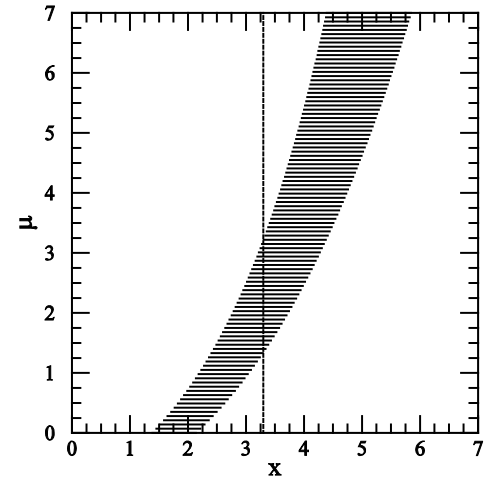
Given $p(x|\mu)$ from a model:
For each value of μ , one draws a horizontal *acceptance interval* $[x_1, x_2]$ such that $p(x \in [x_1, x_2] | \mu) = \text{C.L.} = 1 - \alpha$. (“Ordering principle” for x is used to well-define.)

Upon observing x , obtaining the value x_0 , one draws the vertical line through x_0 .



The vertical *confidence interval* $[\mu_1, \mu_2]$ with Confidence Level $\text{C.L.} = 1 - \alpha$ is the union of all values of μ for which the corresponding acceptance interval is intercepted by the vertical line.

Important note: x and μ need not have the same range, units, or (in generalization to higher dimensions) dimensionality!



I think it is *much* easier to avoid confusion when x and μ are qualitatively different.

Louis Lyons gives the example where x is the flux of solar neutrinos and μ is the temperature at the center of the sun.

I like examples where x and μ have different dimensions: Neyman's original paper has 2D observation space and 1D parameter space – to be discussed later.

Famous confusion re Gaussian $p(x|\mu)$ where μ is mass ≥ 0

It is *crucial* to distinguish between the data x , which *can* be negative (no problem), and the mass parameter μ , 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 $\mathcal{L}(\mu)$).

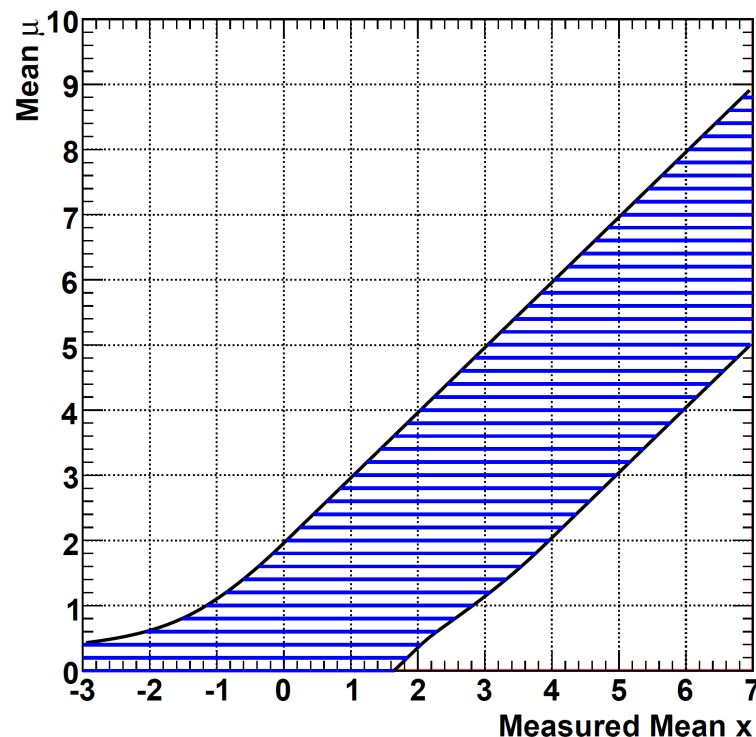
Famous confusion re Gaussian $p(x|\mu)$ where μ is mass ≥ 0

It is *crucial* to distinguish between the data x , which *can* be negative (no problem), and the mass parameter μ , 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 $\mathcal{L}(\mu)$).

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 Neyman construction graph has x of both signs but only non-negative $\mu \geq 0$. Example: Construction on right is LR ordering advocated by Feldman-Cousins



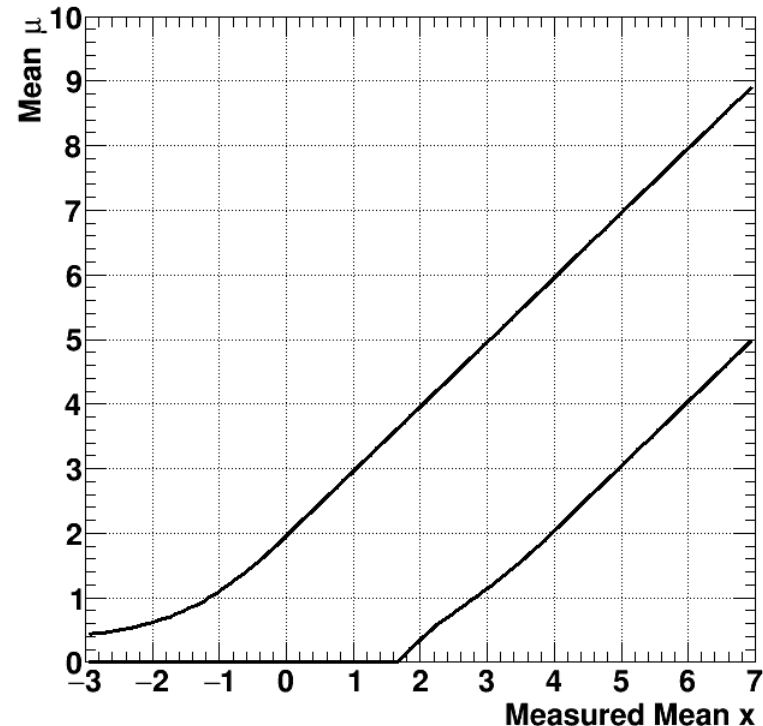
“Confidence Belt”

From the earliest days (as in 1934 Clopper-Pearson example in backup), the horizontal line segments are suppressed and only the envelope (black curves in figure), whose interior is called a *confidence belt*, is typically plotted.

I added the line segments for demonstrating the construction the F-C paper after reading Neyman’s 1937 paper, which had 2D acceptance regions shown. This practice is fortunately spreading.

Earlier works have statements that I found cryptic, such as “Notice that the confidence belt is *constructed horizontally but read vertically.*”

Area inside black curves is “confidence belt”



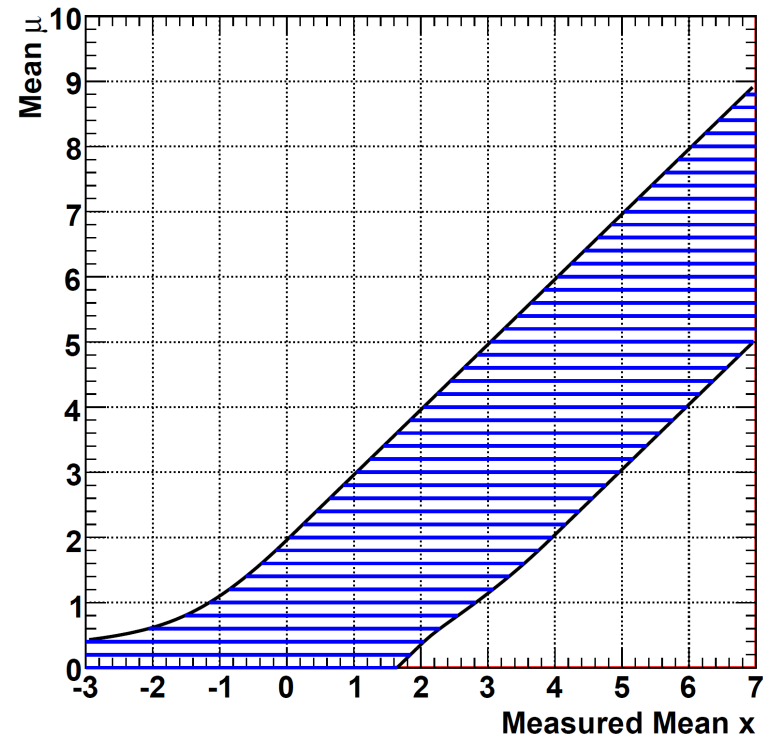
“Confidence Belt”

From the earliest days (as in 1934 Clopper-Pearson example in backup), the horizontal line segments are suppressed and only the envelope (black curves in figure), whose interior is called a *confidence belt*, is typically plotted.

I added the line segments for demonstrating the construction the F-C paper after reading Neyman’s 1937 paper, which had 2D acceptance regions shown. This practice is fortunately spreading.

Earlier works have statements that I found cryptic, such as “Notice that the confidence belt is *constructed horizontally* but *read vertically*.”

Area inside black curves is “confidence belt”



Confidence Intervals and Coverage

Recall: how is a *vector* defined in abstract math class?

Confidence Intervals and Coverage

Recall: how is a *vector* defined in abstract math class?

In math, one defines a *vector space* as a set with certain properties, and then the definition of a *vector* is “an element of a vector space”.

(A vector is not defined in isolation.)

Confidence Intervals and Coverage

Recall: how is a *vector* defined in abstract math class?

In math, one defines a *vector space* as a set with certain properties, and then the definition of a *vector* is “an element of a vector space”.

(A vector is not defined in isolation.)

Similarly, whether constructed in practice by Neyman’s construction or some other technique, a *confidence interval* is defined to be “a element of a confidence set”, where the *confidence set* is a set of intervals defined to have the property of frequentist *coverage* under repeated sampling:

Confidence Intervals and Coverage (cont.)

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 intervals obtained by Neyman's construction 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.

*** For μ_t , the probability that x_0 is in its acceptance region is C.L., by construction. For those x_0 's, the vertical line will intercept μ_t 's acceptance region, and so μ_t will be put into the confidence interval.**

Confidence Intervals and Coverage (cont.)

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

One of the complaints about confidence intervals is that the consumer often forgets (if he or she ever knew) that **the random variables in this equation are μ_1 and μ_2 , and not μ_t , and that coverage is a property of the set, not of an individual interval!**

Please don't forget!

Confidence Intervals and Coverage (cont.)

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

One of the complaints about confidence intervals is that the consumer often forgets (if he or she ever knew) that the random variables in this equation are μ_1 and μ_2 , and not μ_t , and that coverage is a property of the set, not of an individual interval!

Please don't forget!

A lot of confusion might have been avoided if Neyman had chosen the names “coverage intervals” and “coverage level”!

Maybe we can have a summit meeting treaty where frequentists stop saying “confidence” and Bayesians stop saying “noninformative”!

Confidence Intervals and Coverage (cont.)

It *is* true (in precisely the sense defined by the ordering principle used in the Neyman construction) that the confidence interval consists of those values of μ for which the observed x is among the C.L. least extreme values to be observed.

Famous 1934 Construction of Clopper and Pearson: Central Confidence Intervals for a Binomial Parameter

Biometrika, Vol. 26, No. 4. (Dec., 1934), pp. 404-413

This appeared in early days (before Neyman's 1937 paper giving comprehensive discussion of the construction) when the concept of what we now call coverage was discussed from differing points of view of both Fisher and Neyman, with different names and arguments for justifications.

In retrospect, this paper (citing Fisher, Neyman, and Neyman's students) corresponds exactly to a Neyman construction of central confidence intervals. See backup for details.

As I have discussed, for central intervals the answers are the same as those obtained from upper and lower limits.

The same method was applied to confidence intervals for *Poisson mean* by Garwood in his 1934 thesis, and published in 1936.

Also standard in HEP when no background!

Controversy when background – see PDG RRP.

Classical Hypothesis Testing

At this point, we set aside confidence intervals for the moment and consider from the beginning the nominally different topic of *hypothesis testing*.

In fact, we will soon find that in frequentist statistics, certain hypothesis tests will take us immediately back to confidence intervals. But first, we consider the more general framework.

Frequentist hypothesis testing, often called “classical” hypothesis testing, was developed by R.A. Fisher in unfriendly competition with J. Neyman and E. Pearson. Modern testing has a mix of ideas from both.

Classical Hypothesis Testing (cont.)

In Neyman-Pearson (N-P) hypothesis testing (James06), frame discussion in terms of null hypothesis H_0 (e.g. Standard Model) and an alternative H_1 (e.g., some Beyond-SM model).

Then $p(x|\mu)$ is different for H_0 and H_1 , either because parameter μ is different, or $p()$ itself is different.

Classical Hypothesis Testing (cont.)

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 probability α (smallish number).

Then “reject” H_0 if observed x_0 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)

Classical Hypothesis Testing (cont.)

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 probability α (smallish number).

Then “reject” H_0 if observed x_0 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)

To quantify the performance of this test if H_1 is true, we define:

β = probability (with x generated according to H_1) of accepting H_0 when it is false, i.e., not claiming a discovery when there is one (Type II error)

$1-\beta$ is called the *power* of the test.

So a given α will correspond to a threshold value C of the test statistic x (perhaps an output of a neural net). In HEP, events are “background” (H_0) if $x \leq C$ and “signal” (H_1) if $x > C$.

Classical Hypothesis Testing Jargon

There is a tradeoff between Type I and Type 2 errors.

Competing analysis algorithms can be compared by looking at graphs of power $1-\beta$ vs Type 1 error prob α at various μ , and at graphs of $1-\beta$ vs μ at various α (*power function*). See James06, pp. 258, 262.

Testing H_0 vs H_1 , including the binary classification problem of assigning items to one of two classes, is ubiquitous in science.

There are several sets of words used instead of N-P's Type I error prob α and Type II error prob β and their “power” $1-\beta$.

Classical Hypothesis Testing Jargon

Typically, H_0 is the “boring” hypothesis (background event, no disease, no signal of something new), and H_1 is something new (signal for new science, presence of disease, etc.)

Rejecting H_0 is called a “positive” result of the test (even if news of a disease), while not rejecting H_0 is a “negative” test results.

So test results are called true +, true –, false +, or false –.

“ROC curve” jargon for test performance

The acronym ROC (for “receiver operating characteristic”) from the early days of radar is commonly used in many fields of engineering and science.

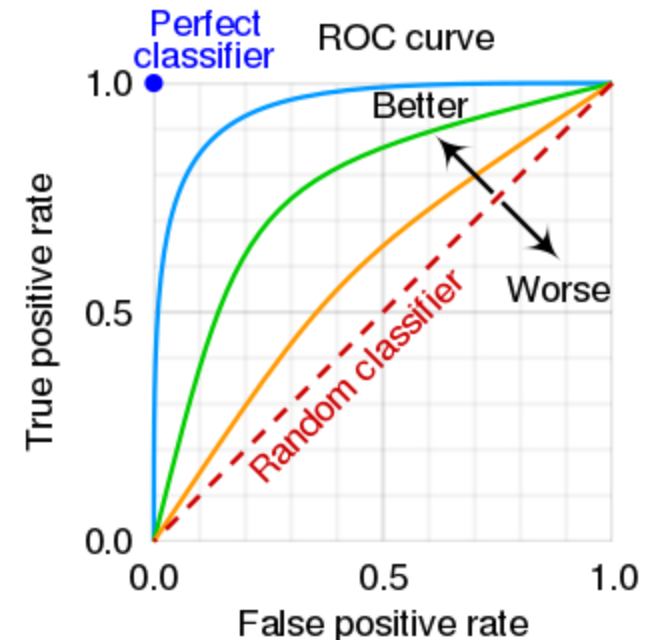
The **ROC curve** is a graph of power $1-\beta$ vs Type 1 error prob α , often with new names: **True Positive Rate vs. False Positive Rate**. This definition uses “rate” as a synonym for “probability” (UGH):

False Positive Rate

= $P(\text{reject } H_0 \mid H_0 \text{ true})$
= α = N-P Type I error prob.
(= $1 - \text{specificity}$ in medical jargon)

True Positive Rate

= $P(\text{reject } H_0 \mid H_0 \text{ false})$
= $1 - \beta$ = N-P power of test.
(= *sensitivity* in medical jargon)

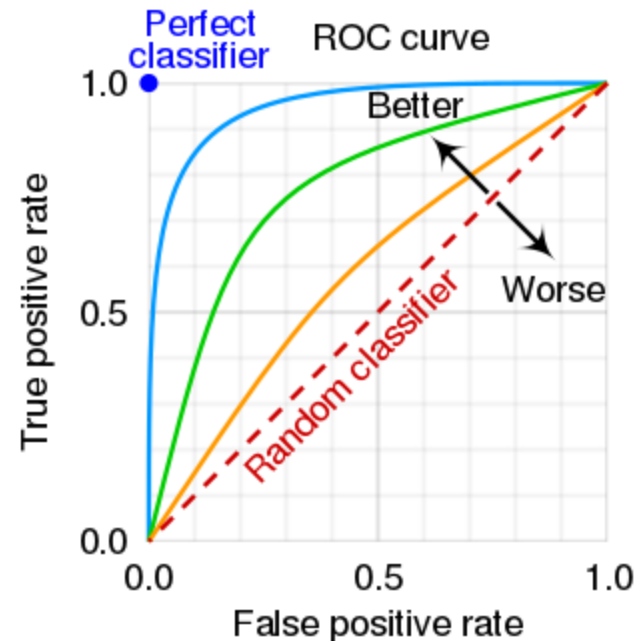


“ROC curve” (cont.)

Again, there is internally a threshold of a test statistic x that determines **tradeoff between α and $1 - \beta$, and hence FPR and TPR.**

See https://en.wikipedia.org/wiki/Receiver_operating_characteristic

Commonly used machine learning software scikit-learn provides ROC curves with TPR vs FPR.



Jargon in HEP for Classical Hypothesis Testing

P(called signal | bkgnd) = α is the “mistag” probability for bkgnd
(object is “tagged” as signal even though it is background)

P(called signal | signal) = $1 - \beta$ is the “efficiency” for signal
(fraction of true signal events that are tagged as signal)

P(called bkgnd | signal) = β is the “inefficiency” for signal

As C is changed, one maps out curves with axes labeled with these terms instead of $1 - \beta$ vs α or TPR vs. FPR.

Note: Assignment of H_0 and H_1 is arbitrary and can be reversed, with corresponding reversal of α and β .

Typically H_0 is the simpler hypothesis.

Warning: “ROC curves” in the TMVA package in ROOT are labeled “background rejection” vs “efficiency”, meaning $1 - \alpha$ vs $1 - \beta$.

Beware of the many definitions and conventions!

Review of jargon for medical tests

Recall: results are true +, true –, false +, or true –.

P(test is + | disease) = $1 - \beta$ is called the “sensitivity” of test:
Probability of correct diagnosis if actually diseased (true +) .

P(test is – | no disease) = $1 - \alpha$ is called the “specificity” of test:
Probability of correct diagnosis if no disease (true –) .

Value of threshold C determines tradeoff between sensitivity and specificity.

Jargon medical jargon (cont.)

Question: Suppose that we know the sensitivity and the specificity. Consider the people with positive test results. What fraction of them have the disease?

I.e., what is $P(\text{disease} \mid \text{test is } +)$? This is called the *positive predictive value (PPV)* in medicine.

Jargon medical jargon (cont.)

Question: Suppose that we know the sensitivity and the specificity. Consider the people with positive test results. What fraction of them have the disease?

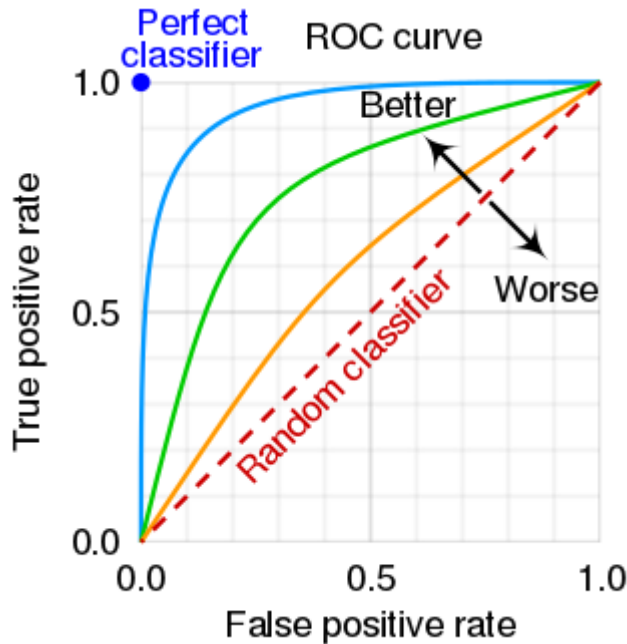
I.e., what is $P(\text{disease} \mid \text{test is } +)$? This is called the *positive predictive value (PPV)* in medicine.

Answer: *Cannot be determined from the given information!*

Need in addition: $P(\text{disease})$, the true fraction of *all* people that have the disease. Then Bayes's Thm inverts the conditionality:
 $P(\text{disease} \mid \text{test is } +) \propto P(\text{test is } + \mid \text{disease}) P(\text{disease})$

Exercise: write sensitivity, specificity, and PPV in terms of numbers of true +, true –, false +, or true –. Check your answers:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8698426/>

Classical Hypothesis Testing (cont.)



Where to live on the ROC curve (choice of “the operating point” $\alpha = \text{FPR}$) is a *long* discussion (even longer when considered as sample size increases, so curve moves toward upper left.)

The N-P language of “accept” or “reject” H_0 should not be mistaken for a complete theory of decision-making:

Decision on whether to declare discovery requires 2 more inputs:

- 1) Prior belief in H_0 vs H_1 . (Can affect choice of α)
- 2) Cost of Type I error (false discovery claim) vs cost of Type II error (missed discovery). (Can also affect choice of α)

A one-size-fits-all criterion of α corresponding to 5σ is without foundation!

Classical Hypothesis Testing: Simple Hypotheses

In idealized cases, a hypothesis may have no floating (unfixed) parameters.

N-P called such hypotheses *simple*, in contrast to *composite* hypotheses that have unfixed parameters.

Examples in HEP where both H_0 and H_1 are *simple* are rare, but we do have a few examples where the quantity of interest is simple in both hypotheses, and the role of unfixed parameters does not spoil the “simplicity”, e.g., H_0 vs H_1 being:

“jet originated from a quark” vs “jet originated from a gluon”

spin-1 vs spin-2 for a new resonance in $\mu^+\mu^-$

$J^P=0^+$ vs $J^P=0^-$ for the Higgs-like boson

Simple Hypotheses Testing: Neyman-Pearson Lemma

IX. *On the Problem of the most Efficient Tests of Statistical Hypotheses.*

By J. NEYMAN, *Nencki Institute, Soc. Sci. Lit. Varsoviensis, and Lecturer at the Central College of Agriculture, Warsaw,* and E. S. PEARSON, *Department of Applied Statistics, University College, London.*

(Communicated by K. PEARSON, F.R.S.)

(Received August 31, 1932.—Read November 10, 1932.)

Phil. Transactions of the Royal Society of London. Vol. 231, (1933), pp. 289-337

If Type I error probability α is specified in a test of simple hypothesis H_0 against simple hypothesis H_1 , then the Type II error probability β is minimized by ordering \mathbf{x} according to the likelihood ratio $\lambda = \mathcal{L}(\mathbf{x} | H_0) / \mathcal{L}(\mathbf{x} | H_1)$.

One finds cutoff $\lambda_{\text{cut},\alpha}$ for that α and rejects H_0 if $\lambda \leq \lambda_{\text{cut},\alpha}$.

The “lemma” applies only to a very special case: no nuisance parameters, not even undetermined parameters of interest! But it has inspired many generalizations, and likelihood ratios are an oft-used component of both frequentist and Bayesian methods.

For an outline of a proof, see Stuart99, p. 176

Nested Hypothesis Testing

In contrast to two disjoint simple hypotheses, it is common in HEP for H_0 to be *nested* in H_1 .

This happens when there is an undetermined parameter μ in H_1 , and H_0 corresponds to a particular parameter value μ_0 (e.g., zero, 1, or ∞). So consider:

$H_0: \mu = \mu_0$ (the “point null”, or “sharp hypothesis”) vs
 $H_1: \mu \neq \mu_0$ (the “continuous alternative”).

Common examples:

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

$H^0 \rightarrow \gamma\gamma$ before discovery of this decay, $\mu =$ signal strength:
null $\mu_0 = 0$, alternative $\mu > 0$

$H^0 \rightarrow \gamma\gamma$ after discovery of this decay:
null $\mu_0 =$ Standard Model prediction, alternative any other $\mu \neq \mu_0$

Nested Hypothesis Testing (cont.)

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

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

In classical/frequentist formalism (but *not* Bayesian formalism), the theory of such nested tests maps to that of confidence intervals!

Nested Hypothesis Testing (cont.)

$H_0: \mu = \mu_0$ (the “point null”, or “sharp hypothesis”) vs
 $H_1: \mu \neq \mu_0$ (the “continuous alternative”).

In classical/frequentist formalism (but *not* Bayesian formalism), the theory of such nested tests maps to that of confidence intervals!

Intuitive argument:

Having observed data x_0 , suppose the 90% C.L. confidence interval for μ is $[\mu_1, \mu_2]$.

This contains all values of μ for which observed x_0 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

$H_0: \mu = \mu_0$ (the “point null”, or “sharp hypothesis”) vs
 $H_1: \mu \neq \mu_0$ (the “continuous alternative”).

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

Intuitive argument:

Having observed data x_0 , suppose the 90% C.L. confidence interval for μ is $[\mu_1, \mu_2]$.

This contains all values of μ for which observed x_0 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_0 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 in 90% C.L. confidence interval $[\mu_1, \mu_2]$.

Nested Hypothesis Testing Duality

Given an ordering:

Test of $\mu = \mu_0$ vs $\mu \neq \mu_0$ at significance level α

\leftrightarrow Is μ_0 in confidence interval for μ with C.L. = $1 - \alpha$?

“There is thus no need to derive optimum properties separately for tests and for intervals; there is a one-to-one correspondence between the problems as in the dictionary in Table 20.1”

Stuart99, p. 175. [Table in backup slides] E.g.,

$$\alpha \leftrightarrow 1 - \text{C.L.}$$

Equal-tailed test \leftrightarrow central confidence intervals

One-tailed tests \leftrightarrow Upper/lower limits

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

Nested Hypothesis Testing (cont.)

Test $\mu = \mu_0$ at $\alpha \leftrightarrow$ Is μ_0 in conf. int. for μ with C.L. = $1 - \alpha$

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):

We emphasized a “new” ordering principle based on LR. While paper was “in proof”, Gary realized that “our” intervals were simply those obtained by “inverting” the LR hypothesis test. In fact it was all on 1¼ pages of “Kendall and Stuart”, plus nuisance parameters! → 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

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

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.

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 is typically converted to Z-value, the equivalent number of Gaussian sigma.*

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

(Z is unfortunately sometimes called “the significance S”),

*Although these lectures are not “statistics in practice”, I mention ROOT commands for *one-tailed* conversions:

```
zvalue = -TMath::NormQuantile(pvalue)
pvalue = 0.5*TMath::Erfc(zvalue/sqrt(2.0))
```

(Thanks, Igor Volobouev.) Note that $p\text{-value} > 0.5$ means $Z\text{-value} < 0$.

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.

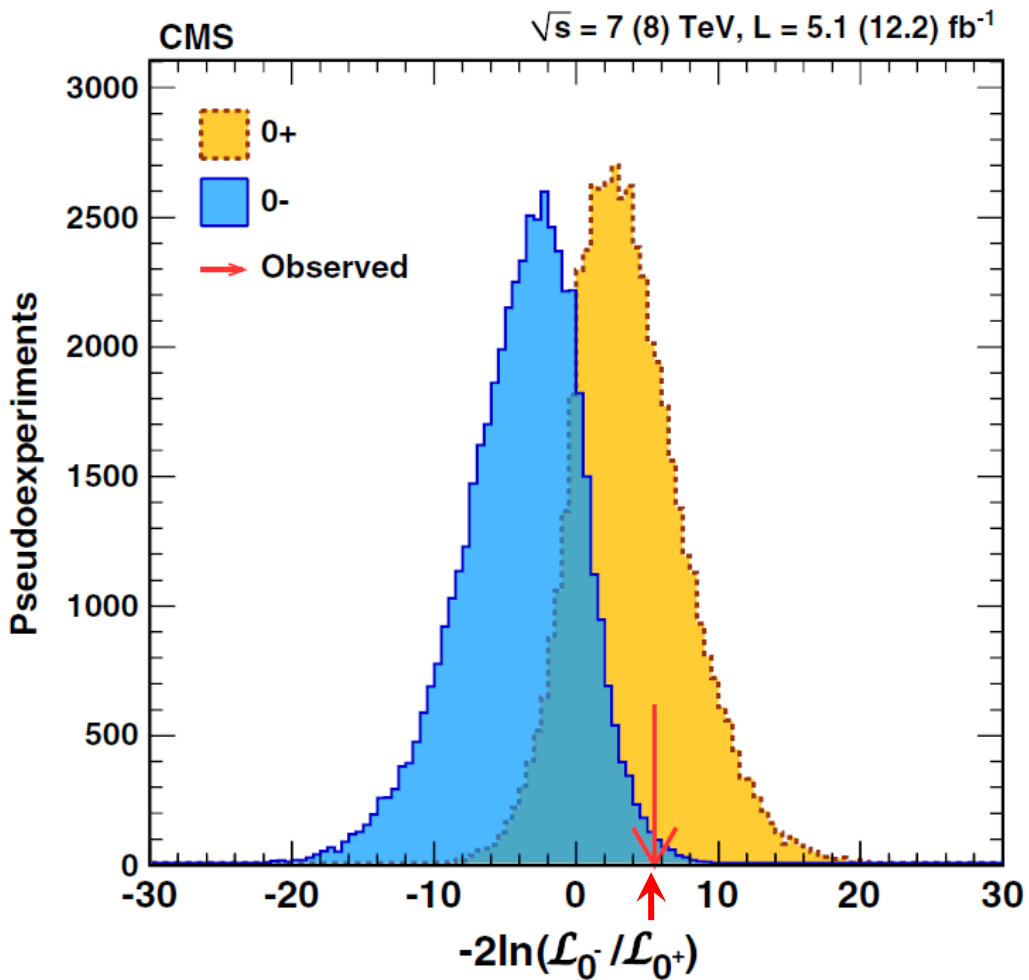
I defend their use in HEP. See <https://arxiv.org/abs/1310.3791>.)

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) !**

Early CMS Higgs spin-parity test of 0^+ vs. 0^-



Paper reported (fixing typo):

- 1) $-2\ln(\mathcal{L}_{0^-} / \mathcal{L}_{0^+}) = 5.5$ favoring 0^+
- 2) for $H_0: 0^-$, p-value = 0.0072
- 3) for $H_0: 0^+$, p-value = 0.7
- 4) $CL_s = (0.0072)/(1-0.7) = 0.024$,
 “a more conservative value for judging whether the observed data are compatible with 0^- ”

FIG. 3 (color online). Expected distribution of $-2 \ln \mathcal{L}_{0^-} / \mathcal{L}_{0^+}$ under the pure pseudoscalar and pure scalar hypotheses (histograms). The arrow indicates the value determined from the observed data.
 CMS, Phys. Rev. Lett. 110 (2013) 081803

N.B. See backup for figure and pointer to paper by Demortier and Lyons discussing two p-values in simple-vs-simple case.

Classical frequentist goodness of fit (g.o.f.)

If H_0 is specified but the alternative H_1 is not, then only the Type I error probability α can be calculated, since the Type II error probability β depends on H_1 .

A test with this feature is called a test for *goodness-of-fit* (to H_0). (Fisher called them significance tests.)

With no alternative specified, the question “Which test is best?” is thus ill-posed.

Despite the popularity of tests with universal maps from test statistics to α (in particular χ^2 and Kolmogorov tests), they may be ill-suited for many problems: they may have poor power $(1 - \beta)$ against relevant alternative H_1).

A plethora of possible tests in 1D are described in the book by D'Agostino and Stephens, a must-read for those wanting to invent a new test.

Classical frequentist goodness of fit (cont.)

As multi-D unbinned ML fits have proliferated in recent decades, there are increasing needs for multi-D unbinned g.o.f. tests.

E.g., is it reasonable that 1000 events scattered in a 5D sample space have been drawn from a particular pdf (which may have parameters that were fit using an unbinned M.L. fit to those 1000 events)?

This is an ill-posed question, but we are looking for good omnibus tests. Then getting the null distribution of the test statistic from simulation is typically doable, it seems.

One can follow an unbinned ML fit with a binned g.o.f. test such as χ^2 , but this brings in its own issues.

At a loss of power but increase in transparency, one can also perform tests on 1D or 2D distributions of the marginalized densities.

Machine learning is also having an impact.

See Appendix B of my writeup on arXiv for more on g.o.f.

Likelihood (Ratio) Intervals for 1 parameter

Recall: Likelihood function $\mathcal{L}(\mu)$ is invariant under reparametrization from μ to $f(\mu)$: $\mathcal{L}(\mu) = \mathcal{L}(f(\mu))$.

So likelihood ratios $\mathcal{L}(\mu_1) / \mathcal{L}(\mu_2)$ and log-likelihood differences $\ln \mathcal{L}(\mu_1) - \ln \mathcal{L}(\mu_2)$ are also invariant.

After using maximum-likelihood method to obtain estimate $\hat{\mu}$ that maximizes either $\mathcal{L}(\mu)$ or $\mathcal{L}(f(\mu))$, one can obtain a likelihood interval $[\mu_1, \mu_2]$ as the union of all μ for which

$$2\ln \mathcal{L}(\hat{\mu}) - 2\ln \mathcal{L}(\mu) \leq Z^2, \text{ for } Z \text{ real.}$$

Likelihood (Ratio) Intervals for 1 parameter

Recall: Likelihood function $\mathcal{L}(\mu)$ is invariant under reparametrization from μ to $f(\mu)$: $\mathcal{L}(\mu) = \mathcal{L}(f(\mu))$.

So *likelihood ratios* $\mathcal{L}(\mu_1) / \mathcal{L}(\mu_2)$ and *log-likelihood differences* $\ln \mathcal{L}(\mu_1) - \ln \mathcal{L}(\mu_2)$ are also invariant.

After using maximum-likelihood method to obtain estimate $\hat{\mu}$ that maximizes either $\mathcal{L}(\mu)$ or $\mathcal{L}(f(\mu))$, one can obtain a likelihood interval $[\mu_1, \mu_2]$ as the union of all μ for which

$$2\ln \mathcal{L}(\hat{\mu}) - 2\ln \mathcal{L}(\mu) \leq Z^2, \text{ for } Z \text{ real.}$$

As sample size increases (under important regularity conditions) this interval approaches a central confidence interval with C.L. corresponding to $\pm Z$ Gaussian standard deviations

But! Regularity conditions, in particular requirement that $\hat{\mu}$ not be on the boundary, need to be carefully checked.

E.g., if $\mu \geq 0$ on physical grounds, then $\hat{\mu} = 0$ requires care.

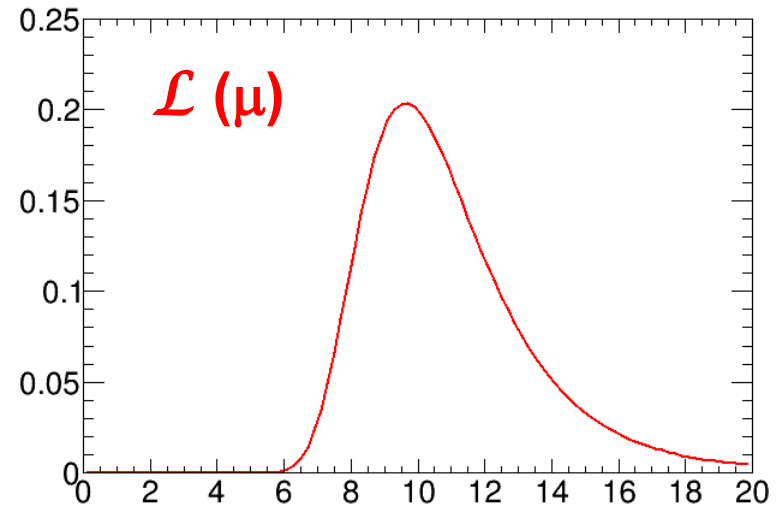
This is a special case of an important theorem by S.S. Wilks, to be discussed later in these lectures.

Gaussian pdf $p(x|\mu,\sigma)$ with σ a function of μ : $\sigma = 0.2 \mu$ Observed $x_0 = 10.0$.

Recall:

$\mathcal{L}(\mu)$ for observed $x_0 = 10.0$.

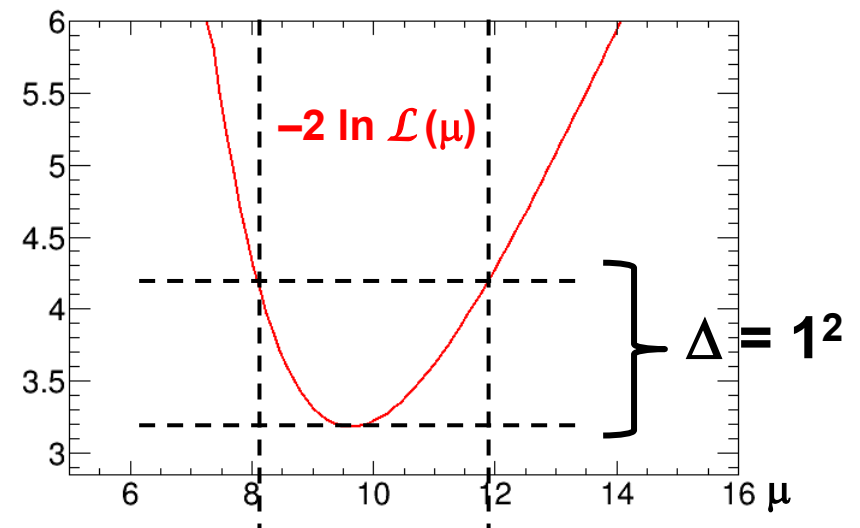
$\mu_{ML} = 9.63$



Likelihood ratio interval for μ at approximate 68% C.L.:

$[\mu_1, \mu_2] = [8.10, 11.9]$.

Compare with exact confidence interval $[8.33, 12.5]$.



Binomial Likelihood-Ratio Interval example

Recall example of $n_{\text{on}}=3$
successes in $n_{\text{tot}}=10$ trials.

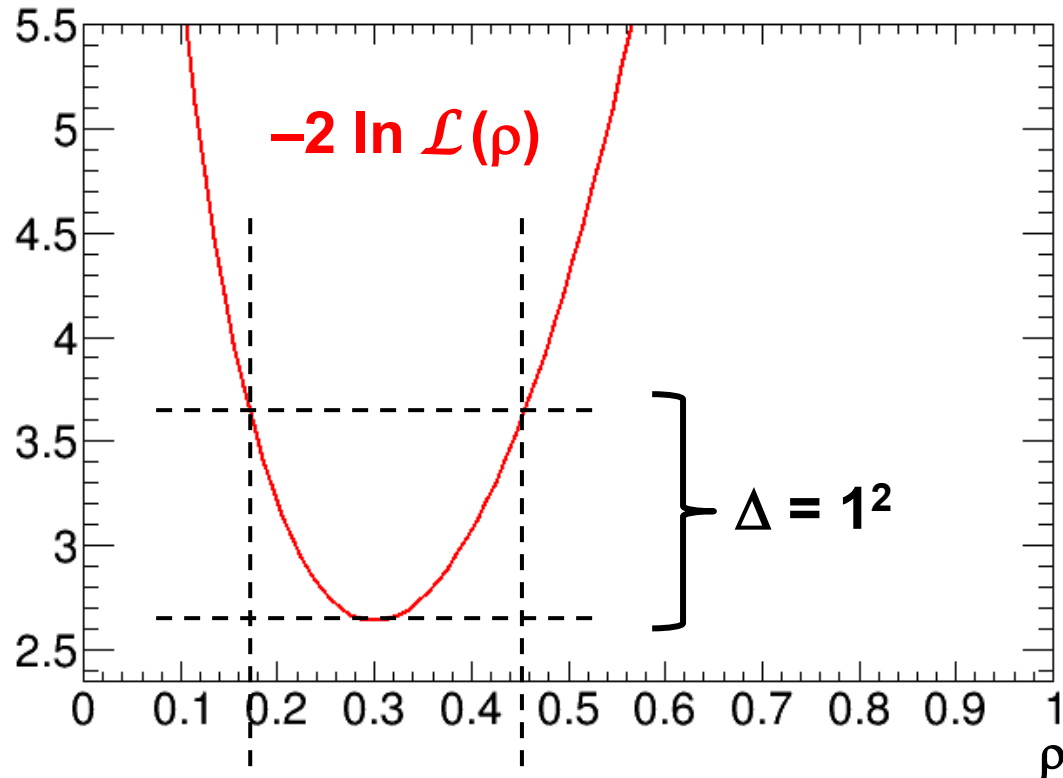
Minimum $-2 \ln \mathcal{L}(\rho) = 2.64$.

Obtain interval from

$-2 \ln \mathcal{L}(\rho) = 2.64 + 1 = 3.64$

\Rightarrow likelihood-ratio interval

$[\rho_1, \rho_2] = [0.17, 0.45]$



Also recall:

Copper-Pearson $[\rho_1, \rho_2] = [0.14, 0.51]$

Wilson $[\rho_1, \rho_2] = [0.18, 0.46]$

Poisson Likelihood-Ratio Interval example

Approx “68% C.L.” likelihood-ratio interval for Poisson process with $n=3$ observed:

$$\mathcal{L}(\mu) = \mu^3 \exp(-\mu)/3!$$

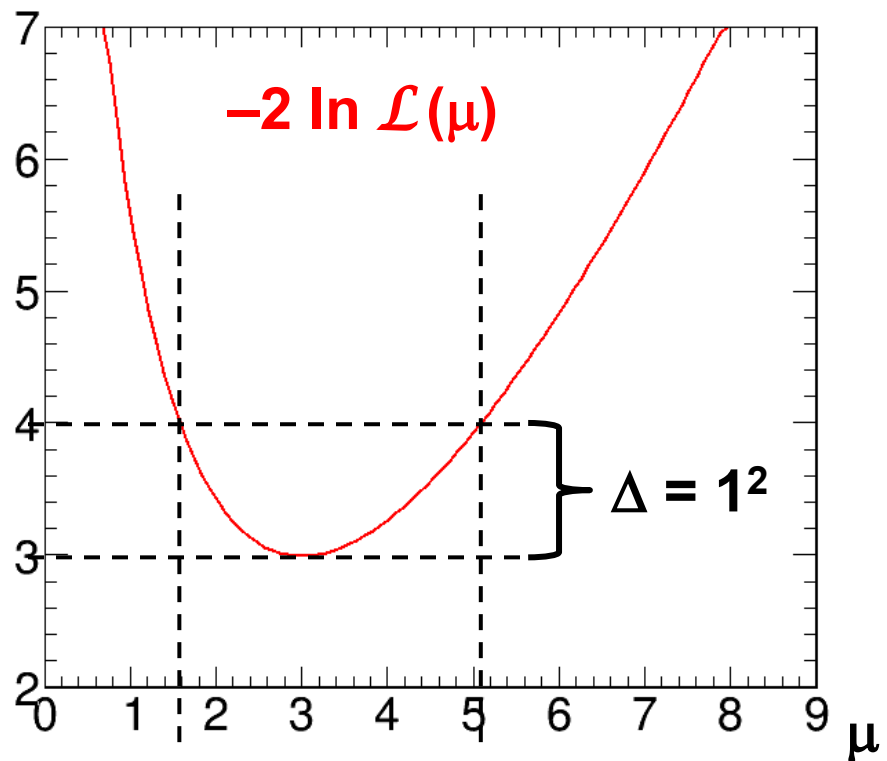
Recall maximum at $\mu = 3$.

$$-2 \ln \mathcal{L}(3) = 2.99$$

$\Delta 2 \ln \mathcal{L} = 1^2$ yields LR interval
 $[\mu_1, \mu_2] = [1.58, 5.08]$

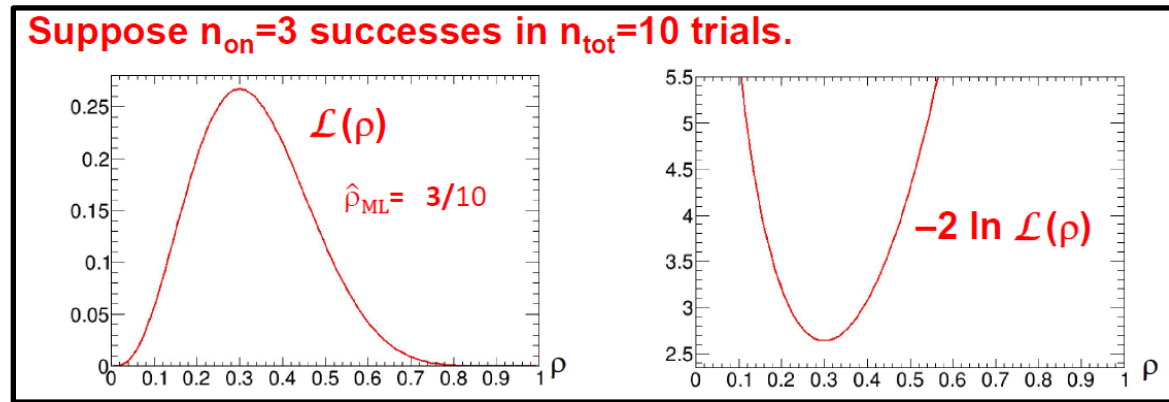
Neyman construction central (Garwood):

$$[\mu_1, \mu_2] = [1.37, 5.92]$$



Recall 3 methods of interval construction for binomial param ρ

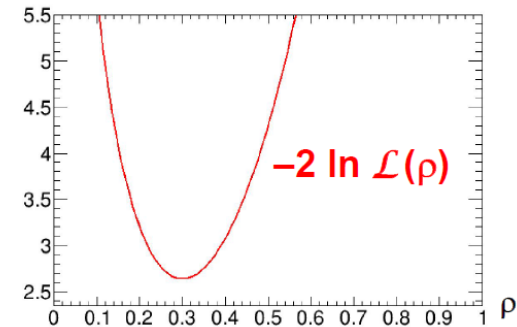
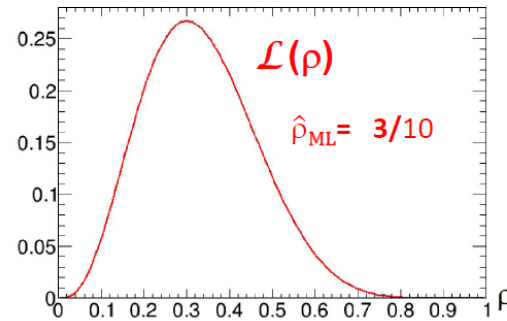
Bayesian and likelihood intervals: $\text{Bi}(n_{\text{on}} | n_{\text{tot}}, \rho)$ is evaluated *only* at observed $n_{\text{on}}=3$.



Recall 3 methods of interval construction for binomial param ρ

Bayesian and likelihood intervals: $\text{Bi}(n_{\text{on}} | n_{\text{tot}}, \rho)$ is evaluated only at observed $n_{\text{on}}=3$.

Suppose $n_{\text{on}}=3$ successes in $n_{\text{tot}}=10$ trials.



Confidence intervals use, in addition, probabilities for values of n_{on} not observed.

$n_{\text{on}} = 3, n_{\text{tot}}=10$.

Find ρ_1 such that

$\text{Bi}(n_{\text{on}} < 3 | \rho_1) = 84\%$

$\text{Bi}(n_{\text{on}} \geq 3 | \rho_1) = 16\%$

(lower limit at 84% C.L.)

Solve: $\rho_1 = 0.142$

And find ρ_2 such that

$\text{Bi}(n_{\text{on}} > 3 | \rho_2) = 84\%$

$\text{Bi}(n_{\text{on}} \leq 3 | \rho_2) = 16\%$

(upper limit at 84% C.L.)

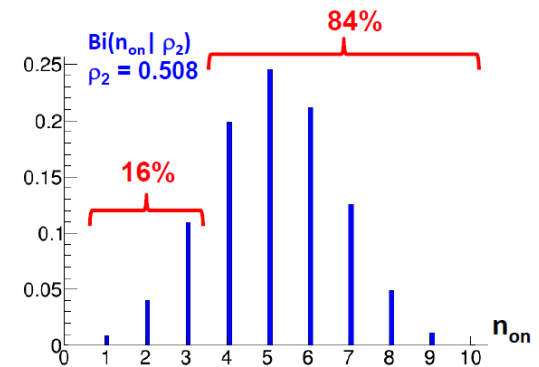
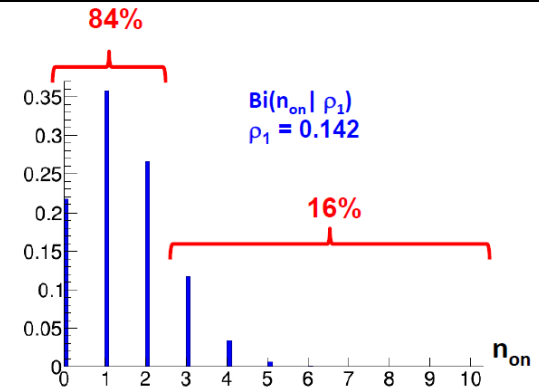
Solve: $\rho_2 = 0.508$

Then $[\rho_1, \rho_2] = (0.142, 0.508)$

is central confidence interval

with 68% C.L. Same as

Clopper and Pearson (1934)



Likelihood Principle

In both Bayesian methods and likelihood-ratio based methods, the probability (density) for obtaining the *data at hand* is used (via the likelihood function), *but probabilities for obtaining other data are not used!*

In contrast, in typical frequentist calculations (confidence intervals, p-values), one also uses probabilities of data that could have been observed but that was *not observed*.

The assertion that only the former is valid is captured by the *Likelihood Principle**:

If two experiments yield likelihood functions that are proportional, then Your inferences from the two experiments should be identical.

*There are various versions of the L.P., strong and weak forms etc. See Stuart99 and book by Berger and Wolpert.

Likelihood Principle (cont.)

L.P. is built into Bayesian inference (except e.g., when Jeffreys prior leads to violation).

L.P. is violated by p-values and confidence intervals.

Jeffreys (Theory of Probability, 1961, p. 385) still seems to be unsurpassed in his ironic criticism of tail probabilities, which include probabilities of data *more extreme* than that observed:

“What the use of [the p-value] implies, therefore, is that a hypothesis that may be true may be rejected because it has not predicted observable results that have not occurred.”

Although practical experience indicates that the L.P. may be too restrictive, it is useful to keep in mind. When frequentist results “make no sense” or “are unphysical”, in my experience the underlying reason can be traced to a bad violation of the L.P.

Likelihood Principle Example #1

The “Karmen Problem”

- You expect background events sampled from a Poisson distribution with mean $b=2.8$, assumed known precisely.
- For signal mean μ , the total number of events n is then sampled from Poisson mean $\mu+b$.
- So $P(n) = (\mu+b)^n \exp(-\mu-b)/n!$
- Then you observe no events at all! I.e., $n=0$.
- $\mathcal{L}(\mu) = (\mu+b)^0 \exp(-\mu-b)/0! = \exp(-\mu) \exp(-b)$

Likelihood Principle Example #1

The “Karmen Problem”

- You expect background events sampled from a Poisson distribution with mean $b=2.8$, assumed known precisely.
- For signal mean μ , the total number of events n is then sampled from Poisson mean $\mu+b$.
- So $P(n) = (\mu+b)^n \exp(-\mu-b)/n!$
- Then you observe no events at all! I.e., $n=0$.
- $\mathcal{L}(\mu) = (\mu+b)^0 \exp(-\mu-b)/0! = \exp(-\mu) \exp(-b)$

Note that changing b from 0 to 2.8 changes $\mathcal{L}(\mu)$ only by the constant factor $\exp(-b)$. This gets renormalized away in any Bayesian calculation, and is irrelevant for likelihood *ratios*.

So for zero events observed, likelihood-based inference about signal mean μ is *independent of expected b* .

For essentially all frequentist confidence interval constructions, the fact that $n=0$ is less likely for $b=2.8$ than for $b=0$ results in *narrower* confidence intervals for μ as b increases.

Clear violation of the L.P.

Likelihood Principle Example #2

Binomial problem, famous among statisticians, translated to HEP

You want to measure the efficiency ε of some trigger selection.

You count until reaching $n_{\text{tot}}=100$ zero-bias events, and note that of these, $m=10$ passed selection.

The probability for m is binomial with binomial parameter ε :

$$\text{Bi}(m \mid n_{\text{tot}}, \varepsilon) = \frac{n_{\text{tot}}!}{m! (n_{\text{tot}}-m)!} \varepsilon^m (1 - \varepsilon)^{(n_{\text{tot}}-m)}$$

Estimate $\varepsilon = 10/100$, compute binomial confidence interval for ε .

Also, plugging in the observed data, the likelihood function is

$$\mathcal{L}(\varepsilon) = \frac{100!}{10! 90!} \varepsilon^{10} (1 - \varepsilon)^{90}$$

Likelihood Principle Example #2 (cont.)

Your colleague *in a different experiment* counts zero-bias events until $m=10$ have passed her trigger. She notes that this requires $n_{\text{tot}}=100$ events (a coincidence).

Intuitively, $\varepsilon = 10/100$ over-estimates ε because she stopped *just* upon reaching 10 passed events, and indeed an unbiased estimate of ε and confidence interval will be slightly different from the binomial case.

Relevant distribution here is (a version of) the *negative binomial*:

$$\text{NBi}(n_{\text{tot}} | m, \varepsilon) = \frac{(n_{\text{tot}}-1)!}{(m-1)!} \varepsilon^m (1 - \varepsilon)^{(n_{\text{tot}}-m)}$$

Also, plugging in the observed data, the likelihood function is

$$\mathcal{L}(\varepsilon) = \frac{99!}{9! 90!} \varepsilon^{10} (1 - \varepsilon)^{90}$$

Likelihood Principle Example #2 (cont.)

So both you and your friend observed 10 successes out of 100 trials, but with different *stopping rules*.

Your likelihood function is based on *binomial* distribution:

$$\mathcal{L}(\varepsilon) = \frac{100!}{10! 90!} \varepsilon^{10} (1 - \varepsilon)^{90}$$

Your friend's is based on *negative binomial* distribution:

$$\mathcal{L}(\varepsilon) = \frac{99!}{9! 90!} \varepsilon^{10} (1 - \varepsilon)^{90}$$

The two likelihoods differ by (only) a constant factor, so the (strong) LP says that inferences should be *identical*.

In contrast, frequentist inferences use probabilities of data not obtained, and result in different confidence intervals and p-values for the different stopping rules.

Likelihood Principle Example #2 (cont.)

The two efficiency measurements had a different *stopping rules*: one stopped after n_{tot} events, and the other stopped after m events passing the trigger.

Frequentist confidence intervals depend on the stopping rule; the likelihood function did not (except for an overall constant).

So Bayesians will get the same answer in both cases, unless the *prior* depends on the stopping rule.

Amusing sidebar: the Jeffreys prior is indeed different for the two distributions, so use of Jeffreys prior violates (strong) L.P.

Stopping Rule Principle

The strong L.P. typically implies, as in this example, that the inference is independent of the stopping rule! This *irrelevance* has been elevated to the “Stopping Rule Principle”.

(It is sometimes amusing to ask a recent Bayesian convert if they know that they just bought the Stopping Rule Principle.)

Concepts that average/sum over the sample space, such as bias and tail probabilities, do not exist in pure Bayesian framework.

Stopping Rule Principle (cont.)

Famous quote by L.J. (Jimmie) Savage, early subjective Bayesian advocate:

“...I learned the stopping-rule principle from Professor Barnard, in conversation in the summer of 1952. Frankly, I then thought it a scandal that anyone in the profession could advance an idea so patently wrong, even as today I can scarcely believe that some people resist an idea so patently right.”

L.J. Savage et al.,

The Foundations of Statistical Inference: A Discussion

Methuen & Co., London, 1962

Scans of “The Complete Savage Forum” on D. Mayo’s web site,

http://www.phil.vt.edu/dmayo/PhilStatistics/supplementary_articles.htm

Summary of Three Ways to Make Intervals

	Bayesian Credible	Frequentist Confidence	Likelihood Ratio
Requires prior pdf?	Yes	No	No
Obeys Likelihood Principle?	Yes (exception re Jeffreys prior)	No	Yes
Random variable in “ $P(\mu_t \in [\mu_1, \mu_2])$ ”:	μ_t	μ_1, μ_2	μ_1, μ_2
Coverage guaranteed?	No	Yes (but over-coverage...)	No
Provides $P(\text{parameter} \text{data})$?	Yes	No	No

Frequentist intervals map to frequentist hypothesis tests, as previously discussed.

Bayesian approach to hypothesis testing is also called *model selection*, and is a whole other “can of worms” (J.O. Berger).

68% intervals for Poisson mean with $n=3$ observed

Method	Prior	Interval	Length	Coverage?
Wald, $n \pm \sqrt{n}$	—	(1.27, 4.73)	3.36	no
Garwood, Frequentist central	—	(1.37, 5.92)	4.55	yes
Bayesian central	1	(2.09, 5.92)	3.83	no
Bayesian central	$1/\mu$	(1.37, 4.64)	3.27	no
Bayesian central Jeffreys	$1/\sqrt{\mu}$	(1.72, 5.27)	3.55	no
Likelihood ratio	—	(1.58, 5.08)	3.50	no

Frequentist intervals over-cover due to discreteness of n .

68% intervals for Poisson mean with $n=3$ observed

Method	Prior	Interval	Length	Coverage?
Wald, $n \pm \sqrt{n}$	—	(1.27, 4.73)	3.36	no
Garwood, Frequentist central	—	(1.37, 5.92)	4.55	yes
Bayesian central	1	(2.09, 5.92)	3.83	no
Bayesian central	$1/\mu$	(1.37, 4.64)	3.27	no
Bayesian central Jeffreys	$1/\sqrt{\mu}$	(1.72, 5.27)	3.55	no
Likelihood ratio	—	(1.58, 5.08)	3.50	no

Bayesian lower limits with $1/\mu$ prior are identical to frequentist lower limits.

68% intervals for Poisson mean with n=3 observed

Method	Prior	Interval	Length	Coverage?
Wald, $n \pm \sqrt{(n)}$	—	(1.27, 4.73)	3.36	no
Garwood, Frequentist central	—	(1.37, 5.92)	4.55	yes
Bayesian central	1	(2.09, 5.92)	3.83	no
Bayesian central	$1/\mu$	(1.37, 4.64)	3.27	no
Bayesian central Jeffreys	$1/\sqrt{\mu}$	(1.72, 5.27)	3.55	no
Likelihood ratio	—	(1.58, 5.08)	3.50	no

Bayesian *upper* limits with flat prior are identical to frequentist upper limits.

Since *upper* limits dominate our field, this is why flat prior for Poisson mean became so well established: it is probability matching prior for upper limits, and when background is added, becomes conservative.

Bayesian Hypothesis Testing (Model Selection)

Typically follows Chapter 5 of Harold Jeffreys's book:
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 “Occam's razor”, etc.

Bayesian Hypothesis Testing (Model Selection)

Typically follows Chapter 5 of Harold Jeffreys's book: 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 “Occam's razor”, etc.

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

Here I mainly just say: Beware! There are posted/published applications HEP that lack foundation, in particular by Bayesian standards.

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

Bayesian Hypothesis Testing (Cont.)

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 becomes negligible.

For testing $H_0: \mu = \mu_0$ vs $H_1: \mu \neq \mu_0$, improper priors for μ that work fine for estimation become a disaster; adding cutoff to make them proper just gives (typically arbitrary) cutoff dependence.

Bayesian Hypothesis Testing (Cont.)

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> .

1D parameter space, 2D observation space

Until now we have considered 1 parameter and 1 observation.

Adding a second observation adds surprising subtleties.

As before, μ is parameter (often called θ by statisticians)

An experiment has two observations x_1, x_2 . These could be:

- two samples from same $p(x|\mu)$, or
- samples of two different quantities from joint density $p(x_1, x_2 | \mu)$.

Neyman construction:

For each μ , use an ordering principle on the sample space (x_1, x_2) to select an acceptance region $\mathcal{A}(\mu)$ in the sample space (x_1, x_2) such that $P((x_1, x_2) \in \mathcal{A}(\mu)) = \text{C.L.}$

In fact, this was the illustration in Neyman's original paper.

X—Outline of a Theory of Statistical Estimation Based on the Classical Theory of Probability

By J. NEYMAN

Reader in Statistics, University College, London

(Communicated by H. JEFFREYS, F.R.S.—Received 20 November, 1936—Read 17 June, 1937)

Original paper has one unknown parameter θ_1 on vertical axis and horizontal planes for 2D vectors of observables $\mathbf{E} = (x_1, x_2)$.

Prior to experiment, acceptance regions $\mathcal{A}(\theta_1)$ in E-space planes are determined for each θ_1 (needs ordering principle) with $P(\mathbf{E} \in \mathcal{A}(\theta_1)) = \text{C.L.}$

\mathbf{E}' is data actually observed in expt. Upon obtaining \mathbf{E}' , confidence interval for θ_1 consists of all values of θ_1 for which \mathbf{E}' is in $\mathcal{A}(\theta_1)$.

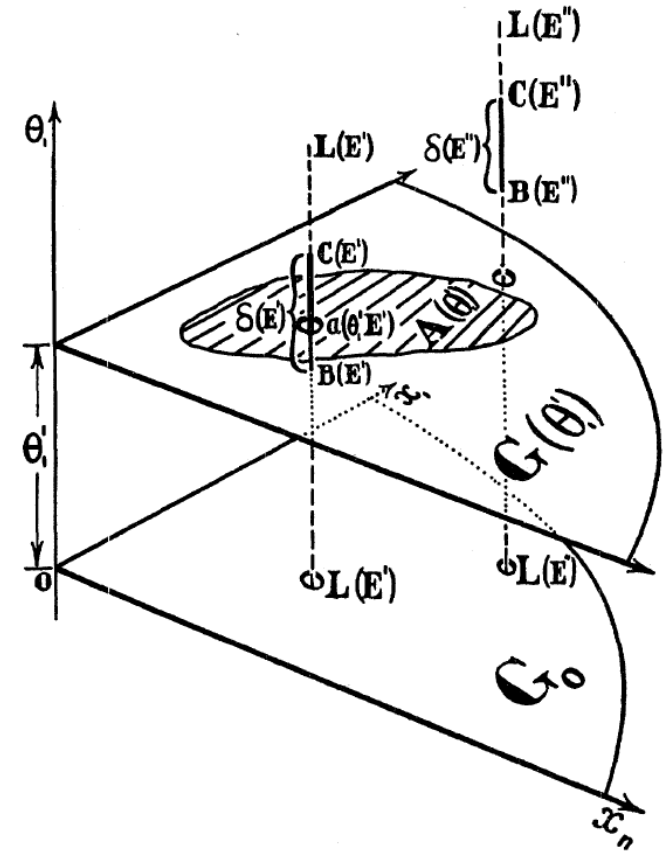


FIG. 1—The general space G .

Recall from Initial Discussion of Probability

For probabilities to be well-defined, the “whole space” needs to be defined. Can be hard for both frequentists and Bayesians!

[...]

Furthermore, it is widely accepted that restricting the “whole space” to a relevant (“conditional”) subspace can sometimes improve the quality of statistical inference. The important topic of such “conditioning” in frequentist inference will be discussed in detail later.

“Later” is now!

Restricting the Sample Space Used by Frequentists

In Neyman's construction in the 2D sample space (x_1, x_2) , the probabilities $P((x_1, x_2) \in \mathcal{A}(\mu))$ associated with each acceptance region $\mathcal{A}(\mu)$ are *unconditional* probabilities with respect to the "whole" sample space of all values of (x_1, x_2) .

In contrast, Bayesian inference is based on a single point in this sample space, the observed (x_1, x_2) , per the Likelihood Principle.

Restricting the Sample Space Used by Frequentists

In Neyman's construction in the 2D sample space (x_1, x_2) , the probabilities $P((x_1, x_2) \in \mathcal{A}(\mu))$ associated with each acceptance region $\mathcal{A}(\mu)$ are *unconditional* probabilities with respect to the “whole” sample space of all values of (x_1, x_2) .

In contrast, Bayesian inference is based on a single point in this sample space, the observed (x_1, x_2) , per the Likelihood Principle.

There can be a middle ground in frequentist inference, in which the probabilities $P((x_1, x_2) \in \mathcal{A}(\mu))$ are *conditional* probabilities conditioned on a function of (x_1, x_2) , in effect restricting the sample space to a “recognizable subset” depending on the observed data.

The function of (x_1, x_2) used for conditional probabilities typically carries information on the *uncertainty* in the point estimate $\hat{\mu}$, but no information on $\hat{\mu}$ itself: called an *ancillary statistic*.

Restricting the Sample Space Used by Frequentists (cont.)

Restricting the sample space in this way is known as *conditioning (on an ancillary)*. Two famous examples:

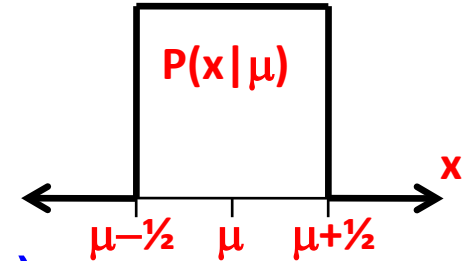
- 1) A somewhat artificial example of Welch where the conditioning arises from mathematical structure**
- 2) A more physical example of Cox where the argument for conditioning seems “obvious”**

Famous example of B.L. Welch (1939)

$$\text{Let } p(x|\mu) = \begin{cases} 1 & \text{if } \mu - 1/2 \leq x \leq \mu + 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

Two values x_1, x_2 are observed.

$$\hat{\mu} = \bar{x} = (x_1 + x_2)/2 \text{ (Only for } n=2! \text{ See Backup)}$$



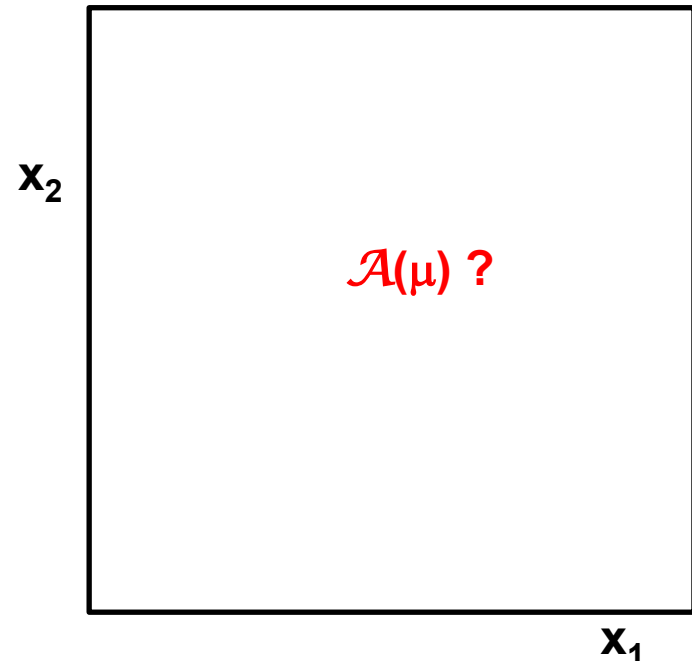
What is 68% C.L. central confidence interval for μ ?

**Neyman construction: Define acceptance region $\mathcal{A}(\mu)$ containing 68% of unit square of (x_1, x_2) centered on μ .
What to use?**

Centrality implies symmetry.

Need something else to rank points in the plane.

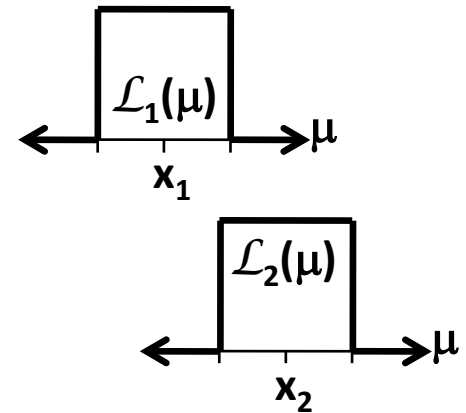
N-P Lemma gives most powerful ranking, but first let's think about some examples.



“Lucky” sample with $|x_1 - x_2|$ close to 1.

$\mathcal{L}(\mu) = \mathcal{L}_1(\mu) \times \mathcal{L}_2(\mu)$ very narrow.

Reasonable to expect small uncertainty in $\hat{\mu}$?



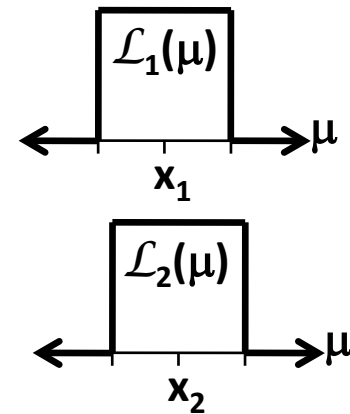
“Unlucky” sample, $|x_1 - x_2|$ close to 0.

$\mathcal{L}(\mu)$ full width close to 1;

Second observation added no useful info.

Expect 68% C.L. conf. interval 0.68 long?

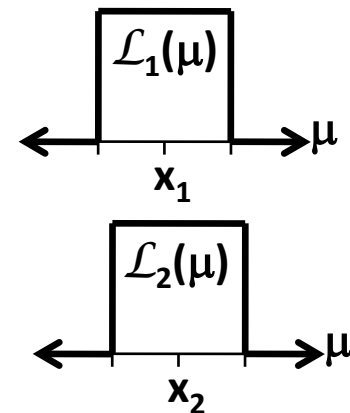
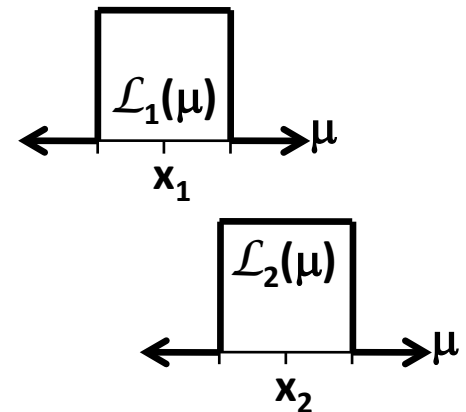
Guess reasonable answer: conf. interval centered on $\hat{\mu}$ with length $0.68(1 - |x_1 - x_2|)$



Seems reasonable for *post-data uncertainty* to depend on $|x_1 - x_2|$.

$|x_1 - x_2|$ is classic example of an ancillary statistic A: has info on *uncertainty* on μ estimate, but no info on μ itself.

Idea dating to Fisher and before: divide the full “*unconditional*” sample space into “recognizable subsets” and report probs using the “relevant” subset rather than the whole space.



See backup slides: Confidence intervals $\hat{\mu} \pm 0.34(1 - |x_1 - x_2|)$, as thought reasonable! Known as “conditioning” on ancillary statistic A: Post-data, proceed as if A had been fixed, rather than randomly sampled!...even though less power!

Brad Efron's talk at PhyStat-2003 gave similar example using Cauchy distribution, where ancillary statistic is curvature of $\mathcal{L}(\mu)$. He called conditional answer "correct".

(<http://www.slac.stanford.edu/econf/C030908/papers/MOAT003.pdf>)

Summary: conditioning on an ancillary statistic A means:

Even though A was randomly sampled in the experimental procedure, after data is obtained proceed as if A had been fixed to the value observed. Ignore sample space with all those other values of A that you could have obtained, but did not.

Famous “weighing machine” example of D.R. Cox (1958)

For measuring the mean μ with Gaussian resolution, one of two devices is selected randomly with equal probability:

- Device #1 with σ_1
- Device #2 with $\sigma_2 \ll \sigma_1$.

So in my notation, x_1 is the index (1 or 2) chosen randomly and specifying device, and x_2 is the single sample from the selected Gaussian measurement.

In Behnke13 (p. 122) Luc Demortier gives an example in HEP: μ is mass of decaying particle with probability p_h to decay hadronically (mass resolution σ_1) and probability $1-p_h$ to decay leptonically with mass resolution σ_2). Thus the “measuring machine” chosen randomly is the detector used to measure the decay mode chosen by QM.

Example of D.R. Cox (cont.)

- Device #1 with σ_1
- Device #2 with $\sigma_2 \ll \sigma_1$.

x_1 is the index (1 or 2) chosen randomly and specifying device;
 x_2 is the single sample from the selected measurement.

So $\hat{\mu} = x_2$. What is the confidence interval?

Example of D.R. Cox (cont.)

- Device #1 with σ_1
- Device #2 with $\sigma_2 \ll \sigma_1$.

x_1 is the index (1 or 2) chosen randomly and specifying device;
 x_2 is the single sample from the selected measurement.

So $\hat{\mu} = x_2$. What is the confidence interval?

The index x_1 is an *ancillary statistic* (gives info on uncertainty but not on point estimate), and it is reasonable (obvious?) to condition on it. I.e., we report confidence interval giving correct coverage in *subspace of measurements that used the same device we used*.

So 68% C.L. confidence interval:

- $\hat{\mu} \pm \sigma_1$ if Device #1 randomly selected
- $\hat{\mu} \pm \sigma_2$ if Device #2 randomly selected

As in Welch example, it turns out that more powerful tests (confidence intervals shorter on average!) can be found.

Example of D.R. Cox (cont.)

Demortier gives details on how average length of intervals optimized in the unconditional sample space is shorter in the HEP example. See backup for details of Cox's example.

In both cases, the idea is to undercover badly (~38% coverage for 68% C.L. intervals) when Device #1 is used, and overcover (~100% coverage) when Device #2 is used. The smaller σ of Device #2 means that average length of interval goes down, *when averaging over the entire unconditional sample space.*

One gives up power with Device #1 and uses it in Device #2.

Cox: "If, however, our object is to say 'what can we learn from the data that we have', the unconditional test is surely no good."

The Welch and Cox (and Efron) examples reveal a real conflict between N-P optimization for power and conditioning to optimize relevance.

Conditionality Principle

The assertion that inference should be conditioned on an ancillary in the Welch example (where it comes out of the math) is often called the “**Conditionality Principle**”.

Conditioning in the Cox example (a “mixture experiment” where the ancillary has physical meaning about which experiment was performed) is then the “**Weak Conditionality Principle**”.

Conditionality Principle (cont.)

But note: in sufficiently complicated cases (for example if there is more than one ancillary statistic), the procedure is less clear.

In many situations, ancillary statistics do not exist, and it is not at all clear how to restrict the “whole space” to the relevant part for frequentist coverage.

The Bayesian answer: restrict the whole space to the point observed! But the price is giving up coverage.

When there are “recognizable subsets” with varying coverage, Buehler has discussed how a “conditional frequentist” can win bets against an “unconditional frequentist” – see Backup.

Conditioning in HEP

A classic example is a *measurement of the branching fraction of a particular decay mode* when the *total* number of decays N can fluctuate because the experiment design is to run for a fixed length of time. Then N is an ancillary statistic.

You perform an experiment and obtain N total decays, and then do a toy M.C. of repetitions of the experiment. Do you let N fluctuate, or do you fix it to the value observed?

It may seem that the toy M.C. should include your *complete* procedure, including fluctuations in N .

Conditioning in HEP

A classic example is a *measurement of the branching fraction of a particular decay mode* when the *total* number of decays N can fluctuate because the experiment design is to run for a fixed length of time. Then N is an ancillary statistic.

You perform an experiment and obtain N total decays, and then do a toy M.C. of repetitions of the experiment. Do you let N fluctuate, or do you fix it to the value observed?

It may seem that the toy M.C. should include your *complete* procedure, including fluctuations in N .

But the above arguments would point toward *conditioning on the value of the ancillary statistic actually obtained*. So your branching fraction measurement is binomial with trials N .

(Originally discussed in HEP by F. James and M. Roos, Nucl. Phys. B 172 (1980) 475. For more complete discussion, see Cousins, Hyme, Tucker, <https://arxiv.org/abs/0905.3831>)

Review: Conditioning and the Likelihood Principle

We have seen that *unconditional frequentists* compute probabilities with respect to the *whole sample space*.

Post-data, *conditional frequentists* try to refer to a *relevant subset* of the whole sample space (typically not easy).

We also saw that pure Bayesians refer only to the probability of the data observed (L.P.). This is literally the ultimate extreme in conditioning, *conditioning (in the continuous case) on a point of measure zero!* (You can't get any more "relevant".)

This is why coverage is not built into Bayesian answers.

Conditioning and the Likelihood Principle (cont.)

It is not surprising that Bayesians argue for the importance of *relevance* of the inference, and criticize frequentists for danger of irrelevance (and difficulty of diagnostic of irrelevance).

And it is not surprising that frequentists argue for the importance of a useful measure of “error rates”, in the sense of Type 1 and Type 2 errors, coverage, etc., which may at best be estimates if L.P. is observed.

(Finally!) More than One Parameter

Generalize to two parameters μ_1 and μ_2 , true values unknown.

Let data x be a multi-D vector, so the model is $p(x|\mu_1, \mu_2)$.

Observed vector value is x_0 .

First consider the desire to obtain a 2D confidence/credible *region* in the parameter space (μ_1, μ_2) . All three methods discussed for intervals handle this in a straightforward (in principle) generalization:

- **Bayesian:** put observed data vector x_0 into $p(x|\mu_1, \mu_2)$ to obtain the likelihood function $\mathcal{L}(\mu_1, \mu_2)$. Multiply by prior pdf $p(\mu_1, \mu_2)$ to obtain 2D posterior pdf $p(\mu_1, \mu_2|x_0)$. Use posterior pdf to obtain credible regions, etc., in (μ_1, μ_2) .

(Finally!) More than One Parameter

Generalize to two parameters μ_1 and μ_2 , true values unknown.

Let data x be a multi-D vector, so the model is $p(x|\mu_1, \mu_2)$.

Observed vector value is x_0 .

First consider the desire to obtain a 2D confidence/credible region in the parameter space (μ_1, μ_2) . All three methods discussed for intervals handle this in a straightforward (in principle) generalization:

- **Bayesian: put observed data vector x_0 into $p(x|\mu_1, \mu_2)$ to obtain the likelihood function $\mathcal{L}(\mu_1, \mu_2)$. Multiply by prior pdf $p(\mu_1, \mu_2)$ to obtain 2D posterior pdf $p(\mu_1, \mu_2|x_0)$. Use posterior pdf to obtain credible regions, etc., in (μ_1, μ_2) .**
- **Confidence intervals: perform Neyman construction: Find acceptance regions for x as a function of (μ_1, μ_2) . The 2D confidence region is union of all (μ_1, μ_2) for which x_0 is in acceptance region.**

(Finally!) More than One Parameter

Generalize to two parameters μ_1 and μ_2 , true values unknown.

Let data x be a multi-D vector, so the model is $p(x|\mu_1, \mu_2)$.

Observed vector value is x_0 .

First consider the desire to obtain a 2D confidence/credible region in the parameter space (μ_1, μ_2) . All three methods discussed for intervals handle this in a straightforward (in principle) generalization:

- **Bayesian:** put observed data vector x_0 into $p(x|\mu_1, \mu_2)$ to obtain the likelihood function $\mathcal{L}(\mu_1, \mu_2)$. Multiply by prior pdf $p(\mu_1, \mu_2)$ to obtain 2D posterior pdf $p(\mu_1, \mu_2|x_0)$. Use posterior pdf to obtain credible regions, etc., in (μ_1, μ_2) .
- **Confidence intervals:** perform Neyman construction: Find acceptance regions for x as a function of (μ_1, μ_2) . The 2D confidence region is union of all (μ_1, μ_2) for which x_0 is in acceptance region.
- **Likelihood regions:** recall 1D method $2\ln\mathcal{L}(\hat{\mu}) - 2\ln\mathcal{L}(\mu) \leq Z^2 \dots$

Likelihood ratio regions in $\geq 2D$

(Recall: *differences* in $\ln\mathcal{L}$ correspond to *ratios* of likelihoods)

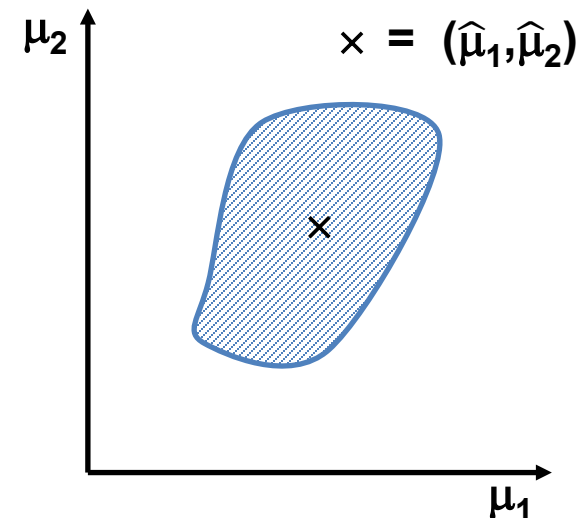
Find global maximum of $\mathcal{L}(\mu_1, \mu_2)$, yielding point estimates $(\hat{\mu}_1, \hat{\mu}_2)$.

Find contour bounded by $2\Delta\ln\mathcal{L} = 2\ln\mathcal{L}(\hat{\mu}_1, \hat{\mu}_2) - 2\ln\mathcal{L}(\mu_1, \mu_2) \leq C$, where C comes from Wilks's Theorem, tabulated in PDG RPP:

Table 39.2: Values of $\Delta\chi^2$ or $2\Delta\ln L$ corresponding to a coverage probability $1 - \alpha$ in the large data sample limit, for joint estimation of m parameters.

$(1 - \alpha)$ (%)	$m = 1$	$m = 2$	$m = 3$
68.27	1.00	2.30	3.53
90.	2.71	4.61	6.25
95.	3.84	5.99	7.82
95.45	4.00	6.18	8.03
99.	6.63	9.21	11.34
99.73	9.00	11.83	14.16

2D joint $\approx 68\%$ C.L.
 $2\Delta\ln\mathcal{L} = 2.3$



<http://pdg.lbl.gov/2018/reviews/rpp2018-rev-statistics.pdf>

As in 1D, Wilks's Theorem is asymptotic (large N) result, with various "regularity conditions" to be satisfied.

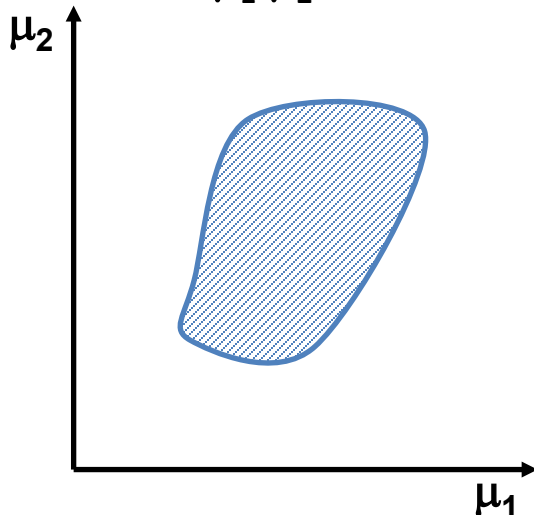
Nuisance Parameters

Frequently one is interested in considering one parameter at a time, irrespective of the value of other parameter(s).

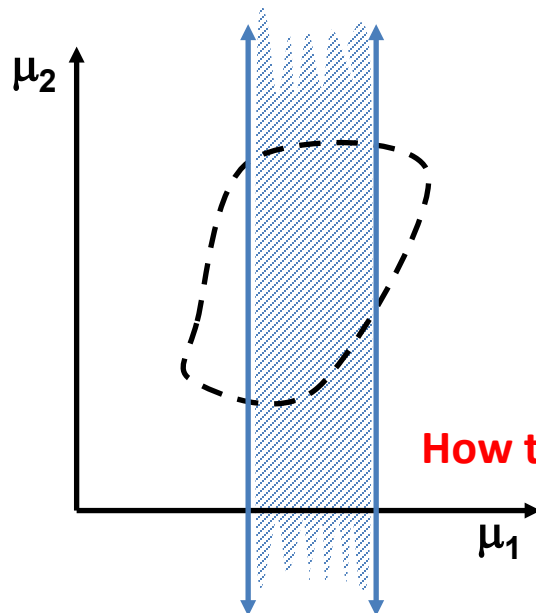
The parameter under consideration at the moment is called the “parameter of interest” and the other parameters (at that moment) are called “nuisance parameters”.

E.g., if μ_1 is of interest and μ_2 is a nuisance, one seeks a 2D confidence region that is a vertical “stripe” in the (μ_1, μ_2) plane.

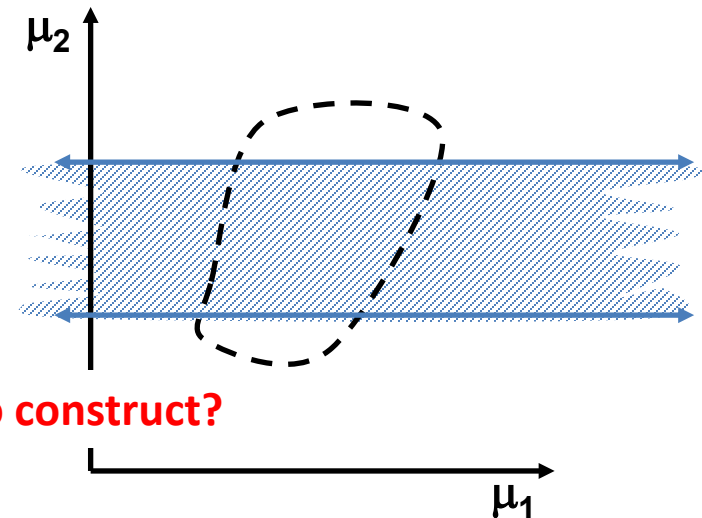
2D joint $\approx 68\%$ C.L. confidence region for (μ_1, μ_2)



2D region to get 1D $\approx 68\%$ C.L. interval for μ_1 (μ_2 is nuisance)



2D region to get 1D $\approx 68\%$ C.L. interval for μ_2 (μ_1 is nuisance)



How to construct?

Aside: Systematic Uncertainties as Nuisance Parameters

I have begun with just one parameter of interest and one nuisance parameter. Analyses in HEP can have hundreds or even thousands of nuisance parameters.

A typical measurement in HEP has many subsidiary measurements of quantities not of direct physics interest, but which enter into the calculation of the physics quantity of particular interest.

Aside: Systematic Uncertainties as Nuisance Parameters

I have begun with just one parameter of interest and one nuisance parameter. Analyses in HEP can have hundreds or even thousands of nuisance parameters.

A typical measurement in HEP has many subsidiary measurements of quantities not of direct physics interest, but which enter into the calculation of the physics quantity of particular interest.

E.g., if an absolute cross section is measured, one will have uncertainty in the integrated luminosity L , in the background level b , the efficiency e of detecting the signal, etc. In HEP, we call these systematic uncertainties, but statisticians (for the obvious reason) refer to L , b , and e as *nuisance parameters*.

Each of the three main classes of constructing intervals (Bayesian, Neyman confidence, likelihood ratio) has a “native” way to incorporate the uncertainty on the nuisance parameters.

But this remains a subject of frontier statistics research.

Nuisance Parameters I: Bayesian Credible Intervals

Construct a multi-D prior pdf $p(\text{parameters})$ for the space spanned by all parameters.

Multiply it by $\mathcal{L}(\text{data}|\text{parameters})$ for the data obtained to yield multi-D posterior pdf.

Integrate over the full subspace of all nuisance parameters **(marginalization)**.

Thus obtain posterior pdf for the parameter of interest.

Math is reduced to the case of no nuisance parameters.

Nuisance Parameters I: Bayesian Credible Intervals

Construct a multi-D prior pdf $p(\text{parameters})$ for the space spanned by all parameters.

Multiply it by $\mathcal{L}(\text{data}|\text{parameters})$ for the data obtained to yield multi-D posterior pdf.

Integrate over the full subspace of all nuisance parameters (**marginalization**).

Thus obtain posterior pdf for the parameter of interest.

Math is reduced to the case of no nuisance parameters.

Problems: The multi-D prior pdf is a problem for both subjective and non-subjective priors. In HEP there has been little use of the favored non-subjective priors (reference priors of Bernardo and Berger). The high-D integral can be a technical problem, more and more overcome by Markov Chain Monte Carlo.

As with all Bayesian analyses, how to interpret probability if default priors are used?

Nuisance Parameters II: Neyman Construction

For each point in the subspace of nuisance parameters, treat them as fixed true values and perform a Neyman construction for multi-D confidence regions in the full space of all parameters. Project these regions onto the subspace of the parameter of interest.

Nuisance Parameters II: Neyman Construction

For each point in the subspace of nuisance parameters, treat them as fixed true values and perform a Neyman construction for multi-D confidence regions in the full space of all parameters. Project these regions onto the subspace of the parameter of interest.

Problem: Typically intractable and causes overcoverage, and therefore rarely attempted.

Tractability can sometimes be recovered by doing the construction in the lower dimensional space of the profile likelihood function, obtaining approximate coverage.

(This is one way to interpret the Kendall and Stuart page on likelihood ratio test with nuisances.)

Typically “elimination” is done in a way technically feasible, including parametric bootstrap later in lectures, and the coverage is studied with toy Monte Carlo.

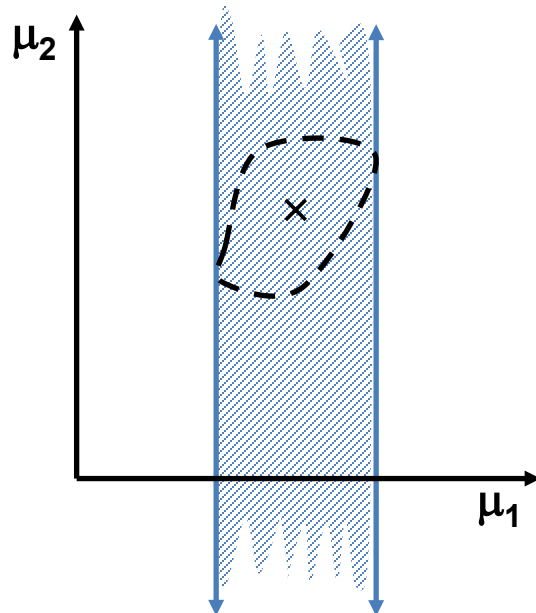
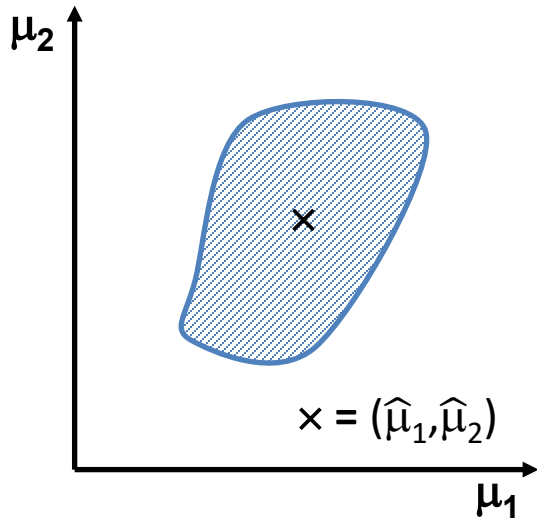
Nuisance Parameters III: Likelihood Ratio intervals

Many of us raised on MINUIT MINOS read F. James, “Interpretation of the Shape of the Likelihood Function around Its Minimum,” *Computer Physics Communications* 20 (1980) 29.

Whereas 2D region has $m=2$ and hence $2\Delta\ln\mathcal{L} = 2.3$ from PDG RPP table, for 1D interval on μ_1 , we first make *2D contour* with $m=1$ value, $2\Delta\ln\mathcal{L} = 1$ (black dashed), and then find extrema in μ_1 :

2D joint $\approx 68\%$ C.L.
 $2\Delta\ln\mathcal{L} = 2.3$

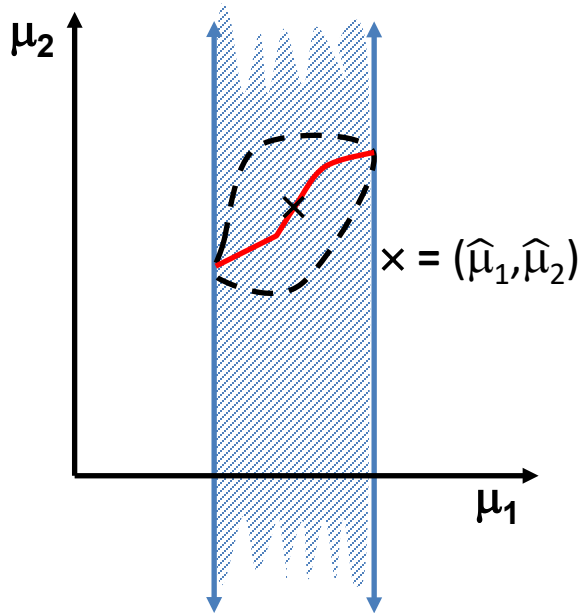
2D region to get 1D
 $\approx 68\%$ C.L. interval for μ_1
 $2\Delta\ln\mathcal{L} = 1$ (dashed)



...and then at the Fermilab Confidence Limits Workshop in 2000, statistician Wolfgang Rolke expressed the construction a different way:

Profile Likelihood Function

2D region to get 1D
68% C.L. interval for μ_1
 $2\Delta\ln\mathcal{L} = 1$ (dashed)



Red curve is path $(\mu_1, \hat{\mu}_2)$
along which profile \mathcal{L} is
evaluated

For each μ_1 , find the value $\hat{\mu}_2$ that
minimizes $-2\ln\mathcal{L}(\mu_1, \hat{\mu}_2)$, shown in red.

Make 1D plot vs μ_1 of this “profile
likelihood function” (of only μ_1).

Use the $m=1$ threshold on $2\Delta\ln\mathcal{L}_{\text{profile}}$.

One obtains the exact same interval as
“MINOS” on the left? Can you see why?

Since 2000, the “profile” terminology
has permeated HEP. The notation is also
used in the “Kendall and Stuart” page
that I showed re F-C.

Warning: Combining profile likelihoods
from two experiments is unreliable.
Apply profiling after combining full
likelihoods.

Likelihood Ratio intervals (cont.)

Problems:

Coverage is not guaranteed, particularly at low N . By using best-fit value of the nuisance parameters corresponding to each value of the parameter of interest, this has an (underserved?) reputation for underestimating the true uncertainties.

In Poisson problems, this is partially compensated by effect due to discreteness of n , and profile likelihood (MINUIT MINOS) gives good performance in many problems. See Rolke et al., NIM A551 (2005) 493.

In some cases (for example when there are spikes in \mathcal{L}), marginalization may give better frequentist performance, according to statisticians.

Later I will talk about the parametric bootstrap to construct more accurate intervals.

Aside: Profile Likelihood Ratio (PLR) test statistic

In more general notation, let μ be the parameter of interest and ν be a vector of nuisance parameters, so the *profile likelihood function* of μ is

$$\mathcal{L}(\mu, \hat{\nu}), \text{ also written as: } \sup_{\nu} \mathcal{L}(\mu, \nu) .$$

A useful quantity for hypothesis testing (and hence for confidence interval construction) is this *profile likelihood function* divided by the *global maximum likelihood*, thus obtaining the *profile likelihood ratio* test statistic,

$$\Lambda = \mathcal{L}(\mu, \hat{\nu}) / \mathcal{L}(\hat{\mu}, \hat{\nu}).$$

We frequently use $-2\ln\Lambda$ as a test statistic.

A Bayesian-inspired alternative to the profile likelihood function, discussed later, is the “*integrated*” (or “*marginalized*”) likelihood,

$$\int \mathcal{L}(\mu, \nu) \pi(\nu) d\nu,$$

where $\pi(\nu)$ is a weight function in the spirit of a prior pdf for ν .

Evaluation of coverage with toy MC

For a single parameter of interest μ , after elimination of nuisance parameters by some approximate method and construction of intervals perhaps involving more approximations, one reports as usual the confidence interval $[\mu_1, \mu_2]$ at some C.L.

It is important to check that the approximations in the whole procedure have not materially altered the claimed coverage:

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

Typically the performance is evaluated with toy MC.

First I describe the most thorough check (very CPU intensive), and then some approximations.

In frequentist statistics, the true values of all parameters are typically *fixed but unknown*. A complete, rigorous check of coverage considers a fine multi-D grid of *all* parameters, and *for each multi-D point in the grid*, generates an ensemble of toy MC pseudo-experiments, runs the full analysis procedure, and finds the fraction of intervals covering μ_t used for that ensemble, i.e., $P(\mu_t \in [\mu_1, \mu_2])$, and compares to C.L.

Evaluation of coverage with toy MC (cont.)

Thus a thorough check of frequentist coverage includes:

- 1) Fix all parameters (of interest and nuisance) to a single set of true values. For this set,
 - a) Loop over “pseudo-experiments”
 - b) For each pseudo-experiment, loop over events, generating each event with toy data generated from the statistical model with parameters set equal to the fixed set.
 - c) Perform the same analysis on the toy events in the pseudo-experiment as was done for the real data.
 - d) Find that fraction of the pseudo-experiments for which parameter(s) of interest are included in stated confidence intervals or regions.
- 2) *Repeat for a various other fixed sets of all parameters. But...the ideal of a fine grid is usually impractical.*

So the issue is what selection of “various other fixed sets” is adequate

Evaluation of coverage with toy MC (cont.)

The ideal of a fine grid is usually impractical.

So the issue is what selection of “various other fixed sets” is adequate.

Obviously one should check coverage if the set of true values is equal to the global best-fit values.

Just as obviously, this may not be adequate. Some exploration is needed, particularly in directions where uncertainty on a parameter depends strongly on the parameter. One can start by varying a few critical parameters by one or two s.d., trying parameters near boundary, and seeing how stable coverage is.

A Bayesian-inspired approach is to calculate a weighted average of coverage over a neighborhood parameter sets for the nuisance parameters. This requires a choice of multi-D prior (Recall problems.) Instead of fixing the true values of nuisance parameters during the toy-tossing, one samples the true parameters from the posterior distribution.

Constructing the intervals with toy MC: the *parametric bootstrap*

This approach is very common at the LHC. It generally improves on the profile likelihood ratio intervals described above. I think that it is best explained using the nested-hypothesis-test view in the duality of hypothesis tests and confidence intervals.

So, we consider the test of a particular value of the parameter of interest μ , and for that fixed value (and using the data observed), we find the best-fit values of all of the nuisance parameters.

We generate toy MC with these *fixed* values, and construct the distribution of the test statistic (typically a *profile likelihood ratio*). For a given C.L. = $1-\alpha$, we see if the μ being tested is rejected or not. By testing various values of μ , we construct the confidence interval for μ .

Hybrid Techniques: Introduction to Pragmatism

Given the difficulties with all three classes of interval estimation, especially when incorporating nuisance parameters, it can sometimes be useful to relax foundational rigor and:

- Treat nuisance parameters in a Bayesian way (marginalization) while treating the parameter of interest in a frequentist way.

Virgil Highland and I were early advocates of this for lumi uncertainty in upper limit calculation (NIM A320 (1992) 331).

Kyle Cranmer exposed problems when used for background mean in 5σ discovery context.

For review of background case and connection to Box's semi-Bayesian “prior predictive p-value”, see NIM A595 (2008) 480, <https://arxiv.org/abs/physics/0702156>

Introduction to Pragmatism (cont.)

- Use the Bayesian framework (even without the priors recommended by statisticians), but evaluate the frequentist performance. In effect (as in profile likelihood) one gets approximate coverage while respecting the L.P. In fact, the statistics literature going back to 1963 has attempts to find prior pdfs that lead to posterior pdfs with good frequentist coverage: *probability matching priors*. (At lowest order in 1D, it is the Jeffreys prior!)

Studies of coverage with toy MC

Numerous studies have been done for elimination of nuisances parameters in the test statistic, many concluding that results are relatively insensitive to profiling vs marginalization, so that choice can be made based on CPU time.

See for example John Conway's talk and writeup at PhyStat-2011, <https://indico.cern.ch/event/107747/timetable/#all.detailed>

Anecdotally, the choice for toys is more important than the choice for test statistic.

Looks at the statistics literature on nuisances

Almost 20 years ago, Luc Demortier and I both looked in the statistics literature regarding nuisance parameters – I thought my note was fairly thorough until I read his! Our writeups:

R.D. Cousins, Treatment of Nuisance Parameters in High Energy Physics, and Possible Justifications and Improvements in the Statistics Literature, with response by statistician Nancy Reid, <http://www.physics.ox.ac.uk/phystat05/proceedings/default.htm>

Luc Demortier, P Values: What They Are and How to Use Them, <https://www-cdf.fnal.gov/~luc/statistics/cdf8662.pdf> (2007)

See also his chapter in Behnke13.

Recently, statistician Larry Wasserman and I posted an informal review of marginalizing vs profiling: <https://arxiv.org/abs/2404.17180>

Priors for nuisance parameters

It used to be (unfortunately) common practice to express, say, a 50% systematic uncertainty on a positive quantity as a Gaussian with 50% rms. Then “truncate” Gaussian by not using non-positive values.

In Bayesian calculation, interaction of uniform prior for Poisson mean and “truncated Gaussian” for systematic uncertainty in efficiency leads to integral that diverges if truncation is at origin...evaluating integral numerically will not even notice! Luc Demortier exposed this:

<http://www.ippp.dur.ac.uk/old/Workshops/02/statistics/proceedings//demortier.pdf>

Recommendation: Use log-normal or (certain) Gamma distributions instead of truncated Gaussian. Recipes at
http://www.physics.ucla.edu/~cousins/stats/cousins_lognormal_prior.pdf

“State of the Art” in HEP

All three main classes of methods are commonly used on the parameter of interest.

- Both marginalization and profiling are used to treat nuisance parameters. At present, I think that profiling is much more common at the LHC.
- As mentioned, the parametric bootstrap is common.
- Many people have the good practice of checking coverage.
- Too little attention is given to priors. But flat prior for Poisson mean is “safe” for *upper* limits (only!).

A serious analysis using any of the main methods requires coding up the likelihood function.

- Doing this (once!) with RooFit modeling language gives access to RooStats techniques for all three classes of calculations, mix-match nuisance parameter treatments.

ATLAS and CMS Conventions

For many years, ATLAS and CMS physicists have collaborated on statistics tools (the RooStats software), and attempted to have some coherence in methods so that results could be compared, and (when worth the effort) combined.

An important development was the paper by Cowan, Cranmer, Gross and Vitells, <https://arxiv.org/abs/1007.1727>, that extended asymptotic formulas to a number of cases where Wilks's theorem was not valid.

As the CCGV asymptotic formulas applied to the “fully frequentist” treatment of nuisance parameters, for consistency we tended to use that in many cases at small N as well. Toy MC is thus done in a frequentist manner (parametric bootstrap).

For upper limits, there was a lot of discussion without convergence, and the two physics coordinators in 2010 decreed that CL_s be used in most cases. (See below.)

ATLAS and CMS Conventions (cont.)

The ATLAS/CMS Higgs results followed these trends. A jointly written description is, “Procedure for the LHC Higgs boson search combination in Summer 2011,” <http://cds.cern.ch/record/1379837>.

Many issues were further discussed and described in the ATLAS-CMS combination papers for mass and couplings, <https://arxiv.org/abs/1503.07589> and <https://arxiv.org/abs/1606.02266> . In particular, a lot of attention was paid to correlations.

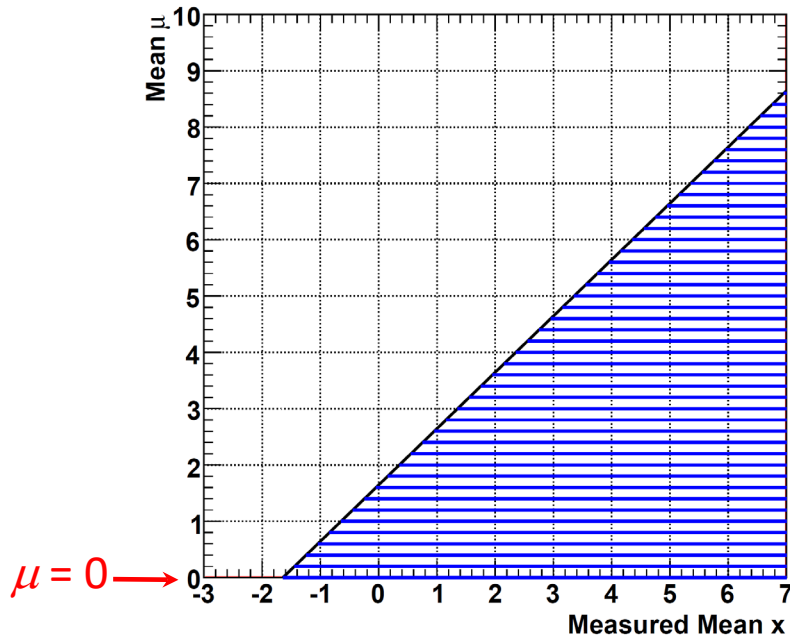
In the last few years, Feldman-Cousins starts to be used more, without my pushing. (Initially some at the LHC were very opposed, evidently because it could return two-sided interval not including zero when they really wanted a strict upper limit.)

CMS recently submitted documentation of its publicly available software, “The CMS statistical analysis and combination tool: COMBINE”, <https://arxiv.org/abs/2404.06614> .

For ATLAS statistics software tools, see backup slide.

Downward fluctuations in searches for excesses

**Classic example: Upper limit on mean μ of Gaussian pdf for x .
Frequentist UL construction if $\mu \geq 0$ in model ($\sigma=1$) :**



**Frequentist 1-sided 95% C.L.
Upper Limits, for $\alpha = 1 - \text{C.L.} = 5\%$.**

As observation x becomes increasingly negative, standard frequentist upper limit becomes small and then null.

For $x < -1.64 \sigma$, the confidence interval is the *null set*!

Downward fluctuations in searches for excesses

Issue acute 20-30 years ago in expts to measure ν_e mass in (tritium β decay): several measured $m_\nu^2 < 0$.

This is a very long story; see my “virtual talk”

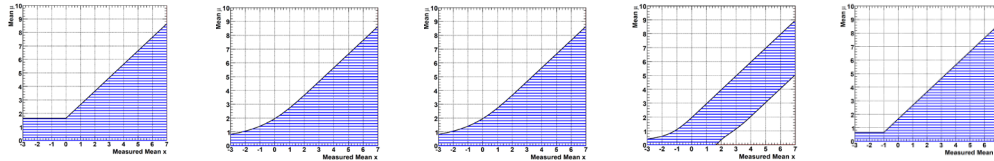
http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf

and related post <https://arxiv.org/abs/1109.2023> .

Contains intro to “Buehler’s betting game”, related to conditioning.



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



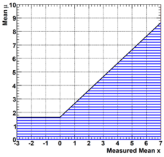
Bob Cousins
Univ. of California, Los Angeles

Virtual Talk

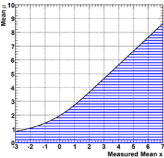
12 September 2011

http://www.physics.ucla.edu/~cousins/stats/cousins_bounded_gaussian_virtual_talk_12sep2011.pdf

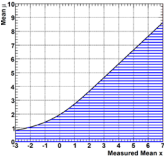
Five methods used for bounded Gaussian mean problem



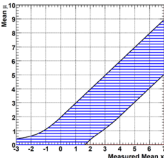
1) 1960's and beyond:
 $UL = \max(x, 0) + 1.64\sigma$



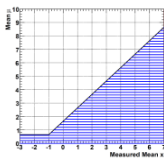
2) 1979 "PDG" (real 1986 PDG) and beyond:
Bayesian with uniform prior



3) 1997: Alex Read et al. (LEP)
 CL_s



4) 1997: Feldman and Cousins (NOMAD)
Unified Approach



5) 2010: Power Constrained Limits;
Cowan, Cranmer, Gross, Vitells (ATLAS):
 $UL = \max(0, \max(x, x_{pCL}) + 1.64\sigma)$

CL_s

The unfortunately named CL_s is the traditional frequentist one-tailed p-value for upper limits divided by another tail probability less than 1. The limits are thus (intentionally) conservative.

Definition, discussion, references in PDG RPP, Section 40.4.2.4 <https://pdg.lbl.gov/2024/reviews/rpp2024-rev-statistics.pdf> .

One way to mitigate the problem of excluding models to which one is not sensitive is the CL_s method, where the measure used to test a parameter is increased for decreasing sensitivity [43, 44]. The procedure is based on a statistic called CL_s, which is defined as

$$\text{CL}_s = \frac{p_\mu}{1 - p_b} , \quad (40.84)$$

where p_b is the p -value of the background-only hypothesis. In the usual formulation of the method, both p_μ and p_b are defined using a single test statistic, and the definition of CL_s above assumes this statistic is continuous; more details can be found in Refs. [43, 44].

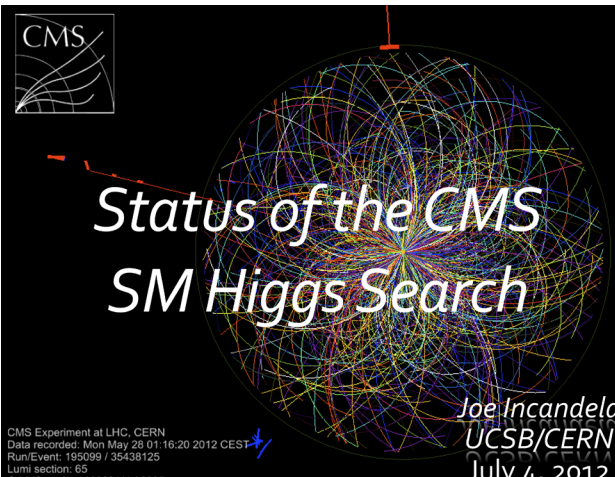
A point in a model's parameter space is regarded as excluded if one finds $\text{CL}_s \leq \alpha$. As the denominator in Eq. (40.84) is always less than or equal to unity, the exclusion criterion based on CL_s is more stringent than the usual requirement $p_\mu \leq \alpha$. In this sense the CL_s procedure is conservative, and the coverage probability of the corresponding intervals will exceed the nominal confidence level $1 - \alpha$. If the experimental sensitivity to a given value of μ is very low, then one finds that as p_μ decreases, so does the denominator $1 - p_b$, and thus the condition $\text{CL}_s \leq \alpha$ is effectively prevented from being satisfied. In this way the exclusion of parameters in the case of low sensitivity is suppressed. [...]

A few more notes/history are in my arxiv post on these lectures.

CL_s inherits all issues of p-values

- **What is new (non-standard statistics) in CL_s is combining the two p-values into 1 quantity.**
 - This step is called “the CL_s criterion” in CMS papers.
- **The p-values themselves have long existed in the statistics literature and should be designated that way (and not with the names that the inventor of named CL_s gave them). All the issues of p-values (choice of test statistic, how to eliminate nuisance parameters) of course still exist.**
- **LEP, Tevatron, and LHC Higgs combination groups differ in choices (!)**
 - What specific likelihood ratio used in test statistic
 - Treatment of nuisance parameters
 - Ensembles used for “Toy M.C.” used to get distribution of test statistic under H_0 (no Higgs) and H_1 (SM with Higgs)

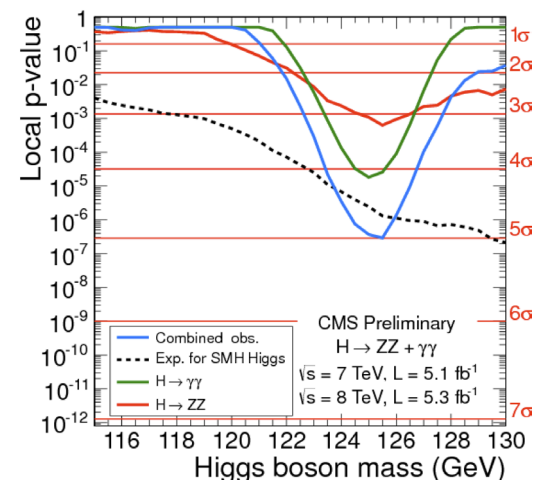
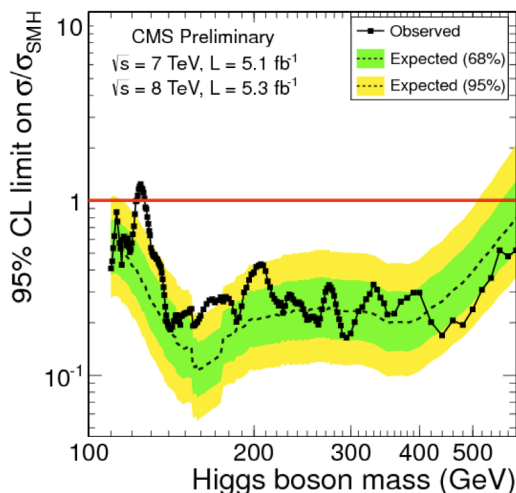
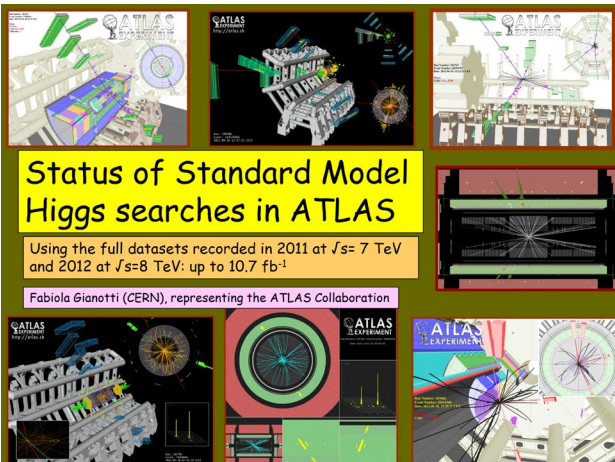
Statistics in practice on July 4, 2012: Imagine being in the audience for the talks on the discovery of the Higgs boson



Status of the CMS SM Higgs Search

Joe Incandela
UCSB/CERN
July 4, 2012

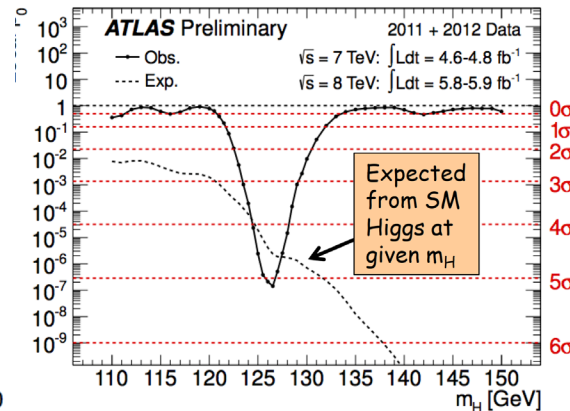
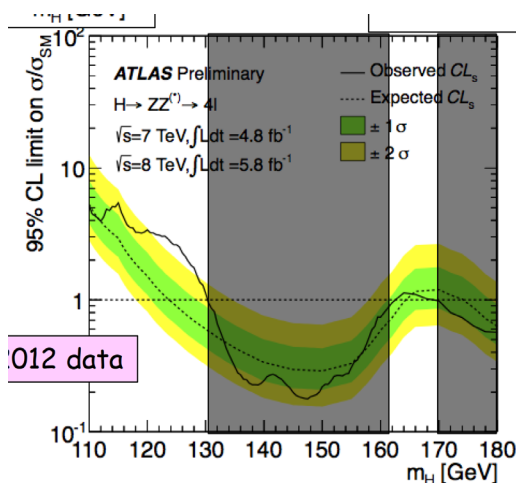
CMS Experiment at LHC, CERN
Data recorded: Mon May 28 01:16:20 2012 CEST
RunEvent: 195099 / 35438125
Lumi section: 65
Orbit/Crossing: 16992111 / 2295

Status of Standard Model Higgs searches in ATLAS

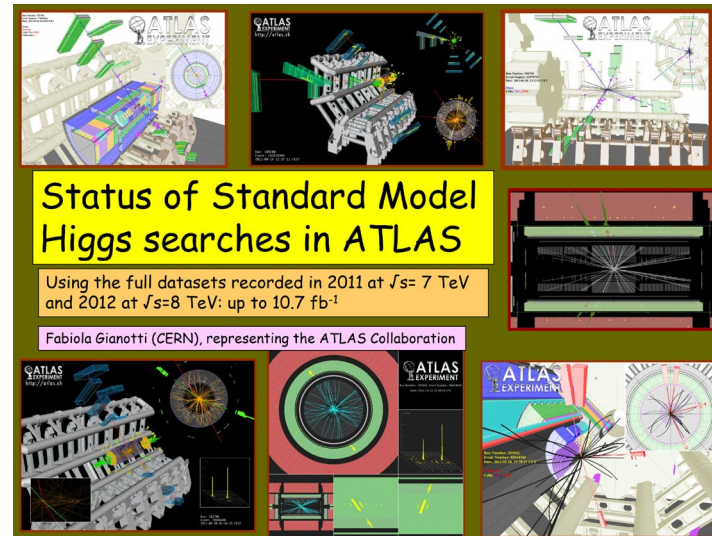
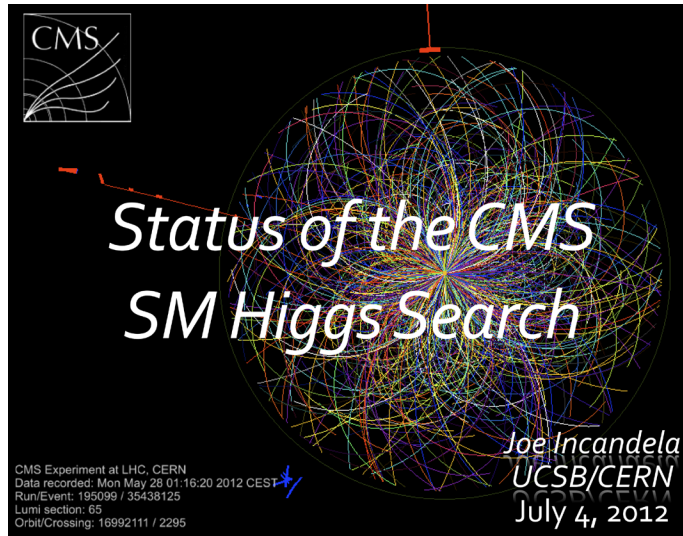
Using the full datasets recorded in 2011 at $\sqrt{s} = 7 \text{ TeV}$ and 2012 at $\sqrt{s} = 8 \text{ TeV}$: up to 10.7 fb^{-1}

Fabiola Gianotti (CERN), representing the ATLAS Collaboration



**Incandela and Gianotti slides are at <https://indico.cern.ch/event/197461/>
 See arxiv writeup for links to journal papers following shortly thereafter.**

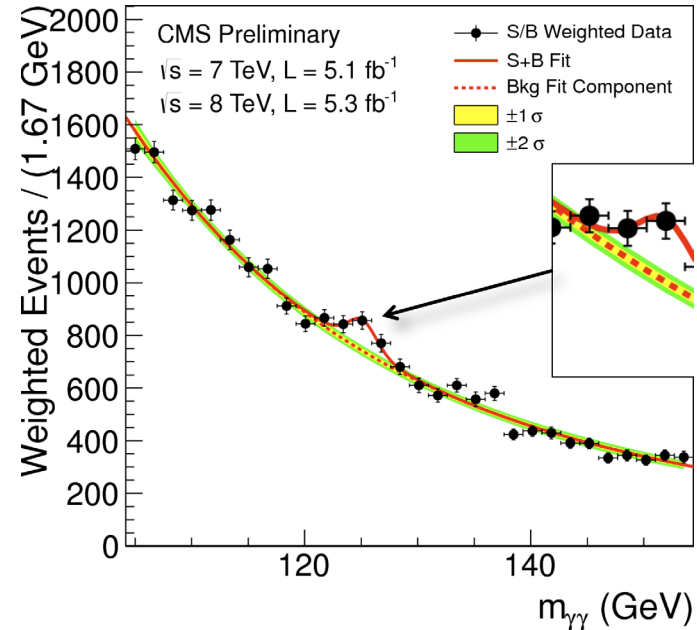
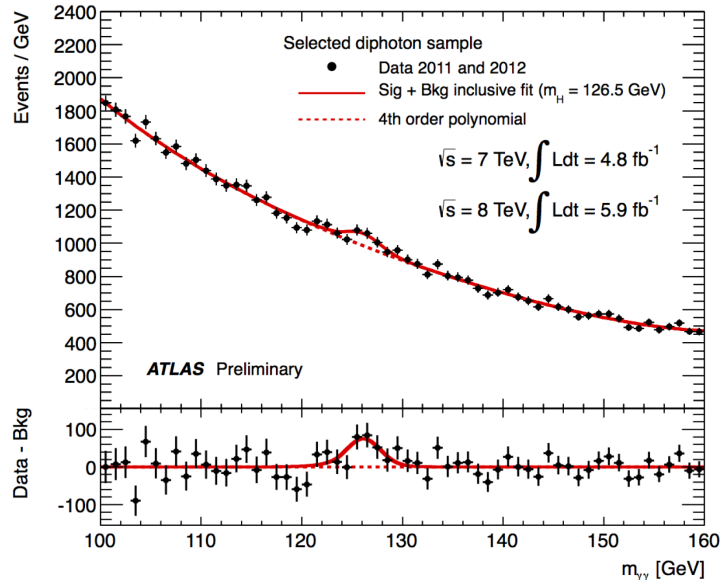
Example of Statistics in Practice



As already mentioned, description by ATLAS and CMS is “Procedure for the LHC Higgs boson search combination in Summer 2011,” <http://cds.cern.ch/record/1379837>, which you can look at for the math details.

The best decay modes for discovering the Higgs boson turned out to be (story too long for today): $H \rightarrow \gamma \gamma$ and $H \rightarrow ZZ^*$, when the Z and Z^* decay to either $e^+ e^-$ or to $\mu^+ \mu^-$ (total of 4 leptons).

Invariant mass of $\gamma\gamma$ final states.

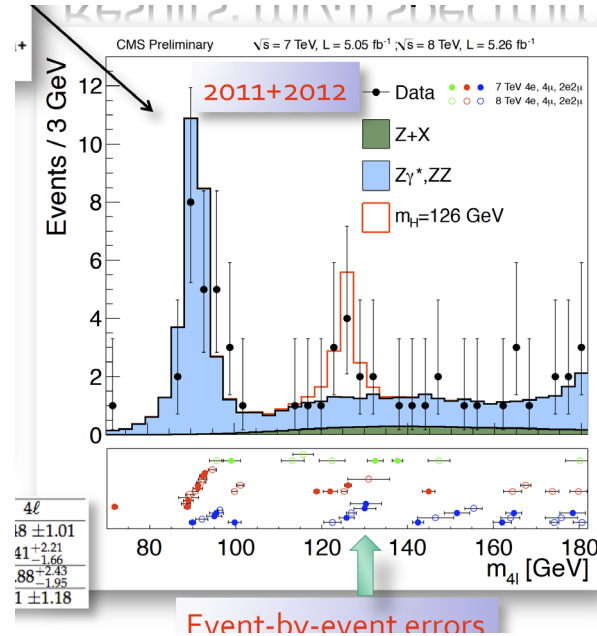
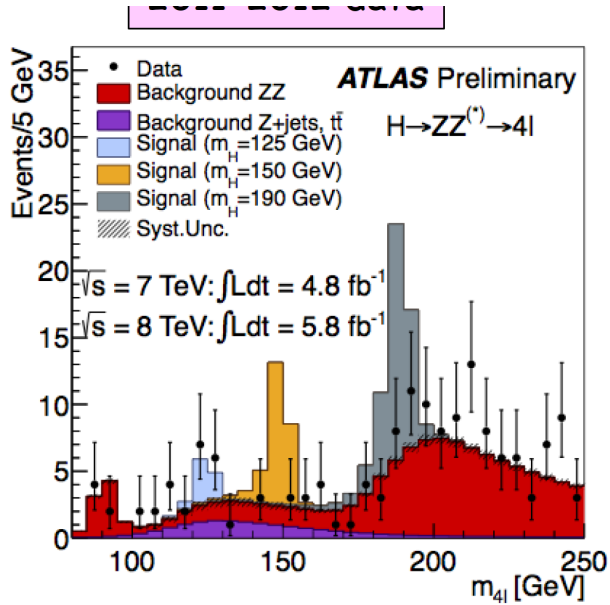


First, consider where we do ***not*** see evidence for a bump.
Consider each mass m separately.

Form profile likelihood function with signal cross section as the parameter of interest, divide by global ML as function of all parameters.

With this PLR as test statistic, calculate 95% C.L. upper limit (using CL_s) on *ratio of cross section to SM cross section* at m .

Invariant mass of 4-lepton final states.

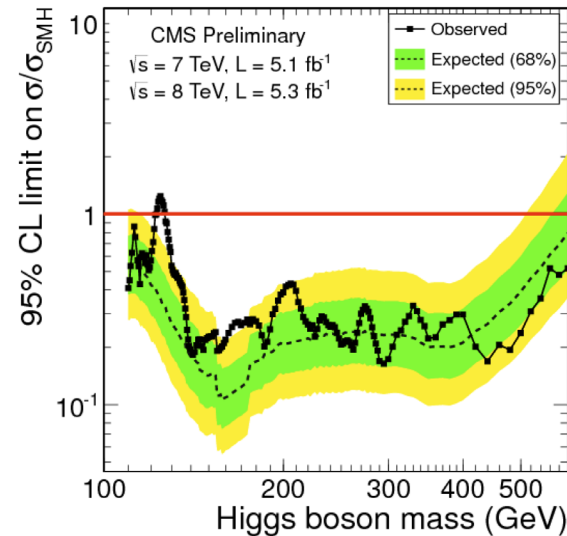
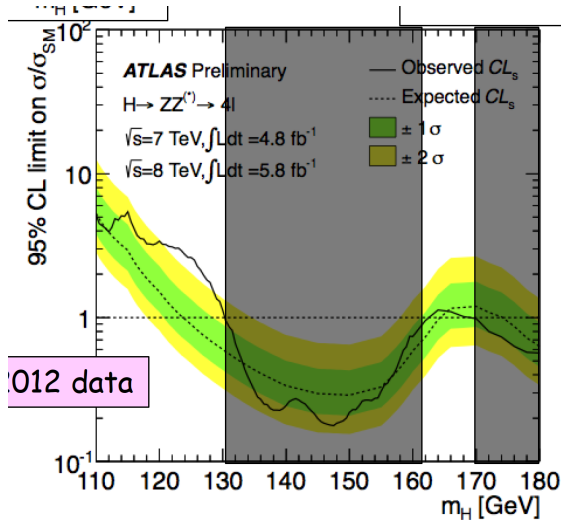


Repeat for 4-lepton final states: $e^+ e^- e^+ e^-$; $e^+ e^- \mu^+ \mu^-$; $\mu^+ \mu^- \mu^+ \mu^-$.

Then, at each mass m , perform one grand combination of all likelihood functions, and compute combined profile likelihood ratio, and 95% C.L. upper limit on cross section ratio.

Repeat for each mass m separately.

Combined upper limits for each experiment



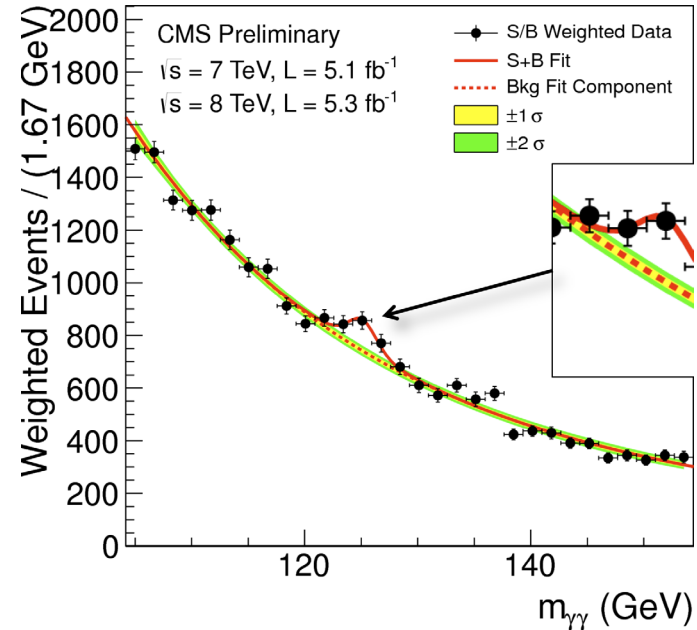
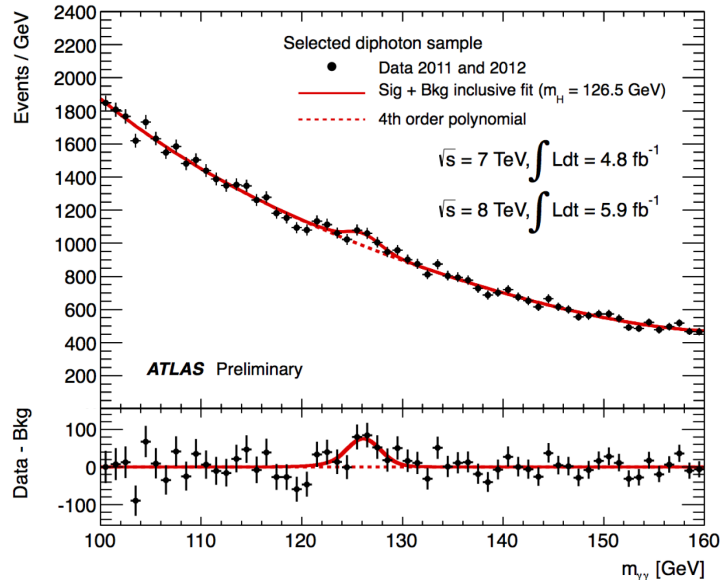
The jagged black curves, based on observed data, are the upper limits (UL) on the cross section ratio as function of m .

The SM Higgs boson is “excluded” at 95% C.L. at masses where the curve drops below 1 (i.e., *not* in a range below 130 GeV).

The dotted black curves show the exclusion if the data had been *equal to* the expected background, *without* any fluctuations.

The green and yellow (“Brazilian flag”) bands are *quantiles* of the expected distribution of ULs if expected background has statistical fluctuations, given no signal. (Not uncertainties!)

Signal searches in invariant mass of $\gamma\gamma$ and 4 leptons



For each final state, start over, again consider each mass m .

Form profile likelihood ratio of

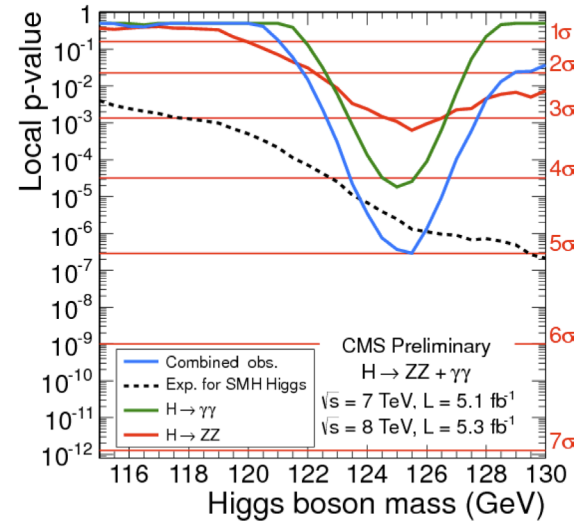
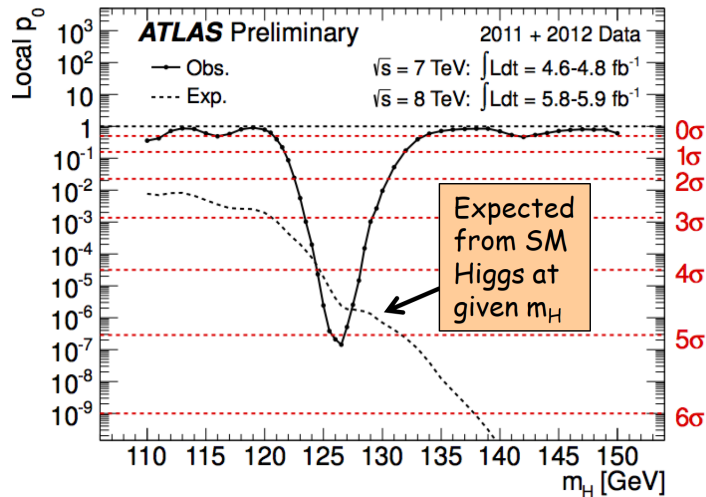
H_0 : background only vs.

H_1 : SM H , except cross section floating (parameter of interest).

From observed counts at mass m , obtain “local” p-value for H_0 .

Then obtain local p-value from PLR using grand combination of likelihood functions of $\gamma\gamma$ and ZZ^* states, convert to Z-values.

Local p-value plots vs mass



In the combinations of $H \rightarrow \gamma\gamma$ and $H \rightarrow ZZ^*$ channels, each experiment had minimum local p-values corresponding to 5σ . (CMS also combined channels with less sensitivity, got 4.9σ .)

This does not yet account for the bias because one scanned a range of masses to look for the smallest p-value.

Correcting for this “look-elsewhere effect” (with somewhat arbitrary choice of range of masses) yielded “global p-values” corresponding to about 4.1 to 4.3σ .

See LLE references in backup.

My advocacy for >15 years:

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

Unsound statements you can now avoid*

- “It makes no sense to talk about the probability density of a constant of nature.”
- “Frequentist confidence intervals for efficiency measurements don’t work when all trials give successes.”
- “We used a uniform prior because this introduces the least bias.”
Or “a noninformative prior since it contained no information.”
- “The total number of events could fluctuate in our experiment, so *obviously* our toy Monte Carlo should let the number of events fluctuate.”
- We used Delta-likelihood contours so there was no Gaussian approximation.”
- “A five-sigma effect constitutes a discovery.”
- “The confidence level tells you how much confidence one has that the true value is in the confidence interval.”
- “We used the tail area under the likelihood function to measure the significance.”
- “Statistics is obvious, so I prefer not to read the literature and just figure it out for myself.”

*Taken from real life

Thanks again!

Thanks to many in HEP (Frederick James, Gary Feldman, Louis Lyons, Luc Demortier, + numerous others) from whom I learned...

... And many statisticians that Louis invited to PhyStat meetings. For Bayesian statistics, that was especially Jim Berger (multiple times), Michael Goldstein, and David van Dyk (multiple times).

...and to CMS Statistics Committee (Olaf Behnke, Igor Volobouev, et al.) for many discussions and comments on earlier versions of the slides...

...and to the authors of numerous papers from which I learned, including early (1980s) Bayesian papers by Harrison Prosper...

...and to Diego Tonelli of the LHCb experiment for encouragement and comments on an earlier version of these slides.

This work was partially supported by the U.S. Department of Energy under Award Number DE–SC0009937.

References Cited in Talk Slides

Behnke13: O. Behnke et al., Data Analysis in High energy Physics, Wiley-VCH, 2013.

James06: Frederick James, Statistical Methods in Experimental Physics, World Scientific, 2006.

Stuart99: A. Stuart, K. Ord, S. Arnold, Kendall's Advanced Theory of Statistics, Vol. 2A, 6th edition, 1999; and earlier editions by Kendall and Stuart.

Recommended reading

Books: I usually recommend the following progression, reading the first three cover-to-cover, and consulting the rest as needed:

- 1) **Philip R. Bevington and D.Keith Robinson**, Data Reduction and Error Analysis for the Physical Sciences (Quick read for undergrad-level review)
 - 2) **Glen Cowan**, Statistical Data Analysis (Solid foundation for HEP)
 - 3) **Frederick James**, Statistical Methods in Experimental Physics, World Scientific, 2006. (This is the second edition of the influential 1971 book by **Eadie et al.**, has more advanced theory, many examples)
 - 4) **A. Stuart, K. Ord, S. Arnold**, Kendall's Advanced Theory of Statistics, Vol. 2A, 6th edition, 1999; and earlier editions of this "**Kendall and Stuart**" series. (Comprehensive old treatise on classical frequentist statistics; anyone contemplating a NIM paper on statistics should look in here first!)
 - 5) **George Casella and R.L. Berger**, Statistical Inference, 2nd, Ed. 2002. A more modern, less dense text on similar topics as Kendall and Stuart.
 - 6) Recent book by HEP "experts": **O. Behnke et al.**, Data Analysis in HEP, 2013
- PhyStat conference series:** From Confidence Limits Workshops in 2000, links: <https://phystat.github.io/Website/> . (Click on past Workshops). See also other links.
- My **Bayesian reading list** is the set of citations in my Comment, Phys. Rev. Lett. 101 029101 (2008), especially refs 2, 8, 9, 10, 11; 7 for model selection)

BACKUP



Workshop on Confidence Limits

27-28 March, 2000
Fermilab 1-West Conference Room

Jim Berger:

M. Kendall, giving the 'old' frequentist viewpoint of Bayesian analysis:

"If they [Bayesians] would only do as he [Bayes] did and publish posthumously, we should all be saved a lot of trouble."

What should be the view today:
Objective Bayesian analysis is the best frequentist tool around.

New Era for non-subjective Bayesian priors in HEP?

PHYSICAL REVIEW D **82**, 034002 (2010)

Reference priors for high energy physics

Luc Demortier

Laboratory of Experimental High Energy Physics, The Rockefeller University, New York, New York 10065, USA

Supriya Jain*

Homer L. Dodge Department of Physics and Astronomy, University of Oklahoma, Norman, Oklahoma 73019, USA

Harrison B. Prosper

Department of Physics, Florida State University, Tallahassee, Florida 32306, USA

(Received 4 February 2010; published 3 August 2010)

Bayesian inferences in high energy physics often use uniform prior distributions for parameters about which little or no information is available before data are collected. The resulting posterior distributions are therefore sensitive to the choice of parametrization for the problem and may even be improper if this choice is not carefully considered. Here we describe an extensively tested methodology, known as reference analysis, which allows one to construct parametrization-invariant priors that embody the notion of minimal informativeness in a mathematically well-defined sense. We apply this methodology to general cross section measurements and show that it yields sensible results. A recent measurement of the single-top quark cross section illustrates the relevant techniques in a realistic situation.

DOI: [10.1103/PhysRevD.82.034002](https://doi.org/10.1103/PhysRevD.82.034002)

PACS numbers: 06.20.Dk, 14.65.Ha

So far, use in HEP has been limited.

Robert Kass's Questions for Classifying Kinds of Bayesians*

1. **Is it important for Bayesian inferences to have good frequentist operating characteristics?**
2. **Does the Bayesian paradigm do anything more than produce candidate procedures, to be judged according to frequentist criteria?**
3. **Is there a useful role for default (a.k.a. “objective”) Bayesian inferences as representing approximately subjective inferences?**
4. **Is it possible to interpret default Bayesian inference as anything other than approximately subjective?**
5. **Assuming that the data analyst has done a thorough and careful job, is it appropriate to interpret default Bayesian inferences as representing, approximately, what any reasonable person ought to think given the data and appropriate background information?**
6. **Is there any useful meaning to the word “objective,” beyond signifying such overwhelming evidence that reasonable people will be forced to agree.**
7. **Is the word “objective” so easy to misunderstand that its utility in the context of Bayesian inference is, on average, negative?**
8. **Is it important to distinguish scientific inference from decision-making?**
9. **Are there scientific settings in which formal elicitation procedures are useful?**

* “Kinds of Bayesians (Comment on articles by Berger and by Goldstein)”, *Bayesian Analysis* 1 437 (2006)

Clopper and Pearson's construction

x = number of successes (here, integer 0-10 out of 10 trials)

Inner corners of the steps give the intervals; traditional to draw the curved "belts" connecting them, but only evaluated at the integers. Tricky to draw, read! (See next slide for details.)

Discreteness of x typically requires horizontal acceptance intervals to contain more than 95% probability, so there is over-coverage in the vertical confidence intervals.

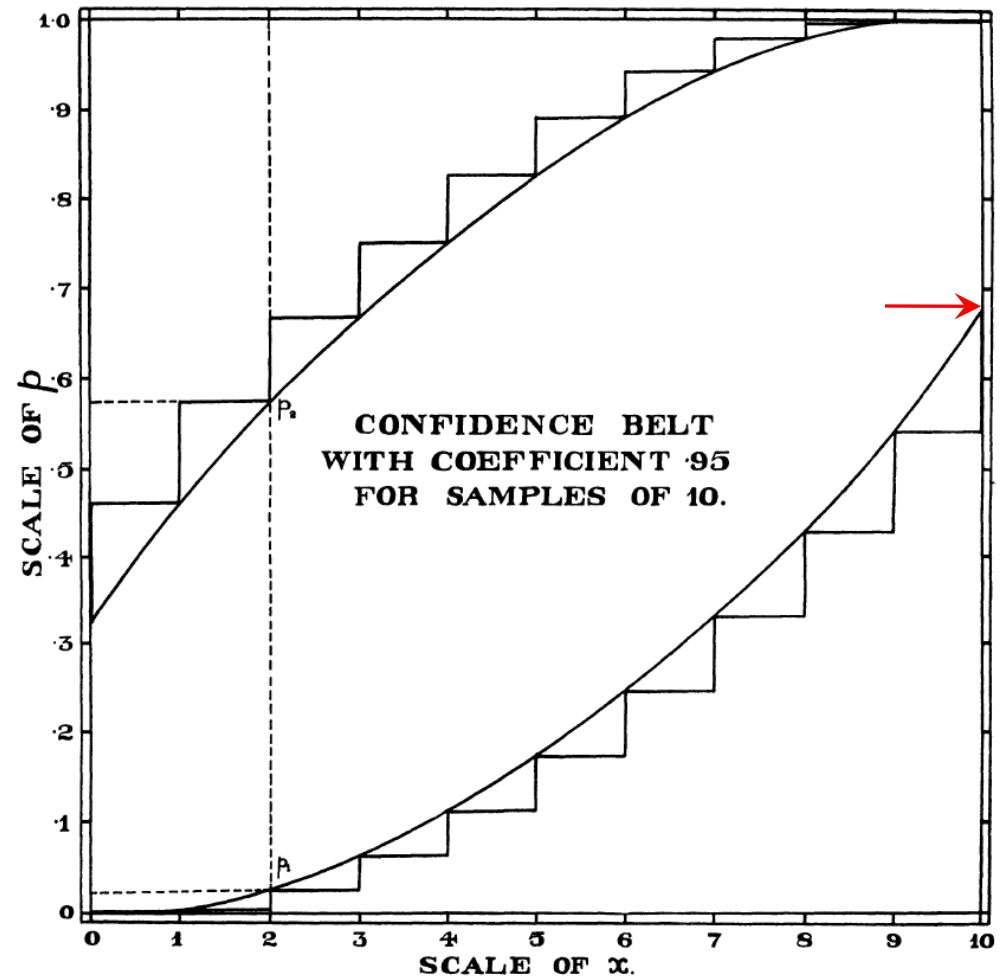


FIG. 1

E.g. 95% C.L. central interval for p if 10/10 successes/trials: (0.69,1.0)

Clopper and Pearson's construction (cont.)

Partial details of construction:

Blue lines are two of the acceptance intervals having central 95% or more probability, at continuous ρ .

Note data x is discrete, so graph is only read at discrete x .

If you stare at it long enough, you will see connection between upper/lower limits and central intervals, for discrete data.

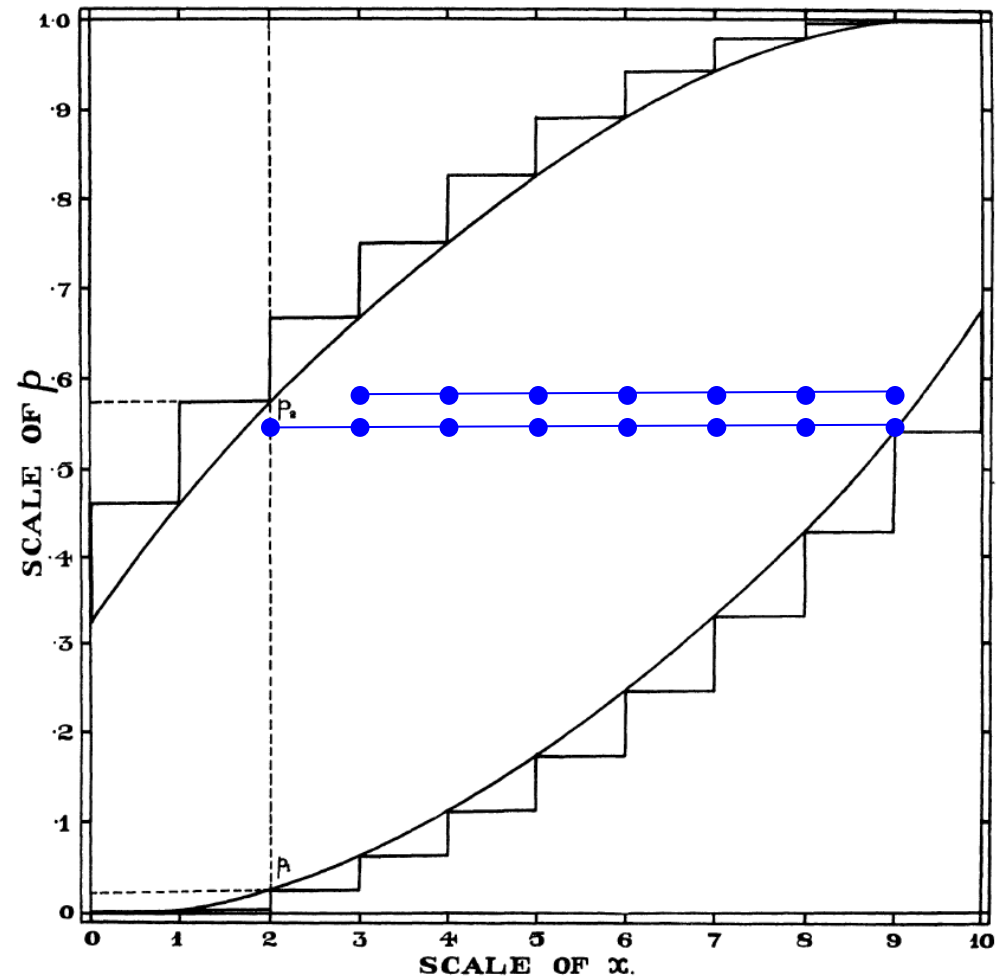


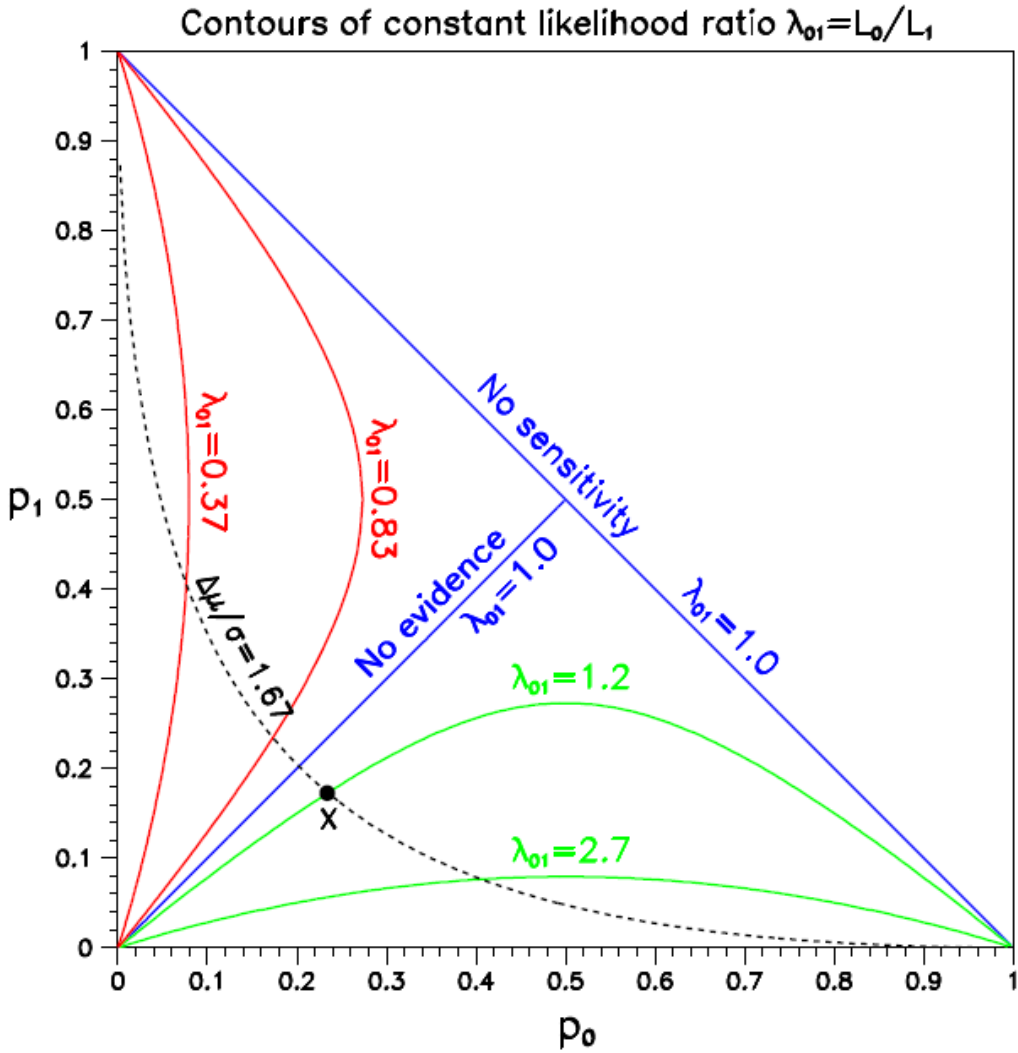
FIG. 1

Luc Demortier and Louis Lyons, “Testing Hypotheses in Particle Physics: Plots of p_0 versus p_1 ”

Test of point null vs point alternative, two Gaussians with same σ , peak separation $\Delta\mu$.

At a glance can see that contours of constant λ_{01} are completely different topology from contours of e.g. p_0 .

For rest of plot, you will have to read their paper or stare at it for a long time.
<http://arxiv.org/abs/1408.6123>



Classical Hypothesis Testing: Duality

Test $\mu=\mu_0$ at $\alpha \leftrightarrow$ Is μ_0 in conf. int. for μ with C.L. = $1- \alpha$

“There is thus no need to derive optimum properties separately for tests and for intervals; there is a one-to-one correspondence between the problems as in the dictionary in Table 20.1”

Stuart99, p. 175.

Table 20.1 Relationships between hypothesis testing and interval estimation

Property of test	Property of corresponding confidence interval
Size = α	Confidence coefficient = $1 - \alpha$
Power = probability of rejecting a false value of $\theta = 1 - \beta$	Probability of not covering a false value of $\theta = 1 - \beta$
Most powerful	Uniformly most accurate
	$\left\{ \begin{array}{l} \text{Unbiased} \\ 1 - \beta \geq \alpha \end{array} \right\}$
Equal-tails test $\alpha_1 = \alpha_2 = \frac{1}{2}\alpha$	Central interval

Referred to as “inverting a test” to obtain intervals; vice versa.

Approximate Confidence Regions Using $\Delta(-\ln\mathcal{L})$

(included in appendix to MINUIT users guide)

Computer Physics Communications 20 (1980) 29–35

INTERPRETATION OF THE SHAPE OF THE LIKELIHOOD FUNCTION AROUND ITS MINIMUM

F. JAMES

Data Handling Division, CERN, Geneva, Switzerland

It often happens that the solution of a minimum problem is itself straightforward, but the calculation or interpretation of the resulting parameter uncertainties, as determined by the shape of the function at the minimum, is considerably more complicated. The purpose of this note is to clarify the most commonly encountered difficulties in parameter error determination. These difficulties may arise in connection with any fitting program, but will be discussed here with the terminology of the program MINUIT for the convenience of MINUIT users.

The most common causes of misinterpretation may be grouped into three categories:

1. Proper normalization of the user-supplied chi-square or likelihood function, and appropriate ERROR DEF.
2. Non-linearities in the problem formulation, leading to different errors being calculated by different techniques, such as MIGRAD, HESSE and MINOS.
3. Multiparameter error definition and interpretation.

All these topics are discussed in some detail by Eadie et al., which may be consulted for further details.

Table 1

Table of UP for multiparameter confidence regions

Number of parameters	Confidence level (probability contents desired inside hypercontour of " $\chi^2 = \chi^2_{\min} + UP$ ")				
	50%	70%	90%	95%	99%
1	0.46	1.07	2.70	3.84	6.63
2	1.39	2.41	4.61	5.99	9.21
3	2.37	3.67	6.25	7.82	11.36
4	3.36	4.88	7.78	9.49	13.28
5	4.35	6.06	9.24	11.07	15.09
6	5.35	7.23	10.65	12.59	16.81
7	6.35	8.38	12.02	14.07	18.49
8	7.34	9.52	13.36	15.51	20.09
9	8.34	10.66	14.68	16.92	21.67
10	9.34	11.78	15.99	18.31	23.21
11	10.34	12.88	17.29	19.68	24.71

If FCN is $-\log(\text{likelihood})$ instead of chi-square, all values of UP should be divided by 2.

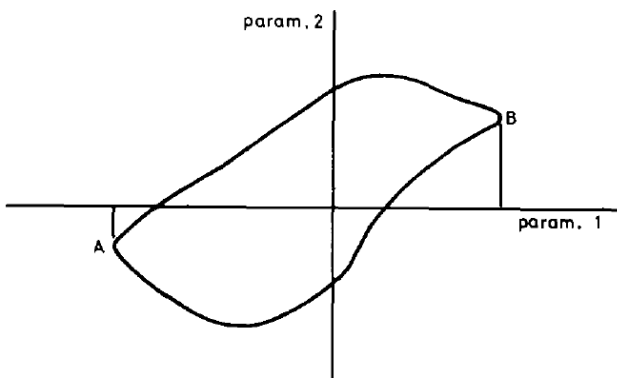


Fig. 1. MINOS errors for parameter 1.

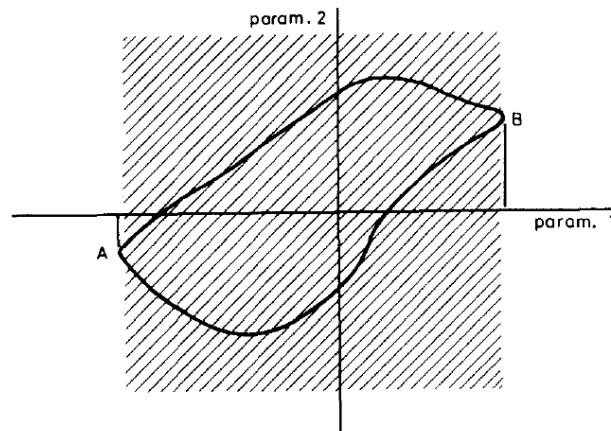


Fig. 2. MINOS error confidence region for parameter 1.

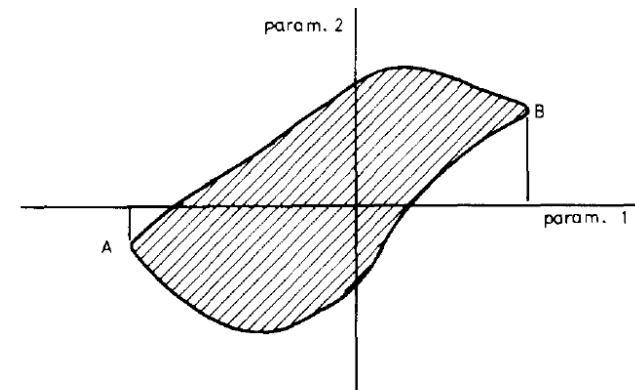


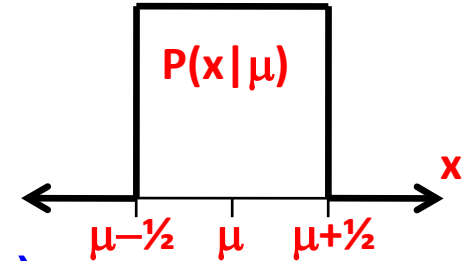
Fig. 4. Optimal confidence region for parameters 1 and 2.

Famous example of B.L. Welch (1939)

$$\text{Let } p(x|\mu) = \begin{cases} 1 & \text{if } \mu - 1/2 \leq x \leq \mu + 1/2 \\ 0 & \text{otherwise.} \end{cases}$$

Two values x_1, x_2 are observed.

$$\hat{\mu} = \bar{x} = (x_1 + x_2)/2 \text{ (Only for } n=2! \text{ See Backup)}$$



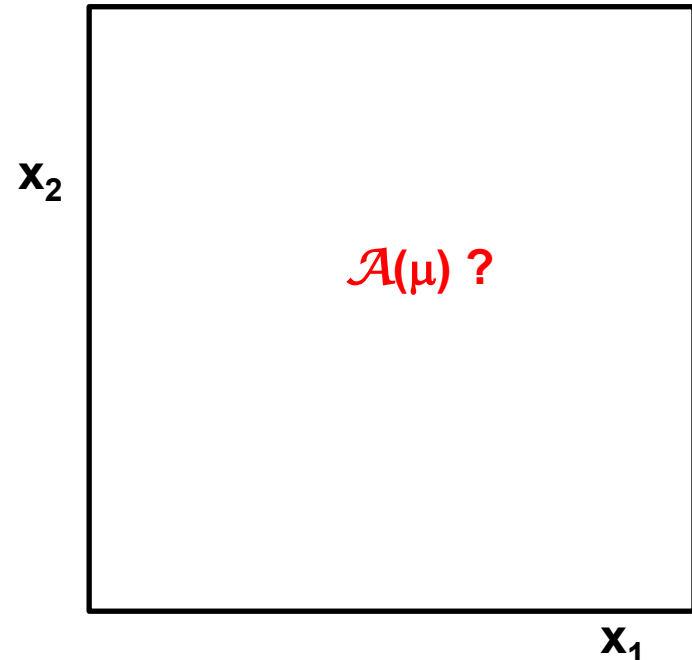
What is 68% C.L. central confidence interval for μ ?

Neyman construction: Define acceptance region $\mathcal{A}(\mu)$ containing 68% of unit square of (x_1, x_2) centered on μ . What to use?

Centrality implies symmetry.

Need something else to rank points in the plane.

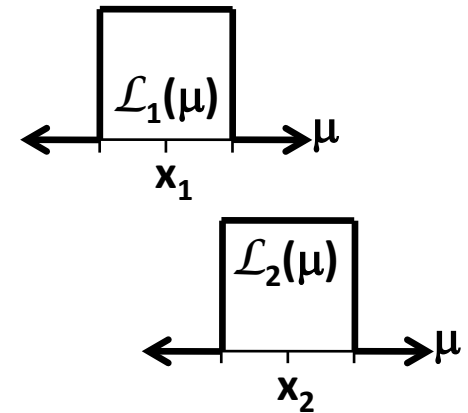
N-P Lemma gives most powerful ranking, but first let's think about some examples.



“Lucky” sample with $|x_1 - x_2|$ close to 1.

$\mathcal{L}(\mu) = \mathcal{L}_1(\mu) \times \mathcal{L}_2(\mu)$ very narrow.

Reasonable to expect small uncertainty in $\hat{\mu}$?



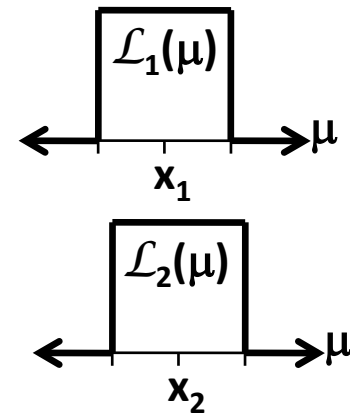
“Unlucky” sample, $|x_1 - x_2|$ close to 0.

$\mathcal{L}(\mu)$ full width close to 1;

Second observation added no useful info.

Expect 68% C.L. conf. interval 0.68 long?

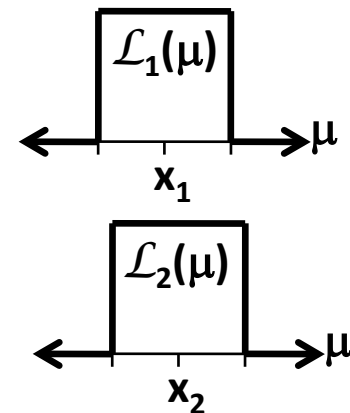
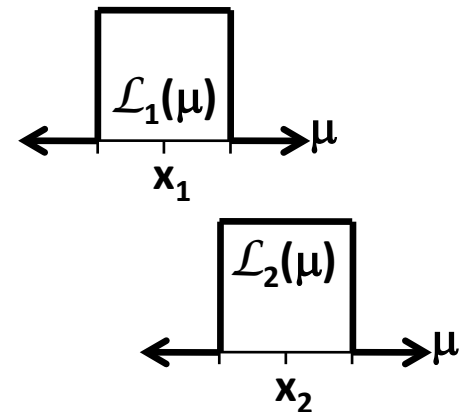
Guess reasonable answer: conf. interval centered on $\hat{\mu}$ with length $0.68(1 - |x_1 - x_2|)$



Seems reasonable for *post-data uncertainty* to depend on $|x_1 - x_2|$.

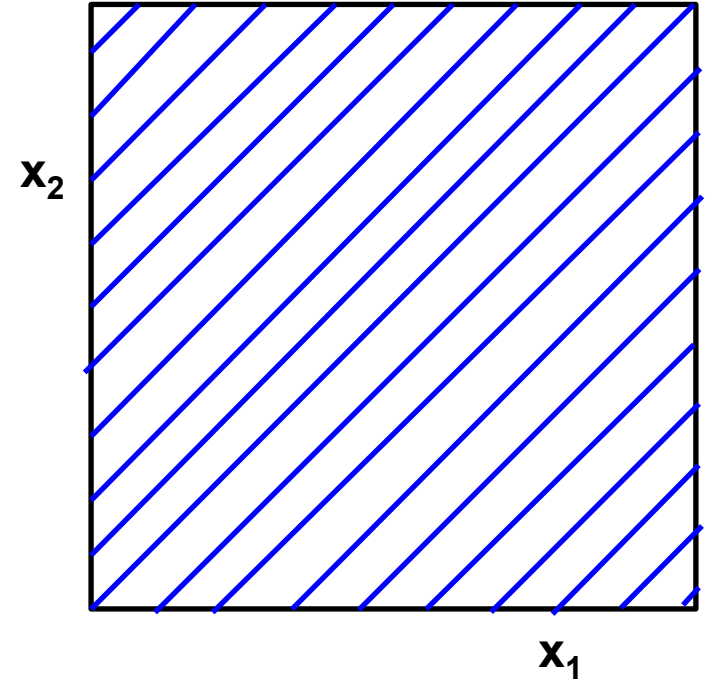
Classic example of an ancillary statistic A: has info on *uncertainty* on μ estimate, but no info on μ itself.

Idea dating to Fisher and before: divide the full “*unconditional*” sample space into “recognizable subsets” and report probs using the “relevant” subset rather than the whole space.



Acceptance regions partitioned by $|x_1 - x_2|$

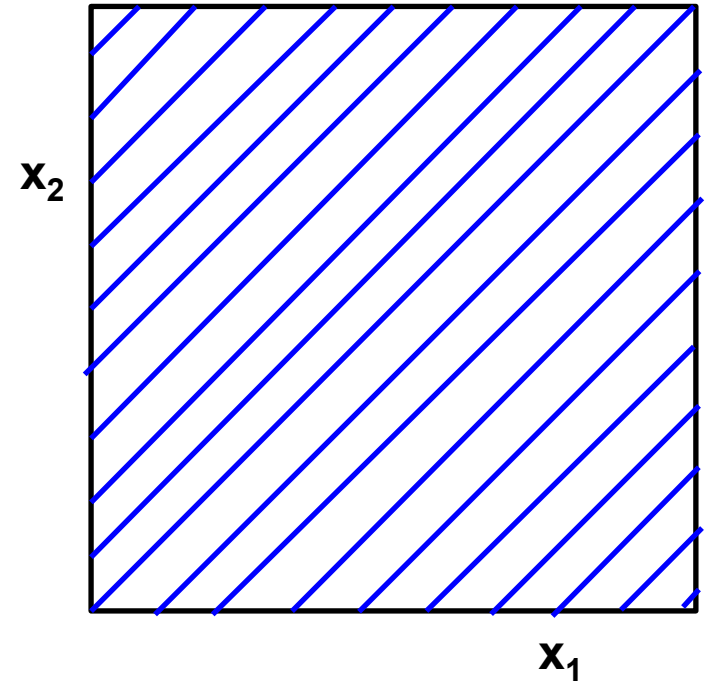
Partition full sample space via ancillary statistic $A = |x_1 - x_2|$ (blue lines)



Acceptance regions partitioned by $|x_1 - x_2|$

Partition full sample space via ancillary statistic $A = |x_1 - x_2|$ (blue lines)

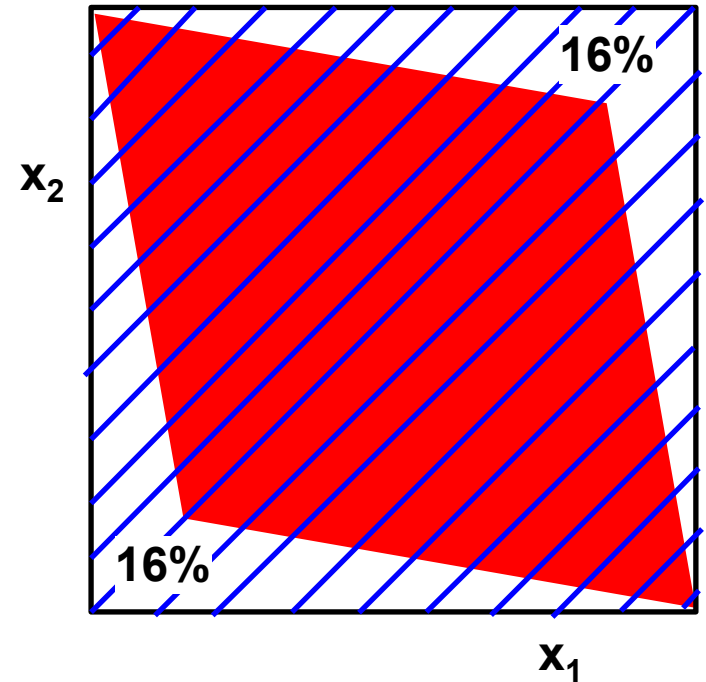
Within each partition, choose central 68% prob acceptance region



Acceptance regions partitioned by $|x_1 - x_2|$

Partition full sample space via ancillary statistic $A = |x_1 - x_2|$ (blue lines)

Within each partition, choose central 68% prob acceptance region (red fill)



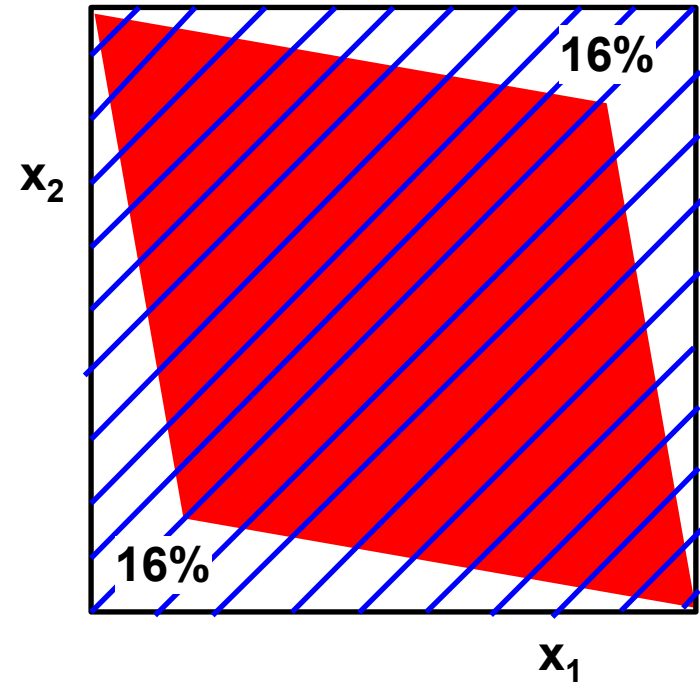
Acceptance regions partitioned by $|x_1 - x_2|$

Partition full sample space via ancillary statistic $A = |x_1 - x_2|$ (blue lines)

Within each partition, choose central 68% prob acceptance region (red fill)

We are thus using *conditional probabilities* (still frequentist) $p(x|A, \mu)$ in Neyman construction, with desired prob 68% *within each partition*.

Resulting $\mathcal{A}(\mu)$ fills 68% of square, so correct unconditional probability as well.



Acceptance regions partitioned by $|x_1 - x_2|$

Partition full sample space via ancillary statistic $A = |x_1 - x_2|$ (blue lines)

Within each partition, choose central 68% prob acceptance region (red fill)

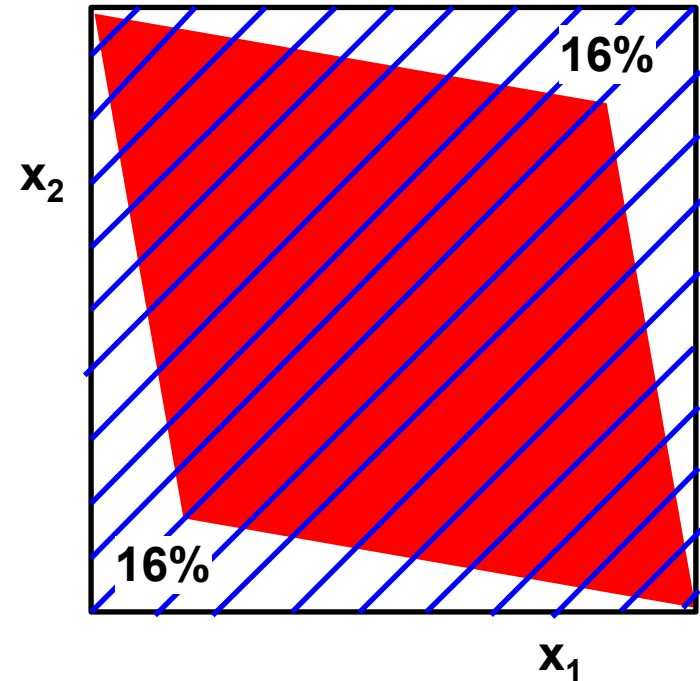
We are thus using *conditional probabilities* (still frequentist) $p(x|A, \mu)$ in Neyman construction, with desired prob 68% *within each partition*.*

Resulting $\mathcal{A}(\mu)$ fills 68% of square, so correct unconditional probability as well.

⇒ **Confidence intervals $\hat{\mu} \pm 0.34(1 - |x_1 - x_2|)$, as thought reasonable!**

Known as “conditioning” on ancillary statistic A : Post-data, *proceed as if A had been fixed, rather than randomly sampled!*

N.B. A set of measure zero has zero prob even if non-zero pdf, so in general care needed in conditioning on exact value of continuous A in $p(x|A, \mu)$.

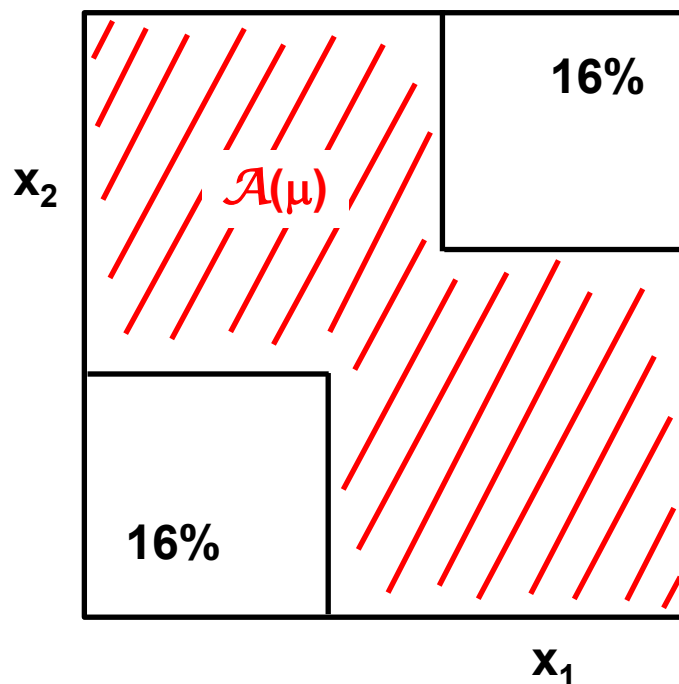
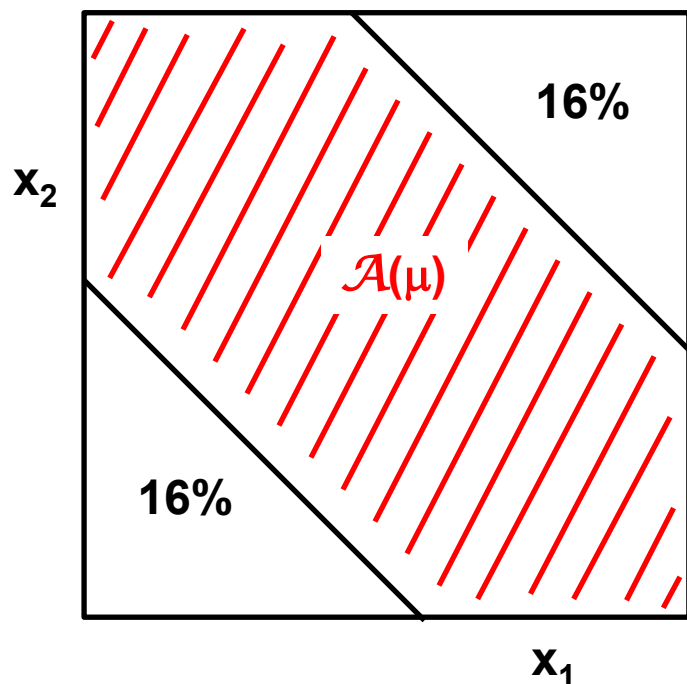


Now the catch: one can find acceptance regions $\mathcal{A}(\mu)$ that correspond to hypothesis tests with *more power* (lower Type 2 error prob β) in *the unconditional sample space!*

These have 100% coverage in the subspace where $|x_1 - x_2| \approx 1$ (narrow likelihood), while badly undercovering when $x_1 \approx x_2$.

Interval length indep of $|x_1 - x_2|$
 $\Rightarrow \hat{\mu} \pm 0.22$ at 68% C.L.

Most powerful?



In fact Welch's 1939 paper argued *against* conditioning because it is less powerful in the unconditional sample space!

Neyman's position is not completely clear but he seems to have been against conditioning on ancillaries (which was Fisher's idea) when it meant an overall loss of power.*

Most modern writers use this example as an “obvious” argument *in favor* of conditioning, unless one is in “industrial” setting where unconditional ensemble is sampled repeatedly and the result for an individual sample is not of much interest.

*See J.O. Berger, <https://projecteuclid.org/euclid.ss/1056397485> and

E.L. Lehmann, <http://www.jstor.org/stable/2291263> (also discussion of Cox example).

Example of D.R. Cox (cont.)

Demortier gives details on how average interval length is shorter in the HEP example. Here I give Cox's discussion.

E.g., if one is testing $\mu=0$ vs $\mu=\mu_1$, with μ_1 roughly the size of σ_1 (the larger σ), consider the following 68% CL intervals:

$\hat{\mu} \pm (0.48)\sigma_1$ if Device #1 used (covers true μ in 37% of uses)

$\hat{\mu} \pm 5\sigma_2$ if Device #2 used (covers true μ nearly 100% of uses)

So true μ is covered in $(37/2 + 100/2)\% = 68\%$ of all intervals!
Unconditional (full sample space) coverage is correct, but conditional coverage is not.

Due to the smallness of σ_2 , average length of all intervals is smaller conditional intervals with independent coverage.

One gives up power with Device #1 and uses it in Device #2.

Cox: "If, however, our object is to say 'what can we learn from the data that we have', the unconditional test is surely no good."

These examples reveal a real conflict between N-P optimization for power and conditioning to optimize relevance.

Look-Elsewhere Effect

In these lectures, I did not have time for the LEE.

A starting point for self-study is the discussion in:

Louis Lyons, “Comments on ‘Look Elsewhere Effect’ ”.

https://users.physics.ox.ac.uk/~lyons/LEE_feb7_2010.pdf .

See also Section 9.2 of my paper on the Jeffreys-Lindley Paradox, <https://arxiv.org/abs/1310.3791> .

An important paper is

Eilam Gross, Ofer Vitells, “Trial factors for the look elsewhere effect in high energy physics,” <https://arxiv.org/abs/1005.1891>

Sufficiency, Conditionality, Likelihood Principles

There is a lot more to the Likelihood Principle than I had time to discuss. I omitted the important (frequentist) concept of a “sufficient statistic”, due to Fisher. This is a way to describe data reduction without loss of relevant information. E.g., for testing a binomial parameter, one needs only the total numbers of successes and trials, and not the information on exactly which trials had successes. See Stuart99 for math definitions.

The Sufficiency Principle says (paraphrasing – there are strong and weak forms) that if the observed values of the sufficient statistic in two experiments are the same, then they constitute equivalence evidence for use in inference.

Birnbaum famously argued (1962) that the Conditionality Principle and the Sufficiency Principle imply the Likelihood Principle. Controversy continues. For recent discussion, see D. Mayo (2014), <https://projecteuclid.org/euclid.ss/1408368565#toc>, with comments by six statisticians and rejoinder.

Point Estimation

Point Estimation

Most of these slides are about intervals – I have not yet much about what to quote as the “measured value”. Statisticians call this the “point estimate”.

- **There is a huge literature on point estimation – see e.g. Ch. 7 and 8 in James06.**
- **If you are an expert on interval estimation, one approach is to use that machinery to get a point estimate.**
 - **E.g., one might take the mid-point of (say) your 68% C.L. central interval. But a better approach is probably to let the C.L. go to 0, so that your interval gets shorter and shorter, and use the limiting point. E.g. for likelihood ratio intervals, this results in the Maximum Likelihood estimate.**
- **But to give you an idea of how rich the subject is, I show a few interesting things from James06.**

Point Estimation: Traditional Desiderata

- **Consistency: Estimate converges toward true value as number of observations N increases**
- **Unbiasedness: Expectation value of estimate is equal to the true value.**
- **Efficiency: Estimate has minimum variance**
- **Minimum loss of information: (technical definition)**
- **Robustness: Insensitivity to departures from the assumed distribution**

One can add:

- **Simplicity: transparent and understandable**
- **Minimum computer time: still relevant in online applications, less relevant otherwise**
- **Minimum loss of physicist's time (how much weight to put on this?)**

Bias and consistency are independent properties

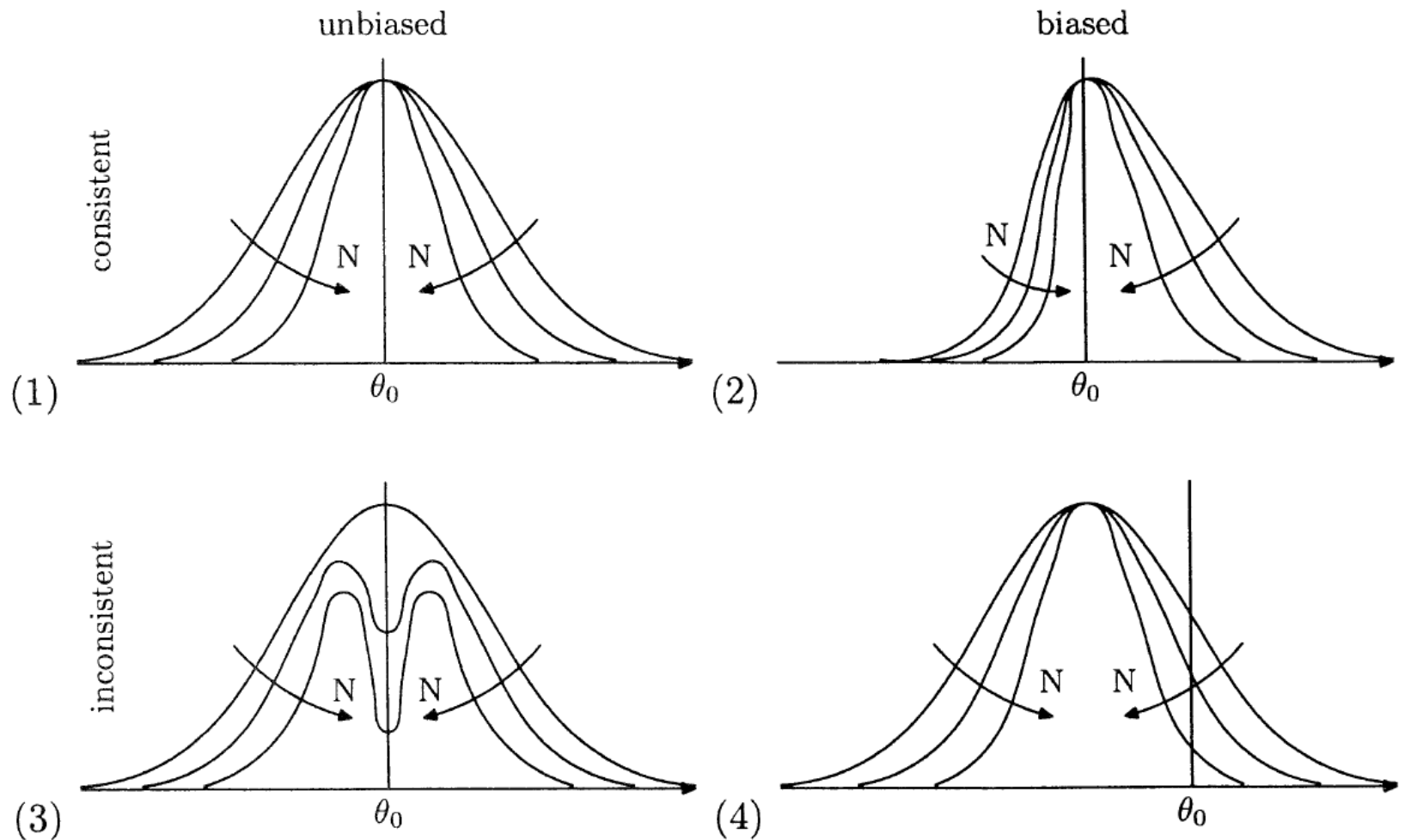


Fig. 7.2. Examples of probability density functions $f(\hat{\theta}, \theta_0)$ (not normalized) with different combinations of consistency and bias. The arrows show the direction of increasing N .

BUT (!) Other desired properties can be impossible to achieve simultaneously

- **How to choose? A thorough analysis requires further input: what are the costs of not incorporating various desiderata? Then formal *decision theory* can be used to choose estimator.**
- **In practice in HEP, *Maximum Likelihood* estimates are often used (even though they are typically not unbiased).**
 - **Consistent**
 - **Other excellent asymptotic properties (estimate is asymptotically normal)**
 - **For finite N , works well in so-called exponential family (includes Poisson, Gaussian, binomial)**
 - **Invariant under reparameterization**

Simple example illustrating diversity of point estimators (James06, p. 209)

- If $p(x|\mu) = f(x-\mu)$, where f is some pdf, then μ is called a *location parameter*. Common examples are:
 - Normal: $p \sim \exp(-(x-\mu)^2/2\sigma^2)$
 - Uniform: $p = \text{constant}$ for $|x-\mu| < a$; $p=0$ otherwise
 - Cauchy: $p \sim 1/(a^2 + (x-\mu)^2)$
 - Double exponential: $p \sim \exp(-a|x-\mu|)$
- These examples are all symmetric about μ :
 $p(\mu+y) = p(\mu-y)$
- Suppose you are given $N=11$ values of x randomly sampled from $p(x|\mu)$. What estimator (function of the 11 values) gives you the “best” estimate of μ ?
- If by “best” you mean minimum variance, it is the M.L. estimate, resulting in a different formula for each!

Minimum Variance Location Estimator

Normal	Sample mean (L_2)
Uniform	Midrange: mean of extreme values (L_∞)
Cauchy	M.L. estimate (no simple formula)
Double-exponential	Median: middle value (L_1)

Three of the four are special cases of L_p , the estimator that minimizes the sum over the observations of $|x_i - \mu|^p$.

Different values of p put different emphasis on observations in the tails.

If true distribution departs from that assumed, estimate of location is no longer optimal. Sensitivity is in tails!

See nice discussion of asymptotic variance and robustness in James06, pp. 211 ff.

ATLAS statistics software tools

Many thanks to Kyle Cranmer for compiling this list of the individual tools in use in the ATLAS Collaboration for the different tasks and stages:

RooFit: (core modelling)

<https://inspirehep.net/literature/621398>

RooStats: (statistical testing)

<https://inspirehep.net/literature/868303>

HistFactory: (specific modelling for histogram based analyses)

<https://inspirehep.net/literature/1236448>

HistFitter: (sits on top of HistFactory offers top-level steering functionality more like CMS's Combine)

<https://inspirehep.net/literature/1320562>

There is also TRexFitter that is like HistFitter and widely used, but no paper currently.

More recently, there is the python based implementation of HistFactory, which also has a different format for saving results, which is being used by HEPData. This is now fairly widely used.

pyhf: <https://inspirehep.net/literature/1845084>

And there is also the python-based approach to tools like HistFitter/TRexFitter called Cabinetry.

<https://inspirehep.net/literature/1911802>