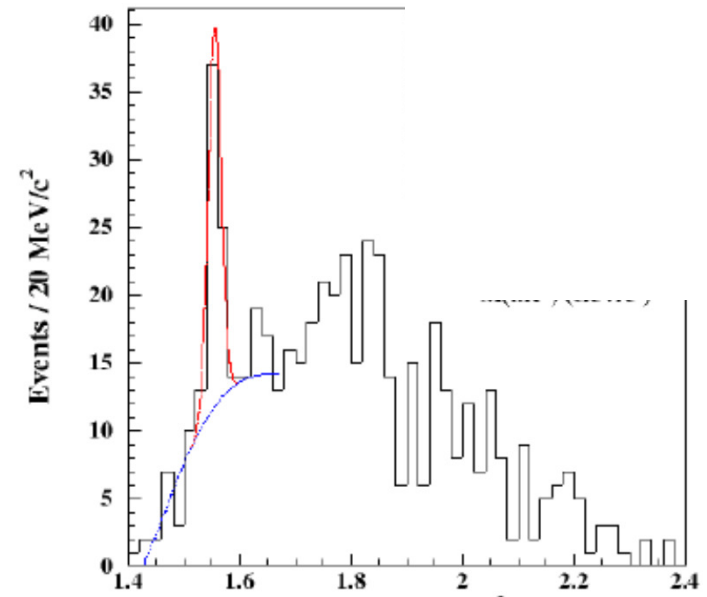
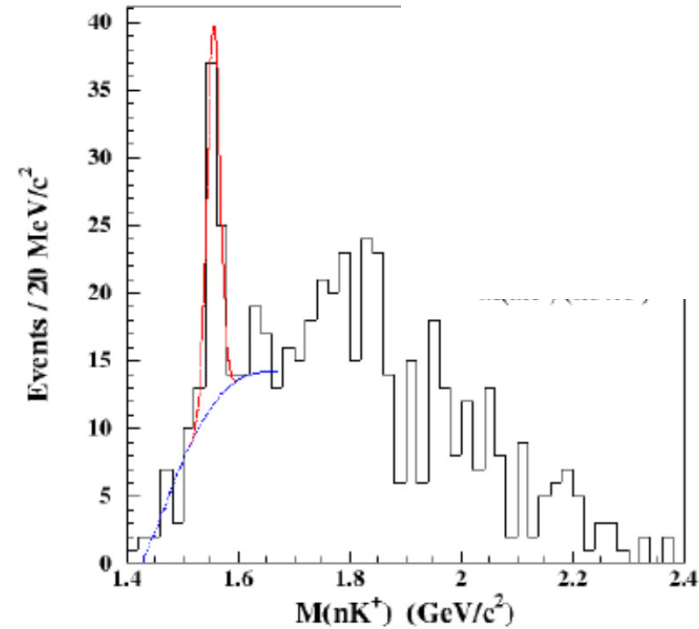


Is there evidence for a peak in this data?



Is there evidence for a peak in this data?



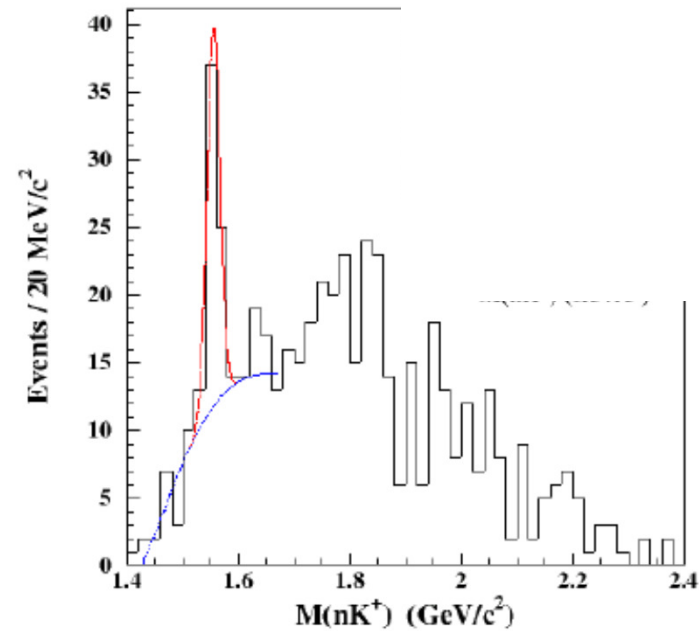
“Observation of an Exotic  $S=+1$

Baryon in Exclusive Photoproduction from the Deuteron”

S. Stepanyan et al, CLAS Collab, Phys.Rev.Lett. 91 (2003) 252001

“The statistical significance of the peak is  $5.2 \pm 0.6 \sigma$ ”

Is there evidence for a peak in this data?



“Observation of an Exotic  $S=+1$

Baryon in Exclusive Photoproduction from the Deuteron”

S. Stepanyan et al, CLAS Collab, Phys.Rev.Lett. 91 (2003) 252001

“The statistical significance of the peak is  $5.2 \pm 0.6 \sigma$ ”

“A Bayesian analysis of pentaquark signals from CLAS data”

D. G. Ireland et al, CLAS Collab, Phys. Rev. Lett. 100, 052001 (2008)

“The  $\ln(\text{RE})$  value for g2a (-0.408) indicates weak evidence in favour of the data model without a peak in the spectrum.”

Comment on “Bayesian Analysis of Pentaquark Signals from CLAS Data”

Bob Cousins, <http://arxiv.org/abs/0807.1330>

# Statistical Issues in Searches for New Physics

Louis Lyons

Imperial College, London  
and  
Oxford

Theme: Using data to make judgements about H1 (New Physics) versus H0 (S.M. with nothing new)

Why?

Experiments are expensive and time-consuming

so

Worth investing effort in statistical analysis

→ better information from data

Topics:

Blind Analysis

LEE = Look Elsewhere Effect

Why  $5\sigma$  for discovery?

Significance

$P(A|B) \neq P(B|A)$

Meaning of p-values

Wilks' Theorem

Background Systematics

Coverage

$p_0$  v  $p_1$  plots

Upper Limyts

Higgs search: Discovery and spin

(N.B. Several of these topics have no unique solutions from Statisticians)

Conclusions

# H0 or H0 versus H1 ?

H0 = null hypothesis

e.g. Standard Model, with nothing new

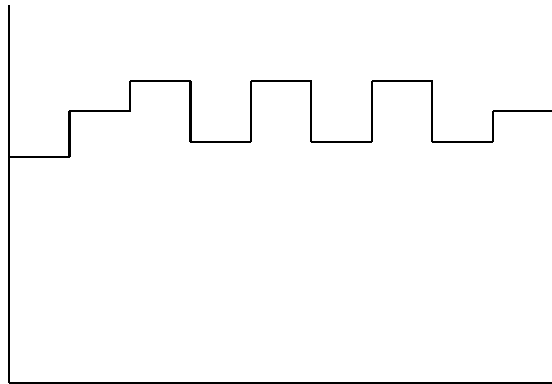
H1 = specific New Physics e.g. Higgs with  $M_H = 125$  GeV

H0: “Goodness of Fit” e.g.  $\chi^2$ , p-values

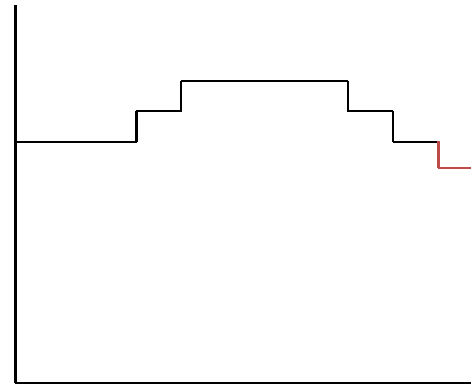
H0 v H1: “Hypothesis Testing” e.g.  $\mathcal{L}$ -ratio

Measures how much data favours one hypothesis wrt other

H0 v H1 likely to be more sensitive



or



# Examples of Hypotheses

## 1) Event selector (Event = particle interaction)

Events produced at CERN LHC at enormous rate

Online 'trigger' to select events for recording ( $\sim 1$  kiloHertz)

e.g. events with many particles

Offline selection based on required features

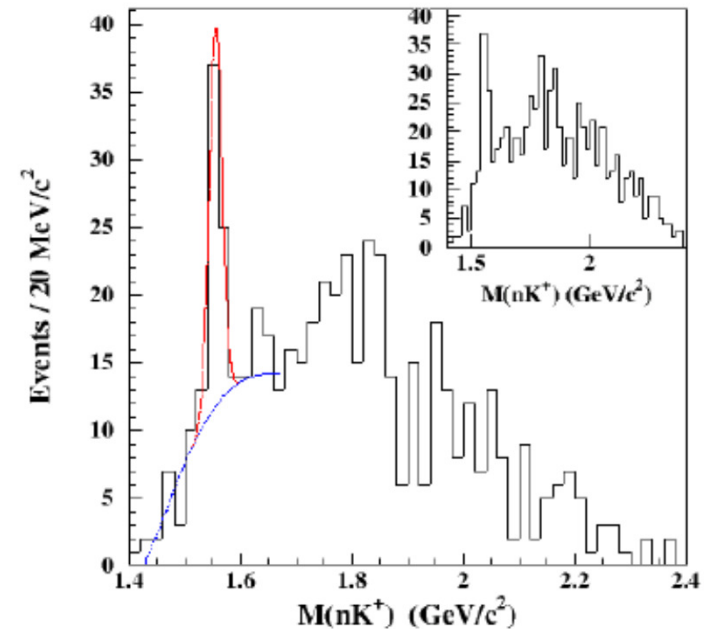
e.g. H0: At least 2 muons H1: 0 or 1 muon

Possible outcomes: Events assigned as H0 or H1

## 2) Result of experiment

e.g. H0 = nothing new, just b  
 H1 = new particle produced as well, b+s  
 (Higgs, SUSY, 4<sup>th</sup> neutrino,.....)

Possible outcomes	H0	H1	
	✓	X	Exclude H1
	X	✓	Discovery
	✓	✓	No decision
	X	X	?

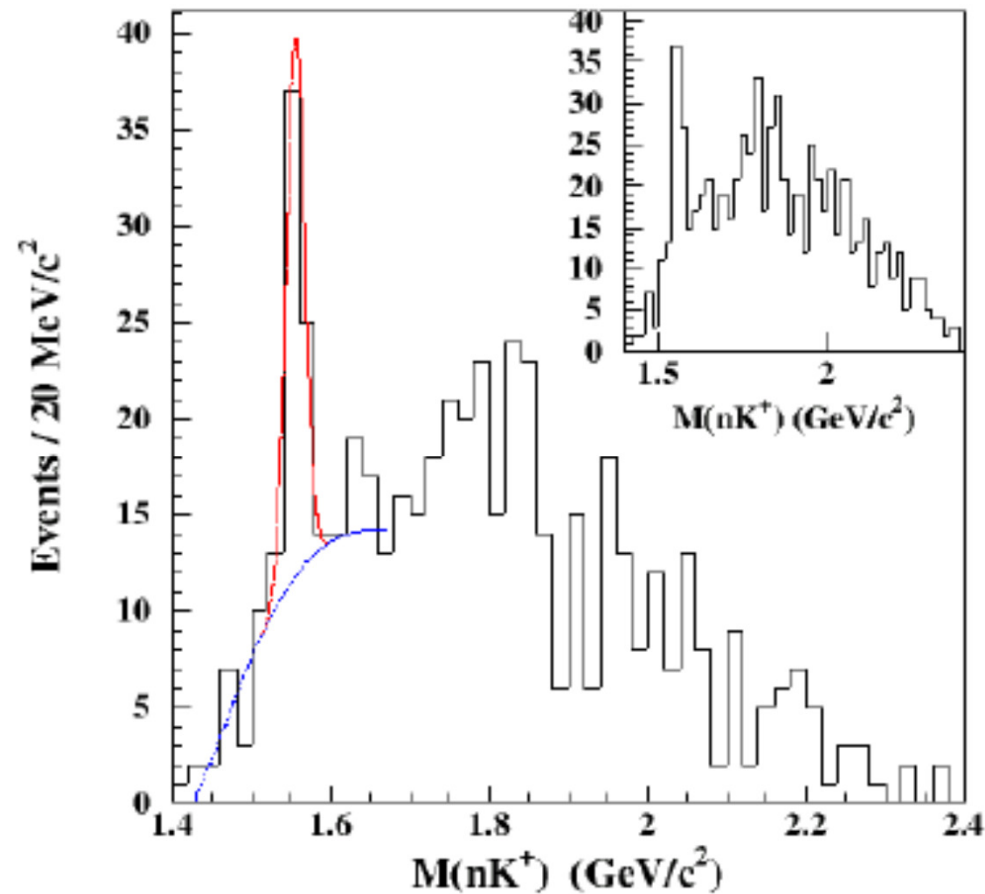


# Choosing between 2 hypotheses

Hypothesis testing: New particle or statistical fluctuation?

$$H_0 = b$$

$$H_1 = b + s$$





# Choosing between 2 hypotheses

Possible methods:

$\Delta\chi^2$

p-value of statistic →

$\ln\mathcal{L}$ -ratio

Bayesian:

Posterior odds

Bayes factor

Bayes information criterion (BIC)

Akaike ..... (AIC)

Minimise “cost”

See ‘Comparing two hypotheses’

<http://www-cdf.fnal.gov/physics/statistics/notes/H0H1.pdf>

# Bayesian methods?

Particle Physicists like Frequentism.

“Avoid personal beliefs”

“Let the data speak for themselves”

For parameter  $\phi$  determination,

a) Range  $\Delta$  of prior for  $\phi$  unimportant, provided.....

b) Bayes' upper limit for Poisson rate (constant prior) agrees with Frequentist UL

c) Easier to incorporate systematics / nuisance parameters

BUT for Hypothesis Testing (e.g. smooth background versus bgd + peak),

$\Delta$  does not cancel, so posterior odds etc. depend on  $\Delta$ .

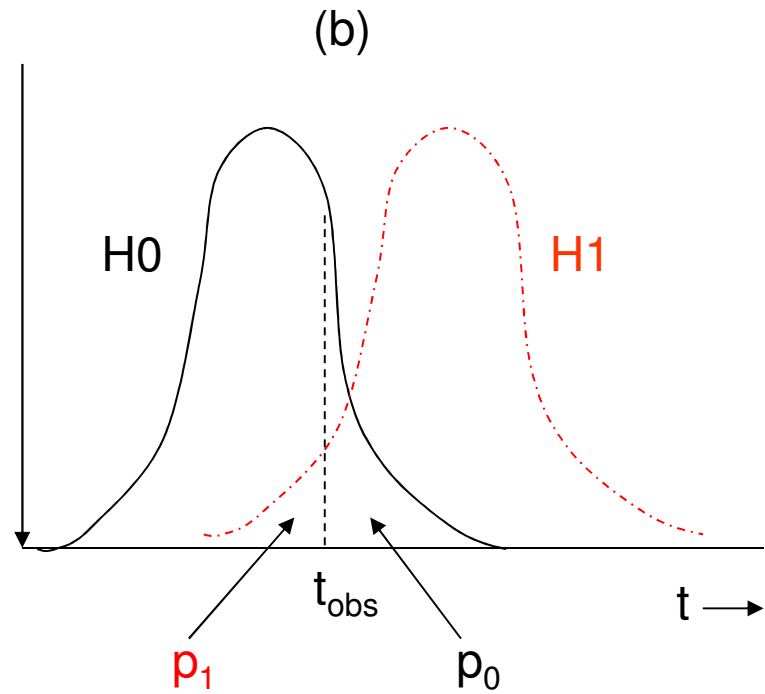
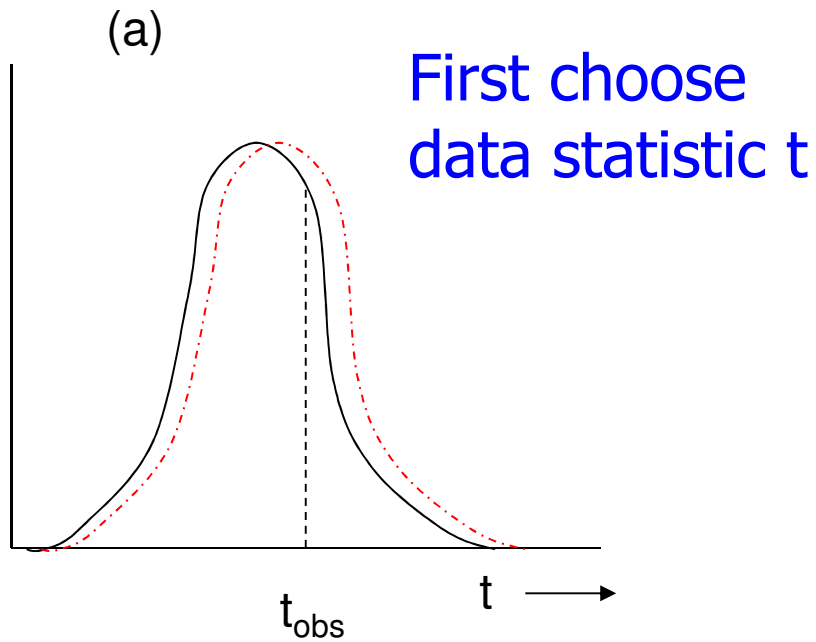
Compare: Frequentist local  $p \rightarrow$  global  $p$  via Look Elsewhere Effect

Effects similar but not same (not surprising):

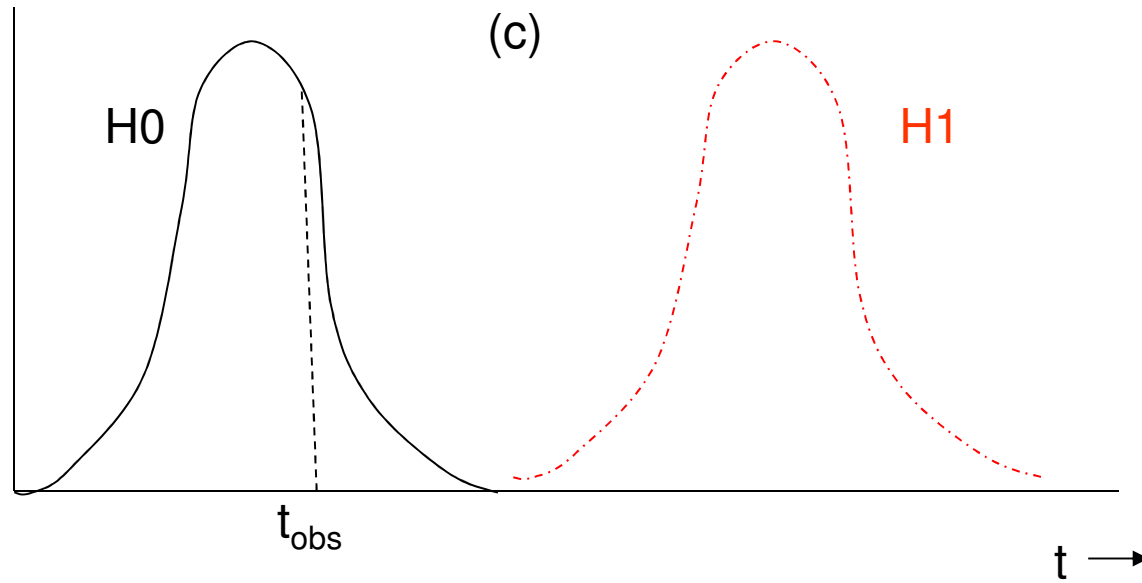
LEE also can allow for selection options, different plots, etc.

Bayesian prior also includes signal strength prior (unless  $\delta$ -function at expected value)

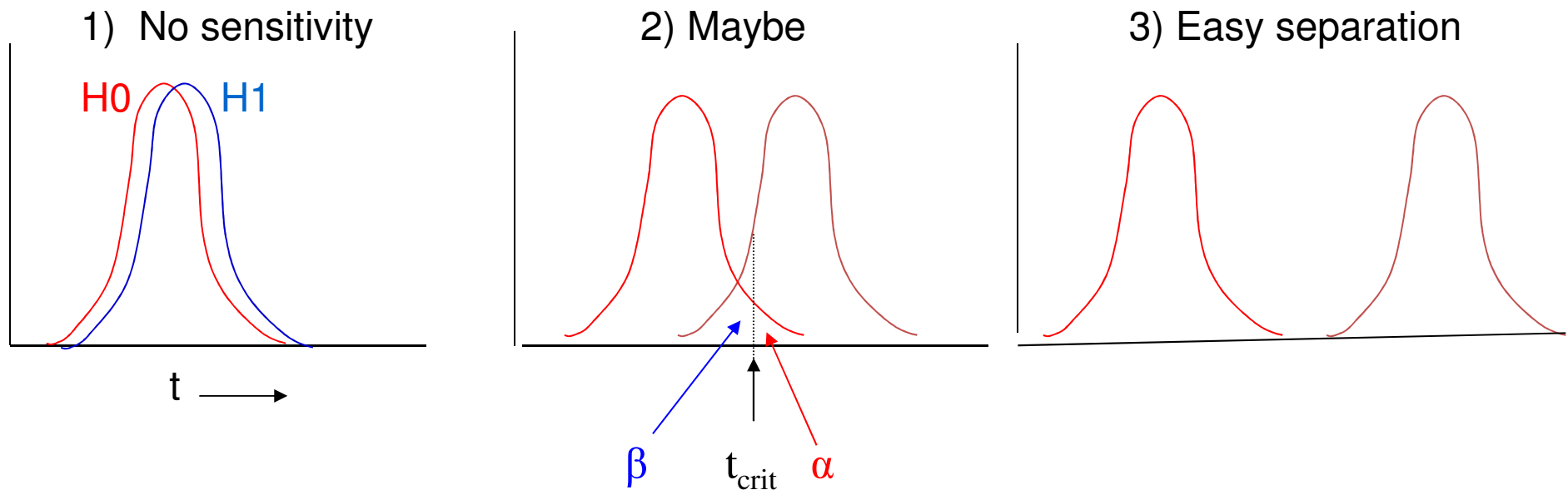
Usually we don't perform prior sensitivity analysis.



With 2 hypotheses, each with own pdf, p-values are defined as tail areas, pointing in towards each other



# Procedure for choosing between 2 hypotheses



Procedure: Obtain expected distributions for data statistic (e.g.  $\mathcal{L}$ -ratio) for H0 and H1

Choose  $\alpha$  (e.g. 95%,  $3\sigma$ ,  $5\sigma$  ?) and CL for  $p_1$  (e.g. 95%)

Given  $b$ ,  $\alpha$  determines  $t_{\text{crit}}$

$b+s$  defines  $\beta$ . For  $s > s_{\text{min}}$ , separation of curves → discovery or excln

$1-\beta$  = Power of test

Now data: If  $t_{\text{obs}} \geq t_{\text{crit}}$  (i.e.  $p_0 \leq \alpha$ ), **discovery at level  $\alpha$**

If  $t_{\text{obs}} < t_{\text{crit}}$ , no discovery. If  $p_1 < 1 - \text{CL}$ , **exclude H1**

**Slide 12**

---

**N1**

NPL, 06/11/2005

# BLIND ANALYSES

**Why blind analysis?** Data statistic, selections, corrections, method  
Dunnington (1932) e/m with detector location hidden

## Methods of blinding

- Add random number to result \*
- Study procedure with simulation only
- Look at only first fraction of data
- Keep the signal box closed
- Keep MC parameters hidden
- Keep unknown fraction visible for each bin

## Disadvantages

- Takes longer time
- Usually not available for searches for unknown

After analysis is unblinded, don't change anything unless .....

\* Luis Alvarez suggestion re “discovery” of free quarks

See Klein and Roodman review: ARNPS 55 (2005) 141

# Look Elsewhere Effect (LEE)

Prob of bgd fluctuation at that place = local p-value

Prob of bgd fluctuation 'anywhere' = global p-value

Global p > Local p

Where is 'anywhere'?

- a) Any location in this histogram in sensible range
  - b) Any location in this histogram
  - c) Also in histogram produced with different cuts, binning, etc.
  - d) Also in other plausible histograms for this analysis
  - e) Also in other searches in this PHYSICS group (e.g. SUSY at CMS)
  - f) In any search in this experiment (e.g. CMS)
  - g) In all CERN expts (e.g. LHC expts + NA62 + OPERA + ASACUSA + ....)
  - h) In all HEP expts
- etc.

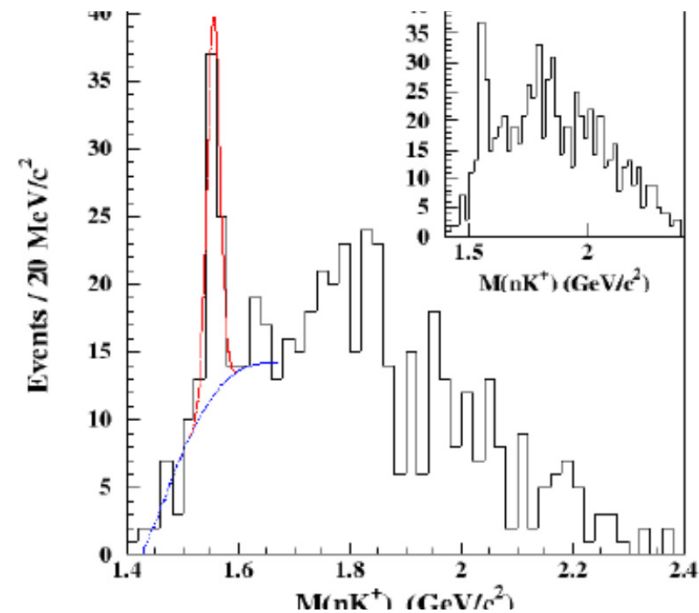
d) relevant for graduate student doing analysis

f) relevant for experiment's Spokesperson

## INFORMAL CONSENSUS:

Quote local p, and global p according to a) above.

Explain which global p



# Example of LEE: Stonehenge





12 is the number of constellations

6 is the number of ages (2160) we spend on each side of the galactic equator

18 number of breaths we take each minute or our life

Missing two large stones in top half. Should be 6 and 6

*Alpha Draconis*

NORTH

*Beta Ursa Minor*

If small stones = 432 years each then the half circle in the center would be  
 $20 \times 432 = 8640$  years  
 $8640 \text{ divided by } 2160 = 4\text{th time.}$

30 Stones in Outer ring =  
 $360 \text{ divided by } 30 = 12$

60 Stones in Second ring =  
 $360 \text{ divided by } 60 = 6$

20 Stones in Center ring =  
 $360 \text{ divided by } 20 = 18$

IF THIS WAS EAST

WINTER SOLSTICE

SUMMER SOLSTICE

WEST  
BALANCED LOCATION IN SPACE

STONEHENGE

The Book of Truth  
A New Perspective on the Hopi Creation Story  
by Thomas O. Mills

Stonehenge from a Hopi point of view.

Doesn't make sense with today's eastward direction.

1 degree = 72 years  
 $360 \times 72 = 25,920$

*Sirius*

TODAY'S EAST

SOUTH

*Zeta Orionis*

$25,920 \text{ divided by } 60 = 432$   
 $432 \times 5 = 2,160$   
Should be 5 stones between each division on the Second ring.

$25,920 \text{ divided by } 12 = 2160$

$25,920 \text{ divided by } 6 = 4320$

$25,920 \text{ divided by } 18 = 1440$

Center Stone in Center Ring would be divided in half by sun rays when Earth in perfect balance.  
Nine on each side + 2 = 20.

# Are alignments significant?

- Atkinson replied with his article "Moonshine on Stonehenge" in [Antiquity](#) in 1966, pointing out that some of the pits which ..... had used for his sight lines were more likely to have been natural depressions, and that he had allowed a margin of error of up to 2 degrees in his alignments. Atkinson found that the probability of so many alignments being visible from 165 points to be close to 0.5 rather than the "one in a million" possibility which ..... had claimed.
- ..... had been examining stone circles since the 1950s in search of astronomical alignments and the [megalithic yard](#). It was not until 1973 that he turned his attention to Stonehenge. He chose to ignore alignments between features within the monument, considering them to be too close together to be reliable. He looked for landscape features that could have marked lunar and solar events. However, one of ..... 's key sites, Peter's Mound, turned out to be a twentieth-century rubbish dump.

# Why 5 $\sigma$ for Discovery?

Statisticians ridicule our belief in extreme tails (esp. for systematics)

Our reasons:

1) Past history (Many 3 $\sigma$  and 4 $\sigma$  effects have gone away)

2) LEE

3) Worries about underestimated systematics

4) Subconscious Bayes calculation

$$\frac{p(H_1|x)}{p(H_0|x)} = \frac{p(x|H_1)}{p(x|H_0)} * \frac{\pi(H_1)}{\pi(H_0)}$$

Posterior      Likelihood      Priors  
ratio

**“Extraordinary claims require extraordinary evidence”**

N.B. Points 2), 3) and 4) are experiment-dependent

Alternative suggestion:

L.L. “Discovering the significance of 5 $\sigma$ ”

<http://arxiv.org/abs/1310.1284>

## How many $\sigma$ 's for discovery?

SEARCH	SURPRISE	IMPACT	LEE	SYSTEMATICS	No. $\sigma$
Higgs search	Medium	Very high	M	Medium	5
Single top	No	Low	No	No	3
SUSY	Yes	Very high	Very large	Yes	7
$B_s$ oscillations	Medium/Low	Medium	$\Delta m$	No	4
Neutrino osc	Medium	High	$\sin^2 2\theta, \Delta m^2$	No	4
$B_s \rightarrow \mu \mu$	No	Low/Medium	No	Medium	3
Pentaquark	Yes	High/V. high	M, decay mode	Medium	7
$(g-2)_\mu$ anom	Yes	High	No	Yes	4
H spin $\neq 0$	Yes	High	No	Medium	5
4 <sup>th</sup> gen q, l, $\nu$	Yes	High	M, mode	No	6
Dark energy	Yes	Very high	Strength	Yes	5
Grav Waves	No	High	Enormous	Yes	8

Suggestions to provoke discussion, rather than 'carved in stone on Mt. Sinai'



Bob Cousins: "2 independent expts each with  $3.5\sigma$  better than one expt with  $5\sigma$ "

David van Dyk: "A calibrated  $3.5\sigma$  experiment is better than a  $5\sigma$  uncalibrated one"

## Significance

$$\text{Significance} = S/\sqrt{B} \quad \text{or similar ?}$$

### Potential Problems:

- Uncertainty in B
- Non-Gaussian behaviour of Poisson, especially in tail
- Number of bins in histogram, no. of other histograms [LEE]
- Choice of cuts, bins (Blind analyses)

### For future experiments:

- Optimising: Could give  $S = 0.1$ ,  $B = 10^{-4}$ ,  $S/\sqrt{B} = 10$

$$P(A | B) \neq P(B | A)$$

Remind Lab or University media contact person that:

Prob[data, given H0] is very small

does not imply that

Prob[H0, given data] is also very small.

e.g. Prob{data | speed of  $v \leq c$ } = very small

does not imply

Prob{speed of  $v \leq c$  | data} = very small

or Prob{speed of  $v > c$  | data}  $\sim 1$

Everyday situation, my granddaughter's example:

$p(\text{bread for breakfast} | \text{murderer}) \sim 95\%$

$p(\text{murderer} | \text{bread for breakfast}) \sim 10^{-6}$

$$P(A | B) \neq P(B | A)$$

Remind Lab or University media contact person that:

Prob[data, given H0] is very small

does not imply that

Prob[H0, given data] is also very small.

e.g. Prob{data | speed of  $v \leq c$ } = very small

does not imply

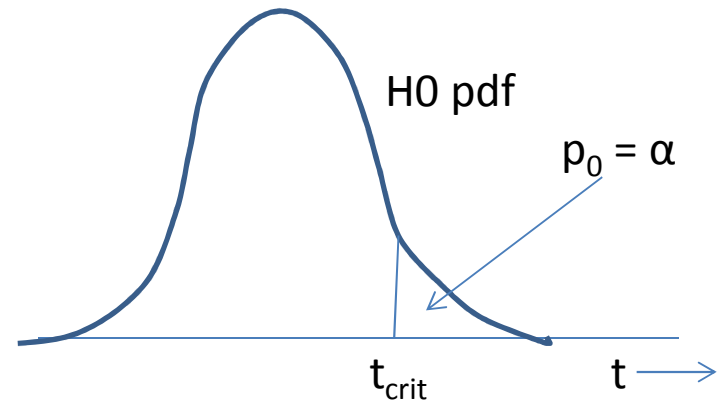
Prob{speed of  $v \leq c$  | data} = very small

or Prob{speed of  $v > c$  | data}  $\sim 1$

Everyday example  $p(\text{pregnant} | \text{female}) \sim 3\%$

$p(\text{female} | \text{pregnant}) \gg 3\%$

# What p-values are (and are not)



Reject  $H_0$  if  $t > t_{\text{crit}}$  ( $p < \alpha$ )

p-value = prob that  $t \geq t_{\text{obs}}$

Small  $p \rightarrow$  data and theory have poor compatibility

Small p-value does **NOT** automatically imply that theory is unlikely

Bayes  $\text{prob}(\text{Theory}|\text{data})$  related to  $\text{prob}(\text{data}|\text{Theory}) = \text{Likelihood}$

by Bayes Th, including Bayesian prior

p-values are misunderstood. e.g. Anti-HEP jibe:

“Particle Physicists don’t know what they are doing, because half their  $p < 0.05$  exclusions turn out to be wrong”

Demonstrates lack of understanding of p-values

[**All** results rejecting energy conservation with  $p < \alpha = .05$  cut will turn out to be ‘wrong’]



# Are p-values useful?

Particle Physicists use p-values for exclusion and for discovery

Have come in for strong criticism:

People think it is  $\text{prob}(\text{theory} | \text{data})$

p-values over-emphasize evidence (much smaller than  $\mathcal{L}$ -ratio)

Over 50% of results with  $p_0 < 5\%$  are wrong

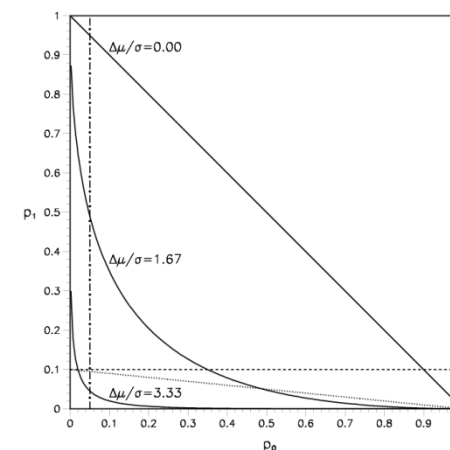
In Particle Physics, we use  $\mathcal{L}$ -ratio as data statistic for p-values

Can regard this as:

p-value method, which just happens to use  $\mathcal{L}$ -ratio as test statistic

or

This is a  $\mathcal{L}$ -ratio method with p-values used as calibration



# Are p-values useful?

Particle Physicists use p-values for exclusion and for discovery

Have come in for strong criticism:

People think it is  $\text{prob}(\text{theory}|\text{data})$  **Stop using relativity, because misunderstood?**  
p-values over-emphasize evidence (much smaller than  $\mathcal{L}$ -ratio)

**Is mass or height `better` for sizes of mice and elephants?**

Over 50% of results with  $p_0 < 5\%$  are wrong **Confusing  $p(A;B)$  with  $p(B;A)$**

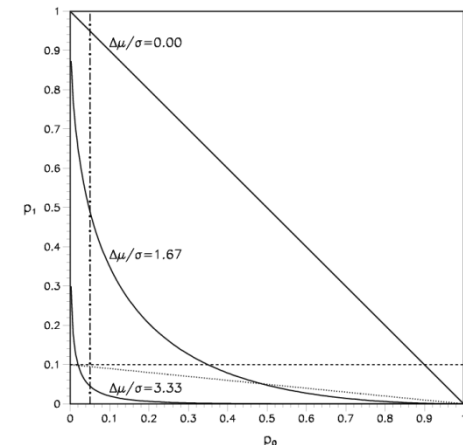
In Particle Physics, we use  $\mathcal{L}$ -ratio as data statistic for p-values

Can regard this as:

p-value method, which just happens to use  $\mathcal{L}$ -ratio as test statistic

or

This is a  $\mathcal{L}$ -ratio method with p-values used as calibration



# $p_0$ v $p_1$ plots

Preprint by Luc Demortier and LL,  
“Testing Hypotheses in Particle Physics:  
Plots of  $p_0$  versus  $p_1$ ”  
<http://arxiv.org/abs/1408.6123>

For hypotheses  $H_0$  and  $H_1$ ,  $p_0$  and  $p_1$   
are the tail probabilities for data  
statistic  $t$

Provide insights on:

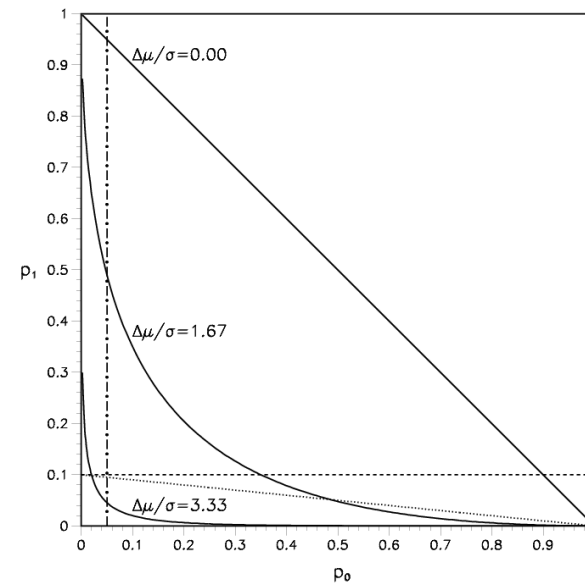
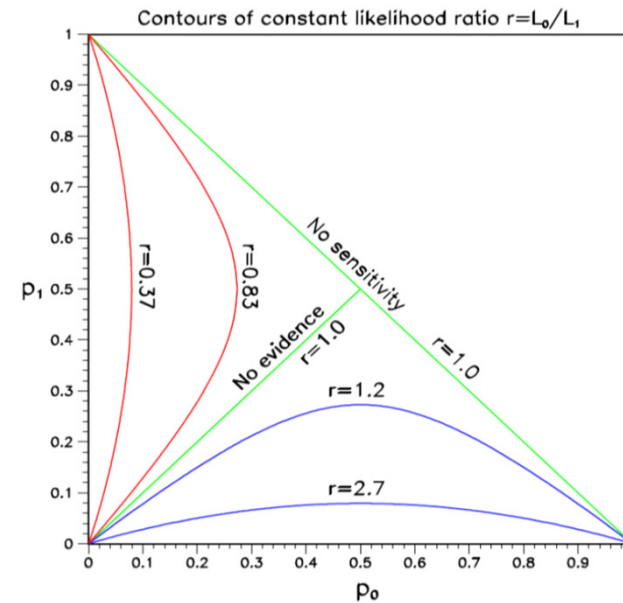
CLs for exclusion

Punzi definition of sensitivity

Relation of p-values and Likelihoods

Probability of misleading evidence

Jeffreys-Lindley paradox

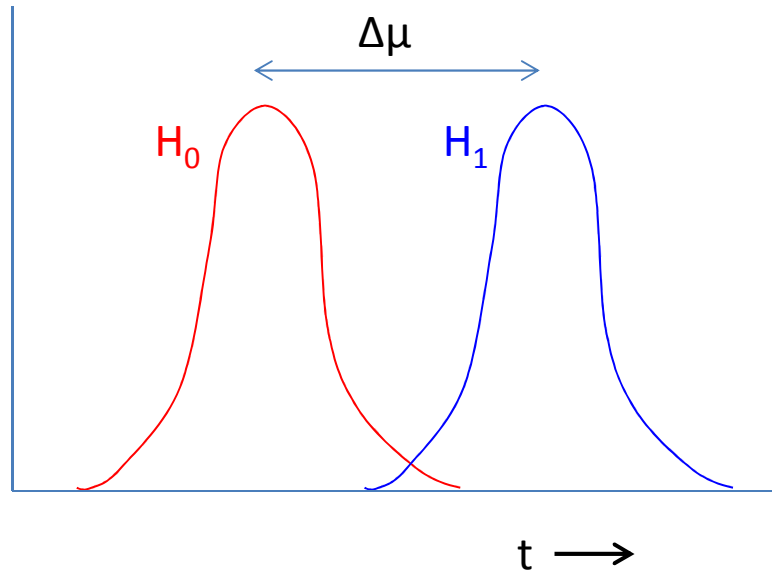


CLs =  $p_1/(1-p_0)$   $\rightarrow$  diagonal line

Provides protection against excluding  $H_1$  when little or no sensitivity

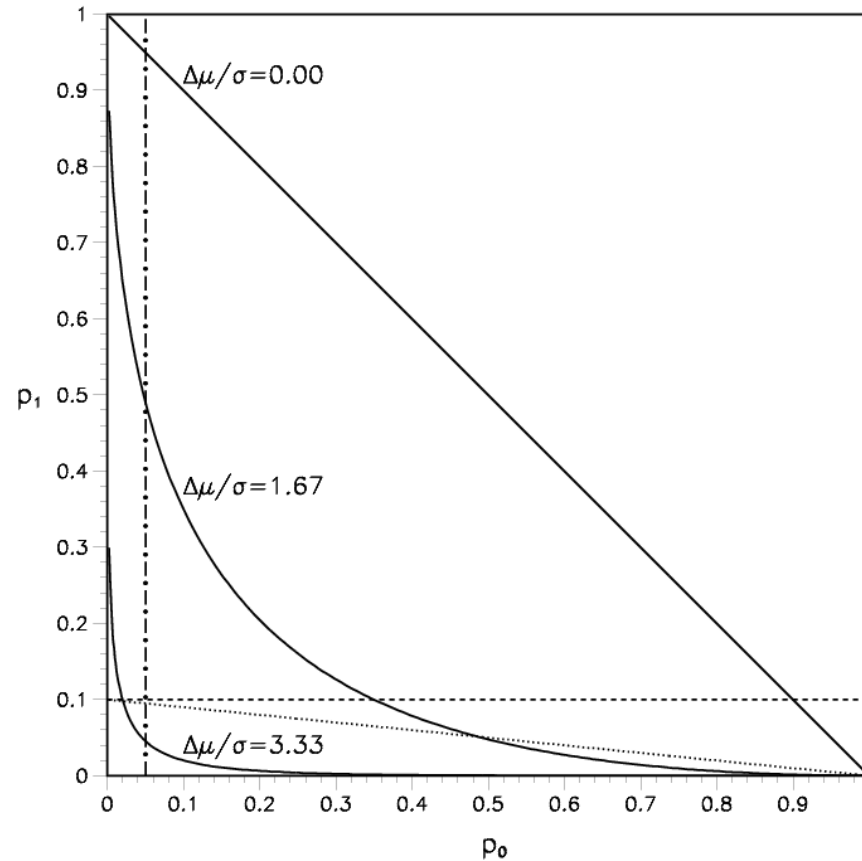
Punzi definition of sensitivity:

Enough separation of pdf's for no chance of ambiguity



Can read off power of test  
e.g. If  $H_0$  is true, what is  
prob of rejecting  $H_1$ ?

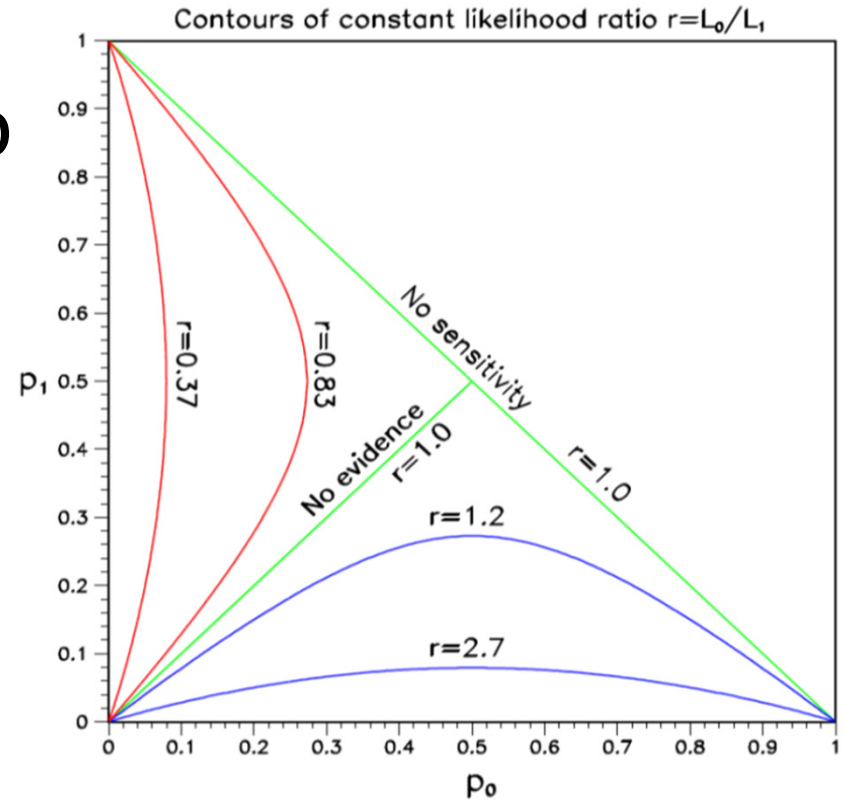
**N.B.  $p_0$  = tail towards  $H_1$**   
 **$p_1$  = tail towards  $H_0$**



# Why $p \neq$ Likelihood ratio

Measure different things:

$p_0$  refers just to  $H_0$ ;  $\mathcal{L}_{01}$  compares  $H_0$  and  $H_1$



Depends on amount of data:

e.g. Poisson counting expt little data:

For  $H_0$ ,  $\mu_0 = 1.0$ . For  $H_1$ ,  $\mu_1 = 10.0$

Observe  $n = 10$   $p_0 \sim 10^{-7}$   $\mathcal{L}_{01} \sim 10^{-5}$

Now with 100 times as much data,  $\mu_0 = 100.0$   $\mu_1 = 1000.0$

Observe  $n = 160$   $p_0 \sim 10^{-7}$   $\mathcal{L}_{01} \sim 10^{+14}$

# Jeffreys-Lindley Paradox

$H_0$  = simple,  $H_1$  has  $\mu$  free  
 $p_0$  can favour  $H_1$ , while  $B_{01}$  can favour  $H_0$   
 $B_{01} = L_0 / \int L_1(s) \pi(s) ds$

Likelihood ratio depends on signal :  
 e.g. Poisson counting expt small signal s:

For  $H_0$ ,  $\mu_0 = 1.0$ . For  $H_1$ ,  $\mu_1 = 10.0$

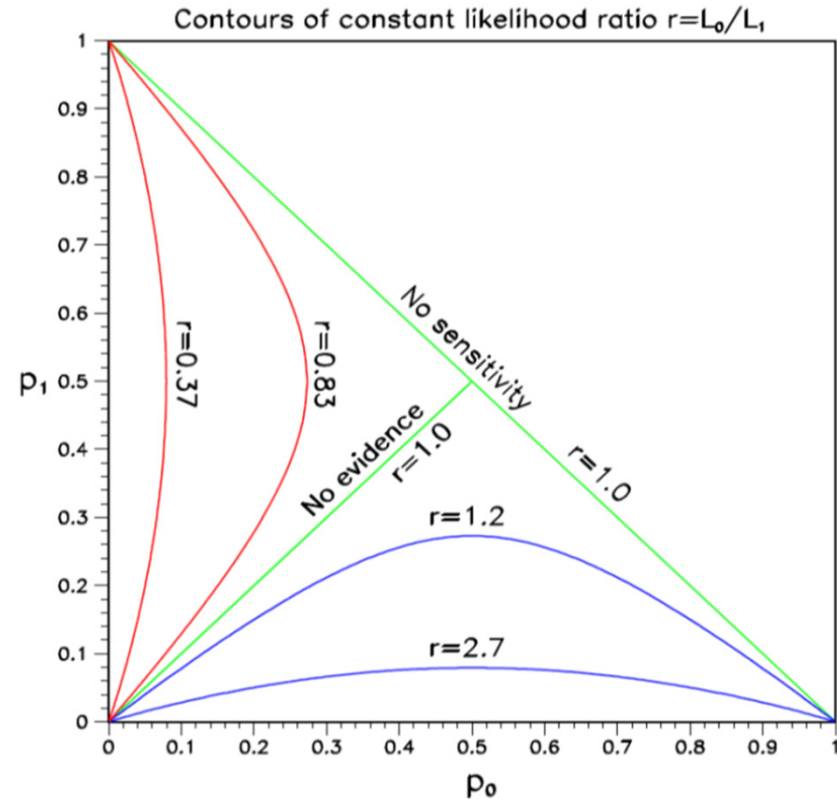
Observe  $n = 10$   $p_0 \sim 10^{-7}$   $L_{01} \sim 10^{-5}$  and favours  $H_1$

Now with 100 times as much signal s,  $\mu_0 = 100.0$   $\mu_1 = 1000.0$

Observe  $n = 160$   $p_0 \sim 10^{-7}$   $L_{01} \sim 10^{+14}$  and favours  $H_0$

$B_{01}$  involves intergration over s in denominator, so a wide enough range  
 will result in favouring  $H_0$

However, for  $B_{01}$  to favour  $H_0$  when  $p_0$  is equivalent to  $5\sigma$ , integration  
 range for s has to be  $O(10^6)$  times Gaussian widths



# Combining different p-values

Several results quote independent p-values for same effect:

$p_1, p_2, p_3, \dots$  e.g. 0.9, 0.001, 0.3 .....

What is combined significance? Not just  $p_1 * p_2 * p_3, \dots$

If 10 expts each have  $p \sim 0.5$ , product  $\sim 0.001$  and is clearly **NOT** correct  
combined p

$$S = z * \sum_{j=0}^{n-1} (-\ln z)^j / j! , \quad z = p_1 p_2 p_3 \dots$$

(e.g. For 2 measurements,  $S = z * (1 - \ln z) \geq z$  )

Problems:

**1) Recipe is not unique (Uniform dist in n-D hypercube  $\rightarrow$  uniform in 1-D)**

**2) Formula is not associative**

Combining  $\{p_1$  and  $p_2\}$ , and then  $p_3\}$  gives different answer  
from  $\{p_3$  and  $p_2\}$ , and then  $p_1\}$  , or all together

Due to different options for “more extreme than  $x_1, x_2, x_3$ ”.

**3) Small p's due to different discrepancies**

\*\*\*\*\* Better to combine data \*\*\*\*\*

# Wilks' Theorem

Data = some distribution e.g. mass histogram

For  $H_0$  and  $H_1$ , calculate best fit weighted sum of squares  $S_0$  and  $S_1$

Examples: 1)  $H_0$  = polynomial of degree 3

$H_1$  = polynomial of degree 5

2)  $H_0$  = background only

$H_1$  = bgd+peak with free  $M_0$  and cross-section

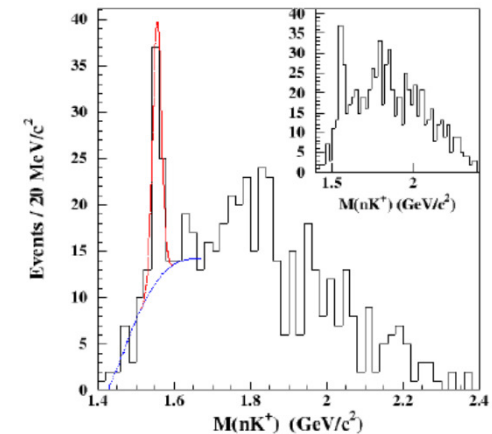
3)  $H_0$  = normal neutrino hierarchy

$H_1$  = inverted hierarchy

If  $H_0$  true,  $S_0$  distributed as  $\chi^2$  with  $\text{ndf} = \nu_0$

If  $H_1$  true,  $S_1$  distributed as  $\chi^2$  with  $\text{ndf} = \nu_1$

If  $H_0$  true, what is distribution of  $\Delta S = S_0 - S_1$ ? Expect not large. Is it  $\chi^2$ ?



**Wilks' Theorem:**  $\Delta S$  distributed as  $\chi^2$  with  $\text{ndf} = \nu_0 - \nu_1$  provided:

a)  $H_0$  is true

b)  $H_0$  and  $H_1$  are nested

c) Params for  $H_1 \rightarrow H_0$  are well defined, and not on boundary

d) Data is asymptotic



# Wilks' Theorem, contd

Examples: Does Wilks' Th apply?

1)  $H_0$  = polynomial of degree 3

$H_1$  = polynomial of degree 5

**YES:  $\Delta S$  distributed as  $\chi^2$  with ndf =  $(d-4) - (d-6) = 2$**

2)  $H_0$  = background only

$H_1$  = bgd + peak with free  $M_0$  and cross-section

**NO:  $H_0$  and  $H_1$  nested, but  $M_0$  undefined when  $H_1 \rightarrow H_0$ .  $\Delta S \neq \chi^2$   
(but not too serious for fixed  $M$ )**

3)  $H_0$  = normal neutrino hierarchy

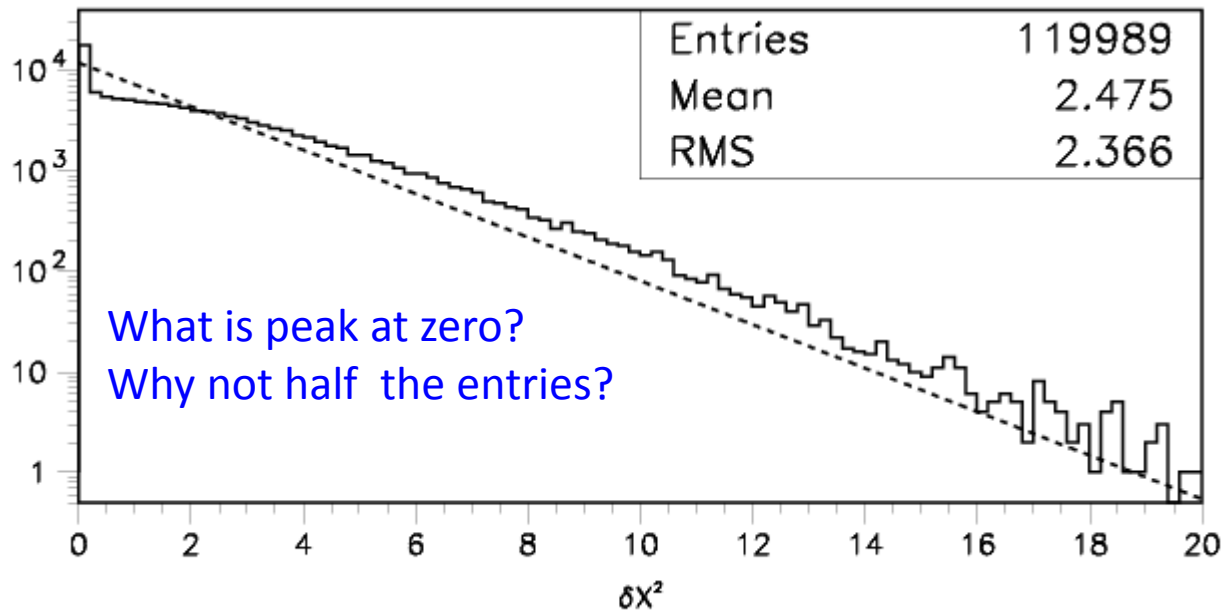
$H_1$  = inverted hierarchy

**NO: Not nested.  $\Delta S \neq \chi^2$  (e.g. can have  $\Delta\chi^2$  negative)**

N.B. 1: Even when **W. Th.** does not apply, it does not mean that  $\Delta S$  is irrelevant, but you cannot use **W. Th.** for its expected distribution.

N.B. 2: For large ndf, better to use  $\Delta S$ , rather than  $S_1$  and  $S_0$  separately

# Is difference in S distributed as $\chi^2$ ?

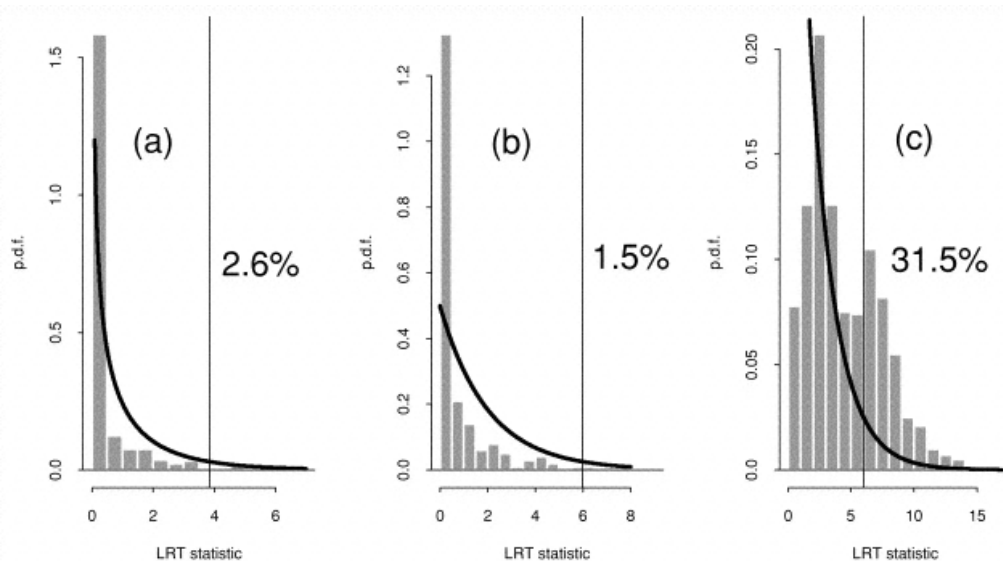


Demortier:

H0 = quadratic bgd

H1 = ..... +

Gaussian of fixed width,  
variable location & ampl



Protassov, van Dyk, Connors, ....

H0 = continuum

(a) H1 = narrow emission line

(b) H1 = wider emission line

(c) H1 = absorption line

Nominal significance level = 5%

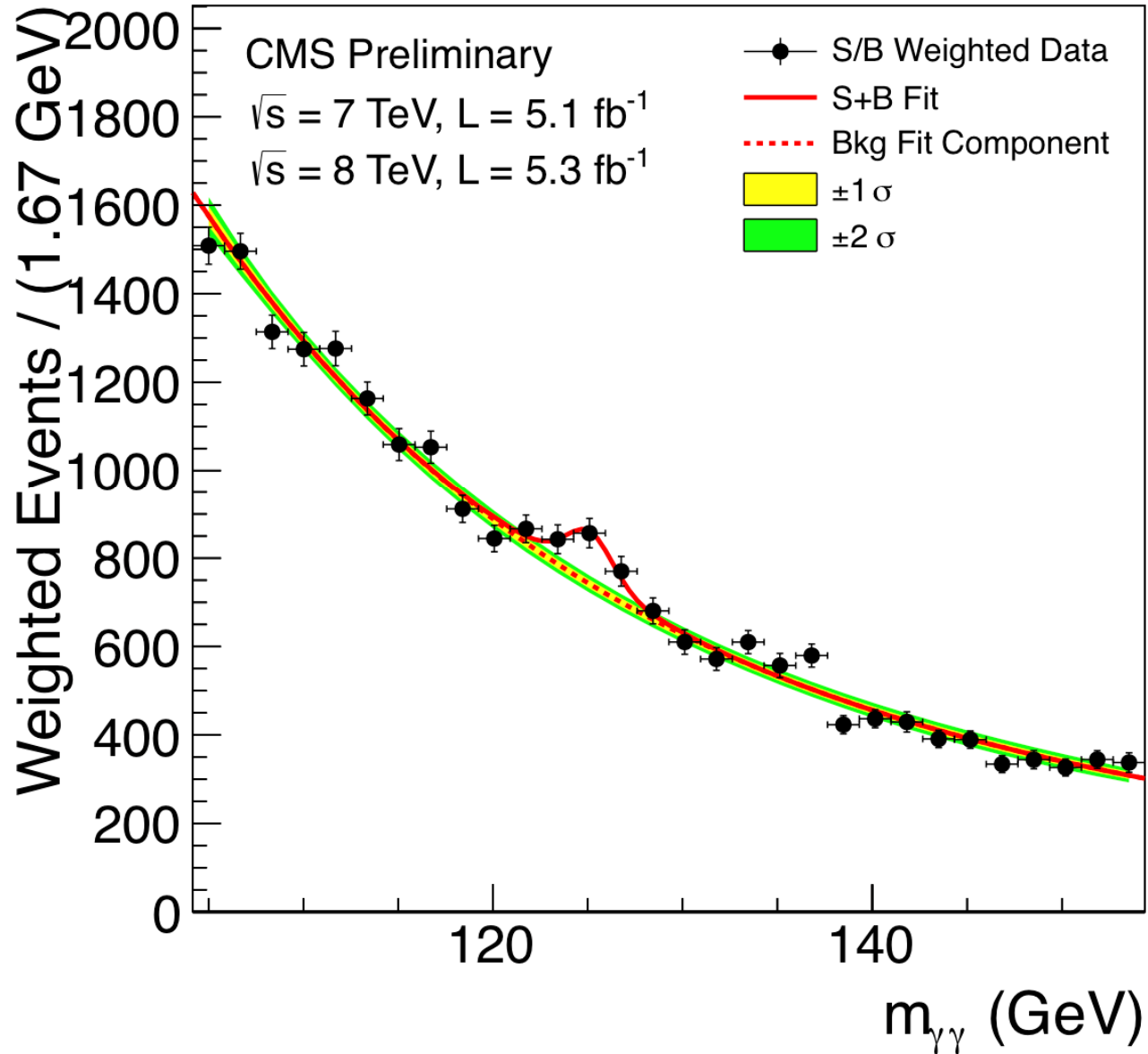
## Is difference in $S$ distributed as $\chi^2$ ?, contd.

So need to determine the  $\Delta S$  distribution by Monte Carlo

N.B.

- 1) For mass spectrum, determining  $\Delta S$  for hypothesis  $H_1$  when data is generated according to  $H_0$  is not trivial, because there will be lots of local minima
- 2) If we are interested in  $5\sigma$  significance level, needs lots of MC simulations (or intelligent MC generation)
- 3) Asymptotic formulae may be useful (see K. Cranmer, G. Cowan, E. Gross and O. Vitells, 'Asymptotic formulae for likelihood-based tests of new physics', <http://link.springer.com/article/10.1140%2Fepjc%2Fs10052-011-1554-0>)

# Background systematics



# Background systematics, contd

Signif from comparing  $\chi^2$ 's for H0 (bgd only) and for H1 (bgd + signal)

Typically, bgd = functional form  $f_a$  with free params

e.g. 4<sup>th</sup> order polynomial

Uncertainties in params included in signif calculation

But what if functional form is different ? e.g.  $f_b$

Typical approach:

If  $f_b$  best fit is bad, not relevant for systematics

If  $f_b$  best fit is ~comparable to  $f_a$  fit, include contribution to systematics

But what is 'comparable'?

Other approaches:

Profile likelihood over different bgd parametric forms

<http://arxiv.org/pdf/1408.6865v1.pdf>

Background subtraction

sPlots

Non-parametric background

Bayes

etc

No common consensus yet among experiments on best approach

{Spectra with multiple peaks are more difficult}

# “Handling uncertainties in background shapes: the discrete profiling method”

Dauncey, Kenzie, Wardle and Davies (Imperial College, CMS)

[arXiv:1408.6865v1](https://arxiv.org/abs/1408.6865v1) [physics.data-an]

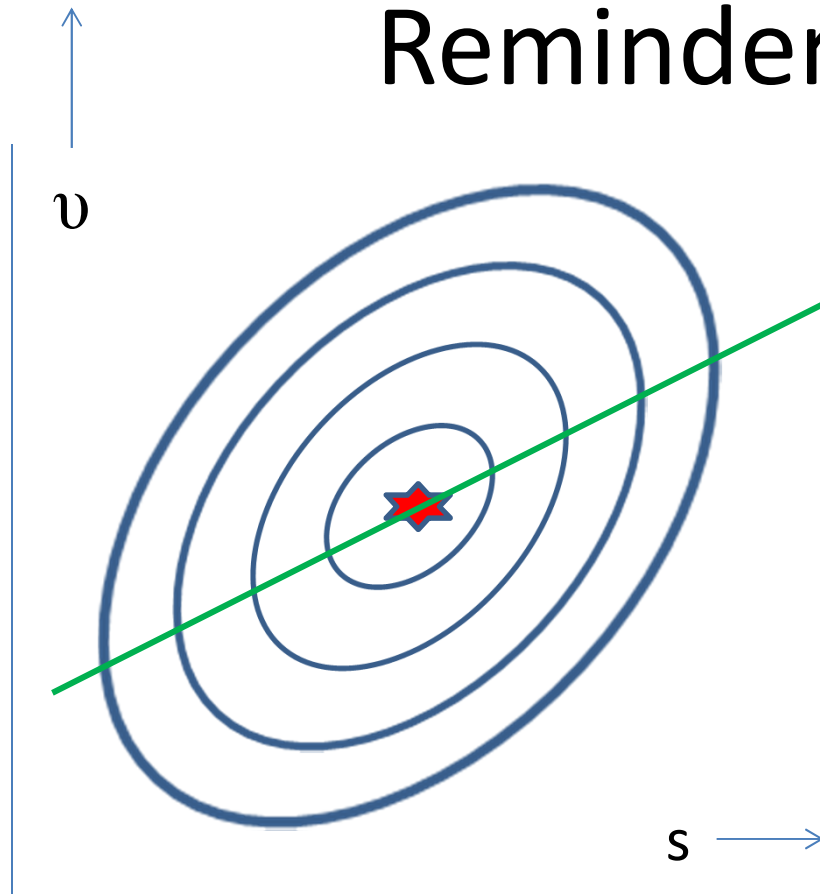
Has been used in CMS analysis of  $H \rightarrow \gamma\gamma$

Problem with ‘Typical approach’: Alternative functional forms do or don’t contribute to systematics by hard cut, so systematics can change discontinuously wrt  $\Delta\chi^2$

Method is like profile  $\mathcal{L}$  for continuous nuisance params

Here ‘profile’ over discrete functional forms

# Reminder of Profile $\mathcal{L}$



Stat uncertainty on  $s$  from width of  $\mathcal{L}$  fixed at  $v_{\text{best}}$

Total uncertainty on  $s$  from width of  $\mathcal{L}(s, v_{\text{prof}(s)}) = \mathcal{L}_{\text{prof}}$

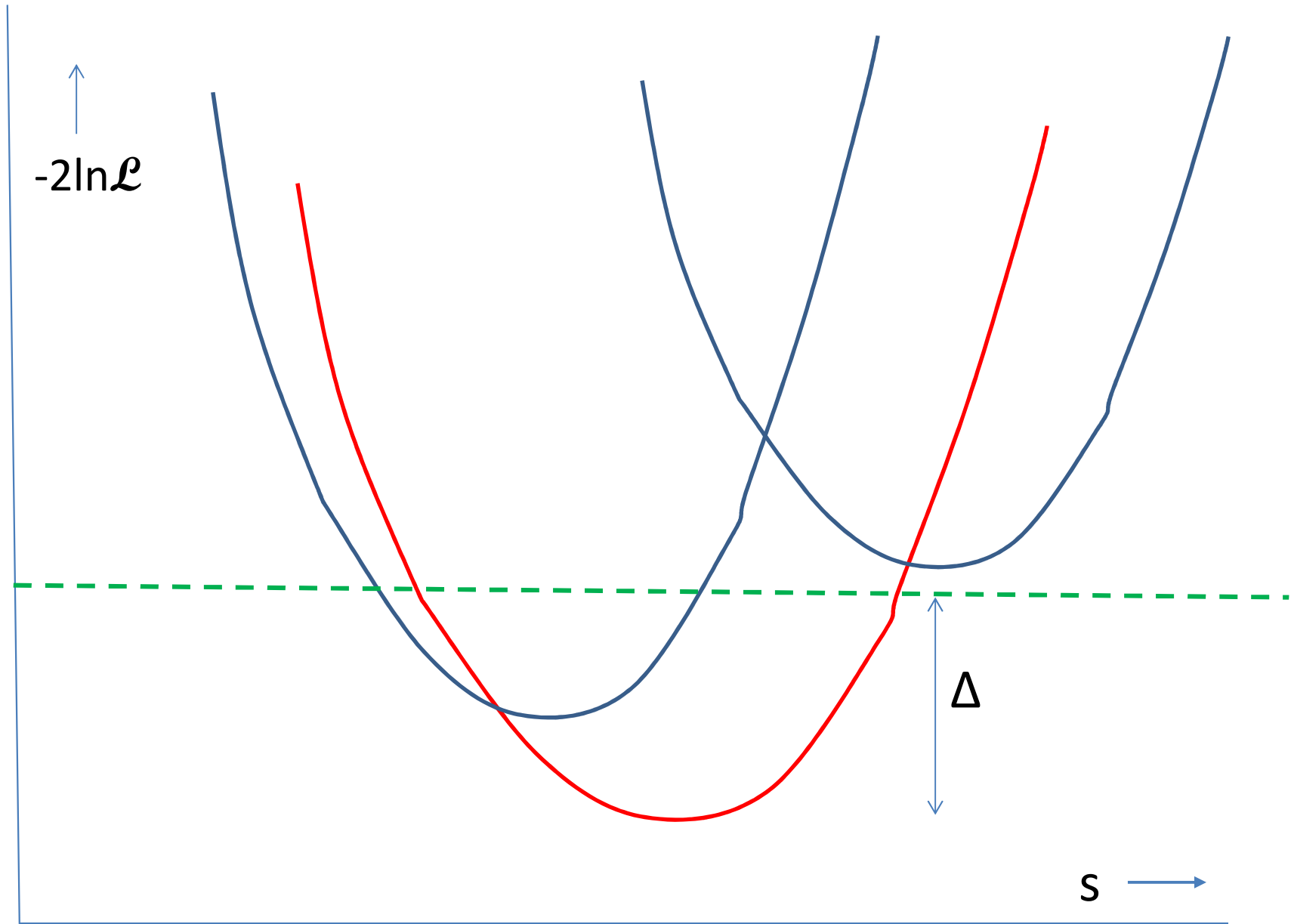
$v_{\text{prof}(s)}$  is best value of  $v$  at that  $s$   
 $v_{\text{prof}(s)}$  as fn of  $s$  lies on green line

Contours of  $\ln \mathcal{L}(s, v)$

$s$  = physics param

$v$  = nuisance param

Total uncert  $\geq$  stat uncertainty





**Red curve:** Best value of nuisance param  $\nu$

**Blue curves:** Other values of  $\nu$

Horizontal line: Intersection with red curve  $\rightarrow$   
statistical uncertainty

‘Typical approach’: Decide which blue curves have small enough  $\Delta$   
Systematic is largest change in minima wrt red curves’.

Profile L: Envelope of lots of blue curves

Wider than red curve, because of systematics ( $\nu$ )

For  $\mathcal{L} =$  multi-D Gaussian, agrees with ‘Typical approach’

Dauncey et al use envelope of finite number of functional forms

Point of controversy!

Two types of 'other functions':

a) Different function types e.g.

$$\sum a_i x_i \text{ versus } \sum a_i / x_i$$

b) Given fn form but different number of terms

DDKW deal with b) by  $-2\ln L \rightarrow -2\ln L + kn$

$n$  = number of extra free params wrt best

$k = 1$ , as in AIC (= Akaike Information Criterion)

Opposition claim choice  $k=1$  is arbitrary.

DDKW agree but have studied different values, and say  $k = 1$  is optimal for them.

Also, any parametric method needs to make such a choice

# Example of misleading inference

Ofer Vitells, Weizmann Institute PhD thesis (2014)

On-off problem (signal + bgd, bgd only)

e.g.  $n_{\text{on}} = 10$ ,  $m_{\text{off}} = 0$

i.e. convincing evidence for signal

Now, to improve analysis, look at spectra of events (e.g. in mass) in “on” and “off” regions

e.g. Use 100 narrow bins  $\rightarrow n_i = 1$  for 10 bins,  $m_i = 0$  for all bins

Assume bins are chosen so that signal expectation  $s_i$  is uniform in all bins  
but bgd  $b_i$  is unknown

$$\text{Likelihood: } \mathcal{L}(s, b_i) = e^{-Ks} e^{-(1+\tau)\sum b_i} \prod_j (s+b_j)$$

$K$  = number of bins (e.g. 100)

$\tau$  = scale factor for bgd (e.g. 1)

$j$  = "on" bins with event (e.g. 1..... 10)

Profile over background nuisance params  $b_i$

$\mathcal{L}_{\text{prof}}(s)$  has largest value at

$s=0$  if  $n_{\text{on}} < K/(1+\tau)$

$s=n_{\text{on}}/K$  if  $n_{\text{on}} \geq K/(1+\tau)$

Similar result for Bayesian marginalisation of  $\mathcal{L}(s, b_i)$  over backgrounds  $b_i$

i.e. With many bins, profile (or marginalised)  $\mathcal{L}$  has largest value at  $s=0$ , even though  $n_{\text{on}} = 10$  and  $m_{\text{off}}=0$

BUT when mass distribution ignored (i.e. just counting experiment), signal+bgd is favoured over just bgd

# WHY?

Background given greater freedom with large number  $K$  of nuisance parameters

Compare:

Neyman and Scott, "Consistent estimates based on partially consistent observations", *Econometrica* 16: 1-32 (1948)

Data =  $n$  pairs  $X_{1i} = G(\mu_i, \sigma^2)$   
 $X_{2i} = G(\mu_i, \sigma^2)$

Param of interest =  $\sigma^2$

Nuisance params =  $\mu_i$ . Number increases with  $n$

Profile  $\mathcal{L}$  estimate of  $\sigma^2$  are biased  $E = \sigma^2/2$   
and inconsistent (bias does not tend to 0 as  $n \rightarrow \infty$ )

## MORAL: Beware!

# WHY LIMITS?

Michelson-Morley experiment → death of aether

HEP experiments: If UL on expected rate for new particle  $<$  expected, exclude particle

CERN CLW (Jan 2000)

FNAL CLW (March 2000)

Heinrich, PHYSTAT-LHC, “Review of Banff Challenge”

# Methods (no systematics)

Bayes (needs priors e.g. const,  $1/\mu$ ,  $1/\sqrt{\mu}$ ,  $\mu$ , .....

Frequentist (needs ordering rule,  
possible empty intervals, F-C)

Likelihood (DON'T integrate your L)

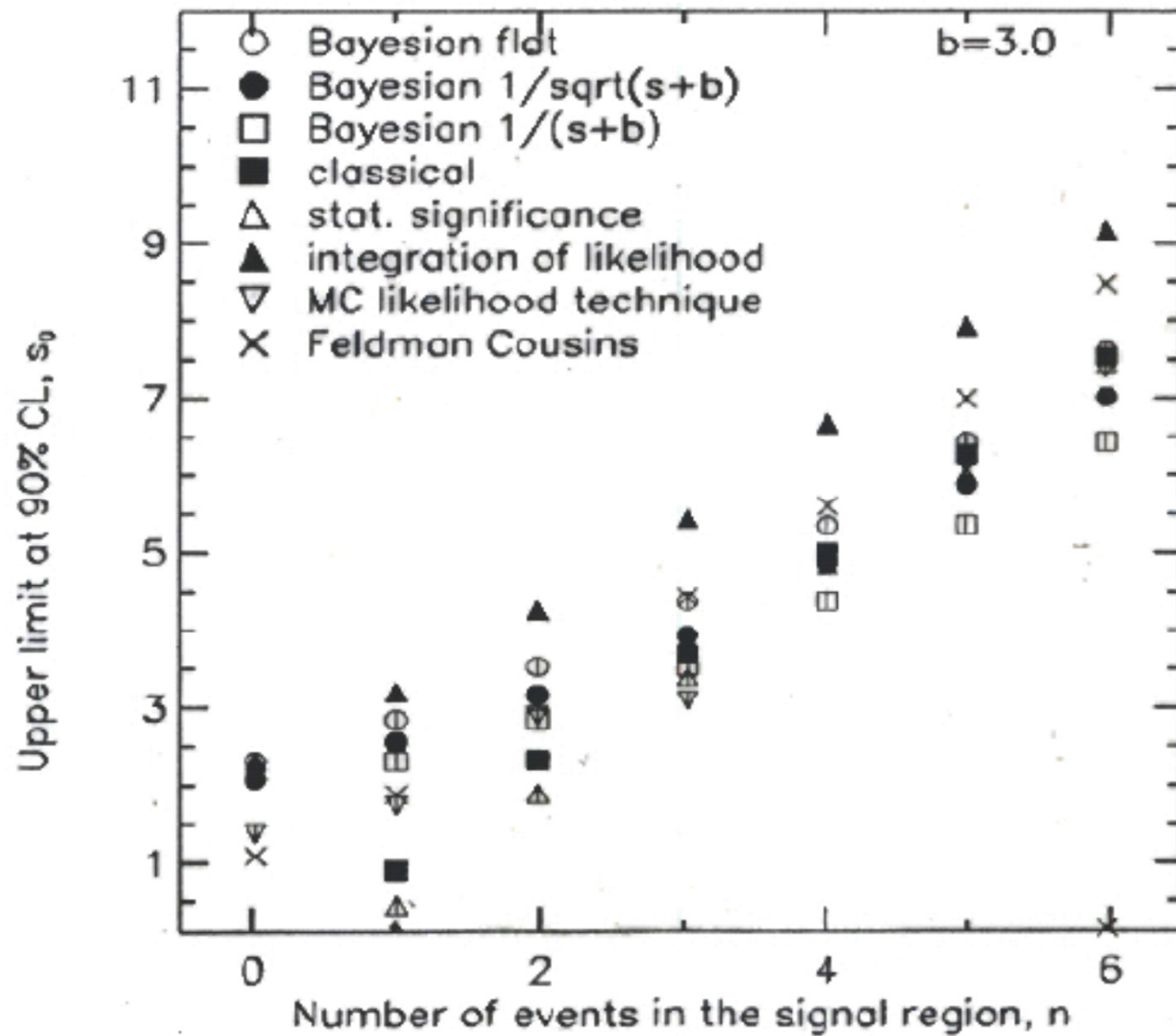
$$\chi^2(\sigma^2 = \mu)$$

$$\chi^2(\sigma^2 = n)$$

Recommendation 7 from CERN CLW: “Show your L”

- 1) Not always practical
- 2) Not sufficient for frequentist methods

Ilya Narsky, FNAL CLW 2000





# DESIRABLE PROPERTIES

- Coverage
- Interval length
- Behaviour when  $n < b$
- Limit increases as  $\sigma_b$  increases
- Unified with discovery and interval estimation

## 90% Classical interval for Gaussian

$$\sigma = 1 \quad \mu \geq 0 \quad \text{e.g. } m^2(v_e)$$

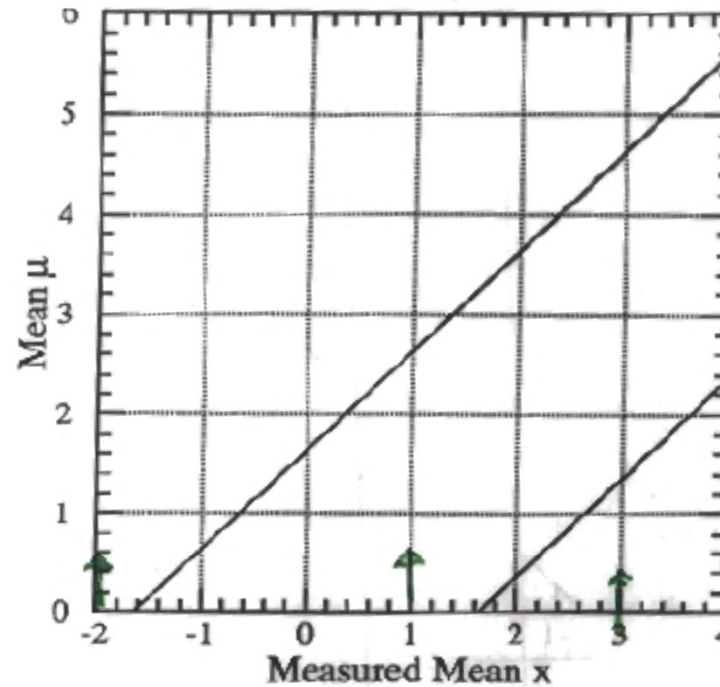


FIG. 3. Standard confidence belt for 90% C.L. central confidence intervals for the mean of a Gaussian, in units of the rms deviation.

- $X_{\text{obs}} = 3$  Two-sided range
- $X_{\text{obs}} = 1$  Upper limit
- $X_{\text{obs}} = -2$  No region for  $\mu$

# FELDMAN - COUSINS

Wants to avoid empty classical intervals →

Uses “ $\mathcal{L}$ -ratio ordering principle” to resolve ambiguity about “which 90% region?”

[Neyman + Pearson say  $\mathcal{L}$ -ratio is best for hypothesis testing]

Unified → No ‘Flip-Flop’ problem

# Classical (Neyman) Confidence Intervals

Uses only  $P(\text{data}|\text{theory})$

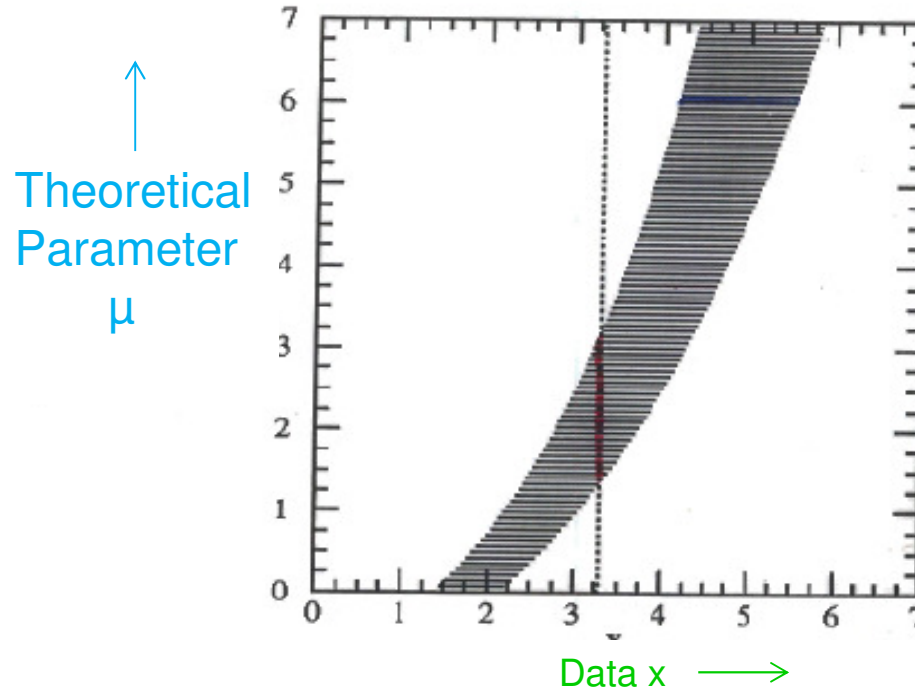


FIG. 1. A generic confidence belt construction and its use. For each value of  $\mu$ , one draws a horizontal acceptance interval  $[x_1, x_2]$  such that  $P(x \in [x_1, x_2] | \mu) = \alpha$ . Upon performing an experiment to measure  $x$  and obtaining the value  $x_0$ , one draws the dashed vertical line through  $x_0$ . The confidence interval  $[\mu_1, \mu_2]$  is the union of all values of  $\mu$  for which the corresponding acceptance interval is intercepted by the vertical line.

Example:

Param = Temp at centre of Sun

Data = Est. flux of solar neutrinos

$$\text{Prob}(\mu_l < \mu < \mu_u) = \alpha$$

$$\mu \geq 0$$

No prior for  $\mu$

# Classical (Neyman) Confidence Intervals

Uses only  $P(\text{data}|\text{theory})$

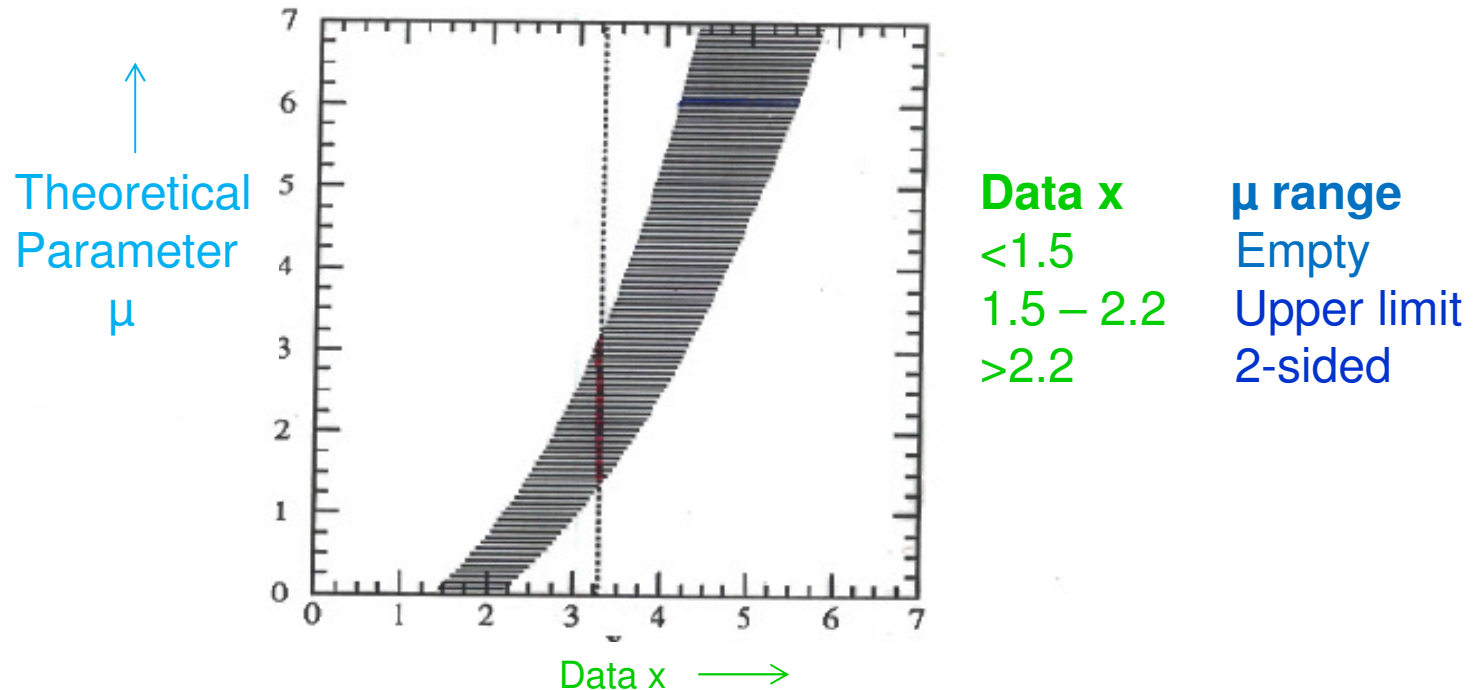


FIG. 1. A generic confidence belt construction and its use. For each value of  $\mu$ , one draws a horizontal acceptance interval  $[x_1, x_2]$  such that  $P(x \in [x_1, x_2] | \mu) = \alpha$ . Upon performing an experiment to measure  $x$  and obtaining the value  $x_0$ , one draws the dashed vertical line through  $x_0$ . The confidence interval  $[\mu_1, \mu_2]$  is the union of all values of  $\mu$  for which the corresponding acceptance interval is intercepted by the vertical line.

Example:

Param = Temp at centre of Sun

Data = est. flux of solar neutrinos

$$\mu \geq 0$$

No prior for  $\mu$

## Feldman-Cousins 90% conf intervals

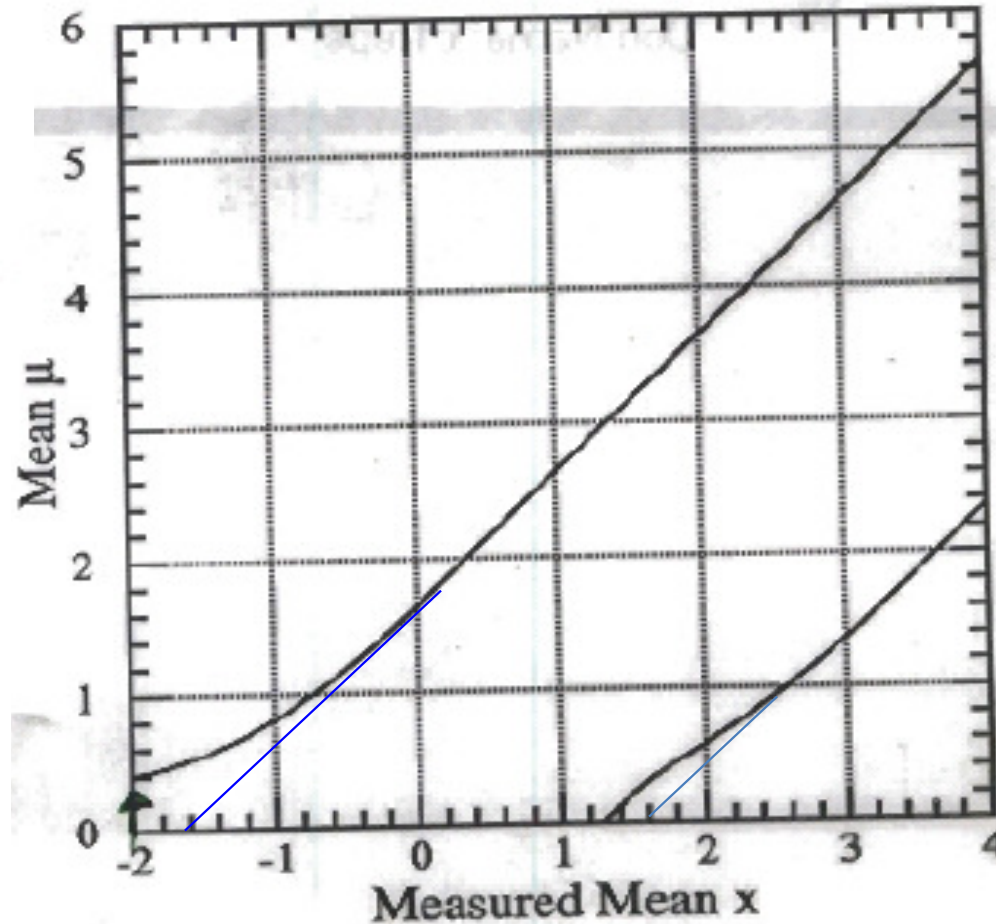


FIG. 10. Plot of our 90% confidence intervals for mean of a Gaussian, constrained to be non-negative, described in the text.

$X_{\text{obs}} = -2$  now gives upper limit

————— central confidence region

# Features of Feldman-Cousins

**Reduces/Eliminates empty intervals**

**Unifies 1-sided and 2-sided intervals**

**Eliminates flip-flop**

**No arbitrariness of intervals**

**'Readily' extends to several dimensions (Other ordering rules have trouble)**

**Less overcoverage than 'no more than 5%' at each end**

Neyman construction is CPU intensive, esp in several dimensions

Problem dealing with systematics consistently with main analysis

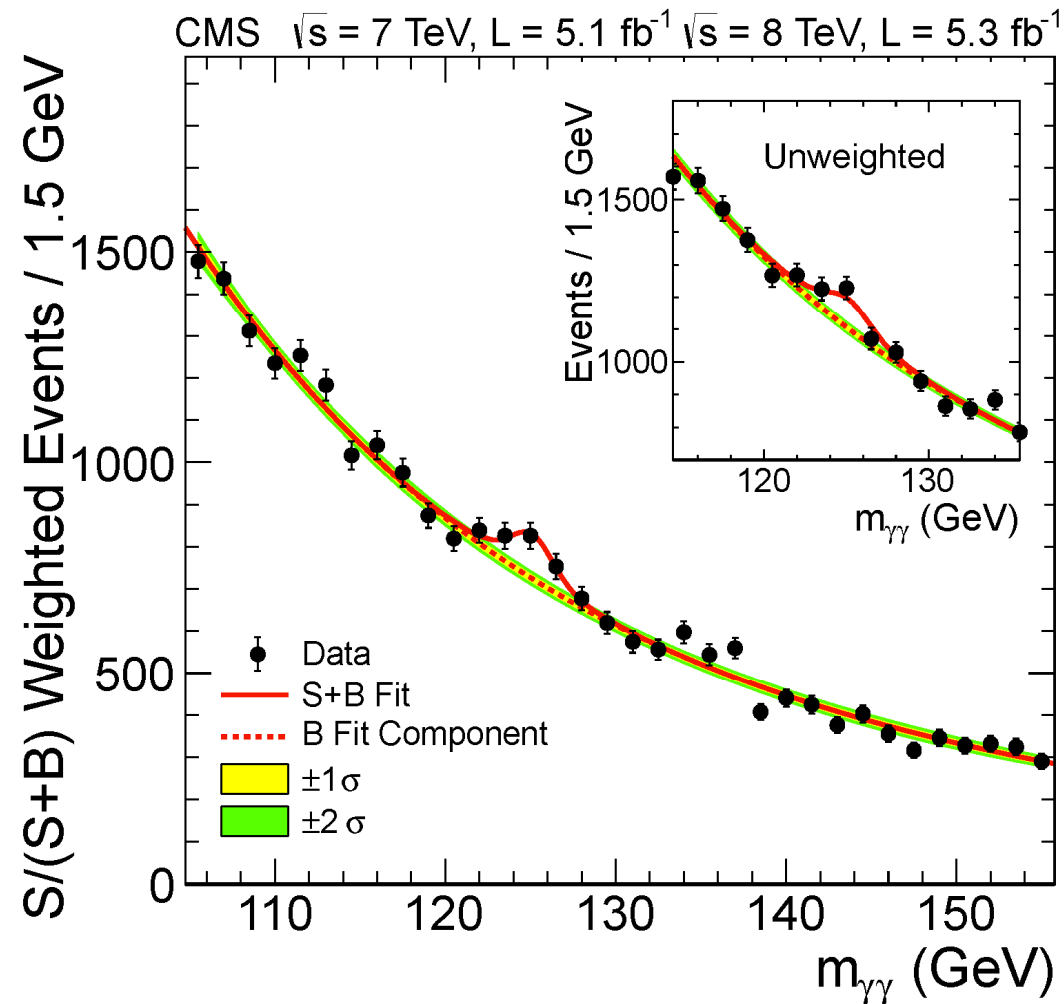
Minor pathologies: Discontinuous intervals

Behaviour wrt background

Tight limits when  $n_{\text{obs}}$  less than  $b$

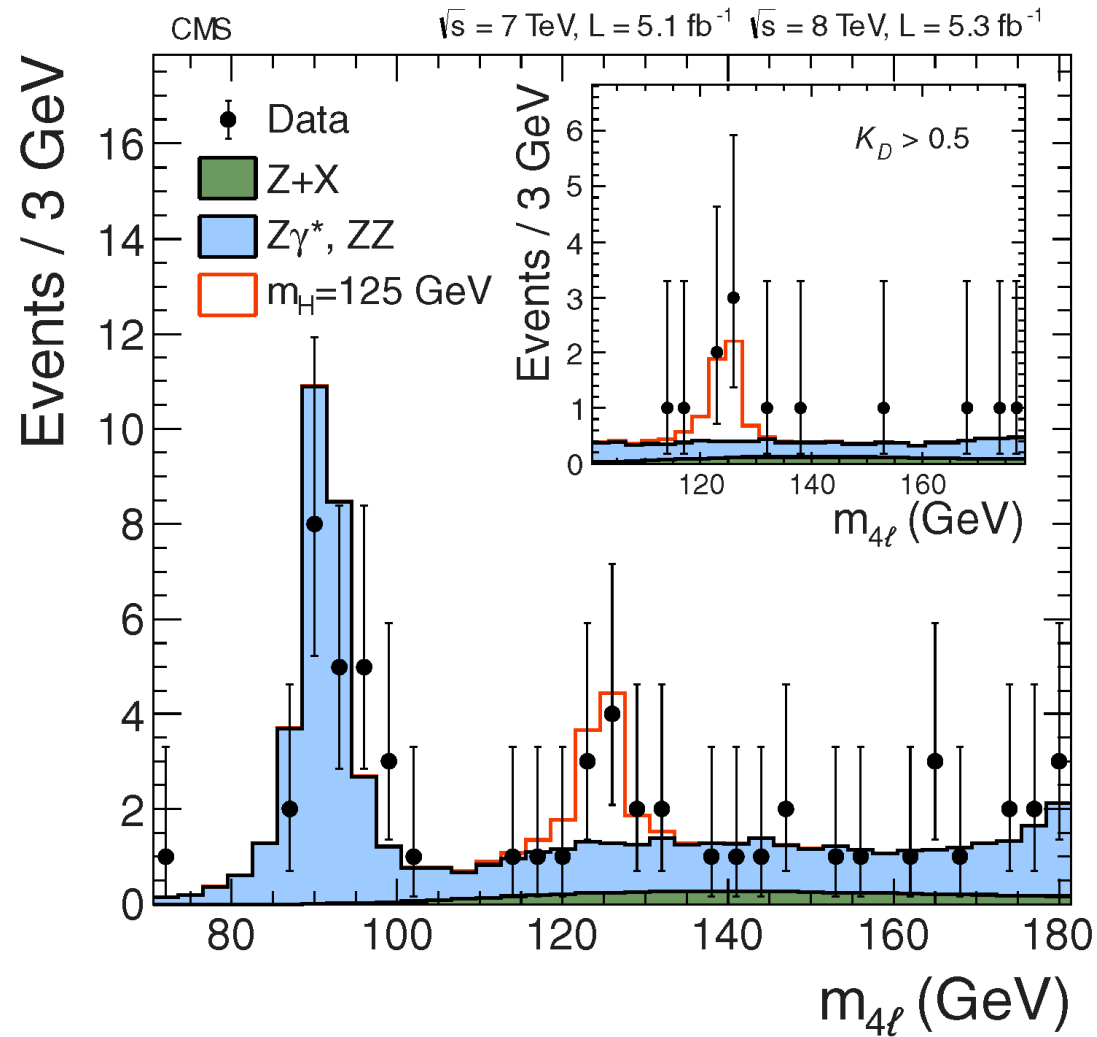
Quicker exclusion of  $s=0$  wrt standard frequentist

# Search for Higgs: $H \rightarrow \gamma\gamma$ : low S/B, high statistics

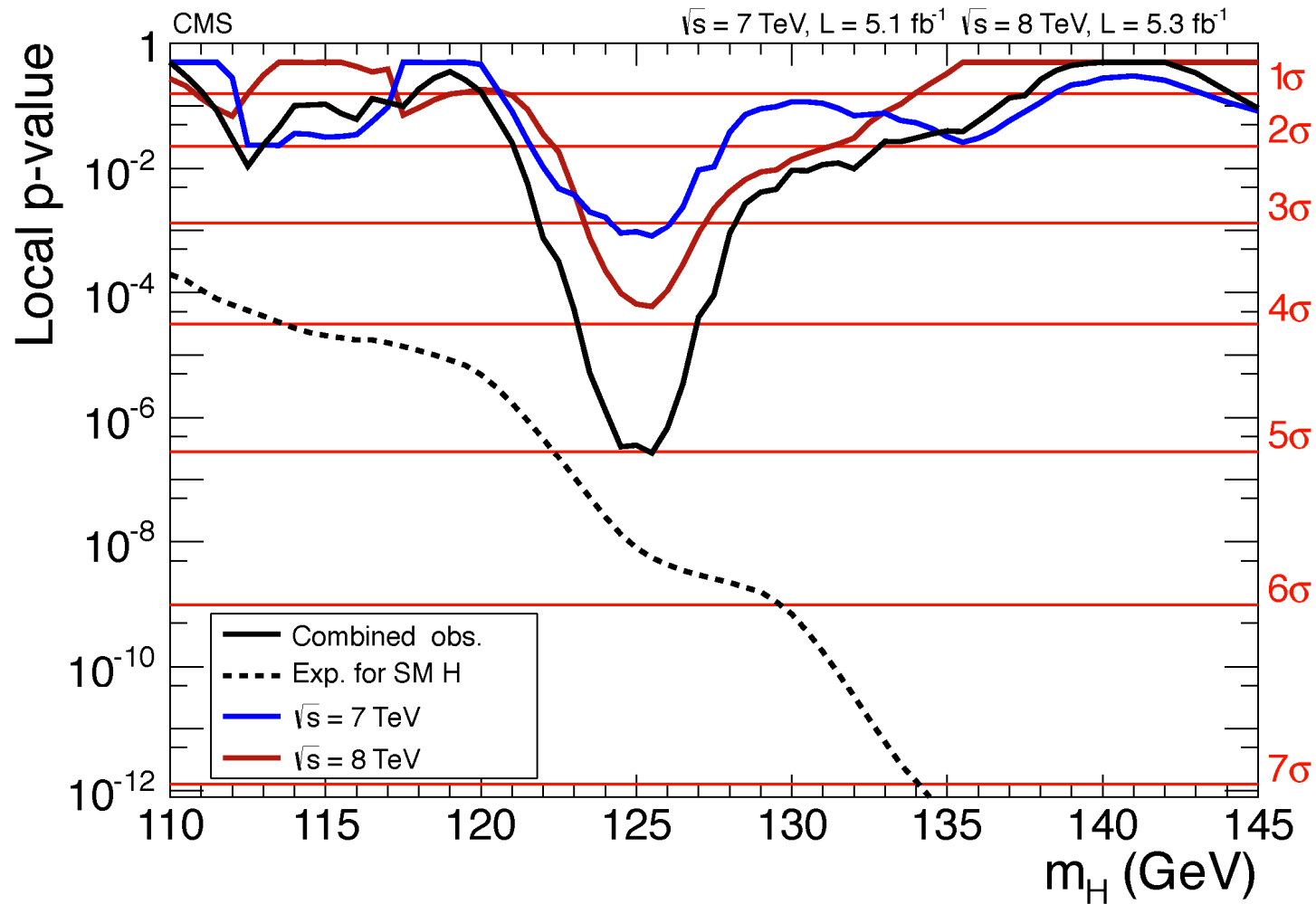




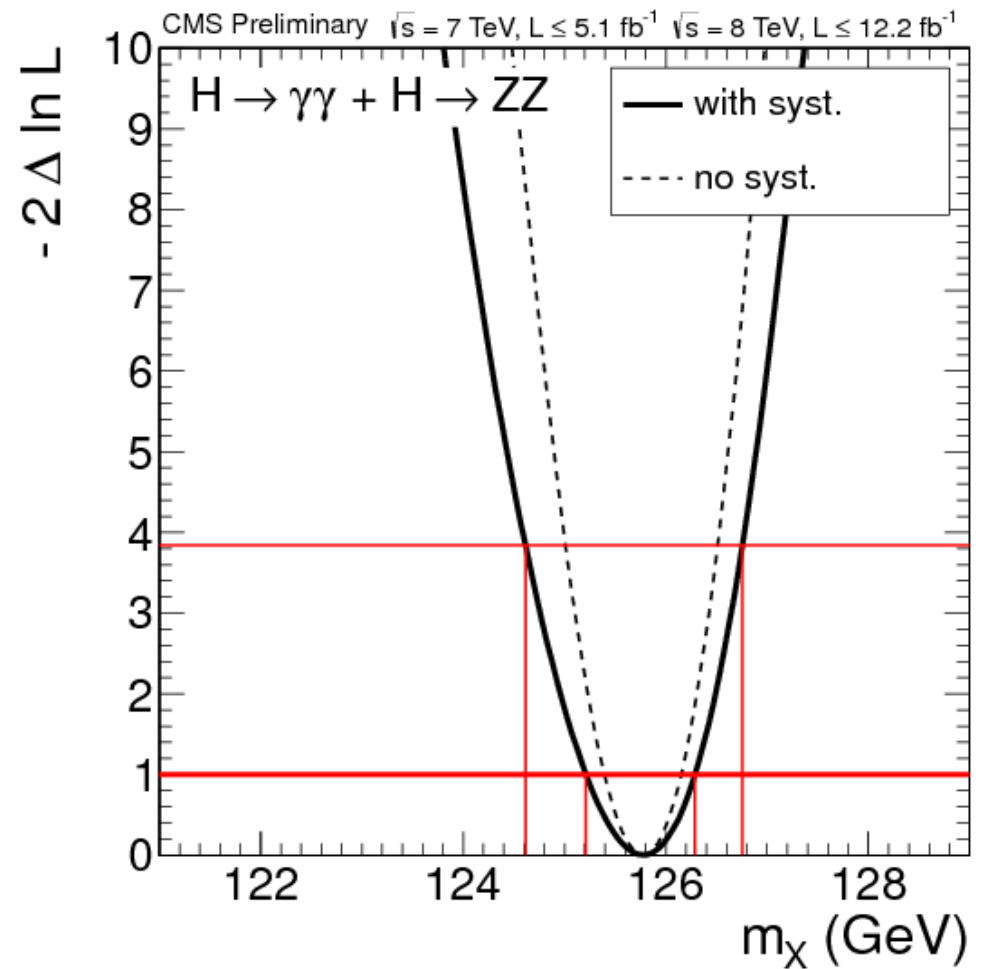
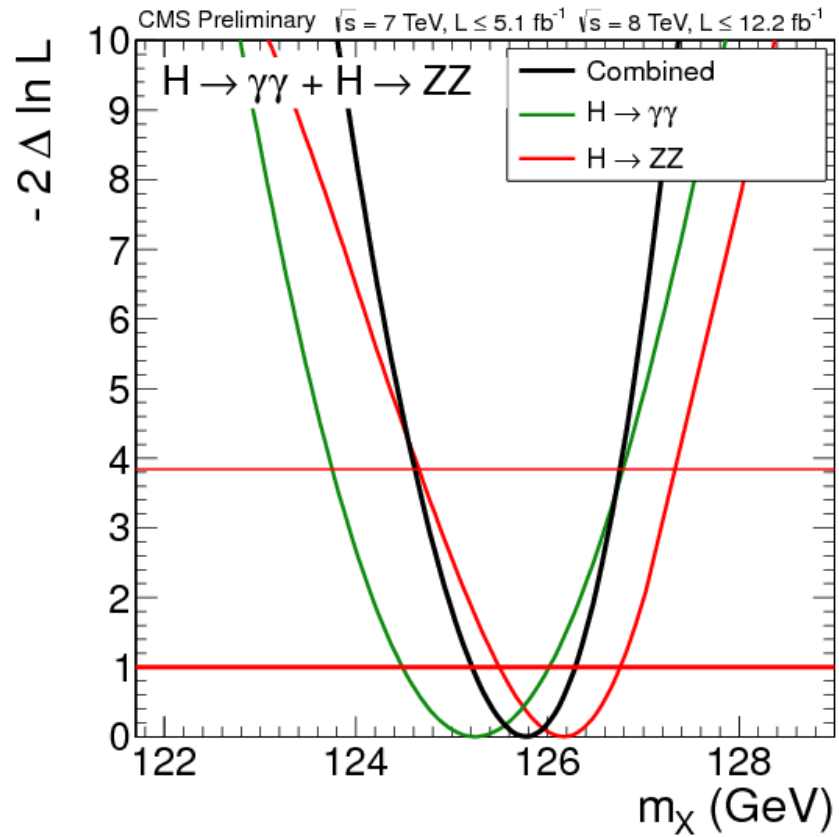
# $H \rightarrow Z Z \rightarrow 4 \ell$ : high S/B, low statistics



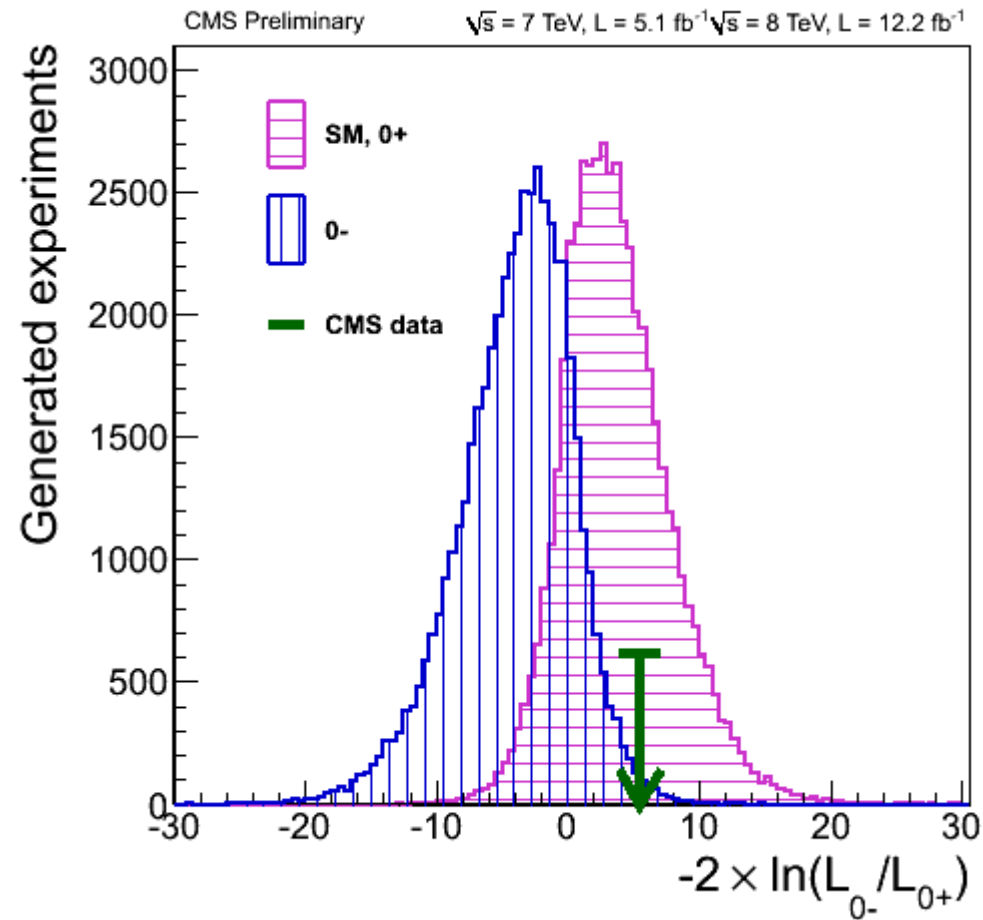
# p-value for 'No Higgs' versus $m_H$



# Mass of Higgs: Likelihood versus mass



# Comparing $0^+$ versus $0^-$ for Higgs (like Neutrino Mass Hierarchy)



<http://cms.web.cern.ch/news/highlights-cms-results-presented-hcp>

# Conclusions

## Resources:

Software exists: e.g. RooStats

Books exist: Barlow, Cowan, James, Lista, Lyons, Roe,.....

New: `Data Analysis in HEP: A Practical Guide to Statistical Methods', Behnke et al.

PDG sections on Prob, Statistics, Monte Carlo

CMS and ATLAS have Statistics Committees (and BaBar and CDF earlier) – see their websites

Before re-inventing the wheel, try to see if Statisticians have already found a solution to your statistics analysis problem.

Don't use a square wheel if a circular one already exists.

**“Good luck”**