

# Introduction to Statistical Issues for Phystat-DM

Louis Lyons

Imperial College, London  
and  
Oxford

PHYSTAT-DM@ Stockholm  
July 2019

## Why bother ?

Experiments are expensive and time-consuming, so:

Worth investing effort in statistical analysis

→ better information from data

# Statistical Procedures

Parameter Determination

**Central value and range (or upper limit)**

**e.g. flux of WIMPs**

Comparing data with Hypotheses → Discoveries, Upper Limits,...

**Just one Hypothesis**

**Goodness of Fit**

**e.g. Just known particles**

**Comparing 2 Hypotheses**

**Hypothesis Testing**

**e.g. Just known particles or Also WIMPs**

# TOPICS

## Introduction

Some issues related to Discovery claims

Choosing between 2 hypotheses

p-values (including CL<sub>s</sub>)

Blind analyses

Look Elsewhere Effect

Why 5-sigma for discovery

Background systematics

## Upper Limits

## Summary

**Other important topics not included:**

Combining results

Likelihoods (including Coverage)

Bayes and Frequentism

MVA: How Neural Networks work

Wilks Theorem

Discovery of Higgs

# 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>

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

## Topics:

Example of Hypotheses

H0 or H0 v H1?

Blind Analysis

Why  $5\sigma$  for discovery?

Significance

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

Meaning of p-values

Wilks' Theorem

LEE = Look Elsewhere Effect

Background Systematics

**Upper Limits**

Higgs search: Discovery and spin

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

# Examples of types of Hypotheses

## 1) Event selector

Selection of event sample based on required features

e.g. H0: Cerenkov ring produced by electron    H1: Produced by other particle

Possible outcomes:    Events assigned as H0 or H1

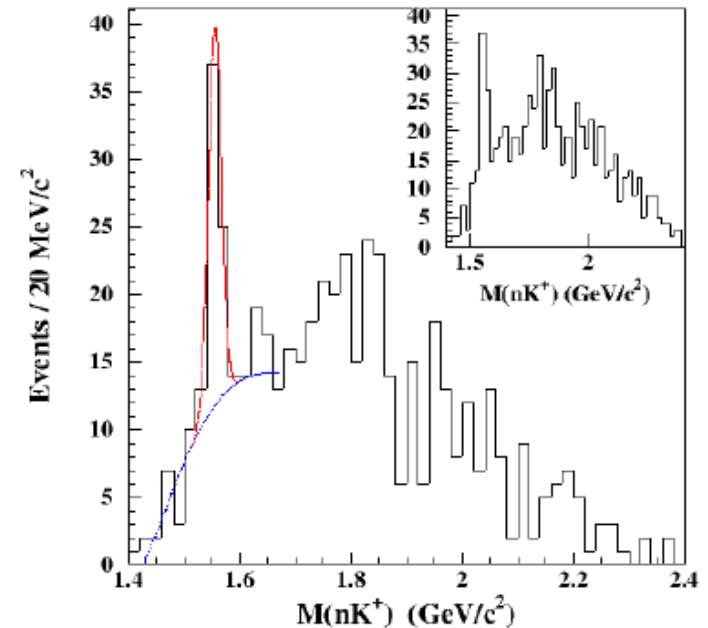
## 2) Result of experiment

e.g. H0 = nothing new

H1 = new particle produced as well

(Sterile neutrino,.....)

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



# Errors of 1<sup>st</sup> and 2<sup>nd</sup> Kind

- 1<sup>st</sup> Kind: Reject H0 when H0 true  
Should happen at rate  $\alpha$
- 2<sup>nd</sup> Kind: Fail to reject H0 when H0 is false  
Rate depends on:
  - How similar H0 and H1 are
  - Relative rates of H0 and H1 (for event selector)

For event selector: E1<sup>st</sup> = Loss of efficiency

E2<sup>nd</sup> = Contamination

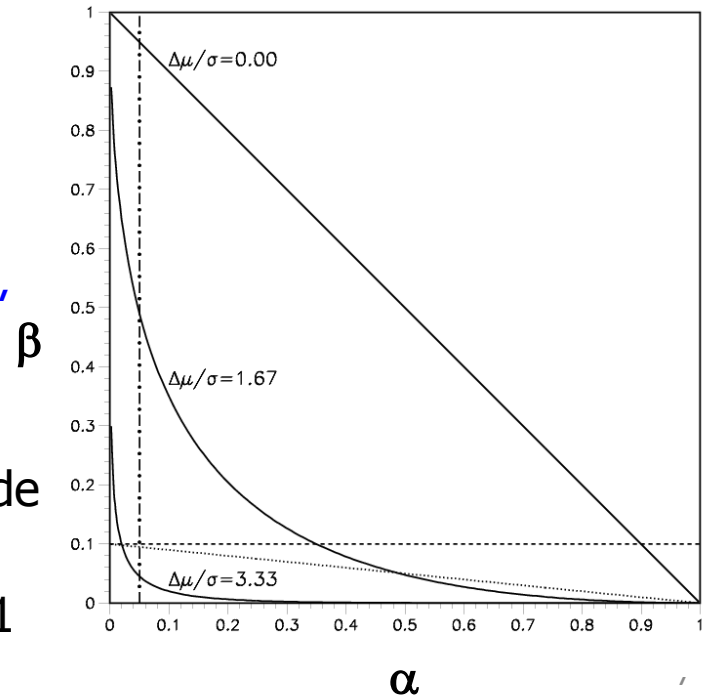
As  $\alpha \downarrow$ , efficiency  $\uparrow$  and contamination  $\downarrow$

For result of expt , E1<sup>st</sup> gives incorrect result  
E2<sup>nd</sup> fails to make discovery

$$\alpha = E1^{st}$$

$\beta$  = Prob of failing to exclude  
H0, if H1 = true

$$1 - \beta = \text{power of test for H1}$$



# H0 or H0 versus H1 ?

H0 = null hypothesis

e.g. known backgrounds, with nothing new

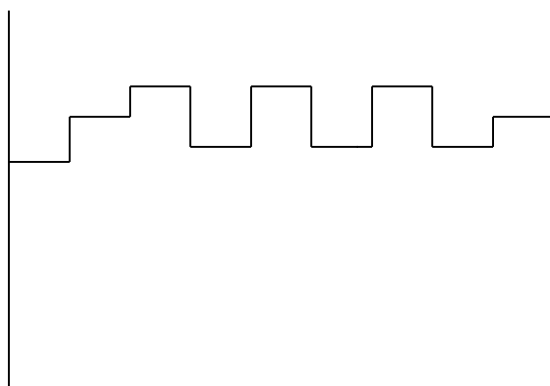
H1 = specific New Physics e.g. WIMP with  $M_W = 70$  GeV

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

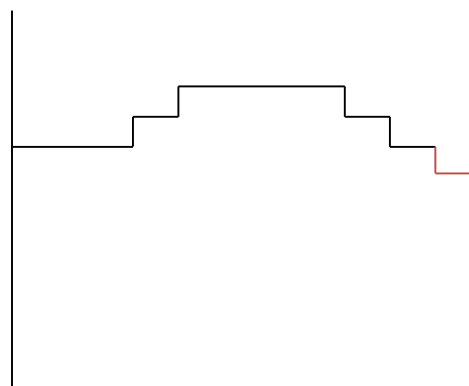
H0 v H1: “Hypothesis Testing” e.g.  $\mathcal{L}$ -ratio

Measures how much data favours one hypothesis wrt other

H0 v H1 likely to be more sensitive for H1



or



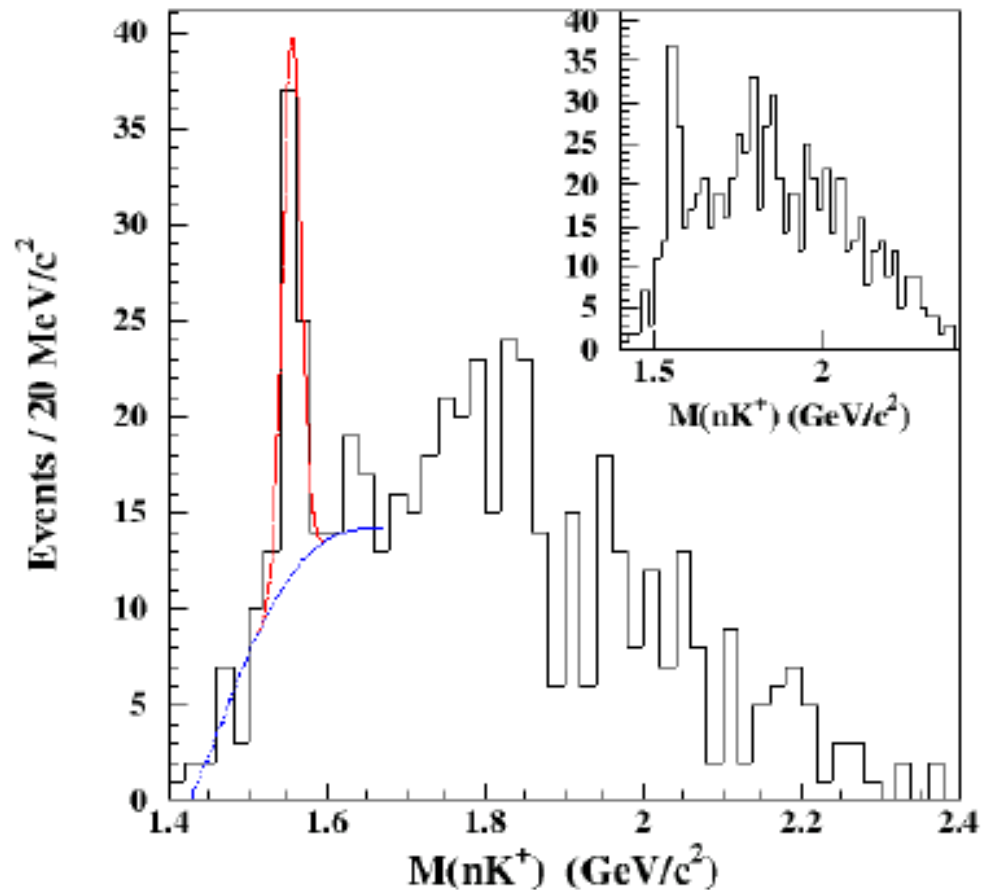


# Choosing between 2 hypotheses

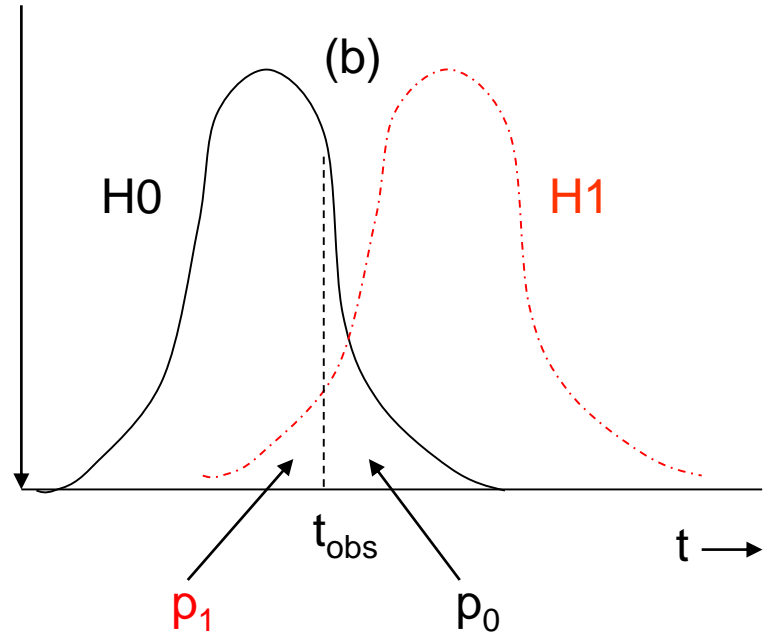
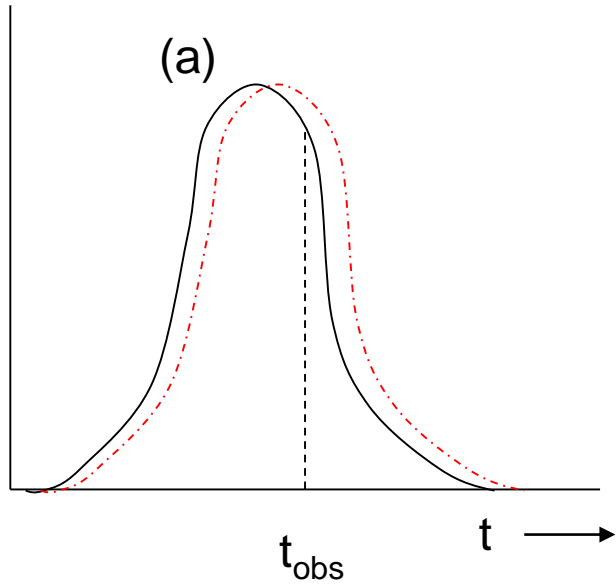
Hypothesis testing: New particle or statistical fluctuation?

$H_0 = b$

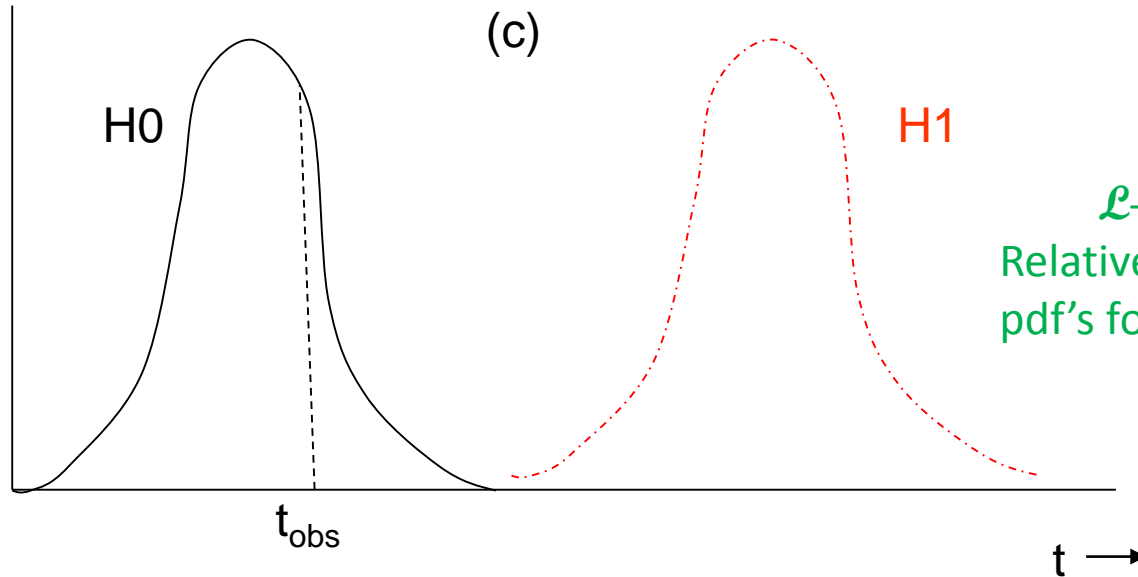
$H_1 = b + s$



First define your data statistic  $t$  ( $n$ ,  $\mathcal{L}$ -ratio, etc.)

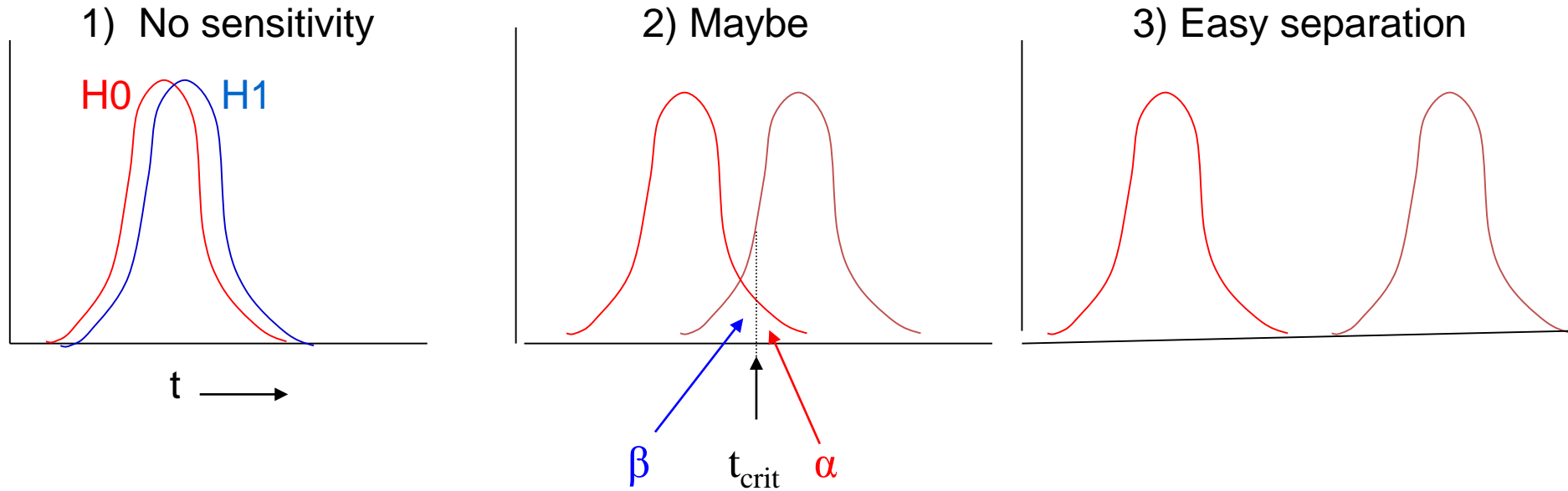


With 2 hypotheses, each with own pdf, p-values are defined as tail areas, pointing in towards each other



$\mathcal{L}$ -ratio:  
Relative heights of pdf's for  $H_0$  and  $H_1$

# 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  $\rightarrow$  discovery or excln

$1-\beta = \text{Power of test}$

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**

# p-values and z-score (number of sigma)

Conventional to convert p-values to number of sigma for one-sided tail of Gaussian

e.g.  $16\% = 1\sigma$

$3 \times 10^{-7} = 5\sigma$

Statisticians call this 'z-score'

Does NOT imply that actual pdf is Gaussian

Just convention

Simply easier to remember than corresponding p-value

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

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

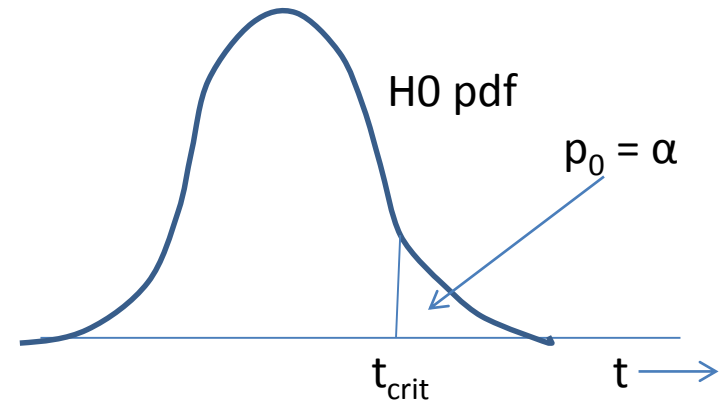
or Prob{speed of  $v > c$  | data}  $\sim 1$

Everyday situation:

$p(\text{eat bread} | \text{murderer}) \sim 99\%$

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

# 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’]

# $p_0$ 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

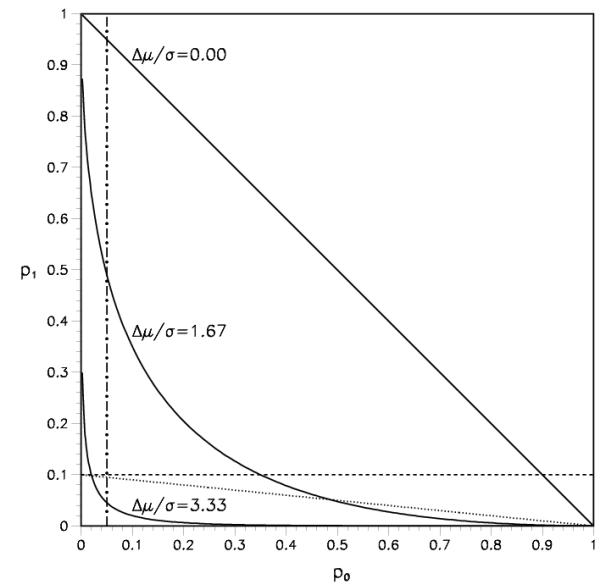
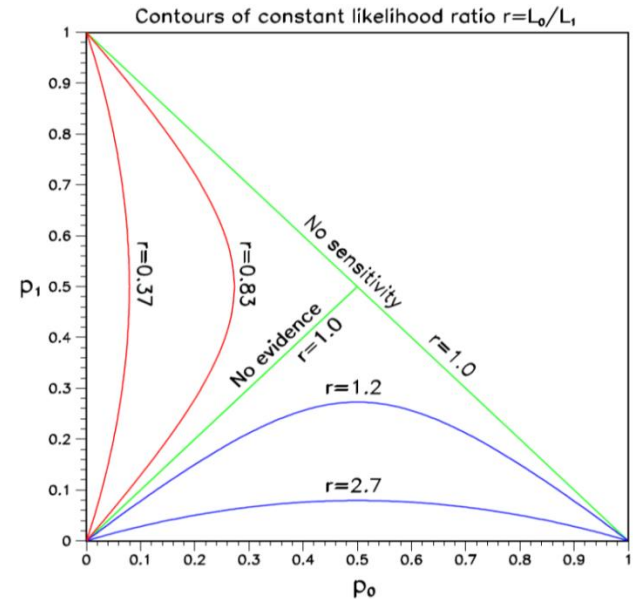
Punzi definition of sensitivity

Relation of p-values and Likelihoods

Probability of misleading evidence

Sampling to foregone conclusion

Jeffreys-Lindley paradox

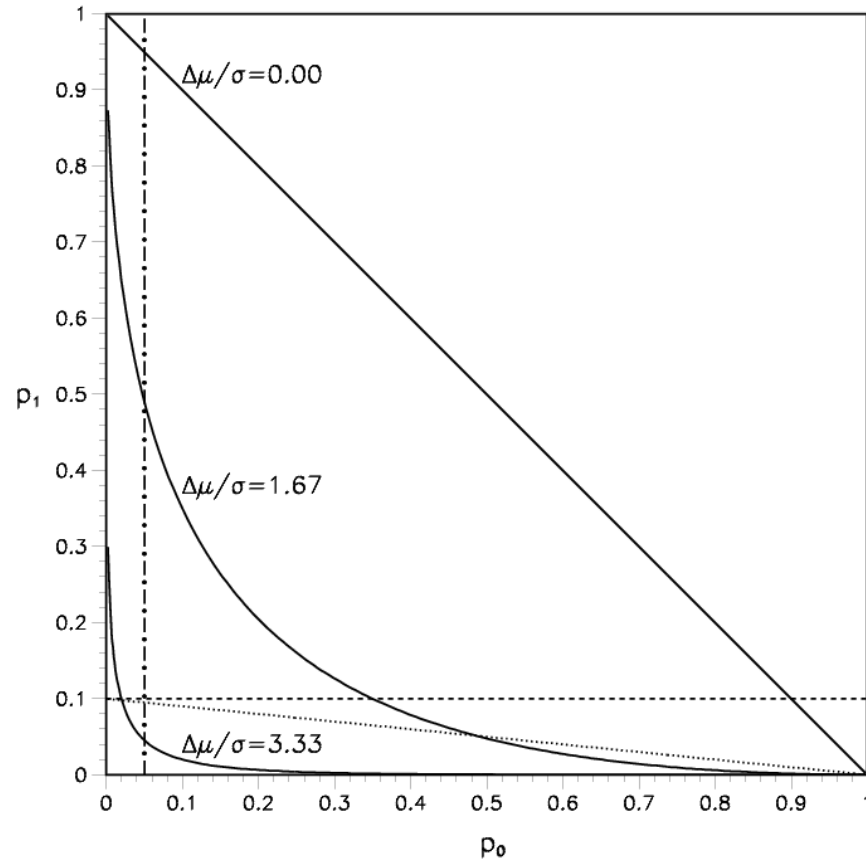
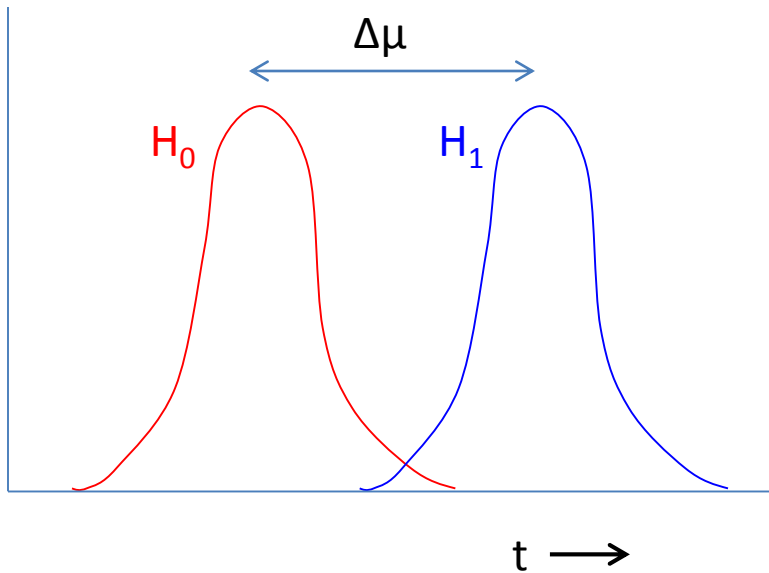


CLs =  $p_1/(1-p_0)$   $\rightarrow$  diagonal line

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:

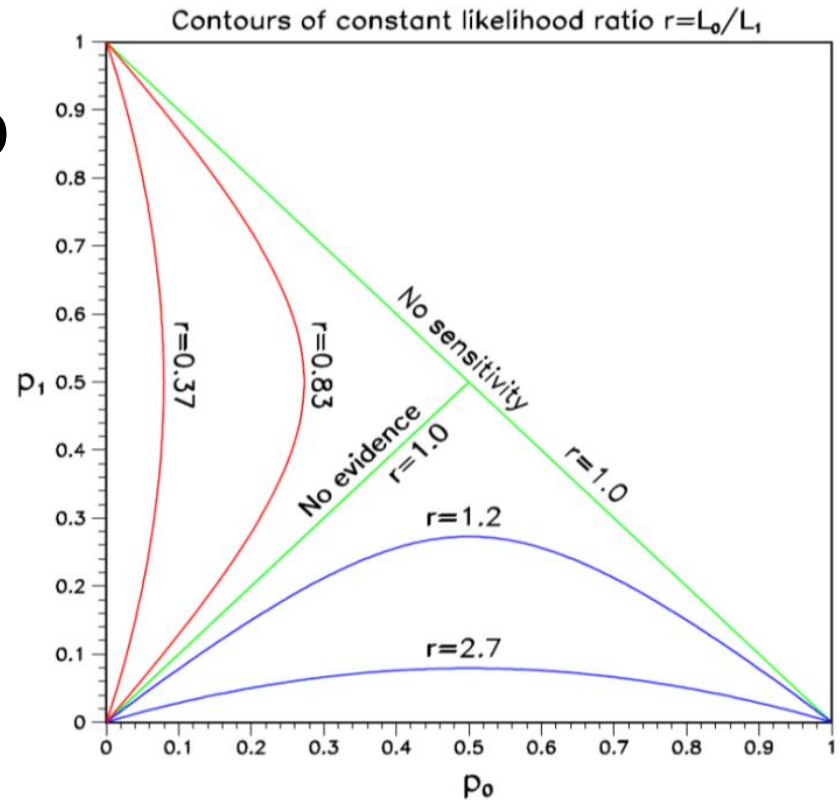
e.g. Poisson counting expt little data:

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  $H_1$ ,  $\mu_1 = 10.0$

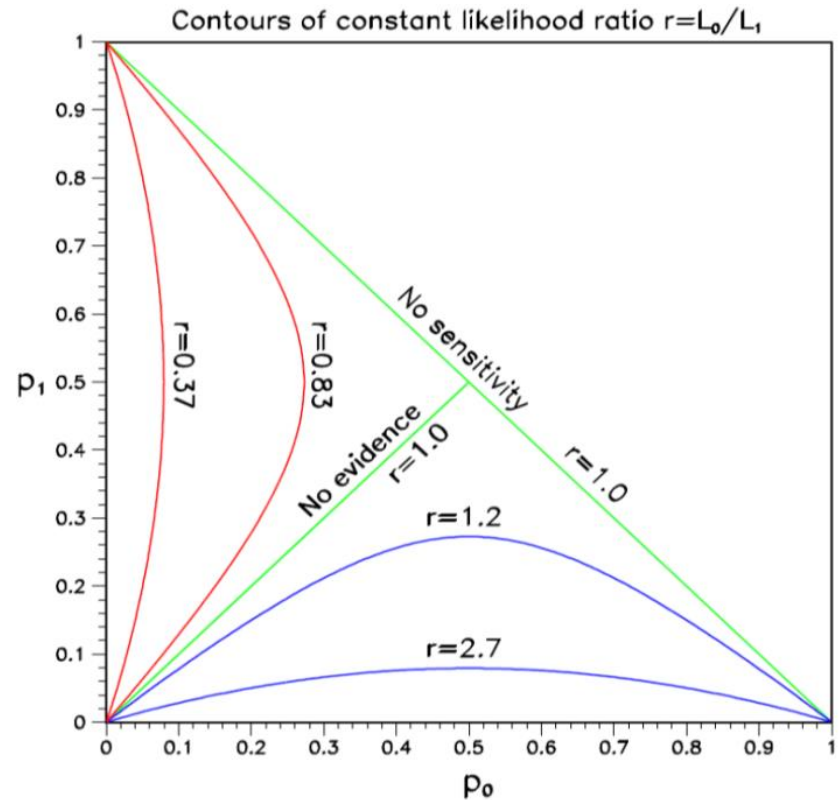
Observe  $n = 10$   $p_0 \sim 10^{-7}$   $L_{01} \sim 10^{-5}$  and favours  $H_1$

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 \*\*\*\*\*

# Significance

(Number of  $\sigma$  = p-value converted to Gaussian one-sided tail)

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

## 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  $S = 0.1$ ,  $B = 10^{-4}$ ,  $S/\sqrt{B} = 10$

# BLIND ANALYSES

Why blind analysis? Data statistic, selections, corrections, method

## 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

# 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

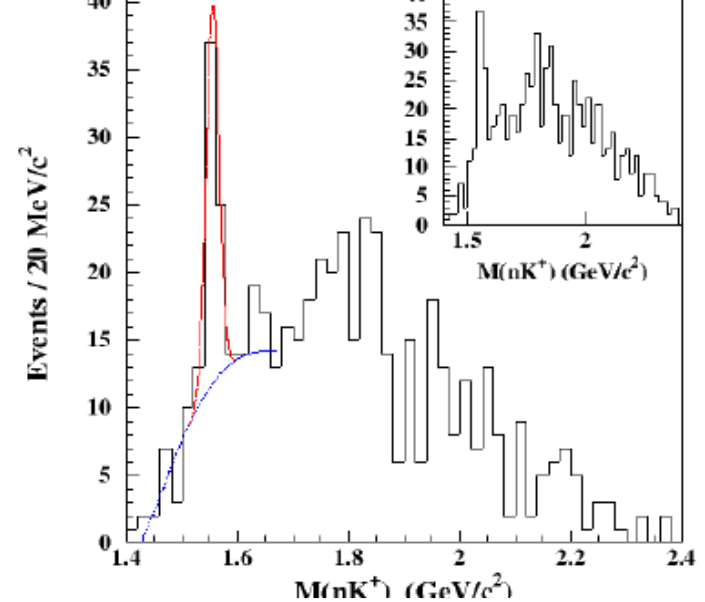
f) relevant for experiment's Spokesperson

**INFORMAL CONSENSUS:**

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

Explain which global p

For DM 'counting' experiments, there is almost no LEE.



# 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  
 $8640$  divided by  $2160 = 4$ th time.

30 Stones in Outer ring =  
 $360$  divided by  $30 = 12$

60 Stones in Second ring =  
 $360$  divided by  $60 = 6$

20 Stones in Center ring =  
 $360$  divided by  $20 = 18$

IF THIS WAS EAST

WINTER SOLSTICE

SUMMER SOLSTICE

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*

WINTER SOLSTICE

SUMMER SOLSTICE

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

$25,920$  divided by  $12 = 2160$

$25,920$  divided by  $6 = 4320$

$25,920$  divided by  $18 = 1440$

WEST  
BALANCED LOCATION IN SPACE

STONEHENGE

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

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 prob      Likelihood ratio      Priors

“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 'delivered on Mt. Sinai'

**How would you rate 'Dark Matter'?**

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

# SYSTEMATICS

- Harder than statistical uncertainties
- Requires much more thought and effort

Different types:

A) On measured quantities to extract answer

B) On implicit assumptions

e.g Simple pendulum expt  $\tau = 2\pi \sqrt{L/g}$

A)  $\tau$  and  $L$

B) Point mass; massless string; small amplitude; no damping

Systematics can be:

- i) Measured in subsidiary (or main) analysis
- ii) Exptl effects not directly measured; or inconsistent results
- iii) Different theories

Just one example here: **BACKGROUND SYSTEMATICS**

# SYSTEMATICS

- Harder than statistical uncertainties
- Requires much more thought and effort

Different types:

A) On measured quantities to extract answer

B) On implicit assumptions

e.g Simple pendulum expt  $\tau = 2\pi \sqrt{L/g}$

A)  $\tau$  and  $L$

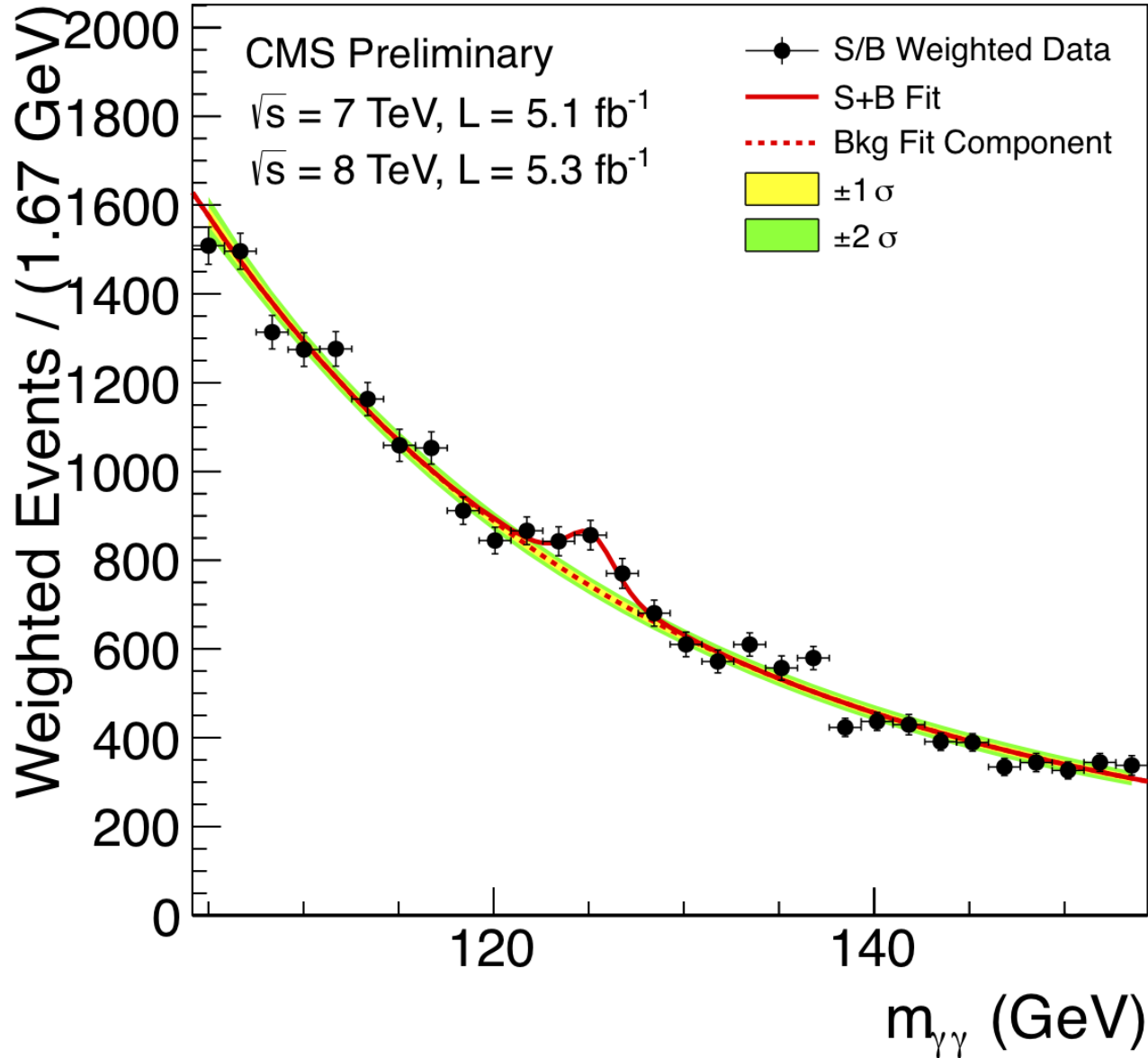
B) Point mass; massless string; small amplitude; no damping

Systematics can be:

- |  |      |
|--|------|
| i) Measured in subsidiary (or main) analysis                     | GOOD |
| ii) Exptl effects not directly measured; or inconsistent results | BAD  |
| iii) Different theories  | UGLY |

Just one example here: **BACKGROUND SYSTEMATICS**

# 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  $\sim$ comparable to  $f_a$  fit, include contribution to systematics

But what is ' $\sim$ comparable'?

Other approaches:

Profile likelihood over different bgd parametric forms

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

Background subtraction

sPlots

Non-parametric background

Bayes

Yellin's Optimal Interval

Cowan's 'Error on the error'

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  $H \rightarrow \gamma\gamma$

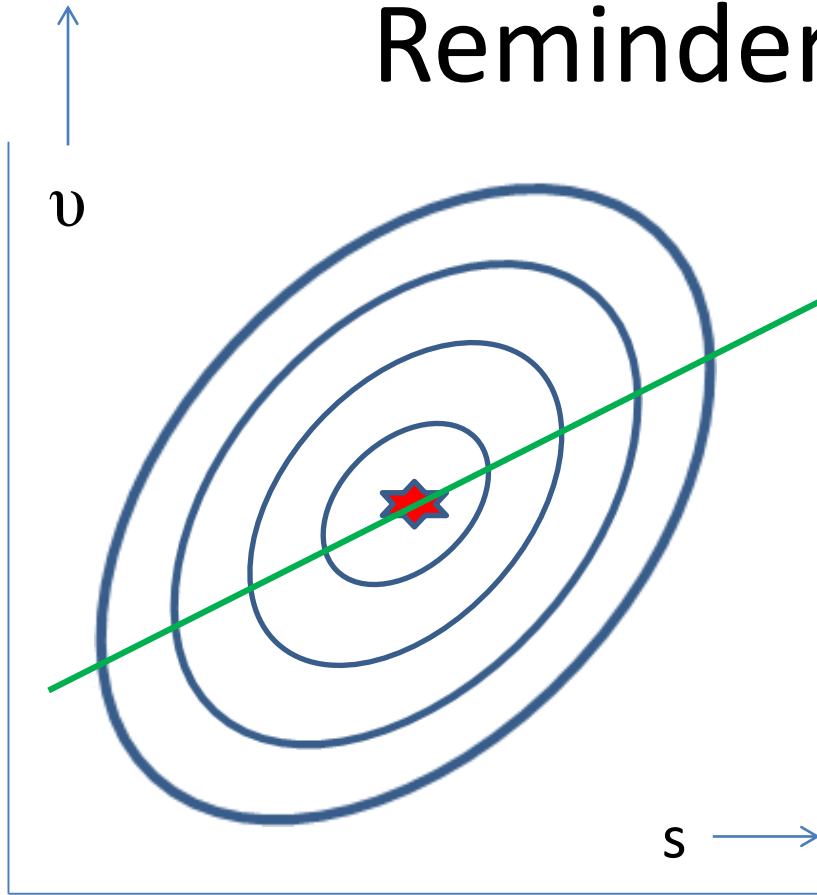
Problem with ‘Typical approach’: Alternative functional forms do or don’t contribute to systematics by hard cut, so systematics can change discontinuously wrt  $\Delta\chi^2$

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

Here ‘profile’ over discrete functional forms



# Reminder of Profile $\mathcal{L}$



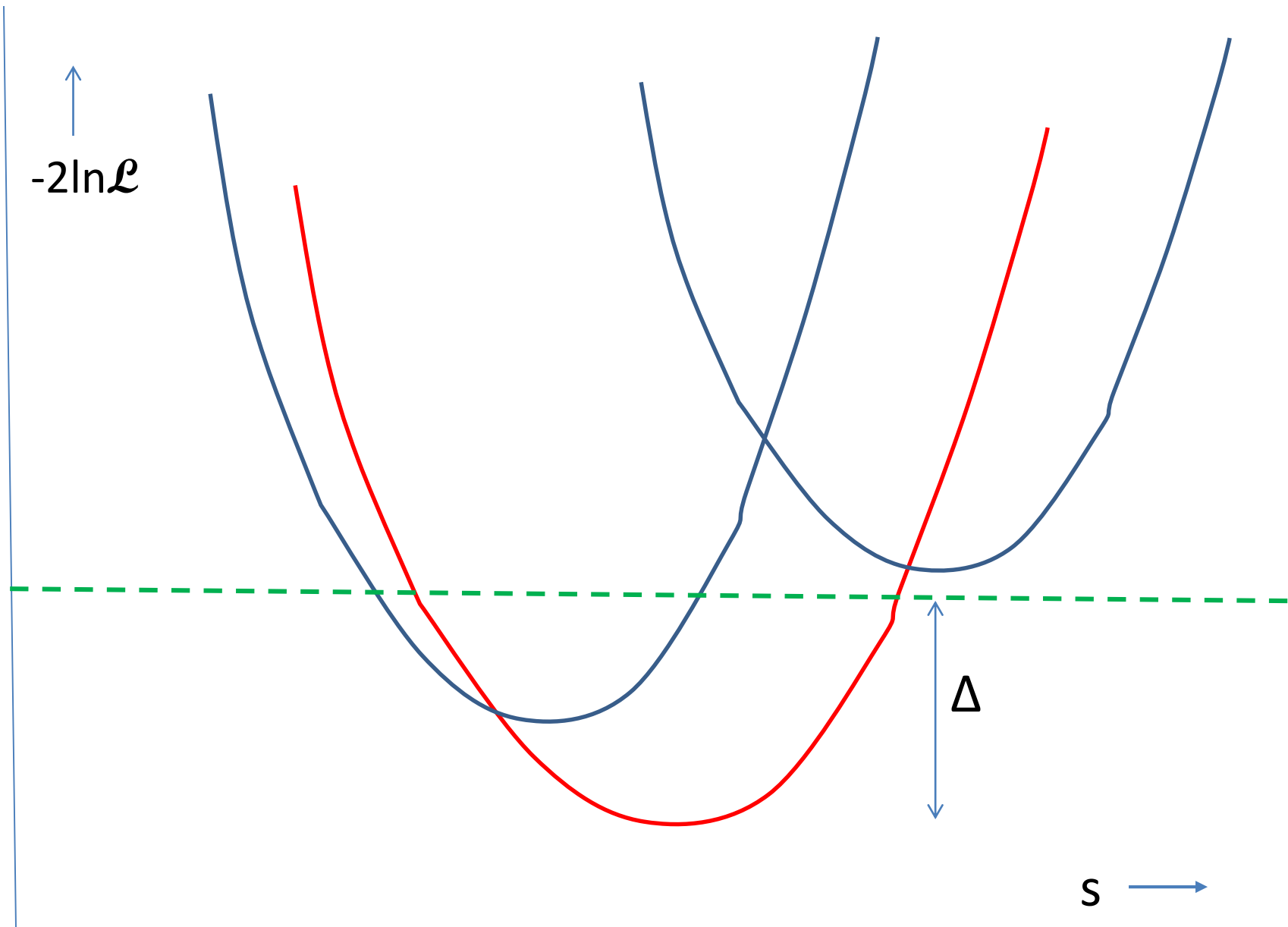
Contours of  $\ln \mathcal{L}(s, v)$   
 $s$  = physics param  
 $v$  = nuisance param

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

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

# WHY LIMITS?

Michelson-Morley experiment → death of aether

HEP experiments:

If UL on rate for new particle < expected,  
exclude particle

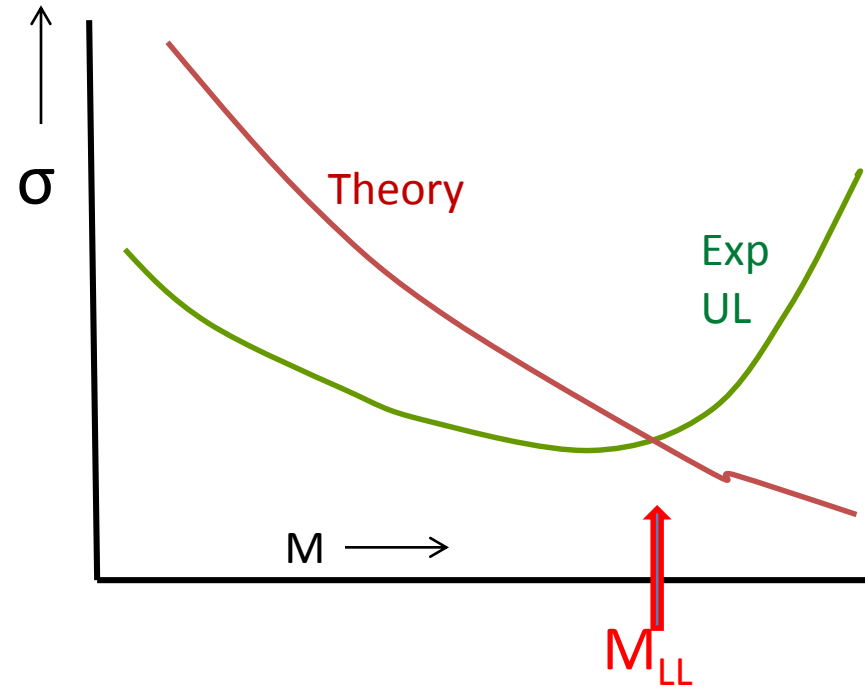
(Almost) all direct DM searches  
quote UL on DM flux, rather than  
claiming a discovery (i.e. flux  $\neq 0$ )

If theory curve below UL on  $\sigma$ ,  
expt not sensitive enough to exclude  
any mass.

CERN CLW (Jan 2000)

FNAL CLW (March 2000)

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



# Methods for ULs (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)

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

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

$CL_s$

Power Constrained Limits

Optimal Interval Method

Recommendation 7 from CERN CLW (2000): “Show your  $\mathcal{L}$ ”

1) Not always practical

2) Not sufficient for frequentist methods

# Power Constrained Limits

When  $n_{\text{obs}} < b$  (expected background), downward fluctuation in data  $\rightarrow$  Tighter than expected limits

Avoid most extreme cases by quoting expectation (or  $\text{exp} - k\sigma$ ) instead of actual limit.

Suggested by Cowan et al (ATLAS), but abandoned and not used.

**NOT RECOMMENDED**

# Optimal Interval

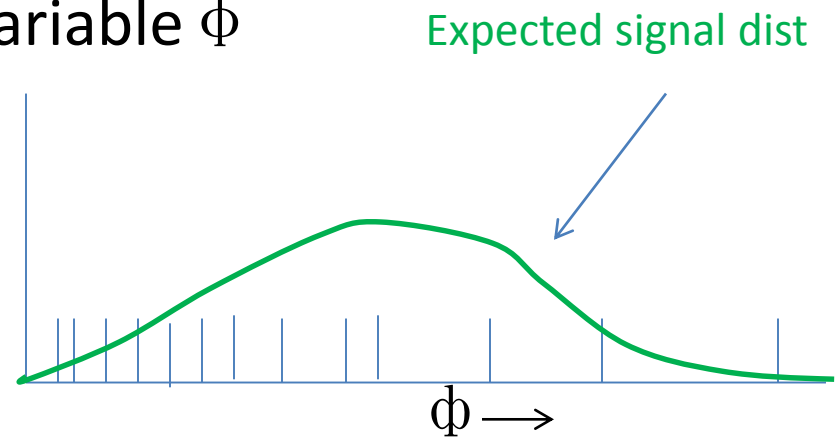
Steven Yellin, PHYSTAT2003

Good for case when shape of background is uncertain – see PRD 66 (2005) 032005

Use distrib of events in some variable  $\phi$

Choose gap with largest expected signal.

Then use intervals with  $n$  events (rather than just zero events)



Extended to deal with larger event numbers (arXiv:0709.2701)

Combining these upper limits (arXiv:1105.2928)

Method used by CDMS, CRESST, Edelweiss

(If bgd is known, better to use different method)



# DESIRABLE PROPERTIES

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

# Ilya Narsky, FNAL CLW 2000

Poisson counting expt

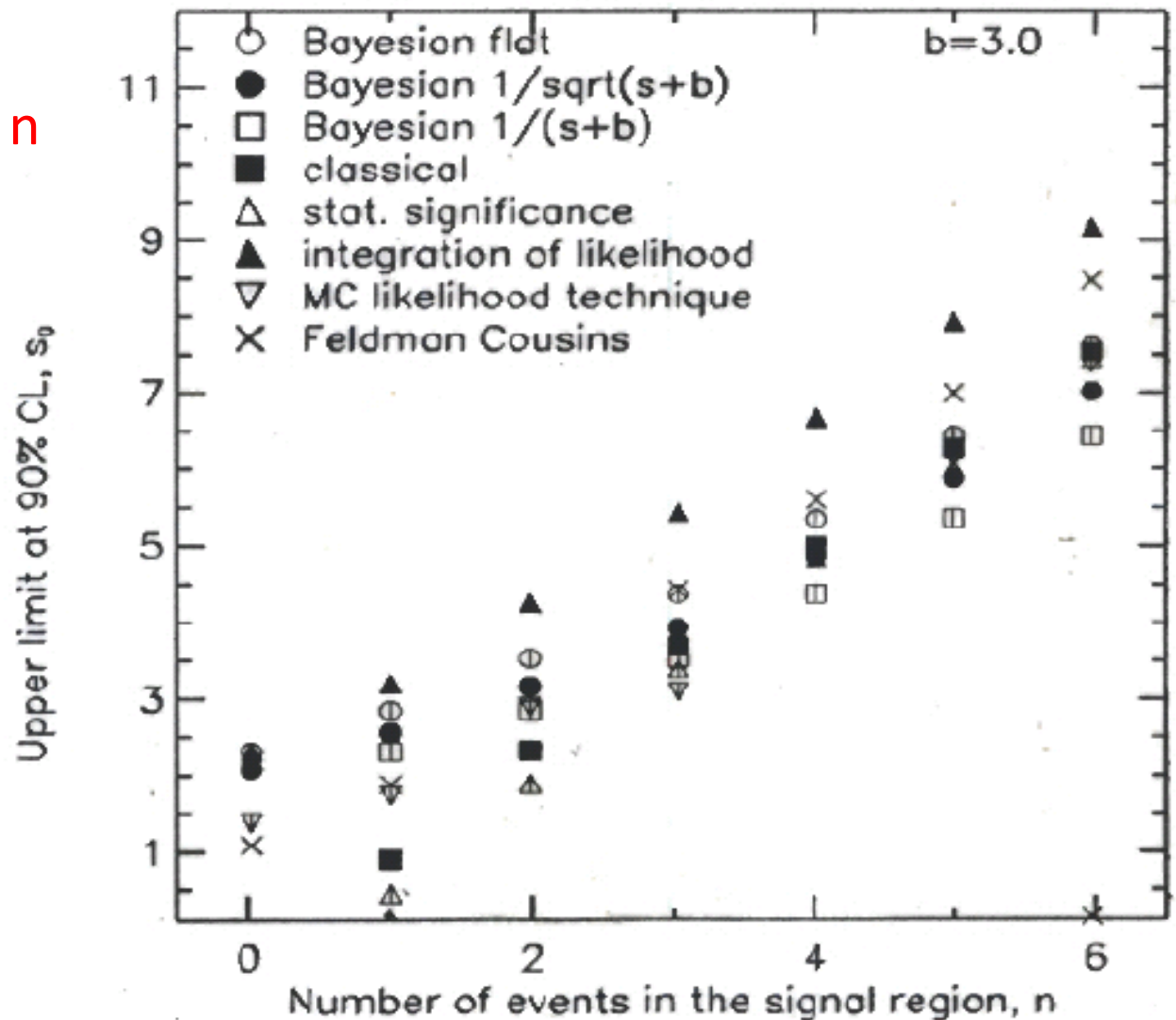
Poisson param =  $s+b$

$b = 3.0$

Observed number =  $n$

What is UL for  $s$ ?

Most dramatic  
differences when  $n < b$



# Conclusions

## Resources:

Software exists: e.g. RooStats

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

‘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.

Neutrino expts might go for combined Statistics Committee. Is that appropriate for direct DM experiments?

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

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

**“Good luck”**



**BACK-UP**

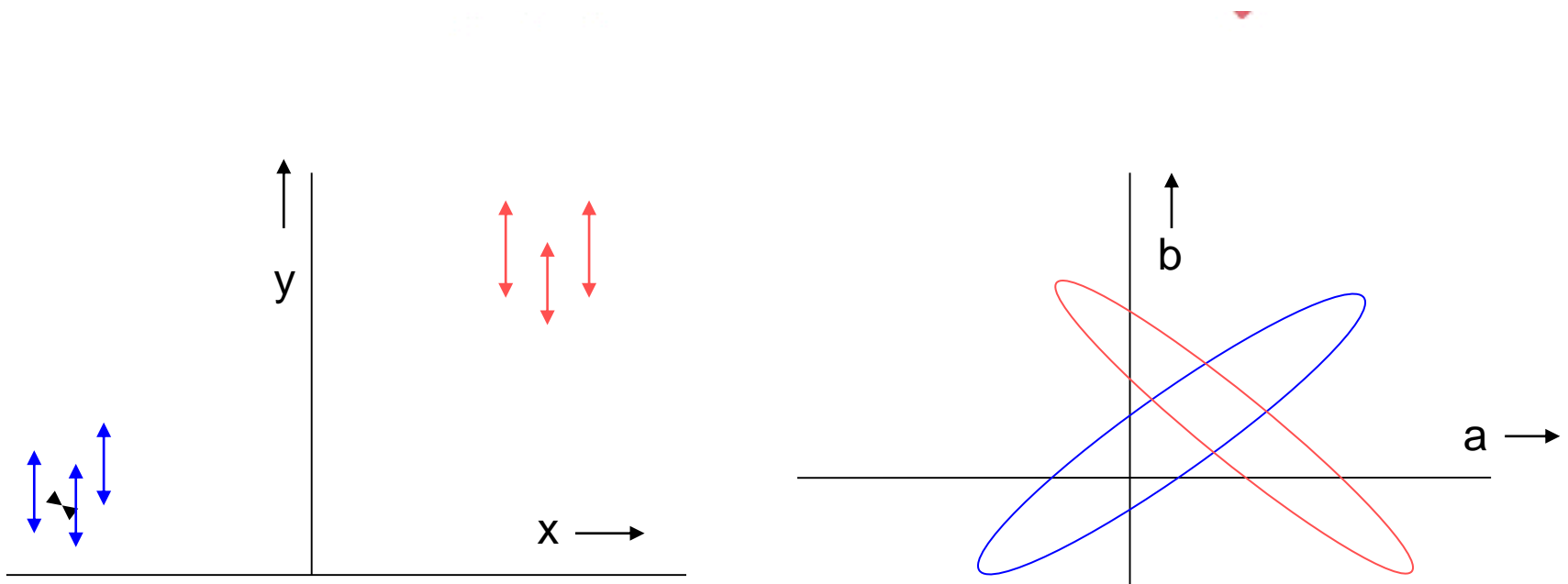
# COMBINING RESULTS

- Better to combine data than combine results  
(Problems with non-Gaussian estimates  
dealing with correlations  
uncertainty estimates)
- Beware of uncertainty estimates that depend on parameter estimate  
e.g.  $n \pm \sqrt{n}$      $100 \pm 10$  and  $80 \pm 9$   
or  $\tau \pm \tau/\sqrt{N}$      $1.00 \pm 0.10$  and  $1.20 \pm 0.12$  (N=100)

# Combining: oddities

- 1 variable :  
Best combination of 2 correlated measurements can be outside range of measurements
- 2 variables,  $\alpha$   $\beta$   
Uncertainties on  $\alpha_{\text{best}}$  and  $\beta_{\text{best}}$  much smaller than individual uncertainties.
- 2 variables,  $\alpha$   $\beta$   
 $\alpha_{\text{best}} > \alpha_1$  and  $\alpha_2$     $\beta_{\text{best}} > \beta_1$  and  $\beta_2$

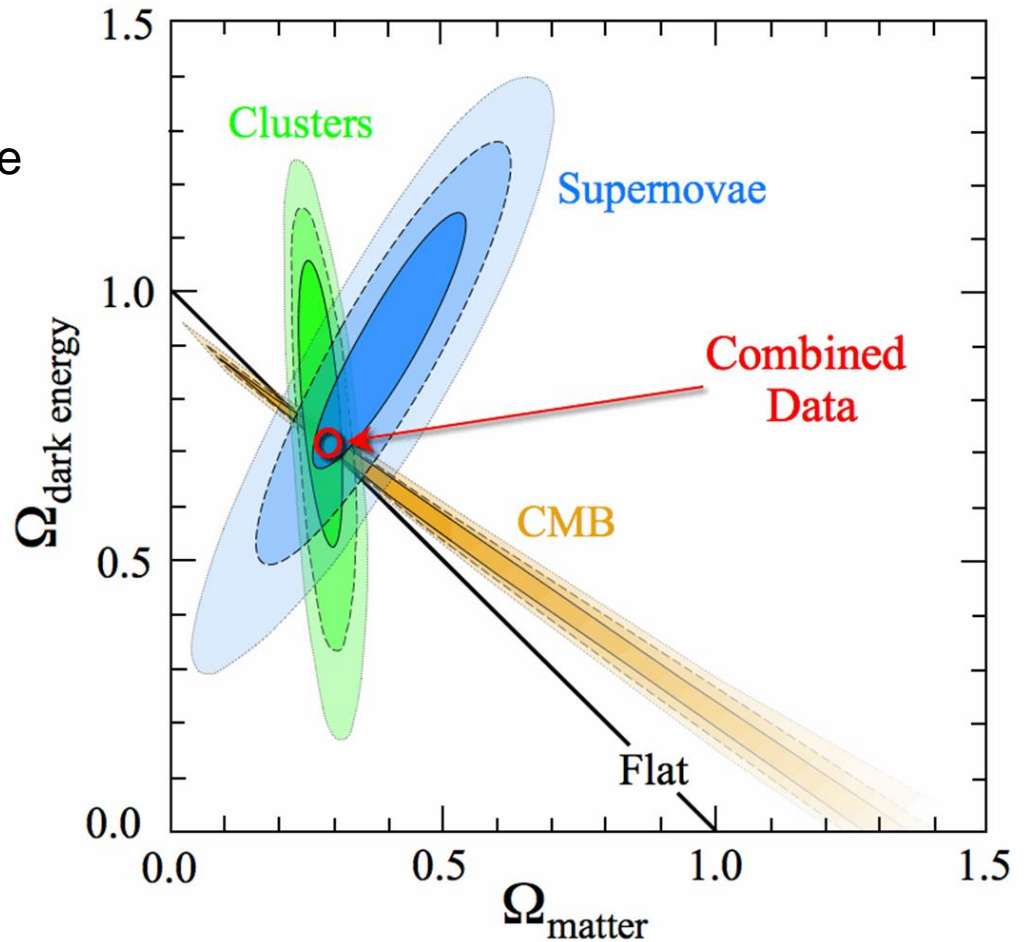
Straight line fit to red points has large uncertainties on intercept and on gradient  
Straight line fit to blue points has large uncertainties on intercept and on gradient  
Combined straight line fit to red and blue points has much smaller uncertainties on intercept and on gradient





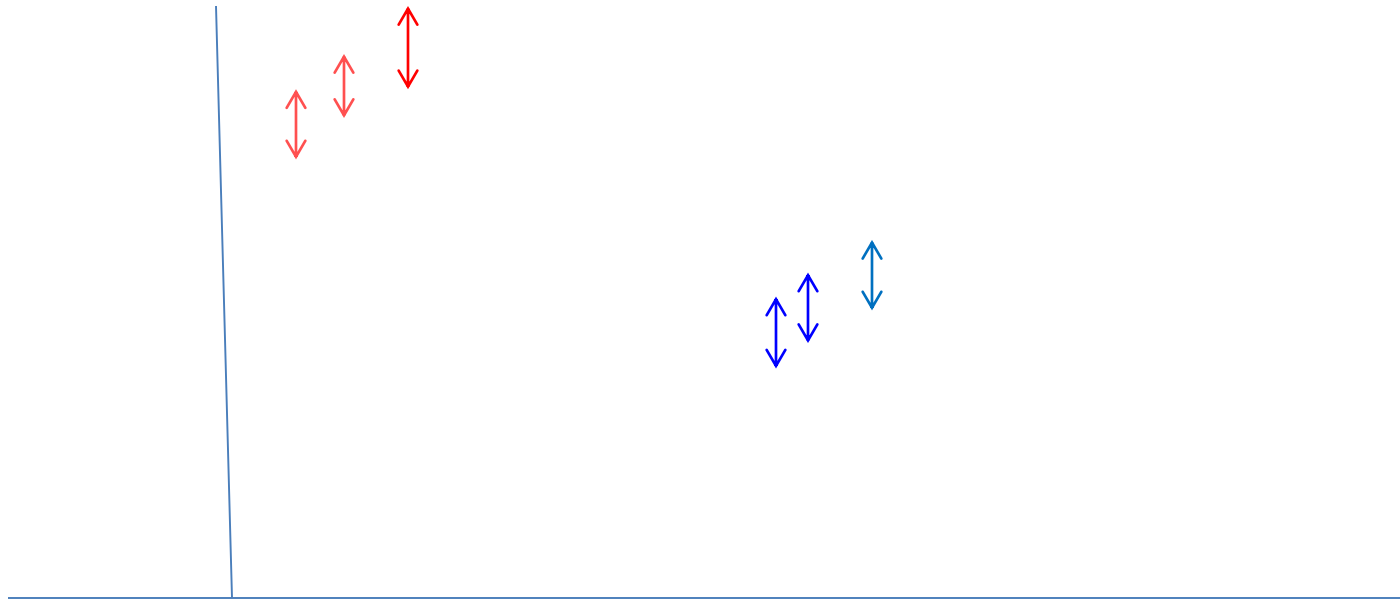
# Uncertainty on $\Omega_{\text{dark energy}}$

When combining pairs of variables, the uncertainties on the **combined parameters** can be **much** smaller than any of **the individual** uncertainties  
e.g.  $\Omega_{\text{dark energy}}$



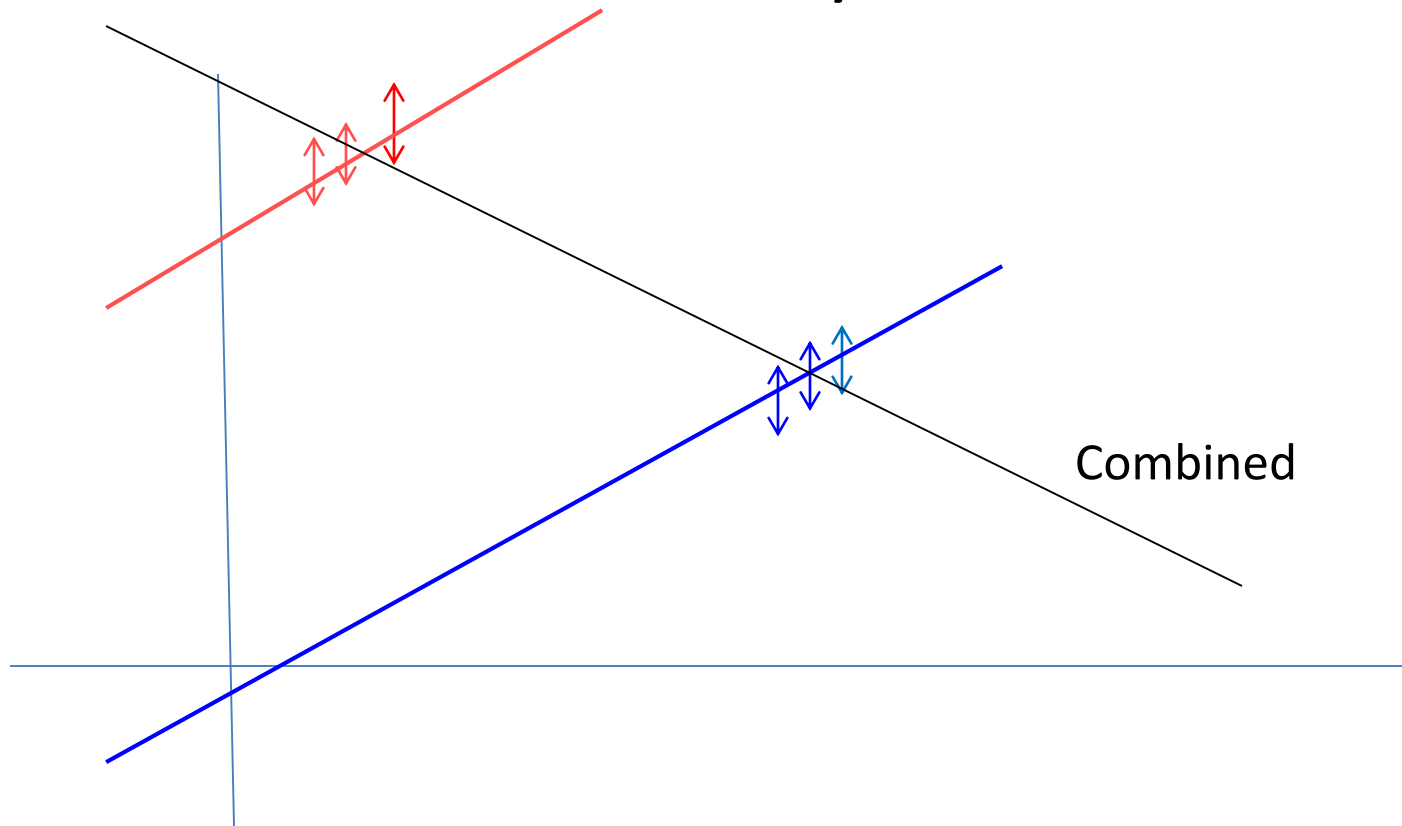
# Best values of params a and b outside range of individual values

$$y = a + bx$$



# Best values of params a and b outside range of individual values

$$y = a + bx$$



# Likelihoods

Here just for parameter determination  
Also very important for Hypothesis Testing,  
in Bayesian and Frequentist approaches

Procedure:

Write down  $P(\text{data} \mid \text{hypothesis}' \text{ param})$

pdf: Regard this as fn of data, for fixed param values

Likelihood: Fn of parameter, for given data

e.g. Poisson  $P(n \mid \mu) = e^{-\mu} \mu^n / n!$

Data:

Can be individual values. Does not have to be a histogram

# Simple example of Likelihood: Angular distribution

$$y = N (1 + \beta \cos^2\theta) \quad N = 1/\{2(1+\beta/3)\}$$

$$y_i = N (1 + \beta \cos^2\theta_i)$$

= probability density of observing  $\theta_i$ , given  $\beta$

$$\mathcal{L}(\beta) = \prod y_i$$

= probability density of observing the data set  $y_i$ , given  $\beta$

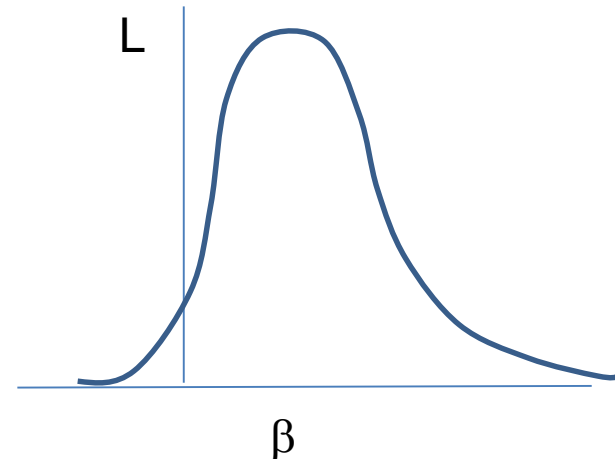
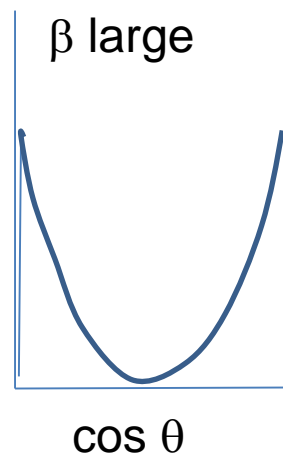
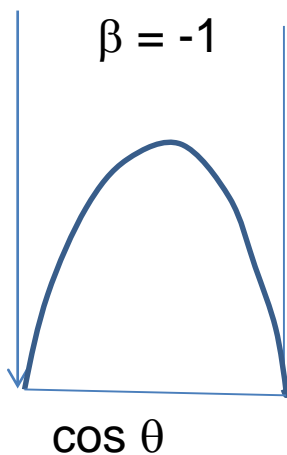
Best estimate of  $\beta$  is that which maximises  $\mathcal{L}$

Values of  $\beta$  for which  $\mathcal{L}$  is very small are ruled out

Precision of estimate for  $\beta$  comes from width of  $\mathcal{L}$  distribution

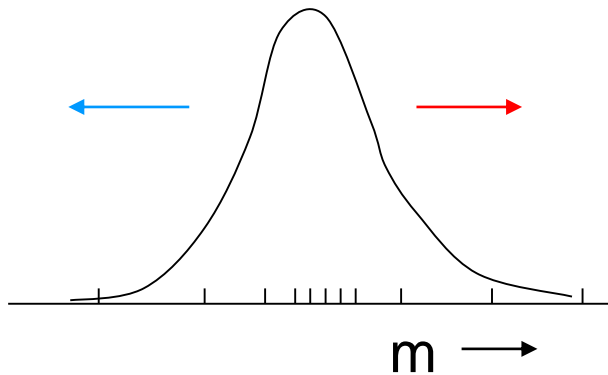
\*\*\*\*\* **CRUCIAL** to normalise  $y$   $N = 1/\{2(1 + \beta/3)\}$

(Information about parameter  $\beta$  comes from **shape** of exptl distribution of  $\cos\theta$ )

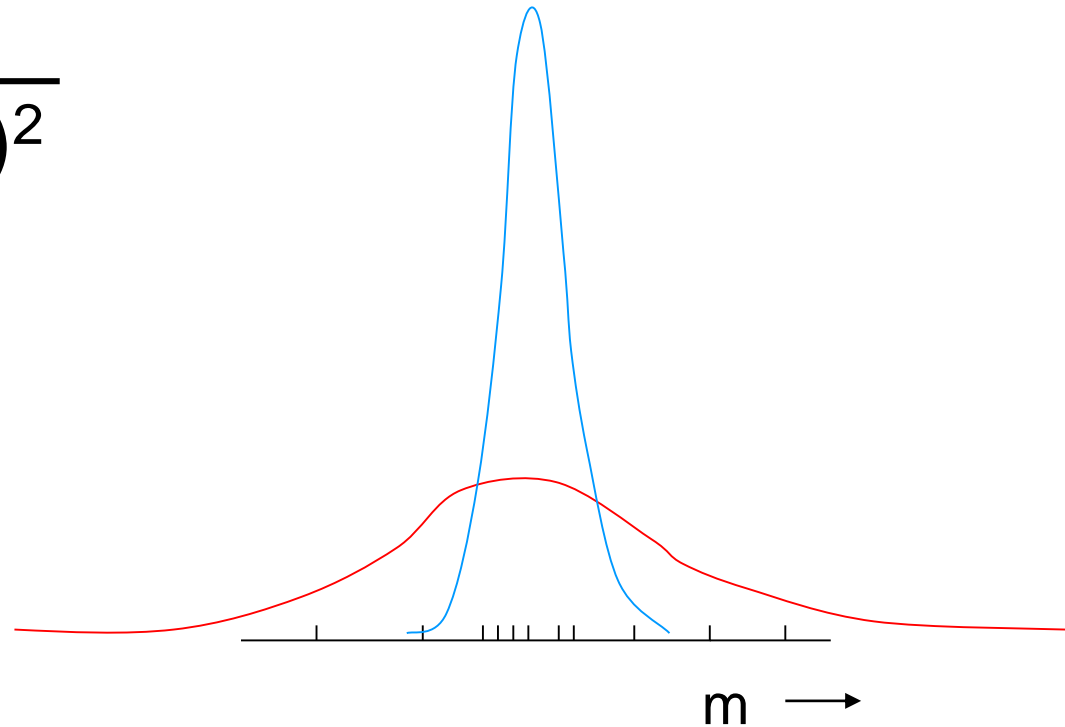


# How it works: Resonance

$$y \sim \frac{\Gamma/2}{(m-M_0)^2 + (\Gamma/2)^2}$$



Vary  $M_0$

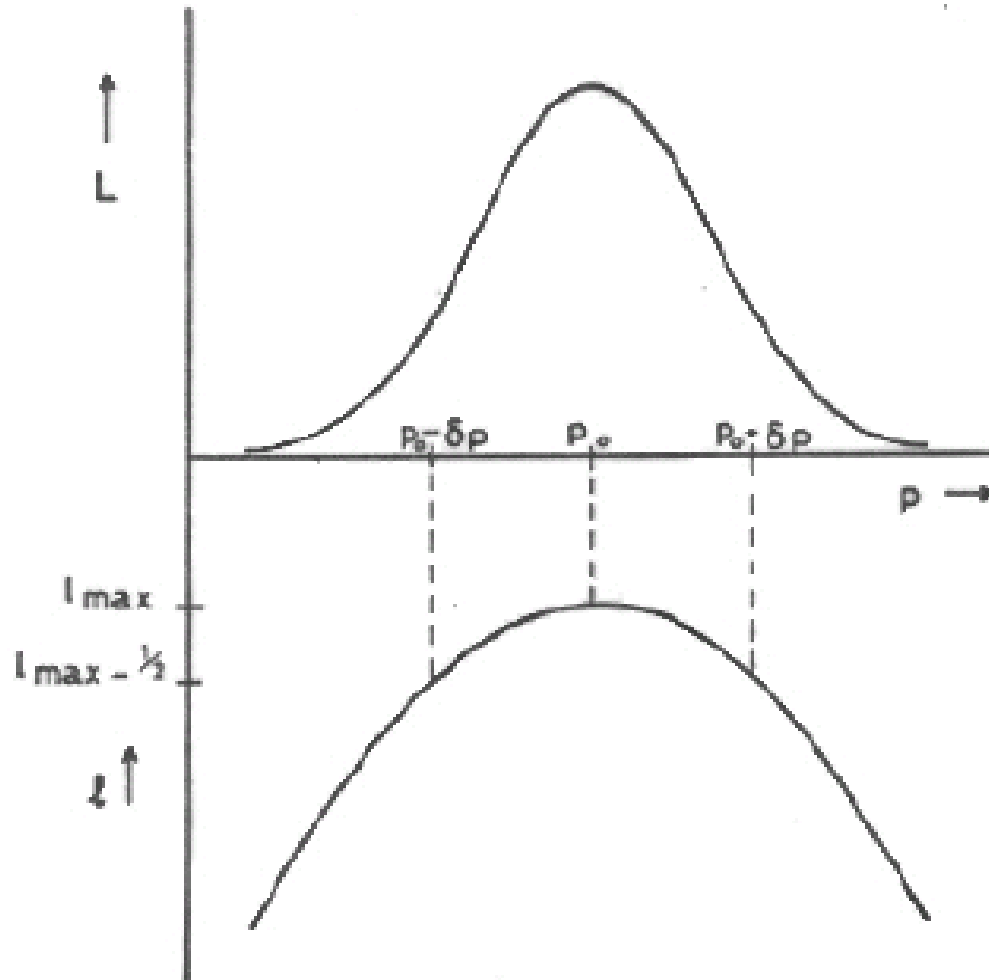


Vary  $\Gamma$

Conventional to consider

$$\ell = \ln(\mathcal{L}) = \sum \ln(y_i)$$

For large  $N$ ,  $\mathcal{L} \rightarrow$  Gaussian



# $\Delta \ln \mathcal{L} = -1/2$ rule

If  $\mathcal{L}(\mu)$  is Gaussian, following definitions of  $\sigma$  are equivalent:

1) RMS of  $\mathcal{L}(\mu)$

2)  $1/\sqrt{-d^2 \ln \mathcal{L} / d\mu^2}$

3)  $\ln(\mathcal{L}(\mu_0 \pm \sigma)) = \ln(\mathcal{L}(\mu_0)) - 1/2$

If  $\mathcal{L}(\mu)$  is non-Gaussian, these are no longer the same

~~“Procedure 3) above still gives interval that contains the true value of parameter  $\mu$  with 68% probability”~~

Heinrich: CDF note 6438 (see CDF Statistics Committee Web-page)

Barlow: Phystat05



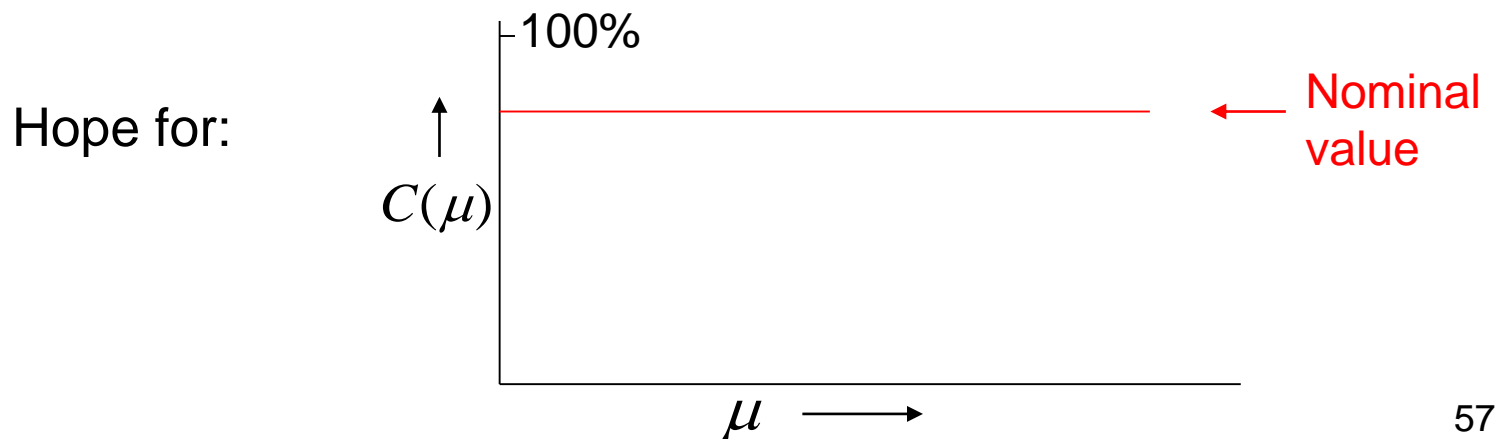
## COVERAGE

How often does quoted range for parameter include param's true value?

N.B. Coverage is a property of **METHOD**, not of a particular exptl result

Coverage can vary with  $\mu$

Study coverage of different methods for Poisson parameter  $\mu$ , from observation of number of events  $n$



# COVERAGE

If true for all  $\mu$  : “correct coverage”

$P < \alpha$  for some  $\mu$  “undercoverage”  
(this is serious !)

$P > \alpha$  for some  $\mu$  “overcoverage”

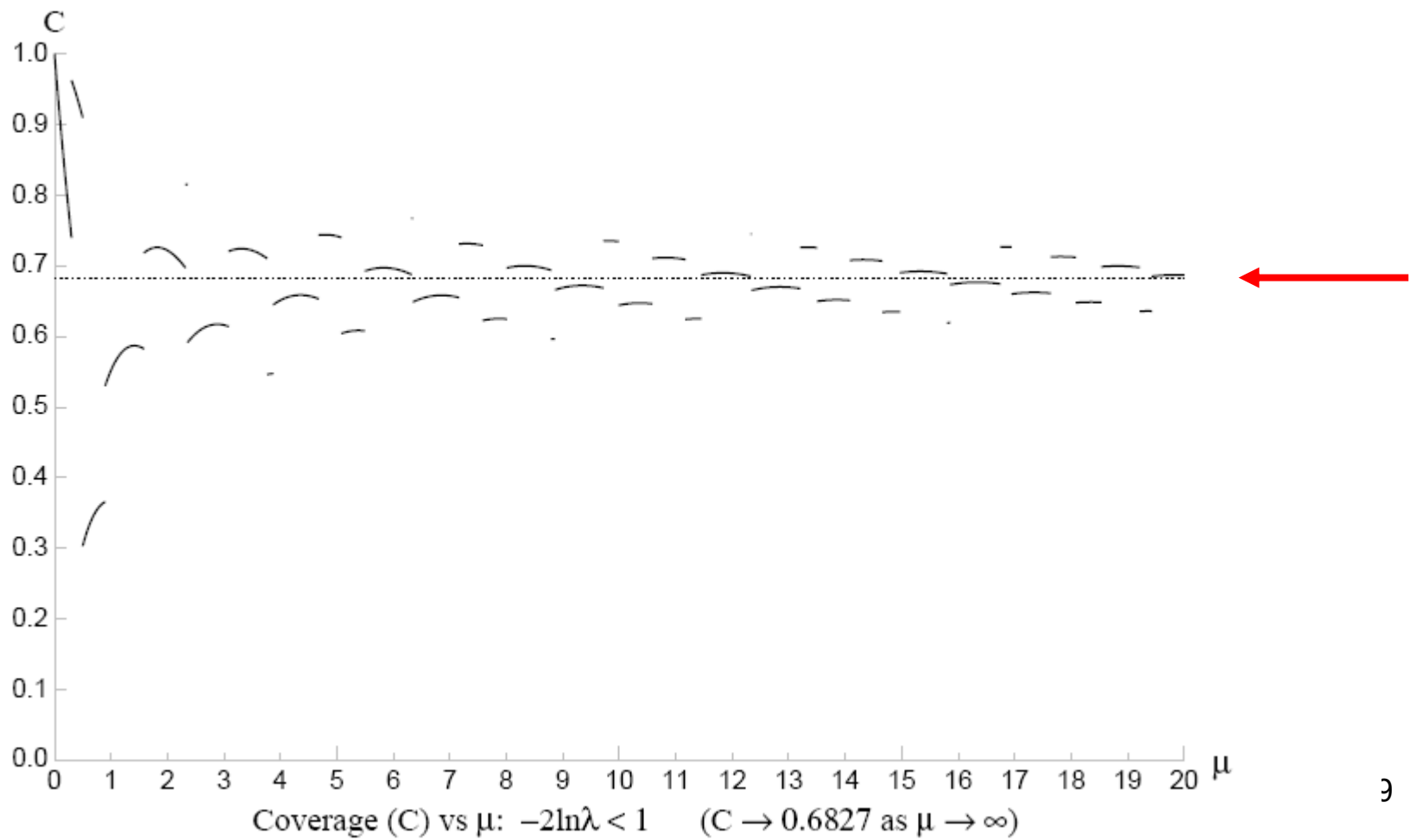
Conservative

Loss of rejection  
power

# Coverage : $\mathcal{L}$ approach (Not frequentist)

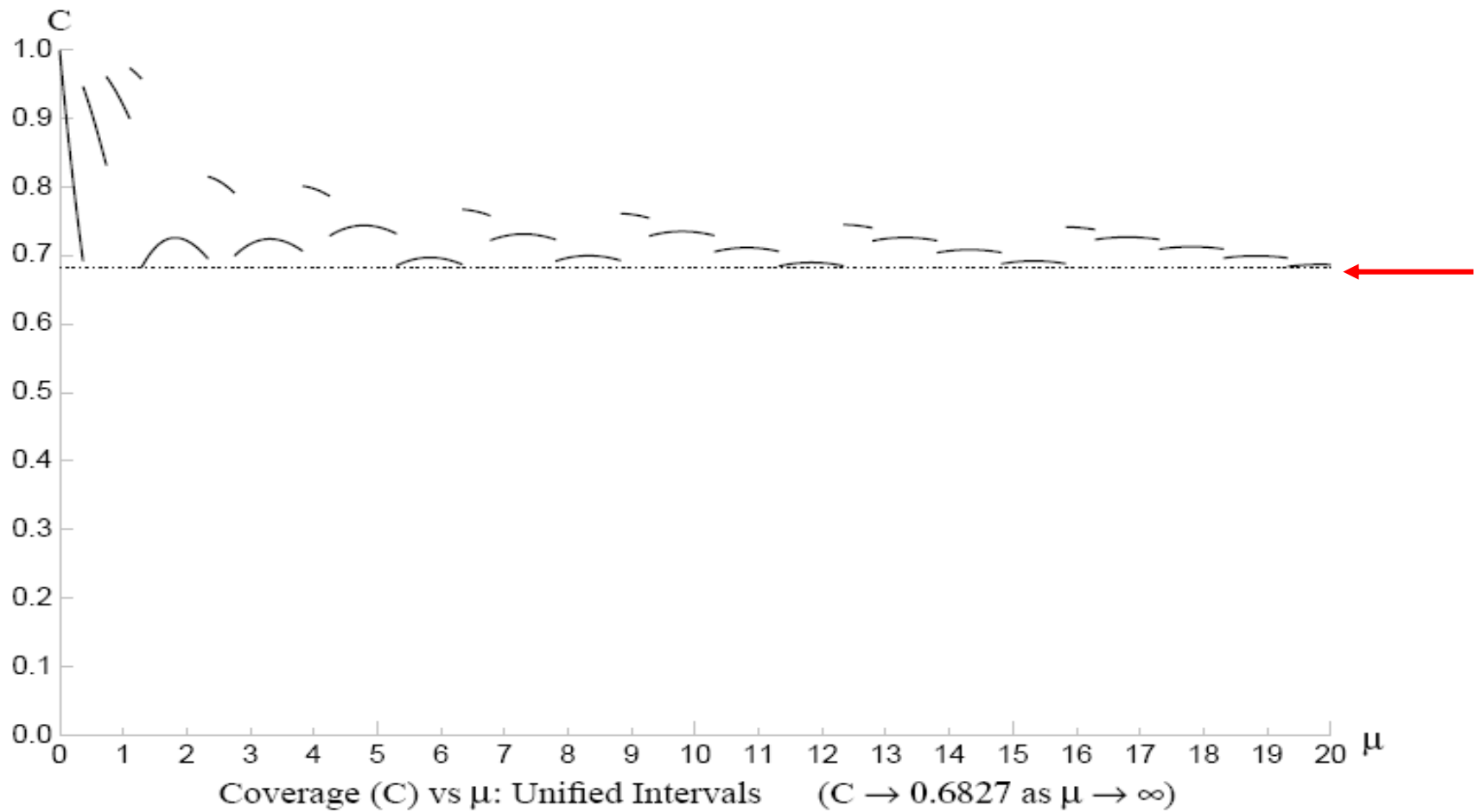
$P(n, \mu) = e^{-\mu} \mu^n / n!$  (Joel Heinrich CDF note 6438)

$-2 \ln \lambda < 1$        $\lambda = P(n, \mu) / P(n, \mu_{\text{best}})$       **UNDERCOVERS**



# Feldman-Cousins Unified intervals

Neyman construction so NEVER undercovers



	Moments	Max Like	Least squares
Easy?	Yes, if...	Normalisation, maximisation messy	Minimisation
Efficient?	Not very	Usually best	Sometimes = Max Like
Input	Separate events	<b>Separate events</b>	Histogram
Goodness of fit	Messy	No (unbinned)	Easy
Constraints	No	Yes	Yes
N dimensions	Easy if ....	Norm, max messier	Easy
Weighted events	Easy	Errors difficult	Easy
Bgd subtraction	Easy	Troublesome	Easy
Uncertainty estimates	Observed spread, or analytic	$\left\{ -\frac{\partial^2 l}{\partial p_i \partial p_j} \right\}$	$\left\{ \frac{\partial^2 S}{2 \partial p_i \partial p_j} \right\}$
Main feature	Easy	Best for params	<b>Goodness of Fit</b>



## BAYES and FREQUENTISM: Different views of probability

# We need to make a statement about Parameters, Given Data

The basic difference between the two:

Bayesian : **Probability (parameter, given data)**  
(an anathema to a Frequentist!)

Frequentist : **Probability (data, given parameter)**  
(a likelihood function)

# PROBABILITY

## MATHEMATICAL

Formal

Based on Axioms

## FREQUENTIST

Ratio of frequencies as  $n \rightarrow$  infinity

Repeated “identical” trials

Not applicable to **single event** or **physical constant**

## BAYESIAN Degree of belief

Can be applied to single event or physical constant

(even though these have unique truth)

Varies from person to person \*\*\*

Quantified by “fair bet”



# Bayesian versus Classical

## Bayesian

$$P(A \text{ and } B) = P(A;B) \times P(B) = P(B;A) \times P(A)$$

{ If A and B independent,  $P(A;B) = P(A) \rightarrow P(A \text{ and } B) = P(A) P(B)$  }

e.g. A = event contains t quark

B = event contains W boson

or A = I am in CERN

B = I am giving a lecture

$$P(A;B) = P(B;A) \times P(A) / P(B) \quad \text{Bayes' Theorem}$$

Completely uncontroversial, provided....

# Bayesian

$$P(A; B) = \frac{P(B; A) \times P(A)}{P(B)}$$

Bayes'  
Theorem

$$p(\text{param} \mid \text{data}) \propto p(\text{data} \mid \text{param}) * p(\text{param})$$

↑  
posterior

↑  
likelihood

↑  
prior

Problems:  $p(\text{param})$  Has particular value  
“Degree of belief”

**Prior** What functional form?

**Coverage**

P(parameter)      **Has specific value**

“Degree of Belief”

**Credible interval**

Prior:      **What functional form?**

Uninformative prior:    flat?

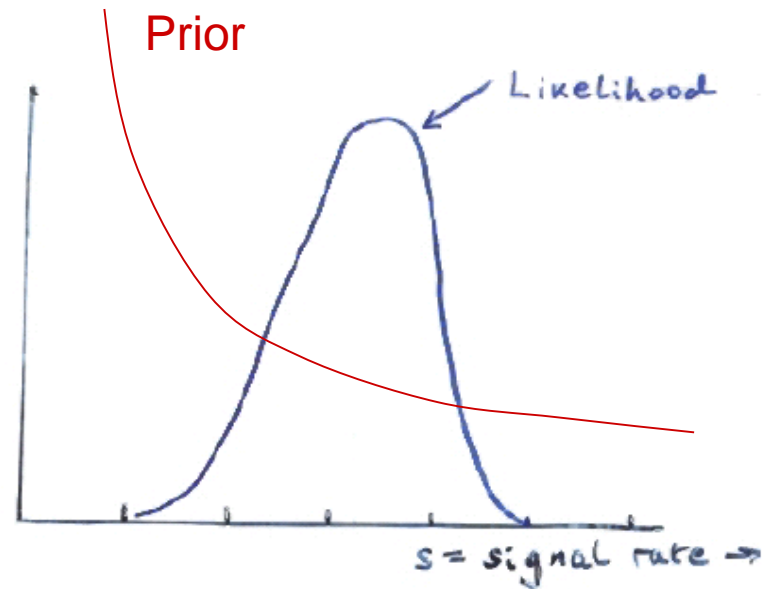
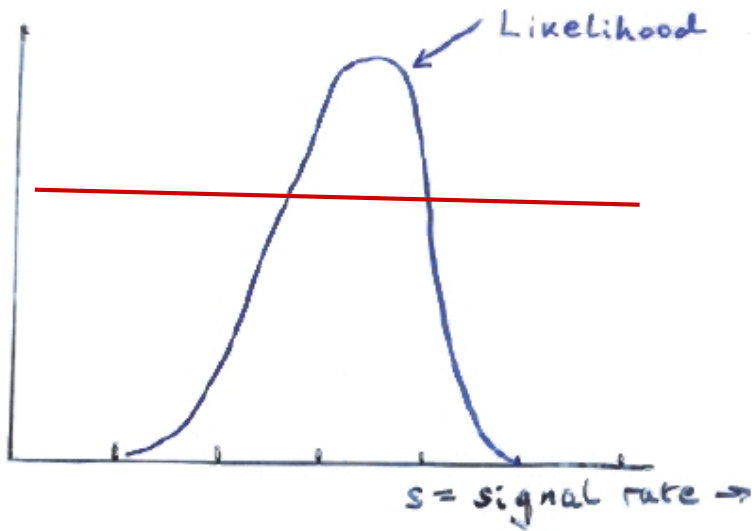
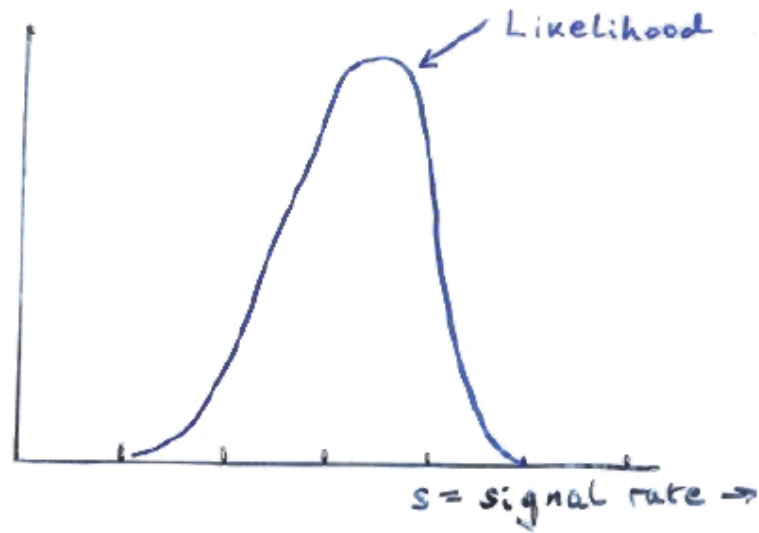
In which variable?    e.g.  $m$ ,  $m^2$ ,  $\ln m$ ,....?

Even more problematic with more params

**Unimportant** if “**data overshadows prior**”

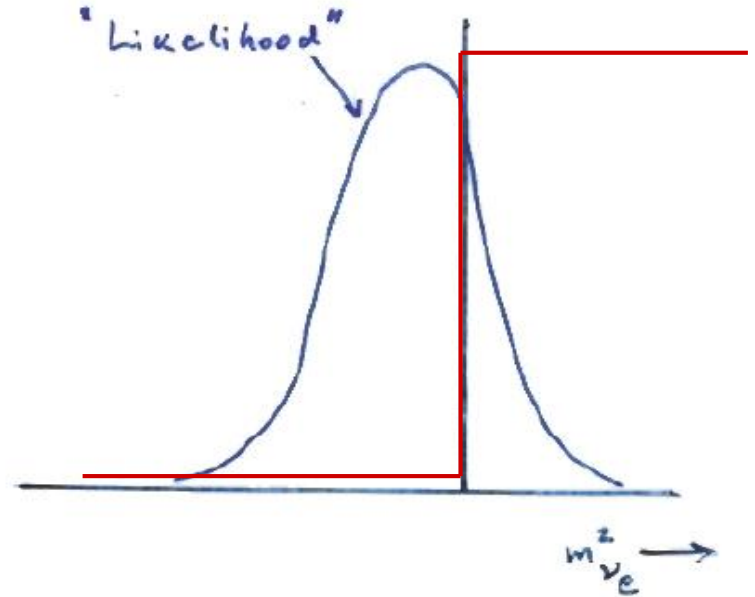
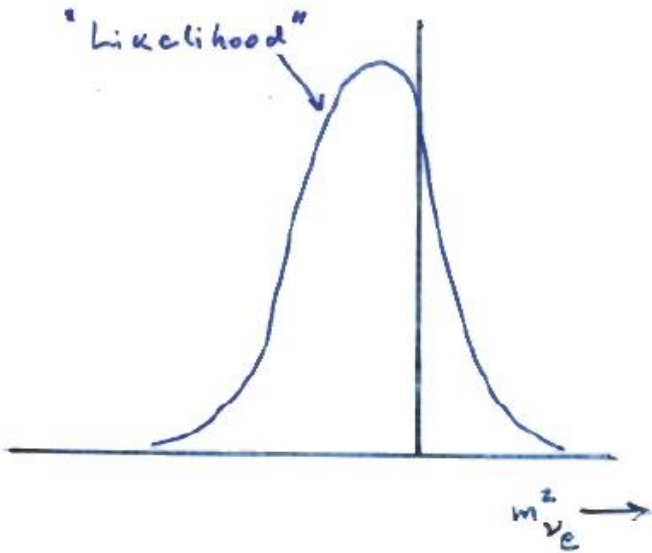
**Important** for limits

**Subjective** or **Objective** prior?



Even more important for **UPPER LIMITS**

# Mass-squared of neutrino



Prior = zero in unphysical region

# Bayes: Specific example

Particle decays exponentially:  $dn/dt = (1/\tau) \exp(-t/\tau)$

Observe 1 decay at time  $t_1$ :  $\mathcal{L}(\tau) = (1/\tau) \exp(-t_1/\tau)$

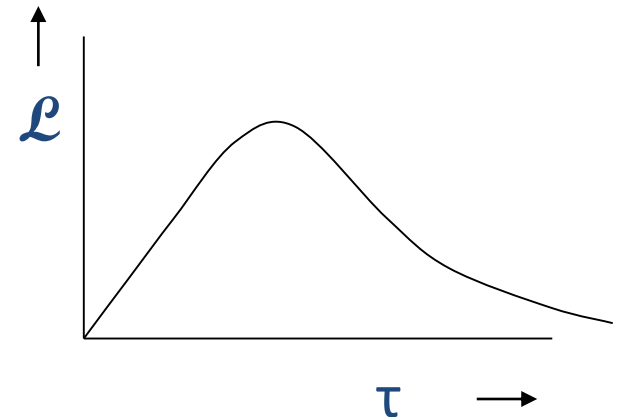
Choose prior  $\pi(\tau)$  for  $\tau$

e.g. constant up to some large  $\tau$

Then posterior  $p(\tau) = \mathcal{L}(\tau) * \pi(\tau)$

has almost same shape as  $\mathcal{L}(\tau)$

Use  $p(\tau)$  to choose interval for  $\tau$  in usual way



Contrast frequentist method for same situation later.

# Classical Approach

Neyman “confidence interval” avoids pdf for  $\mu$

Uses only  $P(x; \mu)$

Confidence interval  $\mu_1 \rightarrow \mu_2$  :

$P(\mu_1 \rightarrow \mu_2 \text{ contains } \mu_t) = \alpha$  True for any  $\mu_t$



Varying intervals  
from ensemble of  
experiments

fixed

Gives range of  $\mu$  for which observed value  $x_0$  was “likely” ( $\alpha$ )

Contrast Bayes : Degree of belief =  $\alpha$  that  $\mu_t$  is in  $\mu_1 \rightarrow \mu_2$

# Classical (Neyman) Confidence Intervals

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

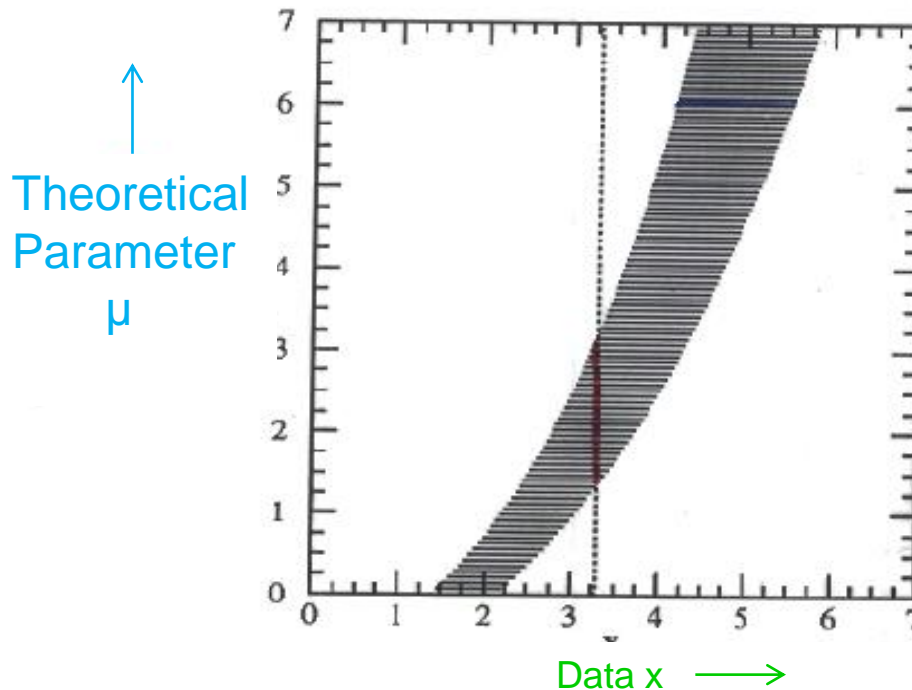


FIG. 1. A generic confidence belt construction and its use. For each value of  $\mu$ , one draws a horizontal acceptance interval  $[x_1, x_2]$  such that  $P(x \in [x_1, x_2] | \mu) = \alpha$ . Upon performing an experiment to measure  $x$  and obtaining the value  $x_0$ , one draws the dashed vertical line through  $x_0$ . The confidence interval  $[\mu_l, \mu_u]$  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  $P(\text{data}|\text{theory})$

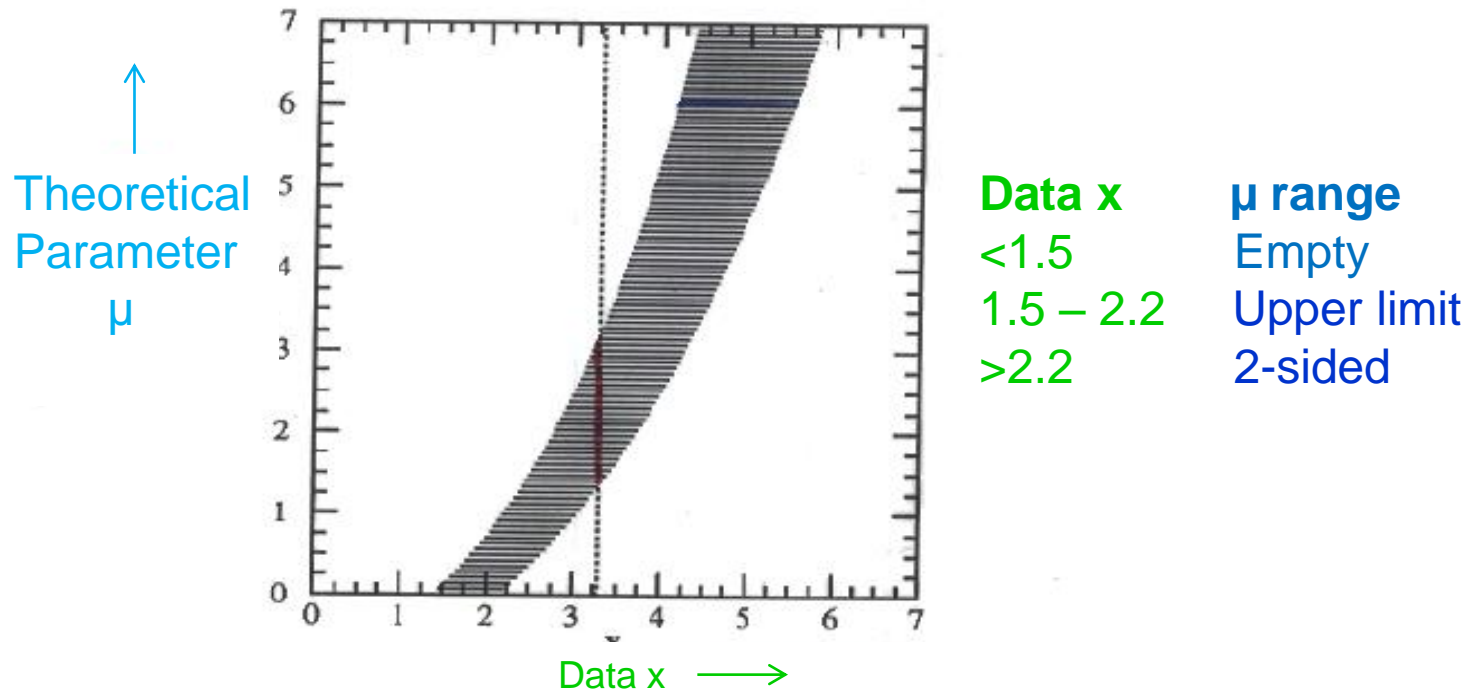


FIG. 1. A generic confidence belt construction and its use. For each value of  $\mu$ , one draws a horizontal acceptance interval  $[x_1, x_2]$  such that  $P(x \in [x_1, x_2] | \mu) = \alpha$ . Upon performing an experiment to measure  $x$  and obtaining the value  $x_0$ , one draws the dashed vertical 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  $\mu$

# 90% Classical interval for Gaussian

$$\sigma = 1 \quad \mu \geq 0$$

e.g.  $m^2(v_e)$ , length of small object

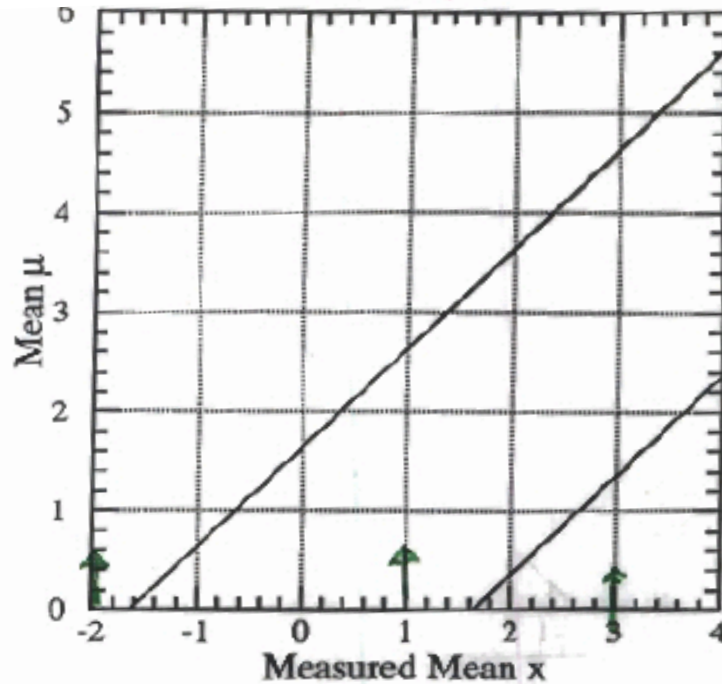


FIG. 3. Standard confidence belt for 90% C.L. central confidence intervals for the mean of a Gaussian, in units of the rms deviation.

$x_{\text{obs}}=3$  Two-sided range

$x_{\text{obs}}=1$  Upper limit

$x_{\text{obs}}=-1$  No region for  $\mu$

Other methods have different behaviour at negative  $x$

# 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

Feldman-Cousins  
90% conf intervals

Uses different  
ordering rule

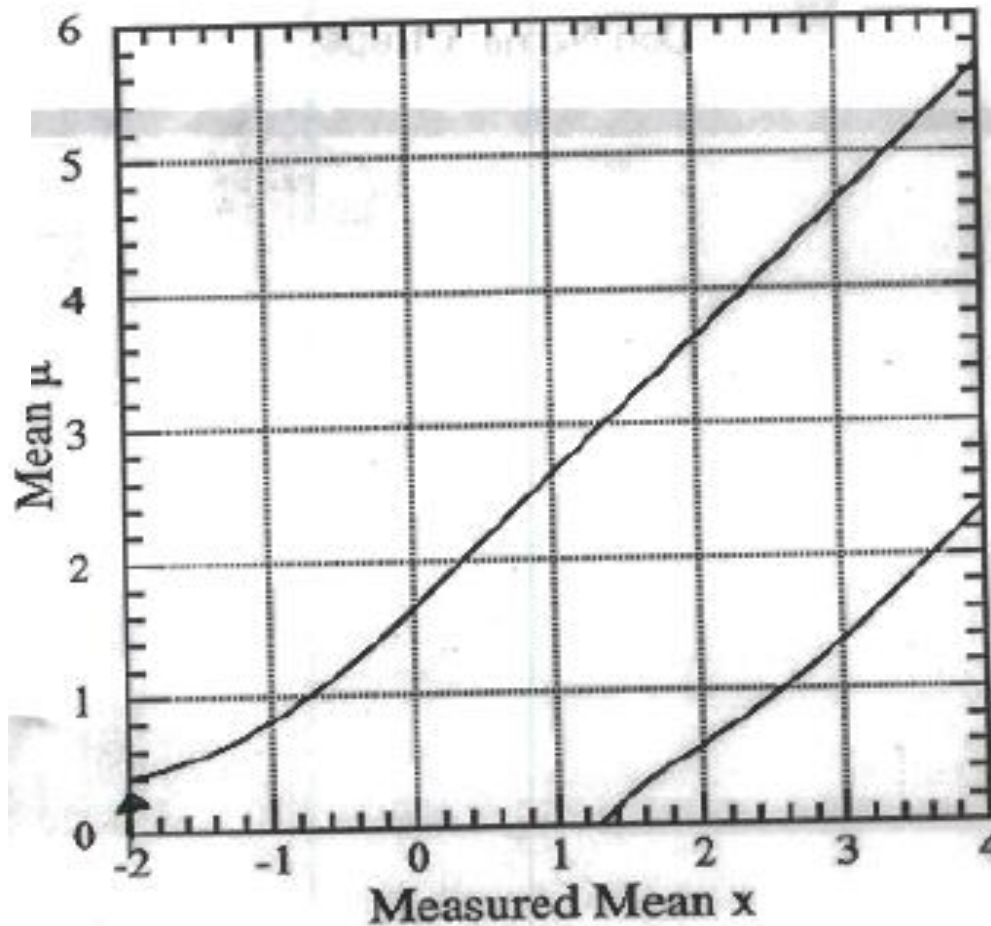


FIG. 10. Plot of our 90% confidence intervals for mean of a Gaussian, constrained to be non-negative, described in the text.

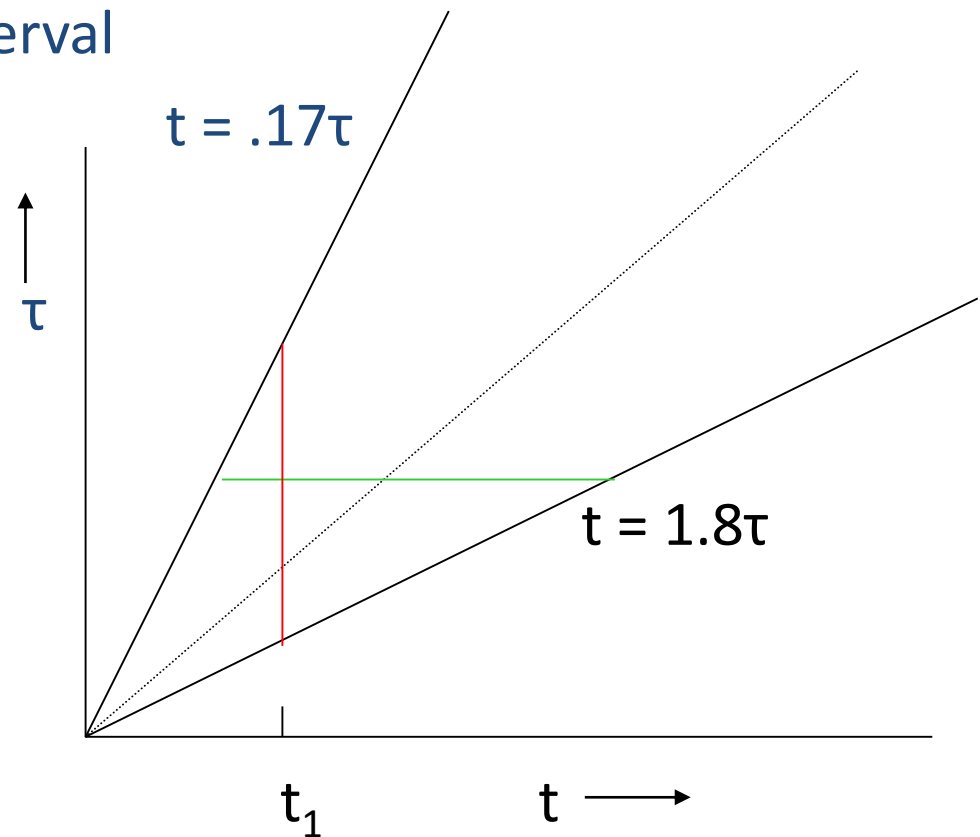
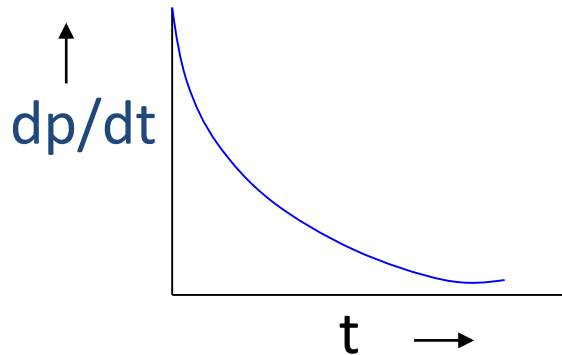
$X_{\text{obs}} = -2$  now gives upper limit

# Frequentism: Specific example

Particle decays exponentially:  $dp/dt = (1/\tau) \exp(-t/\tau)$

Observe 1 decay at time  $t_1$ :  $\mathcal{L}(\tau) = (1/\tau) \exp(-t_1/\tau)$

Construct 68% central interval



68% conf. int. for  $\tau$  from  
 $t_1/1.8 \rightarrow t_1/0.17$

$$\mu_l \leq \mu \leq \mu_u \quad \text{at 90\% confidence}$$

Frequentist

$\mu_l$  and  $\mu_u$  known, but random  
 $\mu$  unknown, but fixed  
Probability statement about  $\mu_l$  and  $\mu_u$

Bayesian

$\mu_l$  and  $\mu_u$  known, and fixed  
 $\mu$  unknown, and random  
Probability/credible statement about  $\mu$

# MULTIVARIATE ANALYSIS

Example: Aim to separate signal from background

Neyman-Pearson Lemma:

Imagine all possible contours that select signal with efficiency  $\varepsilon$  (Loss = Error of 1<sup>st</sup> Kind)

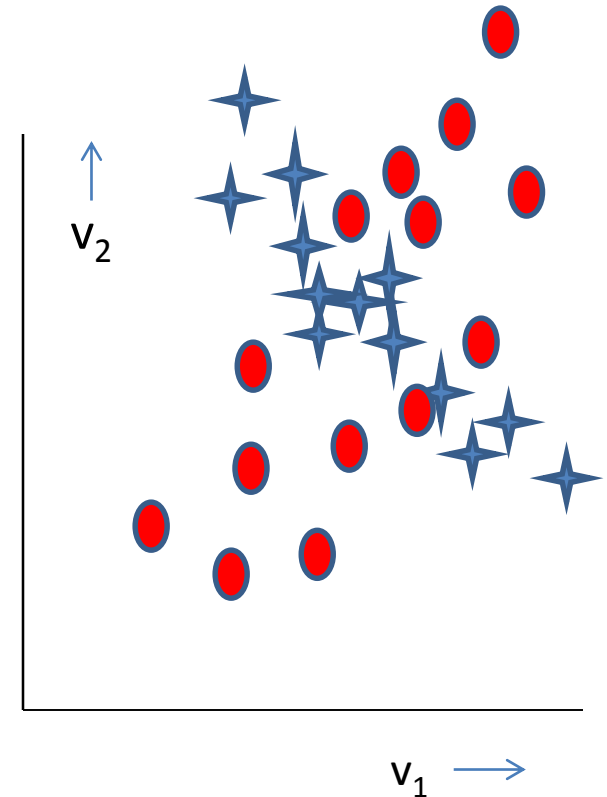
Best is one containing minimal amount of background (Contamination = Error of 2<sup>nd</sup> Kind)

Equivalent to ordering data by

$$\mathcal{L}\text{-ratio} = \mathcal{L}_s(v_1, v_2, \dots) / \mathcal{L}_b(v_1, v_2, \dots)$$

IF variables are independent

$$\mathcal{L}\text{-ratio} = \{\mathcal{L}_s(v_1)/\mathcal{L}_b\{v_1\}\} \times \{\mathcal{L}_s(v_2)/\mathcal{L}_b\{v_2\}\} \times \dots$$



# PROBLEM:

Don't know  $\mathcal{L}$ -ratio exactly because:

- 1) Signal & bkg generated by M.C. with finite statistics
- 2) Nuisance params (systematics) and signal params
- 3) Neglected sources of bkg
- 4) Hard to implement in many dimensions

## METHODS TO DEAL WITH THIS

Cuts

Kernel Density Estimation

Fisher Discriminant

Principal Component Analysis

Boosted Decision Trees

Support Vector Machines

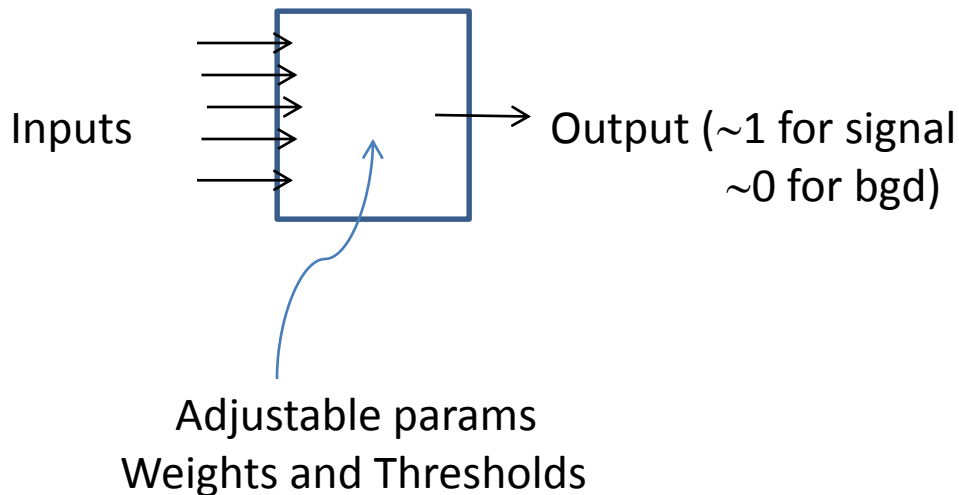
Neural Nets 

Deep Nets

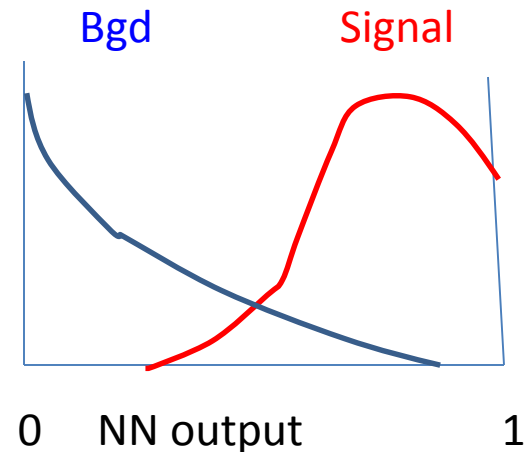


# NEURAL NETWORKS

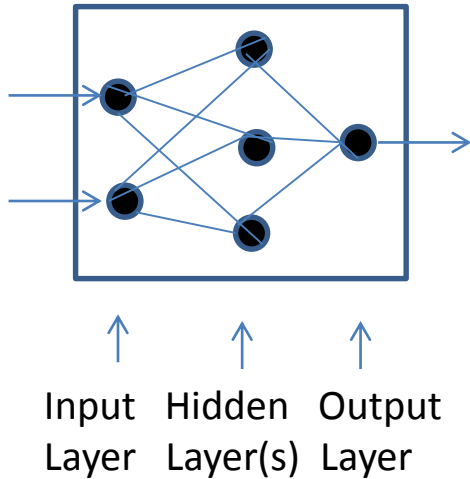
Typical application: Classify events as signal or bgd



- Learning process:  
Input = Known signal & bgd (e.g M.C.)  
Adjust params  $\rightarrow$  'Best' output
- Testing process  
Make sure not 'overtraining'
- Use trained network on actual data  
Classify events as signal if output  $>$  cut



# HOW DOES IT WORK?

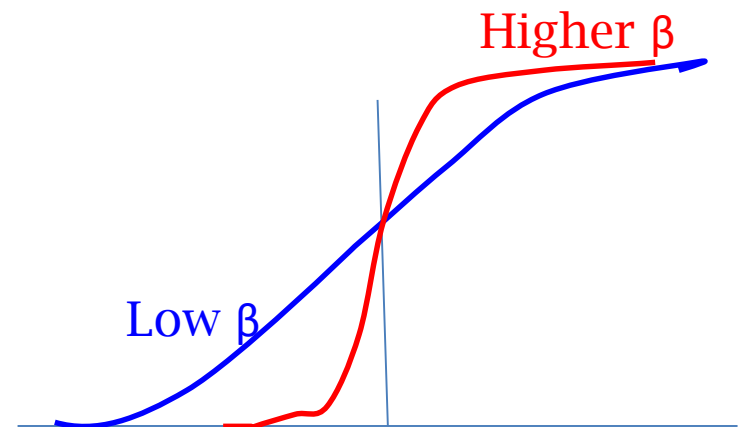


For each hidden or output node  $j$   
$$\text{Output}_j = F [\sum \text{Input}_i * W_{ij} + T_j]$$
  
( $W$  and  $T$  = network params)

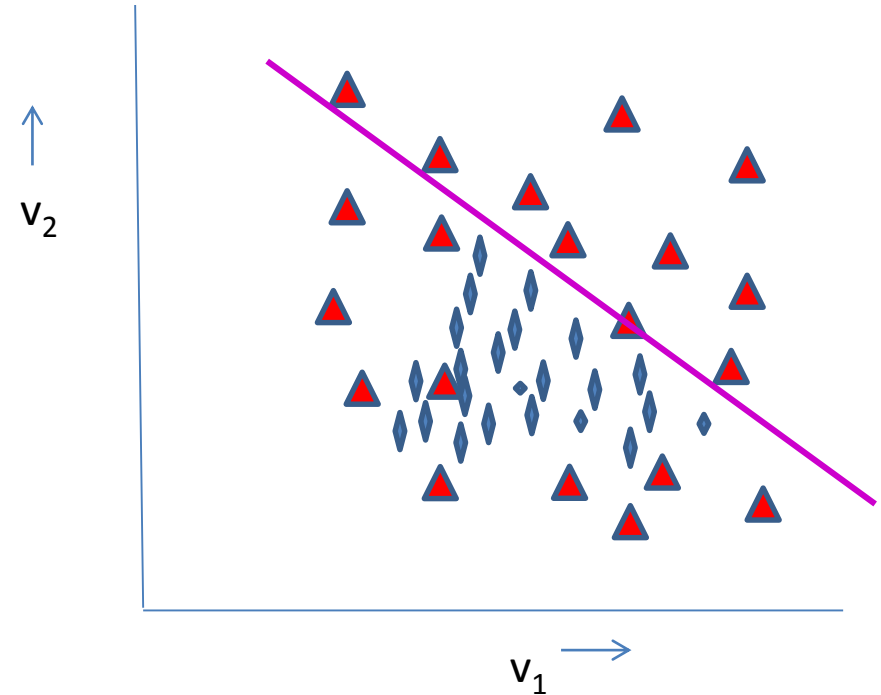
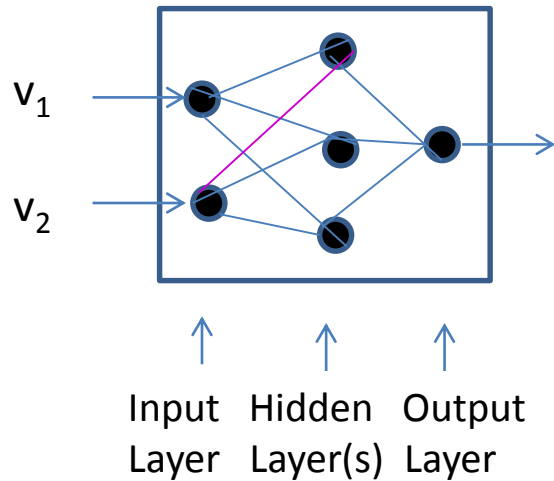
Typical  $F(x) = 1/(1 + e^{-\beta x})$  Sigmoid

For large  $\beta$ , output of node  $j$  is 'ON' if  
$$\sum I_i w_{ij} + T_j > 0$$

This is 'hyper-plane' in  $I$  space ●



# HOW DOES IT WORK?



For First hidden node

Straight line is

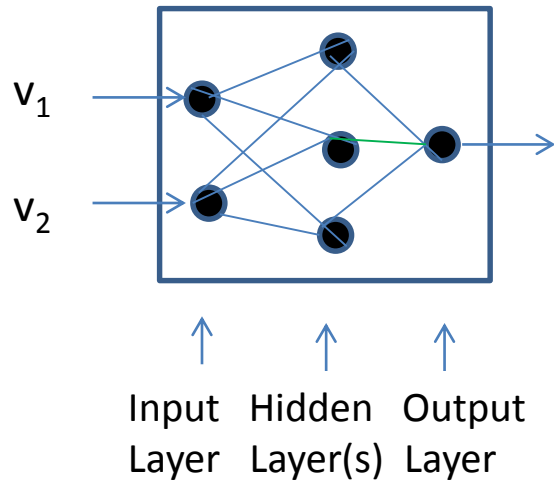
$$w_{11} * v_1 + w_{21} * v_2 + T_{10} = 0$$

where

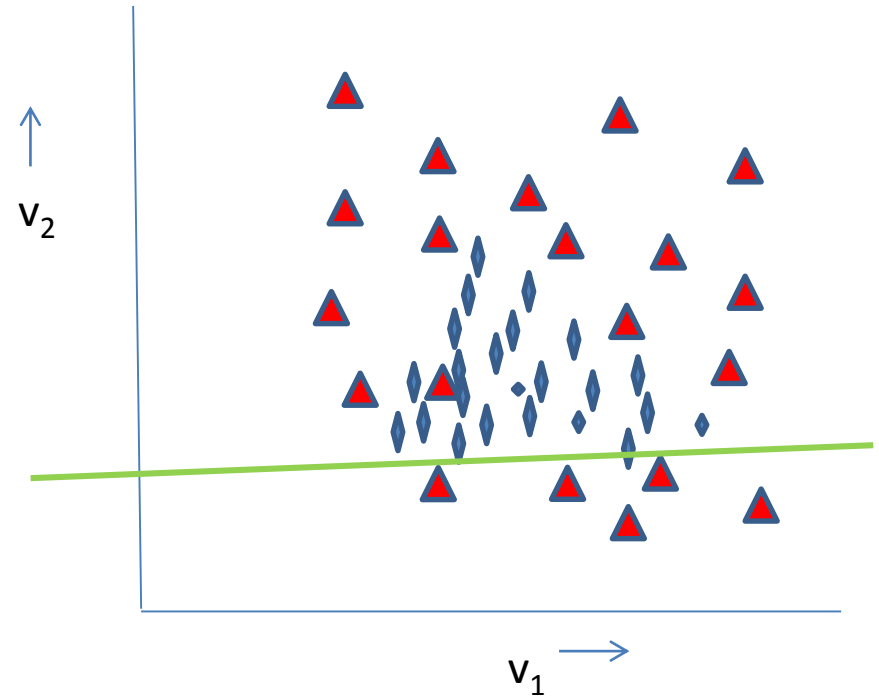
$w_{ij}$  is weight from  $i^{\text{th}}$  input node to  $j^{\text{th}}$  hidden node

$T_{k0}$  is threshold for  $k^{\text{th}}$  hidden node

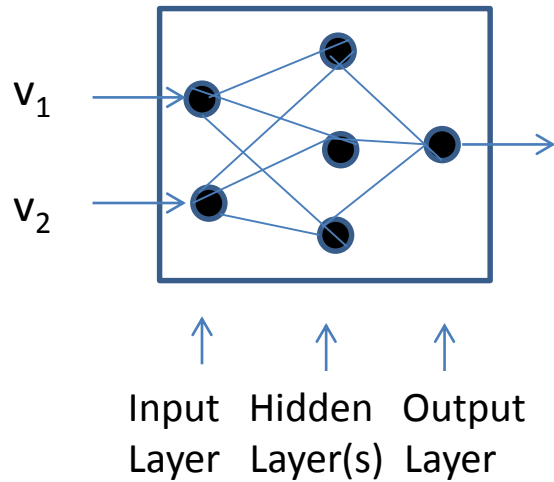
# HOW DOES IT WORK?



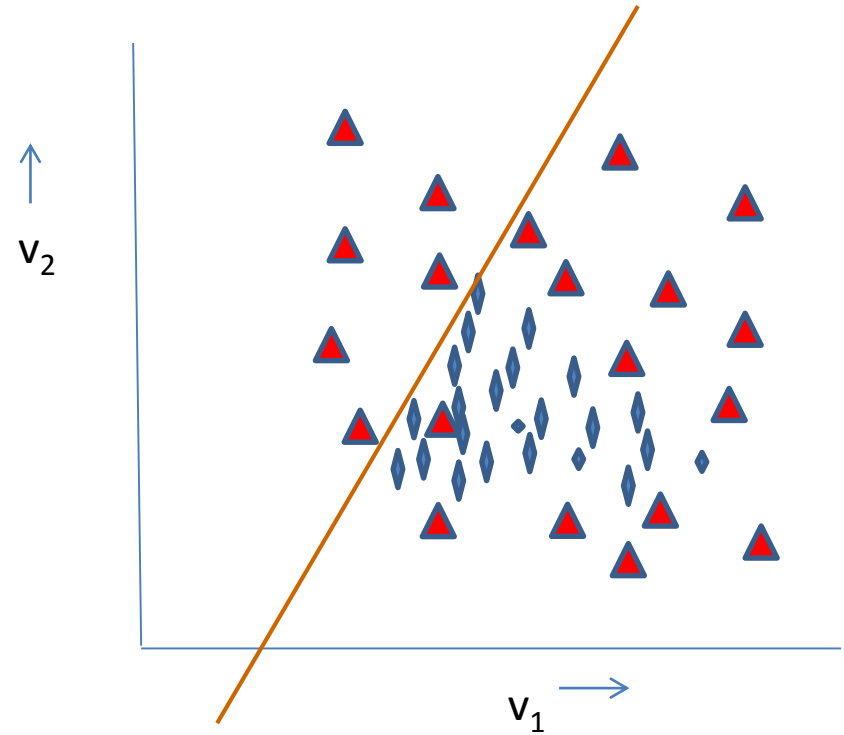
For second hidden node



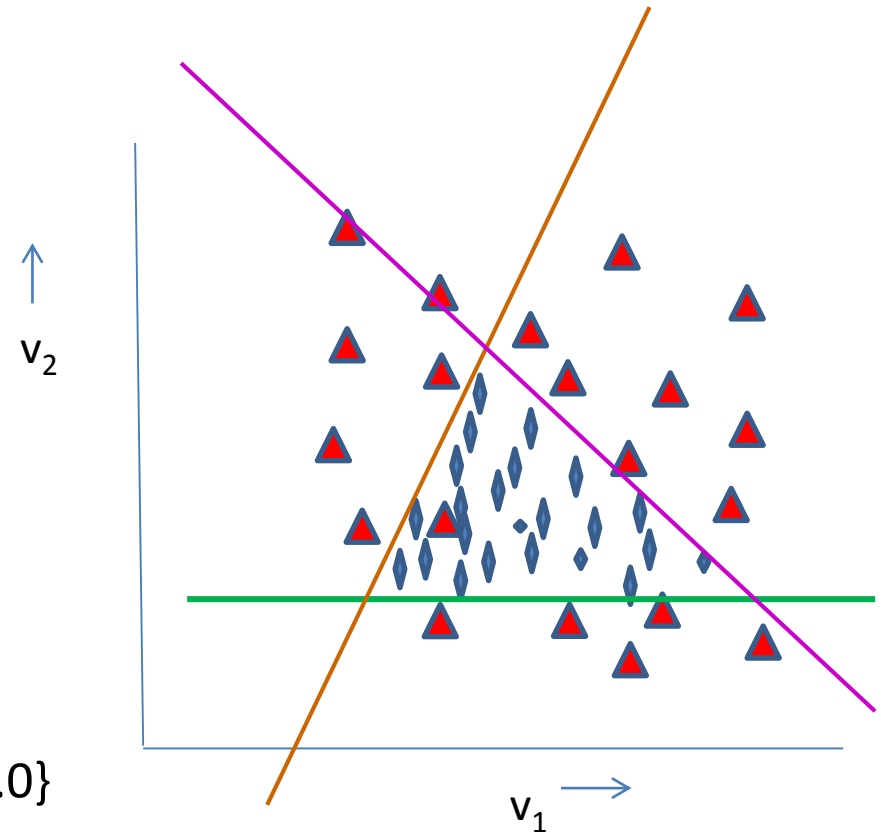
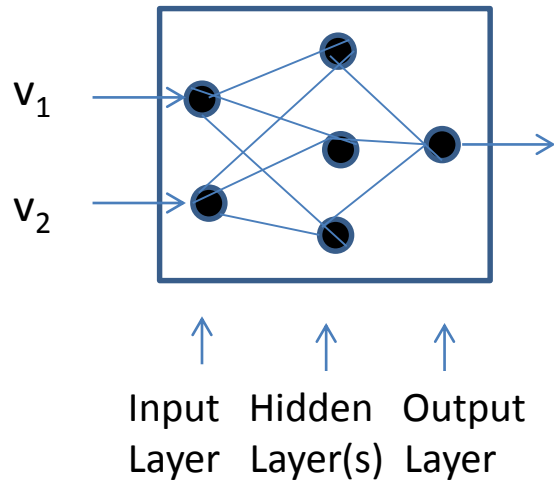
# HOW DOES IT WORK?



For third hidden node



# HOW DOES IT WORK?



Output = Sigmoid $\{0.4H_1 + 0.4H_2 + 0.4H_3 - 1.0\}$   
Output is 'On' only if  $H_1 H_2 H_3$  all are 'On'

- N.B.
- \* Complexity of final region depends on number of hidden nodes.
  - \* Finite  $\beta \rightarrow$  rounded edges for selected region; and contours of constant output in  $(v_1, v_2)$  plane.

# BEWARE

- Training sets are reliable
- Don't train with variable you want to measure
- Data does not extend outside range of training samples (in multi-dimensions)
- Don't overtrain
- Approx equal numbers of signal and bgd

# Is NN better\* than simple cuts?

In principle, NO

Can cut on complicated variable e.g. NN output

In practice: YES

But:

Better NN performance → more work by 'Cuts' analysis to improve performance

\* Better = improved efficiency v mistag rate

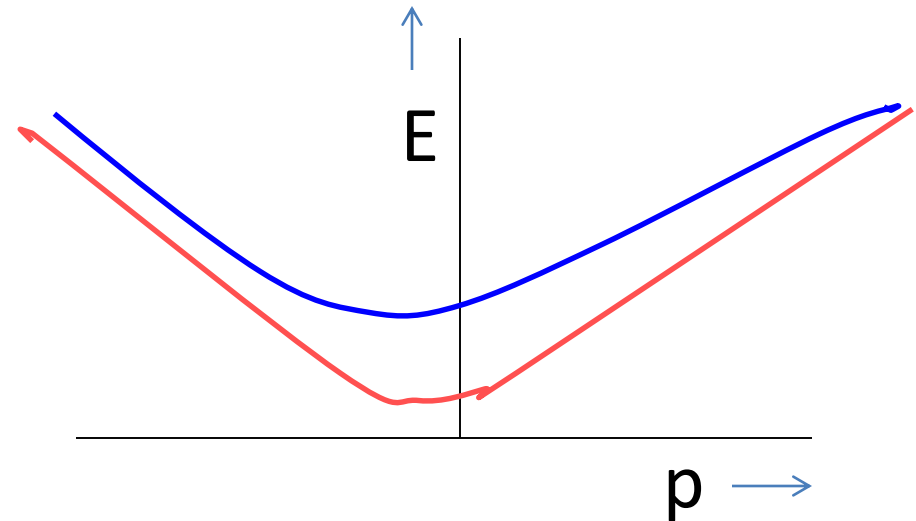


# SIMPLE EXAMPLE

Try to separate  $\pi$  and proton using E and p

$$\pi: E^2 = p^2 + m_\pi^2$$

$$P: E^2 = p^2 + m_p^2$$



Easy:  $p = 0 \rightarrow 2 \text{ GeV}$

Harder:  $p = -4 \rightarrow 4 \text{ GeV}$

Hardest:  $p_x, p_y, p_z = -4 \rightarrow 4 \text{ GeV}$

More realistic: Add expt scatter of data wrt curves

# PHYSICS EXAMPLE

Separate b-jets from light flavour, gluons, W, Z:

Input variables: Track IPs, SV mass, distance, quality, etc.

Output: 0  $\rightarrow$  1

Issues:

Pre NN cuts

Training and testing samples (Where from? How many events? Ratios of different bgds,....)

How many inputs?

Network structure

How many networks?

Single output or several

Systematics (use different sets of testing events}

Stability wrt NN cut

# NN Summary

- **ADVANTAGES:**
  - Very flexible
  - Correlations OK
  - Tunable cut
- **DISADVANTAGES**
  - Training takes time
  - Tendency to include too many variables
  - Treat as black box
- \* Past attitude: Need to convince colleagues NN is sensible  
More recently: Why aren't you using NN?  
Now/future: Why aren't you using a Deep Network?

# Wilks' Theorem

Data = some distribution e.g. mass histogram

For H0 and H1, calculate best fit weighted sum of squares  $S_0$  and  $S_1$

Examples: 1) H0 = polynomial of degree 3

H1 = polynomial of degree 5

2) H0 = background only

H1 = bgd+peak with free  $M_0$  and cross-section

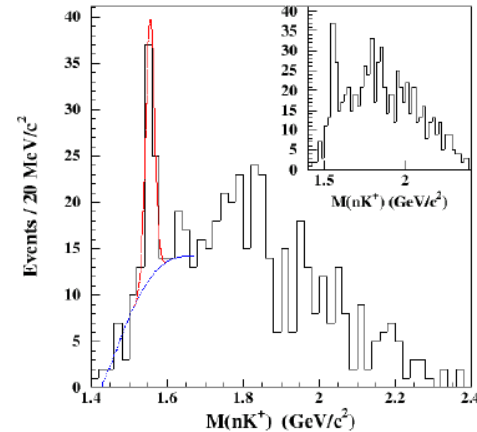
3) H0 = normal neutrino hierarchy

H1 = inverted hierarchy

If H0 true,  $S_0$  distributed as  $\chi^2$  with  $\text{ndf} = \nu_0$

If H1 true,  $S_1$  distributed as  $\chi^2$  with  $\text{ndf} = \nu_1$

If H0 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) H0 is true

b) H0 and H1 are nested

c) Params for H1  $\rightarrow$  H0 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  $\text{ndf} = (d-4) - (d-6) = 2$

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

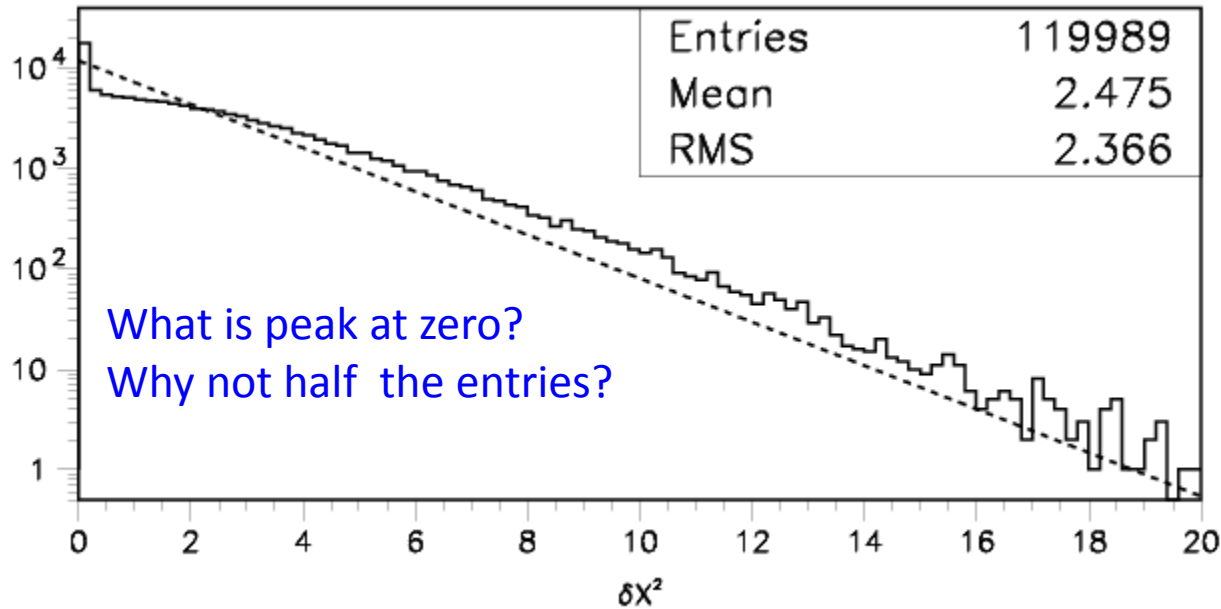
\*\*\*\*\*

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  $\text{ndf}$ , better to use  $\Delta S$ , rather than  $S_1$  and  $S_0$  separately

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

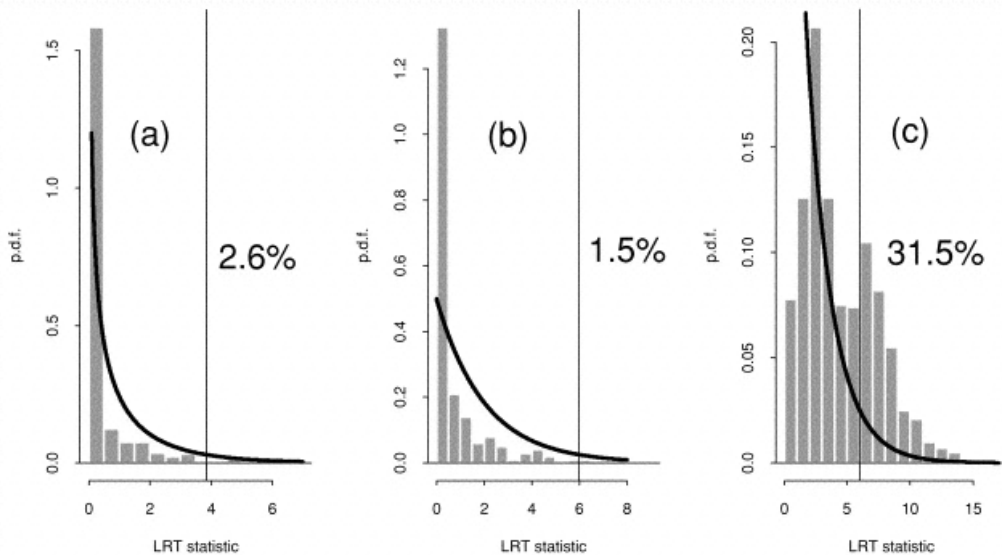


Demortier:

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%

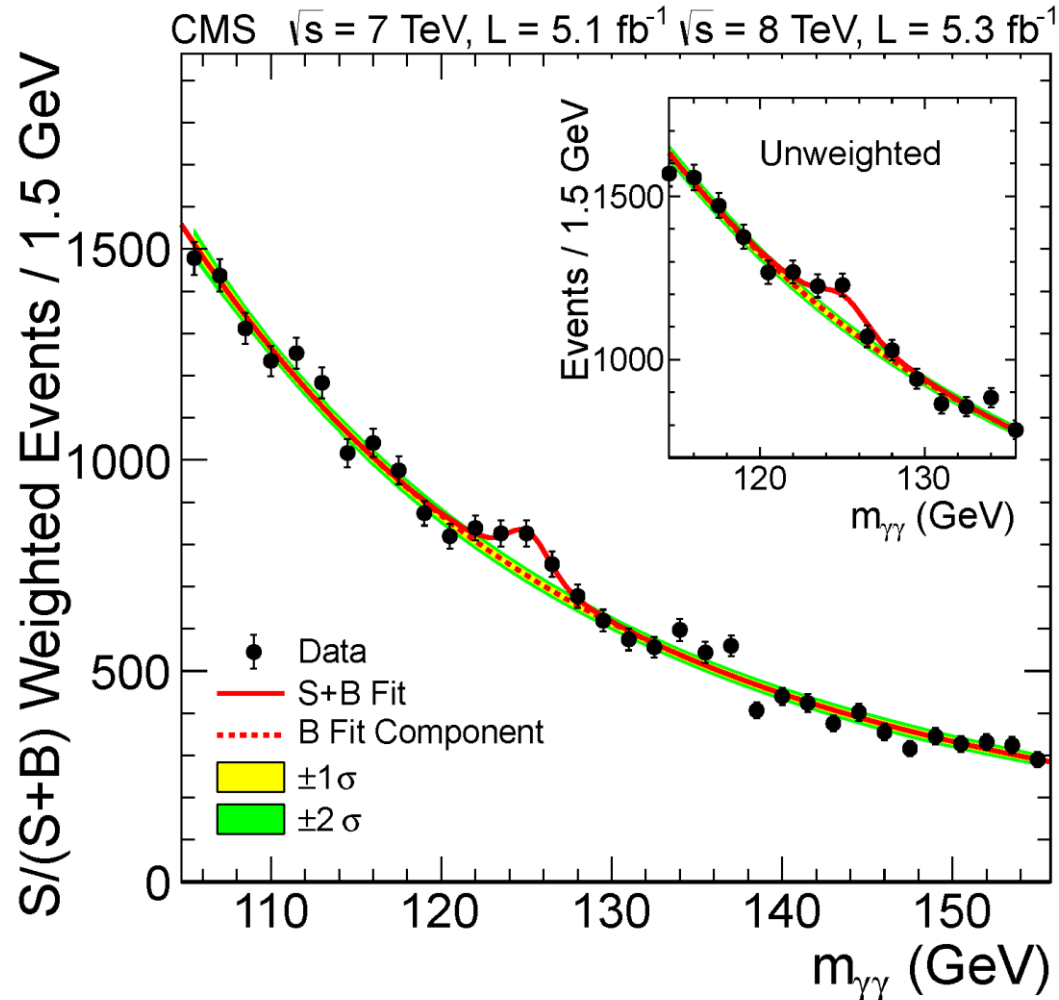
## Is difference in $S$ distributed as $\chi^2$ ?, contd.

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

N.B.

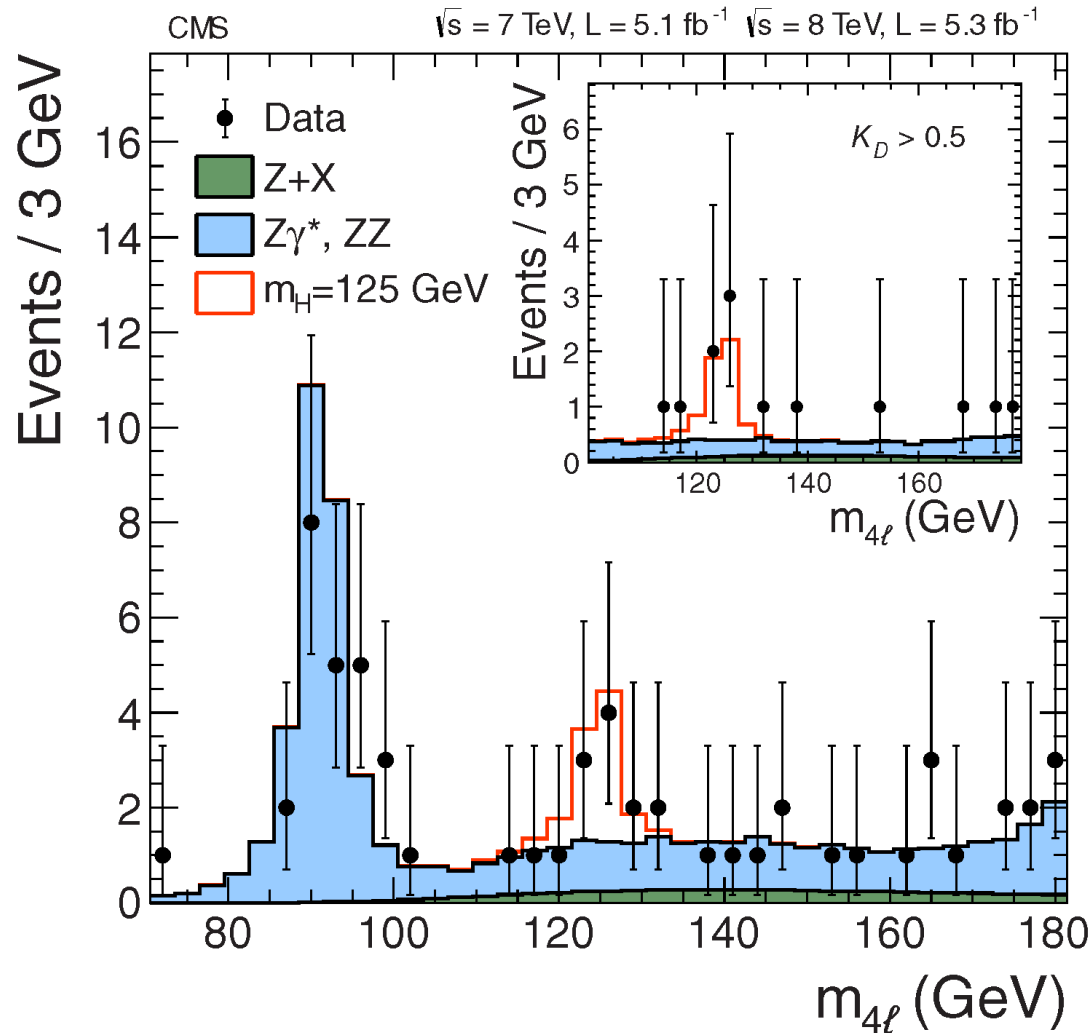
- 1) For mass spectrum, determining  $\Delta S$  for hypothesis  $H_1$  when data is generated according to  $H_0$  is not trivial, because there will be lots of local minima
- 2) If we are interested in  $5\sigma$  significance level, needs lots of MC simulations (or intelligent MC generation)
- 3) 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>)

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

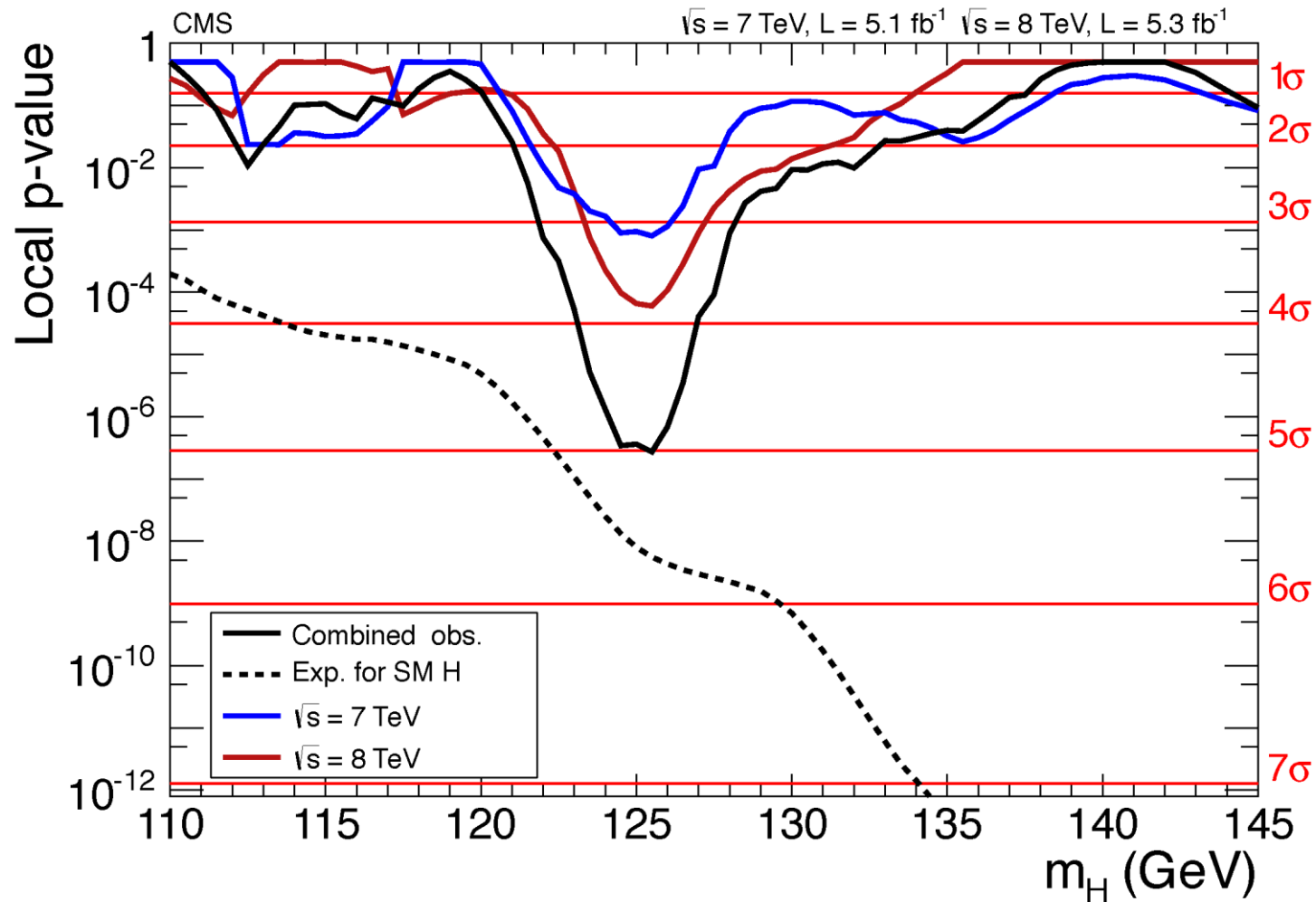


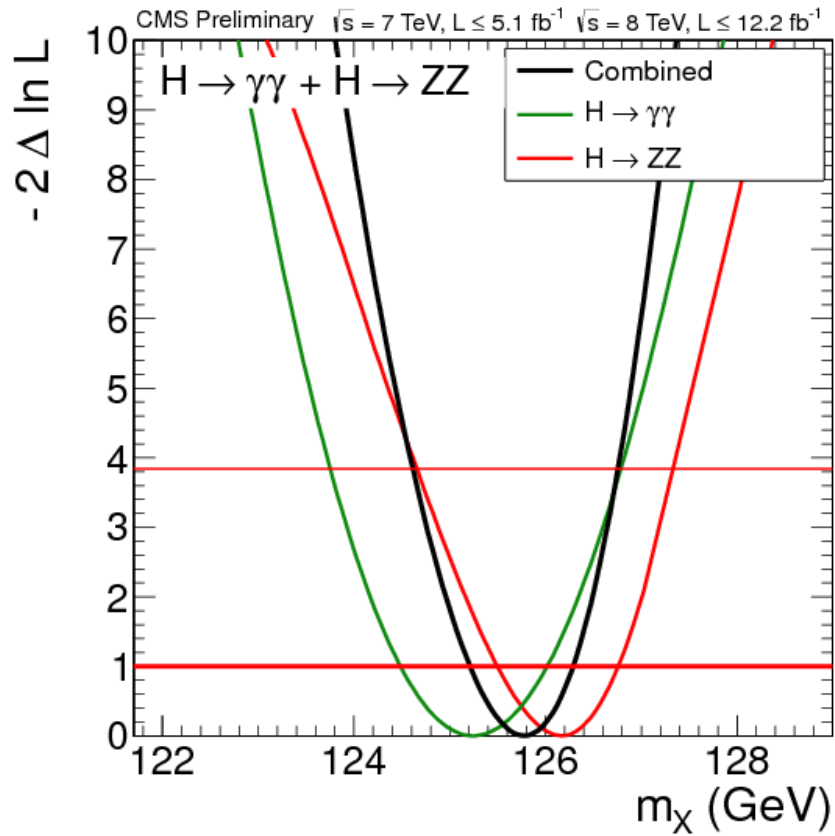


# $H \rightarrow Z Z \rightarrow 4 \ell$ : high S/B, low statistics

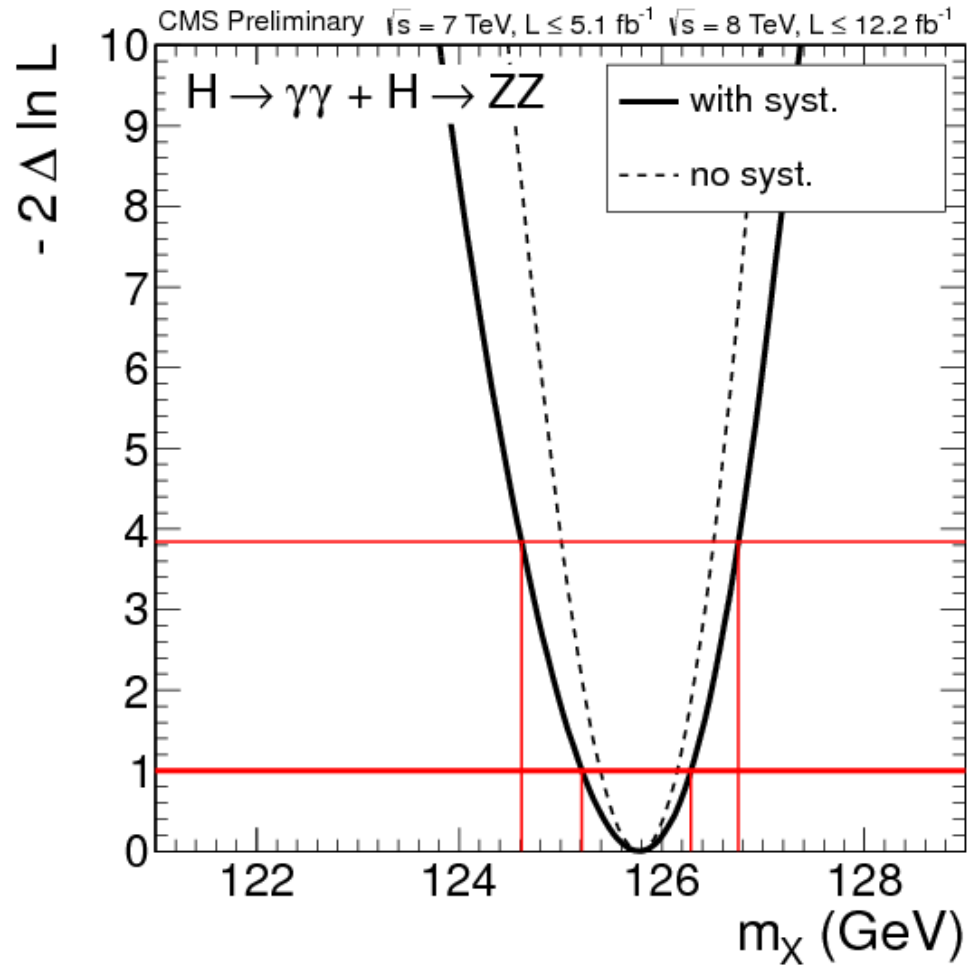


# p-value for 'No Higgs' versus $m_H$

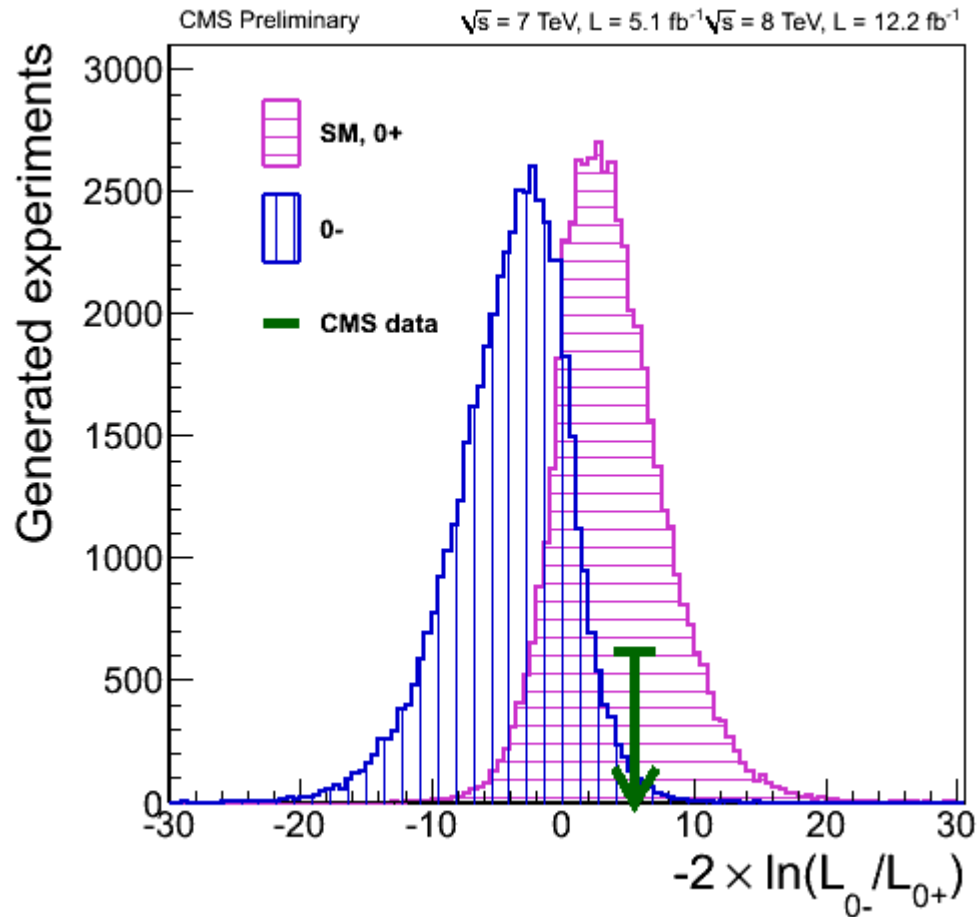




## Mass of Higgs: Likelihood versus mass



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



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