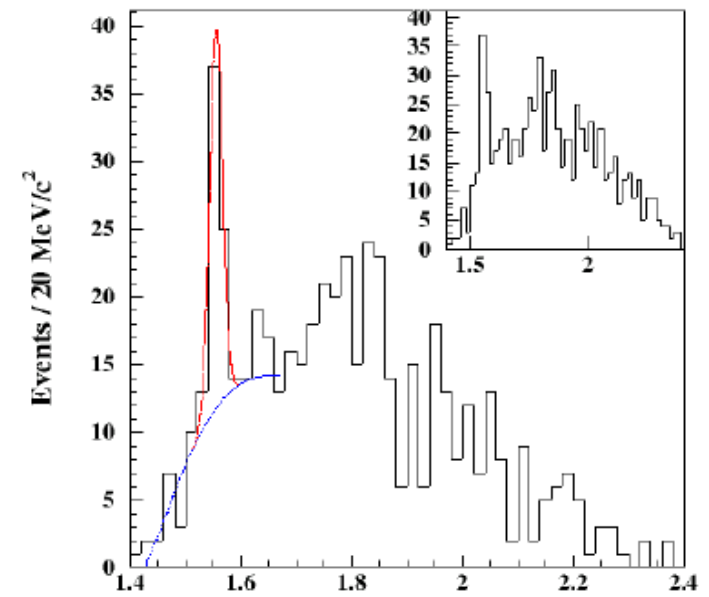
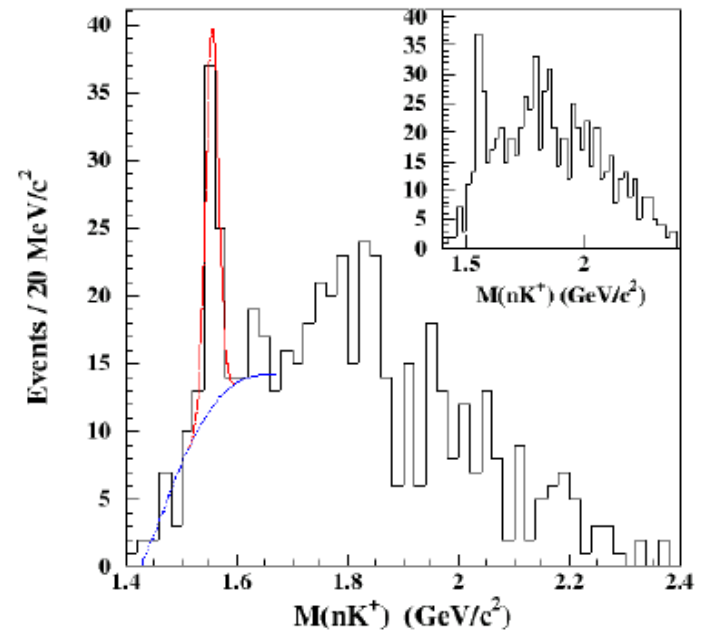


Is there evidence for a peak in this data?



Is there evidence for a peak in this data?



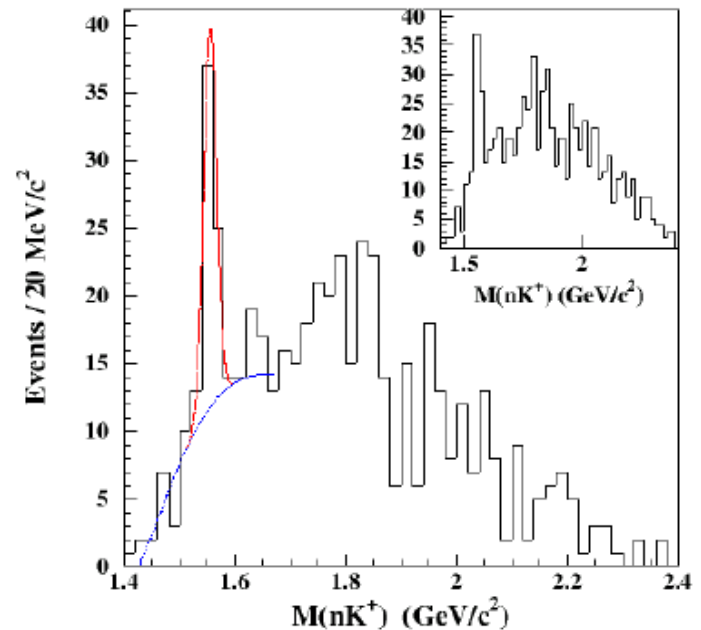
“Observation of an Exotic $S=+1$

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 Phenomena: p-values, Upper Limits and Discovery

Louis Lyons

IC and Oxford

l.lyons@physics.ox.ac.uk

CERN Summer Students,

July 2015

PHYSTAT 2011 Workshop at CERN, Geneva

17-20 January 2011

<http://indico.cern.ch/event/phystat2011>

Contacts:

Albert de Roeck <albert.de.roeck@cern.ch>

Louis Lyons <l.lyons1@physics.ox.ac.uk>

LPCC

LHC Physics Centre at CERN



Jan 17-19: Statistical Issues for Search Experiments

Statistical issues related to discovery claims in search experiments, with emphasis on those at the LHC.

Jan 20: Unfolding

Unfolding of detector effects from experimental distributions

Programme committee

J. Berger (Duke)
V. Blobel (Hamburg)
R. Cousins (UCLA)
D. Cox (Oxford)
G. Cowan (Royal Holloway)
K. Cranmer (NYU)
L. Demortier (Rockefeller)
A. de Roeck (CERN)
B. Efron (Stanford)
G. Flucke (DESY)
E. Gross (Weizmann)
D. Hand (Imperial College)
J. Linnemann (MSU)
R. Lockhart (Simon Fraser)
L. Lyons (Imperial College)
M.L. Mangano (CERN)
S. Schmitt (DESY)
M. Williams (Imperial College)

TOPICS

Discoveries

H_0 or H_0 v H_1

p-values: For Gaussian, Poisson and multi-variate data

What is p good for?

Errors of 1st and 2nd kind

What a p-value is not

Combining p-values

Significance

Look Elsewhere Effect

Blind Analysis

Why 5σ ?

Setting Limits

Case study: Search for Higgs boson

DISCOVERIES

“Recent” history:

Charm	SLAC, BNL	1974
Tau lepton	SLAC	1977
Bottom	FNAL	1977
W, Z	CERN	1983
Top	FNAL	1995
{Pentaquarks	~Everywhere	2002 }
Higgs	CERN	2012
?	CERN	2015?

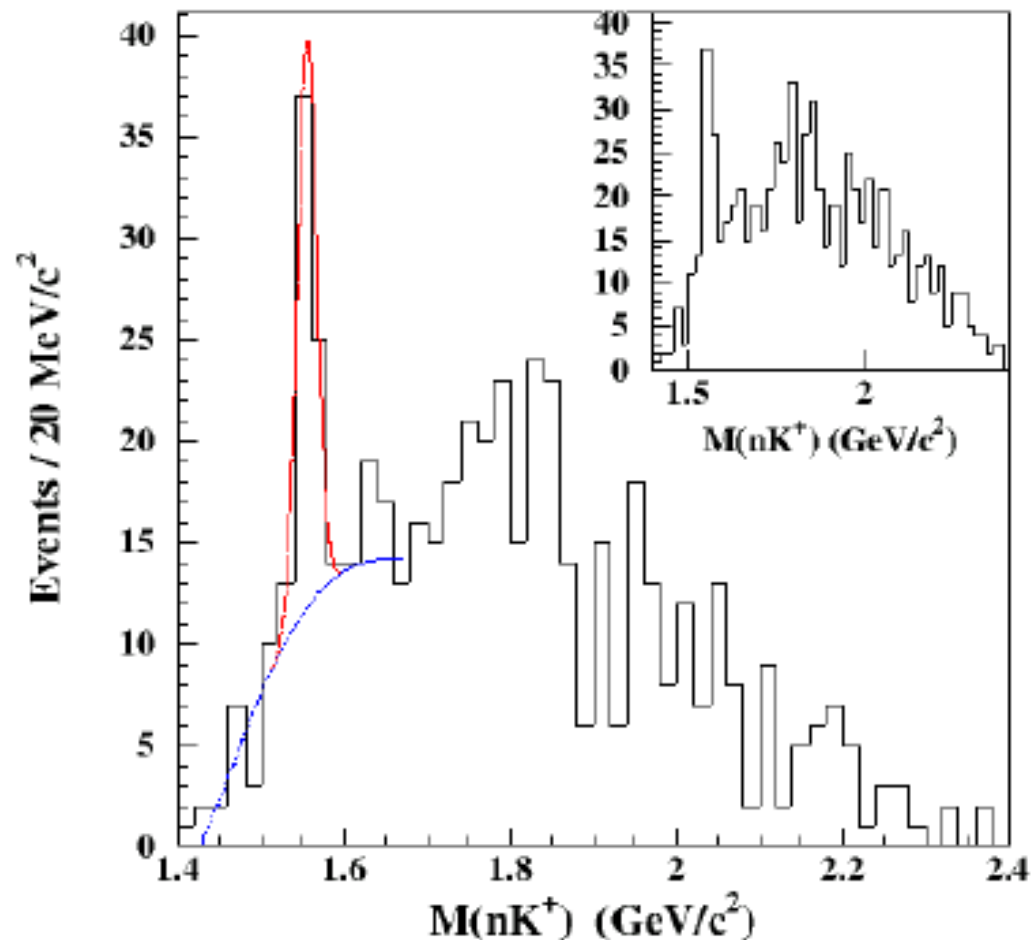
? = SUSY, q and l substructure, extra dimensions,

free q/monopoles, technicolour, 4th generation, black holes,.....

QUESTION: How to distinguish discoveries from fluctuations?

Penta-quarks?

Hypothesis testing: New particle or statistical fluctuation?



H0 or H0 versus H1 ?

H0 = null hypothesis

e.g. Standard Model, with nothing new

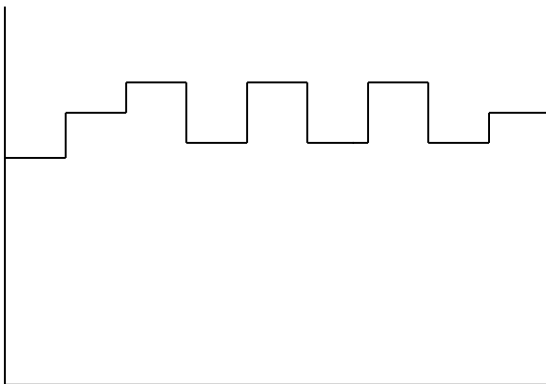
H1 = specific New Physics e.g. Higgs with $M_H = 125$ GeV

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

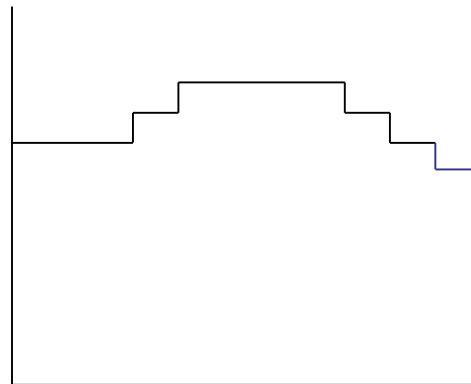
H0 v H1: “Hypothesis Testing” e.g. \mathcal{L} -ratio

Measures how much data favours one hypothesis wrt other

H0 v H1 likely to be more sensitive for H1



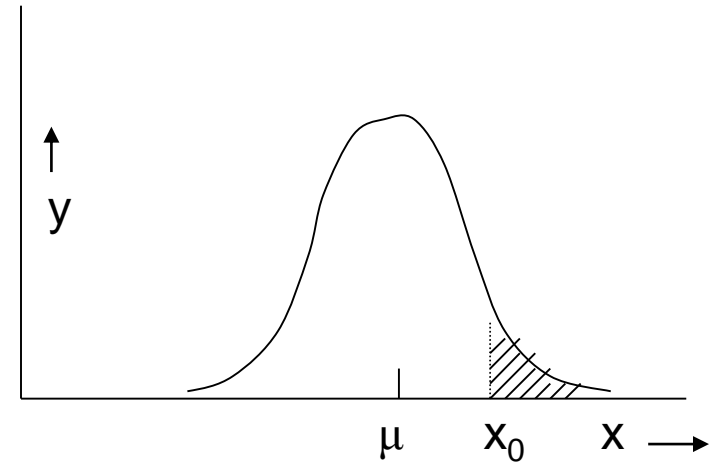
or



p-values

Concept of pdf

Example: Gaussian



y = probability density for measurement x

$$y = \frac{1}{\sqrt{(2\pi)\sigma^2}} \exp\{-0.5*(x-\mu)^2/\sigma^2\}$$

p-value: probability that $x \geq x_0$

Gives probability of “extreme” values of data (in interesting direction)

$(x_0 - \mu)/\sigma$	1	2	3	4	5
p	16%	2.3%	0.13%	0.003%	0.3×10^{-6}

i.e. **Small p = unexpected**

p-values, contd

Assumes:

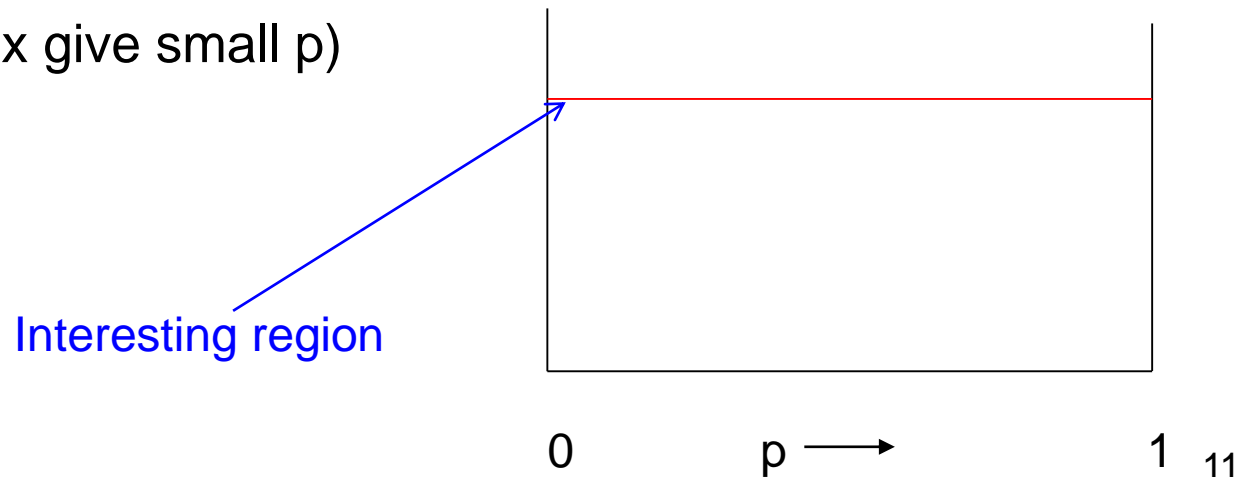
Specific pdf for x (e.g. Gaussian, no long tails)

Data is unbiased

σ is correct

If so, and x is from that pdf \Rightarrow **uniform p-distribution**

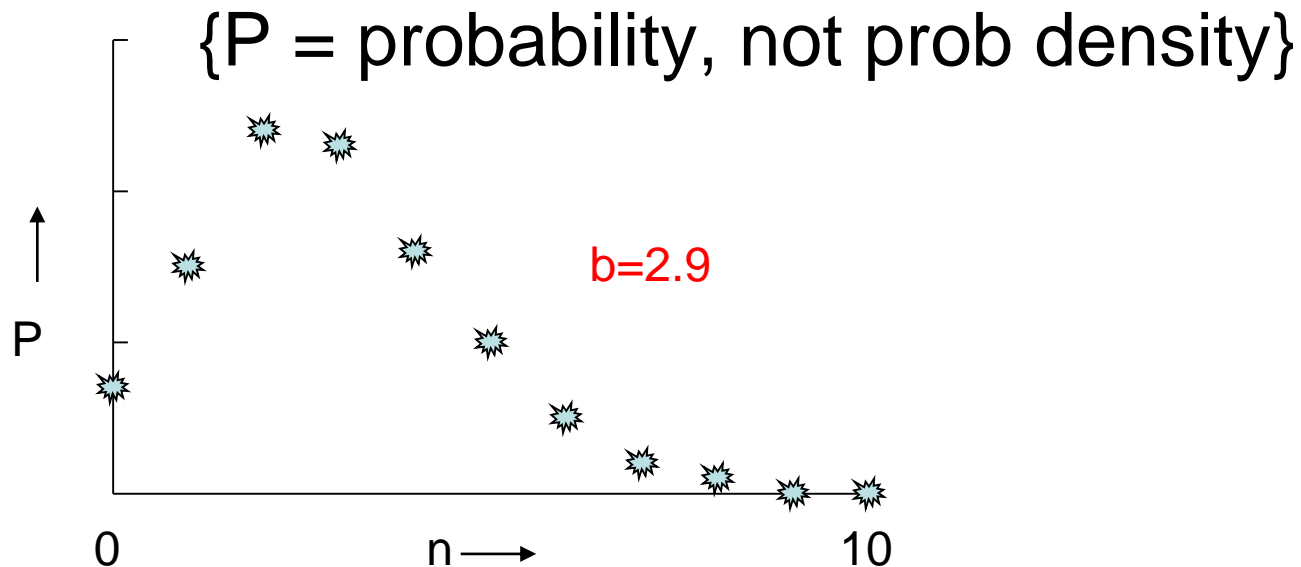
(Events at large x give small p)



p-values for non-Gaussian distributions

e.g. **Poisson** counting experiment, $\text{bgd} = b$

$$P(n) = e^{-b} * b^n / n!$$



For $n=7$, $p = \text{Prob}(\text{ at least 7 events}) = P(7) + P(8) + P(9) + \dots = 0.03$

p-values and σ

p-values often converted into equivalent Gaussian σ

e.g. 3×10^{-7} is “ 5σ ” (one-sided Gaussian tail)

Does NOT imply that pdf = Gaussian

(Simply easier to remember number of σ , than p-value.)

What is p good for?

Used to test whether data is consistent with H_0

Reject H_0 if p is small : $p \leq \alpha$ (How small?)

Sometimes make wrong decision:

Reject H_0 when H_0 is true: **Error of 1st kind**

Should happen at rate α

OR

Fail to reject H_0 when something else

(H_1, H_2, \dots) is true: **Error of 2nd kind**

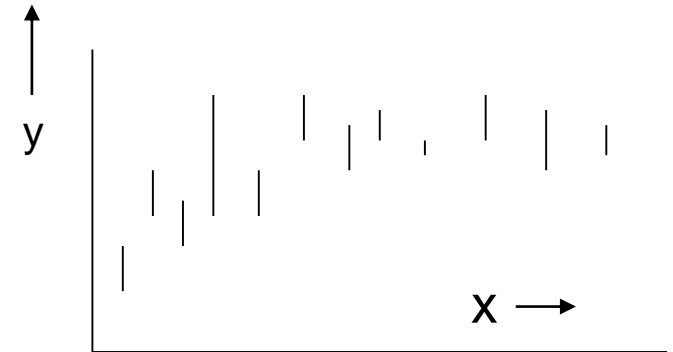
Rate at which this happens depends on.....

Errors of 2nd kind: How often?

e.g.1. Does data lie on straight line?

Calculate S_{\min}

Reject if $S_{\min} \geq 20$



Error of 1st kind: $S_{\min} \geq 20$ Reject H_0 when true

Error of 2nd kind: $S_{\min} < 20$ Accept H_0 when in fact quadratic or..

How often depends on:

- Size of quadratic term

- Magnitude of errors on data, spread in x-values,

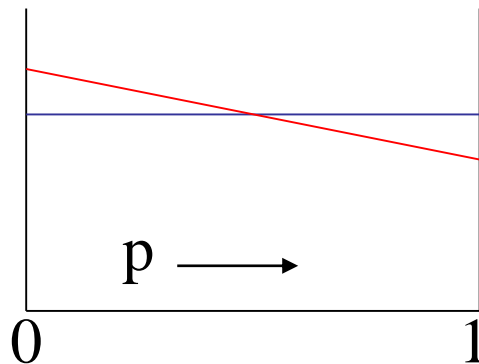
- How frequently quadratic term is present

Errors of 2nd kind: How often?

e.g. 2. Particle identification (TOF, dE/dx , Čerenkov,.....)

Particles are π or μ

Extract p-value for $H_0 = \pi$ from PID information



π and μ have similar masses

Of particles that have $p \sim 1\%$ ('reject H_0 '), fraction that are π is

- a) \sim half, for equal mixture of π and μ
- b) almost all, for "pure" π beam
- c) very few, for "pure" μ beam

p-value is not

Does **NOT** measure $\text{Prob}(H_0 \text{ is true})$

i.e. It is **NOT** $P(H_0|\text{data})$

It is $P(\text{data}|H_0)$

N.B. $P(H_0|\text{data}) \neq P(\text{data}|H_0)$

$P(\text{theory}|\text{data}) \neq P(\text{data}|\text{theory})$

“Of all results with $p \leq 5\%$, half will turn out to be wrong”

N.B. Nothing wrong with this statement

e.g. 1000 tests of energy conservation

~50 should have $p \leq 5\%$, and so reject H_0 = energy conservation

Of these 50 results, all are likely to be “wrong”

Combining different p-values

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

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 ~ 0.001 and is clearly **NOT** correct combined p)

N.B. Problem does not have unique answer

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

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

Significance

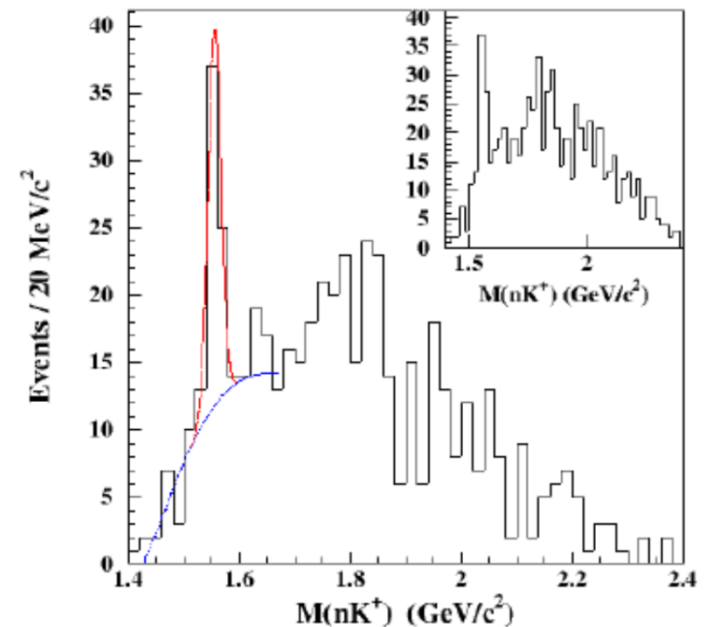
$$\text{Significance} = S/\sqrt{B} ? \text{ \{or } S/\sqrt{(B+S)}?\}$$

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 (Blind analyses)
- Choice of bins (.....)

For future experiments:

- Optimising cuts: Could give $S = 0.1$, $B = 10^{-4}$, $S/\sqrt{B} = 10$
- N.B. $S/\sqrt{(S+B)}$ also has problems
- Best to use proper Poisson p



Look Elsewhere Effect

See 'peak' in bin of histogram

Assuming null hypothesis, p-value is chance of fluctuation at least as significant as observed

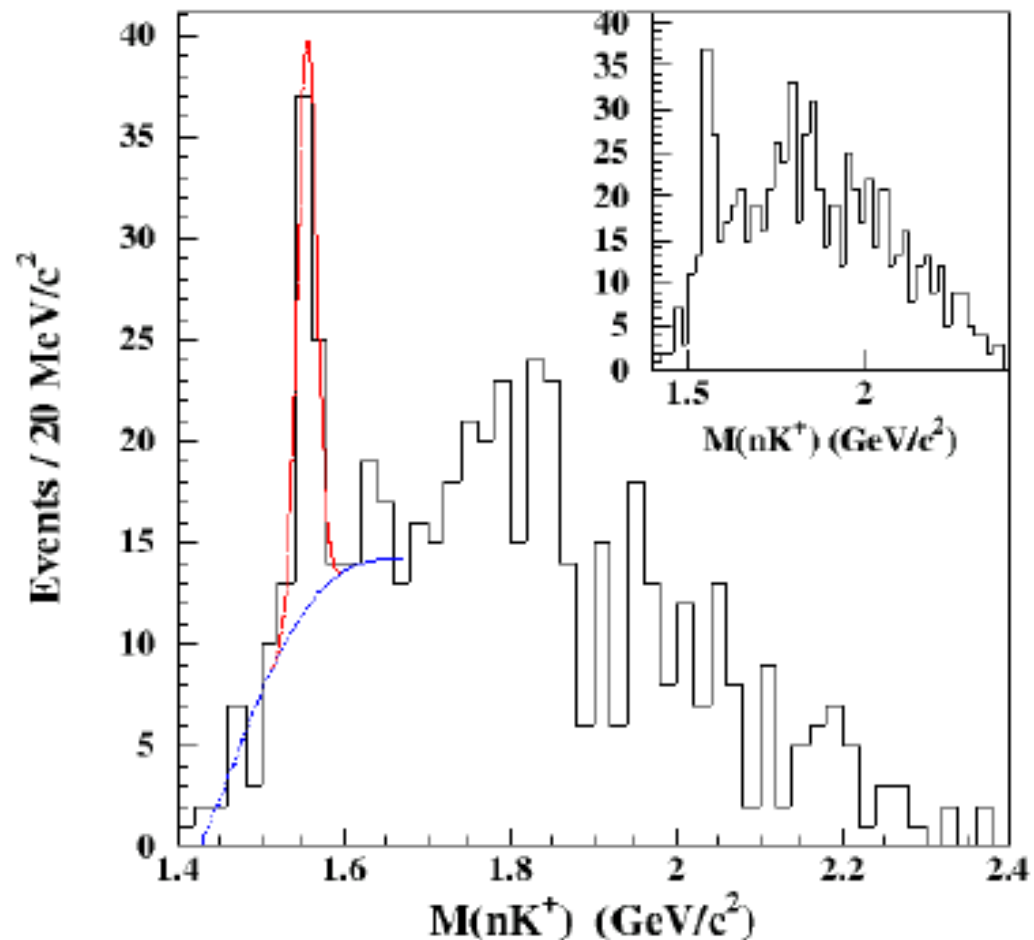
- 1) at the position observed in the data; or
 - 2) anywhere in that histogram; or
 - 3) including other relevant histograms for your analysis; or
 - 4) including other analyses in Collaboration; or
 - 5) in any CERN experiment; or
- etc.

Contrast **local p-value** with 'global' p-value

Specify what is your 'global'

Penta-quarks?

Hypothesis testing: New particle or statistical fluctuation?



Example of LEE: Stonehenge



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-²⁴ century rubbish dump.

BLIND ANALYSES

Why blind analysis?

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

After analysis is unblinded,

* Luis Alvarez suggestion re “discovery” of free quarks

Why 5σ ?

- Past experience with 3σ , 4σ ,... signals
- Look elsewhere effect:

Different cuts to produce data

Different bins (and binning) of this histogram

Different distributions Collaboration did/could look at

Other analyses in Physics subgroup, expt, CERN,...

- Worries about systematics (easily under-estimated?)
- Bayesian priors:

$$\frac{\text{Bayes posteriors}}{\text{Bayes posteriors}} = \frac{\text{Likelihoods} * \text{Priors}}{\text{Likelihoods} * \text{Priors}}$$

↑ ↑ ↑
Bayes posteriors Likelihoods Priors

Prior for $\{H_0 = \text{S.M.}\} \gg \gg$ Prior for $\{H_1 = \text{New Physics}\}$ 26

Why 5σ ?

BEWARE of tails,
especially for nuisance parameters

Same criterion for all searches

Different LEE (contrast muon magnetic moment v. CMS)

Different role of systematics

Different Bayes priors, e.g.

- Single top production

- Higgs

- Highly speculative particle

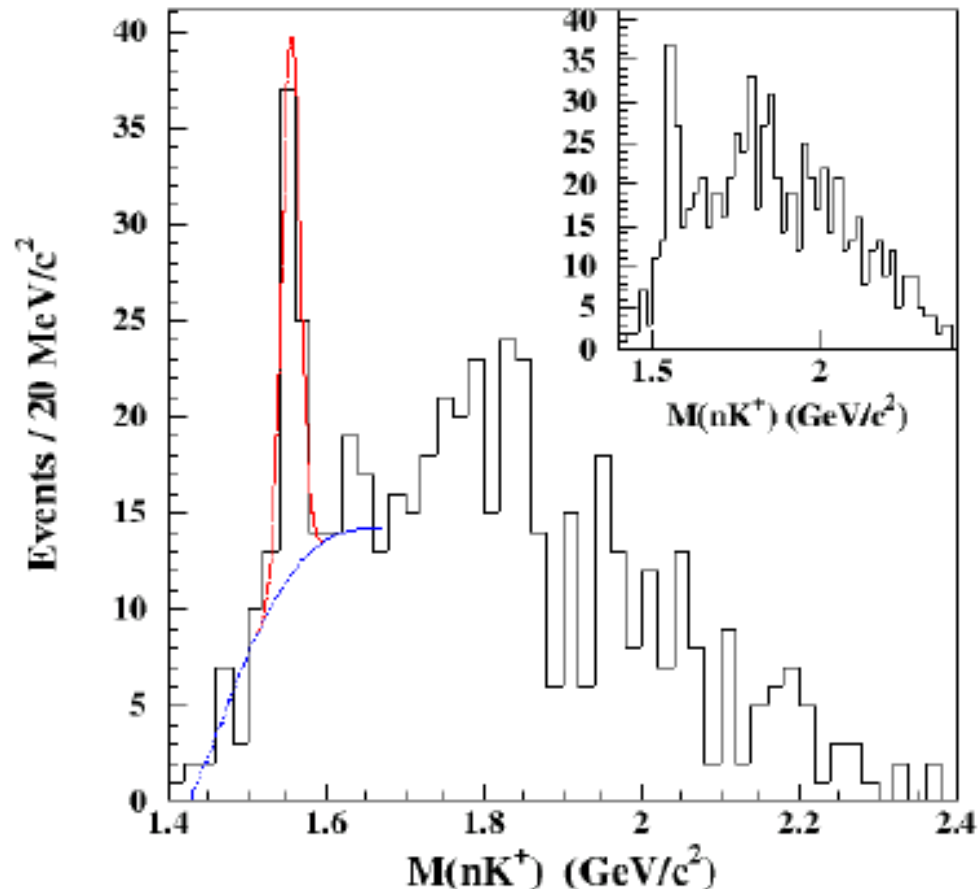
- Energy non-conservation

Blind analysis helps

Choosing between 2 hypotheses

Hypothesis testing: New particle or statistical fluctuation?

$$H_0 = b \quad H_1 = b + s$$



Choosing between 2 hypotheses

Possible methods:

$\Delta\chi^2$

p-value of statistic →

$\ln\mathcal{L}$ -ratio

Bayesian:

Posterior odds

Bayes factor

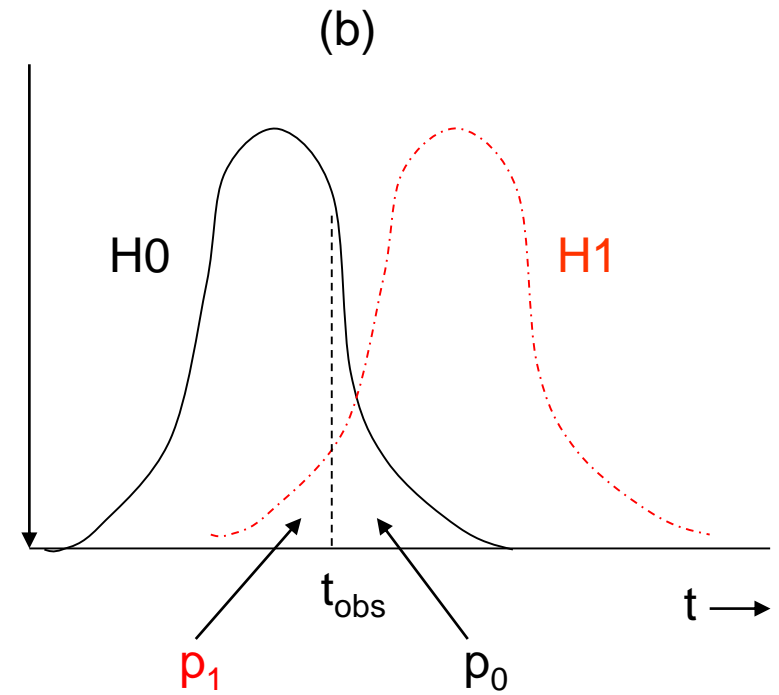
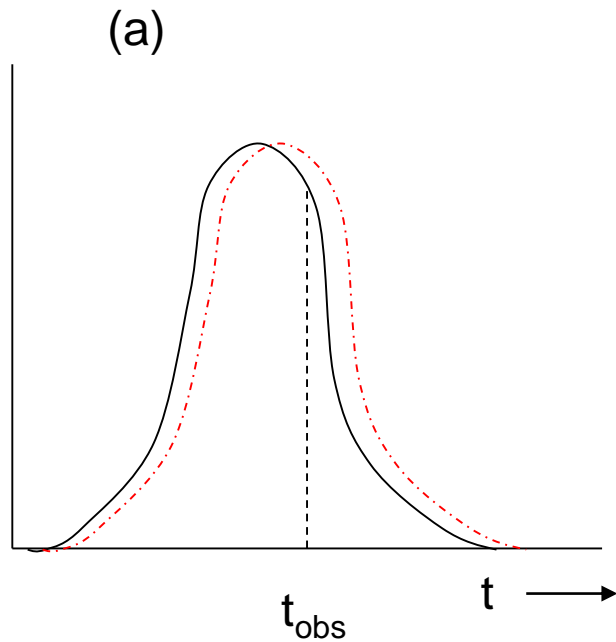
Bayes information criterion (BIC)

Akaike (AIC)

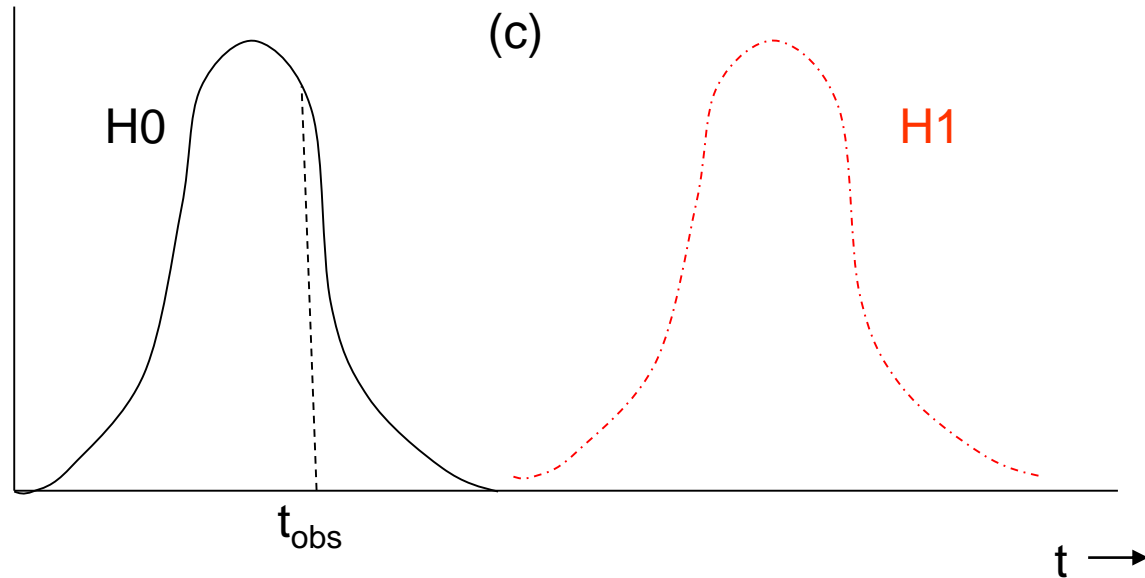
Minimise “cost”

See ‘Comparing two hypotheses’

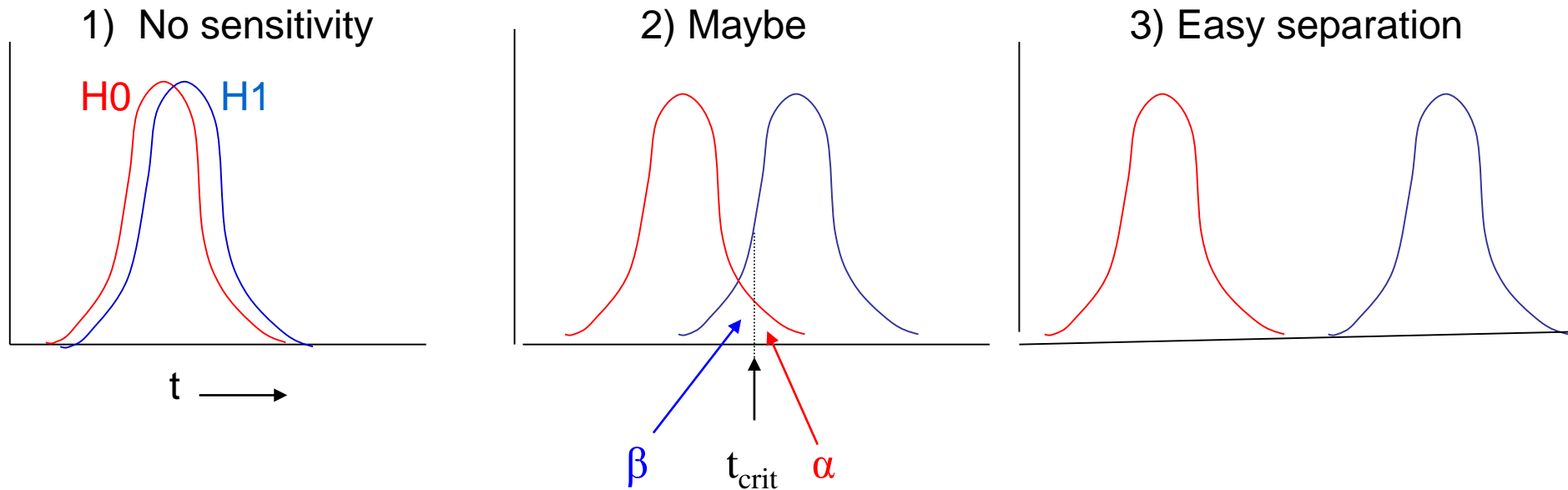
http://www.physics.ox.ac.uk/users/lyons/H0H1_A~1.pdf



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: Choose α (e.g. 95%, 3σ , 5σ ?) and CL for β (e.g. 95%)

Given b , α determines t_{crit}

s defines β . 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 α**

If $t_{\text{obs}} < t_{\text{crit}}$, no discovery. If $p_1 < 1 - \text{CL}$, **exclude H_1**

LIMITS

Look for New Physics s

See no effect. Set upper limit on s

If $s < s_{\text{expected}}$, exclude this sort of New Physics

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

Big industry in Particle Physics

Michelson-Morley experiment \rightarrow death of aether

CERN CLW (Jan 2000)

FNAL CLW (March 2000)

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

SIMPLE PROBLEM?

Gaussian

$\sim \exp\{-0.5*(x-\mu)^2/\sigma^2\}$, with data x_0

No restriction on param of interest μ ; σ known exactly

$$\mu \leq x_0 + k \sigma$$

BUT Poisson $\{\mu = s\varepsilon + b\}$

$$s \geq 0$$

ε and b with uncertainties

Not like : $2 + 3 = ?$

N.B. Actual limit from experiment \neq Expected (median) limit

Methods

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

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

Likelihood (DON'T integrate your L)

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

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

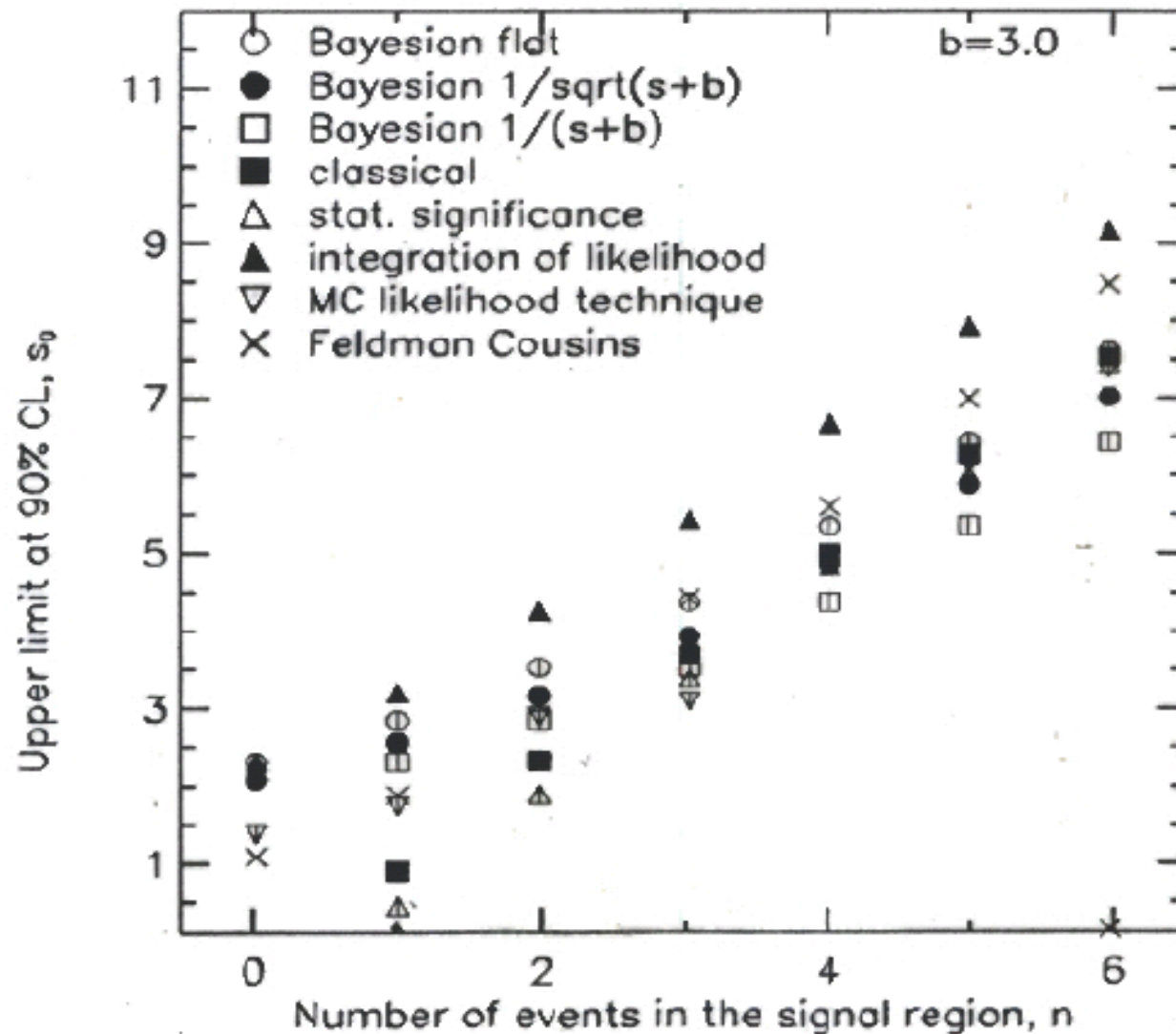
Also have to incorporate systematics

Recommendation 7 from CERN CLW (2000): “Show your L”

1) Not always practical

2) Not sufficient for frequentist methods

Poisson counting expt
Expected bgd = b
Observe n
Set UL for possible signal s



DESIRABLE PROPERTIES

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

INTERVAL LENGTH

Empty \rightarrow Unhappy physicists

Very short \rightarrow False impression of sensitivity

Too long \rightarrow loss of power

(2-sided intervals are more complicated
because 'shorter' is not metric-independent:

e.g. $0 \rightarrow 9$ or $4 \rightarrow 16$ for x^2

cf $0 \rightarrow 3$ or $2 \rightarrow 4$ for x)

Recommendations?

CDF note 7739 (May 2005)

Decide method and procedure in advance

No valid method is ruled out

Bayes is simplest for incorporating nuisance params

Check robustness

Quote coverage

Quote sensitivity

Use same method as other similar expts

Explain method used

Case study: Successful search for Higgs boson

(Meeting of statisticians, atomic physicists, astrophysicists and particle physicist:

“What is value of H_0 ?”))

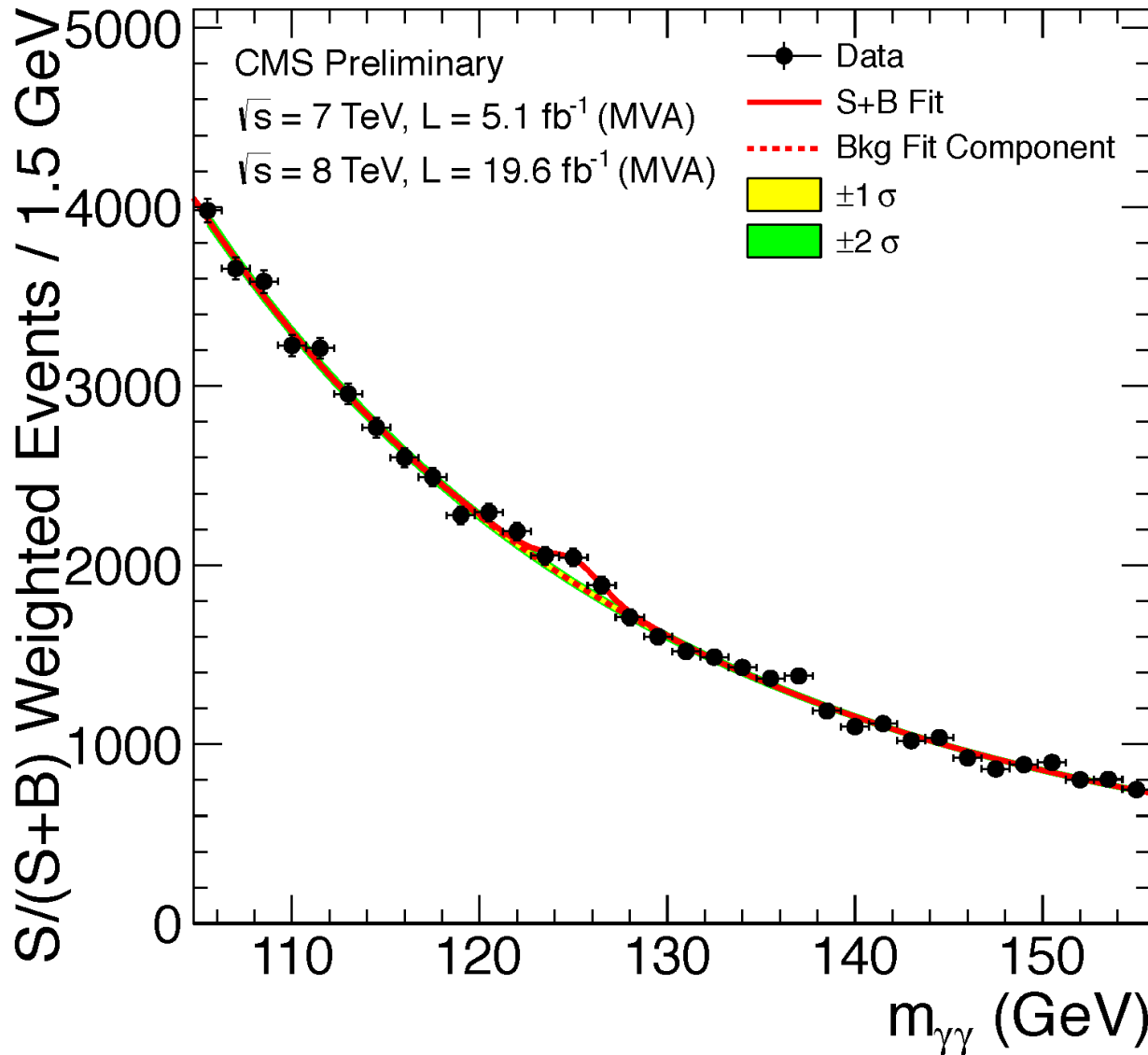
H^0 very fundamental

Wanted to discover Higgs,

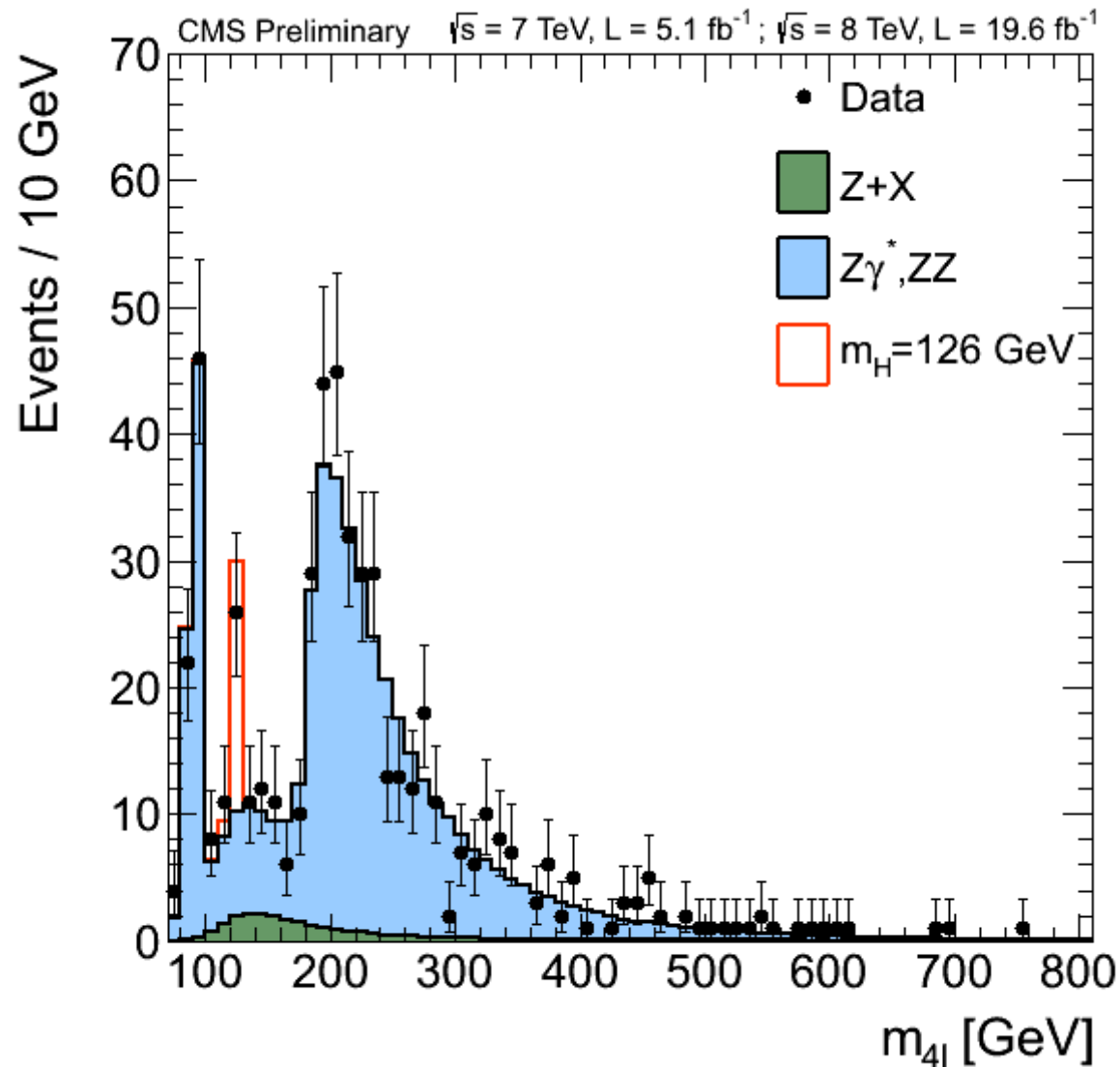
but otherwise exclude

{Other possibility is ‘Not enough data to distinguish’}

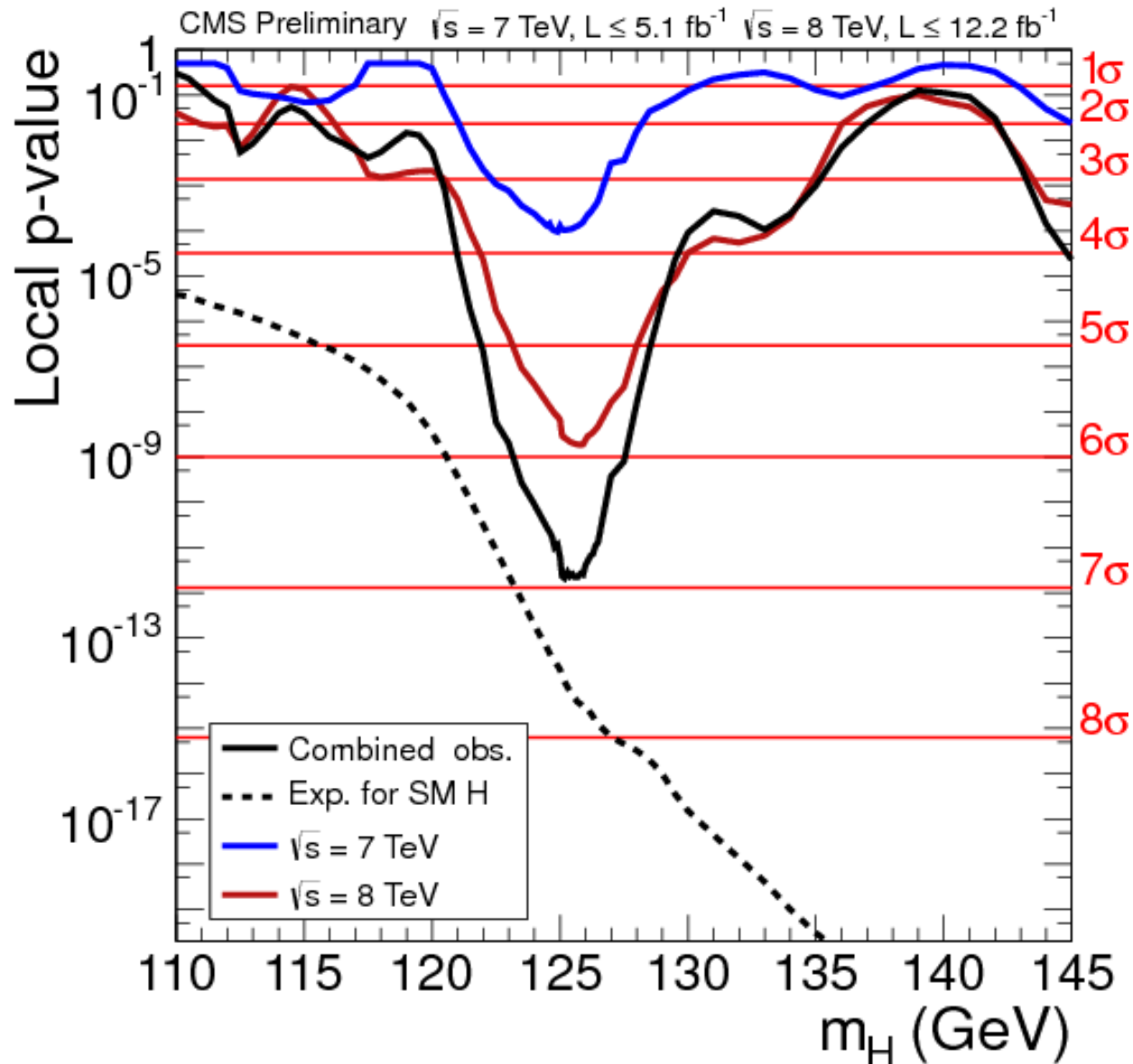
$H \rightarrow \gamma \gamma$: low S/B, high statistics

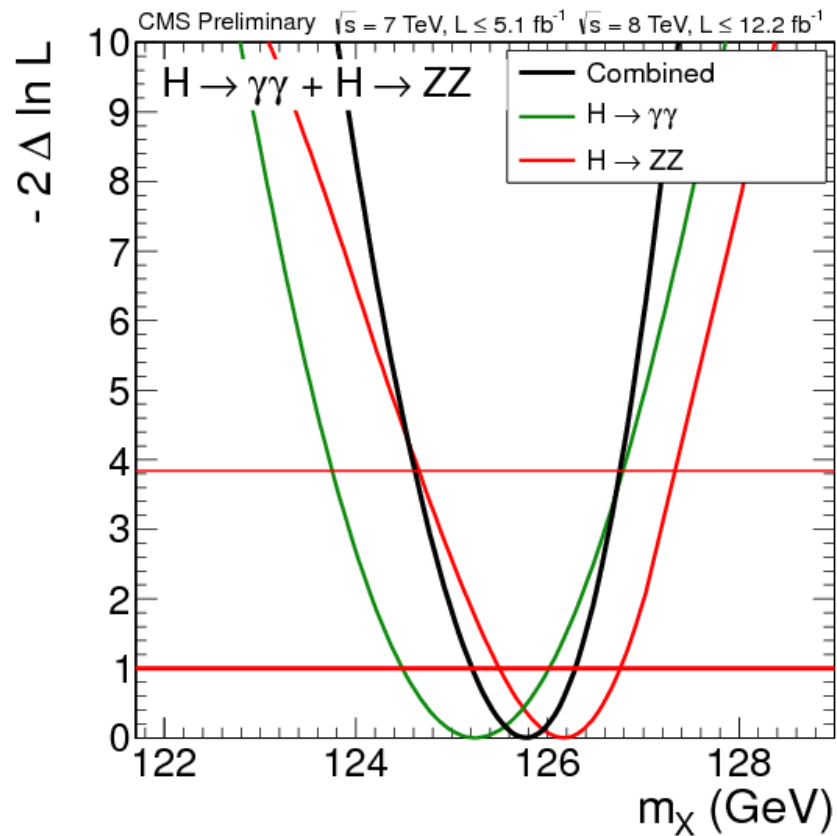


$H \rightarrow Z Z \rightarrow 4 \text{ leptons}$: high S/B, low statistics

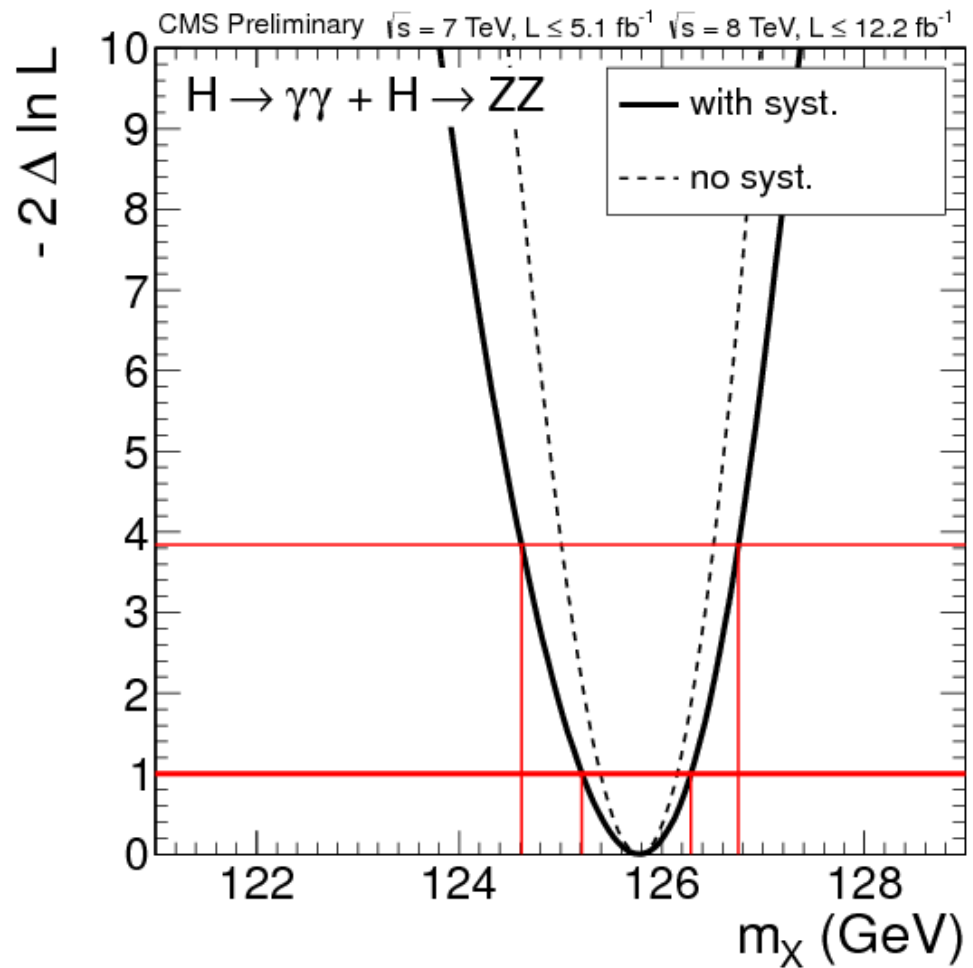


p-value for 'No Higgs' versus m_H

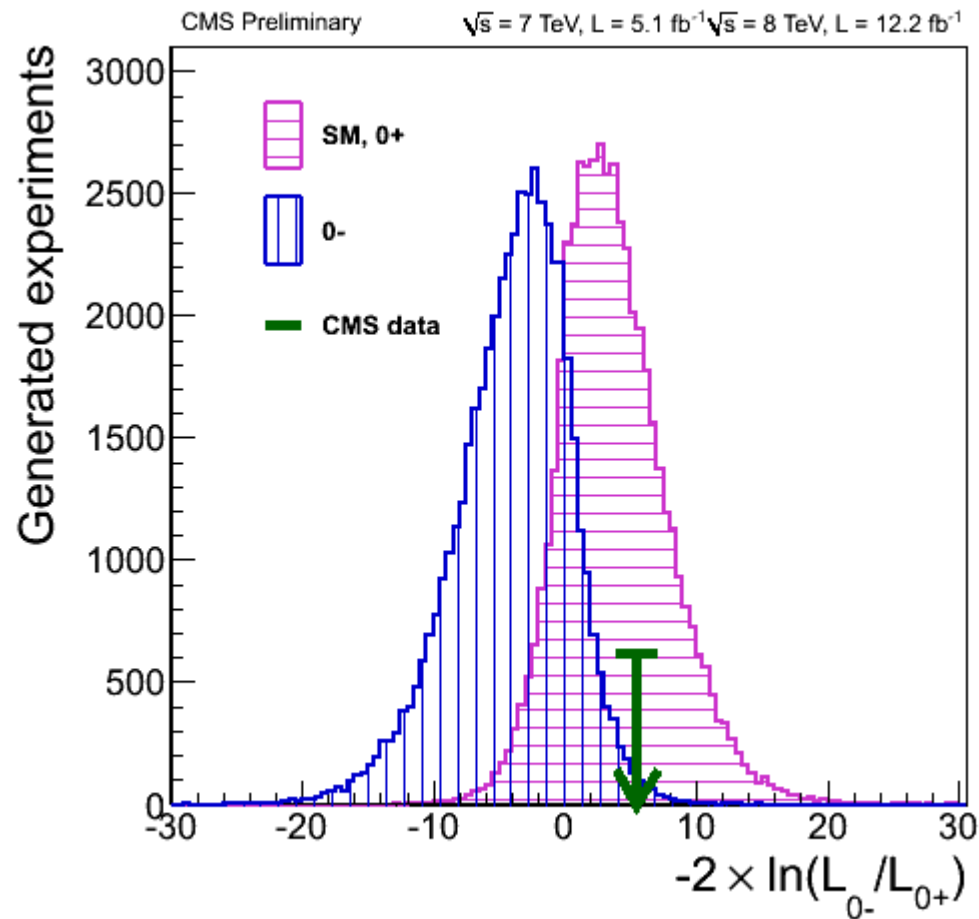




Likelihood versus mass



Comparing 0^+ versus 0^- for Higgs



Summary

- $P(H_0|\text{data}) \neq P(\text{data}|H_0)$
- p-value is NOT probability of hypothesis, given data
- Many different Goodness of Fit tests
 - Most need MC for statistic \rightarrow p-value
- For comparing hypotheses, $\Delta\chi^2$ is better than χ^2_1 and χ^2_2
- Blind analysis avoids personal choice issues
- Different definitions of sensitivity
- Worry about systematics
- H_0 search provides practical example

PHYSTAT2011 Workshop at CERN, Jan 2011 (pre Higgs discovery)

“Statistical issues for search experiments”

Proceedings on website <http://indico.cern.ch/conferenceDisplay.py?confId=107747>

Overall Conclusions

- 1) Best of luck with your statistical analyses
- 2) Your statistical analysis should do justice to your data
- 3) Your problem has probably occurred before, and maybe has been solved

Consult text-books, and statistics information on the web, e.g.
CDF Statistics Committee
CMS Statistics Committee
Particle Data Group Statistics

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

Don't use your own square wheel if a statistician's circular one already exists

- 4) Send me an e-mail (l.lyons@physics.ox.ac.uk)