

# p-values and Discovery

Louis Lyons

Oxford

[l.lyons@physics.ox.ac.uk](mailto:l.lyons@physics.ox.ac.uk)

Kyoto,

May 2007



# PHYSTAT-LHC Workshop



on

## Statistical Issues for LHC Physics

CERN Geneva June 27-29, 2007

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

**Contacts**  
Louis Lyons      l.lyons@physics.ox.ac.uk  
Albert De Roeck      Albert.de.Roeck@cern.ch

**Conference secretary**  
Dorothee Denise      Dorothee.Denise@cern.ch

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

# TOPICS

## Discoveries

$H_0$  or  $H_0$  v  $H_1$

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

Goodness of Fit tests

## Why $5\sigma$ ?

Blind analyses

What is p good for?

Errors of 1<sup>st</sup> and 2<sup>nd</sup> kind

What a p-value is not

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

## THE paradox

Optimising for discovery and exclusion

Incorporating nuisance parameters

# DISCOVERIES

“Recent” history:

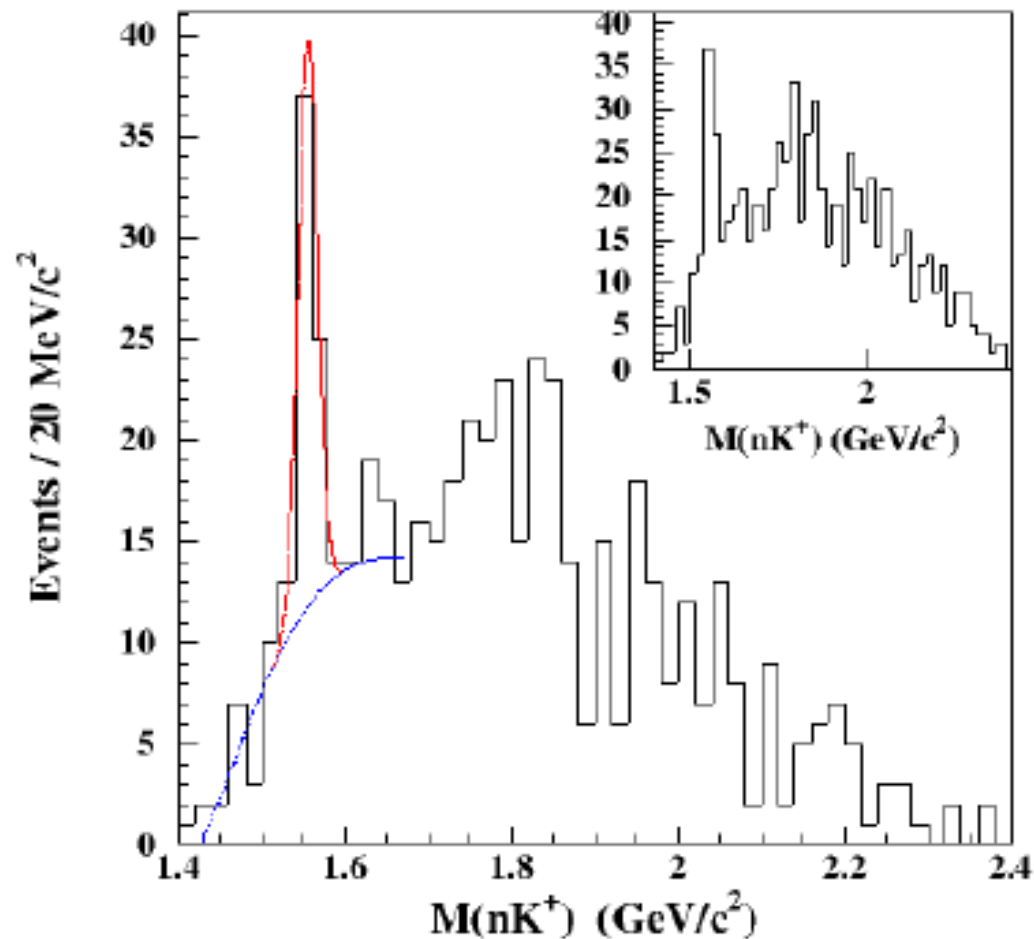
Charm	SLAC, BNL	1974
Tau lepton	SLAC	1977
Bottom	FNAL	1977
W,Z	CERN	1983
Top	FNAL	1995
{Pentaquarks	~Everywhere	2002 }
?	FNAL/CERN	2008?

? = Higgs, SUSY, q and l substructure, extra dimensions,  
free q/monopoles, technicolour, 4<sup>th</sup> generation, black holes,.....

QUESTION: How to distinguish discoveries from fluctuations or goofs?

# Penta-quarks?

Hypothesis testing: New particle or statistical fluctuation?



# H0 or H0 versus H1 ?

H0 = null hypothesis

e.g. Standard Model, with nothing new

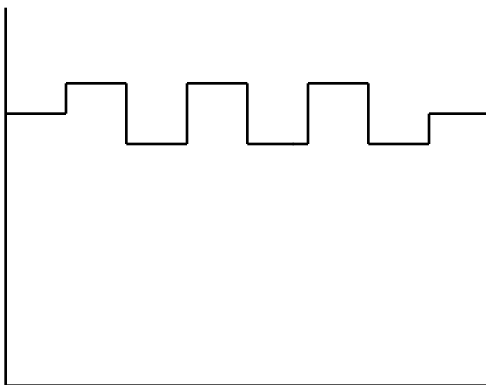
H1 = specific New Physics e.g. Higgs with  $M_H = 120$  GeV

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

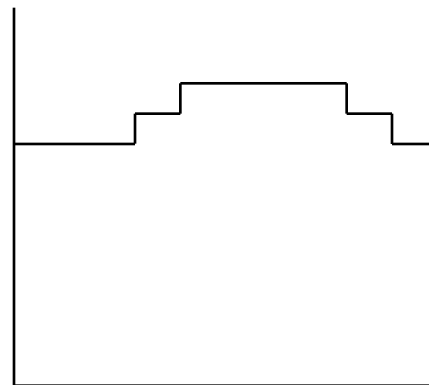
H0 v H1: “Hypothesis Testing” e.g.  $\mathcal{L}$ -ratio

Measures how much data favours one hypothesis wrt other

H0 v H1 likely to be more sensitive



or



# Testing H0:

## Do we have an alternative in mind?

1) Data is number (of observed events)

“H1” usually gives larger number

(smaller number of events if looking for oscillations)

2) Data = distribution. Calculate  $\chi^2$ .

Agreement between data and theory gives  $\chi^2 \sim \text{ndf}$

Any deviations give large  $\chi^2$

So test is independent of alternative?

Counter-example: Cheating undergraduate

3) Data = number or distribution

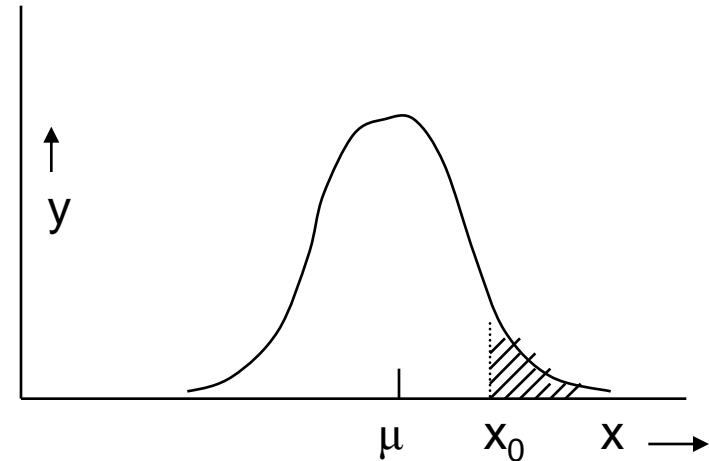
Use  $\mathcal{L}$ -ratio as test statistic for calculating p-value

4) H0 = Standard Model

# p-values

Concept of pdf

Example: **Gaussian**



$y$  = probability density for measurement  $x$

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

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:

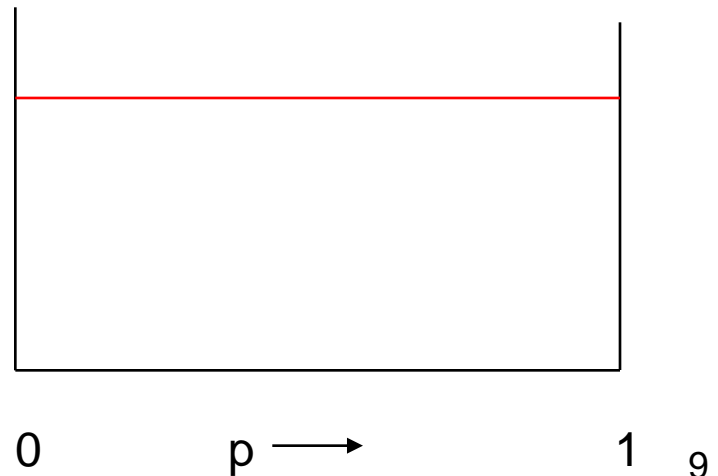
Gaussian pdf (no long tails)

Data is unbiased

$\sigma$  is correct

If so, Gaussian  $x \implies$  **uniform p-distribution**

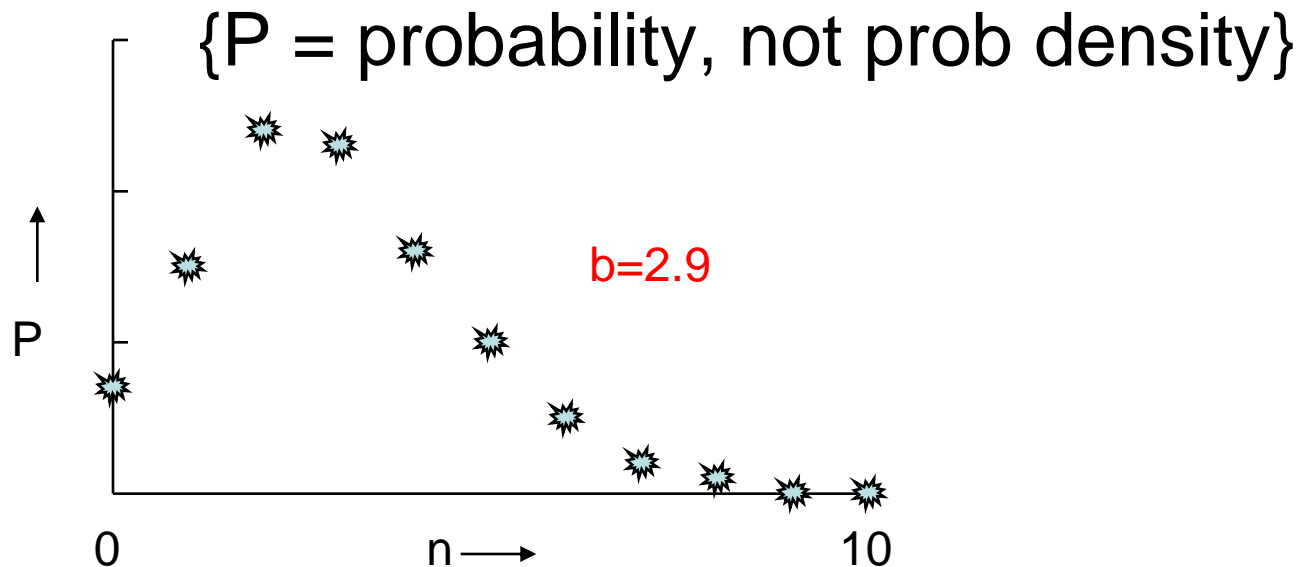
(Events at large  $x$  give small  $p$ )



# p-values for non-Gaussian distributions

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

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



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

# Poisson p-values

$n = \text{integer}$ , so **p has discrete values**

So p distribution cannot be uniform

Replace  $\text{Prob}\{p \leq p_0\} = p_0$ , for continuous p  
by  **$\text{Prob}\{p \leq p_0\} \leq p_0$** , for discrete p  
(equality for possible  $p_0$ )

**p-values often converted into equivalent Gaussian  $\sigma$**   
e.g.  $3 \cdot 10^{-7}$  is “ $5\sigma$ ” (one-sided Gaussian tail)

# Significance

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

## Potential Problems:

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

## For future experiments:

- Optimising  $S / \sqrt{B}$  could give  $S = 0.1$ ,  $B = 10^{-6}$

# Goodness of Fit Tests

Data = individual points, histogram, multi-dimensional,  
multi-channel

$\chi^2$  and number of degrees of freedom

$\Delta\chi^2$  (or  $\ln\mathcal{L}$ -ratio): Looking for a peak

Unbinned  $\mathcal{L}_{\max}$ ?

Kolmogorov-Smirnov

Zech energy test

Combining p-values

Lots of different methods. Software available from:

<http://www.ge.infn.it/statisticaltoolkit>

# $\chi^2$ with $\nu$ degrees of freedom?

1)  $\nu = \text{data} - \text{free parameters}$  ?

Why asymptotic (apart from Poisson  $\rightarrow$  Gaussian) ?

a) Fit flatish histogram with

$$y = N \{ 1 + 10^{-6} \cos(x - \mathbf{x}_0) \} \quad \mathbf{x}_0 = \text{free param}$$

b) Neutrino oscillations: almost degenerate parameters

$$\begin{array}{ll} y \sim 1 - \mathbf{A} \sin^2(1.27 \Delta m^2 L/E) & 2 \text{ parameters} \\ \xrightarrow{\text{Small } \Delta m^2} 1 - \mathbf{A} (1.27 \Delta m^2 L/E)^2 & 1 \text{ parameter} \end{array}$$

# $\chi^2$ with $\nu$ degrees of freedom?

2) Is difference in  $\chi^2$  distributed as  $\chi^2$  ?

H0 is true.

Also fit with H1 with  $k$  extra params

e. g. Look for Gaussian peak on top of smooth background

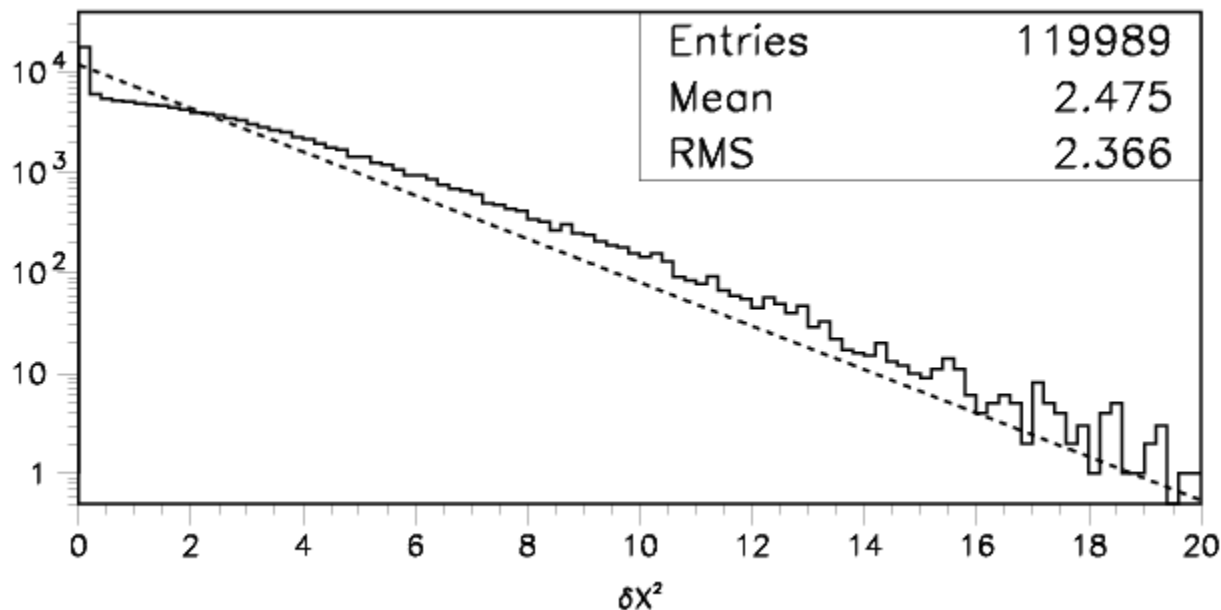
$$y = C(x) + A \exp\{-0.5 ((x-x_0)/\sigma)^2\}$$

Is  $\chi^2_{H0} - \chi^2_{H1}$  distributed as  $\chi^2$  with  $\nu = k = 3$  ?

Relevant for assessing whether enhancement in data is just a statistical fluctuation, or something more interesting

N.B. Under H0 ( $y = C(x)$ ) :  $A=0$  (boundary of physical region)  
 $x_0$  and  $\sigma$  undefined

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

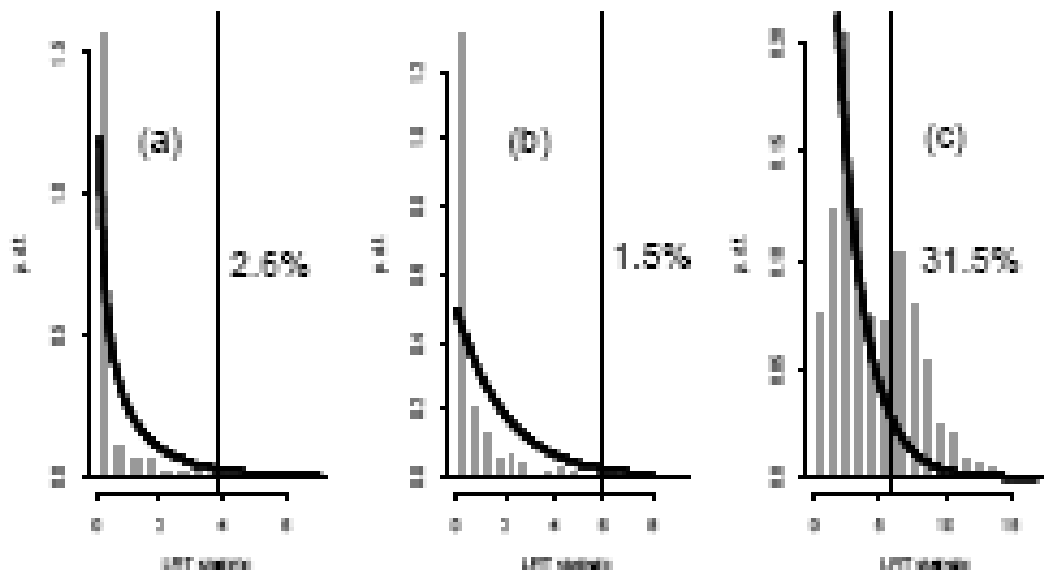


Demortier:

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 $\chi^2$ distributed as $\chi^2$ ?, contd.

So need to determine the  $\Delta\chi^2$  distribution by Monte Carlo

N.B.


- 1) Determining  $\Delta\chi^2$  for hypothesis H1 when data is generated according to H0 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)

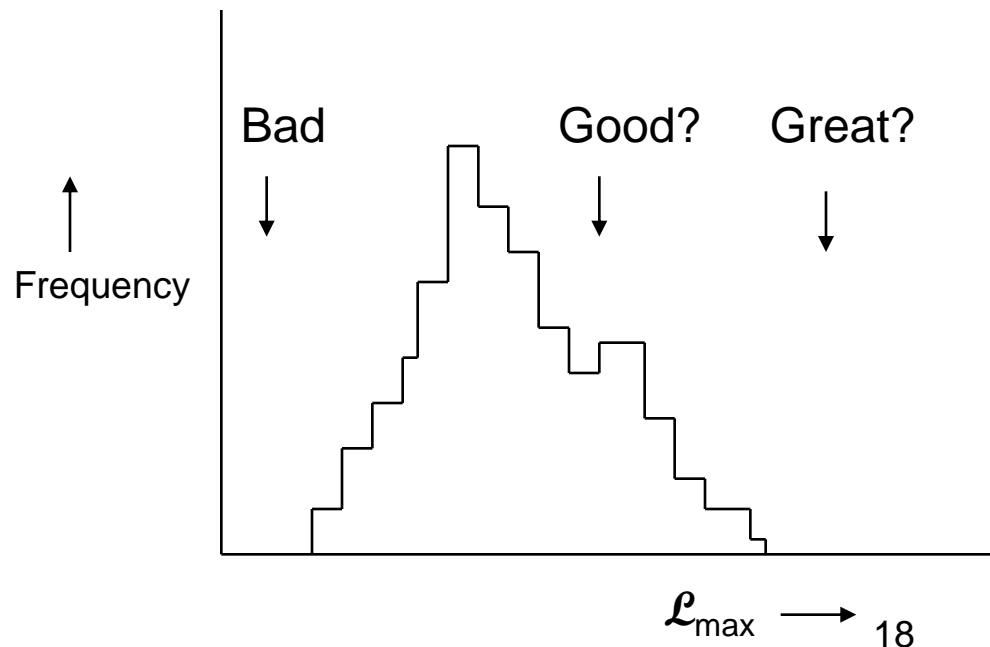
# Unbinned $\mathcal{L}_{\max}$ and Goodness of Fit?

Find params by maximising  $\mathcal{L}$

So larger  $\mathcal{L}$  better than smaller  $\mathcal{L}$

So  $\mathcal{L}_{\max}$  gives Goodness of Fit ??

Monte Carlo distribution  
of unbinned  $\mathcal{L}_{\max}$  



Not necessarily:

$$\mathcal{L}(\text{data}, \text{params})$$

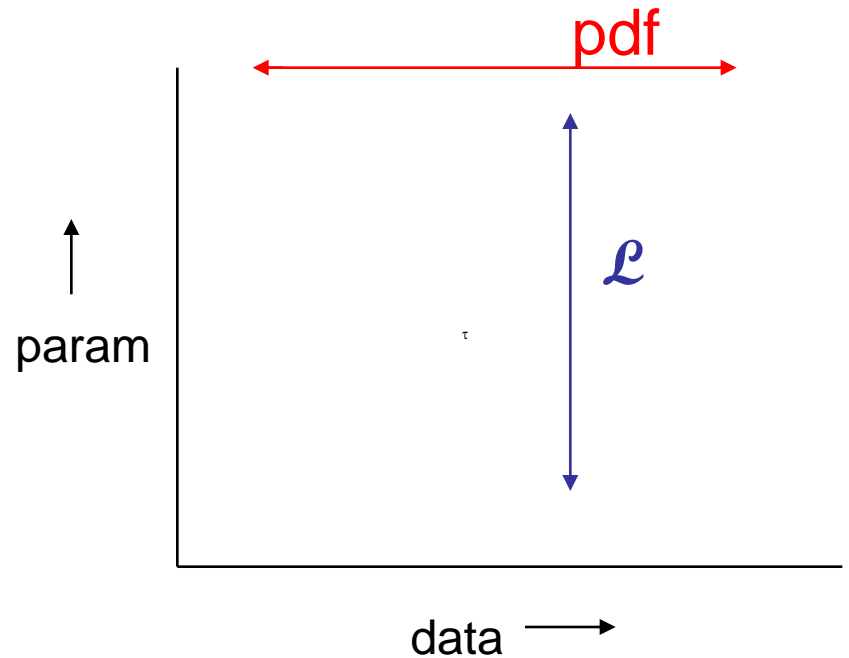


fixed vary

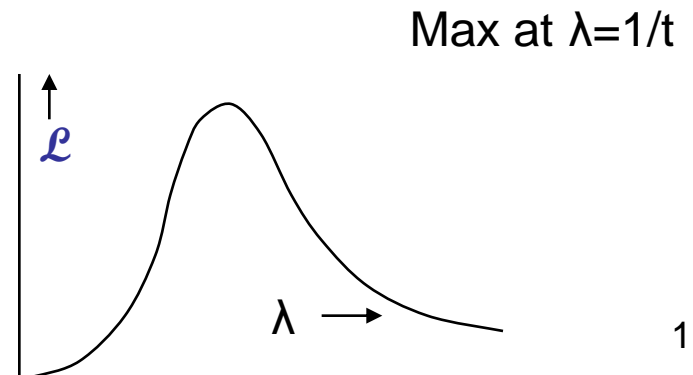
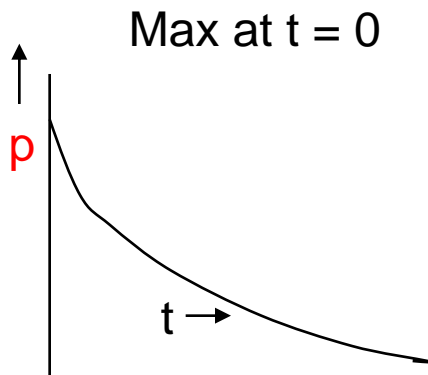
Contrast pdf( $\text{data}, \text{params}$ )



vary fixed



e.g.  $p(t, \lambda) = \lambda \cdot \exp(-\lambda t)$



## Example 1: Exponential distribution

Fit exponential  $\lambda$  to times  $t_1, t_2, t_3, \dots$

[Joel Heinrich, CDF 5639]

$$\mathcal{L} = \prod \lambda e^{-\lambda t}$$

$$\ln \mathcal{L}_{\max}^i = -N(1 + \ln t_{\text{av}})$$

i.e.  $\ln \mathcal{L}_{\max}$  depends only on **AVERAGE**  $t$ , but is

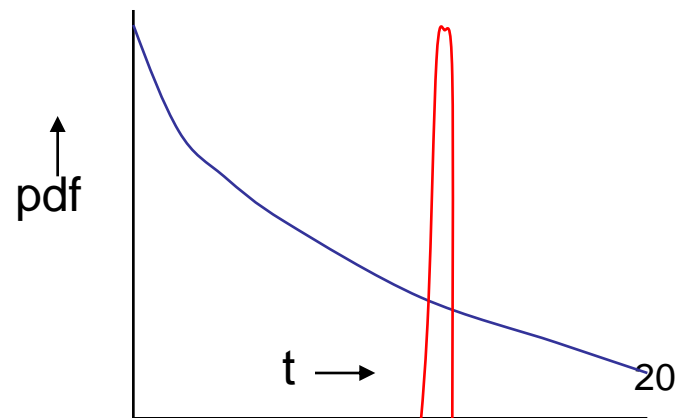
**INDEPENDENT OF DISTRIBUTION OF  $t$**  (except for.....)

(Average  $t$  is a sufficient statistic)

Variation of  $\mathcal{L}_{\max}$  in Monte Carlo is due to variations in samples' average  $t$ , but

**NOT TO BETTER OR WORSE FIT**

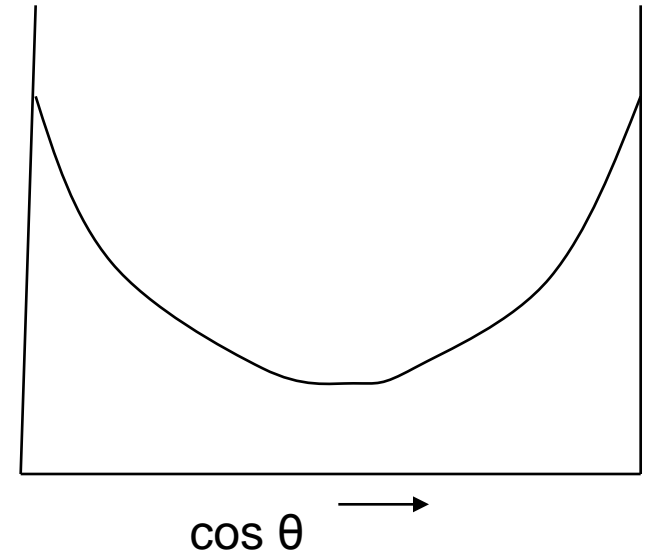
Same average  $t \implies$  same  $\mathcal{L}_{\max}$



## Example 2

$$\frac{dN}{d \cos \theta} = \frac{1 + \alpha \cos^2 \theta}{1 + \alpha / 3}$$

$$\mathcal{L} = \prod_i \frac{1 + \alpha \cos^2 \theta_i}{1 + \alpha / 3}$$



pdf (and likelihood) depends only on  $\cos^2 \theta_i$

Insensitive to **sign** of  $\cos \theta_i$

So data can be in very bad agreement with expected distribution

e.g. all data with  $\cos \theta < 0$ , but  $\mathcal{L}_{\max}$  does not know about it.

Example of general principle

### Example 3

Fit to Gaussian with variable  $\mu$ , fixed  $\sigma$

$$pdf = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right\}$$

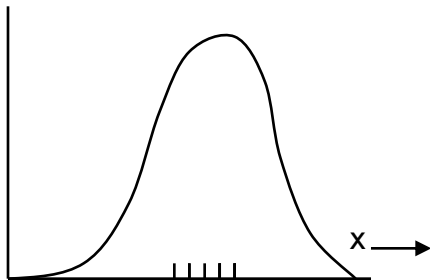
$$\ln \mathcal{L}_{\max} = N(-0.5 \ln 2\pi - \ln \sigma) - 0.5 \sum (x_i - x_{av})^2 / \sigma^2$$

↑
constant
↑
~variance(x)

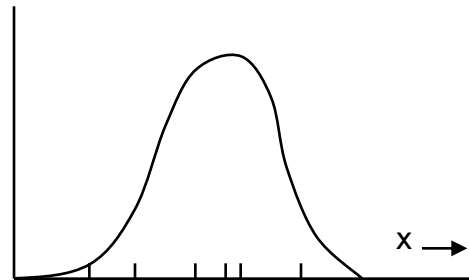
i.e.  $\mathcal{L}_{\max}$  depends only on variance(x),

which is not relevant for fitting  $\mu$  ( $\mu_{\text{est}} = x_{av}$ )

Smaller than expected variance(x) results in larger  $\mathcal{L}_{\max}$



Worse fit, larger  $\mathcal{L}_{\max}$



Better fit, lower  $\mathcal{L}_{\max}$

# $\mathcal{L}_{\max}$ and Goodness of Fit?

## Conclusion:

$\mathcal{L}$  has sensible properties with respect to parameters

**NOT** with respect to data

$\mathcal{L}_{\max}$  within Monte Carlo peak is **NECESSARY**

not **SUFFICIENT**

(‘Necessary’ doesn’t mean that you have to do it!)

# Goodness of Fit: Kolmogorov-Smirnov

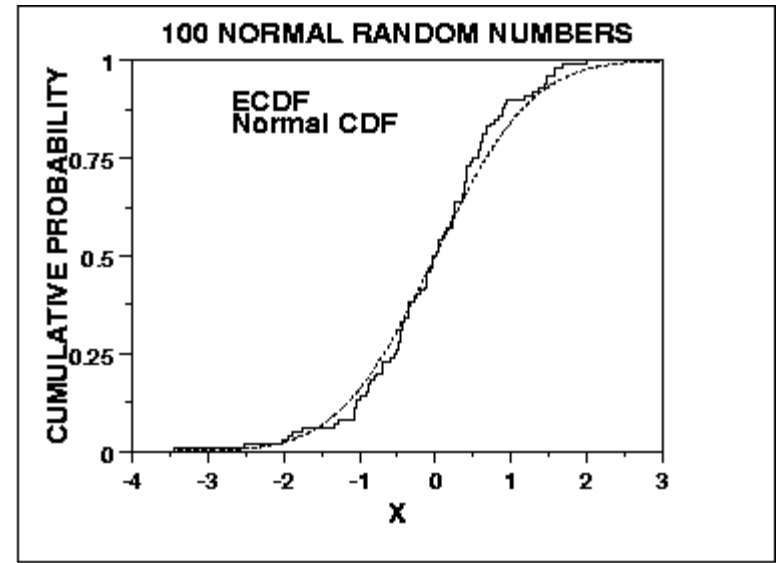
Compares data and model cumulative plots  
Uses largest discrepancy between dists.  
Model can be analytic or MC sample

Uses individual data points

Not so sensitive to deviations in tails  
(so variants of K-S exist)

Not readily extendible to more dimensions

Distribution-free conversion to  $p$ ; depends on  $n$   
(but not when free parameters involved – needs MC)





# Goodness of fit: 'Energy' test

Assign +ve charge to data  $\star$  ; -ve charge to M.C.  $\star$

Calculate 'electrostatic energy E' of charges

If distributions agree,  $E \sim 0$

If distributions don't overlap, E is positive

Assess significance of magnitude of E by MC

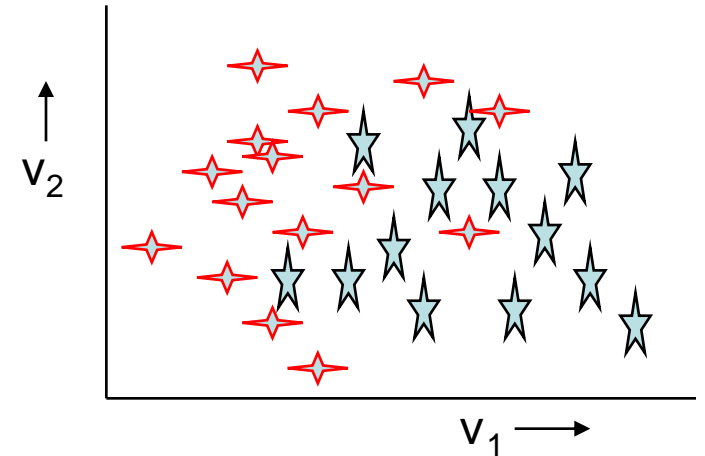
N.B.

- 1) Works in many dimensions
- 2) Needs metric for each variable (make variances similar?)
- 3)  $E \sim \sum q_i q_j f(\Delta r = |r_i - r_j|)$  ,  $f = 1/(\Delta r + \epsilon)$  or  $-\ln(\Delta r + \epsilon)$

Performance insensitive to choice of small  $\epsilon$

See [Aslan and Zech's](#) paper at:

<http://www.ippp.dur.ac.uk/Workshops/02/statistics/program.shtml>



# Combining different p-values

Several results quote 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,  $S = z * (1 - \ln z) \geq z$  )

Slight problem: **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  $x_1, x_2, x_3$ ”.

# Combining different p-values

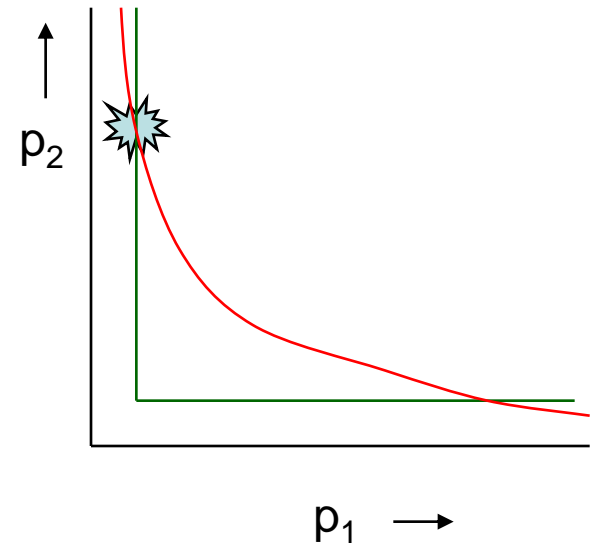
Conventional:

Are set of p-values consistent with H0?

SLEUTH:

How significant is smallest p?

$$1-S = (1-p_{\text{smallest}})^n$$



Combined S	$p_1 = 0.01$		$p_1 = 10^{-4}$	
	$p_2 = 0.01$	$p_2 = 1$	$p_2 = 10^{-4}$	$p_2 = 1$
Conventional	$1.0 \cdot 10^{-3}$	$5.6 \cdot 10^{-2}$	$1.9 \cdot 10^{-7}$	$1.0 \cdot 10^{-3}$
SLEUTH	$2.0 \cdot 10^{-2}$	$2.0 \cdot 10^{-2}$	$2.0 \cdot 10^{-4}$	$2.0 \cdot 10^{-4}$

# Why $5\sigma$ ?

- Past experience with  $3\sigma$ ,  $4\sigma$ ,... signals

- Look elsewhere effect:

Different cuts to produce data

Different bins (and binning) of this histogram

Different distributions Collaboration did/could look at

Defined in SLEUTH

- Bayesian priors:

$$\frac{P(H0|data)}{P(H1|data)} = \frac{P(data|H0) * P(H0)}{P(data|H1) * P(H1)}$$

Bayes posteriors

Likelihoods

Priors

Prior for {H0 = S.M.}  $\gg \gg$  Prior for {H1 = New Physics}



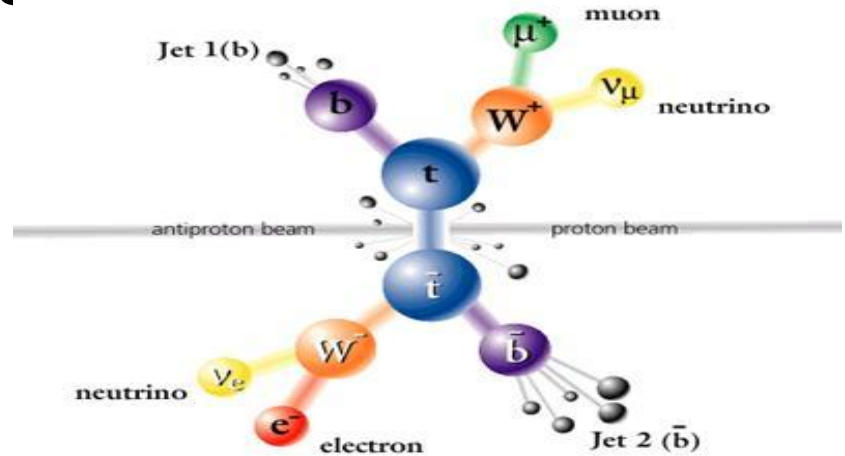
# Sleuth



a quasi-model-independent search strategy for new physics

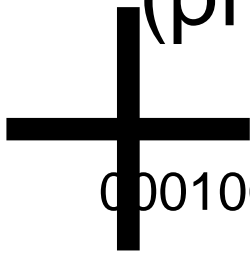
Assumptions:

1. Exclusive final state
2. Large  $\sum p_T$
3. An excess



0608025

(prediction) d(hep-ph)



0001001

Rigorously compute the trials factor associated with looking everywhere

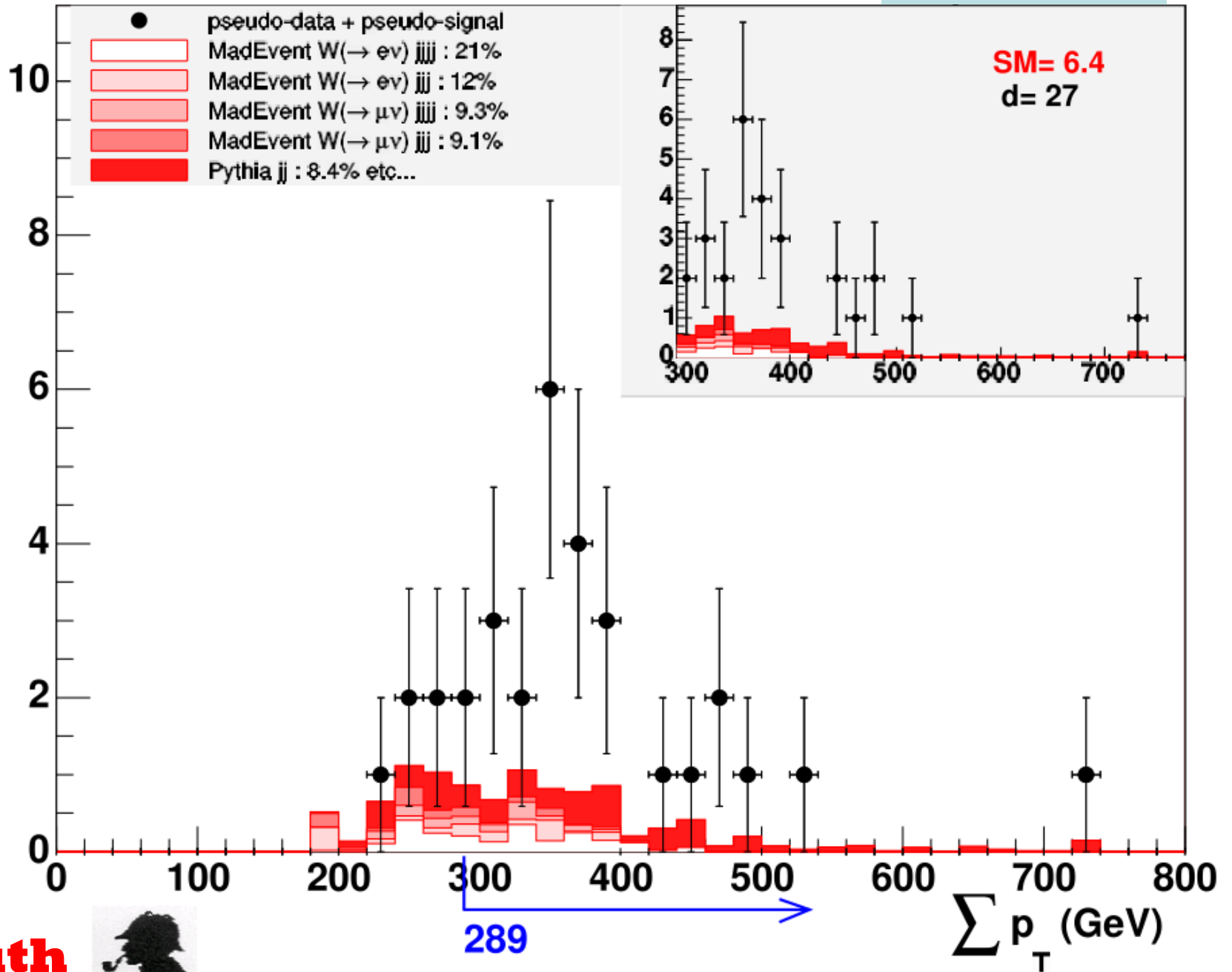
$W b \bar{b} j j$

pseudo discovery

$P_{Wb\bar{b}jj} < 8e-08$

$\tilde{P} < 4e-05$

Number of Events



Sleuth



# 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

# 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 1<sup>st</sup> kind

Should happen at rate  $\alpha$

OR

Fail to reject  $H_0$  when something else

( $H_1, H_2, \dots$ ) is true: Error of 2<sup>nd</sup> kind

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

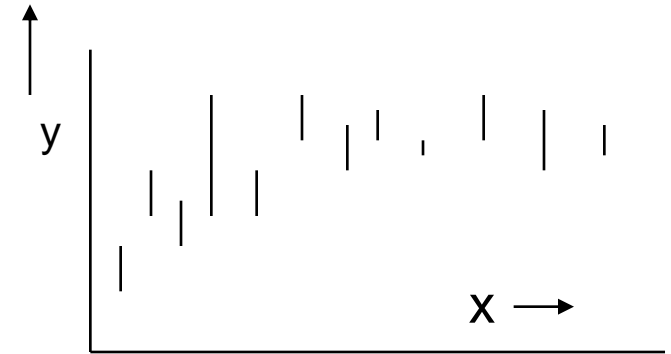


# Errors of 2<sup>nd</sup> kind: How often?

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

Calculate  $\chi^2$

Reject if  $\chi^2 \geq 20$



Error of 1<sup>st</sup> kind:  $\chi^2 \geq 20$  Reject H0 when true

Error of 2<sup>nd</sup> kind:  $\chi^2 \leq 20$  Accept H0 when in fact quadratic or..

How often depends on:

- Size of quadratic term

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

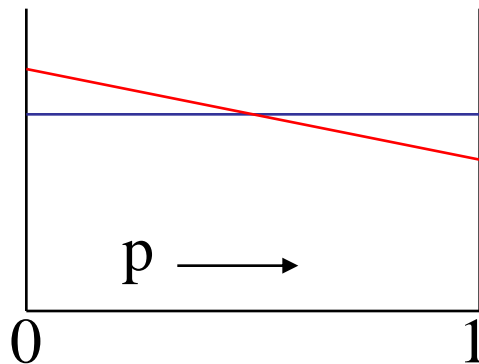
- How frequently quadratic term is present

# Errors of 2<sup>nd</sup> kind: How often?

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

Particles are  $\pi$  or  $\mu$

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



$\pi$  and  $\mu$  have similar masses

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

a)  $\sim$  half, for equal mixture of  $\pi$  and  $\mu$

b) almost all, for "pure"  $\pi$  beam

c) very few, for "pure"  $\mu$  beam

# What is p good for?

## Selecting sample of wanted events

e.g. kinematic fit to select  $t\bar{t}$  events

$t \rightarrow bW, b \rightarrow jj, W \rightarrow \mu\nu$      $\bar{t} \rightarrow \bar{b}W, \bar{b} \rightarrow jj, W \rightarrow jj$

Convert  $\chi^2$  from kinematic fit to p-value

Choose cut on  $\chi^2$  to select  $t\bar{t}$  events

Error of 1<sup>st</sup> kind: Loss of efficiency for  $t\bar{t}$  events

Error of 2<sup>nd</sup> kind: Background from other processes

Loose cut (large  $\chi^2_{\max}$ , small  $p_{\min}$ ): Good efficiency, larger bgd

Tight cut (small  $\chi^2_{\max}$ , larger  $p_{\min}$ ): Lower efficiency, small bgd

Choose cut to optimise analysis:

More signal events: Reduced statistical error

More background: Larger systematic error

# 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  $H_0 = \text{energy conservation}$

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

$P(\text{Data};\text{Theory}) \neq P(\text{Theory};\text{Data})$

Theory = male or female

Data = pregnant or not pregnant

$P(\text{pregnant ; female}) \sim 3\%$

$P(\text{Data};\text{Theory}) \neq P(\text{Theory};\text{Data})$

Theory = male or female

Data = pregnant or not pregnant

$P(\text{pregnant ; female}) \sim 3\%$

but

$P(\text{female ; pregnant}) \gg \gg 3\%$

# Aside: Bayes' Theorem

$$P(A \text{ and } B) = P(A|B) * P(B) = P(B|A) * P(A)$$

$$N(A \text{ and } B)/N_{\text{tot}} = N(A \text{ and } B)/N_B * N_B/N_{\text{tot}}$$

If A and B are independent,  $P(A|B) = P(A)$

Then  $P(A \text{ and } B) = P(A) * P(B)$ , but not otherwise

e.g.  $P(\text{Rainy and Sunday}) = P(\text{Rainy})$

$$\begin{aligned} \text{But } P(\text{Rainy and Dec}) &= P(\text{Rainy}|\text{Dec}) * P(\text{Dec}) \\ 25/365 &= 25/31 * 31/365 \end{aligned}$$

$$\text{Bayes' Th: } P(A|B) = P(B|A) * P(A) / P(B)$$

# More and more data

1) Eventually  $p(\text{data}|\text{H}_0)$  will be small, even if data and  $\text{H}_0$  are very similar.

$p$ -value does not tell you how different they are.

2) Also, beware of multiple (yearly?) looks at data.

“Repeated tests eventually sure to reject  $\text{H}_0$ , independent of value of  $\alpha$ ”

Probably not too serious –  
< ~10 times per experiment.

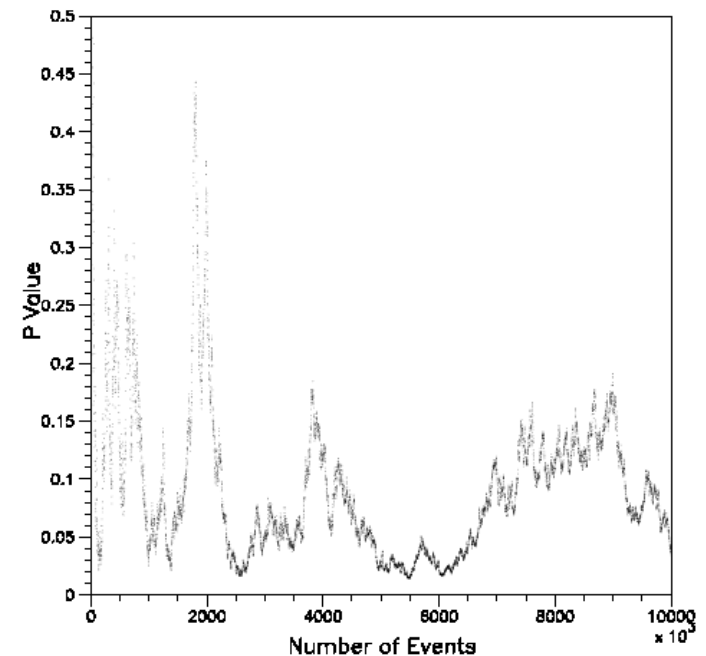
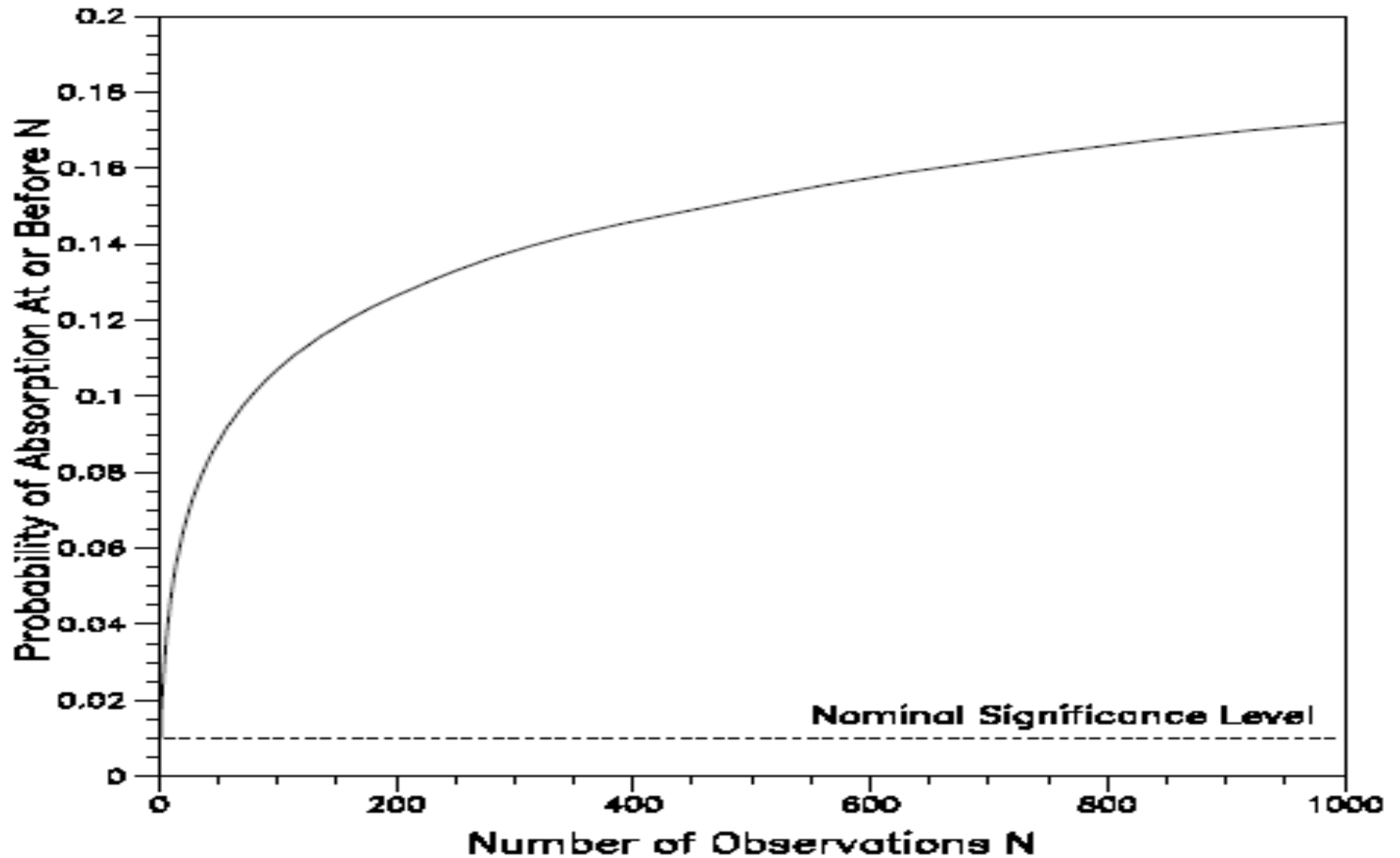


Figure 1:  $P$  value versus sample size.



# More “More and more data”



# PARADOX

Histogram with 100 bins

Fit 1 parameter

$S_{\min}$ :  $\chi^2$  with NDF = 99 (Expected  $\chi^2 = 99 \pm 14$ )

For our data,  $S_{\min}(p_0) = 90$

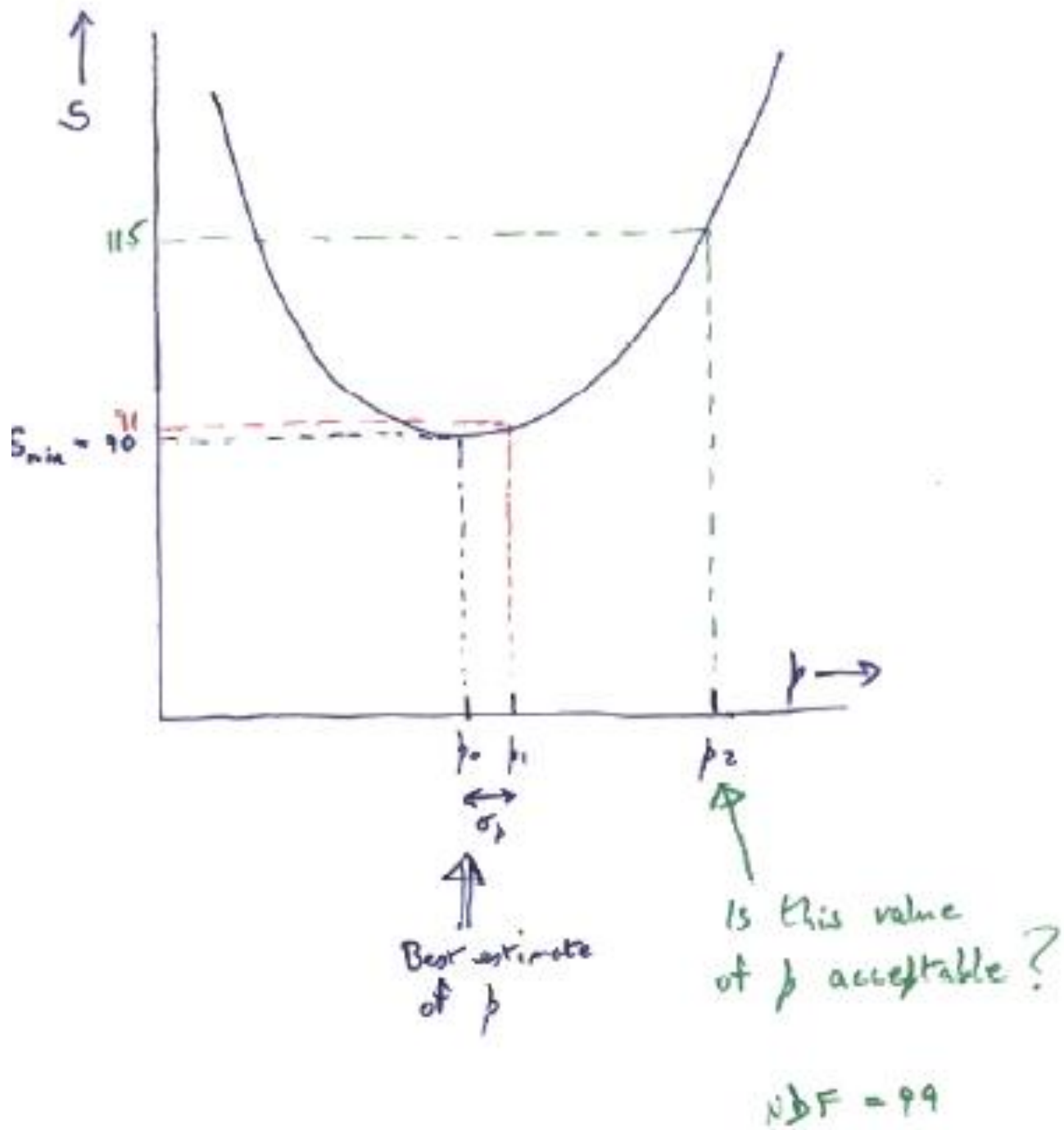
Is  $p_1$  acceptable if  $S(p_1) = 115$ ?

1) YES. Very acceptable  $\chi^2$  probability

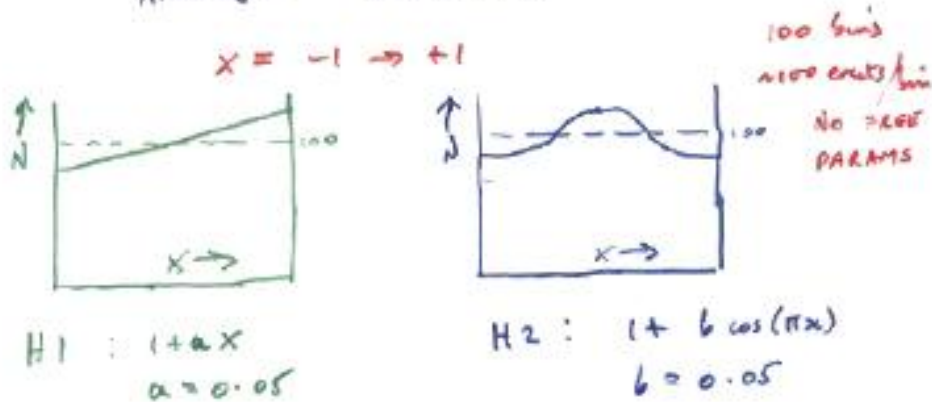
2) NO.  $\sigma_p$  from  $S(p_0 + \sigma_p) = S_{\min} + 1 = 91$

But  $S(p_1) - S(p_0) = 25$

So  $p_1$  is  $5\sigma$  away from best value



ANOTHER EXAMPLE



Generate events according to H1 (+ stat fluct)

Try fitting according to H1 or to H2

$\chi_1^2$                        $\chi_2^2$

Look at dist of  $\chi_1^2$                       As expected for  $NDF=100$

$\chi_2^2$                       Bit bigger. Many #  
"satisfactory"

$\chi_2^2 - \chi_1^2$                       Decision based on  $\Delta\chi^2$   
has much better power

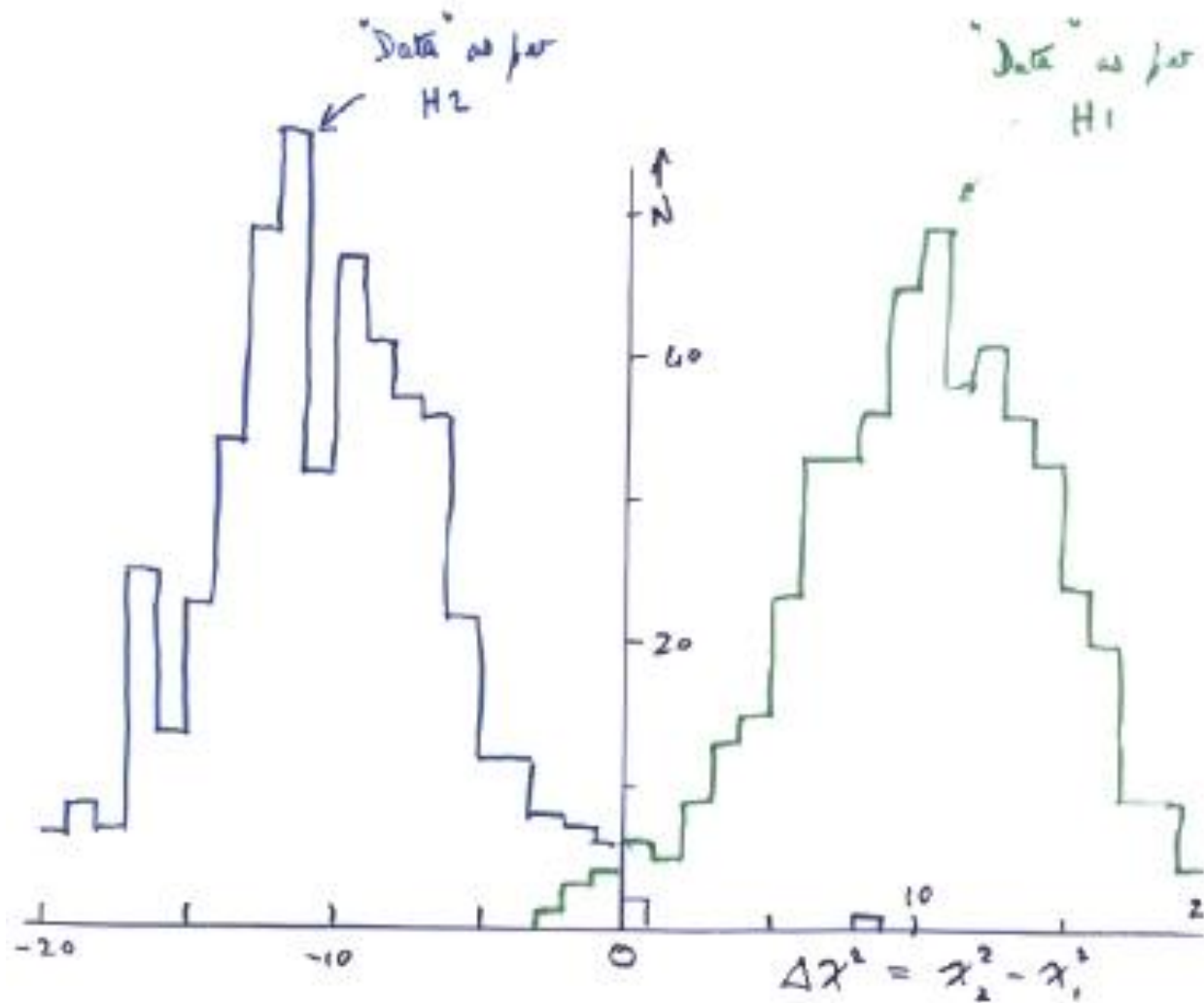
Repeat for events generated according to H2

Look at dist of  $\chi_1^2$   
 $\chi_2^2$   
 $\chi_2^2 - \chi_1^2$

\* 69% have  
 $\chi_2^2 < 130$

# DISTINGUISHING 2 HYPOTHESES ON BASIS OF $\Delta\chi^2$

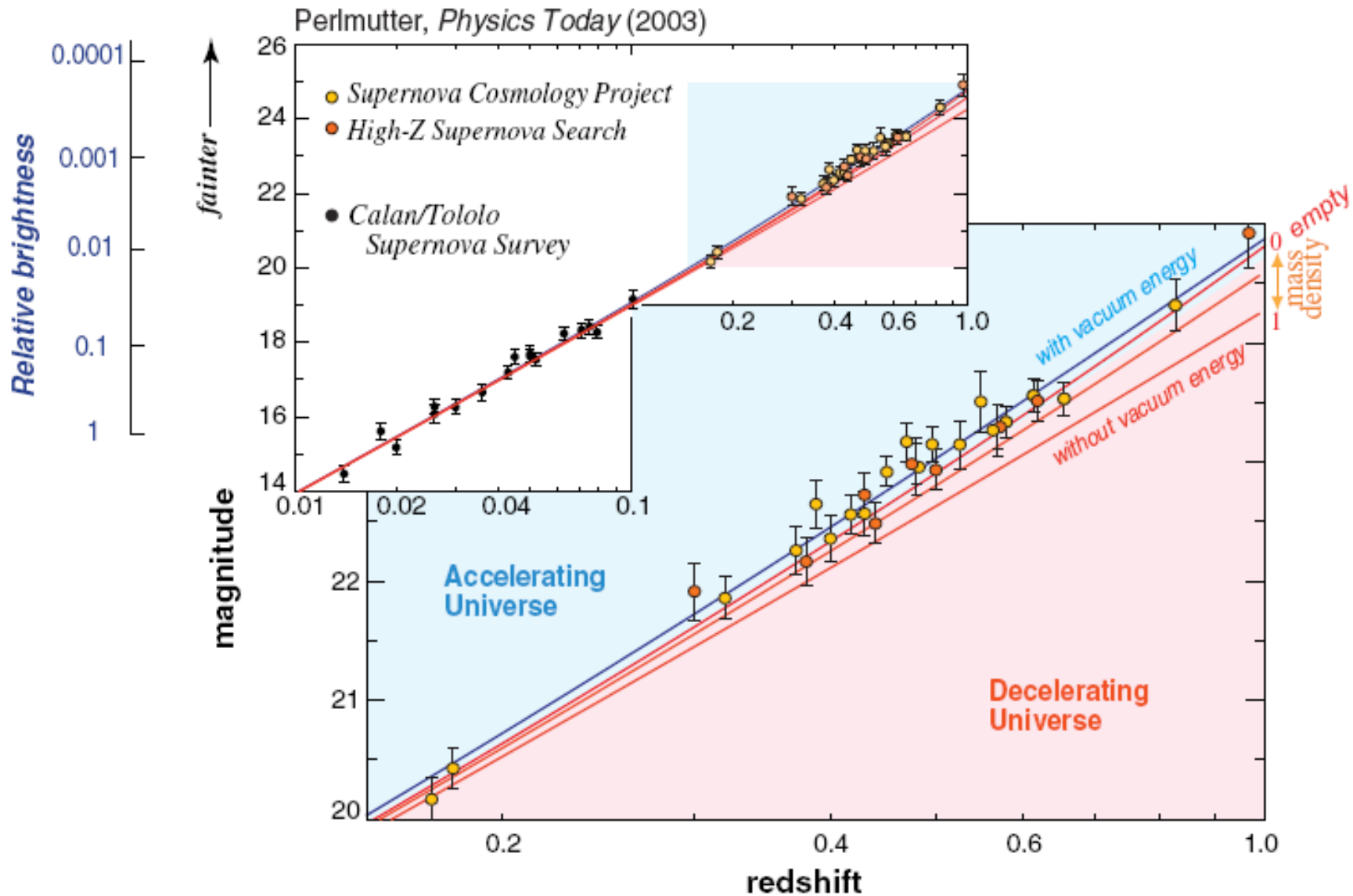
(500 SIMULATIONS)



$$H2 = 1 + 0.05 \cos(\pi x)$$

$$H1 = 1 + 0.05 x$$

# Comparing data with different hypotheses



# Choosing between 2 hypotheses

Possible methods:

$$\Delta\chi^2$$

$\ln\mathcal{L}$ -ratio

Bayesian evidence

Minimise “cost”

# Optimisation for Discovery and Exclusion

Giovanni Punzi, PHYSTAT2003:

“Sensitivity for searches for new signals and its optimisation”

<http://www.slac.stanford.edu/econf/C030908/proceedings.html>

Simplest situation: Poisson counting experiment,

Bgd =  $b$ , Possible signal =  $s$ ,  $n_{\text{obs}}$  counts

(More complex: Multivariate data,  $\ln\mathcal{L}$ -ratio)

Traditional sensitivity:

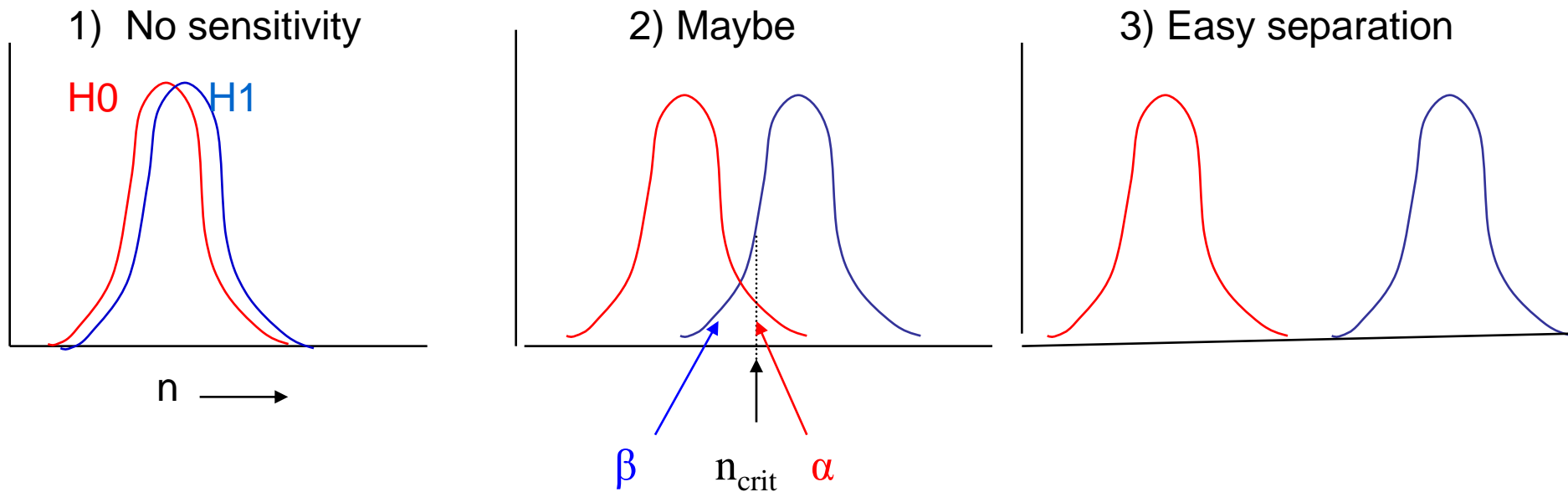
Median limit when  $s=0$

Median  $\sigma$  when  $s \neq 0$  (averaged over  $s$ ?)

Punzi criticism: Not most useful criteria

Separate optimisations





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

Given  $b$ ,  $\alpha$  determines  $n_{crit}$

$s$  defines  $\beta$ . For  $s > s_{min}$ , separation of curves  $\rightarrow$  discovery or excln

$s_{min}$  = Punzi measure of sensitivity For  $s \geq s_{min}$ , 95% chance of  $5\sigma$  discovery

Optimise cuts for smallest  $s_{min}$

Now data: If  $n_{obs} \geq n_{crit}$ , discovery at level  $\alpha$

If  $n_{obs} < n_{crit}$ , no discovery. If  $\beta_{obs} < 1 - CL$ , exclude  $H_1$

# 1) No sensitivity

Data almost always falls in peak

$\beta$  as large as 5%, so 5% chance of H1 exclusion even when no sensitivity. ( $CL_s$ )

# 2) Maybe

If data fall above  $n_{crit}$ , discovery

Otherwise, and  $n_{obs} \rightarrow \beta_{obs}$  small, exclude H1

(95% exclusion is easier than  $5\sigma$  discovery)

But these may not happen  $\rightarrow$  no decision

# 3) Easy separation

Always gives discovery or exclusion (or both!)

Disc	Excl	1)	2)	3)
No	No	□	□	
No	Yes		□	□
Yes	No		(□)	□
Yes	Yes			□!

# Incorporating systematics in p-values

Simplest version:

Observe  $n$  events

Poisson expectation for background only is  $b \pm \sigma_b$

$\sigma_b$  may come from:

acceptance problems

jet energy scale

detector alignment

limited MC or data statistics for backgrounds

theoretical uncertainties

Luc Demortier, “p-values: What they are and how we use them”, CDF memo June 2006

<http://www-cdfd.fnal.gov/~luc/statistics/cdf0000.ps>

Includes discussion of several ways of incorporating nuisance parameters

Desiderata:

Uniformity of p-value (averaged over  $\nu$ , or for each  $\nu$ ?)

p-value increases as  $\sigma_\nu$  increases

Generality

Maintains power for discovery

# Ways to incorporate nuisance params in p-values

- Supremum Maximise  $p$  over all  $v$ . Very conservative
- Conditioning Good, if applicable
- Prior Predictive Box. Most common in HEP  
$$p = \int p(v) \pi(v) dv$$
- Posterior predictive Averages  $p$  over posterior
- Plug-in Uses best estimate of  $v$ , without error
- $\mathcal{L}$ -ratio
- Confidence interval Berger and Boos.  
$$p = \text{Sup}\{p(v)\} + \beta$$
, where  $1-\beta$  Conf Int for  $v$
- Generalised frequentist Generalised test statistic

Performances compared by Demortier

# 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  $\chi^2_1$  and  $\chi^2_2$
- Blind analysis avoids personal choice issues
- Worry about systematics

PHYSTAT Workshop at CERN, June 27  $\rightarrow$  29 2007  
“Statistical issues for LHC Physics Analyses”

# Final message

Send interesting statistical issues to  
[I.lyons@physics.ox.ac.uk](mailto:I.lyons@physics.ox.ac.uk)