# $\mathcal{L}$ikelihoods
## 1) Brief Introduction
## 2) Do's & Dont's

Louis Lyons

Oxford & Imperial College

CMS

1

# Topics

What it is

How it works: Resonance

**Uncertainty estimates**

Coverage

Several Parameters

Do's and Dont's with $\mathcal{L}$ikelihoods:

COMBINING PROFILE  LIKELIHOODS
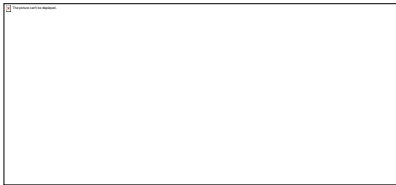NORMALISATION FOR LIKELIHOOD
$\Delta(\ln \mathcal{L}) = 0.5$ RULE
$\mathcal{L}_{max}$  AND GOODNESS OF FIT

Bayes and Frequentism: What is Probability?

2

# DO'S AND DONT'S WITH $\mathcal{L}$

- NORMALISATION FOR LIKELIHOOD

- JUST QUOTE UPPER LIMIT

- $\Delta(\ln \mathcal{L}) = 0.5$ RULE

- $\mathcal{L}_{max}$ AND GOODNESS OF FIT

- 

- BAYESIAN SMEARING OF $\mathcal{L}$

- USE CORRECT $\mathcal{L}$ (PUNZI EFFECT)

3

# Simple example: Parameter for Angular distribution

$$y = N (1 + \beta \cos^2\theta)$$
$$y_i = N (1 + \beta \cos^2\theta_i)$$
$$= \text{probability density of observing } \theta_i, \text{ given } \beta$$

$$\mathcal{L}(\beta) = \Pi \, y_i$$
$$= \text{probability density of observing the data set } y_i, \text{ given } \beta$$

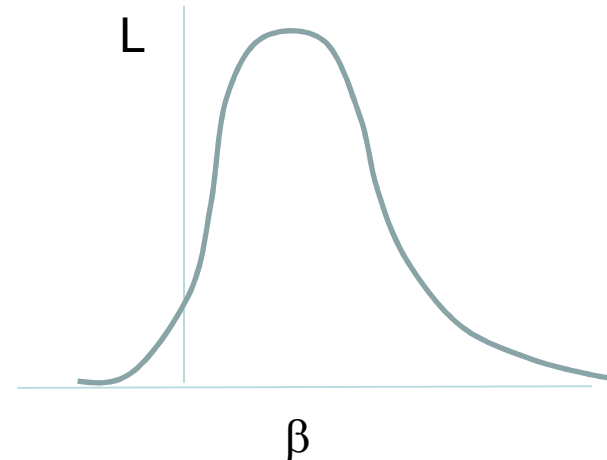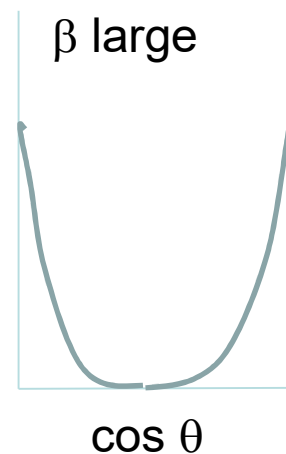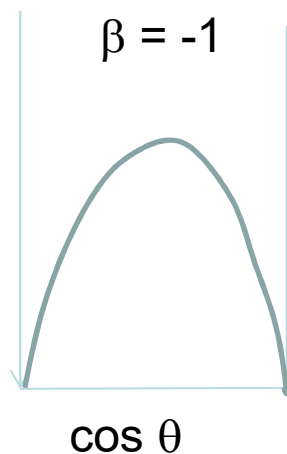Best estimate of $\beta$ is that which maximises $\mathcal{L}$

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

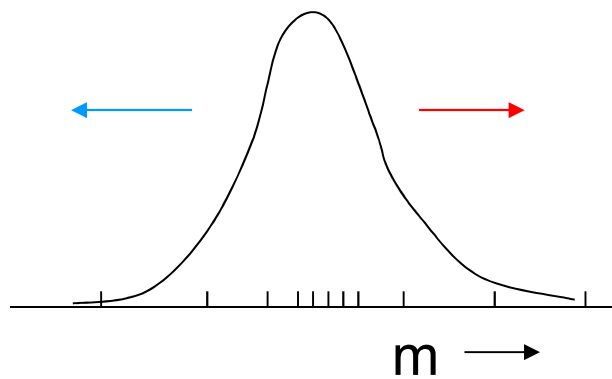**CRUCIAL** to normalise y        $N = 1/\{2(1 + \beta/3)\}$

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



$\beta$ = -1          $\cos \theta$

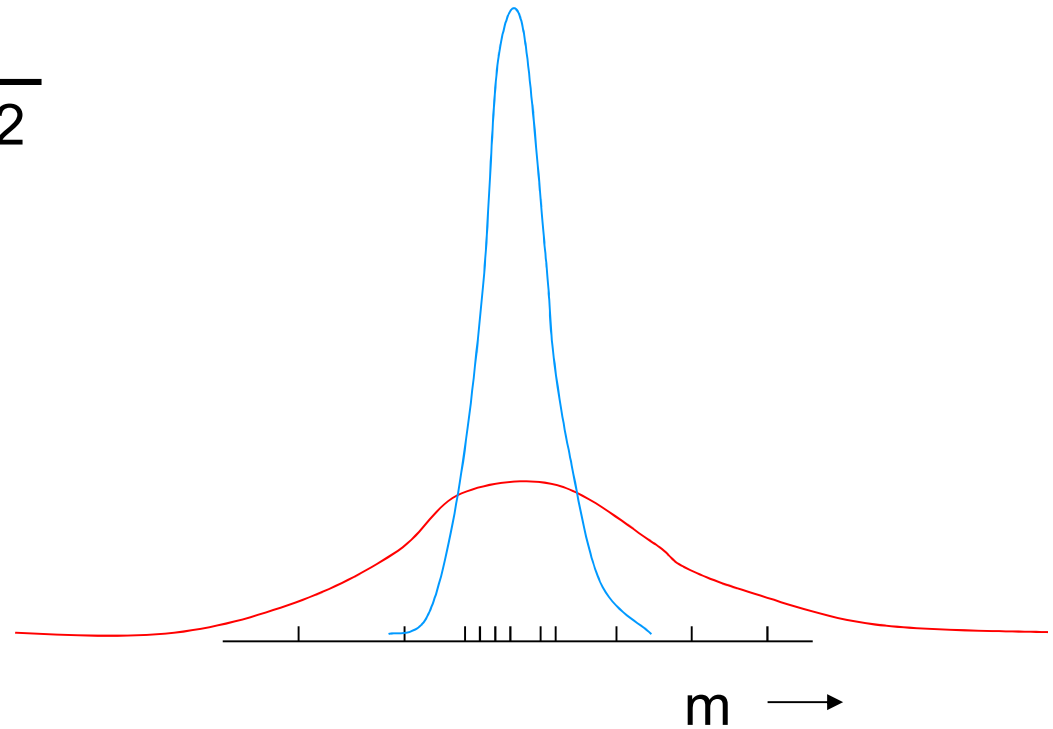$\beta$ large          $\cos \theta$

L          $\beta$

4

# How it works: Resonance

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



Vary $M_0$

Vary $\Gamma$

Find overall optimum by allowing both to vary simultaneously

5

Conventional to consider
$$\ell = \ln(\mathcal{L}) = \Sigma \ln(y_i)$$
For large N, $\mathcal{L} \rightarrow$ Gaussian

"Proof"

Taylor expand $\ell$ about its maximum

$$\ell = \ell_{max} + \tfrac{1}{2!}\ell''\left[\delta\left(\tfrac{\ell}{a}\right)\right]^2 + \cdots$$

$$= \ell_{max} - \tfrac{1}{2c}\delta^2 + \cdots \qquad c = -1/\ell''$$

$$\Rightarrow \mathcal{L} \sim \exp\left(-\tfrac{\delta^2}{2c}\right)$$



6

# Maximum likelihood uncertainty

Range of likely values of param $\mu$ from width of $\mathcal{L}$ or l dists.

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

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

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

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 μ with 68% probability"

Uncertainties from 3) usually asymmetric, and asym uncertainties are messy. So choose param sensibly

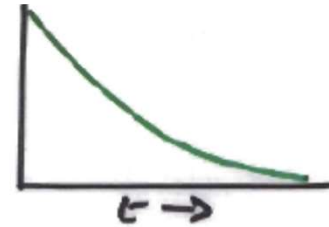e.g 1/p rather than p;        $\tau$ or λ

# Lifetime Determination

Realistic analyses are more complicated than this

$$\frac{dn}{dt} = \frac{1}{\tau} e^{-t/\tau}$$

↳ NORMALISATION

Observe $t_1, t_2 \ldots\ldots t_N$

Use pdf to construct

$$\mathcal{L} = \pi \left(\frac{dn}{dt}\right)_i = \pi \frac{1}{\tau} e^{-t_i/\tau}$$

$$\therefore \quad \ell = \sum_i (-t_i/\tau - \ln \tau)$$

$$\frac{\partial \ell}{\partial \tau} = \sum \left(+ t_i/\tau^2 - \frac{1}{\tau}\right) = 0 = \frac{\sum t_i}{\tau^2} - \frac{N}{\tau}$$

$$\Rightarrow \quad \tau = \sum t_i / N = \bar{t_i} \qquad \text{"Obvious"}$$

$$\frac{\partial^2 \ell}{\partial \tau^2} = -\sum \frac{2t_i}{\tau^3} + \sum \frac{1}{\tau^2} = -2\frac{N}{\tau^2} + \frac{N}{\tau^2} = -\frac{N}{\tau^2}$$

$$\Rightarrow \quad \sigma_\tau = 1 \bigg/ \sqrt{-\frac{\partial^2 \mathcal{L}}{\partial \tau^2}} = \tau / \sqrt{N}$$

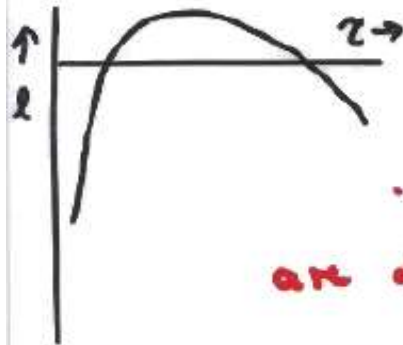N.B. 1) Usual $1/\sqrt{N}$ behaviour

2) $\sigma_\tau \propto \tau_{est}$

BEWARE FOR AVERAGING RESULTS

9

$$\ln \tau - \ln \tau_{max} = \text{Universal Fn of } \tau/\tau_{max}$$

$$\ell(\tau) = \sum -t_i/\tau - N \ln \tau$$

$$\ell(\tau) - \ell(\tau_{max}) = -N \tau_{max}/\tau - N \ln \tau$$
$$+ N + N \ln \tau_{max}$$
$$= N\left[1 + \ln(\tau_{max}/\tau) - \tau_{max}/\tau\right]$$



$$\therefore \text{ For given } N, \ \sigma_+ \approx \sigma_-$$

$$\text{are defined } \left(\sim \frac{\tau_{max}}{\sqrt{N}} \text{ as } N \to \infty\right)$$

For small $N$, $\sigma_+ > \sigma_-$

——"——

$$\ell(\tau_{max}) = -N(1 + \ln \bar{t})$$

N.B. $\ell(\tau_{max})$ depends only on $\bar{t}$, but **not** on distribution of $t_i$

Relevant for whether $\ell_{max}$ is useful for testing goodness of fit

# SEVERAL PARAMETERS

1 param p        $l = \ln\mathcal{L}$
. p from $dl/dp = 0$   i.e from max of $\mathcal{L}$
$$\sigma_p{}^2 = 1/(-d^2l/dp^2)$$
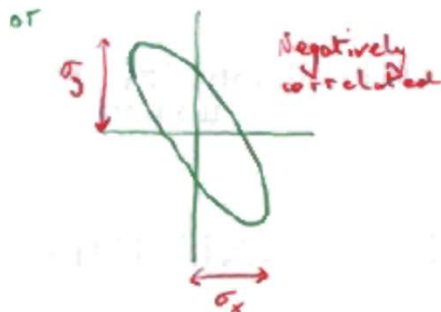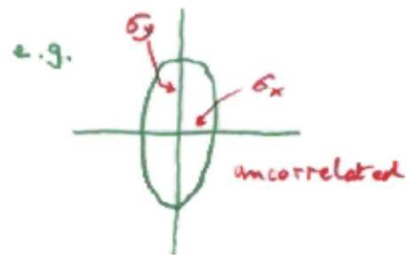
Many dimensions   $l(p_1, p_2,....)$
$p_1$, $p_2$, from $dl/dp_i = 0$  i.e. from max of $\mathcal{L}$
For uncertainties, define
$H_{ij} = d^2l/dp_i\, dp_j$ = Inverse Covariance Matrix
Covariance Matrix $E_{ij} = (H^{-1})_{ij}$
Diagonal Elements for  variances off-diag for covariances



For many params:
N.B1 Ellipsoid with $l=l_{max}$ -0.5 does not have 68% assymptotic coverage.
N.B2 Uncert on x is not given by varying x till $l = l_{max}$ -0.5, while keeping all other params constant

## PROFILE $\mathcal{L}$

$\mathcal{L}_{prof} = \mathcal{L}(\beta, \nu_{best}(\beta))$,  where
$\beta$ = param of interest
$\nu$ = nuisance param(s)
Uncertainty on $\beta$ from decrease in $\ln(\mathcal{L}_{prof})$ by 0.5

11

# ML and EML

ML uses fixed (data) normalisation
EML has normalisation as parameter

Example 1:  Cosmic ray experiment
                          See 96 protons     and     4 heavy nuclei
    ML estimate       96 ± 2% protons       4 ±2% heavy nuclei
    EML estimate      96 ± 10 protons        4 ± 2 heavy nuclei


Example 2:  Decay of resonance
    Use ML  for Branching Ratios
    Use EML for Partial Decay Rates

# Extended Maximum Likelihood

Maximum Likelihood uses shape → parameters
Extended Maximum Likelihood  uses shape and normalisation
i.e. EML uses prob of observing:
    a) sample of N events;    and
    b) given data distribution in x,……
            → shape parameters and normalisation.


Example:   Angular distribution
        Observe N events total            e.g  100
                    F forward                        96
                    B backward                        4

| Rate estimates | ML | EML |
|---|---|---|
| Total | --- | 100±10 |
| Forward | 96±2 | 96±10 |
| Backward | 4±2 | 4± 2 |

# DO'S AND DONT'S WITH $\mathcal{L}$

- COMBINING PROFILE $\mathcal{L}$s

- NORMALISATION FOR $\mathcal{L}$IKELIHOOD

- $\Delta(\ln \mathcal{L}) = 0.5$ RULE

- $\mathcal{L}_{max}$ AND GOODNESS OF FIT

- PDFs and $\mathcal{L}$IKELIHOODS

a) Max Like

Prob for fixed $N$ $=$ Binomial

$\text{Prod of}$ $\text{forwards}$ $\longrightarrow = f^F (1-f)^B \dfrac{N!}{F! \, B!}$ ✦

Maximise $\ln P_a$ wrt $f$ $\Rightarrow \hat{f} = F/N$

$\text{Error on } \hat{f}: \quad 1/\sigma^2 = -\dfrac{\partial^2 \ln P_a}{\partial f^2}$

$\qquad\qquad\qquad = \dfrac{N}{\hat{f}(1-\hat{f})} \qquad f = \hat{f}$

$\Rightarrow$ Estimate of $\hat{F} = Nf = F \pm \sqrt{FB/N}$ ← Completely

$\qquad\qquad\qquad \hat{B} = N(1-f) = B \pm \sqrt{FB/N}$ ← anti-corr

b) EML $\quad P_b = P_a \times \dfrac{e^{-\nu} \nu^N}{N!}$ $\qquad$ Poisson for overall rate

— Expected overall rate (pointing to $\nu$)

Maximise $\ln P_b(\nu, f)$

$\Rightarrow \quad \hat{\nu} = N \pm \sqrt{N}$ ⟩ uncorrelated

$\qquad \hat{f} = F/N \pm \sqrt{\dfrac{f(1-f)}{N}}$

For $\hat{F}$ & $\hat{B}$, either propagate errors for $\hat{F} = \hat{\nu}\hat{f}$

$\qquad\qquad\qquad\qquad\qquad\qquad \hat{B} = \hat{\nu}(1-\hat{f})$

or rewrite eqn ✦ as product of 2 indep Poissons

$\hat{F} = F \pm \sqrt{F}$ ⎫
$\hat{B} = B \pm \sqrt{B}$ ⎭

15

# Danger of combining profile $\mathcal{L}$s

Experiments quote $\mathcal{L}$ikelihood, profiled over nuisance parameters, so that combinations can be performed.
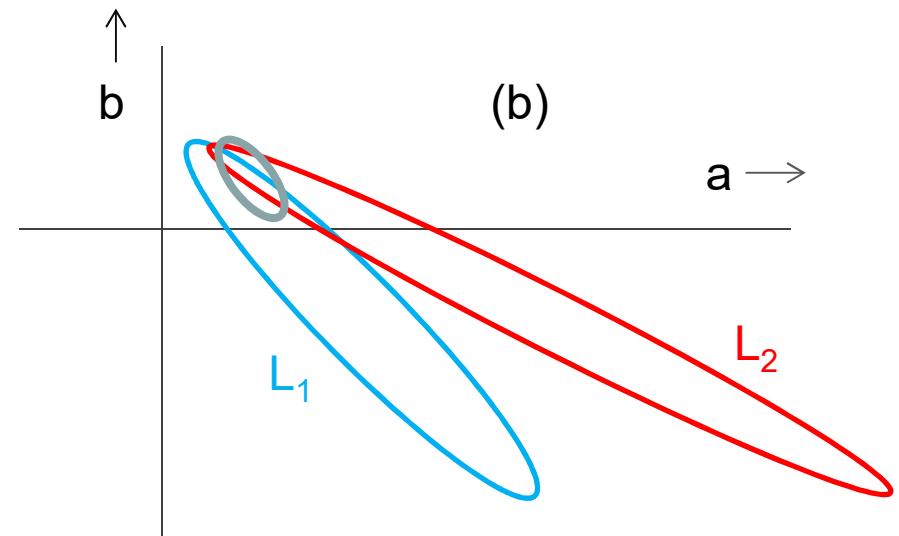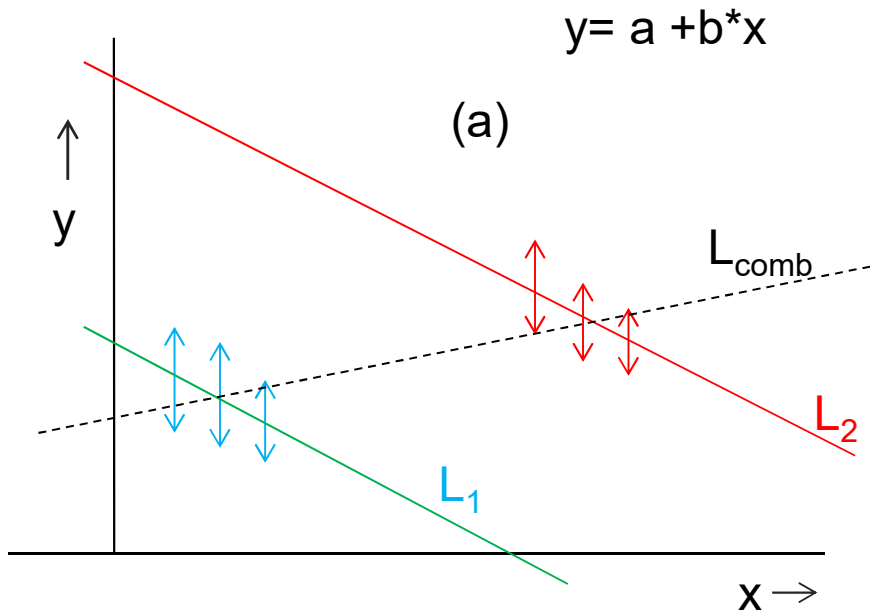
Very simple 'tracking' example:

* No magnetic field

* 2-D fit of straight line y = a + bx

      a = parameter of interest,   b = nuisance param

* Track hits in 2 subdetectors, each of 3 planes

$$y = a + b*x$$

(a)

$y$ ↑

$x$ →

$L_{comb}$

$L_2$

$L_1$

(b)

$b$ ↑

$a$ →

$L_1$

$L_2$

(c)

$\ln\mathcal{L}_{prof}$ ↑
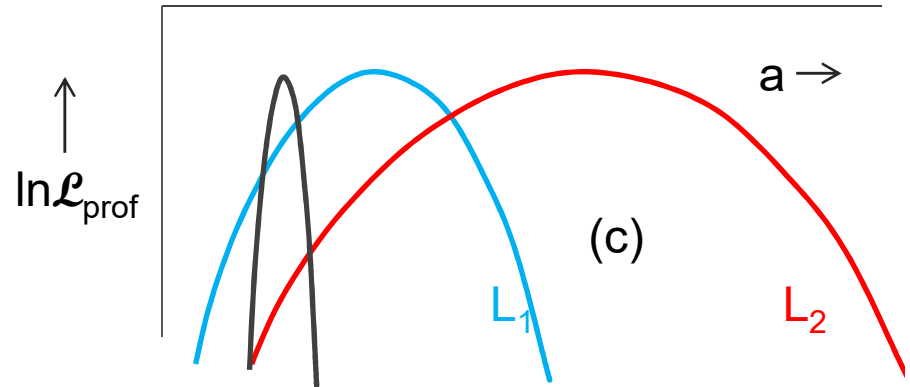
$a$ →

$L_1$

$L_2$

(d)

$b_{best}$ ↑

$a$ →

$L_1$

$L_2$

17

(a) Hits in 2 sub-detectors, each with 3 planes

(b) Covariance ellipses for separate fits $L_1$ and $L_2$, and combined $L_{comb}$

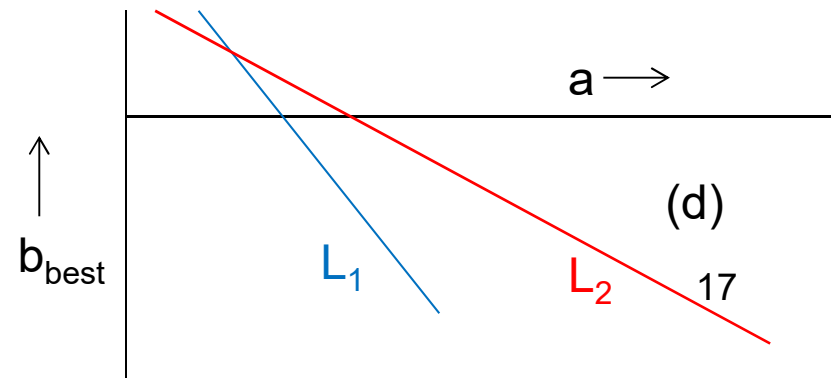(c) $\ln\mathcal{L}_{prof}$ as function of a, for all 3 lines

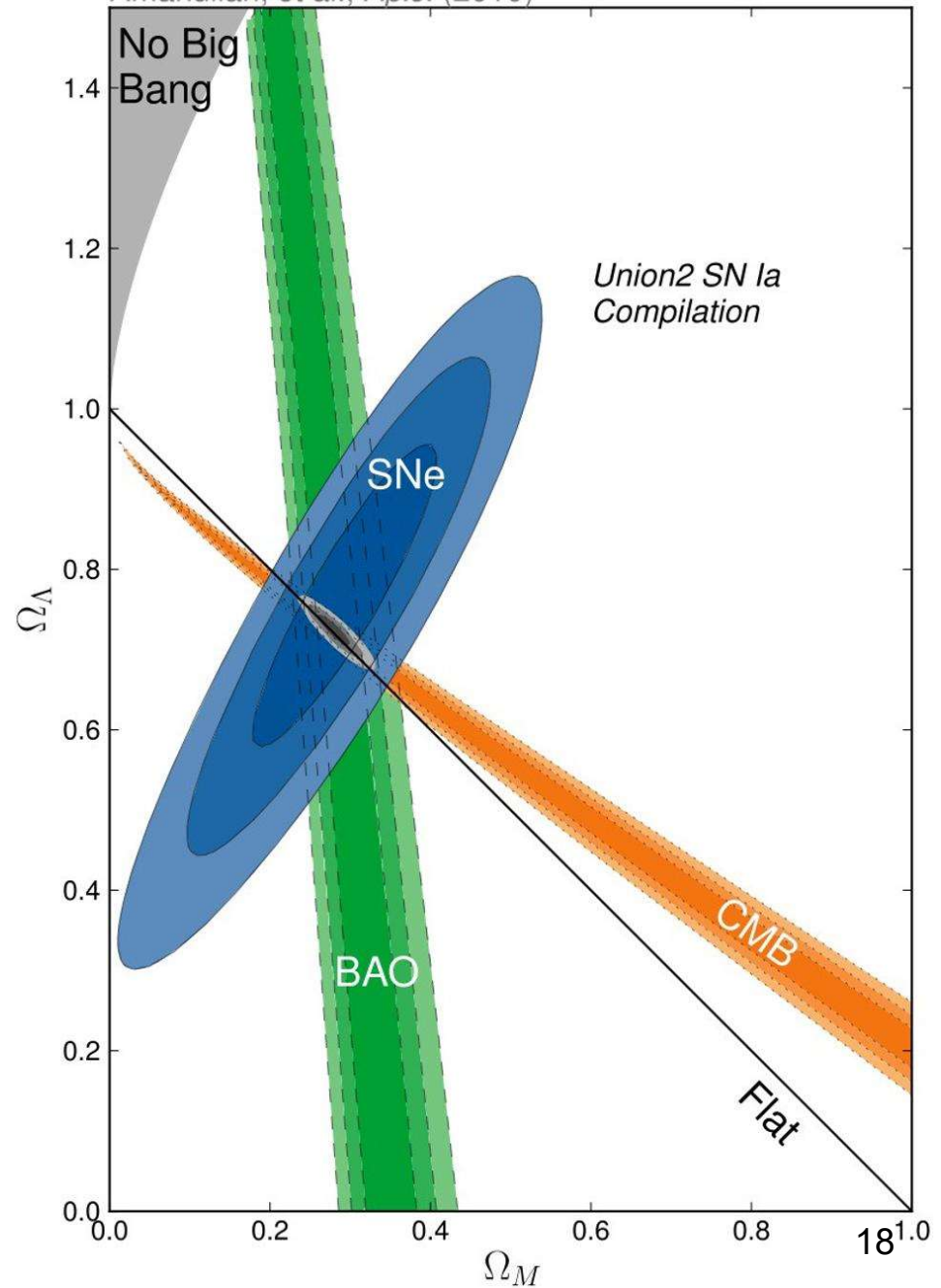(d) $b_{best}$ as a function of a

    N.B. $b_{best}$ for $L_1$ and $L_2$ are the same

*** Combining $\mathcal{L}_{prof}$ for $L_1$ and $L_2$ loses a lot of information, and $a_{best}$ wrong *****

## COSMOLOGY EXAMPLE

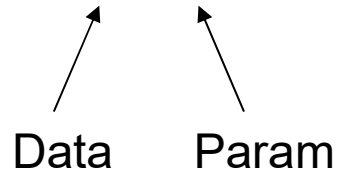Plot of dark energy fraction v dark matter fraction by various methods. Each determines dark energy fraction poorly, but combination is fine, because of different correlations.

Combining Profile Likelihoods would give very large uncertainty on dark energy fraction.



Supernova Cosmology Project
Amanullah, et al., *Ap.J.* (2010)

No Big Bang

Union2 SN Ia Compilation

SNe

BAO

CMB

Flat

$\Omega_\Lambda$

$\Omega_M$

18

# NORMALISATION FOR LIKELIHOOD

$\int P(x \mid \mu) \, dx$ **MUST** be independent of μ
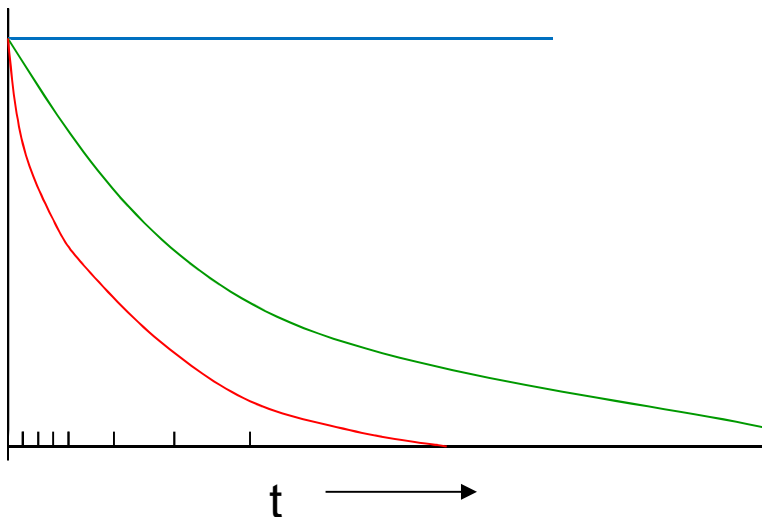
Data     Param

$$[\tau = \sum t_i / N]$$

### Exponential Distribution

**INCORRECT**     $P(t \mid \tau) = \quad e^{-t/\tau}$

Missing $1/\tau$



——— $\tau$ infinite

——— $\tau$ too large

——— $\tau$ about right

t →

19

# QUOTING UPPER LIMIT

**"We observed no significant signal, and our 90% conf upper limit is ….."**

**Need to specify method   e.g.**

$\mathcal{L}$

**Chi-squared (data or theory error)**

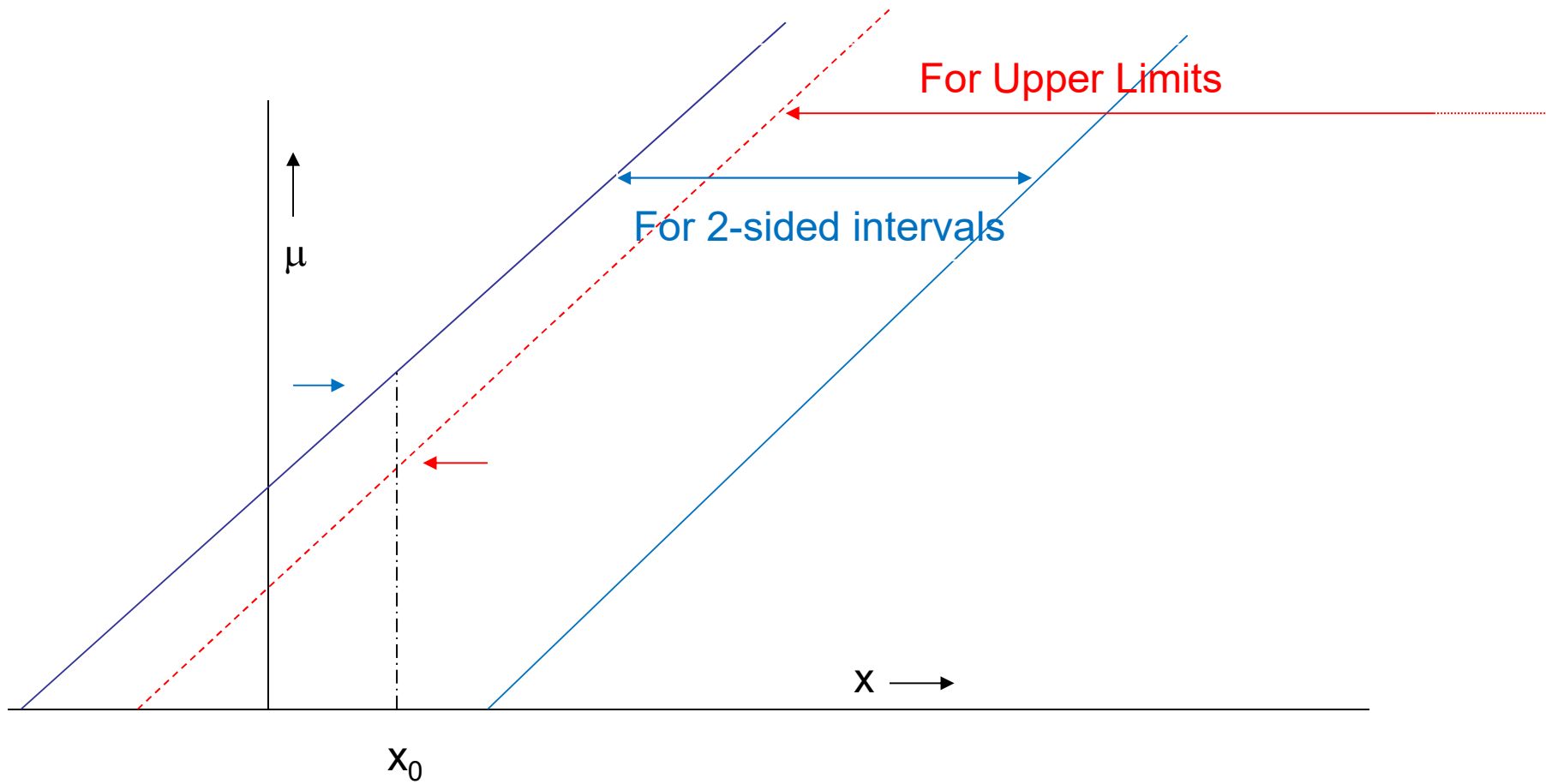**Frequentist  (Central or upper limit)**

**Feldman-Cousins**

**Bayes with prior = const,**

"Show your $\mathcal{L}$"

1) Not always practical

2) Not sufficient for frequentist methods

# 90% C.L. Upper Limits

For Upper Limits

For 2-sided intervals

μ

x

$x_0$

# Δln$\mathcal{L}$ = -1/2 rule

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

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

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

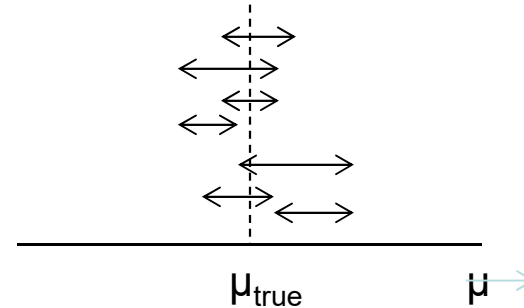3) $\ln(\mathcal{L}(\mu_0\pm\sigma)) = \ln(\mathcal{L}(\mu_0)) -1/2$

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 μ with 68% probability"

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

Barlow: Phystat05

# Coverage



* What it is:

For given statistical method applied to many sets of data to extract confidence intervals for param $\mu$, coverage C is fraction of ranges that contain true value of param.      Can vary with $\mu$

* Does not apply to **your** data:

It is a property of the **statistical method** used

It is **NOT** a probability statement about whether $\mu_{true}$ lies in your confidence range for $\mu$

* Coverage plot for Poisson counting expt

Observe n counts

Estimate $\mu_{best}$ from maximum of likelihood



$\mathcal{L}(\mu) = e^{-\mu}\,\mu^n/n!$    and range of $\mu$ from   $\ln\{\mathcal{L}(\mu_{best})/\mathcal{L}(\mu)\} < 0.5$

For each $\mu_{true}$ calculate coverage $C(\mu_{true})$, and compare with nominal 68%

# Coverage

Fraction of intervals containing true value

Property of method, not of result

Can vary with param

Frequentist concept.  Built in to Neyman construction

Some Bayesians reject idea. Coverage not guaranteed

Integer data (Poisson) → discontinuities

Ideal coverage plot

C ─────────────────────── ←

→ μ

# 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 μ

Study coverage of different methods of Poisson parameter μ, from observation of number of events n

Hope for:



100%

$C(\mu)$

Nominal value

$\mu \longrightarrow$

# 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 (Neyman construction)

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

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



Coverage (C) vs $\mu$: $-2\ln\lambda < 1$    (C $\rightarrow$ 0.6827 as $\mu \rightarrow \infty$)

7

# Neyman central intervals, NEVER undercover

## (Conservative at both ends)



Coverage (C) vs $\mu$: Classical Central Intervals   (C $\rightarrow$ 0.6827 as $\mu \rightarrow \infty$)

28

# Feldman-Cousins Unified intervals

Neyman construction so NEVER undercovers



Coverage (C) vs $\mu$: Unified Intervals     $(C \rightarrow 0.6827$ as $\mu \rightarrow \infty)$

# Probability ordering



Coverage (C) vs $\mu$: Probability Ordering Intervals    (C $\to$ 0.6827 as $\mu \to \infty$)

$\chi^2 = (n-\mu)^2/\mu$      $\Delta \chi^2 = 0.1$ ⟶    24.8% coverage?

?

NOT Neyman :  Coverage = 0% → 100%



Coverage (C) vs $\mu$: $\chi^2 < 0.1$     (C → 0.2482 as $\mu$ → ∞)

# 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 $\mathcal{L}_{max}$ ⟹

Frequency

Bad          Good?     Great?

$\mathcal{L}_{max}$ ⟶

# Difference between $\mathcal{L}$ and pdf

Not necessarily:

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

fixed    vary

Contrast    pdf(data, params)

vary   fixed

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

pdf

$\mathcal{L}$

param

data

Max at t = 0

p

t

Max at $\lambda = 1/t$

$\mathcal{L}$

$\lambda$

33

Example 1

Fit exponential to times $t_1$, $t_2$, $t_3$ …….          [ Joel Heinrich, CDF 5639 ]

$$\mathcal{L} = \prod \lambda \exp(-\lambda t_i)$$

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

i.e. 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          same $\mathcal{L}_{max}$

pdf

34

t

Example 2

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

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

cos θ →

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 $\cos\theta < 0$

and $\mathcal{L}_{max}$ does not know about it.

**Example of general principle**

Example 3

Fit to Gaussian with variable μ, fixed σ

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

constant          ~variance(x)

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

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

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

Worse fit, larger $\mathcal{L}_{max}$                    Better fit, lower $\mathcal{L}_{max}$                    •36

# $\mathcal{L}_{max}$ 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!)

# Binned data and Goodness of Fit using $\mathcal{L}$-ratio

$n_i$

$\mu_i$

$x \rightarrow$

$$\mathcal{L} = \prod_i p_{ni}(\mu_i)$$

$$\mathcal{L}_{best} = \prod_i p_{ni}(\mu_{i,best})$$

$$= \prod_i p_{ni}(n_i)$$

$\ln[\mathcal{L}\text{-ratio}] = \ln[\mathcal{L}/\mathcal{L}_{best}]$

$\xrightarrow{\text{large } \mu_i} -0.5\chi^2$     i.e. Goodness of Fit

$\mathcal{L}_{best}$ is independent of parameters of fit,

and so same parameter values from $\mathcal{L}$ or $\mathcal{L}$-ratio

Baker and Cousins, NIM A221 (1984) 437

# $\mathcal{L}$ and pdf

## Example 1: Poisson

pdf = Probability density function for observing n, given μ

$$P(n;μ) = e^{-μ} μ^n/n!$$

From this, construct $\mathcal{L}$ as

$$\mathcal{L}(μ;n) = e^{-μ} μ^n/n!$$

i.e. use same function of μ and n, but

    for pdf, μ is fixed,   but

    for $\mathcal{L}$,    n is fixed

. . . . . . . . . . . pdf

μ

$\mathcal{L}$

n →

N.B. P(n;μ) exists only at integer non-negative n

    $\mathcal{L}(μ;n)$ exists only as continuous function of non-negative μ

39

Example 2     Lifetime distribution

pdf     $p(t;\lambda) = \lambda e^{-\lambda t}$

So     $\mathcal{L}(\lambda;t) = \lambda e^{-\lambda t}$     (single observed t)

Here both t and $\lambda$ are continuous

pdf maximises at $t = 0$

$\mathcal{L}$ maximises at $\lambda = t$

N.B. Functional form of $p(t)$ and $\mathcal{L}(\lambda)$ are different

Example 3:    Gaussian

$pdf(x;\mu) = \exp\{-(x-\mu)^2/2\sigma^2\} /(\sigma\sqrt{2\pi})$

$\mathcal{L}(\mu;x) \quad = \exp\{-(x-\mu)^2/2\sigma^2\} /(\sigma\sqrt{2\pi})$

N.B. In this case, same functional form for pdf and $\mathcal{L}$

So if you consider just Gaussians, can be confused between pdf and $\mathcal{L}$

So examples 1 and 2 are useful

# Transformation properties of pdf and $\mathcal{L}$

Lifetime example: $dn/dt = \lambda\, e^{-\lambda t}$

Change observable from t to y = $\sqrt{t}$

$$\frac{dn}{dy} = \frac{dn}{dt}\frac{dt}{dy} = 2y\lambda e^{-\lambda y^2}$$

So (a) pdf changes, BUT

(b) $$\int_{t_0}^{\infty} \frac{dn}{dt}\, dt = \int_{\sqrt{t_0}}^{\infty} \frac{dn}{dy}\, dy$$

i.e. corresponding integrals of pdf are
  INVARIANT

Now for $\mathcal{L}$ikelihood

When parameter changes from $\lambda$ to $\tau = 1/\lambda$

(a') $\mathcal{L}$ does not change

dn/dt = $(1/\tau)$ exp$\{-t/\tau\}$

and so $\mathcal{L}(\tau;t) = \mathcal{L}(\lambda=1/\tau;t)$

because identical numbers occur in evaluations of the two $\mathcal{L}$'s

BUT

(b')

So it is NOT meaningful to integrate $\mathcal{L}$

(However,.........)

| | pdf(t;λ) | $\mathcal{L}(\lambda;t)$ |
|---|---|---|
| Value of function | Changes when observable is transformed | INVARIANT wrt transformation of parameter |
| Integral of function | INVARIANT wrt transformation of observable | Changes when param is transformed |
| Conclusion | Max prob density not very sensible | Integrating $\mathcal{L}$ not very sensible |

CONCLUSION:

$$\int_{p_l}^{p_u} \mathcal{L}\, dp = \alpha$$ NOT recognised statistical procedure

[Metric dependent:

$\tau$ range agrees with $\tau_{pred}$

$\lambda$ range inconsistent with $1/\tau_{pred}$ ]

BUT

1) Could regard as "black box"

2) Make respectable by $\mathcal{L}$ $\Longrightarrow$ Bayes' posterior

Posterior($\lambda$) ~ $\mathcal{L}$($\lambda$)* Prior($\lambda$)    [and Prior($\lambda$) can be constant]

**6) BAYESIAN SMEARING OF $\alpha$**

"USE $\ln \mathcal{L}$ FOR $\hat{p}$ & $\sigma_p$

SMEAR IT TO INCORORATE

SYSTEMATIC UNCERTAINTIES"

SCENARIO:

$$n = \text{POISSON}\left(\mu = s\epsilon + b\right)$$

PARAM OF INTEREST ⟶

BACKGROUND

EFFIC/ACCEPTANCE//$\int\alpha$

UNCERTAINTIES

MEASURED IN 'SUBSIDIARY' EXPT

$$P(s, \epsilon \mid n) = \frac{P(n \mid s, \epsilon)\, \pi(s, \epsilon)}{\iint \cdots\cdots\cdots\cdots ds\, d\epsilon}$$

$\alpha$

$$P(s \mid n) = \int P(s, \epsilon \mid n)\, d\epsilon$$

$$= \frac{\int \alpha\, \pi(s)\, \pi(\epsilon)\, d\epsilon}{\iint \cdots\cdots\cdots\cdots ds\, d\epsilon}$$

e.g. $\pi(s) = $ truncated const. $\quad \pi(\epsilon) \sim e^{-\frac{1}{2}\left(\frac{\epsilon - \epsilon_0}{\sigma}\right)}$  BEWARE

i.e. SMEAR $\alpha$ (<u>not</u> $\ln\alpha$) by "prior" for $\epsilon$

46

# Getting $\mathcal{L}$ wrong: Punzi effect

Giovanni Punzi @ PHYSTAT2003
"Comments on $\mathcal{L}$ fits with variable resolution"

Separate two close signals, when resolution σ varies event
 by event, and is different for 2 signals
e.g. 1) Signal 1      $1+\cos^2\theta$
       Signal 2      Isotropic
       and different parts of detector give different σ

       2) M (or τ)
           Different numbers of tracks → different $\sigma_M$ (or $\sigma_\tau$)

Events characterised by $x_i$ and $\sigma_i$

A events centred on x = 0

B events centred on x = 1

$$\mathcal{L}(f)_{wrong} = \Pi\ [f * G(x_i,0,\sigma_i) + (1-f) * G(x_i,1,\sigma_i)]$$

$$\mathcal{L}(f)_{right} = \Pi\ [f*p(x_i,\sigma_i;A) + (1-f) * p(x_i,\sigma_i;B)]$$

$$p(S,T) = p(S|T) * p(T)$$

$$p(x_i,\sigma_i|A) = p(x_i|\sigma_i,A) * p(\sigma_i|A)$$

$$= G(x_i,0,\sigma_i) * p(\sigma_i|A)$$

So

$$\mathcal{L}(f)_{right} = \Pi[f * G(x_i,0,\sigma_i) * p(\sigma_i|A) + (1-f) * G(x_i,1,\sigma_i) * p(\sigma_i|B)]$$

If $p(\sigma|A) = p(\sigma|B)$, $\mathcal{L}_{right} = \mathcal{L}_{wrong}$

but NOT otherwise

Punzi's Monte Carlo for

A : $G(x, 0, \sigma_A)$

B : $G(x, 1, \sigma_B)$

$f_A = 1/3$

| $\sigma_A$ | $\sigma_B$ | $\mathcal{L}_{wrong}$ | | $\mathcal{L}_{right}$ | |
|---|---|---|---|---|---|
| | | $f_A$ | $\sigma_f$ | $f_A$ | $\sigma_f$ |
| 1·0 | 1·0 | 0·336(3) | 0·08 | Same | |
| 1·0 | 1·1 | 0·374(4) | 0·08 | 0·333(0) | 0 |
| 1·0 | 2·0 | 0·645(6) | 0·12 | 0·333(0) | 0 |
| 1 → 2 | 1.5 → 3 | 0·514(7) | 0·14 | 0·335(2) | 0·03 |
| 1.0 | 1 → 2 | 0.482(9) | 0.09 | 0.333(0) | 0 |

1) $\mathcal{L}_{wrong}$ OK for $p(\sigma_A) = p(\sigma_B)$ , but otherwise BIASSED

2) $\mathcal{L}_{right}$ unbiassed, but $\mathcal{L}_{wrong}$ biassed (enormously)!

3) $\mathcal{L}_{right}$ gives smaller $\sigma_f$ than $\mathcal{L}_{wrong}$

# Explanation of Punzi bias

$\sigma_A = 1$      $\sigma_B = 2$



A events with $\sigma = 1$

B events with $\sigma = 2$

x $\rightarrow$

x $\rightarrow$

ACTUAL DISTRIBUTION                  FITTING FUNCTION

[$N_A/N_B$ variable, but same for A and B events]

Fit gives upward bias for $N_A/N_B$ because  (i) that is much better for A events; and

(ii) it does not hurt too much for B events

50

# Another scenario for Punzi problem: PID

A    B                          π    K

M $\longrightarrow$                          TOF $\longrightarrow$

Originally:

Positions of peaks = constant         K-peak $\rightarrow$ π-peak at large momentum

$\sigma_i$ variable, $(\sigma_i)_A \neq (\sigma_i)_B$         $\sigma_i \sim$ constant, $p_K \neq p_\pi$

COMMON FEATURE: Separation/Error $\neq$ Constant

## Where else??

MORAL: Beware of event-by-event variables whose pdf's do not

appear in $\mathcal{L}$

# Avoiding Punzi Bias

BASIC RULE:
Write pdf for ALL observables, in terms of parameters

- Include $p(\sigma|A)$ and $p(\sigma|B)$ in fit
  (But then, for example, particle identification may be determined more by momentum distribution than by PID)

  OR

- Fit each range of $\sigma_i$ separately, and add $(N_A)_i \rightarrow (N_A)_{total}$, and similarly for B

  Incorrect method using $\mathcal{L}_{wrong}$ uses weighted average of $(f_A)_j$, assumed to be independent of j

Talk by Catastini at PHYSTAT05

# What else can we do with £s?

So far mainly parameter determination (also

Baker & Cousins' Goodness of Fit with Likelihood ratio)


Other possibilities:

Frequentist approach:

Construction of parameter confidence intervals

Likelihood ratios for comparing Hypotheses


Bayesian approach:

Together with priors $\rightarrow$ parameter credible intervals;

and Comparing Hypotheses


More in lectures by Olaf Behnke & Glen Cowan

# BAYES and FREQUENTISM
# The Return of an Old Controversy

# Parameter Determination

We need to make a statement about

<span style="color:red">Parameters, Given Data</span>

The basic difference between the two:

Bayesian :     <span style="color:red">Prob(parameter, given data)</span>
                  (an anathema to a Frequentist!)

Frequentist :  <span style="color:red">Prob(data, given parameter)</span>
                  (a likelihood function)

# WHAT IS PROBABILITY?

**MATHEMATICAL**

       **Formal**

       **Based on Axioms**

FREQUENTIST

       Ratio of frequencies as  n→ 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"

       Picture of Bayes

LEGAL PROBABILITY

# Picture of Reverend Bayes



Maybe it isn't Bayes?

"Probability that this is actually a picture of Bayes" is not Frequentist probability.

"Probability of Bayes" is Bayesian probability.

# Bayesian

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

Bayes' Theorem

p(param | data)  α  p(data | param) * p(param)

↑                        ↑                      ↑

posterior              likelihood            prior

Problems:   p(param)   Has particular value

"Degree of belief"

Prior      What functional form?

Coverage

P (Data;Theory) $\neq$ P (Theory;Data)


(Example of P(A;B) $\neq$ P(B;A) )

P (Data;Theory) $\neq$ P (Theory;Data)

Theory = male or female

Data = pregnant or not pregnant

P (pregnant ; female) ~ 3%

P (Data;Theory) $\neq$ P (Theory;Data)

Theory = male or female

Data = pregnant or not pregnant

P (pregnant ; female) ~ 3%

but

P (female ; pregnant) >>>3%

# Classical Approach: Neyman Construction

Neyman "confidence interval" avoids pdf for $\mu$

Uses only P( x; $\mu$ )

Confidence interval $\mu_1 \to \mu_2$ :

P( $\mu_1 \to \mu_2$ 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 $\mu_1 \to \mu_2$

# Classical (Neyman) Confidence Intervals

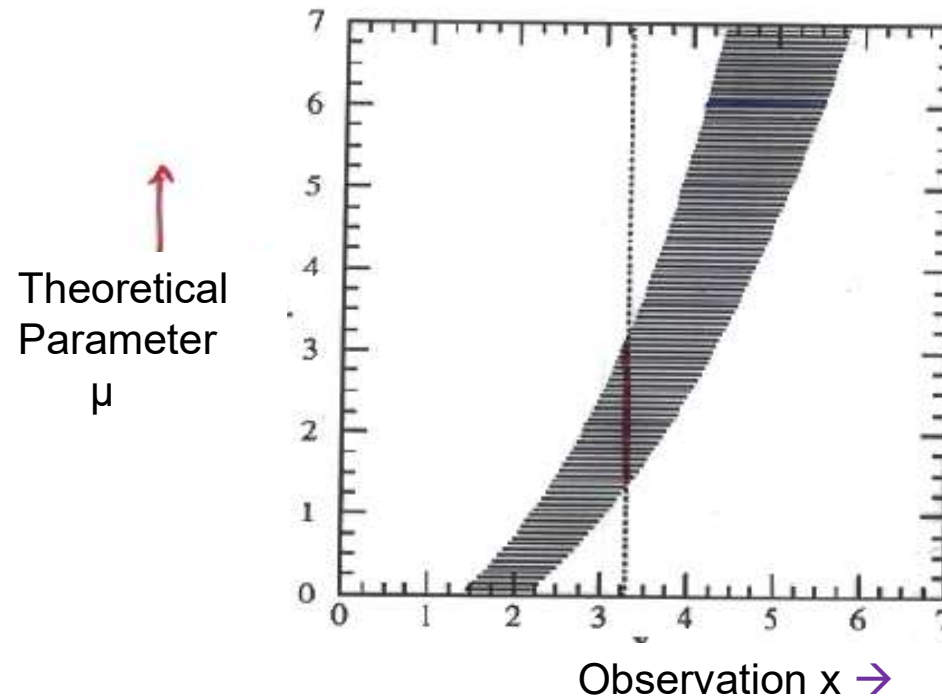## Uses only P(data|theory)

Theoretical
Parameter
μ

Observation x →

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.

μ≥0

## No prior for μ

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

**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 $\mu$

# Conclusions: What you now know

How it works, and how to estimate uncertainties

$\Delta(\ln \mathcal{L}) = 0.5$ rule and coverage

Several Parameters

Commbining Profile $\mathcal{L}$s loses information

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

Intro to Bayes and Frequentism

# FINAL MESSAGE

You cannot become an expert on Statistics by just reading books and listening to lectures.
You have to work at it – solve lots of problems, etc.

Best of luck with Statistics, and with your research and enjoy this School!