

Likelihoods

- 1) Introduction .
- 2) Do's & Dont's

Louis Lyons and Lorenzo Moneta
Imperial College & Oxford
CERN

CERN Academic Training Course

Nov 2016

Topics

(mainly Parameter Determination)

What it is

How it works: Resonance

Uncertainty estimates

Detailed example: Lifetime

Several Parameters

Extended maximum \mathcal{L}

Do's and Dont's with \mathcal{L}

Simple example: Angular distribution

Start with pdf = Prob density fn for data, given param values:

$$y = N (1 + \beta \cos^2\theta)$$

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

= probability density of observing θ_i , given β

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

= probability density of observing the data set y_i , given β

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

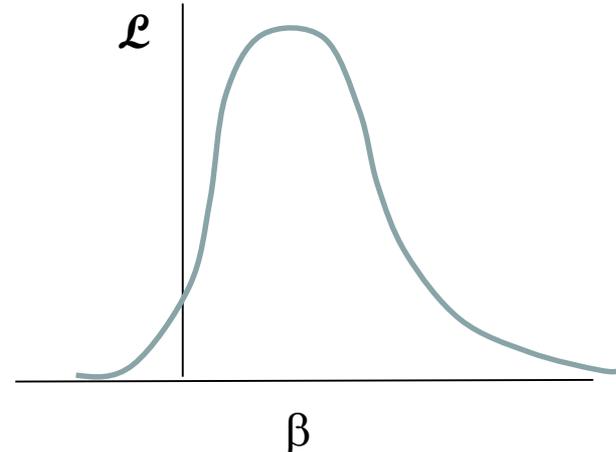
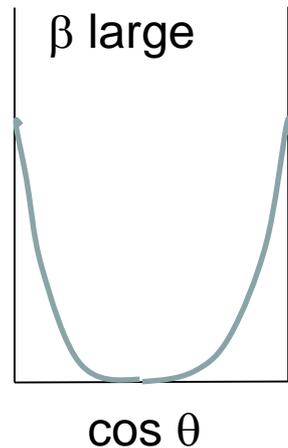
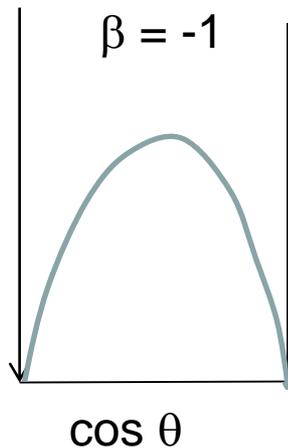
Values of β for which \mathcal{L} is very small are ruled out

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

CRUCIAL to normalise y

$$N = 1/\{2(1 + \beta/3)\}$$

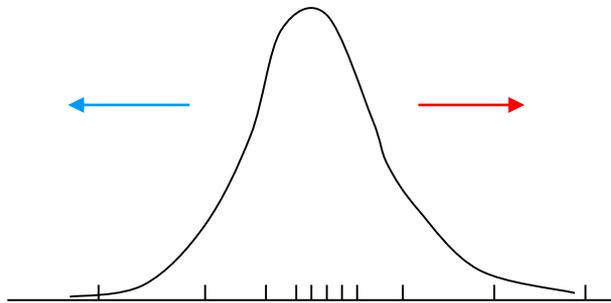
(Information about parameter β comes from **shape** of exptl distribution of $\cos\theta$)



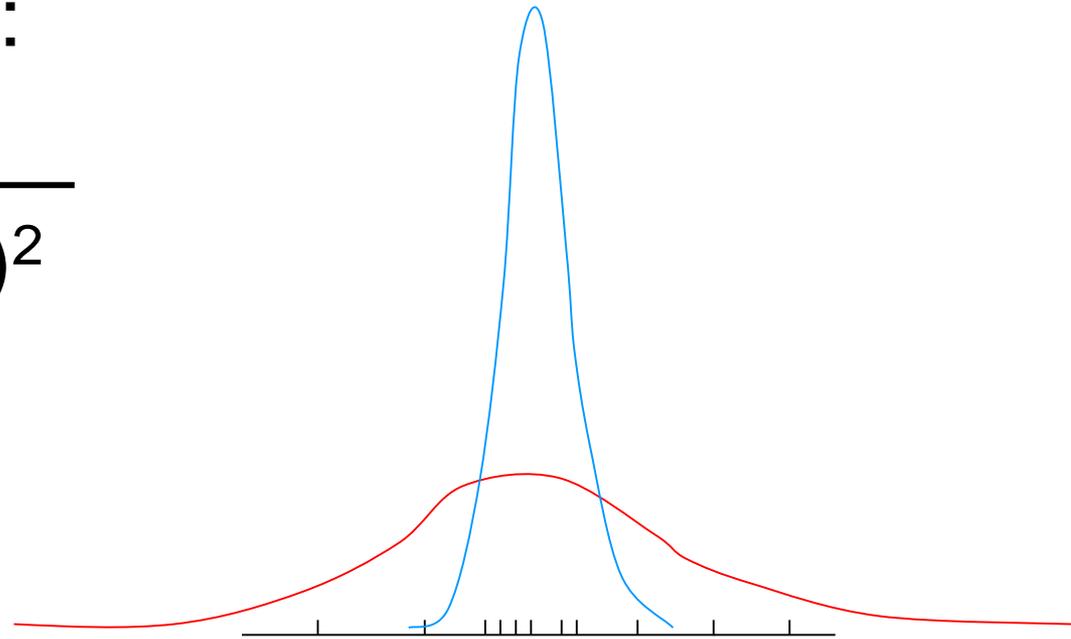
How it works: Resonance

First write down pdf:

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



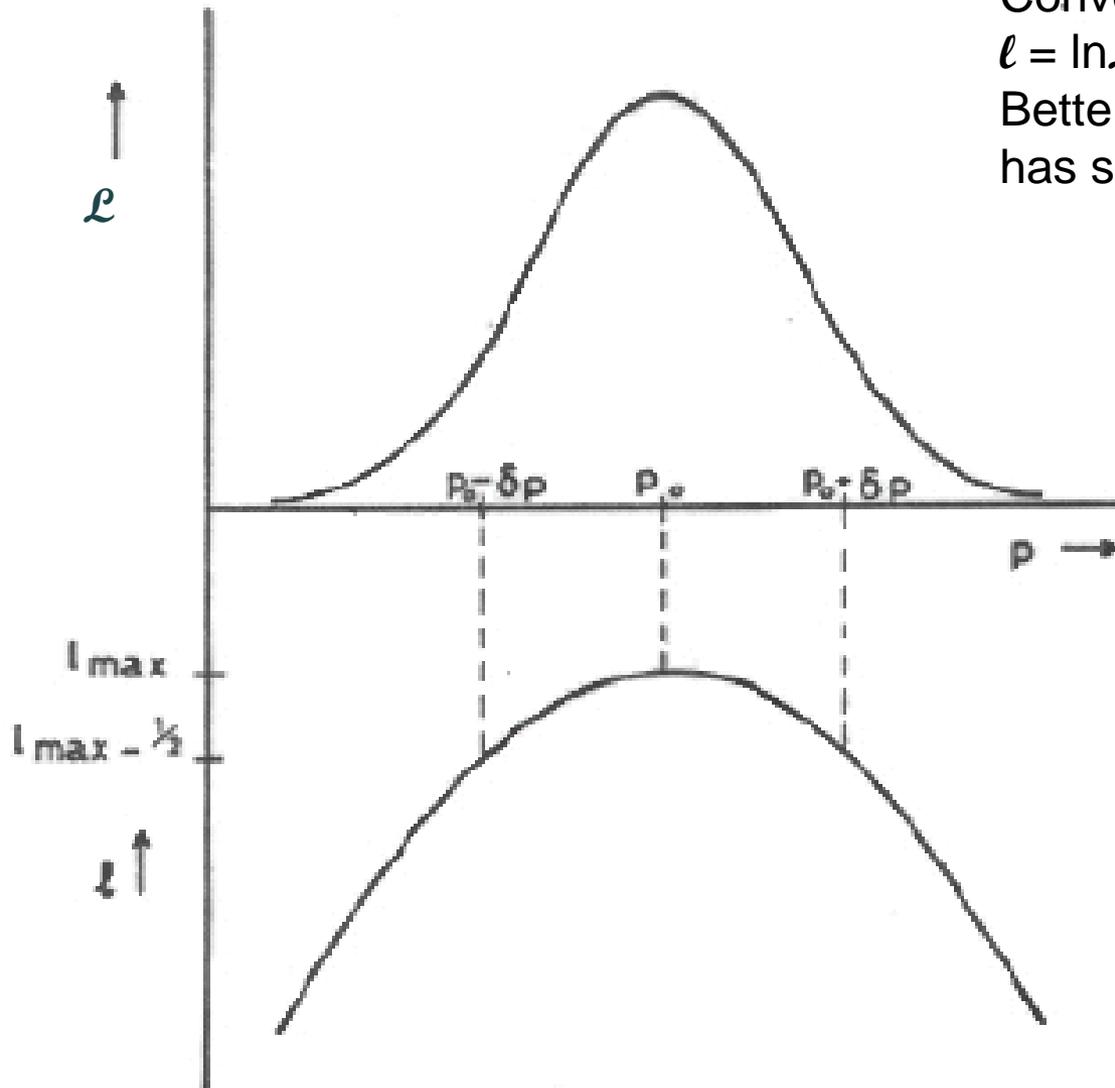
Vary M_0



Vary Γ

**N.B. Can make use
of individual events**

Conventional to consider
 $l = \ln \mathcal{L} = \sum \ln y_i$
Better numerically, and
has some nice properties



Maximum likelihood uncertainty

Range of likely values of param μ from width of \mathcal{L} or 1 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”~~

Return to ‘Coverage’ later

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

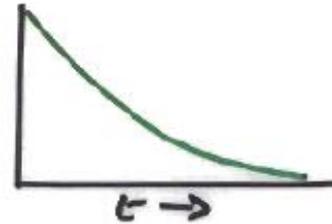
e.g $1/p$ rather than p ; τ or λ

Lifetime Determination

Realistic analyses are more complicated than this

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

↑ NORMALISATION



Observe t_1, t_2, \dots, t_N

Use pdf to construct

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

$$\therefore \mathcal{L} = \sum (-t_i/\tau - \ln \tau)$$

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

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

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

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

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

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

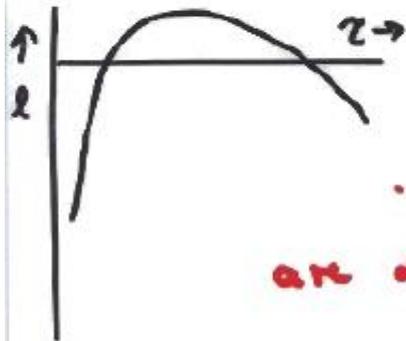
BEWARE FOR AVERAGING RESULTS

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

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

$$l(\tau) - l(\tau_{max}) = -N \tau_{max}/\tau - N \ln \tau$$

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



\therefore For given N , σ_+ & σ_-

are defined ($\sim \frac{\tau_{max}}{\sqrt{N}}$ as $N \rightarrow \infty$)

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

— " —

$$l(\tau_{max}) = -N(1 + \ln \bar{E}) \quad \star \star \star$$

N.B. $l(\tau_{max})$ depends only on \bar{E} ,

but not on distribution of t_i

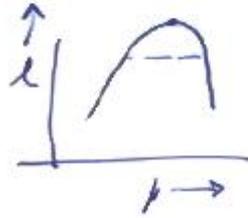
Relevant for whether l_{max} is useful
for testing goodness of fit

Several Parameters

1 param β

$$\beta \text{ from } \frac{\partial \mathcal{L}}{\partial \beta} = 0$$

$$\sigma_{\beta}^2 = 1 / \left(- \frac{\partial^2 \mathcal{L}}{\partial \beta^2} \right)$$

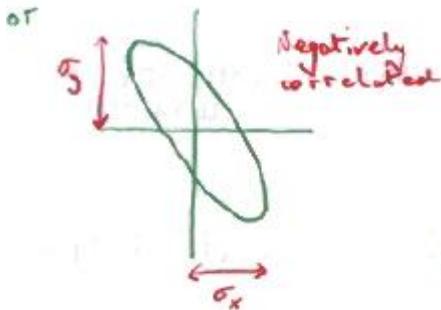
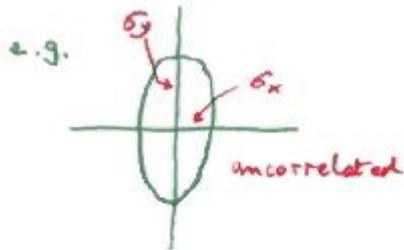


Many dimensions : $\mathcal{L}(\beta_1, \beta_2, \beta_3, \dots)$

$$\beta_1, \beta_2, \beta_3, \dots \text{ from } \frac{\partial \mathcal{L}}{\partial \beta_i} = 0$$

For errors, define $H_{ij} = - \frac{\partial^2 \mathcal{L}}{\partial \beta_i \partial \beta_j} = \text{Inverse Error Matrix}$

$$\text{Error matrix } E_{ij} = (H^{-1})_{ij}$$



N.B. ERROR NOT GIVEN BY

$\mathcal{L} = \mathcal{L}_{\max} - \frac{1}{2}$ WHEN VARYING x
FROM BEST VALUE WHILE
KEEPING y, \dots CONSTANT

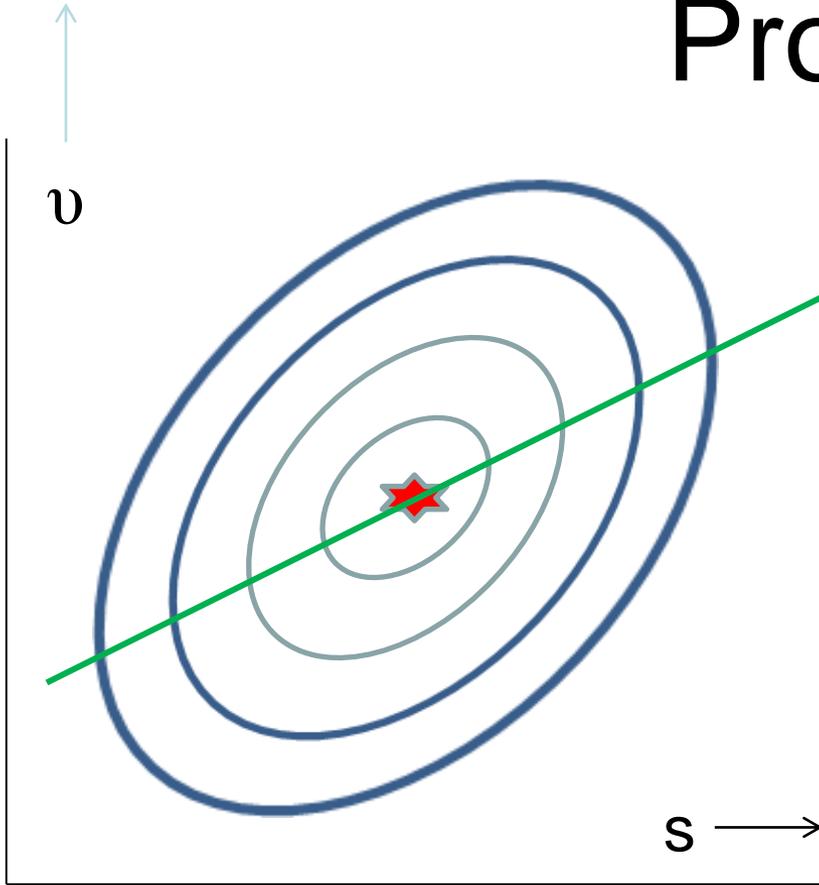
ERROR IS GIVEN BY

$\mathcal{L} = \mathcal{L}_{\max} - \frac{1}{2}$ WHEN VARYING x
FROM BEST VALUE WHILE \dots

PROFILE \mathcal{L}

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

Profile \mathcal{L}



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

s = physics param

v = nuisance param

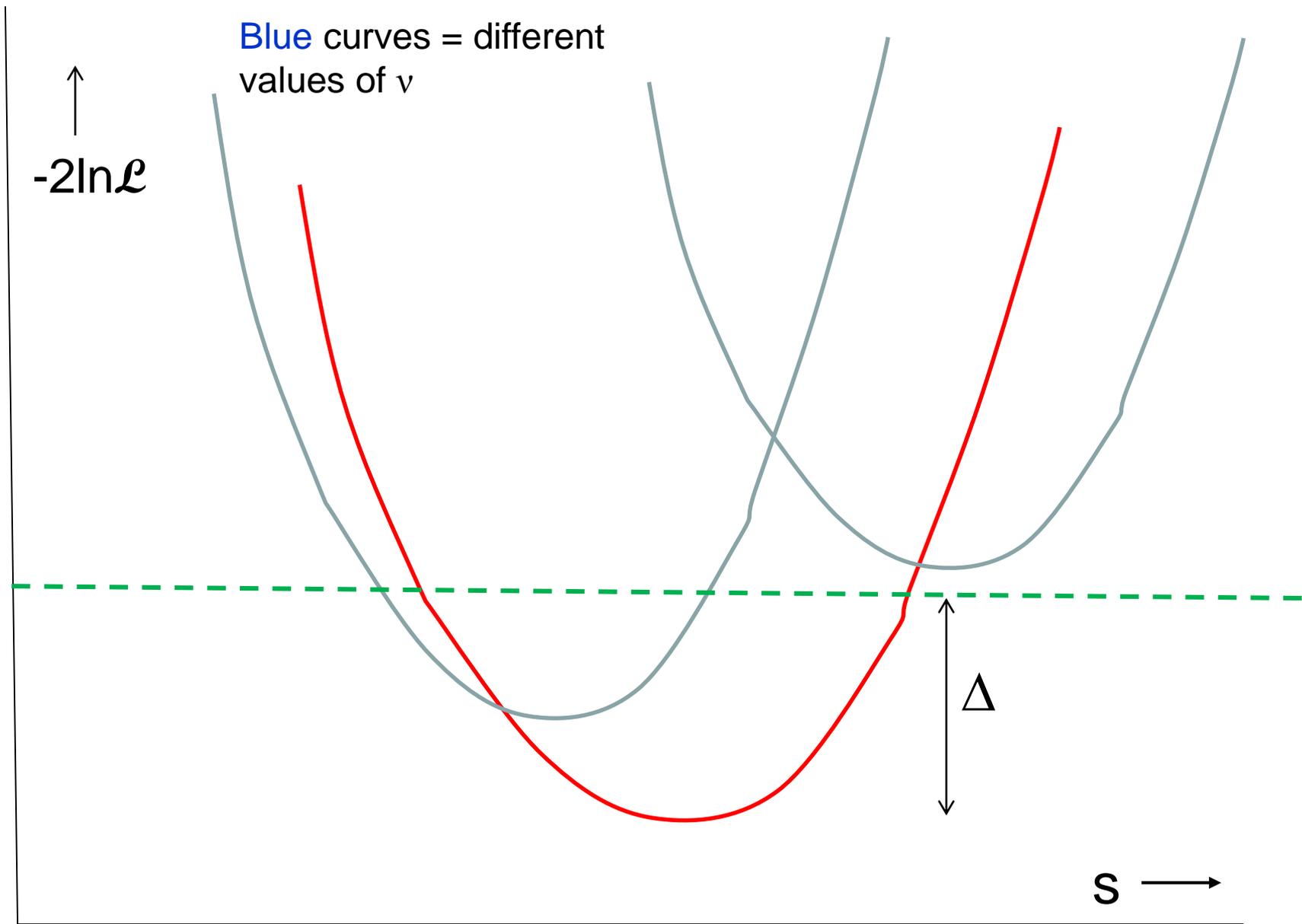
A method for dealing with systematics

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



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 1: 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

ML and EML

ML uses fixed (data) normalisation

EML has normalisation as parameter

Example 2: Cosmic ray experiment

See 96 protons and 4 heavy nuclei

ML estimate $96 \pm 2\%$ protons $4 \pm 2\%$ heavy nuclei

EML estimate 96 ± 10 protons 4 ± 2 heavy nuclei

Example 3: Decay of resonance

Use ML for Branching Ratios

Use EML for Partial Decay Rates

Relation between Poisson and Binomial

N people in lecture, m males and f females (N = m + f)

Assume these are representative of basic rates: ν people νp males $\nu(1-p)$ females

Probability of observing N people = $P_{\text{Poisson}} = e^{-\nu} \nu^N / N!$

Prob of given male/female division = $P_{\text{Binom}} = \frac{N!}{m!f!} p^m (1-p)^f$

Prob of N people, m male and f female = $P_{\text{Poisson}} P_{\text{Binom}}$

$$= \frac{e^{-\nu p} \nu^m p^m}{m!} * \frac{e^{-\nu(1-p)} \nu^f (1-p)^f}{f!}$$

= Poisson prob for males * Poisson prob for females

People	Male	Female
Patients	Cured	Remain ill
Decaying nuclei	Forwards	Backwards
Cosmic rays	Protons	Other particles

	Moments	Max Like	Least squares
Easy?	Yes, if...	Normalisation, maximisation messy	Minimisation
Efficient?	Not very	Usually best	Sometimes = Max Like
Input	Separate events	Separate events	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
Inverse Covariance Matrix	Observed spread, or analytic	$\left\{ -\frac{\partial^2 \ell}{\partial p_i \partial p_j} \right\}$	$\left\{ \frac{\partial^2 S}{2 \partial p_i \partial p_j} \right\}$
Main feature	Easy	Best	Goodness of Fit

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 
- $\int_{p_L}^{p_U} \mathcal{L} dp = 0.90$
- BAYESIAN SMEARING OF \mathcal{L}
- USE CORRECT \mathcal{L} (PUNZI EFFECT)

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

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

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

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

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

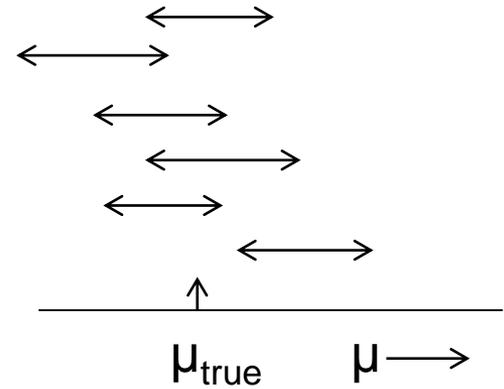
If $\mathcal{L}(\mu)$ is non-Gaussian, these are no longer the same

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

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

Barlow: Phystat05

COVERAGE



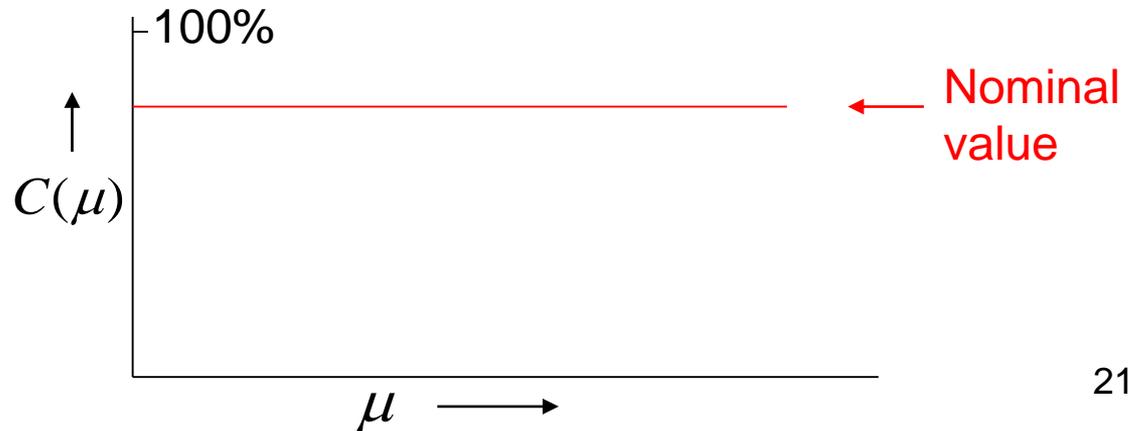
How often does quoted range for parameter include param's true value?

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

Coverage can vary with μ

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

Hope for:



COVERAGE

If true for all μ : “correct coverage”

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

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

Conservative

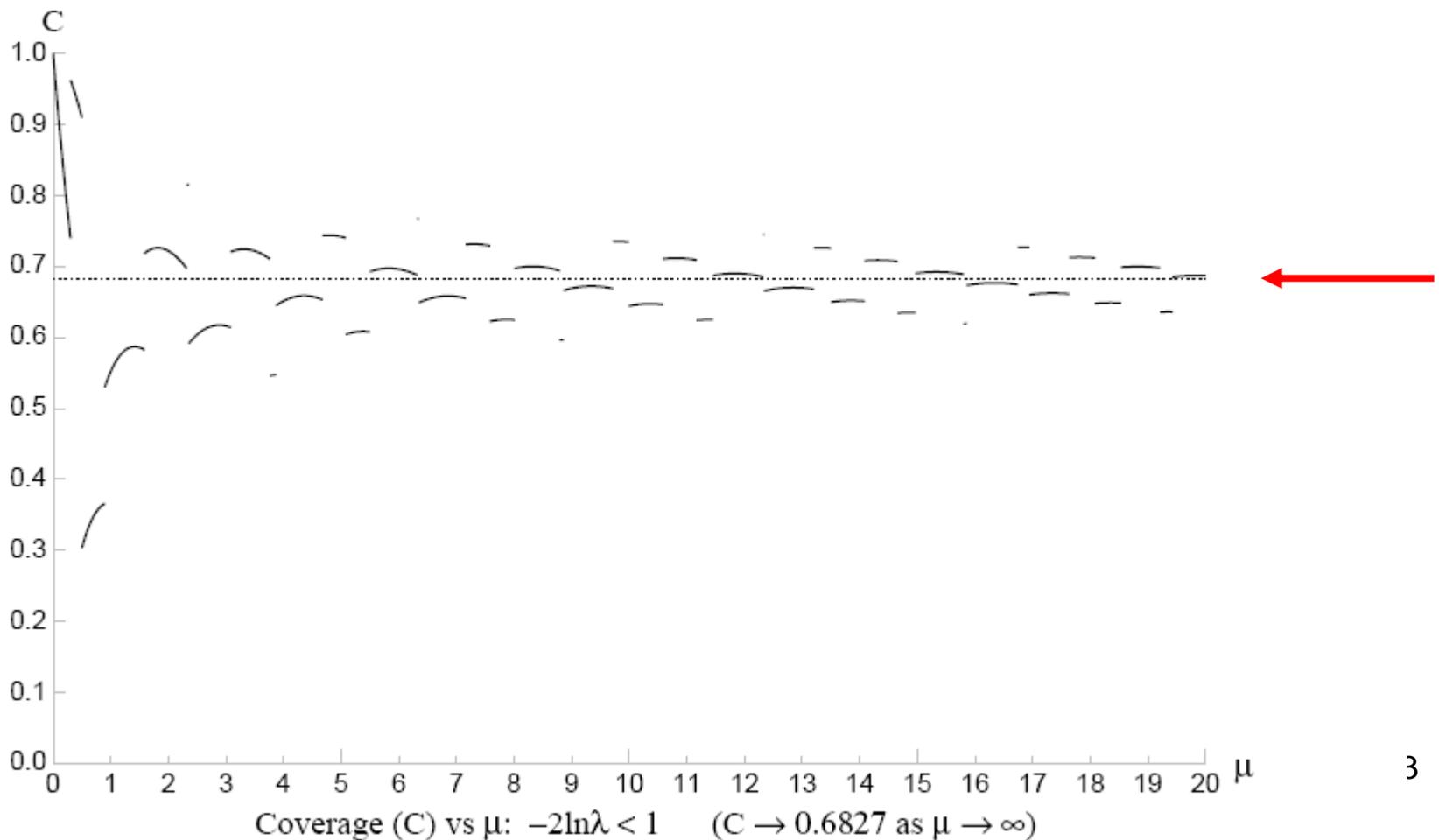
Loss of rejection
power

Some Bayesians regard
Coverage as irrelevant

Coverage : \mathcal{L} approach (Not Neyman construction)

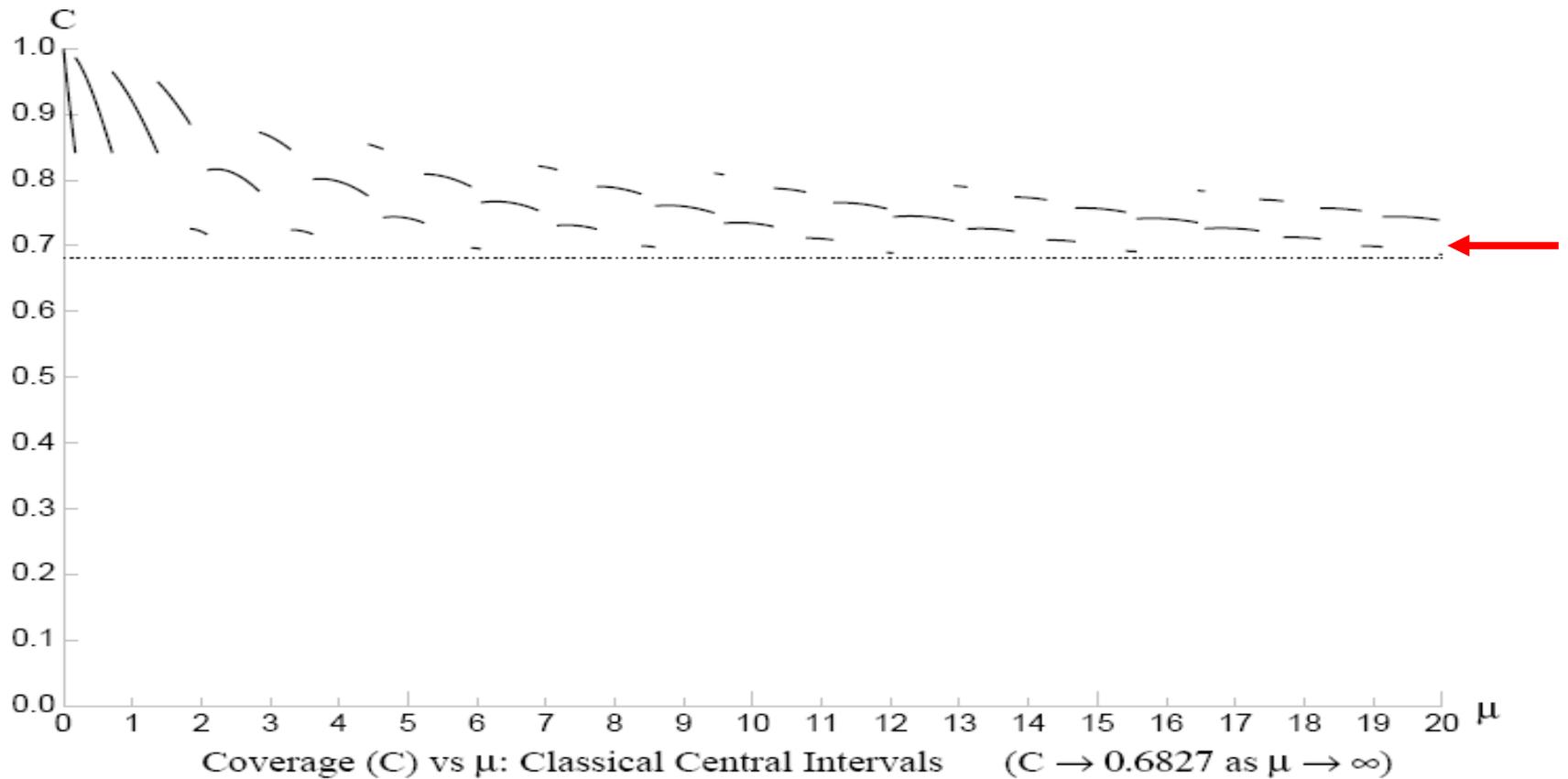
$$P(n, \mu) = e^{-\mu} \mu^n / n! \quad (\text{Joel Heinrich CDF note 6438})$$

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



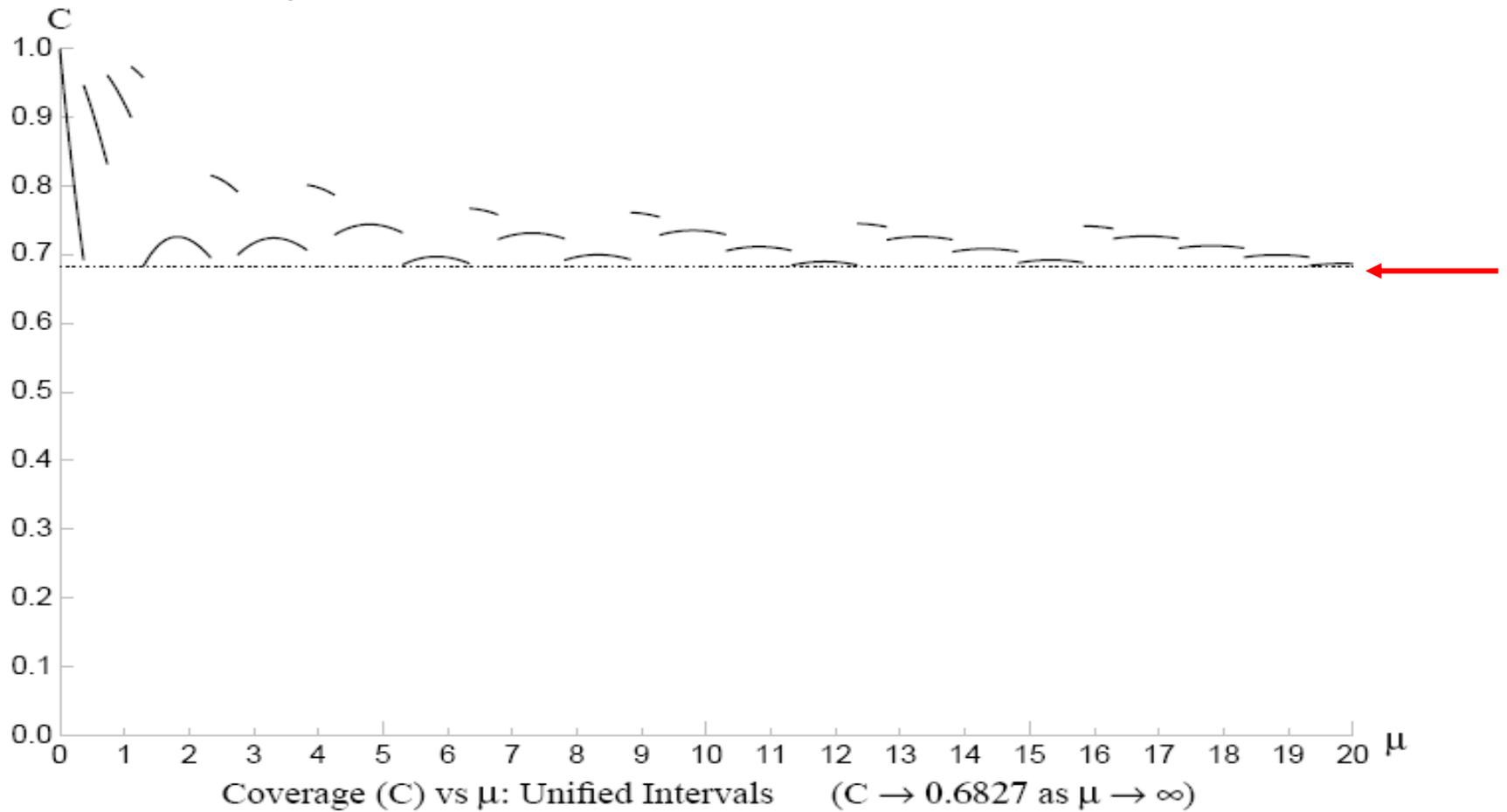
Neyman central intervals, NEVER undercover

(Conservative at both ends)



Feldman-Cousins Unified intervals

Neyman construction so NEVER undercovers



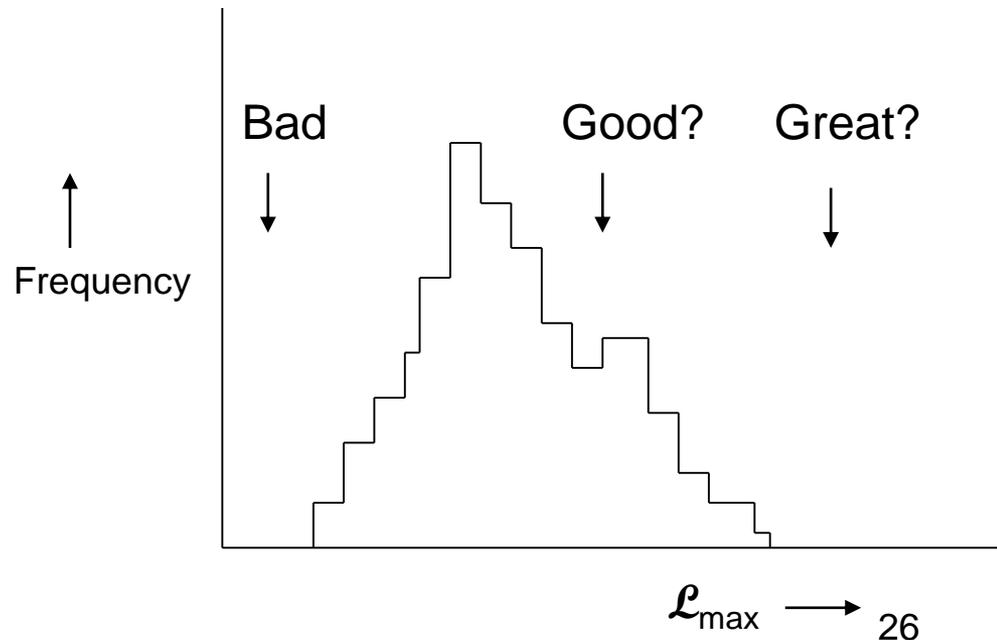
Unbinned \mathcal{L}_{\max} and Goodness of Fit?

Find params by maximising \mathcal{L}

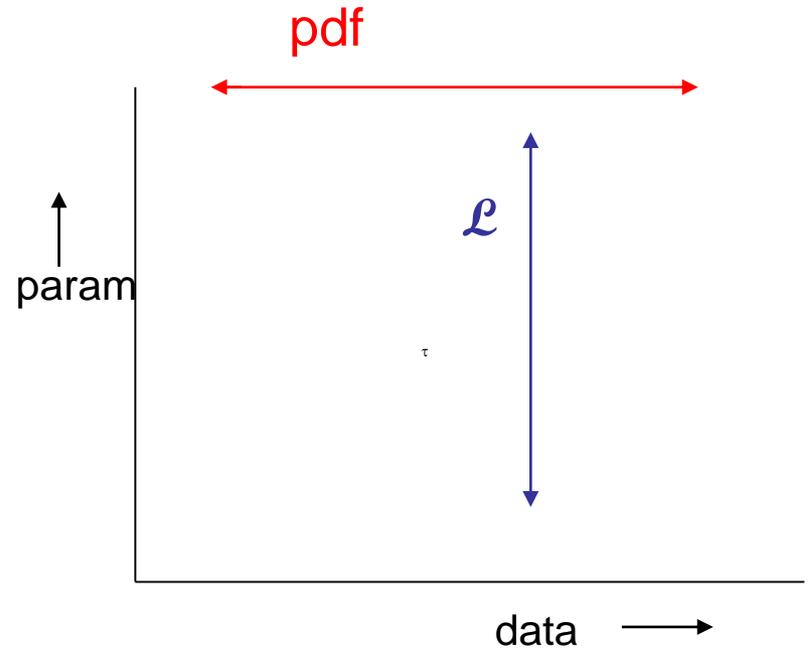
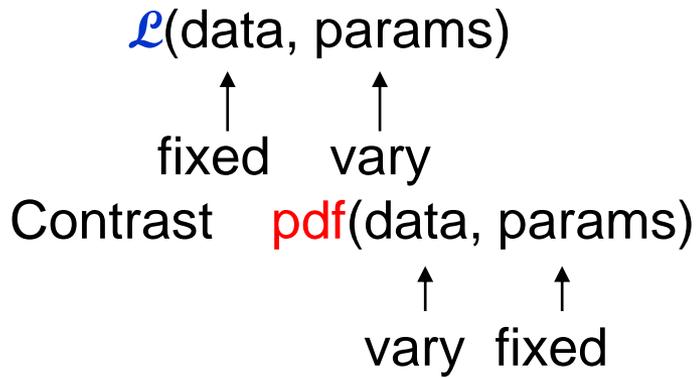
So larger \mathcal{L} better than smaller \mathcal{L}

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

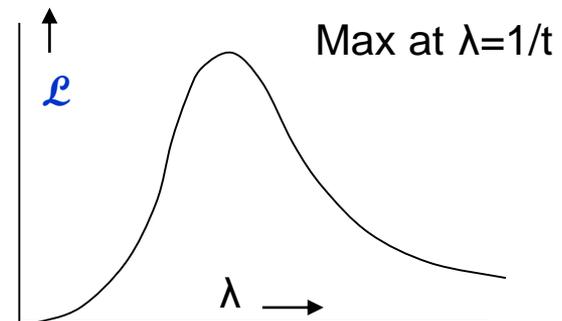
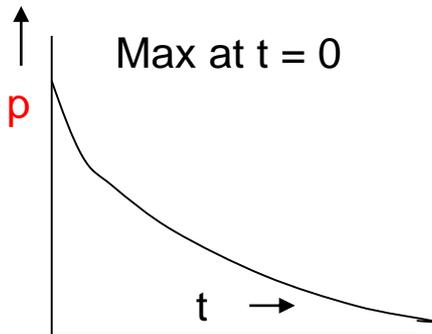
Monte Carlo distribution
of unbinned \mathcal{L}_{\max} 



Not necessarily:



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



Example 1

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

[Joel Heinrich, CDF 5639]

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

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

i.e. Depends only on AVERAGE t , but is

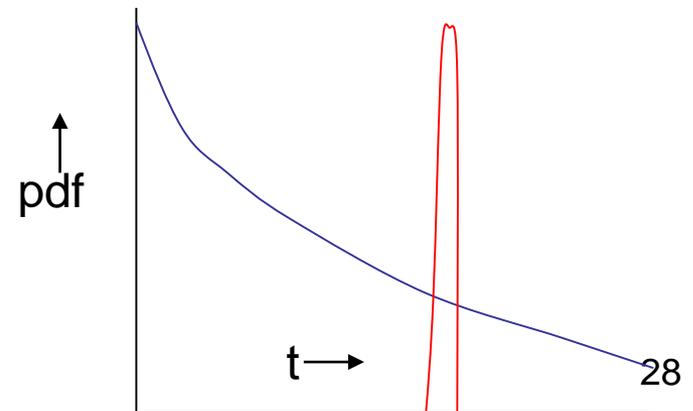
INDEPENDENT OF DISTRIBUTION OF t (except for.....)

(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

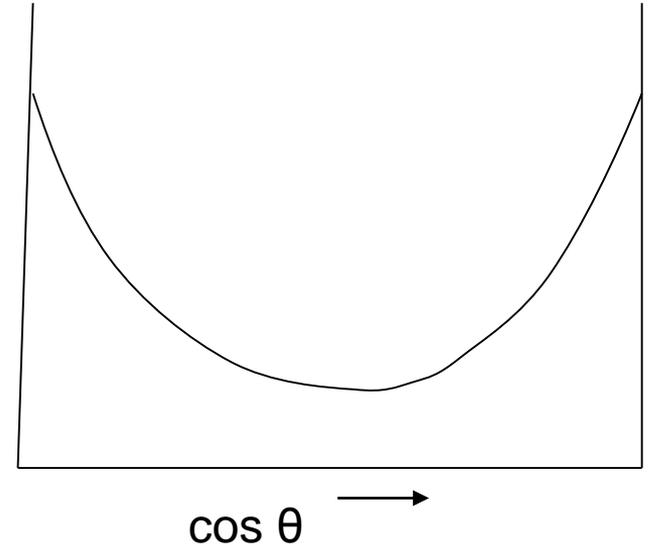
Same average t \longrightarrow same \mathcal{L}_{\max}



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 $\cos \theta < 0$

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

Example of general principle

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

Conclusion:

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

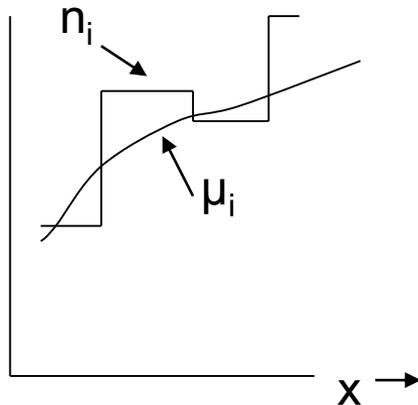
NOT with respect to data

\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



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

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

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

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

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

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

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

Conclusions

How it works, and how to estimate uncertainties

Likelihood or Extended Likelihood

Several Parameters

Likelihood does not guarantee coverage

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

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 \rightarrow different σ_M (or σ_τ)

Events characterised by x_i and σ_i

A events centred on $x = 0$

B events centred on $x = 1$

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

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

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

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

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

So

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

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

but NOT otherwise

Punzi's Monte Carlo for

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

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

$$f_A = 1/3$$

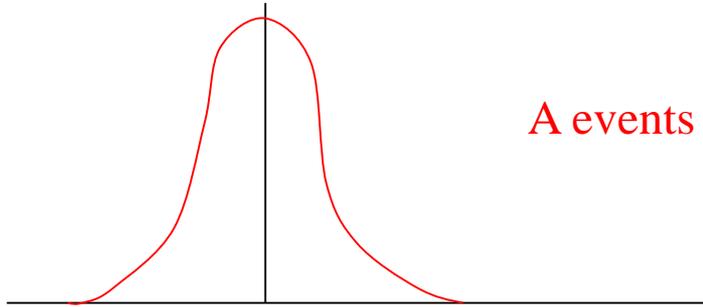
σ_A	σ_B	$\mathcal{L}_{\text{wrong}}$		$\mathcal{L}_{\text{right}}$	
		f_A	σ_f	f_A	σ_f
1.0	1.0	0.336(3)	0.08	Same	
1.0	1.1	0.374(4)	0.08	0.333(0)	0
1.0	2.0	0.645(6)	0.12	0.333(0)	0
1 → 2	1.5 → 3	0.514(7)	0.14	0.335(2)	0.03
1.0	1 → 2	0.482(9)	0.09	0.333(0)	0

- 1) $\mathcal{L}_{\text{wrong}}$ OK for $p(\sigma_A) = p(\sigma_B)$, but otherwise BIASED
- 2) $\mathcal{L}_{\text{right}}$ unbiased, but $\mathcal{L}_{\text{wrong}}$ biased (enormously)!
- 3) $\mathcal{L}_{\text{right}}$ gives smaller σ_f than $\mathcal{L}_{\text{wrong}}$

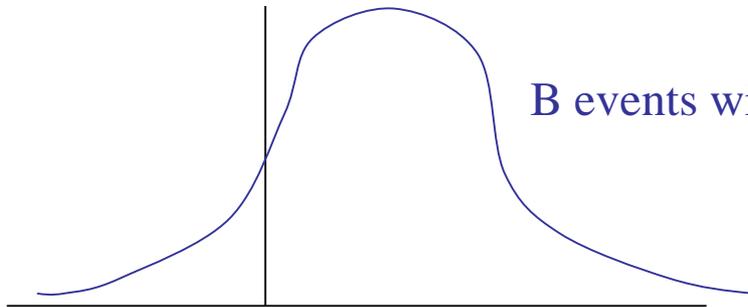
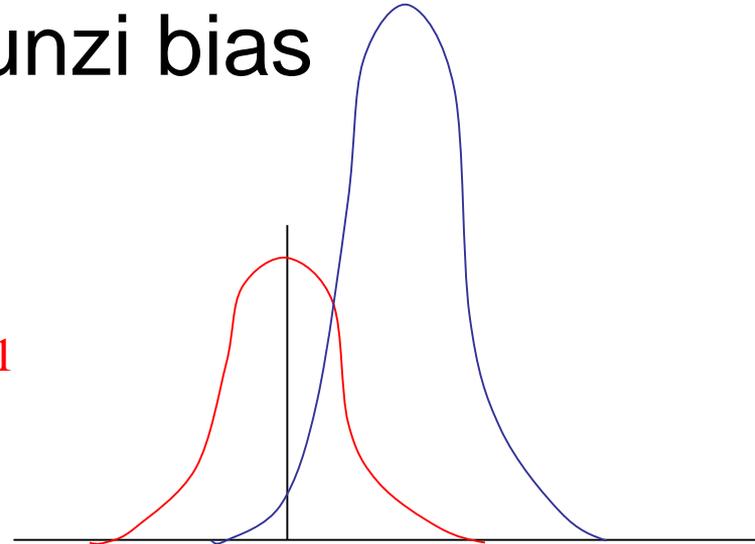
Explanation of Punzi bias

$\sigma_A = 1$

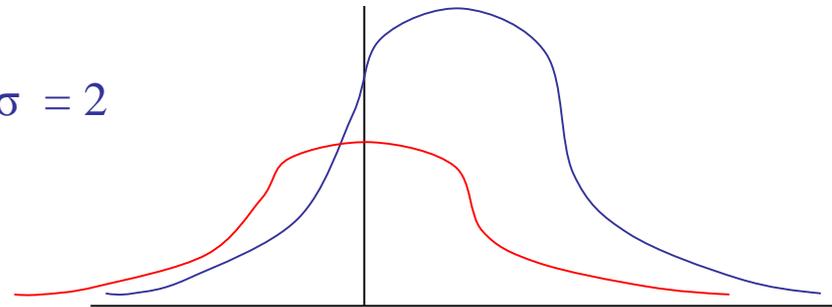
$\sigma_B = 2$



A events with $\sigma = 1$



B events with $\sigma = 2$



x →

x →

ACTUAL DISTRIBUTION

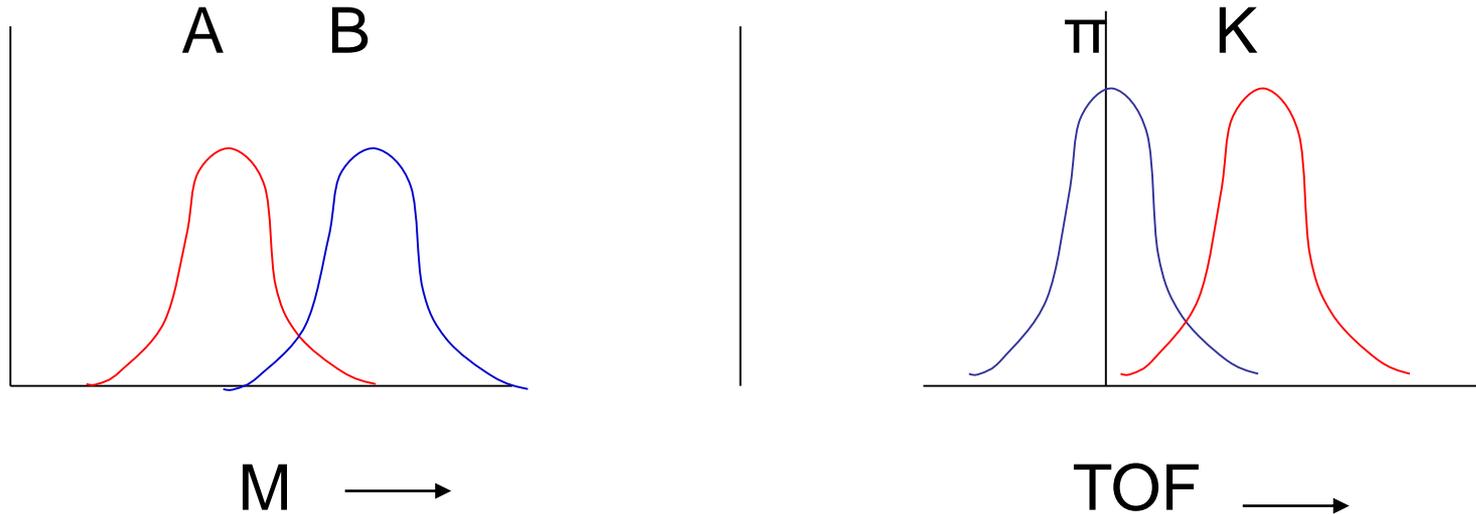
FITTING FUNCTION

[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

Another scenario for Punzi problem: PID



Originally:

Positions of peaks = constant

σ_i variable, $(\sigma_i)_A \neq (\sigma_i)_B$

COMMON FEATURE: Separation/Error \neq Constant

K-peak \rightarrow π -peak at large momentum

$\sigma_i \sim$ constant, $p_K \neq p_\pi$

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 σ_i separately, and add $(N_A)_i \rightarrow (N_A)_{\text{total}}$, and similarly for B

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