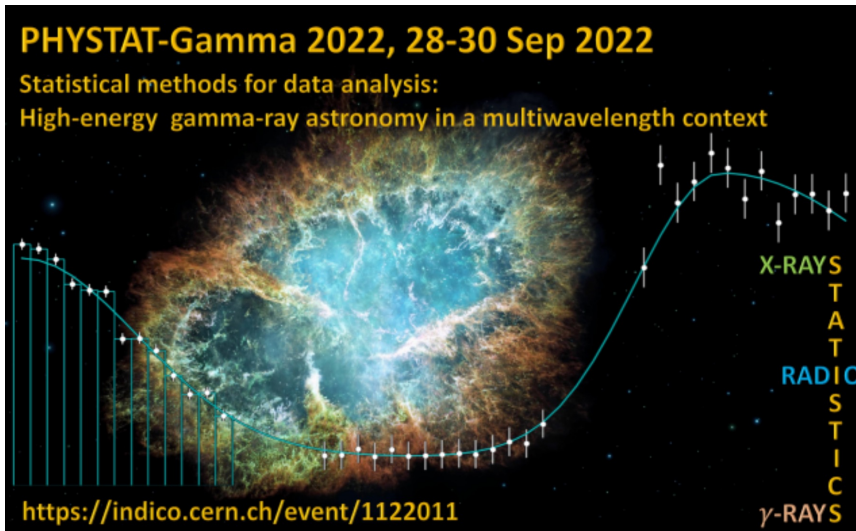


# Statistics 101 – a very fast introduction

## Lecture 1



PHYSTAT-Gamma  
(online)

26-30 September 2022

<https://indico.cern.ch/event/1122011/>



Glen Cowan

Physics Department

Royal Holloway, University of London

[g.cowan@rhul.ac.uk](mailto:g.cowan@rhul.ac.uk)

[www.pp.rhul.ac.uk/~cowan](http://www.pp.rhul.ac.uk/~cowan)

# Outline

- **Lecture 1:** Probability and Bayes theorem, Frequentist and Bayesian statistics, likelihood function, parameter estimation, maximum likelihood, information inequality, properties of MLE
- Lecture 2:** Frequentist hypothesis tests, Neyman-Pearson lemma/likelihood ratio, goodness of fit, p-values and significances, confidence interval from a test, Wilk's theorem and confidence regions

Almost everything is a subset of the University of London course:

[http://www.pp.rhul.ac.uk/~cowan/stat\\_course.html](http://www.pp.rhul.ac.uk/~cowan/stat_course.html)

# A quick review of probability

Frequentist ( $A$  = outcome of repeatable observation)

$$P(A) = \lim_{n \rightarrow \infty} \frac{\text{outcome is in } A}{n}$$

Subjective ( $A$  = hypothesis)

$P(A)$  = degree of belief that  $A$  is true

Conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

E.g. rolling a die,  
outcome  $n = 1, 2, \dots, 6$ :

$$P(n \leq 3 | n \text{ even}) = \frac{P((n \leq 3) \cap n \text{ even})}{P(n \text{ even})} = \frac{1/6}{3/6} = \frac{1}{3}$$

$A$  and  $B$  are independent iff:

$$P(A \cap B) = P(A)P(B)$$

I.e. if  $A, B$  independent, then

$$P(A|B) = \frac{P(A)P(B)}{P(B)} = P(A)$$

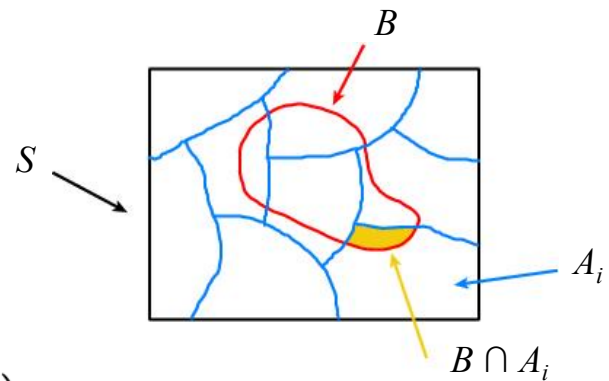
# Bayes' theorem

Use definition of conditional probability and  $P(A \cap B) = P(B \cap A)$

$$\rightarrow P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (\text{Bayes' theorem})$$

If set of all outcomes  $S = \bigcup_i A_i$  with  $A_i$  disjoint, then law of total probability for  $P(B)$  says

$$P(B) = \sum_i P(B \cap A_i) = \sum_i P(B|A_i)P(A_i)$$



so that Bayes' theorem becomes  $P(A|B) = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)}$

Bayes' theorem holds regardless of how probability is interpreted (frequency, degree of belief...).

# Frequentist Statistics – general philosophy

In frequentist statistics, probabilities are associated only with the data, i.e., outcomes of repeatable observations (shorthand:  $x$ ).

Probability = limiting frequency

Probabilities such as

$P$  (string theory is true),

$P(0.117 < \alpha_s < 0.119)$ ,

$P$  (Biden wins in 2024),

etc. are either 0 or 1, but we don't know which.

The tools of frequentist statistics tell us what to expect, under the assumption of certain probabilities, about hypothetical repeated observations.

Preferred theories (models, hypotheses, ...) are those that predict a high probability for data “like” the data observed.

# Bayesian Statistics – general philosophy

In Bayesian statistics, use subjective probability for hypotheses:

probability of the data assuming  
hypothesis  $H$  (the likelihood)

prior probability, i.e.,  
before seeing the data

$$P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) dH}$$

posterior probability, i.e.,  
after seeing the data

normalization involves sum  
over all possible hypotheses

Bayes' theorem has an “if-then” character: **If** your prior probabilities were  $\pi(H)$ , **then** it says how these probabilities should change in the light of the data.

No general prescription for priors (subjective!)

# Hypothesis, likelihood

Suppose the entire result of an experiment (set of measurements) is a collection of numbers  $\mathbf{x}$ .

A (simple) hypothesis is a rule that assigns a probability to each possible data outcome:

$$P(\mathbf{x}|H) = \text{the likelihood of } H$$

Often we deal with a family of hypotheses labeled by one or more undetermined parameters (a composite hypothesis):

$$P(\mathbf{x}|\boldsymbol{\theta}) = L(\boldsymbol{\theta}) = \text{the “likelihood function”}$$

Note:

- 1) For the likelihood we treat the data  $\mathbf{x}$  as fixed.
- 2) The likelihood function  $L(\boldsymbol{\theta})$  is not a pdf for  $\boldsymbol{\theta}$ .

# The likelihood function for i.i.d.\* data

\* i.i.d. = independent and identically distributed

Consider  $n$  independent observations of  $x$ :  $x_1, \dots, x_n$ , where  $x$  follows  $f(x; \theta)$ . The joint pdf for the whole data sample is:

$$f(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

In this case the likelihood function is

$$L(\vec{\theta}) = \prod_{i=1}^n f(x_i; \vec{\theta}) \quad (x_i \text{ constant})$$



# Parameter estimation

The parameters of a pdf are any constants that characterize it,

$$f(x; \theta) = \frac{1}{\theta} e^{-x/\theta}$$

r.v.                      parameter

i.e.,  $\theta$  indexes a set of hypotheses.

Suppose we have a sample of observed values:  $\mathbf{x} = (x_1, \dots, x_n)$

We want to find some function of the data to estimate the parameter(s):

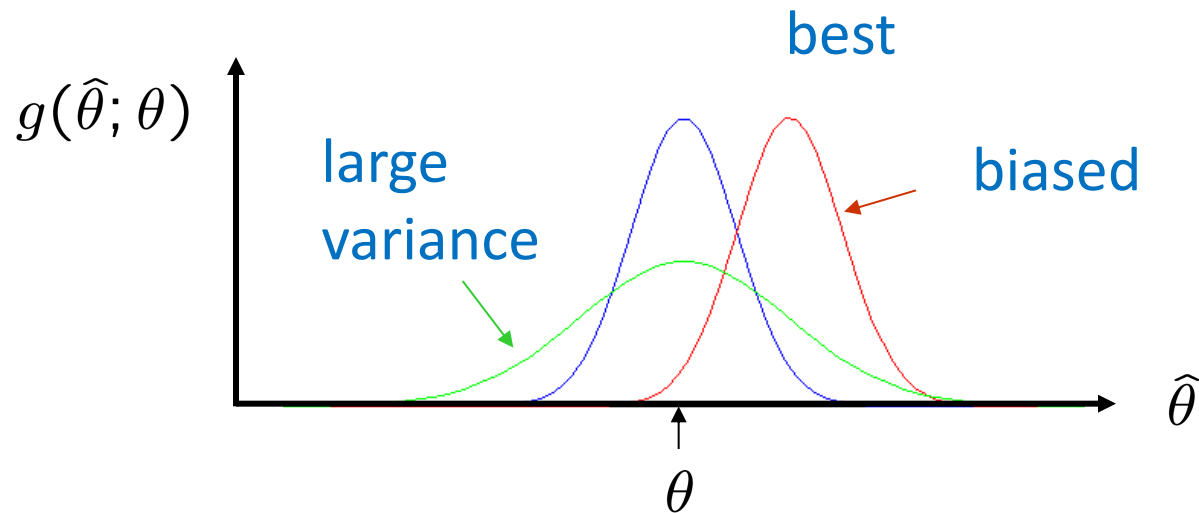
$$\hat{\theta}(\vec{x})$$

← estimator written with a hat

Sometimes we say ‘estimator’ for the function of  $x_1, \dots, x_n$ ; ‘estimate’ for the value of the estimator with a particular data set.

# Properties of estimators

If we were to repeat the entire measurement, the estimates from each would follow a pdf:



We want small (or zero) bias (systematic error):  $b = E[\hat{\theta}] - \theta$

→ average of repeated measurements should tend to true value.

And we want a small variance (statistical error):  $V[\hat{\theta}]$

→ small bias & variance are in general conflicting criteria

# Maximum Likelihood Estimators (MLEs)

We *define* the maximum likelihood estimators or MLEs to be the parameter values for which the likelihood is maximum.

Maximizing  $L$   
equivalent to  
maximizing  $\log L$

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$



Could have multiple maxima (take highest).

MLEs not guaranteed to have any ‘optimal’ properties, (but in practice they’re very good).

# MLE example: parameter of exponential pdf

Consider exponential pdf,  $f(t; \tau) = \frac{1}{\tau} e^{-t/\tau}$

and suppose we have i.i.d. data,  $t_1, \dots, t_n$

The likelihood function is  $L(\tau) = \prod_{i=1}^n \frac{1}{\tau} e^{-t_i/\tau}$

The value of  $\tau$  for which  $L(\tau)$  is maximum also gives the maximum value of its logarithm (the log-likelihood function):

$$\ln L(\tau) = \sum_{i=1}^n \ln f(t_i; \tau) = \sum_{i=1}^n \left( \ln \frac{1}{\tau} - \frac{t_i}{\tau} \right)$$

# MLE example: parameter of exponential pdf (2)

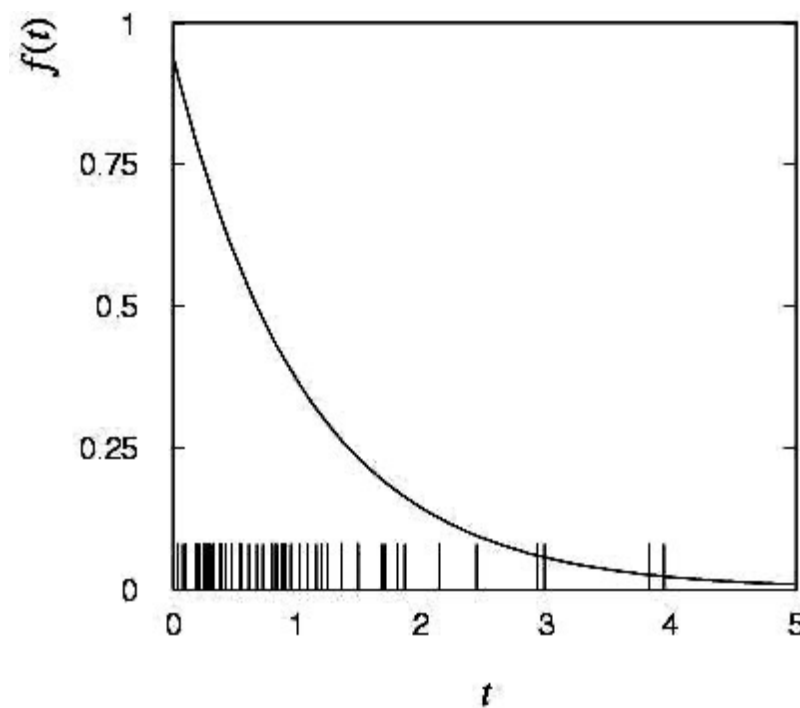
Find its maximum by setting  $\frac{\partial \ln L(\tau)}{\partial \tau} = 0$ ,

$$\rightarrow \hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$$

Monte Carlo test:  
generate 50 values  
using  $\tau = 1$ :

We find the ML estimate:

$$\hat{\tau} = 1.062$$



# MLE example: parameter of exponential pdf (3)

For the exponential distribution one has for mean, variance:

$$E[t] = \int_0^{\infty} t \frac{1}{\tau} e^{-t/\tau} dt = \tau$$

$$V[t] = \int_0^{\infty} (t - \tau)^2 \frac{1}{\tau} e^{-t/\tau} dt = \tau^2$$


For the MLE  $\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$  we therefore find

$$E[\hat{\tau}] = E\left[\frac{1}{n} \sum_{i=1}^n t_i\right] = \frac{1}{n} \sum_{i=1}^n E[t_i] = \tau \quad \longrightarrow \quad b = E[\hat{\tau}] - \tau = 0$$

$$V[\hat{\tau}] = V\left[\frac{1}{n} \sum_{i=1}^n t_i\right] = \frac{1}{n^2} \sum_{i=1}^n V[t_i] = \frac{\tau^2}{n} \quad \longrightarrow \quad \sigma_{\hat{\tau}} = \frac{\tau}{\sqrt{n}}$$

# The information inequality

The information inequality (RCF) sets a lower bound on the variance of any estimator (not only MLE). For a single parameter:

$$(b = E[\hat{\theta}] - \theta)$$


$$V[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 \bigg/ E \left[ -\frac{\partial^2 \ln L}{\partial \theta^2} \right] = \text{MVB} \quad (\text{Minimum Variance Bound})$$

where  $E \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right] = \int \frac{\partial^2 \ln P(\mathbf{x}|\theta)}{\partial \theta^2} P(\mathbf{x}|\theta) d\mathbf{x}$

Proof in Exercise 6.6 of SDA, [http://www.pp.rhul.ac.uk/~cowan/sda/prob/prob\\_6.pdf](http://www.pp.rhul.ac.uk/~cowan/sda/prob/prob_6.pdf)

“Efficiency” of an estimator = MVB / actual variance.

An estimator whose variance equals the MVB is said to be efficient.

# MVB for MLE of exponential parameter

Find 
$$\text{MVB} = - \left( 1 + \frac{\partial b}{\partial \tau} \right)^2 \bigg/ E \left[ \frac{\partial^2 \ln L}{\partial \tau^2} \right]$$

We found for the exponential parameter the MLE  $\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$

and we showed  $b = 0$ , hence  $\partial b / \partial \tau = 0$ .

We find 
$$\frac{\partial^2 \ln L}{\partial \tau^2} = \sum_{i=1}^n \left( \frac{1}{\tau^2} - \frac{2t_i}{\tau^3} \right)$$

and since  $E[t_i] = \tau$  for all  $i$ , 
$$E \left[ \frac{\partial^2 \ln L}{\partial \tau^2} \right] = -\frac{n}{\tau^2},$$

and therefore  $\text{MVB} = \frac{\tau^2}{n} = V[\hat{\tau}]$ . So here the MLE is efficient.



# Large-sample (asymptotic) properties of MLEs

Suppose we have an i.i.d. data sample of size  $n$ :  $x_1, \dots, x_n$

In the large-sample (or “asymptotic”) limit ( $n \rightarrow \infty$ ) and assuming regularity conditions one can show that the likelihood and MLE have several important properties.

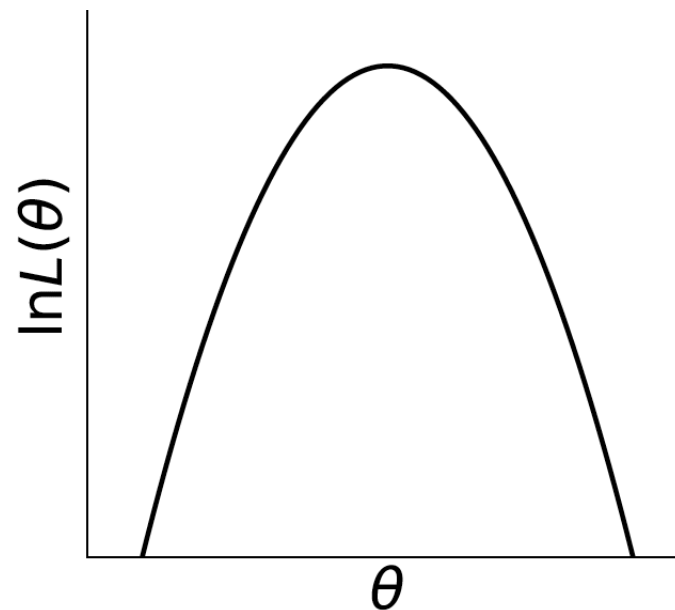
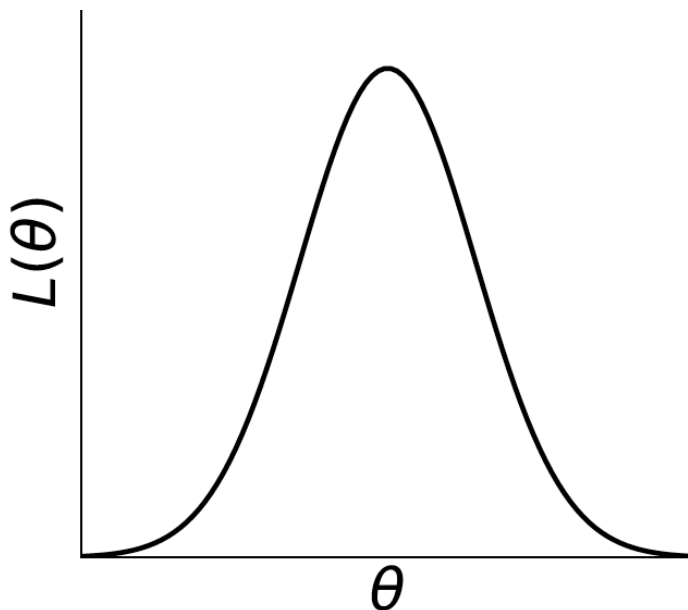
The regularity conditions include:

- the boundaries of the data space cannot depend on the parameter;
- the parameter cannot be on the edge of the parameter space;
- $\ln L(\theta)$  must be differentiable;
- the only solution to  $\partial \ln L / \partial \theta = 0$  is  $\hat{\theta}$ .

In the slides immediately following the properties are shown without proof for a single parameter; the corresponding properties hold also for the multiparameter case,  $\theta = (\theta_1, \dots, \theta_m)$ .

# log-likelihood becomes quadratic

The likelihood function becomes Gaussian in shape, i.e. the log-likelihood becomes quadratic (parabolic).



The MLE becomes increasingly precise as the (log)-likelihood becomes more tightly concentrated about its peak,  
but  $L(\theta) = P(\mathbf{x}|\theta)$  is the probability for  $\mathbf{x}$ , not a pdf for  $\theta$ .

# The MLE converges to the true parameter value

In the large-sample limit, the MLE converges in probability to the true parameter value.

That is, for any  $\varepsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| > \epsilon) = 0$$

The MLE is said to be *consistent*.

# MLE is asymptotically unbiased

In general the MLE can be biased, but in the large-sample limit, this bias goes to zero:

$$\lim_{n \rightarrow \infty} E[\hat{\theta}] - \theta = 0$$

(Recall for the exponential parameter we found the bias was identically zero regardless of the sample size  $n$ .)

# The MLE's variance approaches the MVB

In the large-sample limit, the variance of the MLE approaches the minimum-variance bound, i.e., the information inequality becomes an equality (and bias goes to zero):

$$\lim_{n \rightarrow \infty} V[\hat{\theta}] = - \frac{1}{E \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right]}$$

The MLE is said to be *asymptotically efficient*.

# The MLE's distribution becomes Gaussian

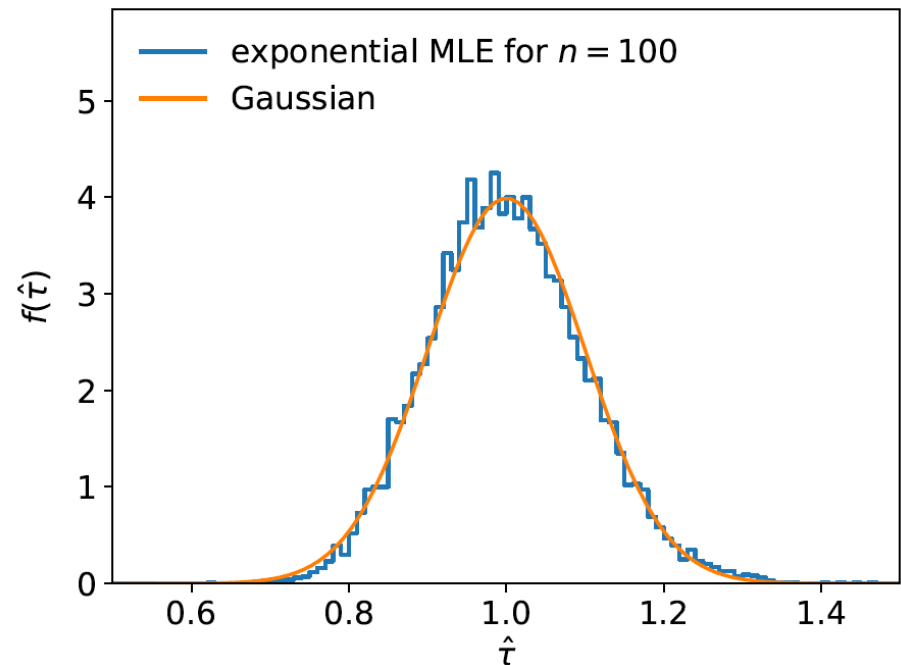
In the large-sample limit, the pdf of the MLE becomes Gaussian,

$$f(\hat{\theta}) = \frac{1}{\sqrt{2\pi}\sigma_{\hat{\theta}}} e^{-(\hat{\theta}-\theta)^2/2\sigma_{\hat{\theta}}^2}$$

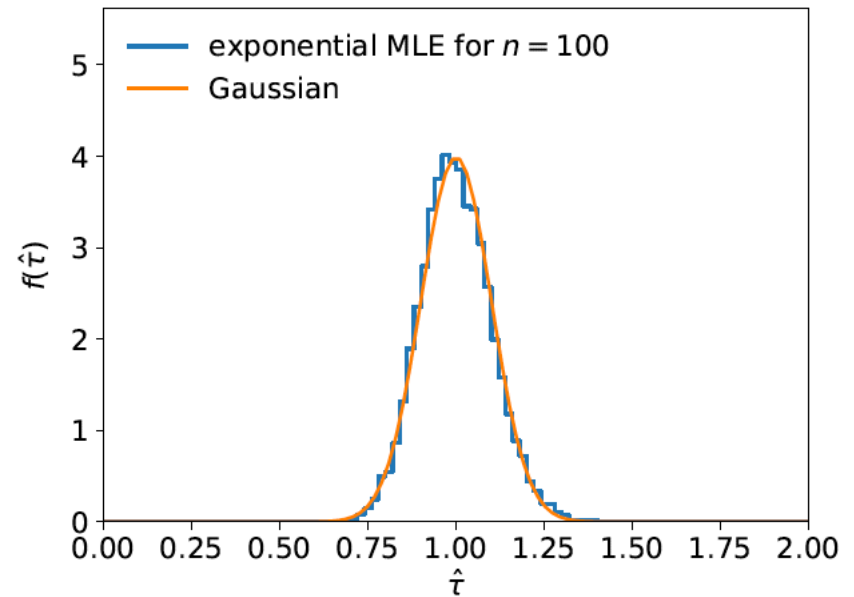
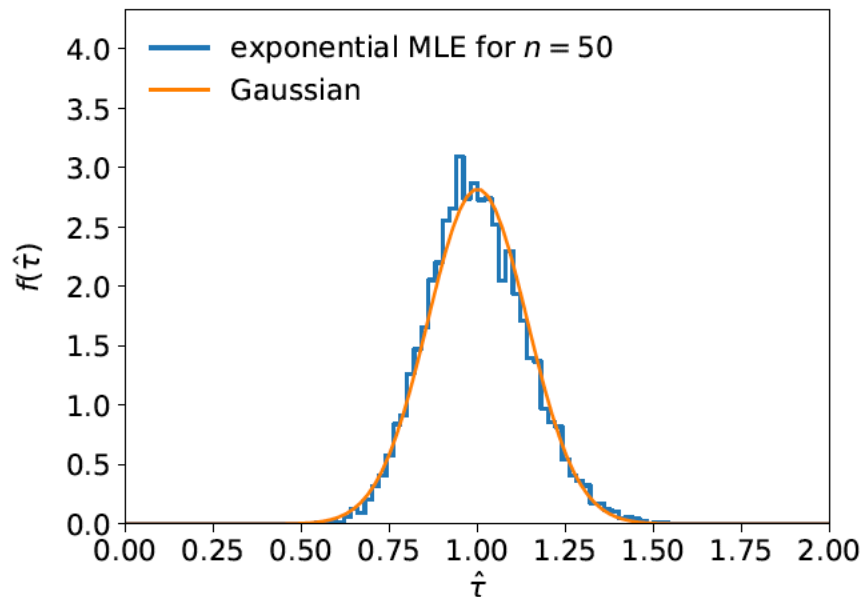
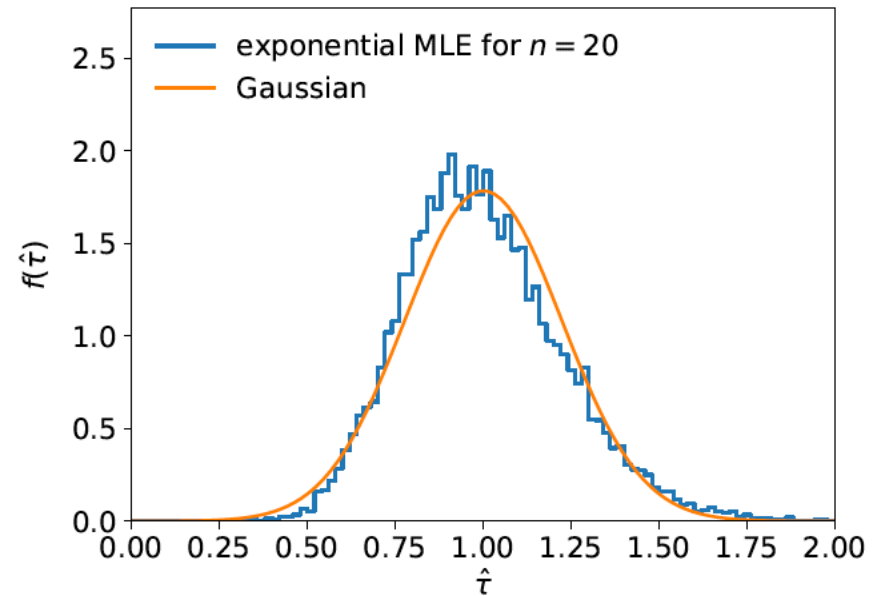
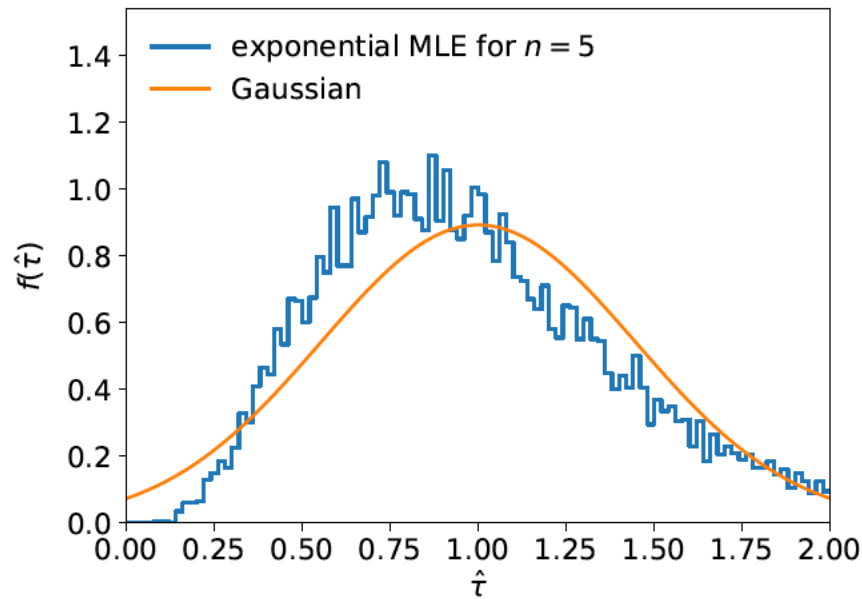
where  $\sigma_{\hat{\theta}}^2$  is the minimum variance bound (note bias is zero).

For example, exponential MLE with sample size  $n = 100$ .

Note that for exponential, MLE is arithmetic average, so Gaussian MLE seen to stem from Central Limit Theorem.



# Distribution of MLE of exponential parameter



# Variance of estimators: graphical method

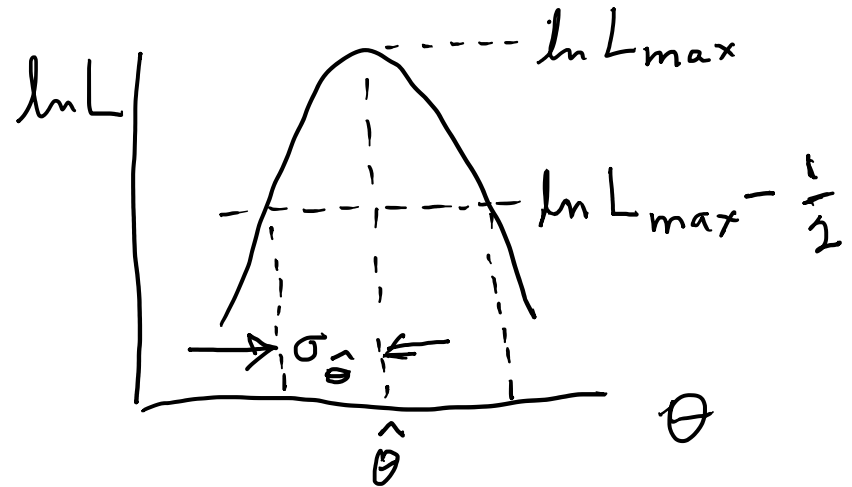
Expand  $\ln L(\theta)$  about its maximum:

$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[ \frac{\partial \ln L}{\partial \theta} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2!} \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta})^2 + \dots$$

First term is  $\ln L_{\max}$ , second term is zero, for third term use information inequality (assume equality):

$$\ln L(\theta) \approx \ln L_{\max} - \frac{(\theta - \hat{\theta})^2}{2\hat{\sigma}_{\hat{\theta}}^2}$$

$$\ln L(\hat{\theta} \pm \hat{\sigma}_{\hat{\theta}}) \approx \ln L_{\max} - \frac{1}{2}$$



→ to get  $\hat{\sigma}_{\hat{\theta}}$ , change  $\theta$  away from  $\hat{\theta}$  until  $\ln L$  decreases by  $1/2$ .



# Example of variance by graphical method

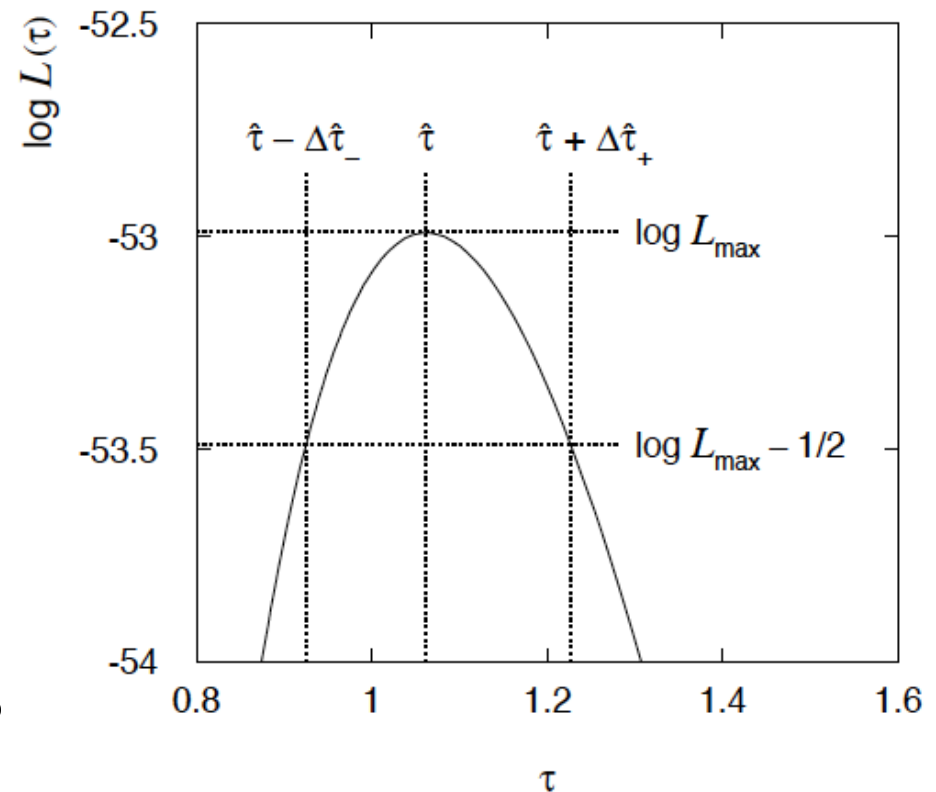
ML example with exponential:

$$\hat{\tau} = 1.062$$

$$\Delta\hat{\tau}_- = 0.137$$

$$\Delta\hat{\tau}_+ = 0.165$$

$$\hat{\sigma}_{\hat{\tau}} \approx \Delta\hat{\tau}_- \approx \Delta\hat{\tau}_+ \approx 0.15$$



Not quite parabolic  $\ln L$  since finite sample size ( $n = 50$ ).

# Information inequality for $N$ parameters

Suppose we have estimated  $N$  parameters  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_N)$

The *Fisher information matrix* is

$$I_{ij} = -E \left[ \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right] = - \int \frac{\partial^2 \ln P(\mathbf{x}|\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} P(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x}$$

and the covariance matrix of estimators  $\hat{\boldsymbol{\theta}}$  is  $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$

The information inequality states that the matrix

$$M_{ij} = V_{ij} - \sum_{k,l} \left( \delta_{ik} + \frac{\partial b_i}{\partial \theta_k} \right) I_{kl}^{-1} \left( \delta_{lj} + \frac{\partial b_l}{\partial \theta_j} \right)$$

is positive semi-definite:

$$\mathbf{z}^T M \mathbf{z} \geq 0 \text{ for all } \mathbf{z} \neq 0, \text{ diagonal elements } \geq 0$$

# Information inequality for $N$ parameters (2)

In practice the inequality is ~always used in the large-sample limit:

bias  $\rightarrow 0$

inequality  $\rightarrow$  equality, i.e,  $M = 0$ , and therefore  $V^{-1} = I$

That is, 
$$V_{ij}^{-1} = -E \left[ \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right]$$

This can be estimated from data using 
$$\hat{V}_{ij}^{-1} = - \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\theta}}$$

Find the matrix  $V^{-1}$  numerically (or with automatic differentiation), then invert to get the covariance matrix of the estimators

$$\hat{V}_{ij} = \widehat{\text{cov}}[\hat{\theta}_i, \hat{\theta}_j]$$

# Multiparameter graphical method for variances

Expand  $\ln L(\boldsymbol{\theta})$  to 2<sup>nd</sup> order about MLE:

$$\ln L(\boldsymbol{\theta}) \approx \ln L(\hat{\boldsymbol{\theta}}) + \sum_i \frac{\partial \ln L}{\partial \theta_i} \Big|_{\hat{\boldsymbol{\theta}}} (\theta_i - \hat{\theta}_i) + \frac{1}{2!} \sum_{i,j} \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \Big|_{\hat{\boldsymbol{\theta}}} (\theta_i - \hat{\theta}_i)(\theta_j - \hat{\theta}_j)$$

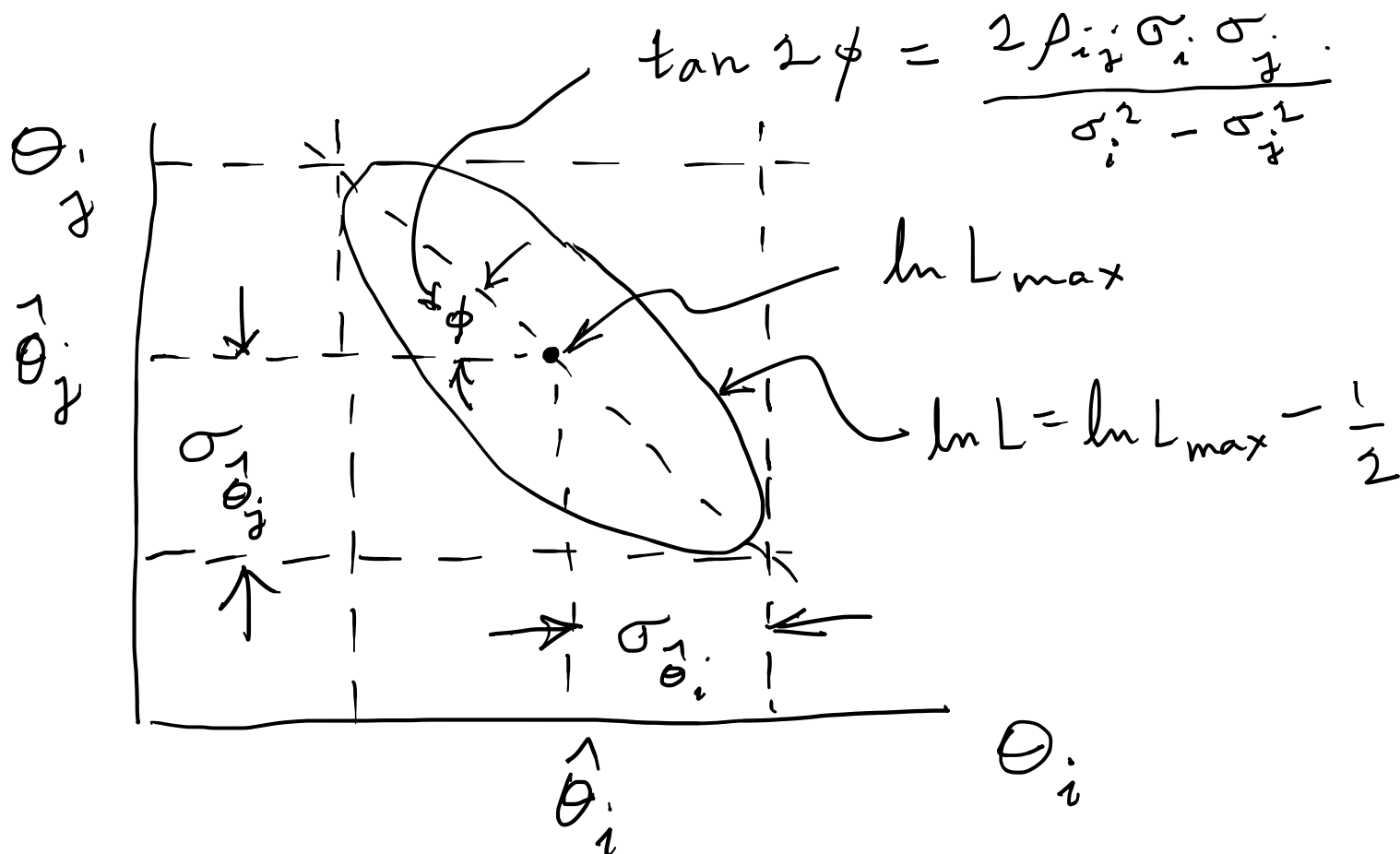
$\ln L_{\max}$                       zero                      relate to covariance matrix of MLEs using information (in)equality.

**Result:**  $\ln L(\boldsymbol{\theta}) = \ln L_{\max} - \frac{1}{2} \sum_{i,j} (\theta_i - \hat{\theta}_i) V_{ij}^{-1} (\theta_j - \hat{\theta}_j)$

So the surface  $\ln L(\boldsymbol{\theta}) = \ln L_{\max} - \frac{1}{2}$  corresponds to

$(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^T V^{-1} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) = 1$  , which is the equation of a (hyper-) ellipse.

## Multiparameter graphical method (2)



Distance from MLE to tangent planes gives standard deviations.

# Extra slides

# Some statistics books, papers, etc.

G. Cowan, *Statistical Data Analysis*, Clarendon, Oxford, 1998

R.J. Barlow, *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences*, Wiley, 1989

Ilya Narsky and Frank C. Porter, *Statistical Analysis Techniques in Particle Physics*, Wiley, 2014.

Luca Lista, *Statistical Methods for Data Analysis in Particle Physics*, Springer, 2017.

L. Lyons, *Statistics for Nuclear and Particle Physics*, CUP, 1986

F. James., *Statistical and Computational Methods in Experimental Physics*, 2nd ed., World Scientific, 2006

S. Brandt, *Statistical and Computational Methods in Data Analysis*, Springer, New York, 1998.

R.L. Workman et al. (Particle Data Group), Prog. Theor. Exp. Phys. 083C01 (2022); [pdg.lbl.gov](https://pdg.lbl.gov) sections on probability, statistics, MC.

# Theory ↔ Statistics ↔ Experiment

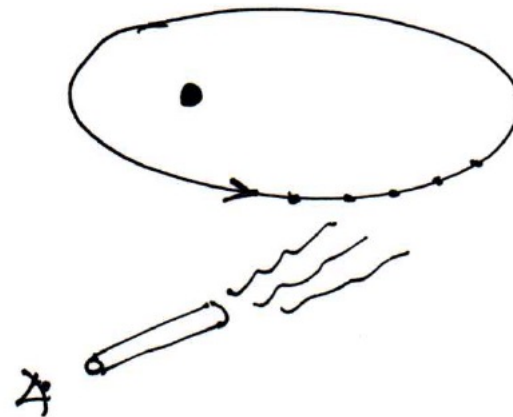
Theory (model, hypothesis):

$$F = -G \frac{m_1 m_2}{r^2}, \dots$$

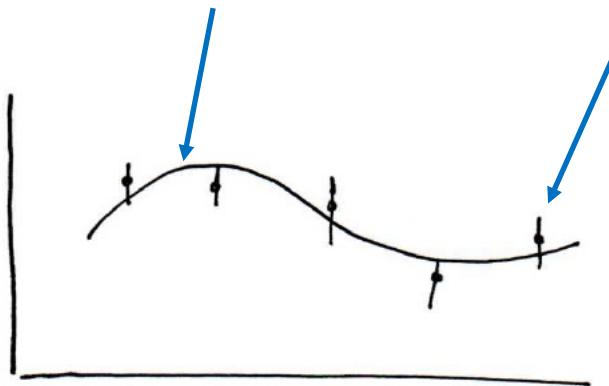
+ response of measurement apparatus

= model prediction

Experiment (observation):



data



Uncertainty enters on many levels

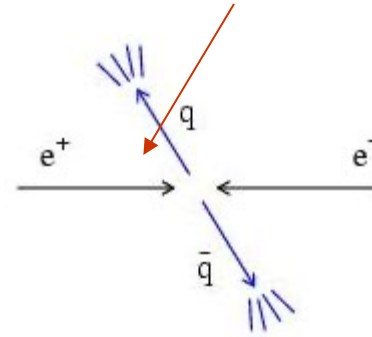
→ quantify with  
**probability**



# Example of ML with 2 parameters

Consider a scattering angle distribution with  $x = \cos \theta$ ,

$$f(x; \alpha, \beta) = \frac{1 + \alpha x + \beta x^2}{2 + 2\beta/3}$$



or if  $x_{\min} < x < x_{\max}$ , need to normalize so that

$$\int_{x_{\min}}^{x_{\max}} f(x; \alpha, \beta) dx = 1 .$$

Example:  $\alpha = 0.5$ ,  $\beta = 0.5$ ,  $x_{\min} = -0.95$ ,  $x_{\max} = 0.95$ ,  
generate  $n = 2000$  events with Monte Carlo.

$$\ln L(\alpha, \beta) = \sum_{i=1}^n \ln f(x_i; \alpha, \beta) \quad \longleftarrow \quad \text{need to find maximum numerically}$$

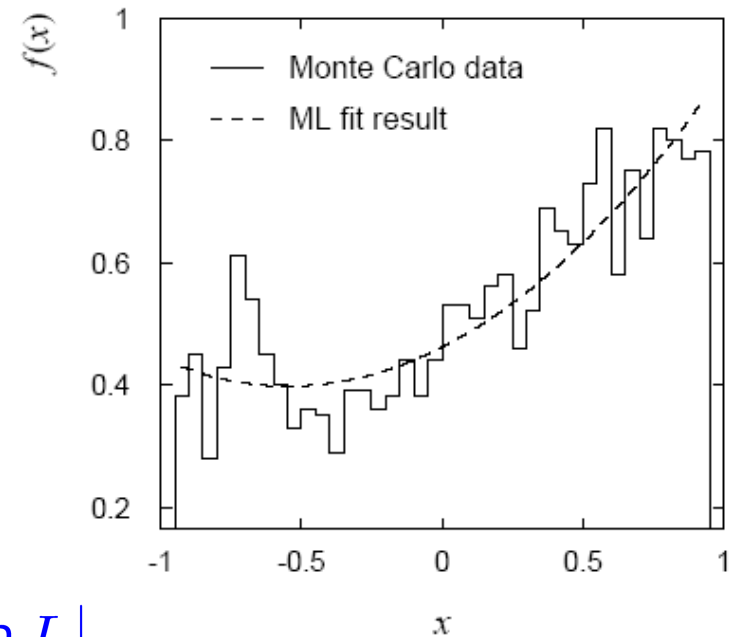
# Example of ML with 2 parameters: fit result

Finding maximum of  $\ln L(\alpha, \beta)$  numerically gives

$$\hat{\alpha} = 0.508$$

$$\hat{\beta} = 0.47$$

N.B. No binning of data for fit,  
but can compare to histogram for  
goodness-of-fit (e.g. 'visual' or  $\chi^2$ ).



(Co)variances from  $(\widehat{V}^{-1})_{ij} = -\left. \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right|_{\vec{\theta} = \vec{\hat{\theta}}}$

$$\hat{\sigma}_{\hat{\alpha}} = 0.052$$

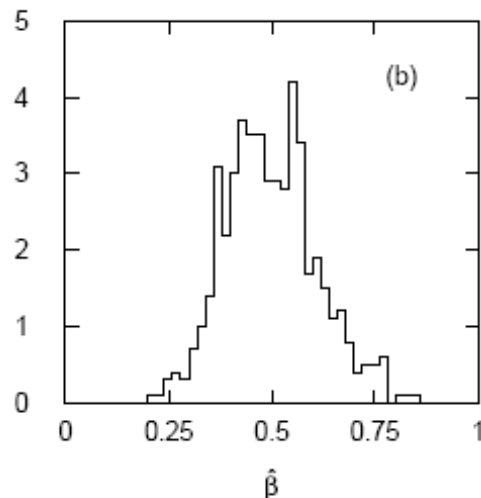
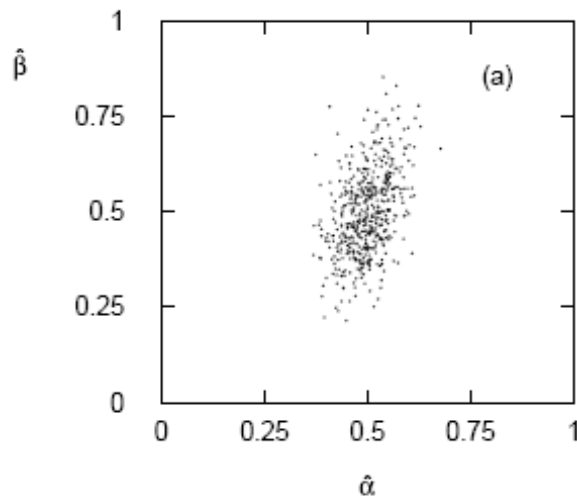
$$\text{cov}[\hat{\alpha}, \hat{\beta}] = 0.0026$$

$$\hat{\sigma}_{\hat{\beta}} = 0.11$$

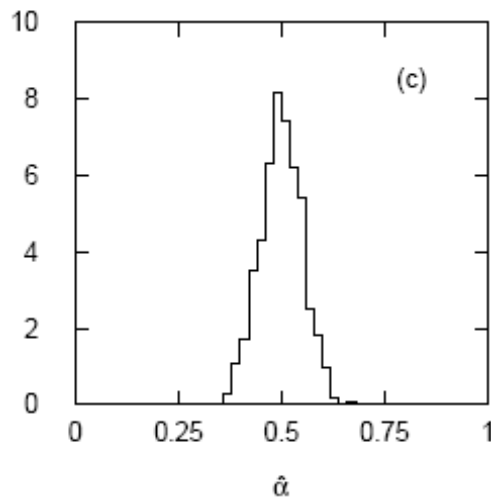
$$r = 0.46 = \text{correlation coefficient}$$

# Two-parameter fit: MC study

Repeat ML fit with 500 experiments, all with  $n = 2000$  events:



$$\begin{aligned}\overline{\hat{\alpha}} &= 0.499 \\ s_{\hat{\alpha}} &= 0.051 \\ \overline{\hat{\beta}} &= 0.498 \\ s_{\hat{\beta}} &= 0.111 \\ \widehat{\text{cov}}[\hat{\alpha}, \hat{\beta}] &= 0.0024 \\ r &= 0.42\end{aligned}$$



Estimates average to  $\sim$ true values;  
(Co)variances close to previous estimates;  
marginal pdfs approximately Gaussian.

# The $\ln L_{\max} - 1/2$ contour for two parameters

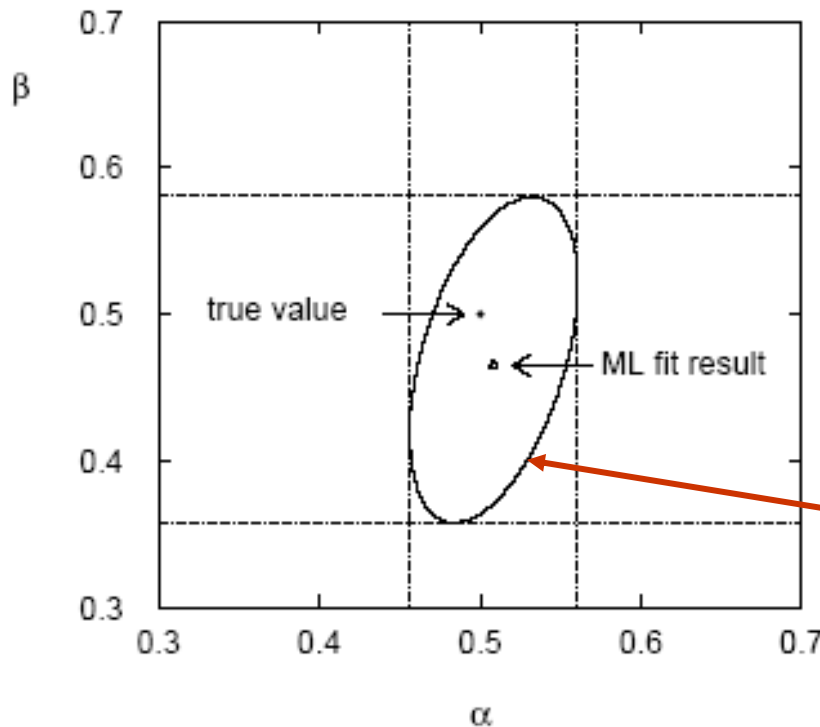
For large  $n$ ,  $\ln L$  takes on quadratic form near maximum:

$$\ln L(\alpha, \beta) \approx \ln L_{\max} - \frac{1}{2(1 - \rho^2)} \left[ \left( \frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right)^2 + \left( \frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right)^2 - 2\rho \left( \frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right) \left( \frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right) \right]$$

The contour  $\ln L(\alpha, \beta) = \ln L_{\max} - 1/2$  is an ellipse:

$$\frac{1}{(1 - \rho^2)} \left[ \left( \frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right)^2 + \left( \frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right)^2 - 2\rho \left( \frac{\alpha - \hat{\alpha}}{\sigma_{\hat{\alpha}}} \right) \left( \frac{\beta - \hat{\beta}}{\sigma_{\hat{\beta}}} \right) \right] = 1$$

## (Co)variances from $\ln L$ contour



The  $\alpha, \beta$  plane for the first MC data set

$$\ln L(\alpha, \beta) = \ln L_{\max} - 1/2$$

→ Tangent lines to contours give standard deviations.

→ Angle of ellipse  $\phi$  related to correlation:  $\tan 2\phi = \frac{2\rho\sigma_{\hat{\alpha}}\sigma_{\hat{\beta}}}{\sigma_{\hat{\alpha}}^2 - \sigma_{\hat{\beta}}^2}$