# Statistics for Particle Physics Lecture day 2





Taller de Altas EnergíasBenasque, Spain (online)5,6 September 2022

http://benasque.org/2022tae/



Glen Cowan Physics Department Royal Holloway, University of London g.cowan@rhul.ac.uk www.pp.rhul.ac.uk/~cowan

		Outline
	Monday 9-11 :	Introduction
		Probability
		Hypothesis tests
		Machine Learning
$\rightarrow$	Tuesday 9-11 :	Parameter estimation
		Confidence limits
		Systematic uncertainties
		Experimental sensitivity
	Tuesday 15:30:	Tutorial on parameter estimation

Almost everything is a subset of the University of London course: http://www.pp.rhul.ac.uk/~cowan/stat\_course.html

#### Parameter estimation

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



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

Suppose we have a sample of observed values:  $x = (x_1, ..., x_n)$ 

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

 $\hat{\theta}(\vec{x}) \leftarrow \text{estimator written with a hat}$ 

Sometimes we say 'estimator' for the function of  $x_1, ..., x_n$ ; '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):  $b = E[\hat{\theta}] - \theta$ 

 $\rightarrow$  average of repeated measurements should tend to true value.

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

 $\rightarrow$  small bias & variance are in general conflicting criteria

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

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

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

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

In this case the likelihood function is

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

## Maximum Likelihood Estimators (MLEs)

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



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,  $t_1, \ldots, t_n$ 

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

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

$$\frac{\partial \ln L(\tau)}{\partial \tau} = 0 ,$$

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

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

TAE 2022 Benasque (online) / Lecture day 2

## Variance of estimators: Monte Carlo method

Having estimated our parameter we now need to report its 'statistical error', i.e., how widely distributed would estimates be if we were to repeat the entire measurement many times.

One way to do this would be to simulate the entire experiment many times with a Monte Carlo program (use ML estimate for MC).

For exponential example, from sample variance of estimates we find:

 $\hat{\sigma}_{\hat{\tau}} = 0.151$ 

Note distribution of estimates is roughly Gaussian – (almost) always true for ML in large sample limit.



## Variance of estimators from information inequality

The information inequality (RCF) sets a lower bound on the variance of any estimator (not only ML):

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

Often the bias *b* is small, and equality either holds exactly or is a good approximation (e.g. large data sample limit). Then,

$$V[\hat{\theta}] \approx -1 \left/ E\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]$$

Estimate this using the 2nd derivative of  $\ln L$  at its maximum:

$$\widehat{V}[\widehat{\theta}] = -\left(\frac{\partial^2 \ln L}{\partial \theta^2}\right)^{-1} \bigg|_{\theta = \widehat{\theta}}$$

#### MVB for MLE of exponential parameter

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

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 MVB  $= \frac{\tau^2}{n} = V[\hat{\tau}]$ . (Here MLE is "efficient").

# Variance of estimators: graphical method

Expand  $lnL(\theta)$  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):

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

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

 $\rightarrow$  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



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

## Confidence intervals by inverting a test

In addition to a 'point estimate' of a parameter we should report an interval reflecting its statistical uncertainty.

Confidence intervals for a parameter  $\theta$  can be found by defining a test of the hypothesized value  $\theta$  (do this for all  $\theta$ ):

Specify values of the data that are 'disfavoured' by  $\theta$ (critical region) such that  $P(\text{data in critical region} | \theta) \le \alpha$ for a prespecified  $\alpha$ , e.g., 0.05 or 0.1.

If data observed in the critical region, reject the value  $\theta$ .

Now invert the test to define a confidence interval as:

set of  $\theta$  values that are not rejected in a test of size  $\alpha$  (confidence level CL is  $1 - \alpha$ ).

Relation between confidence interval and *p*-value

Equivalently we can consider a significance test for each hypothesized value of  $\theta$ , resulting in a *p*-value,  $p_{\theta}$ .

If  $p_{\theta} \leq \alpha$ , then we reject  $\theta$ .

The confidence interval at  $CL = 1 - \alpha$  consists of those values of  $\theta$  that are not rejected.

E.g. an upper limit on  $\theta$  is the greatest value for which  $p_{\theta} > \alpha$ .

In practice find by setting  $p_{\theta} = \alpha$  and solve for  $\theta$ .

For a multidimensional parameter space  $\theta = (\theta_1, \dots, \theta_M)$  use same idea – result is a confidence "region" with boundary determined by  $p_{\theta} = \alpha$ .

## Coverage probability of confidence interval

If the true value of  $\theta$  is rejected, then it's not in the confidence interval. The probability for this is by construction (equality for continuous data):

 $P(\text{reject } \theta | \theta) \leq \alpha = \text{type-I error rate}$ 

Therefore, the probability for the interval to contain or "cover"  $\theta$  is

*P*(conf. interval "covers"  $\theta | \theta \ge 1 - \alpha$ 

This assumes that the set of  $\theta$  values considered includes the true value, i.e., it assumes the composite hypothesis  $P(\mathbf{x}|H,\theta)$ .

#### Frequentist upper limit on Poisson parameter

Consider again the case of observing  $n \sim \text{Poisson}(s + b)$ . Suppose b = 4.5,  $n_{\text{obs}} = 5$ . Find upper limit on s at 95% CL. Relevant alternative is s = 0 (critical region at low n) p-value of hypothesized s is  $P(n \le n_{\text{obs}}; s, b)$ Upper limit  $s_{\text{up}}$  at  $\text{CL} = 1 - \alpha$  found from

$$\alpha = P(n \le n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}} + b)}$$
$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} (1 - \alpha; 2(n_{\text{obs}} + 1)) - b$$
$$= \frac{1}{2} F_{\chi^2}^{-1} (0.95; 2(5 + 1)) - 4.5 = 6.0$$

G. Cowan / RHUL Physics

TAE 2022 Benasque (online) / Lecture day 2

## $n \sim \text{Poisson}(s+b)$ : frequentist upper limit on s

For low fluctuation of *n*, formula can give negative result for  $s_{up}$ ; i.e. confidence interval is empty; all values of  $s \ge 0$  have  $p_s \le \alpha$ .



#### TAE 2022 Benasque (online) / Lecture day 2

### Limits near a boundary of the parameter space

Suppose e.g. b = 2.5 and we observe n = 0.

If we choose CL = 0.9, we find from the formula for  $s_{up}$ 

$$s_{\rm up} = -0.197$$
 (CL = 0.90)

Physicist:

We already knew  $s \ge 0$  before we started; can't use negative upper limit to report result of expensive experiment!

#### Statistician:

The interval is designed to cover the true value only 90% of the time — this was clearly not one of those times.

Not uncommon dilemma when testing parameter values for which one has very little experimental sensitivity, e.g., very small *s*.

## Expected limit for s = 0

Physicist: I should have used CL = 0.95 — then  $s_{up} = 0.496$ 

Even better: for CL = 0.917923 we get  $s_{up} = 10^{-4}$  !

Reality check: with b = 2.5, typical Poisson fluctuation in n is at least  $\sqrt{2.5} = 1.6$ . How can the limit be so low?



## Approximate confidence intervals/regions from the likelihood function

Suppose we test parameter value(s)  $\theta = (\theta_1, ..., \theta_n)$  using the ratio

$$\lambda(\boldsymbol{\theta}) = \frac{L(\boldsymbol{\theta})}{L(\hat{\boldsymbol{\theta}})} \qquad \qquad 0 \le \lambda(\boldsymbol{\theta}) \le 1$$

Lower  $\lambda(\theta)$  means worse agreement between data and hypothesized  $\theta$ . Equivalently, usually define

$$t_{\theta} = -2\ln\lambda(\theta)$$

so higher  $t_{\theta}$  means worse agreement between  $\theta$  and the data.

*p*-value of  $\theta$  therefore

$$p_{\theta} = \int_{t_{\theta, \text{obs}}}^{\infty} f(t_{\theta}|\theta) \, dt_{\theta}$$
need pdf

G. Cowan / RHUL Physics

TAE 2021 / Statistics for PP Lecture 3

## Confidence region from Wilks' theorem

Wilks' theorem says (in large-sample limit and provided certain conditions hold...)

 $f(t_{\theta}|\theta) \sim \chi_n^2 \qquad \begin{array}{l} \text{chi-square dist. with $\#$ d.o.f. =} \\ \# \text{ of components in $\theta = (\theta_1, ..., \theta_n)$.} \end{array}$ 

Assuming this holds, the *p*-value is

$$p_{m{ heta}} = 1 - F_{\chi^2_n}(t_{m{ heta}}) \quad \leftarrow \text{set equal to } lpha$$

To find boundary of confidence region set  $p_{\theta} = \alpha$  and solve for  $t_{\theta}$ :

$$t_{\theta} = F_{\chi_n^2}^{-1}(1-\alpha)$$

Recall also

$$t_{\theta} = -2\ln\frac{L(\theta)}{L(\hat{\theta})}$$

G. Cowan / RHUL Physics

TAE 2021 / Statistics for PP Lecture 3

## Confidence region from Wilks' theorem (cont.)

i.e., boundary of confidence region in  $\theta$  space is where

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}F_{\chi_n^2}^{-1}(1-\alpha)$$

For example, for  $1 - \alpha = 68.3\%$  and n = 1 parameter,

$$F_{\chi_1^2}^{-1}(0.683) = 1$$

and so the 68.3% confidence level interval is determined by

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}$$

Same as recipe for finding the estimator's standard deviation, i.e.,

 $[\hat{\theta} - \sigma_{\hat{\theta}}, \hat{\theta} + \sigma_{\hat{\theta}}]$  is a 68.3% CL confidence interval.

## Example of interval from $\ln L(\theta)$

For n=1 parameter, CL = 0.683,  $Q_{\alpha} = 1$ .



G. Cowan / RHUL Physics

#### TAE 2021 / Statistics for PP Lecture 3

## Multiparameter case

For increasing number of parameters,  $CL = 1 - \alpha$  decreases for confidence region determined by a given

$$Q_{\alpha} = F_{\chi_n^2}^{-1}(1-\alpha)$$

0	1-lpha					
$\mathcal{Q}\alpha$	n = 1	n = 2	n = 3	n = 4	n = 5	
1.0	0.683	0.393	0.199	0.090	0.037	
2.0	0.843	0.632	0.428	0.264	0.151	
4.0	0.954	0.865	0.739	0.594	0.451	
9.0	0.997	0.989	0.971	0.939	0.891	

## Multiparameter case (cont.)

Equivalently,  $Q_{\alpha}$  increases with *n* for a given  $CL = 1 - \alpha$ .

1 0	$\widehat{Q}_{lpha}$					
$1 - \alpha$	n = 1	n = 2	n = 3	n = 4	n = 5	
0.683	1.00	2.30	3.53	4.72	5.89	
0.90	2.71	4.61	6.25	7.78	9.24	
0.95	3.84	5.99	7.82	9.49	11.1	
0.99	6.63	9.21	11.3	13.3	15.1	

## Systematic uncertainties and nuisance parameters In general, our model of the data is not perfect:



Can improve model by including additional adjustable parameters.

$$P(x|\mu) \to P(x|\mu, \theta)$$

Nuisance parameter ↔ systematic uncertainty. Some point in the parameter space of the enlarged model should be "true".

Presence of nuisance parameter decreases sensitivity of analysis to the parameter of interest (e.g., increases variance of estimate).

### **Profile Likelihood**

Suppose we have a likelihood  $L(\mu, \theta) = P(x|\mu, \theta)$  with Nparameters of interest  $\mu = (\mu_1, ..., \mu_N)$  and M nuisance parameters  $\theta = (\theta_1, ..., \theta_M)$ . The "profiled" (or "constrained") values of  $\theta$  are:

$$\hat{\hat{\boldsymbol{\theta}}}(\boldsymbol{\mu}) = \operatorname*{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\mu}, \boldsymbol{\theta})$$

and the profile likelihood is:  $L_{\rm p}(\boldsymbol{\mu}) = L(\boldsymbol{\mu}, \hat{\boldsymbol{\theta}})$ 

The profile likelihood depends only on the parameters of interest; the nuisance parameters are replaced by their profiled values.

The profile likelihood can be used to obtain confidence intervals/regions for the parameters of interest in the same way as one would for all of the parameters from the full likelihood.

## Profile Likelihood Ratio – Wilks theorem

Goal is to test/reject regions of  $\mu$  space (param. of interest).

Rejecting a point  $\mu$  should mean  $p_{\mu} \leq \alpha$  for all possible values of the nuisance parameters  $\theta$ .

Test  $\boldsymbol{\mu}$  using the "profile likelihood ratio":  $\lambda(\boldsymbol{\mu}) = \frac{L(\boldsymbol{\mu}, \hat{\boldsymbol{\theta}})}{L(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\theta}})}$ 

Let  $t_{\mu} = -2 \ln \lambda(\mu)$ . Wilks' theorem says in large-sample limit:  $t_{\mu} \sim \text{chi-square}(N)$ 

where the number of degrees of freedom is the number of parameters of interest (components of  $\mu$ ). So *p*-value for  $\mu$  is

$$p_{\boldsymbol{\mu}} = \int_{t_{\boldsymbol{\mu},\text{obs}}}^{\infty} f(t_{\boldsymbol{\mu}} | \boldsymbol{\mu}, \boldsymbol{\theta}) \, dt_{\boldsymbol{\mu}} = 1 - F_{\chi_N^2}(t_{\boldsymbol{\mu},\text{obs}})$$

G. Cowan / RHUL Physics

TAE 2021 / Statistics for PP Lecture 3

## Profile Likelihood Ratio – Wilks theorem (2)

If we have a large enough data sample to justify use of the asymptotic chi-square pdf, then if  $\mu$  is rejected, it is rejected for any values of the nuisance parameters.

The recipe to get confidence regions/intervals for the parameters of interest at  $CL = 1 - \alpha$  is thus the same as before, simply use the profile likelihood:

$$\ln L_{\rm p}(\boldsymbol{\mu}) = \ln L_{\rm max} - \frac{1}{2} F_{\chi_N^2}^{-1} (1 - \alpha)$$

where the number of degrees of freedom N for the chi-square quantile is equal to the number of parameters of interest.

If the large-sample limit is not justified, then use e.g. Monte Carlo to get distribution of  $t_{\mu}$ .

#### Prototype search analysis

Search for signal in a region of phase space; result is histogram of some variable x giving numbers:

$$\mathbf{n}=(n_1,\ldots,n_N)$$

Assume the  $n_i$  are Poisson distributed with expectation values

$$E[n_i] = \mu s_i + b_i$$

$$frength parameter$$

where

G. Cowan / RHUL Physics

TAE 2021 / Statistics for PP Lecture 3

## Prototype analysis (II)

Often also have a subsidiary measurement that constrains some of the background and/or shape parameters:

$$\mathbf{m} = (m_1, \ldots, m_M)$$

Assume the  $m_i$  are Poisson distributed with expectation values

$$E[m_i] = u_i(\boldsymbol{\theta})$$
nuisance parameters ( $\boldsymbol{\theta}_{s}, \boldsymbol{\theta}_{b}, b_{tot}$ )

Likelihood function is

$$L(\mu, \theta) = \prod_{j=1}^{N} \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)} \quad \prod_{k=1}^{M} \frac{u_k^{m_k}}{m_k!} e^{-u_k}$$

G. Cowan / RHUL Physics

TAE 2021 / Statistics for PP Lecture 3

## The profile likelihood ratio

Base significance test on the profile likelihood ratio:



Define critical region of test of  $\mu$  by the region of data space that gives the lowest values of  $\lambda(\mu)$ .

Important advantage of profile LR is that its distribution becomes independent of nuisance parameters in large sample limit.

## Test statistic for discovery

Suppose relevant alternative to background-only ( $\mu = 0$ ) is  $\mu \ge 0$ . So take critical region for test of  $\mu = 0$  corresponding to high  $a_{\mu}$ 

So take critical region for test of  $\mu = 0$  corresponding to high  $q_0$ and  $\hat{\mu} > 0$  (data characteristic for  $\mu \ge 0$ ).

That is, to test background-only hypothesis define statistic

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \ge 0\\ 0 & \hat{\mu} < 0 \end{cases}$$

i.e. here only large (positive) observed signal strength is evidence against the background-only hypothesis.

Note that even though here physically  $\mu \ge 0$ , we allow  $\hat{\mu}$  to be negative. In large sample limit its distribution becomes Gaussian, and this will allow us to write down simple expressions for distributions of our test statistics.

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

### Distribution of $q_0$ in large-sample limit

Assuming approximations valid in the large sample (asymptotic) limit, we can write down the full distribution of  $q_0$  as

$$f(q_0|\mu') = \left(1 - \Phi\left(\frac{\mu'}{\sigma}\right)\right)\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}\exp\left[-\frac{1}{2}\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)^2\right]$$

The special case  $\mu' = 0$  is a "half chi-square" distribution:

$$f(q_0|0) = \frac{1}{2}\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}e^{-q_0/2}$$

In large sample limit,  $f(q_0|0)$  independent of nuisance parameters;  $f(q_0|\mu')$  depends on nuisance parameters through  $\sigma$ .
# *p*-value for discovery

Large  $q_0$  means increasing incompatibility between the data and hypothesis, therefore *p*-value for an observed  $q_{0,obs}$  is



use e.g. asymptotic formula



From *p*-value get equivalent significance,

$$Z = \Phi^{-1}(1-p)$$

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

# Cumulative distribution of $q_0$ , significance

From the pdf, the cumulative distribution of  $q_0$  is found to be

$$F(q_0|\mu') = \Phi\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)$$

The special case  $\mu' = 0$  is

$$F(q_0|0) = \Phi\left(\sqrt{q_0}\right)$$

The *p*-value of the  $\mu = 0$  hypothesis is

$$p_0 = 1 - F(q_0|0)$$

Therefore the discovery significance Z is simply

$$Z = \Phi^{-1}(1 - p_0) = \sqrt{q_0}$$

G. Cowan / RHUL Physics

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

#### Monte Carlo test of asymptotic formula

- $n \sim \text{Poisson}(\mu s + b)$
- $m \sim \text{Poisson}(\tau b)$
- $\mu =$  param. of interest
- *b* = nuisance parameter

Here take *s* known,  $\tau = 1$ .

Asymptotic formula is good approximation to  $5\sigma$ level ( $q_0 = 25$ ) already for  $b \sim 20$ .



#### How to read the $p_0$ plot

The "local"  $p_0$  means the *p*-value of the background-only hypothesis obtained from the test of  $\mu = 0$  at each individual  $m_{\rm H}$ , without any correct for the Look-Elsewhere Effect.

The "Expected" (dashed) curve gives the median  $p_0$  under assumption of the SM Higgs ( $\mu$  = 1) at each  $m_{\rm H}$ .



The blue band gives the width of the distribution  $(\pm 1\sigma)$  of significances under assumption of the SM Higgs.

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

# Test statistic for upper limits

For purposes of setting an upper limit on  $\mu$  use

$$q_{\mu} = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \end{cases} \quad \text{where} \quad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})}$$

I.e. when setting an upper limit, an upwards fluctuation of the data is not taken to mean incompatibility with the hypothesized  $\mu$ :

From observed  $q_{\mu}$  find *p*-value:  $p_{\mu} = \int_{q_{\mu,\text{obs}}}^{\infty} f(q_{\mu}|\mu) dq_{\mu}$ 

Large sample approximation:

$$p_{\mu} = 1 - \Phi\left(\sqrt{q_{\mu}}\right)$$

To find upper limit at CL =  $1-\alpha$ , set  $p_{\mu} = \alpha$  and solve for  $\mu$ .

G. Cowan / RHUL Physics

#### Monte Carlo test of asymptotic formulae

Consider again  $n \sim \text{Poisson}(\mu s + b)$ ,  $m \sim \text{Poisson}(\tau b)$ Use  $q_{\mu}$  to find *p*-value of hypothesized  $\mu$  values.

E.g.  $f(q_1|1)$  for *p*-value of  $\mu = 1$ . Typically interested in 95% CL, i.e., *p*-value threshold = 0.05, i.e.,  $q_1 = 2.69$  or  $Z_1 = \sqrt{q_1} = 1.64$ . Median[ $q_1|0$ ] gives "exclusion

sensitivity".

Here asymptotic formulae good for s = 6, b = 9.



How to read the green and yellow limit plots For every value of  $m_{\rm H}$ , find the upper limit on  $\mu$ .

Also for each  $m_{\rm H}$ , determine the distribution of upper limits  $\mu_{\rm up}$  one would obtain under the hypothesis of  $\mu$  = 0.

The dashed curve is the median  $\mu_{up}$ , and the green (yellow) bands give the  $\pm 1\sigma$  ( $2\sigma$ ) regions of this distribution.



#### Sensitivity for Poisson counting experiment

Count a number of events  $n \sim \text{Poisson}(s+b)$ , where

- s = expected number of events from signal,
- b = expected number of background events.

To test for discovery of signal compute p-value of s = 0 hypothesis,

$$p = P(n \ge n_{\text{obs}}|b) = \sum_{n=n_{\text{obs}}}^{\infty} \frac{b^n}{n!} e^{-b} = 1 - F_{\chi^2}(2b; 2n_{\text{obs}})$$

Usually convert to equivalent significance:  $Z = \Phi^{-1}(1-p)$ where  $\Phi$  is the standard Gaussian cumulative distribution, e.g., Z > 5 (a 5 sigma effect) means  $p < 2.9 \times 10^{-7}$ .

To characterize sensitivity to discovery, give expected (mean or median) Z under assumption of a given s.

G. Cowan / RHUL Physics

 $s/\sqrt{b}$  for expected discovery significance For large s + b,  $n \to x \sim \text{Gaussian}(\mu, \sigma)$ ,  $\mu = s + b$ ,  $\sigma = \sqrt{(s + b)}$ . For observed value  $x_{\text{obs}}$ , p-value of s = 0 is  $\text{Prob}(x > x_{\text{obs}} | s = 0)$ ,:

$$p_0 = 1 - \Phi\left(\frac{x_{\rm obs} - b}{\sqrt{b}}\right)$$

Significance for rejecting s = 0 is therefore

$$Z_0 = \Phi^{-1}(1 - p_0) = \frac{x_{\text{obs}} - b}{\sqrt{b}}$$

Expected (median) significance assuming signal rate s is

$$\mathrm{median}[Z_0|s+b] = \frac{s}{\sqrt{b}}$$

G. Cowan / RHUL Physics

# Better approximation for significance

Poisson likelihood for parameter s is

 $L(s) = \frac{(s+b)^n}{n!} e^{-(s+b)}$ For now no nuisance params.

To test for discovery use profile likelihood ratio:

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{s} \ge 0 \ , \\ 0 & \hat{s} < 0 \ . \end{cases} \qquad \lambda(s) = \frac{L(s, \hat{\hat{\theta}}(s))}{L(\hat{s}, \hat{\theta})}$$

So the likelihood ratio statistic for testing s = 0 is

$$q_0 = -2\ln\frac{L(0)}{L(\hat{s})} = 2\left(n\ln\frac{n}{b} + b - n\right) \quad \text{for } n > b, \ 0 \text{ otherwise}$$

G. Cowan / RHUL Physics

# Approximate Poisson significance (continued)

For sufficiently large s + b, (use Wilks' theorem),

$$Z = \sqrt{2\left(n\ln\frac{n}{b} + b - n\right)} \quad \text{for } n > b \text{ and } Z = 0 \text{ otherwise.}$$

To find median[Z|s], let  $n \rightarrow s + b$  (i.e., the Asimov data set):

$$Z_{\rm A} = \sqrt{2\left(\left(s+b\right)\ln\left(1+\frac{s}{b}\right) - s\right)}$$

This reduces to  $s/\sqrt{b}$  for  $s \ll b$ .

G. Cowan / RHUL Physics

 $n \sim \text{Poisson}(s+b)$ , median significance, assuming *s*, of the hypothesis s = 0

CCGV, EPJC 71 (2011) 1554, arXiv:1007.1727



"Exact" values from MC, jumps due to discrete data.

Asimov  $\sqrt{q_{0,A}}$  good approx. for broad range of *s*, *b*.

 $s/\sqrt{b}$  only good for  $s \ll b$ .

G. Cowan / RHUL Physics

# Finally

#### Two lectures only enough for a brief introduction to:

- Parameter estimation
- Hypothesis tests ( $\rightarrow$  path to Machine Learning)
- Limits (confidence intervals/regions)
- Systematics (nuisance parameters)
- Experimental sensitivity

Final thought: once the basic formalism is fixed, most of the work focuses on writing down the likelihood, e.g.,  $P(x|\theta)$ , and including in it enough parameters to adequately describe the data (true for both Bayesian and frequentist approaches) so often best to invest most of your time with it.

# Extra slides

Information inequality for N parameters Suppose we have estimated N parameters  $\theta = (\theta_1, ..., \theta_N)$ 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  $\hat{ heta}$  is  $V_{ij} = ext{cov}[\hat{ heta}_i, \hat{ heta}_j]$ 

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:

 $z^{\mathrm{T}}Mz \ge 0$  for all  $z \ne 0$ , diagonal elements  $\ge 0$ 

# 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

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

#### 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

$$\int_{x_{\mathsf{min}}}^{x_{\mathsf{max}}} f(x; lpha, eta) \, dx = \mathbf{1} \; .$$

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

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

G. Cowan / RHUL Physics

TAE 2022 Benasque (online) / Lecture day 2

 $\widehat{x}$ 

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

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

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

0.11



х

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

 $\hat{\sigma}_{\hat{\alpha}} = 0.052 \quad \operatorname{cov}[\hat{\alpha}, \hat{\beta}] = 0.0026$ 

r = 0.46 = correlation coefficient

#### Two-parameter fit: MC study

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



G. Cowan / RHUL Physics

# Multiparameter graphical method for variances

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

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

$$\int_{\ln L_{\max}} zero relate to covariance matrix of$$

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

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

# Multiparameter graphical method (2)



#### Distance from MLE to tangent planes gives standard deviations.

# The $\ln L_{\rm max} - 1/2$ 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:

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

# (Co)variances from In L contour



→ Tangent lines to contours give standard deviations.

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

# Example: fitting a straight line

Data: 
$$(x_i, y_i, \sigma_i), i = 1, ..., n$$
.

Model:  $y_i$  independent and all follow  $y_i \sim \text{Gauss}(\mu(x_i), \sigma_i)$ 

 $\mu(x;\theta_0,\theta_1)=\theta_0+\theta_1x,$ 

assume  $x_i$  and  $\sigma_i$  known.

Goal: estimate  $\theta_0$ 

Here suppose we don't care about  $\theta_1$  (example of a "nuisance parameter")



### Maximum likelihood fit with Gaussian data

In this example, the  $y_i$  are assumed independent, so the likelihood function is a product of Gaussians:

$$L(\theta_0, \theta_1) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] ,$$

Maximizing the likelihood is here equivalent to minimizing

$$\chi^{2}(\theta_{0},\theta_{1}) = -2 \ln L(\theta_{0},\theta_{1}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i};\theta_{0},\theta_{1}))^{2}}{\sigma_{i}^{2}}.$$

i.e., for Gaussian data, ML same as Method of Least Squares (LS)

# $\theta_1$ known a priori

$$L(\theta_0) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right] .$$

 $\chi^{2}(\theta_{0}) = -2 \ln L(\theta_{0}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i}; \theta_{0}, \theta_{1}))^{2}}{\sigma_{i}^{2}}.$ 

For Gaussian  $y_i$ , ML same as LS

Minimize  $\chi^2 \rightarrow \text{estimator } \hat{\theta}_0$ . Come up one unit from  $\chi^2_{\min}$ to find  $\sigma_{\hat{\theta}_0}$ .



#### ML (or LS) fit of $\theta_0$ and $\theta_1$

$$\chi^{2}(\theta_{0},\theta_{1}) = -2 \ln L(\theta_{0},\theta_{1}) + \text{const} = \sum_{i=1}^{n} \frac{(y_{i} - \mu(x_{i};\theta_{0},\theta_{1}))^{2}}{\sigma_{i}^{2}}.$$

Standard deviations from tangent lines to contour

 $\chi^2 = \chi^2_{\rm min} + 1 \; .$ 

Correlation between  $\hat{\theta}_0, \ \hat{\theta}_1$  causes errors to increase.



If we have a measurement  $t_1 \sim \text{Gauss}(\theta_1, \sigma_{t_1})$ 

$$L(\theta_0, \theta_1) = \frac{1}{\sqrt{2\pi\sigma_t}} e^{-(t_1 - \theta_1)^2 / 2\sigma_{t_1}^2} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{1}{2} \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2}\right]$$

$$\chi^2(\theta_0, \theta_1) = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} + \frac{(t_1 - \theta_1)^2}{\sigma_{t_1}^2}$$

The information on  $\theta_1$ improves accuracy of  $\hat{\theta}_0$ .



# **Reminder of Bayesian approach**

In Bayesian statistics we can associate a probability with a hypothesis, e.g., a parameter value  $\theta$ .

Interpret probability of  $\theta$  as 'degree of belief' (subjective). Need to start with 'prior pdf'  $\pi(\theta)$ , this reflects degree of belief about  $\theta$  before doing the experiment.

Our experiment has data x,  $\rightarrow$  likelihood  $L(x|\theta)$ .

Bayes' theorem tells how our beliefs should be updated in light of the data *x*:

$$p(\theta|x) = \frac{L(x|\theta)\pi(\theta)}{\int L(x|\theta')\pi(\theta') d\theta'} \propto L(x|\theta)\pi(\theta)$$

Posterior pdf  $p(\theta|x)$  contains all our knowledge about  $\theta$ .

Bayesian approach:  $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$ We need to associate prior probabilities with  $\theta_0$  and  $\theta_1$ , e.g.,

 $\pi(\theta_0, \theta_1) = \pi_0(\theta_0)\pi_1(\theta_1) \quad \leftarrow \text{suppose knowledge of } \theta_0 \text{ has}$ no influence on knowledge of  $\theta_1$ 

$$\pi_0(\theta_0) = \text{const.} \qquad \leftarrow \text{`non-informative', in any} \\ \text{case much broader than } L(\theta_0)$$

$$\pi_{1}(\theta_{1}) = p(\theta_{1}|t_{1}) \propto p(t_{1}|\theta_{1})\pi_{\mathrm{Ur}}(\theta_{1}) = \frac{1}{\sqrt{2\pi}\sigma_{t}}e^{-(t_{1}-\theta_{1})^{2}/2\sigma_{t}^{2}} \times \mathrm{const.}$$
prior after  $t_{1}$ , Ur = "primordial" Likelihood for control before  $y$  prior measurement  $t_{1}$ 

Bayesian example:  $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$ 

Putting the ingredients into Bayes' theorem gives:



Note here the likelihood only reflects the measurements *y*.

The information from the control measurement  $t_1$  has been put into the prior for  $\theta_1$ .

We would get the same result using the likelihood  $P(y,t|\theta_0,\theta_1)$  and the constant "Ur-prior" for  $\theta_1$ .

# Marginalizing the posterior pdf

We then integrate (marginalize)  $p(\theta_0, \theta_1 | \mathbf{y})$  to find  $p(\theta_0 | \mathbf{y})$ :

$$p(\theta_0|\mathbf{y}) = \int p(\theta_0, \theta_1|\mathbf{y}) \, d\theta_1$$

In this example we can do the integral (rare). We find

$$p(\theta_0|\mathbf{y}) = \frac{1}{\sqrt{2\pi}\sigma_{\theta_0}} e^{-(\theta_0 - \hat{\theta}_0)^2/2\sigma_{\theta_0^2}}$$

 $\hat{\theta}_0 = \text{same as MLE}$ 

 $\sigma_{\theta_0} = \sigma_{\hat{\theta}_0}$  (same as for MLE)

For this example, numbers come out same as in frequentist approach, but interpretation different.

G. Cowan / RHUL Physics

TAE 2022 Benasque (online) / Lecture day 2

# Marginalization with MCMC

Bayesian computations involve integrals like

$$p(\theta_0|x) = \int p(\theta_0, \theta_1|x) d\theta_1$$
.

often high dimensionality and impossible in closed form, also impossible with 'normal' acceptance-rejection Monte Carlo.

Markov Chain Monte Carlo (MCMC) has revolutionized Bayesian computation.

MCMC (e.g., Metropolis-Hastings algorithm) generates correlated sequence of random numbers:

cannot use for many applications, e.g., detector MC; effective stat. error greater than if all values independent .

Basic idea: sample multidimensional  $\theta$  but look only at distribution of parameters of interest.

MCMC basics: Metropolis-Hastings algorithm Goal: given an *n*-dimensional pdf  $p(\theta)$ , generate a sequence of points  $\theta_1, \theta_2, \theta_3, \dots$ 

- 1) Start at some point  $\vec{\theta}_0$
- 2) Generate  $\vec{\theta} \sim q(\vec{\theta}; \vec{\theta}_0)$

Proposal density  $q(\theta; \theta_0)$ e.g. Gaussian centred about  $\theta_0$ 

3) Form Hastings test ratio  $\alpha = \min \left[ 1, \frac{\pi}{n} \right]$ 

$$1, \frac{p(\vec{\theta})q(\vec{\theta}_{0};\vec{\theta})}{p(\vec{\theta}_{0})q(\vec{\theta};\vec{\theta}_{0})} \bigg]$$

- 4) Generate  $u \sim \text{Uniform}[0, 1]$
- 5) If  $u \le \alpha$ ,  $\vec{\theta_1} = \vec{\theta}$ ,  $\leftarrow$  move to proposed point else  $\vec{\theta_1} = \vec{\theta_0} \leftarrow$  old point repeated 6) Iterate

# Metropolis-Hastings (continued)

This rule produces a *correlated* sequence of points (note how each new point depends on the previous one).

Still works if  $p(\theta)$  is known only as a proportionality, which is usually what we have from Bayes' theorem:  $p(\theta|\mathbf{x}) \propto p(\mathbf{x}|\theta)\pi(\theta)$ .

The proposal density can be (almost) anything, but choose so as to minimize autocorrelation. Often take proposal density symmetric:  $q(\theta; \theta_0) = q(\theta_0; \theta)$ 

Test ratio is (*Metropolis*-Hastings):  $\alpha = \min \left[ 1, \frac{p(\theta)}{p(\vec{\theta}_0)} \right]$ 

I.e. if the proposed step is to a point of higher  $p(\theta)$ , take it; if not, only take the step with probability  $p(\theta)/p(\theta_0)$ . If proposed step rejected, repeat the current point.

#### Example: posterior pdf from MCMC

#### Sample the posterior pdf from previous example with MCMC:


## Bayesian method with alternative priors

Suppose we don't have a previous measurement of  $\theta_1$  but rather, an "expert" says it should be positive and not too much greater than 0.1 or so, i.e., something like

$$\pi_1(\theta_1) = \frac{1}{\tau} e^{-\theta_1/\tau} , \quad \theta_1 \ge 0 , \quad \tau = 0.1 .$$

From this we obtain (numerically) the posterior pdf for  $\theta_0$ :

