## **Topics in Statistics for Particle Physics**



Discussion on Statistics LAL Orsay 16 June 2014



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

## Outline

0) Brief review of statistical tests and setting limits.

1) A measure of discovery sensitivity is often used to plan a future analysis, e.g.,  $s/\sqrt{b}$ , gives approximate expected discovery significance (test of s = 0) when counting  $n \sim \text{Poisson}(s+b)$ . A measure of discovery significance is proposed that takes into account uncertainty in the background rate.

2) In many searches for new signal processes, estimates of rates of some background components often based on Monte Carlo with weighted events. Some care (and assumptions) are required to assess the effect of the finite MC sample on the result of the test.

- 3) A few words on the jackknife and bootstrap.
- 4) A few words on Bayesian vs. Frequentist methods.

#### (Frequentist) statistical tests

Consider test of a parameter  $\mu$ , e.g., proportional to cross section. Result of measurement is a set of numbers *x*.

To define test of  $\mu$ , specify *critical region*  $w_{\mu}$ , such that probability to find  $x \in w_{\mu}$  is not greater than  $\alpha$  (the *size* or *significance level*):

 $P(\mathbf{x} \in w_{\mu}|\mu) \le \alpha$ 

(Must use inequality since x may be discrete, so there may not exist a subset of the data space with probability of exactly  $\alpha$ .)

Equivalently define a *p*-value  $p_{\mu}$  such that the critical region corresponds to  $p_{\mu} < \alpha$ .

Often use, e.g.,  $\alpha = 0.05$ .

If observe  $x \in w_{\mu}$ , reject  $\mu$ .

#### Test statistics and *p*-values

Often construct a test statistic,  $q_{\mu}$ , which reflects the level of agreement between the data and the hypothesized value  $\mu$ .

For examples of statistics based on the profile likelihood ratio, see, e.g., CCGV, EPJC 71 (2011) 1554; arXiv:1007.1727.

Usually define  $q_{\mu}$  such that higher values represent increasing incompatibility with the data, so that the *p*-value of  $\mu$  is:



Equivalent formulation of test: reject  $\mu$  if  $p_{\mu} < \alpha$ .

G. Cowan

Confidence interval from inversion of a test

Carry out a test of size  $\alpha$  for all values of  $\mu$ .

The values that are not rejected constitute a *confidence interval* for  $\mu$  at confidence level CL =  $1 - \alpha$ .

The confidence interval will by construction contain the true value of  $\mu$  with probability of at least  $1 - \alpha$ .

The interval depends on the choice of the critical region of the test. Put critical region where data are likely to be under assumption of the relevant alternative to the  $\mu$  that's being tested.

Test  $\mu = 0$ , alternative is  $\mu > 0$ : test for discovery.

Test  $\mu = \mu_0$ , alternative is  $\mu = 0$ : testing all  $\mu_0$  gives upper limit.

#### *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

$$p_0 = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0) \, dq_0$$

will get formula for this later



From *p*-value get equivalent significance,

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

#### Significance from *p*-value

Often define significance Z as the number of standard deviations that a Gaussian variable would fluctuate in one direction to give the same p-value.



$$p = \int_{Z}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \, dx = 1 - \Phi(Z)$$
 1 - TMath::Freq

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

#### 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

$$E[n_i] = \mu s_i + b_i$$
  
strength parameter  
$$f_i(m; \theta_i) dm = b_i = \int_{-\infty}^{-\infty} f_i(m; \theta_i) dm$$

where

$$s_{i} = s_{\text{tot}} \int_{\text{bin } i} f_{s}(x; \boldsymbol{\theta}_{s}) \, dx \,, \quad b_{i} = b_{\text{tot}} \int_{\text{bin } i} f_{b}(x; \boldsymbol{\theta}_{b}) \, dx \,.$$
  
signal background

## 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

## The profile likelihood ratio

Base significance test on the profile likelihood ratio:



The likelihood ratio of point hypotheses gives optimum test (Neyman-Pearson lemma).

The profile LR hould be near-optimal in present analysis with variable  $\mu$  and nuisance parameters  $\theta$ .

#### Test statistic for discovery

Try to reject background-only ( $\mu = 0$ ) hypothesis using

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

i.e. here only regard upward fluctuation of data as 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  $\sigma$ .

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

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

#### Test statistic for upper limits

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

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

95% CL upper limit on  $\mu$  is highest value for which *p*-value is not less than 0.05.

G. Cowan

# Example of a *p*-value ATLAS, Phys. Lett. B 716 (2012) 1-29



Orsay 2014 / Discussion on Statistics

Expected (or median) significance / sensitivity

When planning the experiment, we want to quantify how sensitive we are to a potential discovery, e.g., by given median significance assuming some nonzero strength parameter  $\mu'$ .



So for *p*-value, need  $f(q_0|0)$ , for sensitivity, will need  $f(q_0|\mu')$ ,

G. Cowan

Expected discovery significance for counting experiment with background uncertainty

I. Discovery sensitivity for counting experiment with *b* known:

(a) 
$$\frac{s}{\sqrt{b}}$$

(b) Profile likelihood ratio test & Asimov:

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

II. Discovery sensitivity with uncertainty in b,  $\sigma_b$ :

(a) 
$$\frac{s}{\sqrt{b+\sigma_b^2}}$$

(b) Profile likelihood ratio test & Asimov:

$$\left[2\left((s+b)\ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2}\ln\left[1 + \frac{\sigma_b^2s}{b(b+\sigma_b^2)}\right]\right)\right]^{1/2}$$

Counting experiment with known background Count a number of events  $n \sim Poisson(s+b)$ , 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

 $s/\sqrt{b}$  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

Better approximation for significance Poisson likelihood for parameter *s* is

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

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

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

Orsay 2014 / Discussion on Statistics

params.

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 << b.

 $n \sim \text{Poisson}(s+b)$ , 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  $s \ll b$ .

Orsay 2014 / Discussion on Statistics

## Extending $s/\sqrt{b}$ to case where b uncertain

The intuitive explanation of  $s/\sqrt{b}$  is that it compares the signal, *s*, to the standard deviation of *n* assuming no signal,  $\sqrt{b}$ .

Now suppose the value of *b* is uncertain, characterized by a standard deviation  $\sigma_b$ .

A reasonable guess is to replace  $\sqrt{b}$  by the quadratic sum of  $\sqrt{b}$  and  $\sigma_b$ , i.e.,

$$\operatorname{med}[Z|s] = \frac{s}{\sqrt{b + \sigma_b^2}}$$

This has been used to optimize some analyses e.g. where  $\sigma_b$  cannot be neglected.

#### Profile likelihood with b uncertain

This is the well studied "on/off" problem: Cranmer 2005; Cousins, Linnemann, and Tucker 2008; Li and Ma 1983,...

Measure two Poisson distributed values:

 $n \sim \text{Poisson}(s+b)$ (primary or "search" measurement) $m \sim \text{Poisson}(\tau b)$ (control measurement,  $\tau$  known)

The likelihood function is

$$L(s,b) = \frac{(s+b)^n}{n!} e^{-(s+b)} \frac{(\tau b)^m}{m!} e^{-\tau b}$$

Use this to construct profile likelihood ratio (*b* is nuisance parmeter):  $L(0, \hat{b}(0))$ 

$$\lambda(0) = \frac{L(0, b(0))}{L(\hat{s}, \hat{b})}$$

G. Cowan

#### Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

$$\begin{split} \hat{s} &= n - m/\tau \ , \\ \hat{b} &= m/\tau \ , \\ \hat{b}(s) &= \frac{n + m - (1 + \tau)s + \sqrt{(n + m - (1 + \tau)s)^2 + 4(1 + \tau)sm}}{2(1 + \tau)} \end{split}$$

and in particular to test for discovery (s = 0),

$$\hat{\hat{b}}(0) = \frac{n+m}{1+\tau}$$

#### Asymptotic significance

Use profile likelihood ratio for  $q_0$ , and then from this get discovery significance using asymptotic approximation (Wilks' theorem):

$$Z = \sqrt{q_0}$$
$$= \left[ -2\left(n\ln\left[\frac{n+m}{(1+\tau)n}\right] + m\ln\left[\frac{\tau(n+m)}{(1+\tau)m}\right]\right) \right]^{1/2}$$

for  $n > \hat{b}$  and Z = 0 otherwise.

#### Essentially same as in:

Robert D. Cousins, James T. Linnemann and Jordan Tucker, NIM A 595 (2008) 480– 501; arXiv:physics/0702156.

Tipei Li and Yuqian Ma, Astrophysical Journal 272 (1983) 317–324.

#### Asimov approximation for median significance

To get median discovery significance, replace *n*, *m* by their expectation values assuming background-plus-signal model:

$$n \to s + b$$
  

$$m \to \tau b$$

$$Z_{A} = \left[ -2\left( (s+b) \ln\left[\frac{s+(1+\tau)b}{(1+\tau)(s+b)}\right] + \tau b \ln\left[1+\frac{s}{(1+\tau)b}\right] \right) \right]^{1/2}$$
Or use the variance of  $\hat{b} = m/\tau$ ,  $V[\hat{b}] \equiv \sigma_{b}^{2} = \frac{b}{\tau}$ , to eliminate  $\tau$ :  

$$A = \left[ 2\left( (s+b) \ln\left[\frac{(s+b)(b+\sigma_{b}^{2})}{b^{2}+(s+b)\sigma_{b}^{2}}\right] - \frac{b^{2}}{\sigma_{b}^{2}} \ln\left[1+\frac{\sigma_{b}^{2}s}{b(b+\sigma_{b}^{2})}\right] \right) \right]^{1/2}$$

 $Z_{i}$ 

#### Limiting cases

Expanding the Asimov formula in powers of *s/b* and  $\sigma_b^2/b$  (= 1/ $\tau$ ) gives

$$Z_{\rm A} = \frac{s}{\sqrt{b + \sigma_b^2}} \left( 1 + \mathcal{O}(s/b) + \mathcal{O}(\sigma_b^2/b) \right)$$

So the "intuitive" formula can be justified as a limiting case of the significance from the profile likelihood ratio test evaluated with the Asimov data set. Testing the formulae: s = 5



Orsay 2014 / Discussion on Statistics

#### Using sensitivity to optimize a cut



Figure 1: (a) The expected significance as a function of the cut value  $x_{\text{cut}}$ ; (b) the distributions of signal and background with the optimal cut value indicated.

#### G. Cowan

#### Summary on discovery sensitivity

Simple formula for expected discovery significance based on profile likelihood ratio test and Asimov approximation:

$$Z_{\rm A} = \left[ 2 \left( (s+b) \ln \left[ \frac{(s+b)(b+\sigma_b^2)}{b^2 + (s+b)\sigma_b^2} \right] - \frac{b^2}{\sigma_b^2} \ln \left[ 1 + \frac{\sigma_b^2 s}{b(b+\sigma_b^2)} \right] \right) \right]^{1/2}$$

For large *b*, all formulae OK.

For small *b*,  $s/\sqrt{b}$  and  $s/\sqrt{(b+\sigma_b^2)}$  overestimate the significance.

Could be important in optimization of searches with low background.

Formula maybe also OK if model is not simple on/off experiment, e.g., several background control measurements (checking this).

#### Using MC events in a statistical test

**Prototype analysis** – count *n* events where signal may be present:

 $n \sim \text{Poisson}(\mu s + b)$ 

s = expected events from nominal signal model (regard as known) b = expected background (nuisance parameter)

 $\mu$  = strength parameter (parameter of interest)

**Ideal** – constrain background *b* with a data control measurement *m*, scale factor  $\tau$  (assume known) relates control and search regions:

 $m \sim \text{Poisson}(\tau b)$ 

**Reality** - not always possible to construct data control sample, sometimes take prediction for *b* from MC.

From a statistical perspective, can still regard number of MC events found as  $m \sim \text{Poisson}(\tau b)$  (really should use binomial, but here Poisson good approx.) Scale factor is  $\tau = L_{\text{MC}}/L_{\text{data}}$ .

#### MC events with weights

But, some MC events come with an associated weight, either from generator directly or because of reweighting for efficiency, pile-up. Outcome of experiment is: *n*, *m*, *w*<sub>1</sub>,..., *w*<sub>m</sub>
How to use this info to construct statistical test of μ?
"Usual" (?) method is to construct an estimator for *b*:

$$\hat{b} = \frac{1}{\tau} \sum_{i=1}^{m} w_i$$
  $\hat{\sigma}_{\hat{b}}^2 = \frac{1}{\tau^2} \sum_{i=1}^{m} w_i^2$ 

and include this with a least-squares constraint, e.g., the  $\chi^2$  gets an additional term like

$$\frac{(b-b)^2}{\hat{\sigma}_{\hat{b}}^2}$$

#### Case where *m* is small (or zero)

Using least-squares like this assumes  $\hat{b} \sim$  Gaussian, which is OK for sufficiently large *m* because of the Central Limit Theorem. But  $\hat{b}$  may not be Gaussian distributed if e.g.

*m* is very small (or zero),

the distribution of weights has a long tail.

Hypothetical example:

$$m = 2, w_1 = 0.1307, w_2 = 0.0001605,$$
  
 $\hat{b} = 0.0007 \pm 0.0030$   
 $n = 1$  (!)

Correct procedure is to treat  $m \sim \text{Poisson}$  (or binomial). And if the events have weights, these constitute part of the measurement, and so we need to make an assumption about their distribution.

#### Constructing a statistical test of $\mu$

As an example, suppose we want to test the background-only hypothesis ( $\mu$ =0) using the profile likelihood ratio statistic (see e.g. CCGV, EPJC 71 (2011) 1554, arXiv:1007.1727),

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \ge 0\\ 0 & \hat{\mu} < 0 \end{cases} \quad \text{where} \quad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})} \end{cases}$$

From the observed value of  $q_0$ , the *p*-value of the hypothesis is:

$$p = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0) \, dq_0$$

So we need to know the distribution of the data  $(n, m, w_1, ..., w_m)$ , i.e., the likelihood, in two places:

1) to define the likelihood ratio for the test statistic

2) for  $f(q_0|0)$  to get the *p*-value

#### Normal distribution of weights

Suppose  $w \sim \text{Gauss}(\omega, \sigma_w)$ . The full likelihood function is

$$L(\mu, b, \omega, \sigma_w) = \frac{(\mu s + b)^n}{n!} e^{-(\mu s + b)} \frac{(\tau b/\omega)^m}{m!} e^{-\tau b/\omega} \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma_w} e^{(w_i - \omega)^2/2\sigma_w^2}$$

The log-likelihood can be written:

$$\ln L(\mu, b, \omega, \sigma_w) = n \ln(\mu s + b) - (\mu s + b) + m \ln(\tau b/\omega) - \tau b/\omega$$
$$- m \ln \sigma_w - \frac{m\omega^2}{2\sigma_w^2} + \frac{\omega}{\sigma_w^2} \sum_{i=1}^m w_i - \frac{1}{2\sigma_w^2} \sum_{i=1}^m w_i^2 + C$$

Only depends on weights through: 
$$S_1 = \sum_{i=1}^m w_i$$
,  $S_2 = \sum_{i=1}^m w_i^2$ .
#### Log-normal distribution for weights

Depending on the nature/origin of the weights, we may know:  $w(x) \ge 0$ ,

distribution of w could have a long tail.

So  $w \sim \log$ -normal could be a more realistic model.

I.e, let  $l = \ln w$ , then  $l \sim \text{Gaussian}(\lambda, \sigma_l)$ , and the log-likelihood is

$$\ln L(\mu, b, \lambda, \sigma_l) = n \ln(\mu s + b) - (\mu s + b) + m \ln(\tau b/\omega) - \tau b/\omega$$
$$- m \ln \sigma_l - \frac{m\lambda^2}{2\sigma_l^2} + \frac{\lambda}{\sigma_l^2} \sum_{i=1}^m l_i - \frac{1}{2\sigma_l^2} \sum_{i=1}^m l_i^2.$$

where  $\lambda = E[l]$  and  $\omega = E[w] = \exp(\lambda + \sigma_l^2/2)$ . Need to record *n*, *m*,  $\Sigma_i \ln w_i$  and  $\Sigma_i \ln^2 w_i$ .

# Normal distribution for $\hat{b}$

For m > 0 we can define the estimator for b

$$\hat{b} = \frac{1}{\tau} \sum_{i=1}^{m} w_i$$
  $\hat{\sigma}_{\hat{b}}^2 = \frac{1}{\tau^2} \sum_{i=1}^{m} w_i^2$ 

If we assume  $\hat{b} \sim$  Gaussian, then the log-likelihood is

$$\ln L(\mu, b) = n \ln(\mu s + b) - (\mu s + b) - \frac{1}{2} \frac{(b - \hat{b})^2}{\hat{\sigma}_{\hat{b}}^2}$$

Important simplification: L only depends on parameter of interest  $\mu$  and single nuisance parameter b.

Ordinarily would only use this Ansatz when Prob(*m*=0) negligible.

## Toy weights for test of procedure

Suppose we wanted to generate events according to

$$f(x) = \frac{e^{-x/\xi}}{\xi(1 - e^{-a/\xi})}, \quad 0 \le x \le a.$$

Suppose we couldn't do this, and only could generate x following

$$g(x) = \frac{1}{a} , \quad 0 \le x \le a$$

and for each event we also obtain a weight

$$w(x) = \frac{f(x)}{g(x)} = \frac{a}{\xi} \frac{e^{-x/\xi}}{1 - e^{-a/\xi}}$$

In this case the weights follow:



 $w_{\min} \leq w \leq w_{\max}$ 

G. Cowan

#### Two sample MC data sets

Suppose $n = 17$ , $\tau = 1$ , and		
	weight $w$	$\ln w$
case 1:	0.9684	-0.0320
$a = 5 \ \mathcal{E} = 25$	0.9217	-0.0816
x = 5, 5, 25	1.0238	0.0235
m = 0	1.0063	0.0063
Distribution of w narrow	0.9709	-0.0295
	1.0813	0.0782
	weight $w$	$\ln w$
case 2:	0.1934	-1.6429
$a = 5,  \xi = 1$	0.0561	-2.8809
n = 6	0.7750	-0.2548
Distribution of w broad	0.5039	-0.6853
	0.2059	-1.580
	3.0404	1.1120

#### G. Cowan

## Testing $\mu = 0$ using $q_0$ with n = 17

4	Likelihood used	Distribution of	Significance $Z$
case 1:	to define $q_0$	$w$ for $f(q_0 0)$	to reject $\mu = 0$
$a = 5, \xi = 25$	$w \sim \text{normal}$	normal	2.287
, <u>,</u>	$w \sim \text{normal}$	1/w	2.268
m = 0	$w \sim \log$ -normal	log-normal	2.301
Distribution of	$w \sim \log$ -normal	1/w	2.267
w is narrow	$\hat{b} \sim \mathrm{normal}$	normal	2.289
// 10 HWH 0 W	$\hat{b} \sim \mathrm{normal}$	1/w	2.224

If distribution of weights is narrow, then all methods result in a similar picture: discovery significance  $Z \sim 2.3$ .

Testing $\mu = 0$ using	$q_0$ with $n = 17$	' (cont.)
-------------------------	---------------------	-----------

	Likelihood used	Distribution of	Significance Z
2022 7.	Likelihood used	Distribution of	Significance Z
case 2.	to define $q_0$	w for $f(q_0 0)$	to reject $\mu = 0$
$a = 5, \xi = 1$	$w \sim \text{normal}$	normal	2.163
, , ,	$w \sim \text{normal}$	1/w	1.308
m = 0	$w \sim \log$ -normal	log-normal	0.863
Distribution of	$w \sim \log$ -normal	1/w	0.983
w is broad	$\hat{b} \sim \mathrm{normal}$	normal	1.788
W 15 01000	$\hat{b} \sim \mathrm{normal}$	1/w	1.387

If there is a broad distribution of weights, then:

- 1) If true  $w \sim 1/w$ , then assuming  $w \sim$  normal gives too tight of constraint on *b* and thus overestimates the discovery significance.
- 2) If test statistic is sensitive to tail of *w* distribution (i.e., based on log-normal likelihood), then discovery significance reduced.

Best option above would be to assume  $w \sim \text{log-normal}$ , both for definition of  $q_0$  and  $f(q_0|0)$ , hence Z = 0.863.

#### Case of m = 0

If no MC events found (m = 0) then there is no information with which to estimate the variance of the weight distribution, so the method with  $\hat{b} \sim$  Gaussian  $(b, \sigma_b)$  cannot be used.

For both normal and log-normal distributions of the weights, the likelihood function becomes

$$\ln L(\mu, b, \omega) = n \ln(\mu s + b) - (\mu s + b) - \frac{\tau b}{\omega}$$

If mean weight  $\omega$  is known (e.g.,  $\omega = 1$ ), then the only nuisance parameter is *b*. Use as before profile likelihood ratio to test  $\mu$ .

If  $\omega$  is not known, then maximizing  $\ln L$  gives  $\omega \to \infty$ , no inference on  $\mu$  possible.

If upper bound on  $\omega$  can be used, this gives conservative estimate of significance for test of  $\mu = 0$ .

#### Case of m = 0, test of $\mu = 0$

Asymptotic approx. for test of  $\mu = 0$  (Z =  $\sqrt{q_0}$ ) results in:

$$Z = \sqrt{2n\ln\left(1 + \frac{\tau}{\omega}\right)}$$

Example for n = 5, m = 0,  $\omega = 1$ 



Orsay 2014 / Discussion on Statistics

### Summary on weighted MC

Treating MC data as "real" data, i.e.,  $n \sim$  Poisson, incorporates the statistical error due to limited size of sample.

Then no problem if zero MC events observed, no issue of how to deal with  $0 \pm 0$  for background estimate.

If the MC events have weights, then some assumption must be made about this distribution.

If large sample, Gaussian should be OK,

if sample small consider log-normal.

See draft note for more info and also treatment of weights =  $\pm 1$  (e.g., MC@NLO).

www.pp.rhul.ac.uk/~cowan/stat/notes/weights.pdf

## Jackknife, bootstrap, etc.

To estimate a parameter we have various tools such as maximum likelihood, least squares, etc.



Usually one also needs to know the variance (or the full sampling distribution) of the estimator – this can be more difficult.

Often use asymptotic properties, e.g., sampling distribution of ML estimators becomes Gaussian in large sample limit; std. dev. from curvature of log-likelihood at maximum.

The jackknife and bootstrap are examples of "resampling" methods used to estimate the sampling distribution of statistics.

In HEP we often do this implicitly by using Toy MC to determine sampling properties of statistics (e.g., Brazil plot for  $1\sigma$ ,  $2\sigma$  bands of limits).

G. Cowan

### The Jackknife

Invented by Quenouille (1949) and Tukey (1958).

Suppose data sample consists of *n* events:  $\mathbf{x} = (x_1, \dots, x_n)$ . We have an estimator  $\hat{\theta}(\mathbf{x})$  for a parameter  $\theta$ .

Idea is to produce pseudo data samples  $x_{-i} = (x_1, ..., x_{i-1}, x_{i+1}, ..., x_n)$ by leaving out the *i*th event.

Let  $\hat{\theta}_{-1}$  be the estimator obtained from the data sample  $x_{-i}$ .

Suppose the estimator has a nonzero bias:  $b = E[\hat{\theta}] - \theta$ 

The jackknife estimator  
of the bias is 
$$\hat{b}_{j} = \frac{n-1}{n} \sum_{i=1}^{n} (\hat{\theta}_{-i} - \hat{\theta})$$

See, e.g., Notes on Jackknife and Bootstrap by G. J. Babu: www.iiap.res.in/astrostat/School10/LecFiles/ JBabu\_JackknifeBootstrap\_notes.pdf

#### G. Cowan

#### The Bootstrap (Efron, 1979)

Idea is to produce a set of "bootstrapped" data samples of same size as the original (real) one by sampling from some distribution that approximates the true (unknown) one.

By evaluating a statistic (such as an estimator for a parameter  $\theta$ ) with the bootstrapped-samples, properties of its sampling distribution (often its variance) can be estimated.

If the data consist of *n* events, one way to produce the bootstrapped samples is to randomly select from the original sample *n* events *with replacement* (the non-parametric bootstrap).

That is, some events might get used multiple times, others might not get used at all.

In other cases could generate the bootstrapped samples from a parametric MC model, using parameter values estimated from real data in the MC (parametric bootstrap).

#### The Bootstrap (cont.)

Call the data sample  $x = (x_1, ..., x_n)$ , observed data are  $x_{obs}$ , and the bootstrapped samples are  $x_1^*, x_2^*, ...$ Idea is to use the distribution of

$$\hat{ heta}(\mathbf{x}^*) - \hat{ heta}(\mathbf{x}_{\mathrm{obs}})$$

as an approximation for the distribution of

$$\hat{ heta}(\mathbf{x}) - heta$$

In the first quantity everything is known from the observed data plus bootstrapped samples, so we can use its distribution to estimate bias, variance, etc. of the estimator  $\hat{\theta}$ .

G. Cowan

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



Can improve model by including additional adjustable parameters.

$$L(x|\theta) \to L(x|\theta,\nu)$$

Nuisance parameter  $\leftrightarrow$  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).

*p*-values in cases with nuisance parameters Suppose we have a statistic  $q_{\theta}$  that we use to test a hypothesized value of a parameter  $\theta$ , such that the *p*-value of  $\theta$  is

$$p_{\theta} = \int_{q_{\theta,\text{obs}}}^{\infty} f(q_{\theta}|\theta,\nu) \, dq_{\theta}$$

But what values of *v* to use for  $f(q_{\theta}|\theta, v)$ ? Fundamentally we want to reject  $\theta$  only if  $p_{\theta} < \alpha$  for all *v*.

 $\rightarrow$  "exact" confidence interval

Recall that for statistics based on the profile likelihood ratio, the distribution  $f(q_{\theta}|\theta, v)$  becomes independent of the nuisance parameters in the large-sample limit.

But in general for finite data samples this is not true; one may be unable to reject some  $\theta$  values if all values of v must be considered, even those strongly disfavoured by the data (resulting interval for  $\theta$  "overcovers").

# Profile construction ("hybrid resampling")

K. Cranmer, PHYSTAT-LHC Workshop on Statistical Issues for LHC Physics, 2008. oai:cds.cern.ch:1021125, cdsweb.cern.ch/record/1099969.

Approximate procedure is to reject  $\theta$  if  $p_{\theta} \le \alpha$  where the *p*-value is computed assuming the value of the nuisance parameter that best fits the data for the specified  $\theta$ :

<u>^</u>	"double hat" notation means
$\hat{ u}( heta)$	value of parameter that maximizes
	likelihood for the given $\theta$ .

The resulting confidence interval will have the correct coverage for the points  $(\theta, \hat{v}(\theta))$ .

Elsewhere it may under- or overcover, but this is usually as good as we can do (check with MC if crucial or small sample problem).

# "Hybrid frequentist-Bayesian" method

Alternatively, suppose uncertainty in v is characterized by a Bayesian prior  $\pi(v)$ .

Can use the marginal likelihood to model the data:

$$L_{\rm m}(x|\theta) = \int L(x|\theta,\nu)\pi(\nu) \, d\nu$$

This does not represent what the data distribution would be if we "really" repeated the experiment, since then *v* would not change.

But the procedure has the desired effect. The marginal likelihood effectively builds the uncertainty due to v into the model.

Use this now to compute (frequentist) *p*-values  $\rightarrow$  the model being tested is in effect a weighted average of models.

Example of treatment of nuisance parameters: 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  $\theta_l$  (example of a "nuisance parameter")



Orsay 2014 / Discussion on Statistics

#### 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  $t_1 \sim \text{Gauss}(\theta_1, \sigma_{t_1})$ 

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



G. Cowan

Orsay 2014 / Discussion on Statistics

### Bayesian method

We need to associate prior probabilities with  $\theta_0$  and  $\theta_1$ , e.g.,

$\pi(\theta_0,\theta_1)$	=	$\pi_0(\theta_0)$	$\pi_1(\theta_1)$	'non-i	nformative', in any
$\pi_0(\theta_0)$	=	const.		case n	nuch broader than $L(\theta_0)$
$\pi_1(\theta_1)$	=	$\frac{1}{\sqrt{2\pi}\sigma_{t_1}}$	$-e^{-(\theta_1-t_1)}$	$^{2}/2\sigma_{t_{1}}^{2}$	← based on previous measurement

Putting this into Bayes' theorem gives:

$$p(\theta_0, \theta_1 | \vec{y}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} e^{-(y_i - \mu(x_i; \theta_0, \theta_1))^2 / 2\sigma_i^2} \pi_0 \frac{1}{\sqrt{2\pi\sigma_{t_1}}} e^{-(\theta_1 - t_1)^2 / 2\sigma_{t_1}^2}$$

$$posterior \propto likelihood \times prior$$

# Bayesian method (continued)

We then integrate (marginalize)  $p(\theta_0, \theta_1 | x)$  to find  $p(\theta_0 | x)$ :

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

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

$$p(\theta_0|x) = \frac{1}{\sqrt{2\pi}\sigma_{\theta_0}} e^{-(\theta_0 - \hat{\theta}_0)^2 / 2\sigma_{\theta_0}^2} \text{ with}$$
$$\hat{\theta}_0 = \text{ same as ML estimator}$$
$$\sigma_{\theta_0} = \sigma_{\hat{\theta}_0} \text{ (same as before)}$$

Usually need numerical methods (e.g. Markov Chain Monte Carlo) to do integral.

Digression: 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  $\vec{\theta}$ , look, e.g., only at distribution of parameters of interest. MCMC basics: Metropolis-Hastings algorithm Goal: given an *n*-dimensional pdf  $p(\vec{\theta})$ , generate a sequence of points  $\vec{\theta}_1, \vec{\theta}_2, \vec{\theta}_3, \dots$ 

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

Proposal density  $q(\vec{\theta}; \vec{\theta}_0)$ e.g. Gaussian centred about  $\vec{\theta}_0$ 

3) Form Hastings test ratio  $\alpha = \min \left| 1 \right|$ 

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

- 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).

For our purposes this correlation is not fatal, but statistical errors larger than if points were independent.

The proposal density can be (almost) anything, but choose so as to minimize autocorrelation. Often take proposal density symmetric:  $q(\vec{\theta}; \vec{\theta}_0) = q(\vec{\theta}_0; \vec{\theta})$ 

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

I.e. if the proposed step is to a point of higher  $p(\vec{\theta})$ , take it; if not, only take the step with probability  $p(\vec{\theta})/p(\vec{\theta}_0)$ . If proposed step rejected, hop in place.

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



Summarize pdf of parameter of interest with, e.g., mean, median, standard deviation, etc.

Although numerical values of answer here same as in frequentist case, interpretation is different (sometimes unimportant?)

# Bayesian method with alternative priors

Suppose we don't have a previous measurement of  $\theta_1$  but rather, e.g., a theorist 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  $\theta_0$ :



Bayesian model selection ('discovery')

The probability of hypothesis  $H_0$  relative to its complementary alternative  $H_1$  is often given by the posterior odds:



The Bayes factor is regarded as measuring the weight of evidence of the data in support of  $H_0$  over  $H_1$ .

Interchangeably use  $B_{10} = 1/B_{01}$ 

Lecture 14 page 66

Assessing Bayes factors

One can use the Bayes factor much like a *p*-value (or Z value). The Jeffreys scale, analogous to HEP's  $5\sigma$  rule:

<i>B</i> <sub>10</sub>	Evidence against $H_0$
1 to 3	Not worth more than a bare mention
3 to 20	Positive
20 to 150	Strong
> 150	Very strong

Kass and Raftery, Bayes Factors, J. Am Stat. Assoc 90 (1995) 773.

#### Rewriting the Bayes factor

Suppose we have models  $H_i$ , i = 0, 1, ...,

each with a likelihood  $p(x|H_i, \vec{\theta_i})$ 

and a prior pdf for its internal parameters  $\pi_i(\vec{\theta_i})$ 

so that the full prior is  $\pi(H_i, \vec{\theta}_i) = p_i \pi_i(\vec{\theta}_i)$ 

where  $p_i = P(H_i)$  is the overall prior probability for  $H_i$ .

The Bayes factor comparing  $H_i$  and  $H_i$  can be written

$$B_{ij} = \frac{P(H_i|\vec{x})}{P(H_i)} / \frac{P(H_j|\vec{x})}{P(H_j)}$$

G. Cowan

Orsay 2014 / Discussion on Statistics

Lecture 14 page 68

### Bayes factors independent of $P(H_i)$

For  $B_{ij}$  we need the posterior probabilities marginalized over all of the internal parameters of the models:

$$P(H_i|\vec{x}) = \int P(H_i, \vec{\theta}_i | \vec{x}) d\vec{\theta}_i \qquad \text{Use Bayes} \\ = \frac{\int L(\vec{x} | H_i, \vec{\theta}_i) p_i \pi_i(\vec{\theta}_i) d\vec{\theta}_i}{P(x)}$$

So therefore the Bayes factor is

Ratio of marginal likelihoods

$$B_{ij} = \frac{\int L(\vec{x}|H_i, \vec{\theta}_i) \pi_i(\vec{\theta}_i) d\vec{\theta}_i}{\int L(\vec{x}|H_j, \vec{\theta}_j) \pi_j(\vec{\theta}_j) d\vec{\theta}_j}$$

The prior probabilities  $p_i = P(H_i)$  cancel.

G. Cowan

Orsay 2014 / Discussion on Statistics

Lecture 14 page 69

#### Numerical determination of Bayes factors

Both numerator and denominator of  $B_{ij}$  are of the form

Various ways to compute these, e.g., using sampling of the posterior pdf (which we can do with MCMC).

Harmonic Mean (and improvements) Importance sampling Parallel tempering (~thermodynamic integration) Nested Samplying (MultiNest), ...

Kass and Raftery, Bayes Factors, J. Am. Stat. Assoc. 90 (1995) 773-795.

Cong Han and Bradley Carlin, Markov Chain Monte Carlo Methods for Computing Bayes Factors: A Comparative Review, J. Am. Stat. Assoc. 96 (2001) 1122-1132.

Phil Gregory, *Bayesian Logical Data Analysis for the Physical Sciences*, Cambridge University Press, 2005.

#### G. Cowan

### Priors for Bayes factors

Note that for Bayes factors (unlike Bayesian limits), the prior cannot be improper. If it is, the posterior is only defined up to an arbitrary constant, and so the Bayes factor is ill defined

Possible exception allowed if both models contain *same* improper prior; but having same parameter name (or Greek letter) in both models does not fully justify this step.

If improper prior is made proper e.g. by a cut-off, the Bayes factor will retain a dependence on this cut-off.

In general for Bayes factors, all priors must reflect "meaningful" degrees of uncertainty about the parameters.

### Harmonic mean estimator

E.g., consider only one model and write Bayes theorem as:

$$\frac{\pi(\boldsymbol{\theta})}{m} = \frac{p(\boldsymbol{\theta}|\mathbf{x})}{L(\mathbf{x}|\boldsymbol{\theta})}$$

 $\pi(\theta)$  is normalized to unity so integrate both sides,

posterior expectation

$$m^{-1} = \int \frac{1}{L(\mathbf{x}|\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} = E_p[1/L]$$

Therefore sample  $\theta$  from the posterior via MCMC and estimate *m* with one over the average of 1/L (the harmonic mean of *L*).

M.A. Newton and A.E. Raftery, Approximate Bayesian Inference by the Weighted Likelihood Bootstrap, Journal of the Royal Statistical Society B 56 (1994) 3-48.

G. Cowan

Orsay 2014 / Discussion on Statistics

Lecture 14 page 72
## Improvements to harmonic mean estimator

The harmonic mean estimator is numerically very unstable; formally infinite variance (!). Gelfand & Dey propose variant:

Rearrange Bayes thm; multiply both sides by arbitrary  $pdf f(\theta)$ :

$$\frac{f(\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathbf{x})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} = \frac{f(\boldsymbol{\theta})}{m}$$

Integrate over 
$$\boldsymbol{\theta}$$
:  $m^{-1} = \int \frac{f(\boldsymbol{\theta})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} p(\boldsymbol{\theta}|\mathbf{x}) = E_p \left[ \frac{f(\boldsymbol{\theta})}{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})} \right]$ 

Improved convergence if tails of  $f(\theta)$  fall off faster than  $L(x|\theta)\pi(\theta)$ Note harmonic mean estimator is special case  $f(\theta) = \pi(\theta)$ .

A.E. Gelfand and D.K. Dey, *Bayesian model choice: asymptotics and exact calculations*, Journal of the Royal Statistical Society B 56 (1994) 501-514.

G. Cowan

Orsay 2014 / Discussion on Statistics

Lecture 14 page 73

Importance sampling

Need pdf  $f(\theta)$  which we can evaluate at arbitrary  $\theta$  and also sample with MC.

The marginal likelihood can be written

$$m = \int \frac{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{f(\boldsymbol{\theta})} f(\boldsymbol{\theta}) d\boldsymbol{\theta} = E_f \left[ \frac{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{f(\boldsymbol{\theta})} \right]$$

Best convergence when  $f(\theta)$  approximates shape of  $L(x|\theta)\pi(\theta)$ .

Use for  $f(\theta)$  e.g. multivariate Gaussian with mean and covariance estimated from posterior (e.g. with MINUIT).

## K. Cranmer/R. Trotta PHYSTAT 2011

## The nested sampling algorithm



(animation courtesy of David Parkinson)

An algorithm originally aimed primarily at the Bayesian evidence computation (Skilling, 2006):

$$X(\lambda) = \int_{\mathcal{L}(\theta) > \lambda} P(\theta) d\theta$$

$$P(d) = \int d\theta \mathcal{L}(\theta) P(\theta) = \int_0^1 X(\lambda) d\lambda$$

Feroz et al (2008), arxiv: 0807.4512, Trotta et al (2008), arxiv: 0809.3792

## G. Cowan