

# Statistical Data Analysis 2021/22

## Lecture Week 10



London Postgraduate Lectures on Particle Physics  
University of London MSc/MSci course PH4515



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

Course web page via RHUL moodle (PH4515) and also  
`www.pp.rhul.ac.uk/~cowan/stat_course.html`

# Statistical Data Analysis

## Lecture 10-1

- Upper limits on a Poisson rate parameter
  - Frequentist approach
  - Bayesian approach

# 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  $P(n \leq n_{\text{obs}}; s, b)$

Upper limit  $s_{\text{up}}$  at  $\text{CL} = 1 - \alpha$  found from

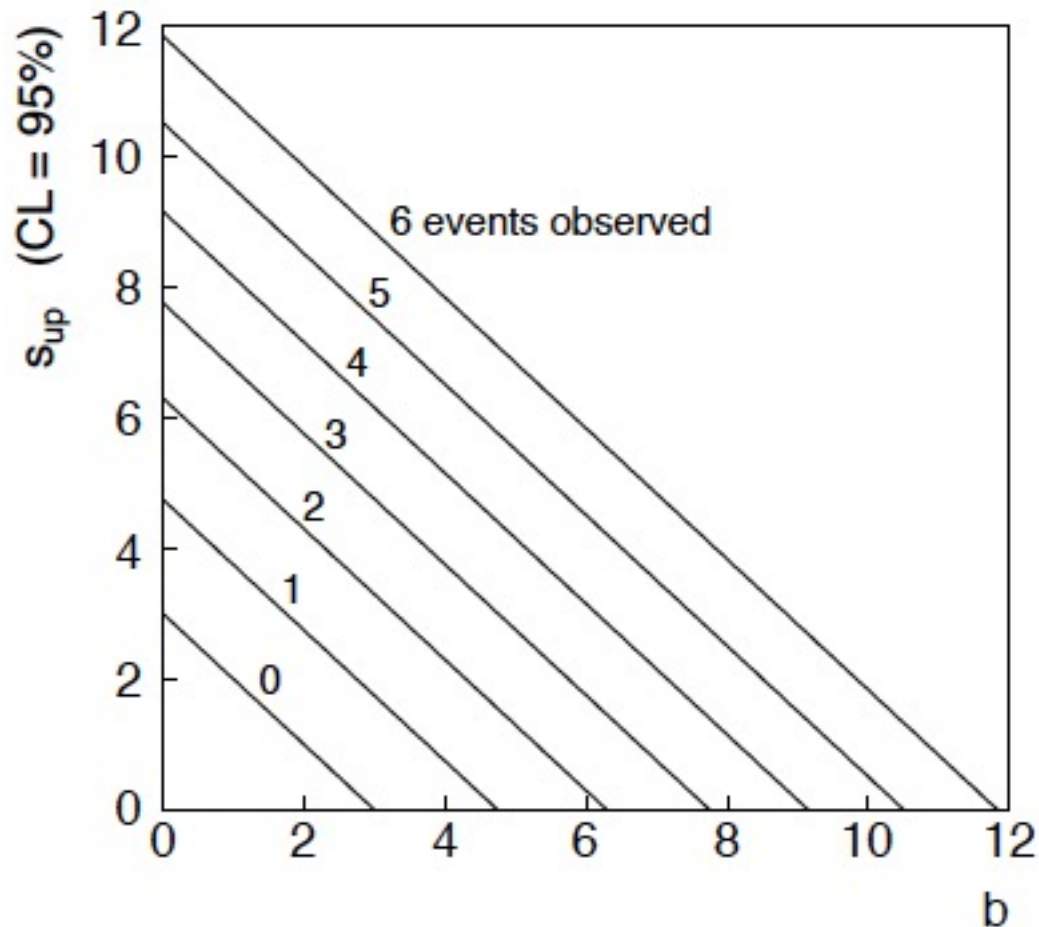
$$\alpha = P(n \leq n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}} + b)}$$

$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1}(1 - \alpha; 2(n_{\text{obs}} + 1)) - b$$

$$= \frac{1}{2} F_{\chi^2}^{-1}(0.95; 2(5 + 1)) - 4.5 = 6.0$$

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

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



# 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  $s_{\text{up}}$

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

Physicist:

We already knew  $s \geq 0$  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 $s = 0$

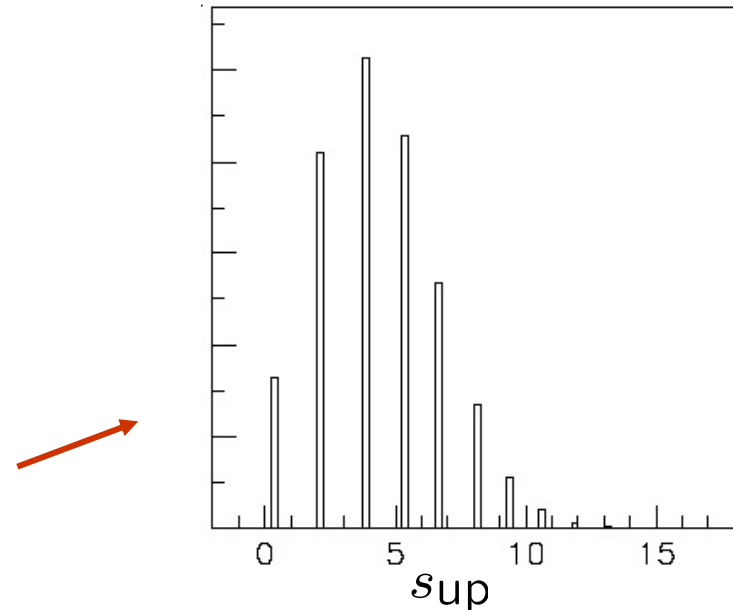
Physicist: I should have used  $CL = 0.95$  — then  $s_{\text{up}} = 0.496$

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

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?

Look at the mean limit for the no-signal hypothesis ( $s = 0$ ) (sensitivity).

Distribution of 95% CL limits  
with  $b = 2.5$ ,  $s = 0$ .  
Mean upper limit = 4.44



# The Bayesian approach to limits

In Bayesian statistics need to start with ‘prior pdf’  $\pi(\theta)$ , this reflects degree of belief about  $\theta$  before doing the experiment.

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

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

Integrate posterior pdf  $p(\theta|x)$  to give interval with any desired probability content.

For e.g.  $n \sim \text{Poisson}(s+b)$ , 95% CL upper limit on  $s$  from

$$0.95 = \int_{-\infty}^{s_{\text{up}}} p(s|n) ds$$

# Bayesian prior for Poisson parameter

Include knowledge that  $s \geq 0$  by setting prior  $\pi(s) = 0$  for  $s < 0$ .

Could try to reflect ‘prior ignorance’ with e.g.

$$\pi(s) = \begin{cases} 1 & s \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Not normalized; can be OK provided  $L(s)$  dies off quickly for large  $s$ .

Not invariant under change of parameter — if we had used instead a flat prior for a nonlinear function of  $s$ , then this would imply a non-flat prior for  $s$ .

Doesn’t really reflect a reasonable degree of belief, but often used as a point of reference; or viewed as a recipe for producing an interval whose frequentist properties can be studied (e.g., coverage probability, which will depend on true  $s$ ).



# Bayesian upper limit with flat prior for $s$

Put Poisson likelihood and flat prior into Bayes' theorem:

$$p(s|n) \propto \frac{(s+b)^n}{n!} e^{-(s+b)} \quad (s \geq 0)$$

Normalize to unit area:

$$p(s|n) = \frac{(s+b)^n e^{-(s+b)}}{\Gamma(b, n+1)}$$

upper incomplete gamma function

Upper limit  $s_{\text{up}}$  determined by requiring

$$1 - \alpha = \int_0^{s_{\text{up}}} p(s|n) ds$$

# Bayesian interval with flat prior for $s$

Solve to find limit  $s_{\text{up}}$ :

$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} [p, 2(n+1)] - b$$

where

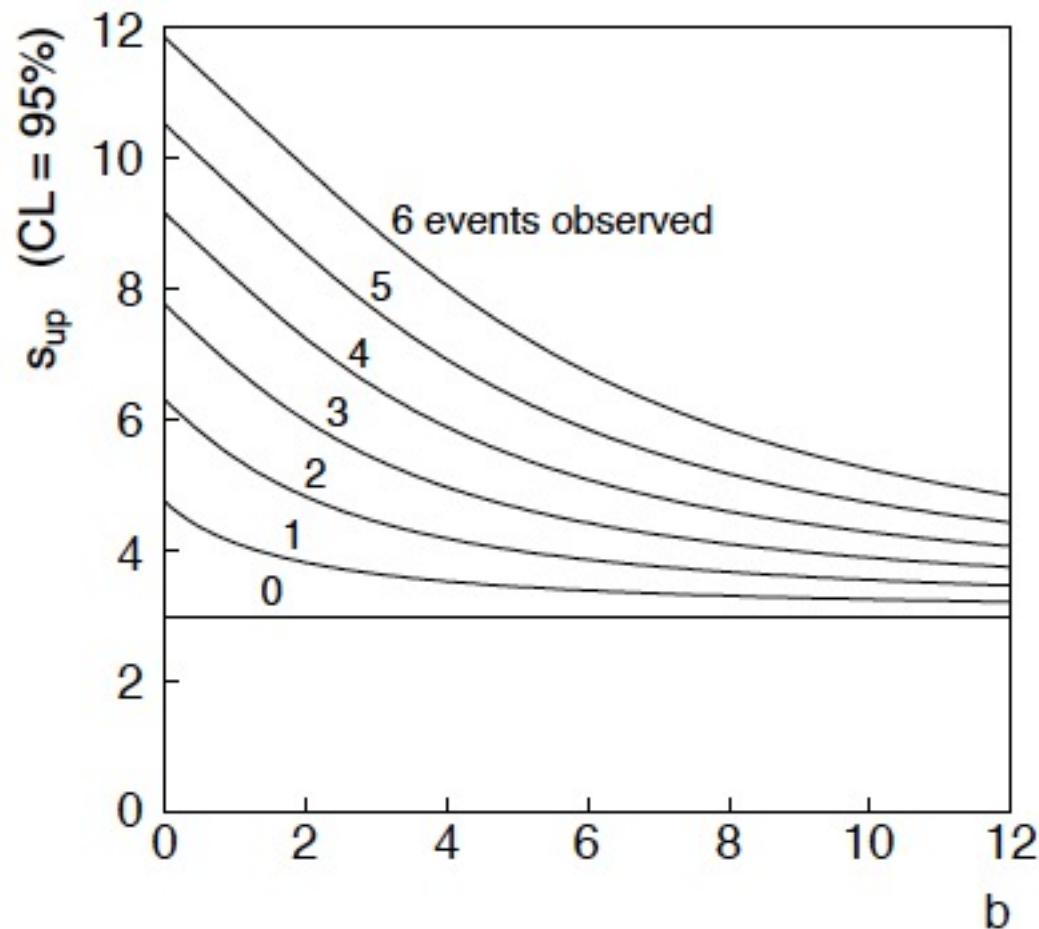
$$p = 1 - \alpha \left( 1 - F_{\chi^2} [2b, 2(n+1)] \right)$$

For special case  $b = 0$ , Bayesian upper limit with flat prior numerically same as one-sided frequentist case ('coincidence').

# Bayesian interval with flat prior for $s$

For  $b > 0$  Bayesian limit is everywhere greater than the (one sided) frequentist upper limit.

Never goes negative. Doesn't depend on  $b$  if  $n = 0$ .



# Statistical Data Analysis

## Lecture 10-2

- Discussion on Bayesian prior probabilities
- Jeffreys' prior
- Example: Poisson mean

# Priors from formal rules

Last time we took the prior for a Poisson mean to be constant to reflect a lack of prior knowledge; we noted this was not invariant under change of parameter.

Because of difficulties in encoding a vague degree of belief in a prior, one often attempts to derive the prior from formal rules, e.g., to satisfy certain invariance principles or to provide maximum information gain for a certain set of measurements.

Often called “objective priors”

Form basis of Objective Bayesian Statistics

The priors do not reflect a degree of belief (but might represent possible extreme cases).

In Objective Bayesian analysis, can use the intervals in a frequentist way, i.e., regard Bayes’ theorem as a recipe to produce an interval with a given coverage probability.

# Priors from formal rules (cont.)

For a review of priors obtained by formal rules see, e.g.,

Robert E. Kass and Larry Wasserman, *The Selection of Prior Distributions by Formal Rules*, J. Am. Stat. Assoc., Vol. 91, No. 435, pp. 1343-1370 (1996).

Formal priors have not been widely used in Particle Physics, but there has been interest in this direction, especially the reference priors of Bernardo and Berger; see e.g.

L. Demortier, S. Jain and H. Prosper, *Reference priors for high energy physics*, Phys. Rev. D 82 (2010) 034002, arXiv:1002.1111.

D. Casadei, *Reference analysis of the signal + background model in counting experiments*, JINST 7 (2012) 01012; arXiv:1108.4270.

# Jeffreys prior

According to *Jeffreys' rule*, take prior according to

$$\pi(\boldsymbol{\theta}) \propto \sqrt{\det(I(\boldsymbol{\theta}))}$$

where

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

is the Fisher information matrix.

One can show that this leads to inference that is invariant under a transformation of parameters in the following sense:

Start with the Jeffreys prior for  $\theta$ :  $\pi_{\theta}(\theta) \sim \sqrt{\det I(\theta)}$

Use it in Bayes' theorem to find:

$$P(\theta|\mathbf{x}) \propto P(\mathbf{x}|\theta)\pi_{\theta}(\theta)$$

## Jeffreys prior (2)

Now consider a function  $\eta(\theta)$ . The posterior for  $\eta$  is

$$P(\eta|\mathbf{x}) = P(\theta|\mathbf{x}) \left| \frac{d\theta}{d\eta} \right|$$

Alternatively, start with  $\eta$  and use its Jeffreys' prior:

$$\pi_{\eta}(\eta) \propto \sqrt{\det I(\eta)}$$

Use this in Bayes' theorem:  $P(\eta|\mathbf{x}) \propto P(\mathbf{x}|\eta)\pi_{\eta}(\eta)$

One can show that Jeffreys' prior results in the same  $P(\eta|\mathbf{x})$  in both cases. For details (single-parameter case) see:

<http://www.pp.rhul.ac.uk/~cowan/stat/notes/JeffreysInvariance.pdf>



# Jeffreys prior for Poisson mean

Suppose  $n \sim \text{Poisson}(\mu)$ . To find the Jeffreys' prior for  $\mu$ ,

$$L(n|\mu) = \frac{\mu^n}{n!} e^{-\mu} \qquad \frac{\partial^2 \ln L}{\partial \mu^2} = -\frac{n}{\mu^2}$$

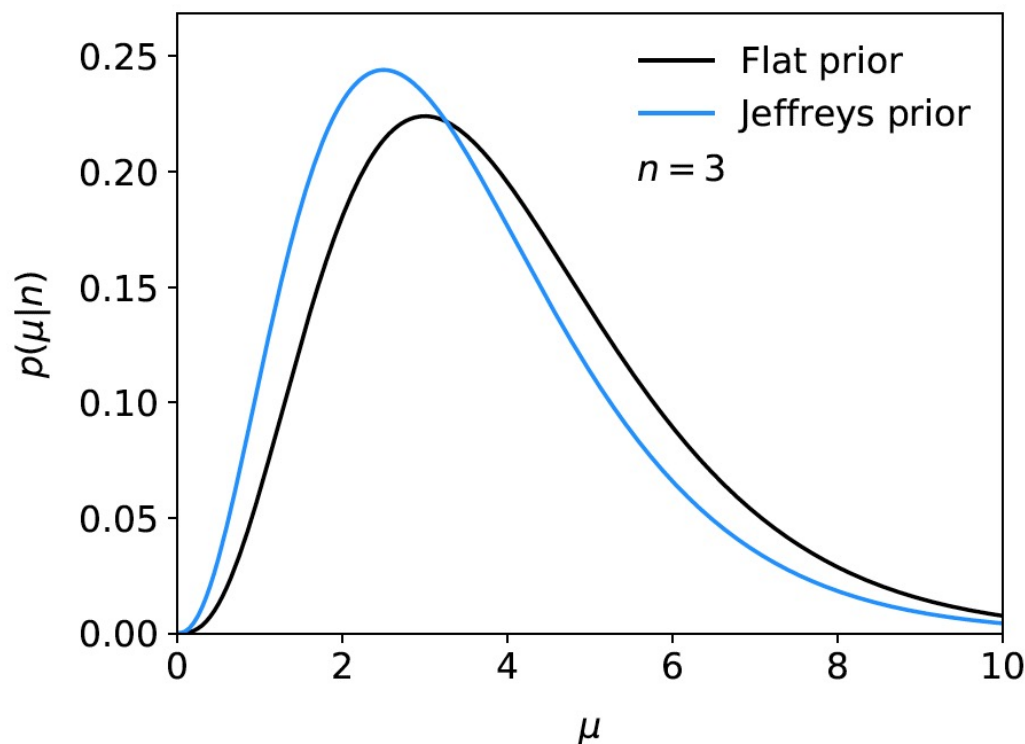
$$I = -E \left[ \frac{\partial^2 \ln L}{\partial \mu^2} \right] = \frac{E[n]}{\mu^2} = \frac{1}{\mu}$$

$$\pi(\mu) \propto \sqrt{I(\mu)} = \frac{1}{\sqrt{\mu}}$$

So e.g. for  $\mu = s + b$ , this means the prior  $\pi(s) \sim 1/\sqrt{s + b}$ , which depends on  $b$ . But this is not designed as a degree of belief about  $s$ .

# Posterior pdf for Poisson mean

From Bayes' theorem,  $p(\mu|n) \propto \mu^n e^{-\mu} \pi(\mu)$



Flat,  $\pi(\mu) = \text{const.}$

$$p(\mu|n) = \frac{\mu^n e^{-\mu}}{\Gamma(n+1)}$$

Jeffreys,  $\pi(\mu) \sim 1/\sqrt{\mu}$

$$p(\mu|n) = \frac{\mu^{n-\frac{1}{2}} e^{-\mu}}{\Gamma(n+\frac{1}{2})}$$

In both cases, posterior is special case of gamma distribution.

# Upper limit for Poisson mean

To find upper limit at  $CL = 1 - \alpha$ , solve

$$1 - \alpha = \int_0^{\mu_{\text{up}}} p(\mu|n) d\mu$$

Jeffreys prior:  $\mu_{\text{up}} = P^{-1}(n + \frac{1}{2}, 1 - \alpha) = 7.03$

Flat prior:  $\mu_{\text{up}} = P^{-1}(n + 1, 1 - \alpha) = 7.75$

$n=3,$   
 $CL=0.95$

where  $P^{-1}$  is the inverse of the normalized lower incomplete gamma function (see `scipy.special`)

$$P(a, \mu_{\text{up}}) = \frac{1}{\Gamma(a)} \int_0^{\mu_{\text{up}}} \mu^{a-1} e^{-\mu} d\mu$$

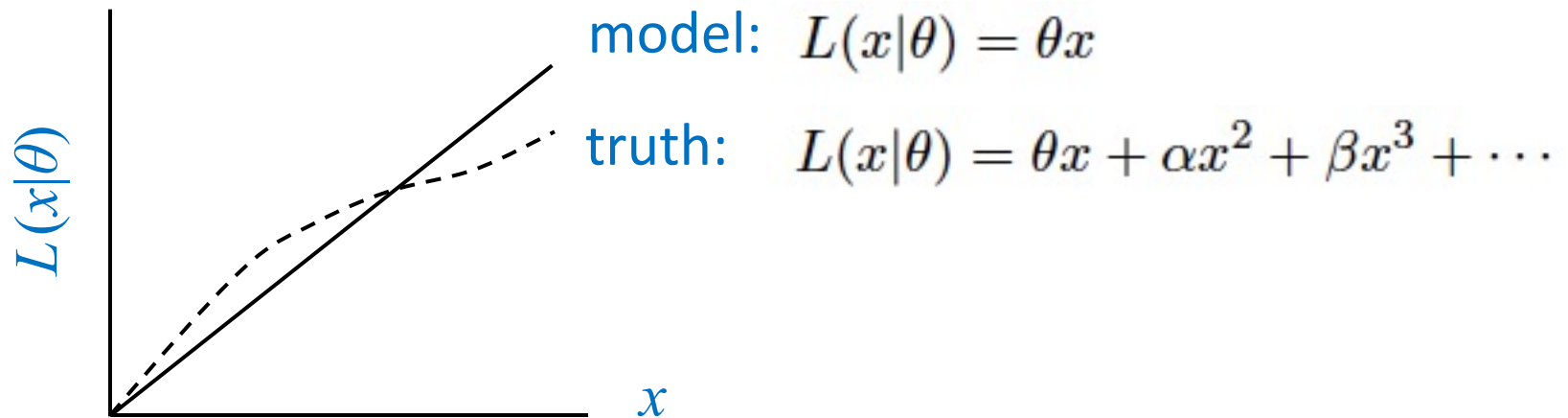
# Statistical Data Analysis

## Lecture 10-3

- Systematic uncertainties and nuisance parameters
- Profile likelihood

# Systematic uncertainties and nuisance parameters

In general, our model of the data is not perfect:



Can improve model by including additional adjustable parameters.

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

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).

# Example: fitting a straight line

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

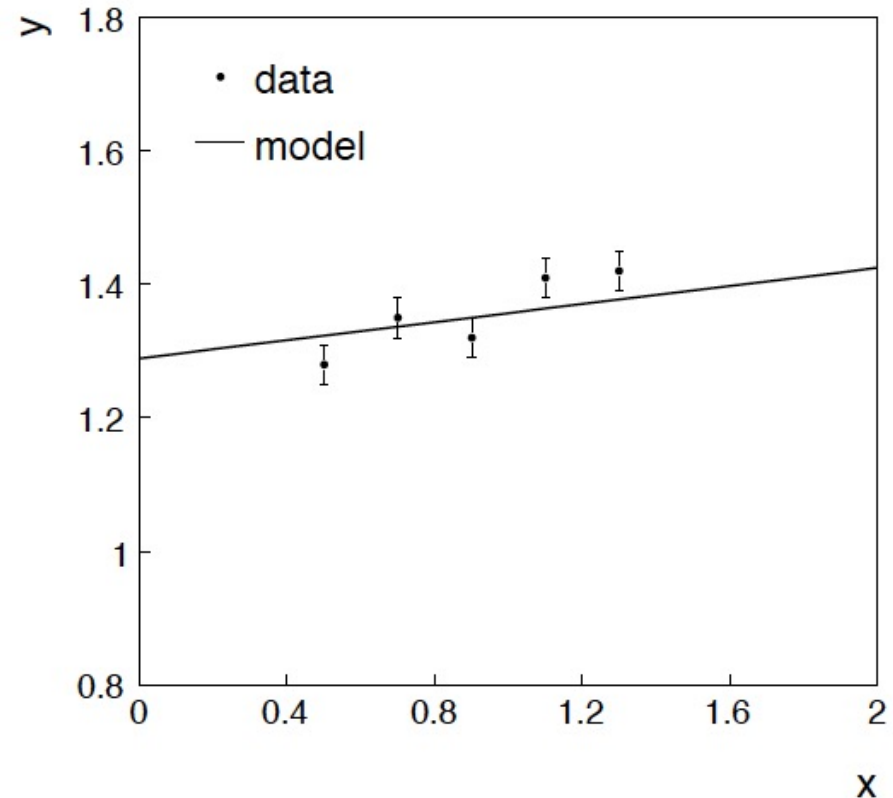
Model:  $y_i$  independent and all follow  $y_i \sim \text{Gauss}(\mu(x_i), \sigma_i)$

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

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

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

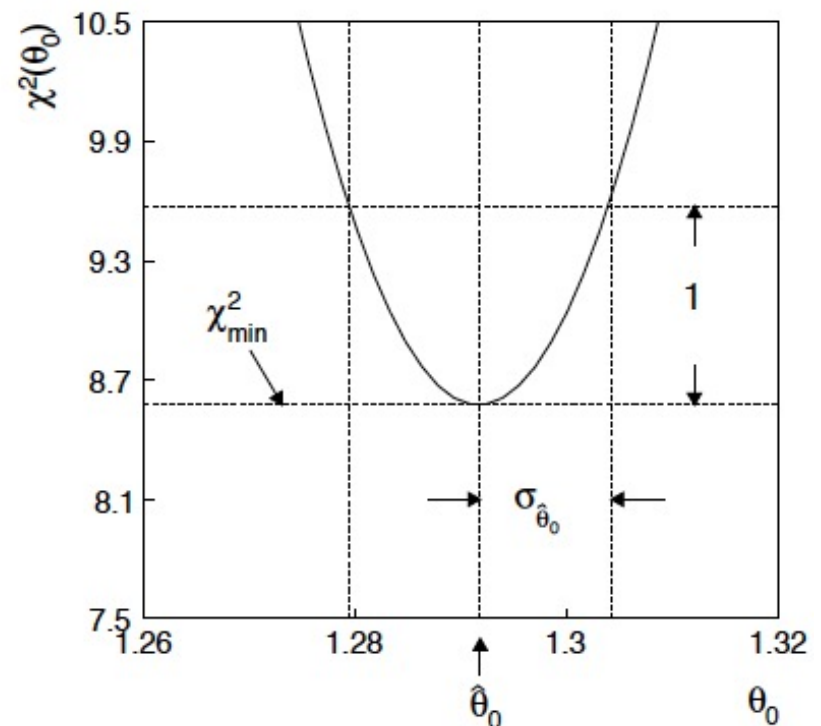
$$\chi^2(\theta_0) = -2 \ln L(\theta_0) + \text{const} = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} .$$

For Gaussian  $y_i$ , ML same as LS

Minimize  $\chi^2 \rightarrow$  estimator  $\hat{\theta}_0$  .

Come up one unit from  $\chi^2_{\min}$

to find  $\sigma_{\hat{\theta}_0}$  .





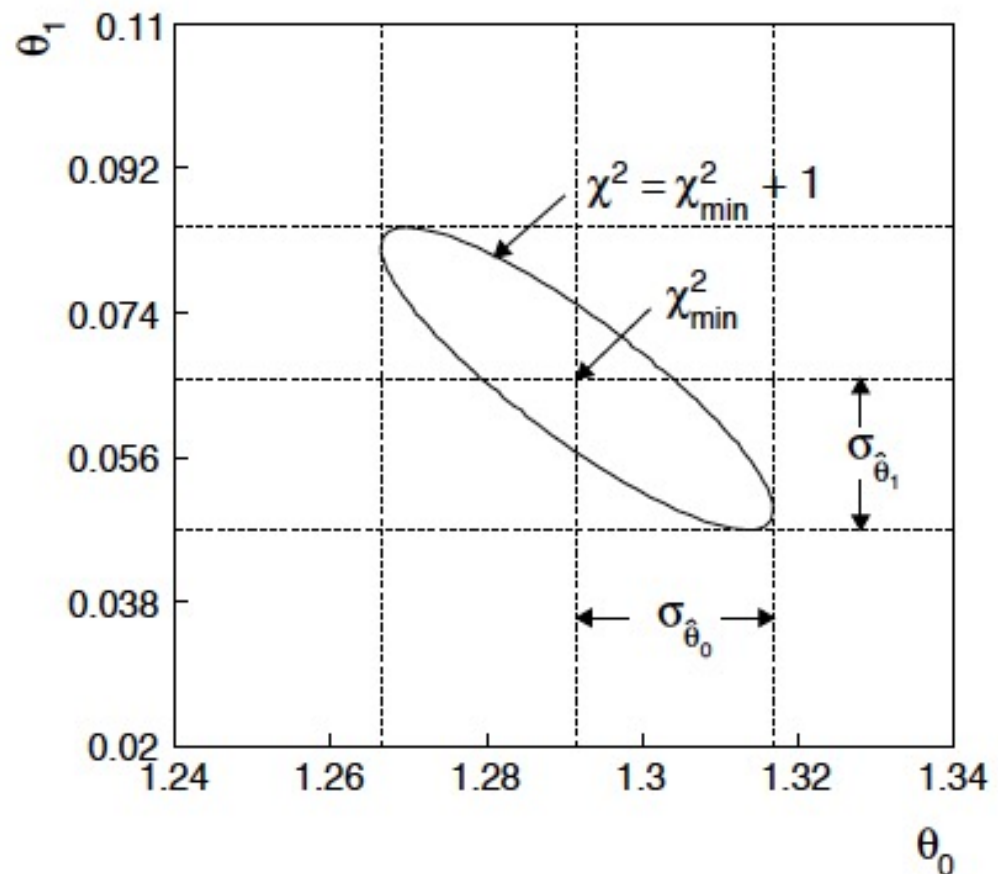
# ML (or LS) fit of $\theta_0$ and $\theta_1$

$$\chi^2(\theta_0, \theta_1) = -2 \ln L(\theta_0, \theta_1) + \text{const} = \sum_{i=1}^n \frac{(y_i - \mu(x_i; \theta_0, \theta_1))^2}{\sigma_i^2} .$$

Standard deviations from  
tangent lines to contour

$$\chi^2 = \chi_{\min}^2 + 1 .$$

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

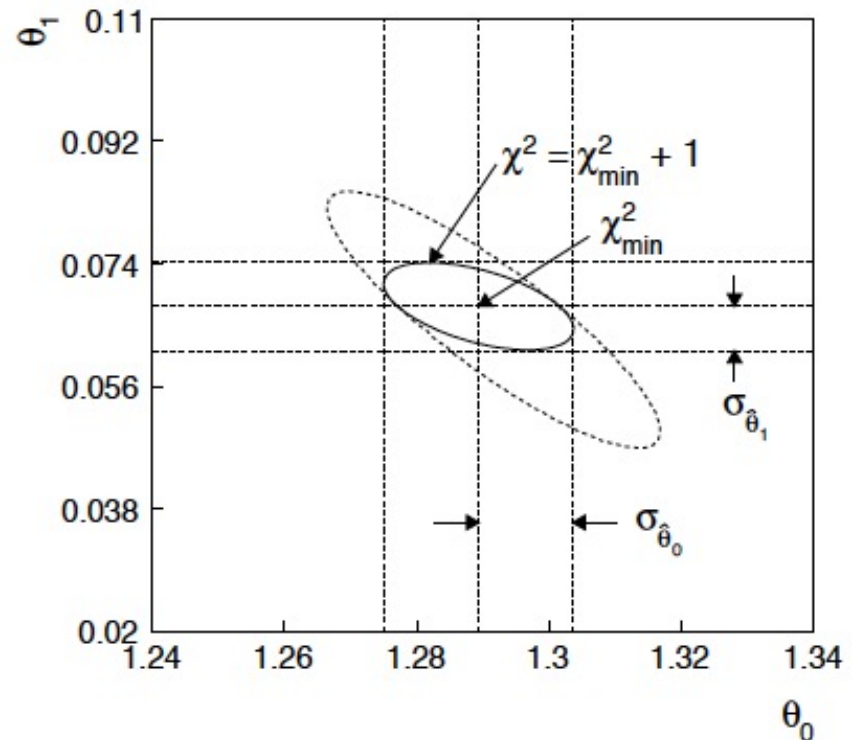


If we have a measurement  $t_1 \sim \text{Gauss}(\theta_1, \sigma_{t_1})$

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

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

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



# Profiling

The  $\ln L = \ln L_{\max} - 1/2$  contour in the  $(\theta_0, \theta_1)$  plane is a confidence region at CL = 39.3%.

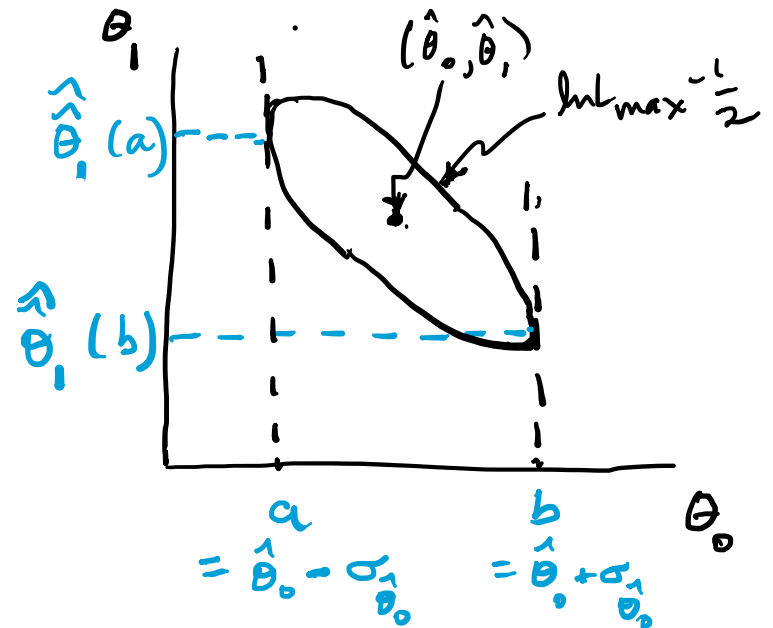
Furthermore if one wants to know only about, say,  $\theta_0$ , then the interval in  $\theta_0$  corresponding to  $\ln L = \ln L_{\max} - 1/2$  is a confidence interval at CL = 68.3% (i.e.,  $\pm 1$  std. dev.).

I.e., form the interval for  $\theta_0$  using

$$\ln L(\theta_0, \hat{\theta}_1(\theta_0)) = \ln L_{\max} - 1/2$$

where  $\theta_1$  is replaced by its “profiled” value

$$\hat{\theta}_1(\theta_0) = \operatorname{argmax}_{\theta_1} L(\theta_0, \theta_1)$$



# Profile Likelihood

Suppose we have a likelihood  $L(\boldsymbol{\mu}, \boldsymbol{\theta}) = P(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\theta})$  with  $N$  parameters of interest  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)$  and  $M$  nuisance parameters  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_M)$ . The “profiled” (or “constrained”) values of  $\boldsymbol{\theta}$  are:

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

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

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

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

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

## 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:

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

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  $t_\mu$ .

# Statistical Data Analysis

## Lecture 10-4

- Bayesian parameter estimation
- Marginalization of posterior pdf
- Markov Chain Monte Carlo

# 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.,

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

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

$$\pi_1(\theta_1) = p(\theta_1|t_1) \propto p(t_1|\theta_1)\pi_{\text{Ur}}(\theta_1) = \frac{1}{\sqrt{2\pi}\sigma_t} e^{-(t_1-\theta_1)^2/2\sigma_t^2} \times \text{const.}$$

prior after  $t_1$ ,  
before  $y$

Ur = “primordial”  
prior

Likelihood for control  
measurement  $t_1$

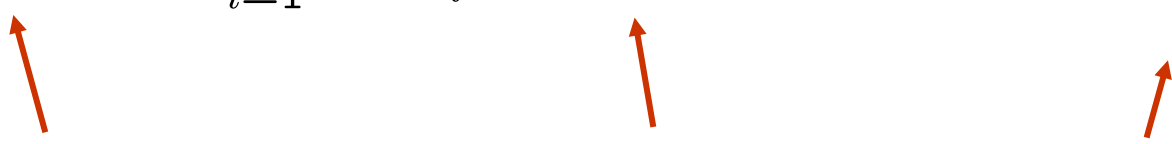
# Bayesian example: $y_i \sim \text{Gauss}(\mu(x_i; \theta_0, \theta_1), \sigma_i)$

Putting the ingredients into Bayes' theorem gives:

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

Diagram illustrating the components of Bayes' theorem:

posterior  $\propto$  likelihood  $\times$  prior



Note here the likelihood only reflects the measurements  $\mathbf{y}$ .

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  $P(\mathbf{y}, t | \theta_0, \theta_1)$  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

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

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

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

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

# 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  $\vec{\theta} \sim q(\vec{\theta}; \vec{\theta}_0)$
- 3) Form Hastings test ratio  $\alpha = \min \left[ 1, \frac{p(\vec{\theta})q(\vec{\theta}_0; \vec{\theta})}{p(\vec{\theta}_0)q(\vec{\theta}; \vec{\theta}_0)} \right]$
- 4) Generate  $u \sim \text{Uniform}[0, 1]$
- 5) If  $u \leq \alpha$ ,  $\vec{\theta}_1 = \vec{\theta}$ ,  move to proposed point  
else  $\vec{\theta}_1 = \vec{\theta}_0$   old point repeated
- 6) Iterate

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



# 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:  $q(\theta; \theta_0) = q(\theta_0; \theta)$

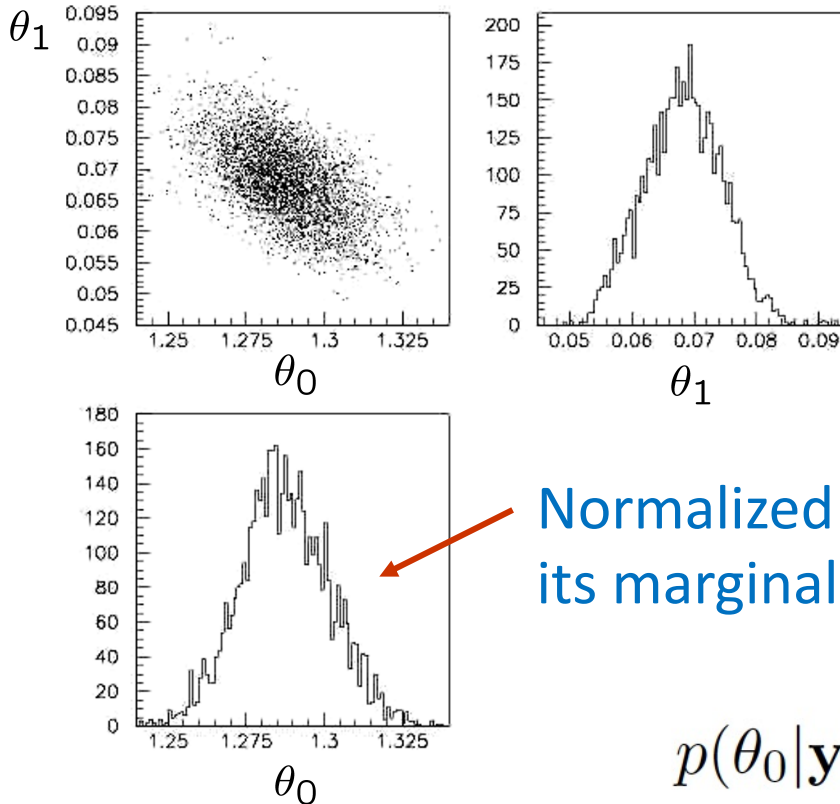
Test ratio is (*Metropolis-Hastings*):  $\alpha = \min \left[ 1, \frac{p(\vec{\theta})}{p(\vec{\theta}_0)} \right]$

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:



Normalized histogram of  $\theta_0$  gives its marginal posterior pdf:

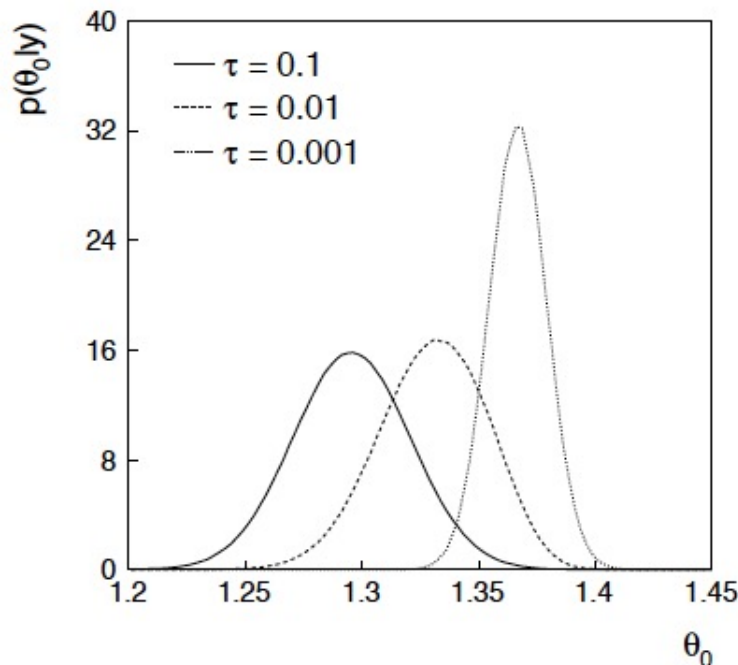
$$p(\theta_0|\mathbf{y}) = \int p(\theta_0, \theta_1|\mathbf{y}) d\theta_1$$

# 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

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

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



This summarizes all knowledge about  $\theta_0$ .

Look also at result from variety of priors.



# Extra slides

# $p$ -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  $\nu$  to use for  $f(q_\theta|\theta, \nu)$ ?

Fundamentally we want to reject  $\theta$  only if  $p_\theta < \alpha$  for all  $\nu$ .

→ “exact” confidence interval

Recall that for statistics based on the profile likelihood ratio, the distribution  $f(q_\theta|\theta, \nu)$  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  $\nu$  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  $p_\theta \leq \alpha$  where the  $p$ -value is computed assuming the value of the nuisance parameter that best fits the data for the specified  $\theta$  :

$$\hat{\hat{\nu}}(\theta)$$

“double hat” notation means value of parameter that maximizes likelihood for the given  $\theta$ .

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

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  $\nu$  is characterized by a Bayesian prior  $\pi(\nu)$ .

Can use the marginal likelihood to model the data:

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

This does not represent what the data distribution would be if we “really” repeated the experiment, since then  $\nu$  would not change.

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

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