# Statistical combinations, etc.

https://indico.desy.de/conferenceDisplay.py?confId=11244



Helmholtz Alliance

Terascale Statistics School DESY, Hamburg March 23-27, 2015



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

# Outline

- 1. Review of some formalism and analysis tasks
- 2. Broad view of combinations & review of parameter estimation
- 3. Combinations of parameter estimates.
- 4. Least-squares averages, including correlations
- 5. Comparison with Bayesian parameter estimation
- 6. Bayesian averages with outliers
- 7. PDG brief overview

Hypothesis, distribution, likelihood, model Suppose the outcome of a measurement is x. (e.g., a number of events, a histogram, or some larger set of numbers). A hypothesis H specifies the probability of the data P(x|H). Often *H* is labeled by parameter(s)  $\theta \rightarrow P(x|\theta)$ . For the probability distribution  $P(x|\theta)$ , variable is x;  $\theta$  is a constant. If e.g. we evaluate  $P(x|\theta)$  with the observed data and regard it as a function of the parameter(s), then this is the likelihood:

 $L(\theta) = P(x|\theta)$  (Data x fixed; treat L as function of  $\theta$ .)

Here use the term 'model' to refer to the full function  $P(x|\theta)$  that contains the dependence both on *x* and  $\theta$ .

(Sometimes write  $L(x|\theta)$  for model or likelihood, depending on context.)

G. Cowan

Bayesian use of the term 'likelihood' We can write Bayes theorem as posterior  $p(\theta|x) = \frac{L(x|\theta)\pi(\theta)}{\int L(x|\theta)\pi(\theta) d\theta}$ 

where  $L(x|\theta)$  is the likelihood. It is the probability for x given  $\theta$ , evaluated with the observed x, and viewed as a function of  $\theta$ .

Bayes' theorem only needs  $L(x|\theta)$  evaluated with a given data set (the 'likelihood principle').

For frequentist methods, in general one needs the full model. For some approximate frequentist methods, the likelihood is enough.

# Theory ↔ Statistics ↔ Experiment



### 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).

# Data analysis in particle physics

Observe events (e.g., pp collisions) and for each, measure a set of characteristics:

particle momenta, number of muons, energy of jets,... Compare observed distributions of these characteristics to predictions of theory. From this, we want to:

Estimate the free parameters of the theory:  $m_{\mu} = 125.4$ 

Quantify the uncertainty in the estimates:  $\pm 0.4$  GeV

Assess how well a given theory stands in agreement with the observed data:  $O^+$  good,  $2^+$  bad

Test/exclude regions of the model's parameter space ( $\rightarrow$  limits)

# Combinations

"Combination of results" can be taken to mean: how to construct a model that incorporates more data.

#### E.g. several experiments,

Experiment 1: data *x*, model  $P(x|\theta) \rightarrow$  upper limit  $\theta_{up,1}$ Experiment 2: data *y*, model  $P(y|\theta) \rightarrow$  upper limit  $\theta_{up,2}$ 

Or main experiment and control measurement(s).

The best way to do the combination is at the level of the data, e.g., (if x, y independent)

 $P(x,y|\theta) = P(x|\theta) P(y|\theta) \rightarrow$  "combined" limit  $\theta_{up,comb}$ 

If the data are not available but rather only the "results" (limits, parameter estimates, *p*-values) then possibilities are more limited.

Usually OK for parameter estimates, difficult/impossible for limits, *p*-values without additional assumptions & information loss.

Quick review of frequentist parameter estimation

Suppose we have a pdf characterized by one or more parameters:

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

random variable

parameter

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

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.

# Maximum likelihood

The most important frequentist method for constructing estimators is to take the value of the parameter(s) that maximize the likelihood:  $\hat{\theta} = \operatorname{argmax} L(x|\theta)$ 

The resulting estimators are functions of the data and thus characterized by a sampling distribution with a given (co)variance:

In general they may have a nonzero bias:

Under conditions usually satisfied in practice, bias of ML estimators is zero in the large sample limit, and the variance is as small as possible for unbiased estimators.

ML estimator may not in some cases be regarded as the optimal trade-off between these criteria (cf. regularized unfolding).

G. Cowan

Terascale Statistics School 2015 / Combination

 $V_{ij} = \operatorname{cov}[\hat{\theta}_i, \hat{\theta}_j]$ 

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

# Ingredients for ML

To find the ML estimate itself one only needs the likelihood  $L(\theta)$ . In principle to find the covariance of the estimators, one requires the full model  $L(x|\theta)$ . E.g., simulate many times independent data sets and look at distribution of the resulting estimates:



G. Cowan

# Ingredients for ML (2)

Often (e.g., large sample case) one can approximate the covariances using only the likelihood  $L(\theta)$ :

$$\widehat{V}_{ij}^{-1} \approx -\frac{\partial^2 \ln L}{\partial \theta_i \, \partial \theta_j} \bigg|_{\theta = \widehat{\theta}}$$



This translates into a simple graphical recipe:

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

 $\rightarrow$  Tangent lines to contours give standard deviations.

 $\rightarrow$  Angle of ellipse  $\phi$  related to correlation: tan  $2\phi$ 

$$=\frac{2\rho\sigma_{\widehat{\alpha}}\sigma_{\widehat{\beta}}}{\sigma_{\widehat{\alpha}}^2-\sigma\widehat{\beta}^2}$$

# The method of least squares

Suppose we measure N values,  $y_1, ..., y_N$ , assumed to be independent Gaussian r.v.s with

$$E[y_i] = \lambda(x_i; \theta)$$
.

Assume known values of the control variable  $x_1, ..., x_N$  and known variances

$$V[y_i] = \sigma_i^2 \, .$$



We want to estimate  $\theta$ , i.e., fit the curve to the data points.

The likelihood function is

$$L(\theta) = \prod_{i=1}^{N} f(y_i; \theta) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{(y_i - \lambda(x_i; \theta))^2}{2\sigma_i^2}\right]$$

G. Cowan

# The method of least squares (2)

The log-likelihood function is therefore

$$\ln L(\theta) = -\frac{1}{2} \sum_{i=1}^{N} \frac{(y_i - \lambda(x_i; \theta))^2}{\sigma_i^2} + \text{ terms not depending on } \theta$$

So maximizing the likelihood is equivalent to minimizing

$$\chi^{2}(\theta) = \sum_{i=1}^{N} \frac{(y_{i} - \lambda(x_{i}; \theta))^{2}}{\sigma_{i}^{2}}$$

Minimum defines the least squares (LS) estimator  $\hat{\theta}$ .

Very often measurement errors are ~Gaussian and so ML and LS are essentially the same.

Often minimize  $\chi^2$  numerically (e.g. program **MINUIT**).

# LS with correlated measurements

If the  $y_i$  follow a multivariate Gaussian, covariance matrix V,

$$g(\vec{y}, \vec{\lambda}, V) = \frac{1}{(2\pi)^{N/2} |V|^{1/2}} \exp\left[-\frac{1}{2}(\vec{y} - \vec{\lambda})^T V^{-1}(\vec{y} - \vec{\lambda})\right]$$

Then maximizing the likelihood is equivalent to minimizing

$$\chi^2(\vec{\theta}) = \sum_{i,j=1}^N (y_i - \lambda(x_i; \vec{\theta}))(V^{-1})_{ij}(y_j - \lambda(x_j; \vec{\theta}))$$

# Linear LS problem

LS has particularly simple properties if  $\lambda(x; \vec{\theta})$  linear in  $\vec{\theta}$ :

$$\lambda(x;ec{ heta}) = \sum\limits_{j=1}^m a_j(x) heta_j$$

where  $a_j(x)$  are any linearly independent functions of x.

Matrix notation: let  $A_{ij} = a_j(x_i)$ ,

$$egin{aligned} \chi^2(ec{ heta}) &= (ec{y} - ec{\lambda})^T \, V^{-1} \, (ec{y} - ec{\lambda}) \ &= (ec{y} - A ec{ heta})^T \, V^{-1} \, (ec{y} - A ec{ heta}) \end{aligned}$$

# Linear LS problem (2)

Set derivitives with respect to  $\theta_i$  to zero,

$$\nabla \chi^2 = -2(A^T V^{-1} \vec{y} - A^T V^{-1} A \vec{\theta}) = 0$$

Solve to get the LS estimators,

$$\hat{\vec{\theta}} = (A^T V^{-1} A)^{-1} A^T V^{-1} \vec{y} \equiv B \vec{y}$$

N.B. estimators  $\hat{\theta}_i$  are linear functions of the measurements  $y_i$ .

# Linear LS problem (3)

Error propagation (exact for linear problem) for  $U_{ij} = \operatorname{cov}[\hat{\theta}_i, \hat{\theta}_j]$ :

 $U = B V B^{T} = (A^{T} V^{-1} A)^{-1}$ 

Equivalently, use

$$(U^{-1})_{ij} = \frac{1}{2} \left[ \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \right]_{\vec{\theta} = \vec{\hat{\theta}}}$$

Equals MVB if  $y_i$  Gaussian)

# Goodness-of-fit with least squares

The value of the  $\chi^2$  at its minimum is a measure of the level of agreement between the data and fitted curve:

$$\chi^2_{\min} = \sum_{i=1}^{N} \frac{(y_i - \lambda(x_i; \hat{\theta}))^2}{\sigma_i^2}$$

It can therefore be employed as a goodness-of-fit statistic to test the hypothesized functional form  $\lambda(x; \theta)$ .

We can show that if the hypothesis is correct, then the statistic  $t = \chi^2_{\text{min}}$  follows the chi-square pdf,

$$f(t; n_{\rm d}) = \frac{1}{2^{n_{\rm d}/2} \Gamma(n_{\rm d}/2)} t^{n_{\rm d}/2 - 1} e^{-t/2}$$

where the number of degrees of freedom is

 $n_{\rm d}$  = number of data points – number of fitted parameters

# Goodness-of-fit with least squares (2)

The chi-square pdf has an expectation value equal to the number of degrees of freedom, so if  $\chi^2_{\rm min} \approx n_{\rm d}$  the fit is 'good'.

More generally, find the *p*-value:  $p = \int_{\chi^2_{\min}}^{\infty} f(t; n_d) dt$ 

This is the probability of obtaining a  $\chi^2_{min}$  as high as the one we got, or higher, if the hypothesis is correct.

# Using LS to combine measurements

Use LS to obtain weighted average of N measurements of  $\lambda$ :

 $y_i$  = result of measurement i, i = 1, ..., N; $\sigma_i^2 = V[y_i]$ , assume known;  $\lambda$  = true value (plays role of  $\theta$ ).

For uncorrelated  $y_i$ , minimize

$$\chi^2(\lambda) = \sum_{i=1}^N rac{(y_i - \lambda)^2}{\sigma_i^2},$$

Set 
$$\frac{\partial \chi^2}{\partial \lambda} = 0$$
 and solve,  
 $\rightarrow \quad \hat{\lambda} = \frac{\sum_{i=1}^N y_i / \sigma_i^2}{\sum_{j=1}^N 1 / \sigma_j^2} \qquad \qquad V[\hat{\lambda}] = \frac{1}{\sum_{i=1}^N 1 / \sigma_i^2}$ 

# Combining correlated measurements with LS

If  $\operatorname{cov}[y_i, y_j] = V_{ij}$ , minimize  $\chi^2(\lambda) = \sum_{i,j=1}^N (y_i - \lambda)(V^{-1})_{ij}(y_j - \lambda),$   $\rightarrow \quad \hat{\lambda} = \sum_{i=1}^N w_i y_i, \quad w_i = \frac{\sum_{j=1}^N (V^{-1})_{ij}}{\sum_{k,l=1}^N (V^{-1})_{kl}}$  $V[\hat{\lambda}] = \sum_{i=1}^N w_i V_{ij} w_j$ 

LS  $\hat{\lambda}$  has zero bias, minimum variance (Gauss–Markov theorem).

That is, if we take the estimator to be a linear form  $\Sigma_i w_i y_i$ , and find the  $w_i$  that minimize its variance, we get the LS solution (= BLUE, Best Linear Unbiased Estimator).

# Example: averaging two correlated measurements

Suppose we have 
$$y_1, y_2$$
, and  $V = \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix}$ 

$$\rightarrow \quad \hat{\lambda} = wy_1 + (1 - w)y_2, \quad w = \frac{\sigma_2^2 - \rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$$
$$V[\hat{\lambda}] = \frac{(1 - \rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} = \sigma^2$$

The increase in inverse variance due to 2nd measurement is

$$\frac{1}{\sigma^2} - \frac{1}{\sigma_1^2} = \frac{1}{1 - \rho^2} \left( \frac{\rho}{\sigma_1} - \frac{1}{\sigma_2} \right)^2 > 0$$

 $\rightarrow$  2nd measurement can only help.

G. Cowan

Negative weights in LS average If  $\rho > \sigma_1/\sigma_2$ ,  $\rightarrow w < 0$ ,

 $\rightarrow$  weighted average is not between  $y_1$  and  $y_2$  (!?) Cannot happen if correlation due to common data, but possible for shared random effect; very unreliable if e.g.  $\rho$ ,  $\sigma_1$ ,  $\sigma_2$  incorrect.

See example in SDA Section 7.6.1 with two measurements at same temperature using two rulers, different thermal expansion coefficients: average is outside the two measurements; used to improve estimate of temperature.

### Covariance, correlation, etc.

For a pair of random variables *x* and *y*, the covariance and correlation are

$$\operatorname{cov}[x,y] = E[xy] - E[x]E[y]$$
  $\rho_{xy} = \frac{\operatorname{cov}[x,y]}{\sigma_x \sigma_y}$ 

One only talks about the correlation of two quantities to which one assigns probability (i.e., random variables).

So in frequentist statistics, estimators for parameters can be correlated, but not the parameters themselves.

In Bayesian statistics it does make sense to say that two parameters are correlated, e.g.,

$$\operatorname{cov}[\theta_i, \theta_j] = \int \theta_i \theta_j p(\theta | x) \, d\theta \, - \, \int \theta_i p(\theta | x) \, d\theta \, \int \theta_j p(\theta | x) \, d\theta$$

# Example of "correlated systematics"

Suppose we carry out two independent measurements of the length of an object using two rulers with different thermal expansion properties.

Suppose the temperature is not known exactly but must be measured (but lengths measured together so *T* same for both),

$$T \sim \text{Gauss}(\tau, \sigma_T)$$

The expectation value of the measured length  $L_i$  (i = 1, 2) is related to true length  $\lambda$  at a reference temperature  $\tau_0$  by

$$E[L_i] = \lambda - \alpha_i (T - \tau_0), \qquad i = 1, 2$$

and the (uncorrected) length measurements are modeled as

$$L_i \sim \text{Gauss}(\lambda - \alpha_i(\tau - \tau_0), \sigma_i)$$

G. Cowan

# Two rulers (2)

The model thus treats the measurements  $T, L_1, L_2$  as uncorrelated with standard deviations  $\sigma_T, \sigma_1, \sigma_2$ , respectively:

$$L(T, L_1, L_2 | \lambda, \tau) = \frac{1}{\sqrt{2\pi}\sigma_T} e^{-(T-\tau)^2/2\sigma_T^2} \prod_{i=1}^2 \frac{1}{\sqrt{2\pi}\sigma_i} e^{-(L_i - \lambda + \alpha_i(\tau - T_0))^2/2\sigma_i^2}$$

Alternatively we could correct each raw measurement:

$$y_i = L_i + \alpha_i (T - \tau_0)$$

which introduces a correlation between  $y_1$ ,  $y_2$  and T

$$\operatorname{cov}[y_1, y_2] = \alpha_1 \alpha_2 \sigma_T^2 \qquad \qquad \operatorname{cov}[y_i, T] = \alpha_i \sigma_T^2$$

But the likelihood function (multivariate Gauss in  $T, y_1, y_2$ ) is the same function of  $\tau$  and  $\lambda$  as before (equivalent!).

Language of  $y_1, y_2$ : temperature gives correlated systematic. Language of  $L_1, L_2$ : temperature gives "coherent" systematic.

# Two rulers (3)

#### Outcome has some surprises:



Estimate of  $\lambda$  does not lie between  $y_1$  and  $y_2$ .

Stat. error on new estimate of temperature substantially smaller than initial  $\sigma_T$ .

These are features, not bugs, that result from our model assumptions.

# Two rulers (4)

We may re-examine the assumptions of our model and conclude that, say, the parameters  $\alpha_1$ ,  $\alpha_2$  and  $\tau_0$  were also uncertain.

We may treat their nominal values as measurements (need a model; Gaussian?) and regard  $\alpha_1$ ,  $\alpha_2$  and  $\tau_0$  as as nuisance parameters.

 $L(L_1, L_2, T, \tilde{\tau}_0, \tilde{\alpha}_1, \tilde{\alpha}_2 | \lambda, \tau, \tau_0, \alpha_1, \alpha_2) =$ 

$$\frac{1}{\sqrt{2\pi}\sigma_T}e^{-(T-\tau)^2/2\sigma_T^2}\prod_{i=1}^2\frac{1}{\sqrt{2\pi}\sigma_i}e^{-(L_i-\lambda+\alpha_i(\tau-\tau_0))^2/2\sigma_i^2}$$

$$\frac{1}{\sqrt{2\pi}\sigma_i}e^{-(\tilde{\alpha}_i-\alpha_i)^2/2\sigma_i^2}\prod_{i=1}^2\frac{1}{\sqrt{2\pi}\sigma_i}e^{-(\tilde{\alpha}_i-\alpha_i)^2/2\sigma_i^2}$$

$$\times \frac{1}{\sqrt{2\pi\sigma_{\tilde{\tau}_0}}} e^{-(\tilde{\tau}_0 - \tau_0)^2/2\sigma_{\tilde{\tau}_0}^2} \prod_{i=1} \frac{1}{\sqrt{2\pi\sigma_{\tilde{\alpha}_i}}} e^{-(\tilde{\alpha}_i - \alpha_i)^2/2\sigma_{\tilde{\alpha}_i}^2}$$

# Two rulers (5)

The outcome changes; some surprises may be "reduced".



# "Related" parameters

Suppose the model for two independent measurements x and y contain the same parameter of interest  $\mu$  and a common nuisance parameter,  $\theta$ , such as the jet-energy scale.

To combine the measurements, construct the full likelihood:

 $L(\mu, \theta) = P(x|\mu, \theta)P(y|\mu, \theta)$ 

Although one may think of  $\theta$  as common to the two measurements, this could be a poor approximation (e.g., the two analyses use jets with different angles/energies, so a single jet-energy scale is not a good model).

Better model: suppose the parameter for x is  $\theta$ , and for y it is

$$\theta' = \theta + \varepsilon$$

where  $\varepsilon$  is an additional nuisance parameter expected to be small.

# Model with nuisance parameters

The additional nuisance parameters in the model may spoil our sensitivity to the parameter of interest  $\mu$ , so we need to constrain them with control measurements.

Often we have no actual control measurements, but some "nominal values" for  $\theta$  and  $\varepsilon$ ,  $\tilde{\theta}$  and  $\tilde{\varepsilon}$ , which we treat as if they were measurements, e.g., with a Gaussian model:

$$p(\tilde{\theta}, \tilde{\varepsilon} | \theta, \varepsilon) = \text{Gauss}(\tilde{\theta} | \theta, \sigma_{\tilde{\theta}}) \text{Gauss}(\tilde{\varepsilon} | \varepsilon, \sigma_{\tilde{\varepsilon}})$$

We started by considering  $\theta$  and  $\theta'$  to be the same, so probably take  $\tilde{\varepsilon} = 0$ .

So we now have an improved model

$$L(\mu, \theta, \varepsilon) = P(x|\mu, \theta) P(y|\mu, \theta, \varepsilon) p(\tilde{\theta}, \tilde{\varepsilon}|\theta, \varepsilon)$$

with which we can estimate  $\mu$ , set limits, etc.

G. Cowan

# Example: fitting a straight line

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

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

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

- 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")



# 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

Terascale Statistics School 2015 / Combination

### 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.

# Marginalization with MCMC

Bayesian computations involve integrals like

$$p(\theta|x) = \int p(\theta, \nu|x) \, d\nu$$

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 naive  $\sqrt{n}$ .

Basic idea: sample full multidimensional parameter space; look, e.g., only at distribution of parameters of interest.

# 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  $\theta_0$ :



The error on the error Some systematic errors are well determined Error from finite Monte Carlo sample

Some are less obvious

Do analysis in *n* 'equally valid' ways and extract systematic error from 'spread' in results.

Some are educated guesses

Guess possible size of missing terms in perturbation series; vary renormalization scale  $(\mu/2 < Q < 2\mu ?)$ 

Can we incorporate the 'error on the error'?

(cf. G. D'Agostini 1999; Dose & von der Linden 1999)

A more general fit (symbolic) Given measurements:  $y_i \pm \sigma_i^{\text{stat}} \pm \sigma_i^{\text{sys}}$ , i = 1, ..., n, and (usually) covariances:  $V_{ij}^{\text{stat}}$ ,  $V_{ij}^{\text{sys}}$ . Predicted value:  $\mu(x_i; \theta)$ , expectation value  $E[y_i] = \mu(x_i; \theta) + b_i$ control variable parameters bias

Often take:  $V_{ij} = V_{ij}^{\text{stat}} + V_{ij}^{\text{sys}}$ Minimize  $\chi^2(\theta) = (\vec{y} - \vec{\mu}(\theta))^T V^{-1} (\vec{y} - \vec{\mu}(\theta))$ 

2 -

Equivalent to maximizing  $L(\theta) \gg e^{-\chi^2/2}$ , i.e., least squares same as maximum likelihood using a Gaussian likelihood function.

G. Cowan

Its Bayesian equivalent Take  $L(\vec{y}|\vec{\theta},\vec{b}) \sim \exp\left[-\frac{1}{2}(\vec{y}-\vec{\mu}(\theta)-\vec{b})^T V_{\text{stat}}^{-1}(\vec{y}-\vec{\mu}(\theta)-\vec{b})\right]$   $\pi_b(\vec{b}) \sim \exp\left[-\frac{1}{2}\vec{b}^T V_{\text{sys}}^{-1}\vec{b}\right]$   $\pi_\theta(\theta) \sim \text{const.}$  Joint probability for all parameters and use Bayes' theorem:  $p(\theta,\vec{b}|\vec{y}) \propto L(\vec{y}|\theta,\vec{b})\pi_\theta(\theta)\pi_b(\vec{b})$ 

To get desired probability for  $\theta$ , integrate (marginalize) over **b**:

$$p(\theta|\vec{y}) = \int p(\theta, \vec{b}|\vec{y}) d\vec{b}$$

→ Posterior is Gaussian with mode same as least squares estimator,  $\sigma_{\theta}$  same as from  $\chi^2 = \chi^2_{\min} + 1$ . (Back where we started!)

G. Cowan

Motivating a non-Gaussian prior  $\pi_b(b)$ 

Suppose now the experiment is characterized by

$$y_i, \quad \sigma_i^{\text{stat}}, \quad \sigma_i^{\text{sys}}, \quad s_i, \quad i = 1, \dots, n$$

where  $s_i$  is an (unreported) factor by which the systematic error is over/under-estimated.

Assume correct error for a Gaussian  $\pi_b(b)$  would be  $s_i \sigma_i^{sys}$ , so

$$\pi_b(b_i) = \int \frac{1}{\sqrt{2\pi} s_i \sigma_i^{\text{Sys}}} \exp\left[-\frac{1}{2} \frac{b_i^2}{(s_i \sigma_i^{\text{Sys}})^2}\right] \pi_s(s_i) \, ds_i$$

Width of  $\pi_s(s_i)$  reflects 'error on the error'.

# Error-on-error function $\pi_s(s)$

A simple unimodal probability density for 0 < s < 1 with adjustable mean and variance is the Gamma distribution:

 $\pi_{s}(s) = \frac{a(as)^{b-1}e^{-as}}{\Gamma(b)}$ Want e.g. expectation value of 1 and adjustable standard deviation  $\sigma_{s}$ , i.e.,  $a = b = 1/\sigma_{s}^{2}$   $\max_{s}(s) = \frac{b/a}{variance} = \frac{b}{a^{2}}$   $\pi_{s}(s) = \frac{1}{\sigma_{s}^{2}}$ 

In fact if we took  $\pi_s(s) \sim inverse \ Gamma$ , we could integrate  $\pi_b(b)$  in closed form (cf. D'Agostini, Dose, von Linden). But Gamma seems more natural & numerical treatment not too painful.

Prior for bias  $\pi_b(b)$  now has longer tails

$$\pi_b(b_i) = \int \frac{1}{\sqrt{2\pi} s_i \sigma_i^{\text{Sys}}} \exp\left[-\frac{1}{2} \frac{b_i^2}{(s_i \sigma_i^{\text{Sys}})^2}\right] \pi_s(s_i) \, ds_i$$







Usually summarize posterior  $p(\mu|y)$  with mode and standard deviation:

 $\sigma_{\rm S} = 0.0$ .  $\mu = 1.000 \pm 0.071$  $\sigma_{\rm S} = 0.5$ :  $\hat{\mu} = 1.000 \pm 0.072$ 

# Simple test with inconsistent data

Case #2: there is an outlier

Posterior  $p(\mu|y)$ :



 $\rightarrow$  Bayesian fit less sensitive to outlier.

 $\rightarrow$  Error now connected to goodness-of-fit.

G. Cowan

### Goodness-of-fit vs. size of error

In LS fit, value of minimized  $\chi^2$  does not affect size of error on fitted parameter.

In Bayesian analysis with non-Gaussian prior for systematics, a high  $\chi^2$  corresponds to a larger error (and vice versa).



### Particle Data Group averages

K.A. Olive et al. (Particle Data Group), Chin. Phys. C, 38, 090001 (2014); pdg.lbl.gov

The PDG needs pragmatic solutions for averages where the reported information may be incomplete/inconsistent.

- Often this involves taking the quadratic sum of statistical and systematic uncertainties for LS averages.
- If asymmetric errors (confidence intervals) are reported, PDG has a recipe to reconstruct a model based on a Gaussian-like function where sigma is a continuous function of the mean.
- Exclusion of inconsistent data "sometimes quite subjective".
- If min. chi-squared is much larger than the number of degrees of freedom  $N_{dof} = N-1$ , scale up the input errors a factor

$$S = \left[\chi^2 / (N-1)\right]^{1/2}$$

so that new  $\chi^2 = N_{dof}$ . Error on the average increased by *S*.

# Summary on combinations

The basic idea of combining measurements is to write down the model that describes all of the available experimental outcomes.

If the original data are not available but only parameter estimates, then one treats the estimates (and their covariances) as "the data". Often a multivariate Gaussian model is adequate for these.

If the reported values are limits, there are few meaningful options.

PDG does not combine limits unless the can be "deconstructed" back into a Gaussian measurement.

ATLAS/CMS 2011 combination of Higgs limits used the histograms of event counts (not the individual limits) to construct a full model (ATLAS-CONF-2011-157, CMS PAS HIG-11-023).

Important point is to publish enough information so that meaningful combinations can be carried out.



G. Cowan

# A quick review of frequentist statistical tests

Consider a hypothesis  $H_0$  and alternative  $H_1$ .

A test of  $H_0$  is defined by specifying a critical region w of the data space such that there is no more than some (small) probability  $\alpha$ , assuming  $H_0$  is correct, to observe the data there, i.e.,

$$P(x \in w \mid H_0) \le \alpha$$

Need inequality if data are discrete.

 $\alpha$  is called the size or significance level of the test.

If x is observed in the critical region, reject  $H_0$ .



# Definition of a test (2)

But in general there are an infinite number of possible critical regions that give the same significance level  $\alpha$ .

So the choice of the critical region for a test of  $H_0$  needs to take into account the alternative hypothesis  $H_1$ .

Roughly speaking, place the critical region where there is a low probability to be found if  $H_0$  is true, but high if  $H_1$  is true:



# Type-I, Type-II errors

Rejecting the hypothesis  $H_0$  when it is true is a Type-I error. The maximum probability for this is the size of the test:

$$P(x \in W \mid H_0) \le \alpha$$

But we might also accept  $H_0$  when it is false, and an alternative  $H_1$  is true.

This is called a Type-II error, and occurs with probability

$$P(x \in \mathbf{S} - W \mid H_1) = \beta$$

One minus this is called the power of the test with respect to the alternative  $H_1$ :

Power = 
$$1 - \beta$$

#### *p*-values

Suppose hypothesis *H* predicts pdf  $f(\vec{x}|H)$  for a set of observations  $\vec{x} = (x_1, \dots, x_n)$ .

We observe a single point in this space:  $\vec{x}_{ODS}$ 

What can we say about the validity of *H* in light of the data?

Express level of compatibility by giving the *p*-value for *H*:

p = probability, under assumption of H, to observe data with equal or lesser compatibility with H relative to the data we got.



This is not the probability that *H* is true!

Requires one to say what part of data space constitutes lesser compatibility with *H* than the observed data (implicitly this means that region gives better agreement with some alternative).

# 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.



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

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

E.g. Z = 5 (a "5 sigma effect") corresponds to  $p = 2.9 \times 10^{-7}$ .

G. Cowan

# Using a *p*-value to define test of $H_0$

One can show the distribution of the *p*-value of H, under assumption of H, is uniform in [0,1].

So the probability to find the *p*-value of  $H_0$ ,  $p_0$ , less than  $\alpha$  is

$$P(p_0 \le \alpha | H_0) = \alpha$$

We can define the critical region of a test of  $H_0$  with size  $\alpha$  as the set of data space where  $p_0 \leq \alpha$ .

Formally the *p*-value relates only to  $H_0$ , but the resulting test will have a given power with respect to a given alternative  $H_1$ .

The Poisson counting experiment

Suppose we do a counting experiment and observe *n* events.

Events could be from *signal* process or from *background* – we only count the total number.

Poisson model:

$$P(n|s,b) = \frac{(s+b)^n}{n!}e^{-(s+b)}$$

s = mean (i.e., expected) # of signal events

b = mean # of background events

Goal is to make inference about *s*, e.g.,

test s = 0 (rejecting  $H_0 \approx$  "discovery of signal process")

test all non-zero *s* (values not rejected = confidence interval)

61

In both cases need to ask what is relevant alternative hypothesis. G. Cowan Terascale Statistics School 2015 / Combination Poisson counting experiment: discovery *p*-value Suppose b = 0.5 (known), and we observe  $n_{obs} = 5$ . Should we claim evidence for a new discovery?

Give *p*-value for hypothesis *s* = 0:

$$p$$
-value =  $P(n \ge 5; b = 0.5, s = 0)$   
=  $1.7 \times 10^{-4} \ne P(s = 0)!$ 



G. Cowan

Terascale Statistics School 2015 / Combination

# Poisson counting experiment: discovery significance Equivalent significance for $p = 1.7 \times 10^{-4}$ : $Z = \Phi^{-1}(1-p) = 3.6$ Often claim discovery if Z > 5 ( $p < 2.9 \times 10^{-7}$ , i.e., a "5-sigma effect")



In fact this tradition should be revisited: *p*-value intended to quantify probability of a signallike fluctuation assuming background only; not intended to cover, e.g., hidden systematics, plausibility signal model, compatibility of data with signal, "look-elsewhere effect" (~multiple testing), etc.

# Confidence intervals by inverting a test

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  $P(\text{data in critical region}) \le \alpha$  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 would not be rejected in a test of size  $\alpha$  (confidence level is  $1 - \alpha$ ).

The interval will cover the true value of  $\theta$  with probability  $\geq 1 - \alpha$ . Equivalently, the parameter values in the confidence interval have *p*-values of at least  $\alpha$ .

#### 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  $s_{up}$  at  $CL = 1 - \alpha$  found from

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

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

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

G. Cowan

### Frequentist upper limit on Poisson parameter







G. Cowan

Terascale Statistics School 2015 / Combination

# Frequentist treatment of nuisance parameters in a test

Suppose we test a value of  $\theta$  with the profile likelihood ratio:

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

We want a *p*-value of  $\theta$ :

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

Wilks' theorem says in the large sample limit (and under some additional conditions)  $f(t_{\theta}|\theta, v)$  is a chi-square distribution with number of degrees of freedom equal to number of parameters of interest (number of components in  $\theta$ ).

Simple recipe for *p*-value; holds regardless of the values of the nuisance parameters!

# Frequentist treatment of nuisance parameters in a test (2)

But for a finite data sample,  $f(t_{\theta}|\theta, v)$  still depends on *v*.

So what is the rule for saying whether we reject  $\theta$ ?

Exact approach is to reject  $\theta$  only if  $p_{\theta} < \alpha$  (5%) for all possible *v*.

Some values of  $\theta$  might not be excluded for a value of v known to be disfavoured.

Less values of  $\theta$  rejected  $\rightarrow$  larger interval  $\rightarrow$  higher probability to cover true value ("over-coverage").

But why do we say some values of v are disfavoured? If this is because of other measurements (say, data y) then include y in the model:

 $P(x, y|\theta, \nu) = P_x(x|\theta, \nu)P_y(y|\nu)$ 

Now v better constrained, new interval for  $\theta$  smaller.

# 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  $\theta$  (the profiled values):

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

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

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

#### Bayesian treatment of nuisance parameters

Conceptually straightforward: all parameters have a prior:  $\pi(\theta, \nu)$ 

Often  $\pi(\theta, \nu) = \pi_{\theta}(\theta)\pi_{\nu}(\nu)$ 

Often  $\pi_{\theta}(\theta)$  "non-informative" (broad compared to likelihood). Usually  $\pi_{\nu}(\nu)$  "informative", reflects best available info. on  $\nu$ . Use with likelihood in Bayes' theorem:

 $p(\theta, \nu | x) \propto L(x | \theta, \nu) \pi(\theta, \nu)$ 

To find  $p(\theta|x)$ , marginalize (integrate) over nuisance param.:

$$p(\theta|x) = \int p(\theta, \nu|x) \, d\nu$$

# Prototype analysis in HEP

Each event yields a collection of numbers  $\vec{x} = (x_1, \dots, x_n)$ 

 $x_1$  = number of muons,  $x_2 = p_t$  of jet, ...

 $\vec{x}$  follows some *n*-dimensional joint pdf, which depends on the type of event produced, i.e., signal or background.

1) What kind of decision boundary best separates the two classes?



2) What is optimal test of hypothesis that event sample contains only background?

G. Cowan

#### Test statistics

The boundary of the critical region for an *n*-dimensional data space  $x = (x_1, ..., x_n)$  can be defined by an equation of the form

$$t(x_1,\ldots,x_n)=t_{\rm cut}$$

where  $t(x_1, ..., x_n)$  is a scalar test statistic.

We can work out the pdfs  $g(t|H_0), g(t|H_1), \ldots$ 

Decision boundary is now a single 'cut' on *t*, defining the critical region.

So for an *n*-dimensional problem we have a corresponding 1-d problem.



Terascale Statistics School 2015 / Combination
Test statistic based on likelihood ratio

For multivariate data x, not obvious how to construct best test.

Neyman-Pearson lemma states:

To get the highest power for a given significance level in a test of  $H_0$ , (background) versus  $H_1$ , (signal) the critical region should have

$$\frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)} > c$$

inside the region, and  $\leq c$  outside, where *c* is a constant which depends on the size of the test  $\alpha$ .

Equivalently, optimal scalar test statistic is

$$t(\mathbf{x}) = \frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)}$$

N.B. any monotonic function of this is leads to the same test.

Terascale Statistics School 2015 / Combination

## Ingredients for a frequentist test

In general to carry out a test we need to know the distribution of the test statistic t(x), and this means we need the full model P(x|H).

Often one can construct a test statistic whose distribution approaches a well-defined form (almost) independent of the distribution of the data, e.g., likelihood ratio to test a value of  $\theta$ :

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

In the large sample limit  $t_{\theta}$  follows a chi-square distribution with number of degrees of freedom = number of components in  $\theta$  (Wilks' theorem).

So here one doesn't need the full model  $P(x|\theta)$ , only the observed value of  $t_{\theta}$ .

Terascale Statistics School 2015 / Combination

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  $\alpha = \min | 1,$ 

$$\left[\frac{p(\vec{\theta})q(\vec{\theta}_{0};\vec{\theta})}{p(\vec{\theta}_{0})q(\vec{\theta};\vec{\theta}_{0})}\right]$$

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

## 6) Iterate

G. Cowan

Terascale Statistics School 2015 / Combination

Lecture 5 page 75

## 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 naive  $\sqrt{n}$  .

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.

G. Cowan

Terascale Statistics School 2015 / Combination