Statistical Methods for Particle Physics Day 1: Fundamentals and Parameter Estimation https://indico.desy.de/indico/event/19085/



Helmholtz Alliance



Terascale Statistics School DESY, 19-23 February, 2018

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

Outline

Day 1: Introduction and parameter estimation Probability, random variables, pdfs Parameter estimation maximum likelihood least squares Bayesian parameter estimation Introduction to unfolding

Day 2: Discovery and Limits

Comments on multivariate methods (brief)

p-values

Testing the background-only hypothesis: discovery Testing signal hypotheses: setting limits Experimental sensitivity

Some statistics books, papers, etc.

- G. Cowan, *Statistical Data Analysis*, Clarendon, Oxford, 1998 R.J. Barlow, *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences*, Wiley, 1989
- Ilya Narsky and Frank C. Porter, *Statistical Analysis Techniques in Particle Physics*, Wiley, 2014.
- L. Lyons, Statistics for Nuclear and Particle Physics, CUP, 1986
- F. James., *Statistical and Computational Methods in Experimental Physics*, 2nd ed., World Scientific, 2006
- S. Brandt, *Statistical and Computational Methods in Data Analysis*, Springer, New York, 1998 (with program library on CD) K.A. Olive et al. (Particle Data Group), *Review of Particle Physics*, Chin. Phys. C, 38, 090001 (2014).; see also pdg.lbl.gov sections on probability, statistics, Monte Carlo

Theory ↔ Statistics ↔ Experiment



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: ± 0.4 GeV

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

To do this we need a clear definition of PROBABILITY

G. Cowan

A definition of probability

Consider a set S with subsets A, B, ...

For all $A \subset S, P(A) \ge 0$ P(S) = 1If $A \cap B = \emptyset, P(A \cup B) = P(A) + P(B)$



Kolmogorov axioms (1933)

Also define conditional probability of *A* given *B*:

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

Subsets A, B independent if: $P(A \cap B) = P(A)P(B)$

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

G. Cowan

Interpretation of probability

I. Relative frequency

A, B, ... are outcomes of a repeatable experiment

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

cf. quantum mechanics, particle scattering, radioactive decay...

- II. Subjective probability

 A, B, ... are hypotheses (statements that are true or false)
 P(A) = degree of belief that A is true

 Both interpretations consistent with Kolmogorov axioms.
- In particle physics frequency interpretation often most useful, but subjective probability can provide more natural treatment of non-repeatable phenomena:

systematic uncertainties, probability that Higgs boson exists,...

Bayes' theorem

From the definition of conditional probability we have,

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

but $P(A \cap B) = P(B \cap A)$, so

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

First published (posthumously) by the Reverend Thomas Bayes (1702–1761)

Bayes' theorem



An essay towards solving a problem in the doctrine of chances, Philos. Trans. R. Soc. 53 (1763) 370; reprinted in Biometrika, 45 (1958) 293.

G. Cowan



G. Cowan

An example using Bayes' theorem

Suppose the probability (for anyone) to have a disease D is:

 $P(D) = 0.001 \leftarrow \text{prior probabilities, i.e.,}$ $P(\text{no } D) = 0.999 \leftarrow \text{before any test carried out}$

Consider a test for the disease: result is + or -

P(+|D) = 0.98 P(-|D) = 0.02 \leftarrow probabilities to (in)correctly identify a person with the disease

$$P(+|\text{no D}) = 0.03 \leftarrow \text{probabilities to (in)correctly}$$

 $P(-|\text{no D}) = 0.97 \leftarrow \text{probabilities to (in)correctly}$

Suppose your result is +. How worried should you be?

G. Cowan

Bayes' theorem example (cont.)

The probability to have the disease given a + result is

$$p(\mathbf{D}|+) = \frac{P(+|\mathbf{D})P(\mathbf{D})}{P(+|\mathbf{D})P(\mathbf{D}) + P(+|\mathrm{no} \ \mathbf{D})P(\mathrm{no} \ \mathbf{D})}$$

$= \frac{0.98 \times 0.001}{0.98 \times 0.001 + 0.03 \times 0.999}$

 $= 0.032 \leftarrow \text{posterior probability}$

i.e. you're probably OK!

Your viewpoint: my degree of belief that I have the disease is 3.2%. Your doctor's viewpoint: 3.2% of people like this have the disease.

G. Cowan

Frequentist Statistics – general philosophy

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

Probability = limiting frequency

Probabilities such as

P (Higgs boson exists), *P* (0.117 < $\alpha_{\rm s}$ < 0.121),

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

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

A hypothesis is is preferred if the data are found in a region of high predicted probability (i.e., where an alternative hypothesis predicts lower probability).

G. Cowan

Bayesian Statistics – general philosophy

In Bayesian statistics, use subjective probability for hypotheses:

probability of the data assuming hypothesis *H* (the likelihood) $P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) dH}$ posterior probability, i.e., after seeing the data $P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) dH}$ normalization involves sum over all possible hypotheses

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

No general prescription for priors (subjective!)

Random variables and probability density functions A random variable is a numerical characteristic assigned to an element of the sample space; can be discrete or continuous.

Suppose outcome of experiment is continuous value *x*

$$P(x \text{ found in } [x, x + dx]) = f(x) dx$$

 $\rightarrow f(x) =$ probability density function (pdf)

$$\int_{-\infty}^{\infty} f(x) \, dx = 1 \qquad x \text{ must be somewhere}$$

Or for discrete outcome x_i with e.g. i = 1, 2, ... we have

$$P(x_i) = p_i$$
 probability mass function
 $\sum_i P(x_i) = 1$ x must take on one of its possible values

Other types of probability densities

Outcome of experiment characterized by several values, e.g. an *n*-component vector, $(x_1, ..., x_n)$

$$\rightarrow$$
 joint pdf $f(x_1, \ldots, x_n)$

Sometimes we want only pdf of some (or one) of the components \rightarrow marginal pdf $f_1(x_1) = \int \cdots \int f(x_1, \dots, x_n) dx_2 \dots dx_n$ x_1, x_2 independent if $f(x_1, x_2) = f_1(x_1) f_2(x_2)$

Sometimes we want to consider some components as constant

$$\rightarrow$$
 conditional pdf $g(x_1|x_2) = \frac{f(x_1, x_2)}{f_2(x_2)}$

G. Cowan

Expectation values

Consider continuous r.v. x with pdf f(x). Define expectation (mean) value as $E[x] = \int x f(x) dx$ Notation (often): $E[x] = \mu$ ~ "centre of gravity" of pdf. For a function y(x) with pdf g(y),

$$E[y] = \int y g(y) dy = \int y(x) f(x) dx$$
 (equivalent)

Variance: $V[x] = E[x^2] - \mu^2 = E[(x - \mu)^2]$

Notation: $V[x] = \sigma^2$

Standard deviation: $\sigma = \sqrt{\sigma^2}$

 σ ~ width of pdf, same units as *x*.



G. Cowan

Covariance and correlation

Define covariance cov[x,y] (also use matrix notation V_{xy}) as

$$cov[x, y] = E[xy] - \mu_x \mu_y = E[(x - \mu_x)(y - \mu_y)]$$

Correlation coefficient (dimensionless) defined as

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

If x, y, independent, i.e., $f(x, y) = f_x(x)f_y(y)$, then

$$E[xy] = \iint xy f(x, y) \, dx \, dy = \mu_x \mu_y$$

$$\rightarrow \operatorname{COV}[x, y] = 0 \qquad x \text{ and } y, \text{`uncorrelated'}$$

N.B. converse not always true.

G. Cowan

Correlation (cont.)



G. Cowan

Parameter estimation

The parameters of a pdf are constants that characterize its shape, e.g.



r.v. 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.

G. Cowan

Properties of estimators

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



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

→ average of repeated measurements should tend to true value.
 And we want a small variance (statistical error): V[θ̂]
 → small bias & variance are in general conflicting criteria

G. Cowan

An estimator for the mean (expectation value)

Parameter:
$$\mu = E[x] = \langle x \rangle = \int_{-\infty}^{\infty} x f(x) dx$$

Estimator:
$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i \equiv \overline{x}$$
 ('sample mean')

We find:
$$b = E[\hat{\mu}] - \mu = 0$$

$$V[\hat{\mu}] = \frac{\sigma^2}{n} \qquad \left(\sigma_{\hat{\mu}} = \frac{\sigma}{\sqrt{n}} \right)$$

G. Cowan

An estimator for the variance

Parameter:
$$\sigma^2 = V[x] = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

Estimator:
$$\widehat{\sigma^2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \overline{x})^2 \equiv s^2$$
 ('sample variance')

We find:

$$b = E[\widehat{\sigma^2}] - \sigma^2 = 0$$
 (factor of *n*-1 makes this so)

$$V[\widehat{\sigma^2}] = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \mu_2 \right) , \text{ where}$$
$$\mu_4 = \int (x-\mu)^k f(x) \, dx$$

$$\mu_k = \int (x - \mu)^k f(x) \, dx$$

G. Cowan

The likelihood function

Suppose the entire result of an experiment (set of measurements) is a collection of numbers x, and suppose the joint pdf for the data x is a function that depends on a set of parameters θ :

$$P(\mathbf{x}|\boldsymbol{\theta})$$

Now evaluate this function with the data obtained and regard it as a function of the parameter(s). This is the likelihood function:

$$L(\boldsymbol{\theta}) = P(\mathbf{x}|\boldsymbol{\theta})$$

(*x* constant)

G. Cowan

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

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

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

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

In this case the likelihood function is

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

G. Cowan

Maximum likelihood estimators

If the hypothesized θ is close to the true value, then we expect a high probability to get data like that which we actually found.



So we define the maximum likelihood (ML) estimator(s) to be the parameter value(s) for which the likelihood is maximum.

ML estimators not guaranteed to have any 'optimal' properties, (but in practice they're very good).

G. Cowan

ML example: parameter of exponential pdf

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

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

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

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

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

G. Cowan

ML example: parameter of exponential pdf (2) Find its maximum by setting $\frac{\partial \ln L(\tau)}{\partial \tau} = 0$,

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

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

We find the ML estimate:

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



ML example: parameter of exponential pdf (3) For the exponential distribution one has for mean, variance:

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

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

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

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

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

G. Cowan

Functions of ML estimators

Suppose we had written the exponential pdf as $f(t; \lambda) = \lambda e^{-\lambda t}$, i.e., we use $\lambda = 1/\tau$. What is the ML estimator for λ ?

For a function (with unique inverse) $\lambda(\tau)$ of a parameter τ , it doesn't matter whether we express *L* as a function of λ or τ .

The ML estimator of a function $\lambda(\tau)$ is simply $\hat{\lambda} = \lambda(\hat{\tau})$

So for the decay constant we have
$$\hat{\lambda} = \frac{1}{\hat{\tau}} = \left(\frac{1}{n}\sum_{i=1}^{n} t_i\right)^{-1}$$

Caveat: $\hat{\lambda}$ is biased, even though $\hat{\tau}$ is unbiased.

Can show
$$E[\hat{\lambda}] = \lambda \frac{n}{n-1}$$
. (bias $\rightarrow 0$ for $n \rightarrow \infty$)

G. Cowan

Example of ML: parameters of Gaussian pdf

Consider independent $x_1, ..., x_n$, with $x_i \sim \text{Gaussian}(\mu, \sigma^2)$

$$f(x;\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

The log-likelihood function is

$$\ln L(\mu, \sigma^2) = \sum_{i=1}^n \ln f(x_i; \mu, \sigma^2)$$
$$= \sum_{i=1}^n \left(\ln \frac{1}{\sqrt{2\pi}} + \frac{1}{2} \ln \frac{1}{\sigma^2} - \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

G. Cowan

Example of ML: parameters of Gaussian pdf (2)

Set derivatives with respect to μ , σ^2 to zero and solve,

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i, \qquad \widehat{\sigma^2} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \widehat{\mu})^2.$$

We already know that the estimator for μ is unbiased.

But we find, however, $E[\widehat{\sigma^2}] = \frac{n-1}{n}\sigma^2$, so ML estimator for σ^2 has a bias, but $b \rightarrow 0$ for $n \rightarrow \infty$. Recall, however, that

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \hat{\mu})^{2}$$

is an unbiased estimator for σ^2 .

G. Cowan

Variance of estimators: Monte Carlo method

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

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

For exponential example, from sample variance of estimates we find: $\hat{a} = 0.151$

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

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



Variance of estimators from information inequality

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

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

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

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

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

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

Variance of estimators: graphical method Expand $\ln L(\theta)$ about its maximum:

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

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

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

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

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

G. Cowan

Example of variance by graphical method



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

G. Cowan

Information inequality for *n* parameters Suppose we have estimated *n* parameters $\vec{\theta} = (\theta_1, \dots, \theta_n)$. The (inverse) minimum variance bound is given by the

Fisher information matrix:

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

The information inequality then states that $V - I^{-1}$ is a positive semi-definite matrix, where $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$. Therefore

$$V[\hat{\theta}_i] \ge (I^{-1})_{ii}$$

Often use I^{-1} as an approximation for covariance matrix, estimate using e.g. matrix of 2nd derivatives at maximum of L.

G. Cowan
Example of ML with 2 parameters

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

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



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

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

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

G. Cowan

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

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

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



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

(MINUIT routine HESSE)

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

 $\hat{\sigma}_{\hat{\beta}} = 0.11$ r = 0.46

G. Cowan

Two-parameter fit: MC study Repeat ML fit with 500 experiments, all with n = 2000 events:



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

2

0

0

0.25

0.5

â

0.75

The $\ln L_{\rm max}$ – 1/2 contour

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

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

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

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

G. Cowan

(Co)variances from ln L contour



 \rightarrow Tangent lines to contours give standard deviations.

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

Correlations between estimators result in an increase in their standard deviations (statistical errors).

ML with binned data

Often put data into a histogram: $\vec{n} = (n_1, \dots, n_N), n_{tot} = \sum_{i=1}^N n_i$

Hypothesis is
$$\vec{\nu} = (\nu_1, \dots, \nu_N), \ \nu_{tot} = \sum_{i=1}^N \nu_i$$
 where

$$\nu_i(\vec{\theta}) = \nu_{\text{tot}} \int_{\text{bin } i} f(x; \vec{\theta}) \, dx$$

If we model the data as multinomial (n_{tot} constant),

$$f(\vec{n};\vec{\nu}) = \frac{n_{\text{tot}}!}{n_1! \dots n_N!} \left(\frac{\nu_1}{n_{\text{tot}}}\right)^{n_1} \cdots \left(\frac{\nu_N}{n_{\text{tot}}}\right)^{n_N}$$

then the log-likelihood function is: $\ln L(\vec{\theta}) = \sum_{i=1}^{N} n_i \ln \nu_i(\vec{\theta}) + C$

G. Cowan

ML example with binned data

Previous example with exponential, now put data into histogram:



 $\hat{\tau} = 1.07 \pm 0.17$ (1.06 \pm 0.15 for unbinned ML with same sample)

Limit of zero bin width \rightarrow usual unbinned ML.

If n_i treated as Poisson, we get extended log-likelihood:

$$\ln L(\nu_{\text{tot}}, \vec{\theta}) = -\nu_{\text{tot}} + \sum_{i=1}^{N} n_i \ln \nu_i(\nu_{\text{tot}}, \vec{\theta}) + C$$

G. Cowan

Relationship between ML and Bayesian estimators

In Bayesian statistics, both θ and x are random variables:

 $L(\theta) = L(\vec{x}|\theta) = f_{\text{joint}}(\vec{x}|\theta)$

Recall the Bayesian method:

Use subjective probability for hypotheses (θ); before experiment, knowledge summarized by prior pdf $\pi(\theta)$; use Bayes' theorem to update prior in light of data:

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

Posterior pdf (conditional pdf for θ given x)

G. Cowan

ML and Bayesian estimators (2) Purist Bayesian: $p(\theta | x)$ contains all knowledge about θ . Pragmatist Bayesian: $p(\theta | x)$ could be a complicated function, \rightarrow summarize using an estimator $\hat{\theta}_{Bayes}$ Take mode of $p(\theta | x)$, (could also use e.g. expectation value)

What do we use for $\pi(\theta)$? No golden rule (subjective!), often represent 'prior ignorance' by $\pi(\theta)$ = constant, in which case

$$\hat{\theta}_{\mathsf{Bayes}} = \hat{\theta}_{\mathsf{ML}}$$

But... we could have used a different parameter, e.g., $\lambda = 1/\theta$, and if prior $\pi_{\theta}(\theta)$ is constant, then $\pi_{\lambda}(\lambda) = \pi_{\theta}(\theta(\lambda)) |d\theta/d\lambda|$ is not!

'Complete prior ignorance' is not well defined.

G. Cowan

Priors from formal rules

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 a Subjective Bayesian analysis, using objective priors can be an important part of the sensitivity analysis.

Priors from formal rules (cont.)

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 certain coverage properties. For a review see:

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 HEP, but there is recent interest in this direction; see e.g.

L. Demortier, S. Jain and H. Prosper, *Reference priors for high energy physics*, arxiv:1002.1111 (Feb 2010)

Jeffreys' prior

According to Jeffreys' rule, take prior according to

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

where

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

is the Fisher information matrix.

One can show that this leads to inference that is invariant under a transformation of parameters.

For a Gaussian mean, the Jeffreys' prior is constant; for a Poisson mean μ it is proportional to $1/\sqrt{\mu}$.

"Invariance of inference" with Jeffreys' prior Suppose we have a parameter θ , to which we assign a prior $\pi_{\theta}(\theta)$. An experiment gives data *x*, modeled by $L(\theta) = P(x|\theta)$. Bayes' theorem then tells us the posterior for θ :

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

Now consider a function $\eta(\theta)$, and we want the posterior $P(\eta|x)$. This must follow from the usual rules of transformation of random variables:

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

G. Cowan

"Invariance of inference" with Jeffreys' prior (2)

Alternatively, we could have just starting with η as the parameter in our model, and written down a prior pdf $\pi_n(\eta)$.

Using it, we express the likelihood as $L(\eta) = P(x|\eta)$ and write Bayes' theorem as

 $P(\eta|x) \propto P(x|\eta)\pi_{\eta}(\eta)$

If the priors really express our degree of belief, then they must be related by the usual laws of probability $\pi_{\eta}(\eta) = \pi_{\theta}(\theta(\eta)) |d\theta/d\eta|$, and in this way the two approaches lead to the same result.

But if we choose the priors according to "formal rules", then this is not guaranteed. For the Jeffrey's prior, however, it does work!

Using $\pi_{\theta}(\theta) \propto \sqrt{I(\theta)}$ and transforming to find $P(\eta|x)$ leads to the same as using $\pi_{\eta}(\eta) \propto \sqrt{I(\eta)}$ directly with Bayes' theorem.

Jeffreys' prior for Poisson mean

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

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

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

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

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*.

G. Cowan

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 θ , 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 χ^2 numerically (e.g. program **MINUIT**).

G. Cowan

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

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

G. Cowan

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 $A_{ij} = a_j(x_i)$,

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

G. Cowan

Linear LS problem (2)

Set derivitives with respect to θ_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 .

G. Cowan

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)

G. Cowan

Example of least squares fit

Fit a polynomial of order *p*: $\lambda(x; \theta_0, \dots, \theta_p) = \sum_{n=0}^{p} \theta_n x^n$



Variance of LS estimators

In most cases of interest we obtain the variance in a manner similar to ML. E.g. for data ~ Gaussian we have

$$\chi^2(\theta) = -2\ln L(\theta) + C$$

and so

$$\widehat{\sigma^2}_{\widehat{\theta}} \approx 2 \left[\frac{\partial^2 \chi^2}{\partial \theta^2} \right]_{\theta = \widehat{\theta}}^{-1}$$

or for the graphical method we take the values of θ where

$$\chi^2(\theta) = \chi^2_{\min} + 1$$



Two-parameter LS fit

2-parameter case (line with nonzero slope):



Angle of ellipse \rightarrow correlation (same as for ML)

G. Cowan

Goodness-of-fit with least squares

The value of the χ^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 χ^2_{min} as high as the one we got, or higher, if the hypothesis is correct.

E.g. for the previous example with 1st order polynomial (line),

 $\chi^2_{\rm min} = 3.99$, $n_{\rm d} = 5 - 2 = 3$, p = 0.263

whereas for the 0th order polynomial (horizontal line),

$$\chi^2_{\rm min} = 45.5$$
, $n_{\rm d} = 5 - 1 = 4$, $p = 3.1 \times 10^{-9}$

G. Cowan

Goodness-of-fit vs. statistical errors

Small statistical error does not mean a good fit (nor vice versa).

Curvature of χ^2 near its minimum \rightarrow statistical errors $(\sigma_{\hat{\theta}})$ Value of $\chi^2_{\min} \rightarrow$ goodness-of-fit

Horizontal line fit, move the data points, keep errors on points same:



G. Cowan

Goodness-of-fit vs. stat. errors (2)

 $\rightarrow \chi^2(\theta_0)$ shifted down, same curvature as before.

Variance of estimator (statistical error) tells us:

if experiment repeated many times, how wide is the distribution of the estimates $\hat{\theta}$. (Doesn't tell us whether hypothesis correct.) P-value tells us:

if hypothesis is correct and experiment repeated many times, what fraction will give equal or worse agreement between data and hypothesis according to the statistic χ^2_{\min} .

Low P-value \rightarrow hypothesis may be wrong \rightarrow systematic error.

LS with binned data



We have

 $y_i =$ number of entries in bin i,

$$\lambda_i(ec{ heta}) = n \int_{x_i^{\min}}^{x_i^{\max}} f(x;ec{ heta}) dx = n p_i(ec{ heta})$$

G. Cowan

LS with binned data (2)

LS fit: minimize

$$\chi^2(ec{ heta}) = \sum\limits_{i=1}^N rac{(y_i - \lambda_i(ec{ heta}))^2}{\sigma_i^2}$$

where $\sigma_i^2 = V[y_i]$, here not known a priori.

Treat the y_i as Poisson r.v.s, in place of true variance take either

 $\sigma_i^2 = \lambda_i(\vec{\theta})$ (LS method)

 $\sigma_i^2 = y_i$ (Modified LS method)

MLS sometimes easier computationally, but χ^2_{min} no longer follows chi-square pdf (or is undefined) if some bins have few (or no) entries.

G. Cowan

LS with binned data — normalization Do not 'fit the normalization':

$$\lambda_i(ec{ heta},
u) =
u \int_{x_i^{ ext{min}}}^{x_i^{ ext{max}}} f(x;ec{ heta}) dx =
u p_i(ec{ heta})$$

i.e. introduce adjustable ν , fit along with θ .

 $\hat{\nu}$ is a bad estimator for n (which we know, anyway!)

$$\hat{\nu}_{\rm LS} = n + \frac{\chi^2_{\rm min}}{2}$$

$$\hat{
u}_{ ext{MLS}} = n - \chi^2_{ ext{min}}$$

LS normalization example

Example with n = 400 entries, N = 20 bins:



Expect χ^2_{\min} around N-m,

 \rightarrow relative error in $\hat{\nu}$ large when N large, n small Either get n directly from data for LS (or better, use ML).

G. Cowan

Using LS to combine measurements

Use LS to obtain weighted average of N measurements of λ :

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

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

G. Cowan

Combining correlated measurements with LS

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

LS λ has zero bias, minimum variance (Gauss–Markov theorem).

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. ρ , σ_1 , σ_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.

G. Cowan, Statistical Data Analysis, Oxford University Press, 1998.

G. Cowan
Unfolding: formulation of the problem

Consider a random variable *y*, goal is to determine pdf f(y). If parameterization $f(y;\theta)$ known, find e.g. ML estimators $\hat{\theta}$. If no parameterization available, construct histogram:



G. Cowan

Migration

Effect of measurement errors: y = true value, x = observed value, migration of entries between bins,

f(y) is 'smeared out', peaks broadened.

 $f_{\text{meas}}(x) = \int R(x|y) f_{\text{true}}(y) \, dy$ $\downarrow \qquad \text{discretize: data are } \mathbf{n} = (n_1, \dots, n_N)$ $\nu_i = E[n_i] = \sum_{j=1}^M R_{ij} \mu_j , \quad i = 1, \dots, N$ $R_{ij} = P(\text{observed in bin } i \mid \text{true in bin } j)$ response matrix

Note μ , v are constants; n subject to statistical fluctuations.

G. Cowan

Efficiency, background

Sometimes an event goes undetected:

 $\sum_{i=1}^{N} R_{ij} = \sum_{i=1}^{N} P(\text{observed in bin } i \,|\, \text{true value in bin } j)$

= P(observed anywhere | true value in bin j)

 $= \varepsilon_j$ \leftarrow efficiency

Sometimes an observed event is due to a background process:

$$\nu_i = \sum_{j=1}^M R_{ij}\mu_j + \beta_i$$

 β_i = expected number of background events in *observed* histogram. For now, assume the β_i are known.

G. Cowan

The basic ingredients



Summary of ingredients 'true' histogram: $\mu = (\mu_1, \dots, \mu_M), \quad \mu_{\text{tot}} = \sum_{i=1}^{N} \mu_i$ probabilities: $\mathbf{p} = (p_1, \dots, p_M) = \boldsymbol{\mu} / \boldsymbol{\mu}_{\text{tot}}$ expectation values for observed histogram: $\nu = (\nu_1, \dots, \nu_N)$ observed histogram: $\mathbf{n} = (n_1, \dots, n_N)$ response matrix: $R_{ij} = P(\text{observed in bin } i \mid \text{true in bin } j)$ efficiencies: $\varepsilon_j = \sum_{i=1}^{N} R_{ij}$ expected background: $\beta = (\beta_1, \dots, \beta_N)$ These are related by: $E[\mathbf{n}] = \boldsymbol{\nu} = R\boldsymbol{\mu} + \boldsymbol{\beta}$

G. Cowan

Maximum likelihood (ML) estimator from inverting the response matrix

Assume $\nu = R\mu + \beta$ can be inverted: $\mu = R^{-1}(\nu - \beta)$

Suppose data are independent Poisson: $P(n_i; \nu_i) = \frac{\nu_i^{n_i}}{n_i!} e^{-\nu_i}$

So the log-likelihood is
$$\ln L(\boldsymbol{\mu}) = \sum_{i=1}^{N} (n_i \ln \nu_i - \nu_i)$$

ML estimator is $\hat{\boldsymbol{\nu}} = \mathbf{n}$

$$\longrightarrow \hat{\mu} = R^{-1}(\mathbf{n} - \boldsymbol{\beta})$$

Example with ML solution



Catastrophic failure???

G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 1

What went wrong?



Applying R^{-1} to $\vec{\nu}$ puts the fine structure back: $\vec{\mu} = R^{-1}\vec{\nu}$. But we don't have ν , only n. R^{-1} "thinks" fluctuations in n are the residual of original fine structure, puts this back into $\hat{\mu}$.

G. Cowan

ML solution revisited

For Poisson data the ML estimators are unbiased:

$$E[\hat{\boldsymbol{\mu}}] = R^{-1}(E[\mathbf{n}] - \boldsymbol{\beta}) = \boldsymbol{\mu}$$

Their covariance is:

$$U_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j] = \sum_{k,l=1}^N (R^{-1})_{ik} (R^{-1})_{jl} \operatorname{cov}[n_k, n_l]$$
$$= \sum_{k=1}^N (R^{-1})_{ik} (R^{-1})_{jk} \nu_k$$

(Recall these statistical errors were huge for the example shown.)

G. Cowan

ML solution revisited (2)

The information inequality gives for unbiased estimators the minimum (co)variance bound:

$$(U^{-1})_{kl} = -E\left[\frac{\partial^2 \log L}{\partial \mu_k \partial \mu_l}\right] = \sum_{i=1}^N \frac{R_{ik} R_{il}}{\nu_i}$$

invert
$$\rightarrow \quad U_{ij} = \sum_{k=1}^{N} (R^{-1})_{ik} (R^{-1})_{jk} \nu_k$$

This is the same as the actual variance! I.e. ML solution gives smallest variance among all unbiased estimators, even though this variance was huge.

In unfolding one must accept some bias in exchange for a (hopefully large) reduction in variance.

Correction factor method

Use equal binning for $\vec{\mu}$, $\vec{\nu}$ and take $\hat{\mu}_i = C_i(n_i - \beta_i)$, where

$$C_i = \frac{\mu_i^{\text{MC}}}{\nu_i^{\text{MC}}} \qquad \begin{array}{l} \nu_i^{\text{MC}} \text{ and } \mu_i^{\text{MC}} \text{ from Monte Carlo} \\ \text{simulation (no background).} \end{array}$$

$$U_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j] = C_i^2 \operatorname{cov}[n_i, n_j]$$

Often $C_i \sim O(1)$ so statistical errors far smaller than for ML.

But the bias
$$b_i = E[\hat{\mu}_i] - \mu_i$$
 is $b_i = \left(\frac{\mu_i^{\text{MC}}}{\nu_i^{\text{MC}}} - \frac{\mu_i}{\nu_i^{\text{sig}}}\right)$
Nonzero bias unless MC = Nature.
 $\nu_i^{\text{sig}} = \nu_i - \beta_i$

G. Cowan

Reality check on the statistical errors

Suppose for some bin *i* we have:

$$C_i = 0.1 \qquad \qquad \beta_i = 0 \qquad \qquad n_i = 100$$

$$\rightarrow \hat{\mu}_i = C_i n_i = 10$$
 $\sigma_{\hat{\mu}_i} = C_i \sqrt{n_i} = 1.0$ (10% stat.
error)

But according to the estimate, only 10 of the 100 events found in the bin belong there; the rest spilled in from outside.

How can we have a 10% measurement if it is based on only 10 events that really carry information about the desired parameter?

Discussion of correction factor method

As with all unfolding methods, we get a reduction in statistical error in exchange for a bias; here the bias is difficult to quantify (difficult also for many other unfolding methods).

The bias should be small if the bin width is substantially larger than the resolution, so that there is not much bin migration.

So if other uncertainties dominate in an analysis, correction factors may provide a quick and simple solution (a "first-look").

Still the method has important flaws and it would be best to avoid it.

Regularized unfolding

Consider 'reasonable' estimators such that for some $\Delta \log L$,

 $\log L(\vec{\mu}) \ge \log L_{\max} - \Delta \log L$

Out of these estimators, choose the 'smoothest', by maximizing

 $\Phi(\vec{\mu}) = \alpha \, \log L(\vec{\mu}) \, + \, S(\vec{\mu}),$

 $S(\vec{\mu}) =$ regularization function (measure of smoothness),

 α = regularization parameter (choose to give desired $\Delta \log L$)

Regularized unfolding (2)

In addition require $\sum_{i=1}^{N} \nu_i = \sum_{i,j} R_{ij} \mu_j = n_{\text{tot}}$, i.e. maximize

$$\varphi(\vec{\mu}, \lambda) = \alpha \log L(\vec{\mu}) + S(\vec{\mu}) + \lambda \left[n_{\text{tot}} - \sum_{i=1}^{N} \nu_i \right]$$

where λ is a Lagrange multiplier, $\partial \varphi / \partial \lambda = 0 \rightarrow \sum_{i=1}^{N} \nu_i = n_{\text{tot}}$.

 $\alpha = 0$ gives smoothest solution (ignores data!),

 $\alpha \to \infty$ gives ML solution (variance too large).

We need: regularization function $S(\vec{\mu})$, a prescription for setting α .

G. Cowan

Tikhonov regularization

Take measure of smoothness = mean square of kth derivative,

$$S[f_{ ext{true}}(y)] = - \int \left(rac{d^k f_{ ext{true}}(y)}{dy^k}
ight)^2 dy \;, ext{ where } k=1,2,\dots$$

If we use Tikhonov (k = 2) with $\log L = -\frac{1}{2}\chi^2$,

$$S(\boldsymbol{\mu}) = -\sum_{i=1}^{M-2} (-\mu_i + 2\mu_{i+1} - \mu_{i+2})^2$$

 $\varphi(\vec{\mu}, \lambda) = -\frac{\alpha}{2}\chi^2(\vec{\mu}) + S(\vec{\mu})$ quadratic in μ_i ,

 \rightarrow setting derivatives of φ equal to zero gives linear equations. Solution using Singular Value Decomposition (SVD).

G. Cowan

SVD implementation of Tikhonov unfolding

A. Hoecker, V. Kartvelishvili, NIM A372 (1996) 469; (TSVDUnfold in ROOT).

Minimizes
$$\chi^2(\mu) + \tau \sum_{i} \left[(\mu_{i+1} - \mu_i) - (\mu_i - \mu_{i-1})^2 \right]$$

Numerical implementation using Singular Value Decomposition. Recommendations for setting regularization parameter τ :

> Transform variables so errors ~ Gauss(0,1); number of transformed values significantly different from zero gives prescription for τ ; or base choice of τ on unfolding of test distributions.

SVD example

A. Höcker, V. Kartvelishvili, NIM A372 (1996) 469.



G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 1

Regularization function based on entropy

Shannon entropy of a set of probabilities is

$$H = -\sum_{i=1}^{M} p_i \log p_i$$

All p_i equal \rightarrow maximum entropy (maximum smoothness)

One $p_i = 1$, all others $= 0 \rightarrow \text{minimum entropy}$

Use entropy as regularization function,

$$S(\vec{\mu}) = H(\vec{\mu}) = -\sum_{i=1}^{M} \frac{\mu_i}{\mu_{\text{tot}}} \log \frac{\mu_i}{\mu_{\text{tot}}}$$

 $\propto \log(\text{number of ways to arrange } \mu_{\text{tot}} \text{ entries in } M \text{ bins})$

Can have Bayesian motivation: $S(\vec{\mu}) \rightarrow \text{prior pdf for } \vec{\mu}$ G. Cowan DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 1

Example of entropy-based unfolding



G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 1

Estimating bias and variance

In general, the equations determining $\hat{\vec{\mu}}(\vec{n})$ are nonlinear. Expand $\hat{\vec{\mu}}(\vec{n})$ about \vec{n}_{obs} (observed data set),

Use error propagation to get covariance $U_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j],$

and estimators for the bias, $b_i = E[\hat{\mu}_i] - \mu_i$,

$$\hat{b}_{i} = \sum_{j=1}^{N} \frac{\partial \hat{\mu}_{i}}{\partial n_{j}} (\hat{\nu}_{j} - n_{j}),$$
where $\hat{\vec{\nu}} = R\hat{\vec{\mu}} + \vec{\beta}$. (N.B. $\hat{\vec{\nu}} \neq \vec{n}$.)

G. Cowan

Choosing the regularization parameter

 $\alpha = 0 \rightarrow \vec{\mu}$ maximally smooth (ignores data).

 $\alpha \to \infty \to ML$ solution (no bias, very large variance).

Possible criteria for best trade-off between bias and variance:

Minimize mean squared error,

$$MSE = \frac{1}{M} \sum_{i=1}^{M} (U_{ii} + \hat{b}_i^2), \text{ or}$$
$$MSE' = \frac{1}{M} \sum_{i=1}^{M} \frac{U_{ii} + \hat{b}_i^2}{\hat{\mu}_i}.$$

G. Cowan

G. Cowan, Statistical Data Analysis, OUP (1998) Ch. 11

Choosing the regularization parameter (2)

Or look at changes in χ^2 from unregularized (ML) solution,

$$\Delta \chi^2 = 2\Delta \log L = N$$

Or require that bias be consistent with zero to within its own error,

$$\chi_b^2 = \sum_{i=1}^M \frac{\hat{b}_i^2}{W_{ii}} = M \text{ where } W_{ij} = \operatorname{cov}[\hat{b}_i, \hat{b}_j].$$

i.e. if bias significantly different from zero, we would subtract it; \rightarrow equivalent to going to smaller $\Delta \log L$ or larger α (less bias). G. Cowan, Statistical Data Analysis, OUP (1998) Ch. 11

Some examples with Tikhonov regularization



G. Cowan

Some examples with entropy regularization



G. Cowan

Stat. and sys. errors of unfolded solution

In general the statistical covariance matrix of the unfolded estimators is not diagonal; need to report full

$$U_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j]$$

But unfolding necessarily introduces biases as well, corresponding to a systematic uncertainty (also correlated between bins).

This is more difficult to estimate. Suppose, nevertheless, we manage to report both U_{stat} and U_{sys} .

To test a new theory depending on parameters θ , use e.g.

$$\chi^2(\boldsymbol{\theta}) = (\boldsymbol{\mu}(\boldsymbol{\theta}) - \hat{\boldsymbol{\mu}})^T (U_{\text{stat}} + U_{\text{sys}})^{-1} (\boldsymbol{\mu}(\boldsymbol{\theta}) - \hat{\boldsymbol{\mu}})$$

Mixes frequentist and Bayesian elements; interpretation of result can be problematic, especially if U_{sys} itself has large uncertainty.

Folding

Suppose a theory predicts $f(y) \rightarrow \mu$ (may depend on parameters θ). Given the response matrix *R* and expected background β , this predicts the expected numbers of observed events:

$$\boldsymbol{\nu} = R\boldsymbol{\mu} + \boldsymbol{\beta}$$

From this we can get the likelihood, e.g., for Poisson data,

$$L(\mathbf{n}|\boldsymbol{\nu}) = \prod_{i=1}^{N} \frac{\nu_i^{n_i}}{n_i!} e^{-\nu_i}$$

And using this we can fit parameters and/or test, e.g., using the likelihood ratio statistic

$$q = -2\ln\frac{L(\mathbf{n}|\boldsymbol{\nu})}{L(\mathbf{n}|\hat{\boldsymbol{\nu}})} \sim \chi_N^2$$

G. Cowan

Versus unfolding

If we have an unfolded spectrum and full statistical and systematic covariance matrices, to compare this to a model μ compute likelihood

$$L(\hat{\boldsymbol{\mu}}|\boldsymbol{\mu}) \sim e^{-\chi^2/2}$$

where

$$\chi^2 = (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})^T (U_{\text{stat}} + U_{\text{sys}})^{-1} (\boldsymbol{\mu} - \hat{\boldsymbol{\mu}})$$

Complications because one needs estimate of systematic bias U_{sys} .

If we find a gain in sensitivity from the test using the unfolded distribution, e.g., through a decrease in statistical errors, then we are exploiting information inserted via the regularization (e.g., imposed smoothness).

ML solution again

From the standpoint of testing a theory or estimating its parameters, the ML solution, despite catastrophically large errors, is equivalent to using the uncorrected data (same information content).

There is no bias (at least from unfolding), so use

$$\chi^2(\boldsymbol{\theta}) = (\boldsymbol{\mu}(\boldsymbol{\theta}) - \hat{\boldsymbol{\mu}}_{\mathrm{ML}})^T U_{\mathrm{stat}}^{-1}(\boldsymbol{\mu}(\boldsymbol{\theta}) - \hat{\boldsymbol{\mu}}_{\mathrm{ML}})$$

The estimators of θ should have close to optimal properties: zero bias, minimum variance.

The corresponding estimators from any unfolded solution cannot in general match this.

Crucial point is to use full covariance, not just diagonal errors.

Unfolding discussion

Unfolding can be a minefield and is not necessary if goal is to compare measured distribution with a model prediction.

Even comparison of uncorrected distribution with *future* theories not a problem, as long as it is reported together with the expected background and response matrix.

> In practice complications because these ingredients have uncertainties, and they must be reported as well.

Unfolding useful for getting an actual estimate of the distribution we think we've measured; can e.g. compare ATLAS/CMS.

Model test using unfolded distribution should take account of the (correlated) bias introduced by the unfolding procedure.

Finally...

Estimation of parameters is usually the "easy" part of statistics:

Frequentist: maximize the likelihood.

Bayesian: find posterior pdf and summarize (e.g. mode).

Standard tools for quantifying precision of estimates: Variance of estimators, confidence intervals,...

But there are many potential stumbling blocks:

bias versus variance trade-off (how many parameters to fit?);goodness of fit (usually only for LS or binned data);choice of prior for Bayesian approach;unexpected behaviour in LS averages with correlations,...

Extra slides

G. Cowan

Some distributions

Distribution/pdf **Binomial** Multinomial Poisson Uniform Exponential Gaussian Chi-square Cauchy Landau Beta Gamma Student's t

Example use in HEP **Branching** ratio Histogram with fixed NNumber of events found Monte Carlo method Decay time Measurement error Goodness-of-fit Mass of resonance **Ionization energy loss** Prior pdf for efficiency Sum of exponential variables Resolution function with adjustable tails

G. Cowan

Binomial distribution

Consider *N* independent experiments (Bernoulli trials): outcome of each is 'success' or 'failure', probability of success on any given trial is *p*.

Define discrete r.v. n = number of successes ($0 \le n \le N$).

Probability of a specific outcome (in order), e.g. 'ssfsf' is $pp(1-p)p(1-p) = p^n(1-p)^{N-n}$ N!

But order not important; there are

 $\frac{1}{n!(N-n)!}$

ways (permutations) to get *n* successes in *N* trials, total probability for *n* is sum of probabilities for each permutation.

G. Cowan

Binomial distribution (2)

The binomial distribution is therefore

$$f(n; N, p) = \frac{N!}{n!(N-n)!}p^n(1-p)^{N-n}$$
random parameters
variable

For the expectation value and variance we find:

$$E[n] = \sum_{n=0}^{N} nf(n; N, p) = Np$$
$$V[n] = E[n^{2}] - (E[n])^{2} = Np(1 - p)$$

G. Cowan

Binomial distribution (3)

Binomial distribution for several values of the parameters:



Example: observe *N* decays of W^{\pm} , the number *n* of which are $W \rightarrow \mu \nu$ is a binomial r.v., *p* = branching ratio.

G. Cowan
Multinomial distribution

Like binomial but now *m* outcomes instead of two, probabilities are

$$\vec{p} = (p_1, \dots, p_m)$$
, with $\sum_{i=1}^m p_i = 1$.

For N trials we want the probability to obtain:

 n_1 of outcome 1, n_2 of outcome 2, \vdots n_m of outcome *m*.

This is the multinomial distribution for $\vec{n} = (n_1, \ldots, n_m)$

$$f(\vec{n}; N, \vec{p}) = \frac{N!}{n_1! n_2! \cdots n_m!} p_1^{n_1} p_2^{n_2} \cdots p_m^{n_m}$$

G. Cowan

Multinomial distribution (2)

Now consider outcome *i* as 'success', all others as 'failure'.

 \rightarrow all n_i individually binomial with parameters N, p_i

$$E[n_i] = Np_i, \quad V[n_i] = Np_i(1-p_i) \quad \text{for all } i$$

One can also find the covariance to be

$$V_{ij} = Np_i(\delta_{ij} - p_j)$$

Example: $\vec{n} = (n_1, \dots, n_m)$ represents a histogram with *m* bins, *N* total entries, all entries independent.

G. Cowan

Poisson distribution

Consider binomial *n* in the limit

 $N \to \infty, \qquad p \to 0, \qquad E[n] = Np \to \nu.$

 \rightarrow *n* follows the Poisson distribution:

$$f(n;\nu) = \frac{\nu^n}{n!}e^{-\nu} \quad (n \ge 0)$$

$$E[n] = \nu, \quad V[n] = \nu.$$

Example: number of scattering events *n* with cross section σ found for a fixed integrated luminosity, with $\nu = \sigma \int L dt$.



10

n

15

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 1

0

0

5

20

Uniform distribution

Consider a continuous r.v. *x* with $-\infty < x < \infty$. Uniform pdf is:



N.B. For any r.v. *x* with cumulative distribution F(x), y = F(x) is uniform in [0,1].

Example: for $\pi^0 \to \gamma \gamma$, E_{γ} is uniform in $[E_{\min}, E_{\max}]$, with $E_{\min} = \frac{1}{2} E_{\pi} (1 - \beta)$, $E_{\max} = \frac{1}{2} E_{\pi} (1 + \beta)$

G. Cowan

Exponential distribution

The exponential pdf for the continuous r.v. *x* is defined by:



Example: proper decay time *t* of an unstable particle

 $f(t;\tau) = \frac{1}{\tau}e^{-t/\tau}$ (τ = mean lifetime)

Lack of memory (unique to exponential): $f(t - t_0 | t \ge t_0) = f(t)$

G. Cowan

Gaussian distribution

The Gaussian (normal) pdf for a continuous r.v. *x* is defined by:

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} e^{-(x-\mu)^2/2\sigma^2}$$

$$E[x] = \mu$$
(N.B. often μ, σ^2 denote mean, variance of any

$$V[x] = \sigma^2$$
r.v., not only Gaussian.)



х

Special case: $\mu = 0$, $\sigma^2 = 1$ ('standard Gaussian'):

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} , \quad \Phi(x) = \int_{-\infty}^x \varphi(x') \, dx'$$

If $y \sim$ Gaussian with μ , σ^2 , then $x = (y - \mu) / \sigma$ follows $\varphi(x)$.

G. Cowan

Gaussian pdf and the Central Limit Theorem

The Gaussian pdf is so useful because almost any random variable that is a sum of a large number of small contributions follows it. This follows from the Central Limit Theorem:

For *n* independent r.v.s x_i with finite variances σ_i^2 , otherwise arbitrary pdfs, consider the sum

$$y = \sum_{i=1}^{n} x_i$$

In the limit $n \to \infty$, y is a Gaussian r.v. with

$$E[y] = \sum_{i=1}^{n} \mu_i \qquad V[y] = \sum_{i=1}^{n} \sigma_i^2$$

Measurement errors are often the sum of many contributions, so frequently measured values can be treated as Gaussian r.v.s.

G. Cowan

Central Limit Theorem (2)

The CLT can be proved using characteristic functions (Fourier transforms), see, e.g., SDA Chapter 10.

For finite *n*, the theorem is approximately valid to the extent that the fluctuation of the sum is not dominated by one (or few) terms.



Beware of measurement errors with non-Gaussian tails.

Good example: velocity component v_x of air molecules.

OK example: total deflection due to multiple Coulomb scattering. (Rare large angle deflections give non-Gaussian tail.)

Bad example: energy loss of charged particle traversing thin gas layer. (Rare collisions make up large fraction of energy loss, cf. Landau pdf.)

Multivariate Gaussian distribution

Multivariate Gaussian pdf for the vector $\vec{x} = (x_1, \dots, x_n)$:

$$f(\vec{x};\vec{\mu},V) = \frac{1}{(2\pi)^{n/2}|V|^{1/2}} \exp\left[-\frac{1}{2}(\vec{x}-\vec{\mu})^T V^{-1}(\vec{x}-\vec{\mu})\right]$$

 $\vec{x}, \vec{\mu}$ are column vectors, $\vec{x}^T, \vec{\mu}^T$ are transpose (row) vectors,

$$E[x_i] = \mu_i, \quad \operatorname{cov}[x_i, x_j] = V_{ij}.$$

For n = 2 this is $f(x_1, x_2; \mu_1, \mu_2, \sigma_1, \sigma_2, \rho) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}}$ $\times \exp\left\{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x_1 - \mu_1}{\sigma_1}\right)^2 + \left(\frac{x_2 - \mu_2}{\sigma_2}\right)^2 - 2\rho \left(\frac{x_1 - \mu_1}{\sigma_1}\right) \left(\frac{x_2 - \mu_2}{\sigma_2}\right) \right] \right\}$

where $\rho = \operatorname{cov}[x_1, x_2]/(\sigma_1 \sigma_2)$ is the correlation coefficient.

G. Cowan

Chi-square (χ^2) distribution

The chi-square pdf for the continuous r.v. $z \ (z \ge 0)$ is defined by

$$f(z;n) = \frac{1}{2^{n/2} \Gamma(n/2)} z^{n/2-1} e^{-z/2} \left\{ \begin{array}{c} 0.5 \\ 0.4 \\ \dots & n=2 \\ \dots & n=5 \\ 0.3 \\ \dots & n=10 \end{array} \right\}$$

$$n = 1, 2, \dots = \text{ number of 'degrees of freedom' (dof)}$$

$$E[z] = n, \quad V[z] = 2n.$$

For independent Gaussian x_i , i = 1, ..., n, means μ_i , variances σ_i^2 ,

$$z = \sum_{i=1}^{n} \frac{(x_i - \mu_i)^2}{\sigma_i^2} \quad \text{follows } \chi^2 \text{ pdf with } n \text{ dof.}$$

Example: goodness-of-fit test variable especially in conjunction with method of least squares.

G. Cowan

Cauchy (Breit-Wigner) distribution

The Breit-Wigner pdf for the continuous r.v. x is defined by

$$f(x; \Gamma, x_0) = \frac{1}{\pi} \frac{\Gamma/2}{\Gamma^2/4 + (x - x_0)^2}$$

$$(\Gamma = 2, x_0 = 0 \text{ is the Cauchy pdf.})$$

$$E[x] \text{ not well defined, } V[x] \to \infty.$$

$$x_0 = \text{ mode (most probable value)}$$

$$\Gamma = \text{ full width at half maximum}$$

Example: mass of resonance particle, e.g. ρ , K^{*}, ϕ^0 , ... Γ = decay rate (inverse of mean lifetime)

G. Cowan

Landau distribution

For a charged particle with $\beta = v/c$ traversing a layer of matter of thickness *d*, the energy loss Δ follows the Landau pdf:



L. Landau, J. Phys. USSR **8** (1944) 201; see also W. Allison and J. Cobb, Ann. Rev. Nucl. Part. Sci. **30** (1980) 253.

G. Cowan

Landau distribution (2)



Beta distribution

$$f(x;\alpha,\beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$



Often used to represent pdf of continuous r.v. nonzero only between finite limits.



Gamma distribution

$$f(x; \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^{\alpha}} x^{\alpha-1} e^{-x/\beta}$$

$$V[x] = \alpha \beta^2$$

 $E[r] = \alpha \beta$

Often used to represent pdf of continuous r.v. nonzero only in $[0,\infty]$.

Also e.g. sum of *n* exponential r.v.s or time until *n*th event in Poisson process ~ Gamma



Student's t distribution

$$f(x;\nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\,\Gamma(\nu/2)} \left(1 + \frac{x^2}{\nu}\right)^{-\left(\frac{\nu+1}{2}\right)}$$



G. Cowan