Statistical Methods for Particle Physics (I)

https://agenda.infn.it/conferenceDisplay.py?confId=14407





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

G. Cowan

Outline

→ Lecture 1:

Quick review of probability Parameter estimation, maximum likelihood Statistical tests for discovery and limits

Lecture 2:

Limits for the Poisson counting experiment Nuisance parameters and systematic uncertainties Tests from profile likelihood ratio More parameter estimation, Bayesian methods (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)
- C. Patrignani et al. (Particle Data Group), *Review of Particle Physics*, Chin. Phys. C, 40, 100001 (2016); 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)

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

Bayes' theorem



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) prior probability, i.e., before seeing the data $P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) dH}$ posterior probability, i.e., after seeing the data 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!)

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

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.

G. Cowan

Properties of estimators

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

Want small variance and small bias, but in general cannot optimize with respect to both; some trade-off necessary.

G. Cowan

Maximum Likelihood (ML) estimators

The most important frequentist method for constructing estimators is to take the value of the parameter(s) that maximize the likelihood (or equivalently the log-likelihod):



 $\theta = \operatorname{argmax} L(x|\theta)$

In some cases we can find the ML estimator as a closed-form function of the data; more often it is found numerically.

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

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.



G. Cowan

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:

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

G. Cowan

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:

$$I_{ij} = E\left[-\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j}\right] = -n \int f(x; \vec{\theta}) \frac{\partial^2 \ln f(x; \vec{\theta})}{\partial \theta_i \partial \theta_j} dx$$

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

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

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

G. Cowan

Prelude to statistical tests: A simulated SUSY event



G. Cowan

Background events



This event from Standard Model ttbar production also has high $p_{\rm T}$ jets and muons, and some missing transverse energy.

→ can easily mimic a SUSY event.

Frequentist statistical tests

Suppose a measurement produces data x; consider a hypothesis H_0 we want to test and alternative H_1

 H_0, H_1 specify probability for \mathbf{x} : $P(\mathbf{x}|H_0), P(\mathbf{x}|H_1)$

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 α , assuming H_0 is correct, to observe the data there, i.e.,

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

Need inequality if data are discrete.

 α is called the size or significance level of the test.

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



G. Cowan

Definition of a test (2)

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

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:



G. Cowan

Classification viewed as a statistical test

Probability to reject H_0 if true (type I error): $\alpha = \int_W f(\mathbf{x}|H_0) d\mathbf{x}$

 α = size of test, significance level, false discovery rate

Probability to accept H_0 if H_1 true (type II error) $\beta = \int_{\overline{W}} f(\mathbf{x}|H_1) d\mathbf{x}$ $1 - \beta = \text{power of test with respect to } H_1$

Equivalently if e.g. H_0 = background, H_1 = signal, use efficiencies:

$$\varepsilon_{\rm b} = \int_W f(\mathbf{x}|H_0) = \alpha$$

$$\varepsilon_{\mathbf{s}} = \int_{W} f(\mathbf{x}|H_1) = 1 - \beta = \text{power}$$

G. Cowan

Purity / misclassification rate

Consider the probability that an event of signal (s) type classified correctly (i.e., the event selection purity),



Note purity depends on the prior probability for an event to be signal or background as well as on s/b efficiencies.

G. Cowan

Physics context of a statistical test

Event Selection: data = individual event; goal is to classify

Example: separation of different particle types (electron vs muon) or known event types (ttbar vs QCD multijet). E.g. test H_0 : event is background vs. H_1 : event is signal. Use selected events for further study.

Search for New Physics: data = a sample of events. Test null hypothesis

 H_0 : all events correspond to Standard Model (background only), against the alternative

 H_1 : events include a type whose existence is not yet established (signal plus background)

Many subtle issues here, mainly related to the high standard of proof required to establish presence of a new phenomenon. The optimal statistical test for a search is closely related to that used for event selection.

G. Cowan

Statistical tests for event selection

Suppose the result of a measurement for an individual event is a collection of numbers $\vec{x} = (x_1, \dots, x_n)$

 x_1 = number of muons,

 $x_2 = \text{mean } p_T \text{ of jets,}$

 $x_3 = missing energy, ...$

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

 $\mathsf{pp} o t\overline{t} \;, \quad \mathsf{pp} o \widetilde{g}\widetilde{g} \;, \ldots$

For each reaction we consider we will have a hypothesis for the pdf of x, e.g., p(x|b), p(x|s)

E.g. here call H_0 the background hypothesis (the event type we want to reject); H_1 is signal hypothesis (the type we want).

G. Cowan

Selecting events

Suppose we have a data sample with two kinds of events, corresponding to hypotheses H_0 and H_1 and we want to select those of type H_1 .

Each event is a point in \vec{x} space. What 'decision boundary' should we use to accept/reject events as belonging to event types H_0 or H_1 ?

Perhaps select events with 'cuts':

 $\begin{array}{ll} x_i & < c_i \\ x_j & < c_j \end{array}$



Other ways to select events

Or maybe use some other sort of decision boundary:

linear

or nonlinear



How can we do this in an 'optimal' way?

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.



Test statistic based on likelihood ratio

How can we choose a test's critical region in an 'optimal way'?

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

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

inside the region, and $\leq c$ outside, where c is a constant chosen to give a test of the desired size.

Equivalently, optimal scalar test statistic is

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

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

G. Cowan

Neyman-Pearson doesn't usually help

We usually don't have explicit formulae for the pdfs f(x|s), f(x|b), so for a given x we can't evaluate the likelihood ratio

$$t(\mathbf{x}) = \frac{f(\mathbf{x}|s)}{f(\mathbf{x}|b)}$$

Instead we may have Monte Carlo models for signal and background processes, so we can produce simulated data:

generate
$$\mathbf{x} \sim f(\mathbf{x}|\mathbf{s}) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_N$$

generate $\mathbf{x} \sim f(\mathbf{x}|\mathbf{b}) \rightarrow \mathbf{x}_1, \dots, \mathbf{x}_N$

This gives samples of "training data" with events of known type. Can be expensive (1 fully simulated LHC event ~ 1 CPU minute).

G. Cowan

Approximate LR from histograms

Want t(x) = f(x|s)/f(x|b) for x here



One possibility is to generate MC data and construct histograms for both signal and background.

Use (normalized) histogram values to approximate LR:

$$t(x) \approx \frac{N(x|s)}{N(x|b)}$$

Can work well for single variable.

Statistics 1 / JENNIFER, Trieste, 30 Jul - 3 Aug 2018

Approximate LR from 2D-histograms

Suppose problem has 2 variables. Try using 2-D histograms:



Approximate pdfs using N(x,y|s), N(x,y|b) in corresponding cells. But if we want *M* bins for each variable, then in *n*-dimensions we have M^n cells; can't generate enough training data to populate.

 \rightarrow Histogram method usually not usable for n > 1 dimension.

Strategies for multivariate analysis

Neyman-Pearson lemma gives optimal answer, but cannot be used directly, because we usually don't have f(x|s), f(x|b).

Histogram method with M bins for n variables requires that we estimate M^n parameters (the values of the pdfs in each cell), so this is rarely practical.

A compromise solution is to assume a certain functional form for the test statistic t(x) with fewer parameters; determine them (using MC) to give best separation between signal and background.

Alternatively, try to estimate the probability densities f(x|s) and f(x|b) (with something better than histograms) and use the estimated pdfs to construct an approximate likelihood ratio.

Multivariate methods

Many new (and some old) methods esp. from Machine Learning: Fisher discriminant (Deep) neural networks Kernel density methods Support Vector Machines Decision trees Boosting Bagging

This is a large topic -- see e.g. lectures

http://www.pp.rhul.ac.uk/~cowan/stat/stat_2.pdf (from around p 38) and references therein.

G. Cowan

Testing significance / goodness-of-fit 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?

Decide what part of the data space represents less compatibility with H than does the point \vec{x}_{ODS} .

This region therefore has greater compatibility with some alternative *H*'.



p-values

Express 'goodness-of-fit' 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!

In frequentist statistics we don't talk about P(H) (unless H represents a repeatable observation). In Bayesian statistics we do; use Bayes' theorem to obtain

$$P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H) \, dH}$$

where $\pi(H)$ is the prior probability for *H*.

For now stick with the frequentist approach; result is *p*-value, regrettably easy to misinterpret as P(H).

Significance from *p*-value

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



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

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

G. Cowan

Test statistics and *p*-values

Consider a parameter μ proportional to rate of signal process.

Often define a function of the data (test statistic) q_{μ} that reflects level of agreement between the data and the hypothesized value μ .

Usually define q_{μ} so that higher values increasingly incompatibility with the data (more compatible with a relevant alternative).

We can define critical region of test of μ by $q_{\mu} \ge \text{const.}$, or equivalently define the *p*-value of μ as:

 $p_{\mu} = \int_{q_{\mu,obs}}^{\infty} f(q_{\mu}|\mu) \, dq_{\mu}$ observed value of q_{μ} pdf of q_{μ} assuming μ Equivalent formulation of test: reject μ if $p_{\mu} < \alpha$.

G. Cowan

Confidence interval from inversion of a test

Carry out a test of size α for all values of μ .

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

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

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

To find edge of interval (the "limit"), set $p_{\mu} = \alpha$ and solve for μ .

Coming up...

We now have most of the ingredients to carry out:

Parameter estimation

Statistical tests

Setting limits (confidence intervals)

In the second lecture we will apply these to some examples and extend the concepts e.g. to deal with systematic uncertainties.