Statistics for Particle Physics Lecture 2





Taller de Altas EnergíasBenasque, Spain6 September 2021

http://benasque.org/2021tae/



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

Outline

Lecture 1: Introduction, probability, parameter estimation

→ Lecture 2: Hypothesis tests, limits

Lecture 3: Systematic uncertainties, experimental sensitivity

Frequentist hypothesis 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.,

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

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 .



Definition of a test (2)

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

Use the alternative hypothesis H_1 to motivate where to place the critical region.

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



Classification viewed as a statistical test

Suppose events come in two possible types:

s (signal) and b (background)

For each event, test hypothesis that it is background, i.e., $H_0 = b$.

Carry out test on many events, each is either of type s or b, i.e., here the hypothesis is the "true class label", which varies randomly from event to event, so we can assign to it a frequentist probability.

Select events for which where H_0 is rejected as "candidate events of type s". Equivalent Particle Physics terminology:

background efficiency
$$arepsilon_{
m b} = \int_W f(\mathbf{x}|H_0) \, d\mathbf{x} = lpha$$

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

signal efficiency

G. Cowan / RHUL Physics

Example of a test for classification



For each event in a mixture of signal (s) and background (b) test

 H_0 : event is of type b

using a critical region W of the form: $W = \{x : x \le x_c\}$, where x_c is a constant that we choose to give a test with the desired size α .

G. Cowan / RHUL Physics

Classification example (2)

Suppose we want $\alpha = 10^{-4}$. Require:

$$\alpha = P(x \le x_{c}|b) = \int_{0}^{x_{c}} f(x|b) \, dx = \frac{4x^{4}}{4} \Big|_{0}^{x_{c}} = x_{c}^{4}$$

and therefore $x_{\rm c} = \alpha^{1/4} = 0.1$

For this test (i.e. this critical region W), the power with respect to the signal hypothesis (s) is

$$M = P(x \le x_{\rm c}|{\rm s}) = \int_0^{x_{\rm c}} f(x|{\rm s}) \, dx = 2x_{\rm c} - x_{\rm c}^2 = 0.19$$

Note: the optimal size and power is a separate question that will depend on goals of the subsequent analysis.

G. Cowan / RHUL Physics

Classification example (3)

Suppose that the prior probabilities for an event to be of type s or b are:

 $\pi_{\rm s} = 0.001$ $\pi_{\rm b} = 0.999$

The "purity" of the selected signal sample (events where b hypothesis rejected) is found using Bayes' theorem:

$$P(\mathbf{s}|x \le x_{\mathbf{c}}) = \frac{P(x \le x_{\mathbf{c}}|\mathbf{s})\pi_{\mathbf{s}}}{P(x \le x_{\mathbf{c}}|\mathbf{s})\pi_{\mathbf{s}} + P(x \le x_{\mathbf{c}}|\mathbf{b})\pi_{\mathbf{b}}}$$

= 0.655

G. Cowan / RHUL Physics

Test statistic based on likelihood ratio

How can we choose a test's critical region in an 'optimal way', in particular if the data space is multidimensional?

Neyman-Pearson lemma states:

For a test of H_0 of size α , to get the highest power with respect to the alternative H_1 we need for all x in the critical region W

"likelihood
$$\frac{P(\mathbf{x}|H_1)}{P(\mathbf{x}|H_0)} \ge c_{\alpha}$$

inside W and $\leq c_{\alpha}$ outside, where c_{α} is a constant chosen to give a test of the desired size.

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.

G. Cowan / RHUL Physics

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

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

This gives samples of "training data" with events of known type.

Use these to construct a statistic that is as close as possible to the optimal likelihood ratio (\rightarrow Machine Learning).

G. Cowan / RHUL Physics

Testing significance / goodness-of-fit

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

We observe a single point in this space: x_{obs} .

 X_i

How can we quantify the level of compatibility between the data and the predictions of *H*?

Decide what part of the data space represents equal or less compatibility with H than does the point x_{obs} . (Not unique!)



p-values

Express level of compatibility between data and hypothesis (sometimes 'goodness-of-fit') by giving the *p*-value for *H*:

 $p = P(\mathbf{x} \in \omega_{\leq}(\mathbf{x}_{obs})|H)$

- probability, under assumption of H, to observe data
 with equal or lesser compatibility with H relative to the
 data we got.
- probability, under assumption of H, to observe data as discrepant with H as the data we got or more so.

Basic idea: if there is only a very small probability to find data with even worse (or equal) compatibility, then *H* is "disfavoured by the data".

If the *p*-value is below a user-defined threshold α (e.g. 0.05) then *H* is rejected (equivalent to hypothesis test as discussed previously).



The *p*-value of H is not the probability that *H* is true!

In frequentist statistics we don't talk about P(H) (unless H represents a repeatable observation).

If we do define P(H), e.g., in Bayesian statistics as a degree of belief, then we need to 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).

p-value from test statistic



If e.g. we define the region of less or eq. compatibility to be $t(x) \ge t_{obs}$ then the *p*-value of *H* is

$$p_H = \int_{t_{\text{obs}}}^{\infty} f(t|H) \, dt = \int_{\{\mathbf{x}: t(\mathbf{x}) \ge t_{\text{obs}}\}} f(\mathbf{x}|H) \, d\mathbf{x}$$

G. Cowan / RHUL Physics

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

One can show the distribution of a *p*-value p_0 of H_0 under assumption of H_0 is uniform between 0 and 1.

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

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



Therefore we can define the critical region of a test of H_0 with size α 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)

In both cases need to ask what is relevant alternative hypothesis.

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 / RHUL Physics

TAE 2021 / Statistics for PP Lecture 2

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.



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

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

in ROOT: p = 1 - TMath::Freq(Z) Z = TMath::NormQuantile(1-p)

in python (scipy.stats):
p = 1 - norm.cdf(Z) = norm.sf(Z)
Z = norm.ppf(1-p)

Result Z is a "number of sigmas". Note this does not mean that the original data was Gaussian distributed.

G. Cowan / RHUL Physics

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

In addition to a 'point estimate' of a parameter we should report an interval reflecting its statistical uncertainty.

Confidence intervals for a parameter θ can be found by defining a test of the hypothesized value θ (do this for all θ):

Specify values of the data that are 'disfavoured' by θ (critical region) such that $P(\text{data in critical region} | \theta) \le \alpha$ for a prespecified α , e.g., 0.05 or 0.1.

If data observed in the critical region, reject the value θ .

Now invert the test to define a confidence interval as:

set of θ values that are not rejected in a test of size α (confidence level CL is $1 - \alpha$).

Relation between confidence interval and *p*-value

Equivalently we can consider a significance test for each hypothesized value of θ , resulting in a *p*-value, p_{θ} .

If $p_{\theta} \leq \alpha$, then we reject θ .

The confidence interval at $CL = 1 - \alpha$ consists of those values of θ that are not rejected.

E.g. an upper limit on θ is the greatest value for which $p_{\theta} > \alpha$.

In practice find by setting $p_{\theta} = \alpha$ and solve for θ .

For a multidimensional parameter space $\theta = (\theta_1, \dots, \theta_M)$ use same idea – result is a confidence "region" with boundary determined by $p_{\theta} = \alpha$.

Coverage probability of confidence interval

If the true value of θ is rejected, then it's not in the confidence interval. The probability for this is by construction (equality for continuous data):

 $P(\text{reject } \theta | \theta) \leq \alpha = \text{type-I error rate}$

Therefore, the probability for the interval to contain or "cover" θ is

P(conf. interval "covers" $\theta | \theta \ge 1 - \alpha$

This assumes that the set of θ values considered includes the true value, i.e., it assumes the composite hypothesis $P(\mathbf{x}|H,\theta)$.

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 $\text{CL} = 1 - \alpha$ found from

$$\alpha = P(n \le n_{\text{obs}}; s_{\text{up}}, b) = \sum_{n=0}^{n_{\text{obs}}} \frac{(s_{\text{up}} + b)^n}{n!} e^{-(s_{\text{up}} + b)}$$
$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} (1 - \alpha; 2(n_{\text{obs}} + 1)) - b$$
$$= \frac{1}{2} F_{\chi^2}^{-1} (0.95; 2(5 + 1)) - 4.5 = 6.0$$

G. Cowan / RHUL Physics

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

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



Limits near a boundary of the parameter space

Suppose e.g. b = 2.5 and we observe n = 0.

If we choose CL = 0.9, we find from the formula for s_{up}

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

Physicist:

We already knew $s \ge 0$ before we started; can't use negative upper limit to report result of expensive experiment!

Statistician:

The interval is designed to cover the true value only 90% of the time — this was clearly not one of those times.

Not uncommon dilemma when testing parameter values for which one has very little experimental sensitivity, e.g., very small *s*.

Expected limit for s = 0

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

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

Reality check: with b = 2.5, typical Poisson fluctuation in n is at least $\sqrt{2.5} = 1.6$. How can the limit be so low?



Next time...

Lecture 1: Introduction, probability, parameter estimation Lecture 2: Hypothesis tests, limits

→ Lecture 3: Systematic uncertainties, experimental sensitivity



Classification example (4)

Suppose an individual event is observed at x = 0.1. What is the probability that this event is background?

$$P(\mathbf{b}|x) = \frac{f(x|\mathbf{b})\pi_{\mathbf{b}}}{f(x|\mathbf{b})\pi_{\mathbf{b}} + f(x|\mathbf{s})\pi_{\mathbf{s}}}$$
$$= \frac{4x^{3}\pi_{\mathbf{b}}}{4x^{3}\pi_{\mathbf{b}} + 2(1-x)\pi_{\mathbf{s}}}$$

= 0.689

(Here nothing to do with the test using $x \le x_c$, just an illustration of Bayes' theorem.)

G. Cowan / RHUL Physics

Compatibility with ${\cal H}$

What does it mean for a region of data space to be less compatible with the predictions of *H*?

It must mean that that region of data space is more compatible with some relevant alternative H'.

So although the definition of the *p*-value does not refer explicitly to an alternative, this enters implicitly through its role in determining the partitioning of the data space into more and less-or-equally compatible regions.

As in the case of hypothesis tests, there may be more than one relevant alternative.

Obvious where to put *W*?

In the 1930s there were great debates as to the role of the alternative hypothesis.

Fisher held that one could test a hypothesis H_0 without reference to an alternative.

Suppose, e.g., H_0 predicts that x (suppose positive) usually comes out low. High values of x are less characteristic of H_0 , so if a high value is observed, we should reject H_0 , i.e., we put W at high x:



G. Cowan / RHUL Physics

Or not so obvious where to put W?

But what if the only relevant alternative to H_0 is H_1 as below:



Here high x is more characteristic of H_0 and not like what we expect from H_1 . So better to put W at low x.

Neyman and Pearson argued that "less characteristic of H_0 " is well defined only when taken to mean "more characteristic of some relevant alternative H_1 ".

G. Cowan / RHUL Physics

Example of *p*-value: exponential decay time

A nuclear sample contains two radioactive isotopes with mean lifetimes $\tau = 0.2$ s and $\tau = 1.0$ s.

For either isotope we expect the decay time to follow $f(t|\tau) = \frac{1}{\tau}e^{-t/\tau}$

A nucleus is observed to decay after a time t_{obs} = 0.6 s.

The *p*-value of the hypothesis *H* that the nucleus is of the type with $\tau = 0.2$ s is

$$p_H = P(t \ge t_{\rm obs} | \tau = 0.2 \,\mathrm{s}) = 0.0498$$

Here we take $t \ge t_{obs}$ as being less compatible with $\tau = 0.2$ s , because greater t is more characteristic of $\tau = 1.0$ s.

If the relevant alternative had been $\tau = 0.1$ s, then one would define the *p*-value as

$$p_H = P(t \le t_{\rm obs} | \tau = 0.2 \,\mathrm{s}) = 0.9502$$



Distribution of the *p*-value

The *p*-value is a function of the data, and is thus itself a random variable with a given distribution. Suppose the *p*-value of *H* is found from a test statistic t(x) as

$$p_H = \int_t^\infty f(t'|H)dt'$$

The pdf of p_H under assumption of H is

$$g(p_H|H) = \frac{f(t|H)}{|\partial p_H/\partial t|} = \frac{f(t|H)}{f(t|H)} = 1 \quad (0 \le p_H \le 1)$$

In general for continuous data, under assumption of H, $p_H \sim \text{Uniform}[0,1]$ and is concentrated toward zero for some (broad) class of alternatives.



G. Cowan / RHUL Physics

The Bayesian approach to limits

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

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

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

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

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

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

Bayesian prior for Poisson parameter

Include knowledge that $s \ge 0$ by setting prior $\pi(s) = 0$ for s < 0. Could try to reflect 'prior ignorance' with e.g.

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

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

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

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

Bayesian upper limit with flat prior for s

Put Poisson likelihood and flat prior into Bayes' theorem:

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

Normalize to unit area:

$$p(s|n) = \frac{(s+b)^n e^{-(s+b)}}{\Gamma(b, n+1)} \longleftarrow \text{ upper incomplete gamma function}$$

Upper limit s_{up} determined by requiring

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

G. Cowan / RHUL Physics

Bayesian interval with flat prior for *s*

Solve to find limit s_{up} :

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

where

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

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

Bayesian interval with flat prior for s

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

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



G. Cowan / RHUL Physics

Low sensitivity to μ

It can be that the effect of a given hypothesized μ is very small relative to the background-only (μ = 0) prediction.

This means that the distributions $f(q_{\mu}|\mu)$ and $f(q_{\mu}|0)$ will be almost the same:



G. Cowan / RHUL Physics

Having sufficient sensitivity

In contrast, having sensitivity to μ means that the distributions $f(q_{\mu}|\mu)$ and $f(q_{\mu}|0)$ are more separated:



That is, the power (probability to reject μ if $\mu = 0$) is substantially higher than α . Use this power as a measure of the sensitivity.

G. Cowan / RHUL Physics

Spurious exclusion

Consider again the case of low sensitivity. By construction the probability to reject μ if μ is true is α (e.g., 5%).

And the probability to reject μ if $\mu = 0$ (the power) is only slightly greater than α .

critical region

This means that with probability of around $\alpha = 5\%$ (slightly higher), one excludes hypotheses to which one has essentially no sensitivity (e.g., $m_{\rm H} = 1000$ TeV).

"Spurious exclusion"

Ways of addressing spurious exclusion

The problem of excluding parameter values to which one has no sensitivity known for a long time; see e.g.,

Virgil L. Highland, *Estimation of Upper Limits from Experimental Data*, July 1986, Revised February 1987, Temple University Report C00-3539-38.

In the 1990s this was re-examined for the LEP Higgs search by Alex Read and others

T. Junk, Nucl. Instrum. Methods Phys. Res., Sec. A 434, 435 (1999); A.L. Read, J. Phys. G 28, 2693 (2002).

and led to the "CL_s" procedure for upper limits.

Unified intervals also effectively reduce spurious exclusion by the particular choice of critical region.

The CL_s procedure

In the usual formulation of CL_s , one tests both the $\mu = 0$ (*b*) and $\mu > 0$ ($\mu s+b$) hypotheses with the same statistic $Q = -2\ln L_{s+b}/L_b$:



G. Cowan / RHUL Physics

The CL_s procedure (2)

As before, "low sensitivity" means the distributions of Q under b and s+b are very close:



G. Cowan / RHUL Physics

The CL_s procedure (3)

The CL_s solution (A. Read et al.) is to base the test not on the usual *p*-value (CL_{s+b}), but rather to divide this by CL_b (~ one minus the *p*-value of the *b*-only hypothesis), i.e.,



Choice of test for limits (2)

In some cases $\mu = 0$ is no longer a relevant alternative and we want to try to exclude μ on the grounds that some other measure of incompatibility between it and the data exceeds some threshold.

If the measure of incompatibility is taken to be the likelihood ratio with respect to a two-sided alternative, then the critical region can contain both high and low data values.

→ unified intervals, G. Feldman, R. Cousins, Phys. Rev. D 57, 3873–3889 (1998)

The Big Debate is whether to use one-sided or unified intervals in cases where small (or zero) values of the parameter are relevant alternatives. Professional statisticians have voiced support on both sides of the debate.

Unified (Feldman-Cousins) intervals

We can use directly

$$t_{\mu} = -2\ln\lambda(\mu)$$
 where

$$\lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})}$$

as a test statistic for a hypothesized μ .

Large discrepancy between data and hypothesis can correspond either to the estimate for μ being observed high or low relative to μ .

This is essentially the statistic used for Feldman-Cousins intervals (here also treats nuisance parameters).

G. Feldman and R.D. Cousins, Phys. Rev. D 57 (1998) 3873.

Lower edge of interval can be at μ = 0, depending on data.

Upper/lower edges of F-C interval for μ versus bfor $n \sim \text{Poisson}(\mu+b)$



Lower edge may be at zero, depending on data.

For n = 0, upper edge has (weak) dependence on b.

G. Cowan / RHUL Physics