Statistics for Particle Physics Lecture 2



Helmholtz Alliance



Introduction to the Terascale https://indico.desy.de/event/46666/

DESY, Hamburg 19 March 2025



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

Outline

Wednesday 9:00

Quick review of probability Hypothesis testing

 \rightarrow Wednesday 9:45

p-values
Confidence intervals / limits

More resources in the University of London course:

https://www.pp.rhul.ac.uk/~cowan/stat_course.html

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} .

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 of size α as seen earlier).



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). 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.

G. Cowan / RHUL Physics

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, suppose relevant alt. is 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

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: i p = 1 - TMath::Freq(Z) p Z = TMath::NormQuantile(1-p) 2

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" $\theta | \theta \ge 1 \square \alpha$

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

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

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

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

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

G. Cowan / RHUL Physics

n ~ 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}

 $s_{\rm up} = -0.197$ (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?



Approximate confidence intervals/regions from the likelihood function

Suppose we test parameter value(s) $\theta = (\theta_1, ..., \theta_n)$ using the ratio

$$\lambda(\theta) = \frac{L(\theta)}{L(\hat{\theta})} \qquad \qquad 0 \le \lambda(\theta) \le 1$$

Lower $\lambda(\theta)$ means worse agreement between data and hypothesized θ . Equivalently, usually define

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

so higher t_{θ} means worse agreement between θ and the data.

p-value of θ therefore $p_{\theta} = \int_{t_{\theta,\text{obs}}}^{\infty} f(t_{\theta}|\theta) dt_{\theta}$ need p

G. Cowan / RHUL Physics

Confidence region from Wilks' theorem

Wilks' theorem says (in large-sample limit and provided certain conditions hold...)

 $f(t_{\theta}|\theta) \sim \chi_n^2 \qquad \begin{array}{l} \text{chi-square dist. with $\#$ d.o.f. =} \\ \# \text{ of components in $\theta = (\theta_1, ..., \theta_n)$.} \end{array}$

Assuming this holds, the *p*-value is

$$p_{m{ heta}} = 1 - F_{\chi^2_n}(t_{m{ heta}}) \quad \leftarrow \text{set equal to } lpha$$

To find boundary of confidence region set $p_{\theta} = \alpha$ and solve for t_{θ} :

$$t_{\theta} = F_{\chi_n^2}^{-1}(1-\alpha)$$

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

Recall also

G. Cowan / RHUL Physics

Confidence region from Wilks' theorem (cont.) i.e., boundary of confidence region in θ space is where

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}F_{\chi_n^2}^{-1}(1-\alpha)$$

For example, for $1 - \alpha = 68.3\%$ and n = 1 parameter,

$$F_{\chi_1^2}^{-1}(0.683) = 1$$

and so the 68.3% confidence level interval is determined by

$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}$$

Same as recipe for finding the estimator's standard deviation, i.e.,

 $[\hat{\theta} - \sigma_{\hat{\theta}}, \hat{\theta} + \sigma_{\hat{\theta}}]$ is a 68.3% CL confidence interval.

Example of interval from $\ln L(\theta)$

For n=1 parameter, CL = 0.683, $Q_{\alpha} = 1$.



Multiparameter case

For increasing number of parameters, $CL = 1 - \alpha$ decreases for confidence region determined by a given

$$Q_{\alpha} = F_{\chi_n^2}^{-1}(1-\alpha)$$

Q_{lpha}	1-lpha					
	n = 1	n = 2	n = 3	n = 4	n = 5	
1.0	0.683	0.393	0.199	0.090	0.037	
2.0	0.843	0.632	0.428	0.264	0.151	
4.0	0.954	0.865	0.739	0.594	0.451	
9.0	0.997	0.989	0.971	0.939	0.891	

Multiparameter case (cont.)

Equivalently, Q_{α} increases with *n* for a given $CL = 1 - \alpha$.

$1 - \alpha$	$ar{Q}_{lpha}$						
	n = 1	n = 2	n = 3	n = 4	n = 5		
0.683	1.00	2.30	3.53	4.72	5.89		
0.90	2.71	4.61	6.25	7.78	9.24		
0.95	3.84	5.99	7.82	9.49	11.1		
0.99	6.63	9.21	11.3	13.3	15.1		

Summary

Here only time to talk about a limited number of topics:

Hypothesis tests p-values Confidence intervals, regions, limits...

Going further:

Machine Learning Bayesian methods Systematic uncertainties

For some exercises and software, see:

https://www.pp.rhul.ac.uk/~cowan/stat/exercises/cowan_stat_exercises.pdf

Extra slides

MLE 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):

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

MLE example: parameter of exponential pdf (2)

Find its maximum by setting

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

$$\frac{\partial \ln L(\tau)}{\partial \tau} = 0 \; ,$$

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

We find the ML estimate:

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



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

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

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:

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

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,

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

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

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

G. Cowan / RHUL Physics

MVB for MLE of exponential parameter

Find MVB =
$$-\left(1 + \frac{\partial b}{\partial \tau}\right)^2 / E\left[\frac{\partial^2 \ln L}{\partial \tau^2}\right]$$

We found for the exponential parameter the MLE

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

and we showed b = 0, hence $\partial b / \partial \tau = 0$.

We find
$$\frac{\partial^2 \ln L}{\partial \tau^2} = \sum_{i=1}^n \left(\frac{1}{\tau^2} - \frac{2t_i}{\tau^3}\right)$$

and since $E[t_i] = \tau$ for all i , $E\left[\frac{\partial^2 \ln L}{\partial \tau^2}\right] = -\frac{n}{\tau^2}$,
and therefore $MVB = \frac{\tau^2}{n} = V[\hat{\tau}]$. (Here MLE is "efficient")

G. Cowan / RHUL Physics

•

Variance of estimators: graphical method

Expand $lnL(\theta)$ 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):

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

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.

Example of variance by graphical method



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