Basics of Statistical Data Analysis



Quantum Universe Lecture DESY / U. Hamburg Zoom / 13 April 2021

https://indico.desy.de/event/29561/



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

Outline

- Basic introduction to frequentist vs Bayesian approaches
- Basics of fitting, hypothesis tests and asymptotics
- Treating systematic uncertainties with nuisance parameters

Everything here is a subset of the University of London course: http://www.pp.rhul.ac.uk/~cowan/stat_course.html

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.
- Luca Lista, *Statistical Methods for Data Analysis in Particle Physics*, Springer, 2017.
- 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.
- P.A. Zyla et al. (Particle Data Group), Prog. Theor. Exp. Phys. 2020, 083C01 (2020); pdg.1b1.gov sections on probability, statistics, MC.

Theory ↔ Statistics ↔ Experiment

Theory (model, hypothesis):

Experiment (observation):

$$F = -G \frac{m_1 m_2}{r^2}, \dots$$

+ response of measurement apparatus

= model prediction



Quick review of probablility

Frequentist (*A* = outcome of repeatable observation):

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

Subjective (*A* = hypothesis):

$$P(A) =$$
degree of belief that A is true

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

Bayes' theorem:
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)}$$

Frequentist Statistics – general philosophy

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

Probability = limiting frequency

Probabilities such as

P (string theory is true), P (0.117 < α_s < 0.119), P (Biden wins in 2024),

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.

Preferred theories (models, hypotheses, ...) are those that predict a high probability for data "like" the data observed.

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

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!)

Parameter estimation

The parameters of a pdf are any constants that characterize it,



i.e., θ indexes a set of hypotheses.

Suppose we have a sample of observed values: $x = (x_1, ..., x_n)$

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.

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

 \rightarrow average of repeated measurements should tend to true value.

And we want a small variance (statistical error): $V[\hat{\theta}]$

 \rightarrow small bias & variance are in general conflicting criteria

Hypothesis, likelihood

Suppose the entire result of an experiment (set of measurements) is a collection of numbers *x*.

A (simple) hypothesis is a rule that assigns a probability to each possible data outcome:

 $P(\mathbf{x}|H)$ = the likelihood of H

Often we deal with a family of hypotheses labeled by one or More undetermined parameters (a composite hypothesis):

$$P(\mathbf{x}|oldsymbol{ heta}) = L(oldsymbol{ heta})$$
 = the "likelihood function"

Note:

- 1) For the likelihood we treat the data x as fixed.
- 2) The likelihood function $L(\theta)$ is not a pdf for θ .

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

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

Maximum Likelihood Estimators (MLEs)

We *define* the maximum likelihood estimators or MLEs to be the parameter values for which the likelihood is maximum.



Could have multiple maxima (take highest).

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

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

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

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

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

We find the ML estimate:

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



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 .



G. Cowan / RHUL Physics

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
$$\varepsilon_{\mathbf{b}} = \int_{W} f(\mathbf{x}|H_0) \, d\mathbf{x} = \alpha$$

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

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

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

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

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). 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

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

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}

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

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} :

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

Recall also

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

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}	$1-\alpha$					
	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 $CL = 1 - \alpha$.

$1 - \alpha$	\widehat{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		

Systematic uncertainties and nuisance parameters In general, our model of the data is not perfect:



Can improve model by including additional adjustable parameters.

$$P(x|\mu) \to P(x|\mu, \theta)$$

Nuisance parameter ↔ systematic uncertainty. Some point in the parameter space of the enlarged model should be "true".

Presence of nuisance parameter decreases sensitivity of analysis to the parameter of interest (e.g., increases variance of estimate).

Prototype search analysis

Search for signal in a region of phase space; result is histogram of some variable x giving numbers:

$$\mathbf{n}=(n_1,\ldots,n_N)$$

Assume the n_i are Poisson distributed with expectation values

$$E[n_i] = \mu s_i + b_i$$
strength parameter

where

G. Cowan / RHUL Physics

Prototype analysis (II)

Often also have a subsidiary measurement that constrains some of the background and/or shape parameters:

$$\mathbf{m} = (m_1, \ldots, m_M)$$

Assume the m_i are Poisson distributed with expectation values

$$E[m_i] = u_i(\boldsymbol{\theta})$$
nuisance parameters ($\boldsymbol{\theta}_{s}, \boldsymbol{\theta}_{b}, b_{tot}$)

Likelihood function is

$$L(\mu, \theta) = \prod_{j=1}^{N} \frac{(\mu s_j + b_j)^{n_j}}{n_j!} e^{-(\mu s_j + b_j)} \quad \prod_{k=1}^{M} \frac{u_k^{m_k}}{m_k!} e^{-u_k}$$

G. Cowan / RHUL Physics

The profile likelihood ratio

Base significance test on the profile likelihood ratio:



Define critical region of test of μ by the region of data space that gives the lowest values of $\lambda(\mu)$.

Important advantage of profile LR is that its distribution becomes independent of nuisance parameters in large sample limit.

Test statistic for discovery

Suppose relevant alternative to background-only ($\mu = 0$) is $\mu \ge 0$. So take critical region for test of $\mu = 0$ corresponding to high q_0

and $\hat{\mu} > 0$ (data characteristic for $\mu \ge 0$).

That is, to test background-only hypothesis define statistic

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{\mu} \ge 0\\ 0 & \hat{\mu} < 0 \end{cases}$$

i.e. here only large (positive) observed signal strength is evidence against the background-only hypothesis.

Note that even though here physically $\mu \ge 0$, we allow $\hat{\mu}$ to be negative. In large sample limit its distribution becomes Gaussian, and this will allow us to write down simple expressions for distributions of our test statistics.

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Distribution of q_0 in large-sample limit

Assuming approximations valid in the large sample (asymptotic) limit, we can write down the full distribution of q_0 as

$$f(q_0|\mu') = \left(1 - \Phi\left(\frac{\mu'}{\sigma}\right)\right)\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}\exp\left[-\frac{1}{2}\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)^2\right]$$

The special case $\mu' = 0$ is a "half chi-square" distribution:

$$f(q_0|0) = \frac{1}{2}\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}e^{-q_0/2}$$

In large sample limit, $f(q_0|0)$ independent of nuisance parameters; $f(q_0|\mu')$ depends on nuisance parameters through σ .

p-value for discovery

Large q_0 means increasing incompatibility between the data and hypothesis, therefore *p*-value for an observed $q_{0,obs}$ is



use e.g. asymptotic formula



From *p*-value get equivalent significance,

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

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Cumulative distribution of q_0 , significance

From the pdf, the cumulative distribution of q_0 is found to be

$$F(q_0|\mu') = \Phi\left(\sqrt{q_0} - \frac{\mu'}{\sigma}\right)$$

The special case $\mu' = 0$ is

$$F(q_0|0) = \Phi\left(\sqrt{q_0}\right)$$

The *p*-value of the $\mu = 0$ hypothesis is

$$p_0 = 1 - F(q_0|0)$$

Therefore the discovery significance Z is simply

$$Z = \Phi^{-1}(1 - p_0) = \sqrt{q_0}$$

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Monte Carlo test of asymptotic formula

- $n \sim \text{Poisson}(\mu s + b)$
- $m \sim \text{Poisson}(\tau b)$
- $\mu =$ param. of interest
- *b* = nuisance parameter

Here take *s* known, $\tau = 1$.

Asymptotic formula is good approximation to 5σ level ($q_0 = 25$) already for $b \sim 20$.



How to read the p_0 plot

The "local" p_0 means the *p*-value of the background-only hypothesis obtained from the test of $\mu = 0$ at each individual $m_{\rm H}$, without any correct for the Look-Elsewhere Effect.

The "Expected" (dashed) curve gives the median p_0 under assumption of the SM Higgs (μ = 1) at each $m_{\rm H}$.



The blue band gives the width of the distribution $(\pm 1\sigma)$ of significances under assumption of the SM Higgs.

Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554

Test statistic for upper limits

For purposes of setting an upper limit on μ use

$$q_{\mu} = \begin{cases} -2\ln\lambda(\mu) & \hat{\mu} \leq \mu \\ 0 & \hat{\mu} > \mu \end{cases} \quad \text{where} \quad \lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})}$$

I.e. when setting an upper limit, an upwards fluctuation of the data is not taken to mean incompatibility with the hypothesized μ :

From observed
$$q_{\mu}$$
 find p -value: $p_{\mu} = \int_{q_{\mu, obs}}^{\infty} f(q_{\mu}|\mu) \, dq_{\mu}$

Large sample approximation:

$$p_{\mu} = 1 - \Phi\left(\sqrt{q_{\mu}}\right)$$

To find upper limit at CL = $1-\alpha$, set $p_{\mu} = \alpha$ and solve for μ .

G. Cowan / RHUL Physics

Monte Carlo test of asymptotic formulae

Consider again $n \sim \text{Poisson}(\mu s + b)$, $m \sim \text{Poisson}(\tau b)$ Use q_{μ} to find *p*-value of hypothesized μ values.

E.g. $f(q_1|1)$ for *p*-value of $\mu = 1$. Typically interested in 95% CL, i.e., *p*-value threshold = 0.05, i.e., $q_1 = 2.69$ or $Z_1 = \sqrt{q_1} = 1.64$. Median[$q_1|0$] gives "exclusion

sensitivity".

Here asymptotic formulae good for s = 6, b = 9.



How to read the green and yellow limit plots For every value of $m_{\rm H}$, find the upper limit on μ .

Also for each $m_{\rm H}$, determine the distribution of upper limits $\mu_{\rm up}$ one would obtain under the hypothesis of μ = 0.

The dashed curve is the median μ_{up} , and the green (yellow) bands give the $\pm 1\sigma$ (2σ) regions of this distribution.



53

Finally

One lecture only enough for a brief introduction to:

Parameter estimation, maximum likelihood
Hypothesis tests, *p*-values
Limits (confidence intervals/regions)
Systematics (nuisance parameters)
Asymptotics (Wilks' theorem)

Final thought: once the basic formalism is fixed, most of the work focuses on writing down the likelihood, e.g., $P(x|\theta)$, and including in it enough parameters to adequately describe the data (true for both Bayesian and frequentist approaches).



Expected discovery significance for counting experiment with background uncertainty

I. Discovery sensitivity for counting experiment with *b* known:

(a)
$$\frac{s}{\sqrt{b}}$$

(b) Profile likelihood ratio test & Asimov:

$$\sqrt{2\left((s+b)\ln\left(1+\frac{s}{b}\right)-s\right)}$$

II. Discovery sensitivity with uncertainty in b, σ_b :

(a)
$$\frac{s}{\sqrt{b+\sigma_b^2}}$$

(b) Profile likelihood ratio test & Asimov:

$$\left[2\left((s+b)\ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2}\ln\left[1 + \frac{\sigma_b^2s}{b(b+\sigma_b^2)}\right]\right)\right]^{1/2}$$

G. Cowan / RHUL Physics

Counting experiment with known background

Count a number of events $n \sim \text{Poisson}(s+b)$, where

- s = expected number of events from signal,
- b = expected number of background events.

To test for discovery of signal compute p-value of s = 0 hypothesis,

$$p = P(n \ge n_{\text{obs}}|b) = \sum_{n=n_{\text{obs}}}^{\infty} \frac{b^n}{n!} e^{-b} = 1 - F_{\chi^2}(2b; 2n_{\text{obs}})$$

Usually convert to equivalent significance: $Z = \Phi^{-1}(1-p)$ where Φ is the standard Gaussian cumulative distribution, e.g., Z > 5 (a 5 sigma effect) means $p < 2.9 \times 10^{-7}$.

To characterize sensitivity to discovery, give expected (mean or median) Z under assumption of a given s.

G. Cowan / RHUL Physics

 s/\sqrt{b} for expected discovery significance For large s + b, $n \to x \sim \text{Gaussian}(\mu, \sigma)$, $\mu = s + b$, $\sigma = \sqrt{(s + b)}$. For observed value x_{obs} , p-value of s = 0 is $\text{Prob}(x > x_{\text{obs}} | s = 0)$,:

$$p_0 = 1 - \Phi\left(\frac{x_{\rm obs} - b}{\sqrt{b}}\right)$$

Significance for rejecting s = 0 is therefore

$$Z_0 = \Phi^{-1}(1 - p_0) = \frac{x_{\text{obs}} - b}{\sqrt{b}}$$

Expected (median) significance assuming signal rate s is

$$\mathrm{median}[Z_0|s+b] = \frac{s}{\sqrt{b}}$$

G. Cowan / RHUL Physics

Better approximation for significance

Poisson likelihood for parameter s is

 $L(s) = \frac{(s+b)^n}{n!} e^{-(s+b)}$ For now no nuisance params.

To test for discovery use profile likelihood ratio:

$$q_0 = \begin{cases} -2\ln\lambda(0) & \hat{s} \ge 0 ,\\ 0 & \hat{s} < 0 . \end{cases} \qquad \lambda(s) = \frac{L(s, \hat{\hat{\theta}}(s))}{L(\hat{s}, \hat{\theta})}$$

So the likelihood ratio statistic for testing s = 0 is

$$q_0 = -2\ln\frac{L(0)}{L(\hat{s})} = 2\left(n\ln\frac{n}{b} + b - n\right) \quad \text{for } n > b, \ 0 \text{ otherwise}$$

G. Cowan / RHUL Physics

Approximate Poisson significance (continued)

For sufficiently large s + b, (use Wilks' theorem),

$$Z = \sqrt{2\left(n\ln\frac{n}{b} + b - n\right)} \quad \text{for } n > b \text{ and } Z = 0 \text{ otherwise.}$$

To find median[Z|s], let $n \rightarrow s + b$ (i.e., the Asimov data set):

$$Z_{\rm A} = \sqrt{2\left(\left(s+b\right)\ln\left(1+\frac{s}{b}\right) - s\right)}$$

This reduces to s/\sqrt{b} for $s \ll b$.

G. Cowan / RHUL Physics

 $n \sim \text{Poisson}(s+b)$, median significance, assuming *s*, of the hypothesis s = 0

CCGV, EPJC 71 (2011) 1554, arXiv:1007.1727



"Exact" values from MC, jumps due to discrete data.

Asimov $\sqrt{q_{0,A}}$ good approx. for broad range of *s*, *b*.

 s/\sqrt{b} only good for $s \ll b$.



Extending s/\sqrt{b} to case where b uncertain

The intuitive explanation of s/\sqrt{b} is that it compares the signal, s, to the standard deviation of n assuming no signal, \sqrt{b} .

Now suppose the value of b is uncertain, characterized by a standard deviation σ_b .

A reasonable guess is to replace \sqrt{b} by the quadratic sum of \sqrt{b} and σ_b , i.e.,

$$\operatorname{med}[Z|s] = \frac{s}{\sqrt{b + \sigma_b^2}}$$

This has been used to optimize some analyses e.g. where σ_b cannot be neglected.

Profile likelihood with *b* uncertain

This is the well studied "on/off" problem: Cranmer 2005; Cousins, Linnemann, and Tucker 2008; Li and Ma 1983,...

Measure two Poisson distributed values:

 $n \sim \text{Poisson}(s+b)$ (primary or "search" measurement) $m \sim \text{Poisson}(\tau b)$ (control measurement, τ known) The likelihood function is

$$L(s,b) = \frac{(s+b)^n}{n!} e^{-(s+b)} \frac{(\tau b)^m}{m!} e^{-\tau b}$$

Use this to construct profile likelihood ratio (*b* is nuisance parameter): $I(0, \hat{k}(0))$

$$\lambda(0) = \frac{L(0, b(0))}{L(\hat{s}, \hat{b})}$$

G. Cowan / RHUL Physics

Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

$$\begin{aligned} \hat{s} &= n - m/\tau ,\\ \hat{b} &= m/\tau ,\\ \hat{b}(s) &= \frac{n + m - (1 + \tau)s + \sqrt{(n + m - (1 + \tau)s)^2 + 4(1 + \tau)sm}}{2(1 + \tau)} \end{aligned}$$

and in particular to test for discovery (s = 0),

$$\hat{\hat{b}}(0) = \frac{n+m}{1+\tau}$$

G. Cowan / RHUL Physics

Asymptotic significance

Use profile likelihood ratio for q_0 , and then from this get discovery significance using asymptotic approximation (Wilks' theorem):

$$Z = \sqrt{q_0}$$
$$= \left[-2\left(n\ln\left[\frac{n+m}{(1+\tau)n}\right] + m\ln\left[\frac{\tau(n+m)}{(1+\tau)m}\right]\right) \right]^{1/2}$$

for $n > \hat{b}$ and Z = 0 otherwise.

Essentially same as in:

Robert D. Cousins, James T. Linnemann and Jordan Tucker, NIM A 595 (2008) 480– 501; arXiv:physics/0702156.

Tipei Li and Yuqian Ma, Astrophysical Journal 272 (1983) 317–324.

G. Cowan / RHUL Physics

Asimov approximation for median significance

To get median discovery significance, replace *n*, *m* by their expectation values assuming background-plus-signal model:

$$\begin{split} n &\to s + b \\ m &\to \tau b \end{split}$$

$$Z_{\rm A} = \left[-2\left((s+b) \ln\left[\frac{s+(1+\tau)b}{(1+\tau)(s+b)}\right] + \tau b \ln\left[1+\frac{s}{(1+\tau)b}\right] \right) \right]^{1/2} \end{aligned}$$
Or use the variance of $\hat{b} = m/\tau$, $V[\hat{b}] \equiv \sigma_b^2 = \frac{b}{\tau}$, to eliminate τ :
$$Z_{\rm A} = \left[2\left((s+b) \ln\left[\frac{(s+b)(b+\sigma_b^2)}{b^2+(s+b)\sigma_b^2}\right] - \frac{b^2}{\sigma_b^2} \ln\left[1+\frac{\sigma_b^2 s}{b(b+\sigma_b^2)}\right] \right) \right]^{1/2}$$

G. Cowan / RHUL Physics

Limiting cases

Expanding the Asimov formula in powers of s/b and σ_b^2/b (= $1/\tau$) gives

$$Z_{\rm A} = \frac{s}{\sqrt{b + \sigma_b^2}} \left(1 + \mathcal{O}(s/b) + \mathcal{O}(\sigma_b^2/b) \right)$$

So the "intuitive" formula can be justified as a limiting case of the significance from the profile likelihood ratio test evaluated with the Asimov data set. Testing the formulae: s = 5



G. Cowan / RHUL Physics

Using sensitivity to optimize a cut



Figure 1: (a) The expected significance as a function of the cut value x_{cut} ; (b) the distributions of signal and background with the optimal cut value indicated.

G. Cowan / RHUL Physics

Summary on discovery sensitivity

Simple formula for expected discovery significance based on profile likelihood ratio test and Asimov approximation:

$$Z_{\rm A} = \left[2 \left((s+b) \ln \left[\frac{(s+b)(b+\sigma_b^2)}{b^2 + (s+b)\sigma_b^2} \right] - \frac{b^2}{\sigma_b^2} \ln \left[1 + \frac{\sigma_b^2 s}{b(b+\sigma_b^2)} \right] \right) \right]^{1/2}$$

For large *b*, all formulae OK.

For small b, s/\sqrt{b} and $s/\sqrt{(b+\sigma_b^2)}$ overestimate the significance.

Could be important in optimization of searches with low background.

Formula maybe also OK if model is not simple on/off experiment, e.g., several background control measurements (check this).

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

Large-sample (asymptotic) properties of MLEs

Suppose we have an i.i.d. data sample of size $n: x_1, ..., x_n$

In the large-sample (or "asymptotic") limit $(n \rightarrow \infty)$ and assuming regularity conditions one can show that the likelihood and MLE have several important properties.

The regularity conditions include:

- the boundaries of the data space cannot depend on the parameter;
- the parameter cannot be on the edge of the parameter space;
- $\ln L(\theta)$ must be differentiable;
- the only solution to $\partial \ln L / \partial \theta = 0$ is $\hat{\theta}$.

In the slides immediately following the properties are shown without proof for a single parameter; the corresponding properties hold also for the multiparameter case, $\theta = (\theta_1, ..., \theta_m)$.
log-likelihood becomes quadratic

The likelihood function becomes Gaussian in shape, i.e. the log-likelihood becomes quadratic (parabolic).



The MLE becomes increasingly precise as the (log)-likelihood becomes more tightly concentrated about its peak, but $L(\theta) = P(\mathbf{x}|\theta)$ is the probability for \mathbf{x} , not a pdf for θ .

The MLE converges to the true parameter value

In the large-sample limit, the MLE converges in probability to the true parameter value.

That is, for any $\varepsilon > 0$,

$$\lim_{n \to \infty} P(|\hat{\theta} - \theta| > \epsilon) = 0$$

The MLE is said to be *consistent*.

MLE is asymptotically unbiased

In general the MLE can be biased, but in the large-sample limit, this bias goes to zero:

$$\lim_{n \to \infty} E[\hat{\theta}] - \theta = 0$$

(Recall for the exponential parameter we found the bias was identically zero regardless of the sample size *n*.)

The MLE's variance approaches the MVB

In the large-sample limit, the variance of the MLE approaches the minimum-variance bound, i.e., the information inequality becomes an equality (and bias goes to zero):

$$\lim_{n \to \infty} V[\hat{\theta}] = -\frac{1}{E\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]}$$

The MLE is said to be *asymptotically efficient*.

The MLE's distribution becomes Gaussian

In the large-sample limit, the pdf of the MLE becomes Gaussian,

$$f(\hat{\theta}) = \frac{1}{\sqrt{2\pi\sigma_{\hat{\theta}}}} e^{-(\hat{\theta}-\theta)^2/2\sigma_{\hat{\theta}}^2}$$

where $\sigma^2_{\hat{ heta}}$ is the minimum variance bound (note bias is zero).

For example, exponential MLE with sample size n = 100.

Note that for exponential, MLE is arithmetic average, so Gaussian MLE seen to stem from Central Limit Theorem.



Distribution of MLE of exponential parameter



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.



DESY / U. Hamburg Basics of Statistical Data Analysis