Statistical Methods for Particle Physics Day 2: Statistical Tests and Limits https://indico.desy.de/indico/event/19085/



Heimnoitz Alliance



Terascale Statistics School DESY, 19-23 February, 2018

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

Outline

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

→ Day 2: Discovery and Limits

Comments on multivariate methods (brief)

p-values

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

Frequentist statistical tests

Consider a hypothesis H_0 and alternative 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.,

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

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 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

Type-I, Type-II errors

Rejecting the hypothesis H_0 when it is true is a Type-I error. The maximum probability for this is the size of the test:

$$P(x \in W \mid H_0) \le \alpha$$

But we might also accept H_0 when it is false, and an alternative H_1 is true.

This is called a Type-II error, and occurs with probability

$$P(x \in \mathbf{S} - W | H_1) = \beta$$

One minus this is called the power of the test with respect to the alternative H_1 :

Power =
$$1 - \beta$$

A simulated SUSY event



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.

Physics context of a statistical test

Event Selection: the event types in question are both known to exist.

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: the null hypothesis is

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

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 \vec{x} , e.g., $f(\vec{x}|H_0)$, $f(\vec{x}|H_1)$, etc.

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

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

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

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.

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.

G. Cowan

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: Fisher discriminant (Deep) neural networks Kernel density methods Support Vector Machines Decision trees Boosting Bagging

Resources on multivariate methods

C.M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006

T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning, 2nd ed., Springer, 2009

R. Duda, P. Hart, D. Stork, Pattern Classification, 2nd ed., Wiley, 2001

A. Webb, Statistical Pattern Recognition, 2nd ed., Wiley, 2002.

Ilya Narsky and Frank C. Porter, *Statistical Analysis Techniques in Particle Physics*, Wiley, 2014.

朱永生(编著),实验数据多元统计分析,科学出版社, 北京,2009。

Software

Rapidly growing area of development – two important resources:

TMVA, Höcker, Stelzer, Tegenfeldt, Voss, Voss, physics/0703039 From tmva.sourceforge.net, also distributed with ROOT Variety of classifiers Good manual, widely used in HEP scikit-learn

> Python-based tools for Machine Learning scikit-learn.org Large user community

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}_{obs} . Note – "less compatible with H" means "more compatible with some alternative H".



G. Cowan

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).

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, $p_H \sim$ Uniform[0,1] and is concentrated toward zero for Some class of relevant alternatives.



G. Cowan

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

One can show the distribution of the *p*-value of H, under assumption of H, is uniform in [0,1].

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

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

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 .

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.



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

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

E.g. Z = 5 (a "5 sigma effect") corresponds to $p = 2.9 \times 10^{-7}$.

G. Cowan

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 DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2 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 = $P(n \ge 5; b = 0.5, s = 0)$
= $1.7 \times 10^{-4} \neq P(s = 0)$



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 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 $P(\text{data in critical region}) \le \alpha$ 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 would not be rejected in a test of size α (confidence level is $1 - \alpha$).

The interval will cover the true value of θ 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 $p_{\theta} = \alpha$ and solve for θ . G. Cowan DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

Frequentist upper limit on Poisson parameter

Consider again the case of observing $n \sim \text{Poisson}(s + b)$.

Suppose b = 4.5, $n_{obs} = 5$. Find upper limit on *s* at 95% CL.

When testing *s* values to find upper limit, relevant alternative is s = 0 (or lower *s*), so critical region at low *n* and *p*-value of hypothesized *s* is $P(n \le n_{obs}; s, b)$.

Upper limit s_{up} at $CL = 1 - \alpha$ from setting $\alpha = p_s$ and solving for s:

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

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

G. Cowan

2

Frequentist upper limit on Poisson parameter

Upper limit s_{up} at $CL = 1 - \alpha$ found from $p_s = \alpha$.



 $n_{\rm obs} = 5,$ b = 4.5

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

 $n \sim \text{Poisson}(s+b)$: frequentist upper limit on *s* For low fluctuation of *n* formula can give negative result for s_{up} ; i.e. confidence interval is empty.



G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

Limits near a physical boundary

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?



G. Cowan
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

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

G. Cowan

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 but this is OK as long as L(s) dies off for large s.

Not invariant under change of parameter — if we had used instead a flat prior for, say, the mass of the Higgs boson, this would imply a non-flat prior for the expected number of Higgs events.

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 (coverage will depend on true *s*).

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').

G. Cowan

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

Priors from formal rules

Because of difficulties in encoding a vague degree of belief in a prior, one often attempts to derive the prior from formal rules, e.g., to satisfy certain invariance principles or to provide maximum information gain for a certain set of measurements.

> Often called "objective priors" Form basis of Objective Bayesian Statistics

The priors do not reflect a degree of belief (but might represent possible extreme cases).

In Objective Bayesian analysis, can use the intervals in a frequentist way, i.e., regard Bayes' theorem as a recipe to produce an interval with certain coverage properties.

Priors from formal rules (cont.)

For a review of priors obtained by formal rules see, e.g.,

Robert E. Kass and Larry Wasserman, *The Selection of Prior Distributions by Formal Rules*, J. Am. Stat. Assoc., Vol. 91, No. 435, pp. 1343-1370 (1996).

Formal priors have not been widely used in HEP, but there is recent interest in this direction, especially the reference priors of Bernardo and Berger; see e.g.

L. Demortier, S. Jain and H. Prosper, *Reference priors for high energy physics*, Phys. Rev. D 82 (2010) 034002, arXiv:1002.1111.

D. Casadei, *Reference analysis of the signal + background model in counting experiments*, JINST 7 (2012) 01012; arXiv:1108.4270.

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
$$\theta$$
 therefore $p_{\theta} = \int_{t_{\theta,obs}}^{\infty} f(t_{\theta}|\theta) dt_{\theta}$
G. Cowan DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2 need pdf

Confidence region from Wilks' theorem Wilks' theorem says (in large-sample limit and providing certain conditions hold...)

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

Assuming this holds, the *p*-value is

 $p_{\theta} = 1 - F_{\chi_n^2}(t_{\theta})$ where $F_{\chi_n^2}(t_{\theta}) \equiv \int_0^{t_{\theta}} f_{\chi_n^2}(t'_{\theta}) t'_{\theta}$

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

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

G. Cowan

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.

G. Cowan

Example of interval from $\ln L$

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

Exponential example, now with only 5 events:



Parameter estimate and approximate 68.3% CL confidence interval:

 $\hat{\tau} = 0.85^{+0.52}_{-0.30}$



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 $CL = 1 - \alpha$.

$1 - \alpha$	$ar{Q}_{lpha}$						
	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		

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

$$s_{i} = s_{\text{tot}} \int_{\text{bin } i} f_{s}(x; \boldsymbol{\theta}_{s}) \, dx \,, \quad b_{i} = b_{\text{tot}} \int_{\text{bin } i} f_{b}(x; \boldsymbol{\theta}_{b}) \, dx \,.$$

signal background

G. Cowan

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

The profile likelihood ratio

Base significance test on the profile likelihood ratio:

 $\lambda(\mu) = \frac{L(\mu, \hat{\hat{\theta}})}{L(\hat{\mu}, \hat{\theta})}$ maximize L for specified μ maximize L

The likelihood ratio of point hypotheses gives optimum test (Neyman-Pearson lemma). In practice the profile LR is near-optimal.

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

G. Cowan

Test statistic for discovery

Try to reject background-only ($\mu = 0$) hypothesis using

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

i.e. here only regard upward fluctuation of data as 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 .

G. Cowan

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

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

G. Cowan

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

Monte Carlo test of asymptotic formula

 $n \sim \text{Poisson}(\mu s + b)$ $m \sim \text{Poisson}(\tau b)$

Here take $\tau = 1$.

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



Example of discovery: 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 ($\mu = 1$) 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.

Return to interval estimation

Suppose a model contains a parameter μ ; we want to know which values are consistent with the data and which are disfavoured.

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 probability that the true value of μ will be rejected is not greater than α , so by construction the confidence interval will contain the true value of μ with probability $\geq 1 - \alpha$.

The interval depends on the choice of the test (critical region).

If the test is formulated in terms of a *p*-value, p_{μ} , then the confidence interval represents those values of μ for which $p_{\mu} > \alpha$.

To find the end points of the interval, set $p_{\mu} = \alpha$ and solve for μ . G. Cowan

Test statistic for upper limits

cf. Cowan, Cranmer, Gross, Vitells, arXiv:1007.1727, EPJC 71 (2011) 1554. For purposes of setting an upper limit on μ one can 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: $p_{\mu} = \int_{q_{\mu,\text{obs}}}^{\infty} f(q_{\mu}|\mu) dq_{\mu}$

Large sample approximation:

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

95% CL upper limit on μ is highest value for which *p*-value is not less than 0.05.

G. Cowan

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

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".



for s = 6, b = 9.

Low sensitivity to μ

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

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



G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

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

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 .



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.

G. Cowan

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

The CL_s procedure (2)

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



G. Cowan

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.,



 $CL_s \leq \alpha$

Increases "effective" *p*-value when the two distributions become close (prevents exclusion if sensitivity is low).

G. Cowan

Setting upper limits on $\mu = \sigma / \sigma_{\rm SM}$

Carry out the CLs procedure for the parameter $\mu = \sigma/\sigma_{SM}$, resulting in an upper limit μ_{up} .

In, e.g., a Higgs search, this is done for each value of $m_{\rm H}$.

At a given value of $m_{\rm H}$, we have an observed value of $\mu_{\rm up}$, and we can also find the distribution $f(\mu_{\rm up}|0)$:



 $\pm 1\sigma$ (green) and $\pm 2\sigma$ (yellow) bands from toy MC;

Vertical lines from asymptotic formulae.

G. Cowan

How to read the green and yellow limit plots For every value of $m_{\rm H}$, find the CLs 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 $\mu = 0$.

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



G. Cowan

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

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

Counting experiment with known background Count a number of events $n \sim Poisson(s+b)$, 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

 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
Better approximation for significance Poisson likelihood for parameter *s* is

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

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

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

params.

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 << b.

G. Cowan

 $n \sim \text{Poisson}(s+b)$, 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 parmeter): $L(0, \hat{b}(0))$

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

G. Cowan

Ingredients for profile likelihood ratio

To construct profile likelihood ratio from this need estimators:

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

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

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

G. Cowan

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

Asimov approximation for median significance

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

$$n \to s + b$$

$$m \to \tau b$$

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

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

 Z_{i}

Limiting cases

Expanding the Asimov formula in powers of *s/b* and σ_b^2/b (= 1/ τ) 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

DESY Terascale School of Statistics / 19-23 Feb 2018 / Day 2

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

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 (checking this).

Finally

Four lectures only enough for a brief introduction to:

Parameter estimation

Unfolding

Statistical tests for discovery and limits

Experimental sensitivity

Many other important topics; some covered in rest of week:

Bayesian methods, MCMC

Multivariate methods, Machine Learning

The look-elsewhere effect, etc., etc.

Final thought: once the basic formalism is understood, most of the work focuses on building the model, i.e., 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).

G. Cowan



G. Cowan

Goodness of fit from the likelihood ratio

Suppose we model data using a likelihood $L(\mu)$ that depends on N parameters $\mu = (\mu_1, ..., \mu_N)$. Define the statistic

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

Value of t_{μ} reflects agreement between hypothesized μ and the data.

Good agreement means $\hat{\mu} \approx \mu$, so t_{μ} is small;

Larger t_{μ} means less compatibility between data and μ .

Quantify "goodness of fit" with *p*-value:
$$p_{\mu} = \int_{t_{\mu,\text{obs}}}^{\infty} f(t_{\mu}|\mu) dt_{\mu}$$

G. Cowan

Likelihood ratio (2)

Now suppose the parameters $\boldsymbol{\mu} = (\mu_1, ..., \mu_N)$ can be determined by another set of parameters $\boldsymbol{\theta} = (\theta_1, ..., \theta_M)$, with M < N.

E.g. in LS fit, use $\mu_i = \mu(x_i; \theta)$ where x is a control variable.

Define the statistic

fit *M* parameters

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

fit *N* parameters

Use q_{μ} to test hypothesized functional form of $\mu(x; \theta)$. To get *p*-value, need pdf $f(q_{\mu}|\mu)$.

G. Cowan

Wilks' Theorem (1938)

Wilks' Theorem: if the hypothesized parameters $\mu = (\mu_1, ..., \mu_N)$ are true then in the large sample limit (and provided certain conditions are satisfied) t_{μ} and q_{μ} follow chi-square distributions.

For case with $\boldsymbol{\mu} = (\mu_1, ..., \mu_N)$ fixed in numerator:



S.S. Wilks, The large-sample distribution of the likelihood ratio for testing composite hypotheses, Ann. Math. Statist. 9 (1938) 60-2.

G. Cowan

Goodness of fit with Gaussian data

Suppose the data are *N* independent Gaussian distributed values:

$$y_i \sim \text{Gauss}(\mu_i, \sigma_i)$$
, $i = 1, \dots, N$
want to estimate known

Likelihood:
$$L(\mu) = \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma_i}} e^{-(y_i - \mu_i)^2/2\sigma_i^2}$$

Log-likelihood:
$$\ln L(\boldsymbol{\mu}) = -\frac{1}{2} \sum_{i=1}^{N} \frac{(y_i - \mu_i)^2}{\sigma_i^2} + C$$

ML estimators: $\hat{\mu}_i = y_i$ $i = 1, \dots, N$

G. Cowan

Likelihood ratios for Gaussian data

The goodness-of-fit statistics become

$$t_{\mu} = -2\ln\frac{L(\mu)}{L(\hat{\mu})} = \sum_{i=1}^{N} \frac{(y_i - \mu_i)^2}{\sigma_i^2} \qquad f(t_{\mu}|\mu) \sim \chi_N^2$$

$$q_{\mu} = -2\ln\frac{L(\mu(\hat{\theta}))}{L(\hat{\mu})} = \sum_{i=1}^{N} \frac{(y_i - \mu_i(\hat{\theta}))^2}{\sigma_i^2} \qquad f(q_{\mu}|\mu) \sim \chi^2_{N-M}$$

So Wilks' theorem formally states the well-known property of the minimized chi-squared from an LS fit.

G. Cowan

Likelihood ratio for Poisson data

Suppose the data are a set of values $n = (n_1, ..., n_N)$, e.g., the numbers of events in a histogram with *N* bins.

Assume $n_i \sim \text{Poisson}(v_i)$, i = 1, ..., N, all independent. Goal is to estimate $v = (v_1, ..., v_N)$.

Likelihood:
$$L(\nu) = \prod_{i=1}^{N} \frac{\nu_i^{n_i}}{n_i!} e^{-\nu_i}$$

Log-likelihood:
$$\ln L(\boldsymbol{\nu}) = \sum_{i=1}^{N} [n_i \ln \nu_i - \nu_i] + C$$

ML estimators: $\hat{\nu}_i = n_i$, $i = 1, \dots, N$

G. Cowan

Goodness of fit with Poisson data

The likelihood ratio statistic (all parameters fixed in numerator):

$$t_{oldsymbol{
u}} = -2\lnrac{L(oldsymbol{
u})}{L(\hat{oldsymbol{
u}})}$$

$$= -2\sum_{i=1}^{N} \left[n_i \ln \frac{\nu_i}{\hat{\nu}_i} - \nu_i + \hat{\nu}_i \right]$$

$$= -2\sum_{i=1}^{N} \left[n_i \ln \frac{\nu_i}{n_i} - \nu_i + n_i \right]$$

Wilks' theorem: $f(t_{\nu}|\nu) \sim \chi_N^2$

G. Cowan

Goodness of fit with Poisson data (2)

Or with *M* fitted parameters in numerator:

$$q_{\nu} = -2\ln\frac{L(\nu(\hat{\theta}))}{L(\hat{\nu})} = -2\sum_{i=1}^{N} \left[n_i \ln\frac{\nu_i(\hat{\theta})}{n_i} - \nu_i(\hat{\theta}) + n_i\right]$$

Wilks' theorem: $f(q_{\nu}|\nu) \sim \chi^2_{N-M}$

Use t_{μ} , q_{μ} to quantify goodness of fit (*p*-value). Sampling distribution from Wilks' theorem (chi-square). Exact in large sample limit; in practice good approximation for surprisingly small n_i (~several).

G. Cowan

Goodness of fit with multinomial data

Similar if data $\mathbf{n} = (n_1, ..., n_N)$ follow multinomial distribution:

$$P(\mathbf{n}|\mathbf{p}, n_{\text{tot}}) = \frac{n_{\text{tot}}!}{n_1! n_2! \dots n_N!} p_1^{n_1} p_2^{n_2} \dots p_N^{n_N}$$

E.g. histogram with N bins but fix: $n_{\text{tot}} = \sum_{i=1}^{N} n_i$

Log-likelihood:
$$\ln L(\nu) = \sum_{i=1}^{N} n_i \ln \frac{\nu_i}{n_{\text{tot}}} + C$$
 $(\nu_i = p_i n_{\text{tot}})$

ML estimators: $\hat{\nu}_i = n_i$ (Only *N*-1 independent; one is n_{tot} minus sum of rest.)

Goodness of fit with multinomial data (2)

The likelihood ratio statistics become:

$$t_{\nu} = -2\sum_{i=1}^{N} n_i \ln \frac{\nu_i}{n_i} \qquad f(t_{\nu}|\nu) \sim \chi^2_{N-1}$$
$$q_{\nu} = -2\sum_{i=1}^{N} n_i \ln \frac{\nu_i(\hat{\theta})}{n_i} \qquad f(q_{\nu}|\nu) \sim \chi^2_{N-M-1}$$

One less degree of freedom than in Poisson case because effectively only *N*–1 parameters fitted in denominator.

Estimators and g.o.f. all at once

Evaluate numerators with θ (not its estimator):

$$\chi_{\rm P}^2(\theta) = -2\sum_{i=1}^N \left[n_i \ln \frac{\nu_i(\theta)}{n_i} - \nu_i(\theta) + n_i \right]$$
(Poisson)
$$\chi_{\rm M}^2(\theta) = -2\sum_{i=1}^N n_i \ln \frac{\nu_i(\theta)}{n_i}$$
(Multinomial)

These are equal to the corresponding $-2 \ln L(\theta)$, so minimizing them gives the usual ML estimators for θ .

The minimized value gives the statistic q_{μ} , so we get goodness-of-fit for free.

Steve Baker and Robert D. Cousins, *Clarification of the use of the chi-square and likelihood functions in fits to histograms*, NIM **221** (1984) 437.

G. Cowan