# Statistical Methods for Particle Physics
## Lecture 1: introduction & statistical tests

`https://indico.fnal.gov/conferenceTimeTable.py?confId=11505`

Lectures on Statistics
HCPSS – Fermilab
11,12 August 2016

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

# Outline

➡ Lecture 1:  Introduction and review of fundamentals
    Probability, random variables, pdfs
    Parameter estimation, maximum likelihood
    Statistical tests for discovery and limits

Lecture 2:  Further topics
    Brief overview of multivariate methods
    Nuisance parameters and systematic uncertainties
    Experimental sensitivity

# Some statistics books, papers, etc.

G. Cowan, *Statistical Data Analysis*, Clarendon, Oxford, 1998

R.J. Barlow, *Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences*, Wiley, 1989

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

L. Lyons, *Statistics for Nuclear and Particle Physics*, CUP, 1986

F. James., *Statistical and Computational Methods in Experimental Physics*, 2nd ed., World Scientific, 2006
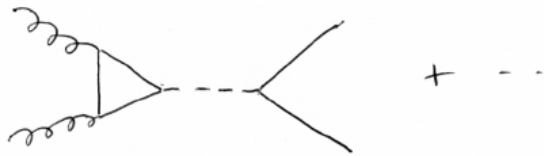
S. Brandt, *Statistical and Computational Methods in Data Analysis*, Springer, New York, 1998 (with program library on CD)

K.A. Olive et al. (Particle Data Group), *Review of Particle Physics*, Chin. Phys. C, 38, 090001 (2014); see also `pdg.lbl.gov` sections on probability, statistics, Monte Carlo
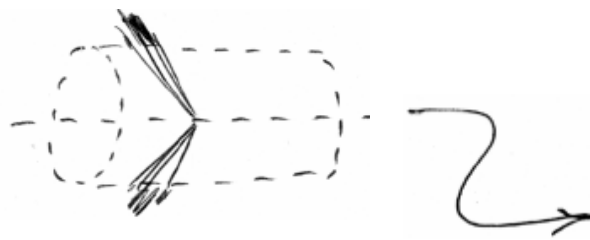
# Theory ↔ Statistics ↔ Experiment

## Theory (model, hypothesis):

## Experiment:

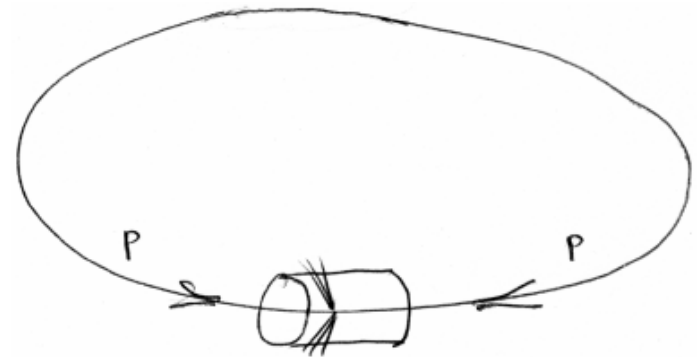$$\mathcal{L} = -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} + i \bar{\psi} \not{D} \psi + \cdots$$

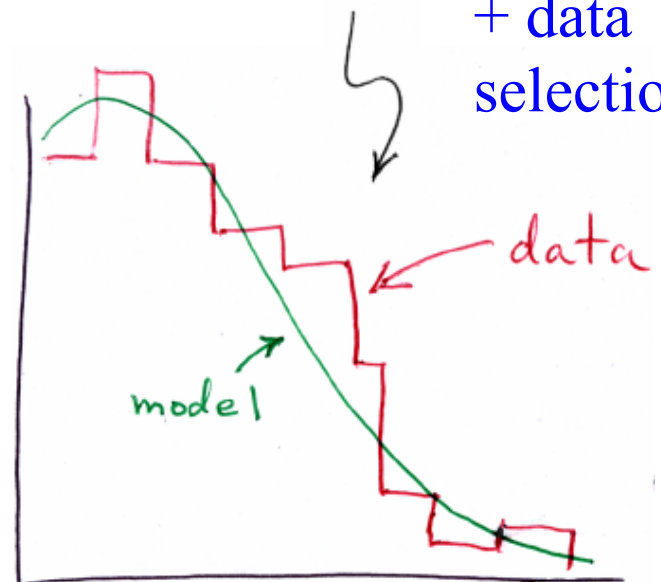$$\sigma = \frac{G_F \alpha_s^2 m_H^2}{288 \sqrt{2\pi}} \times \sim$$

+ simulation
of detector
and cuts

+ data
selection

data

model

# 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):

$$P(A) = \text{degree of belief that } A \text{ is true}$$

Conditional probability:

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

Bayes' theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\Sigma_i P(B|A_i)P(A_i)}$$

# Frequentist Statistics − general philosophy

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

Probability = limiting frequency

Probabilities such as

$P$ (Higgs boson exists),
$P$ $(0.117 < \alpha_s < 0.121)$,

etc. are either 0 or 1, but we don't know which.

The tools of frequentist statistics tell us what to expect, under the assumption of certain probabilities, about hypothetical repeated observations.

A hypothesis is is preferred if the data are found in a region of high predicted probability (i.e., where an alternative hypothesis predicts lower probability).

# Bayesian Statistics − general philosophy

In Bayesian statistics, use subjective probability for hypotheses:

probability of the data assuming
hypothesis $H$ (the likelihood)

prior probability, i.e.,
before seeing the data

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

posterior probability, i.e.,
after seeing the data

normalization involves sum
over all possible hypotheses

Bayes' theorem has an "if-then" character:  If your prior
probabilities were $\pi(H)$, then it says how these probabilities
should change in the light of the data.

No general prescription for priors (subjective!)

# Quick review of frequentist parameter estimation

Suppose we have a pdf characterized by one or more parameters:

$$f(x; \theta) = \frac{1}{\theta} e^{-x/\theta}$$

random variable      parameter

Suppose we have a sample of observed values: $\vec{x} = (x_1, \dots, x_n)$

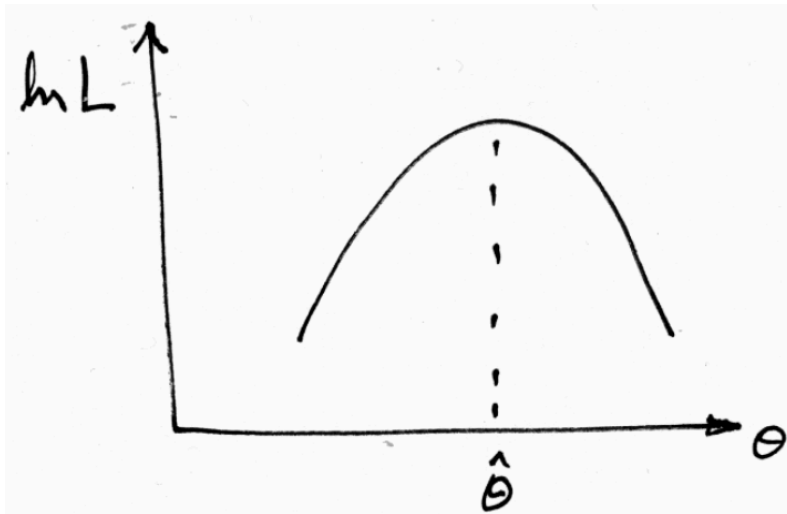We want to find some function of the data to estimate the parameter(s):

$$\hat{\theta}(\vec{x})$$   $\leftarrow$ 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.

# Maximum Likelihood (ML) estimators

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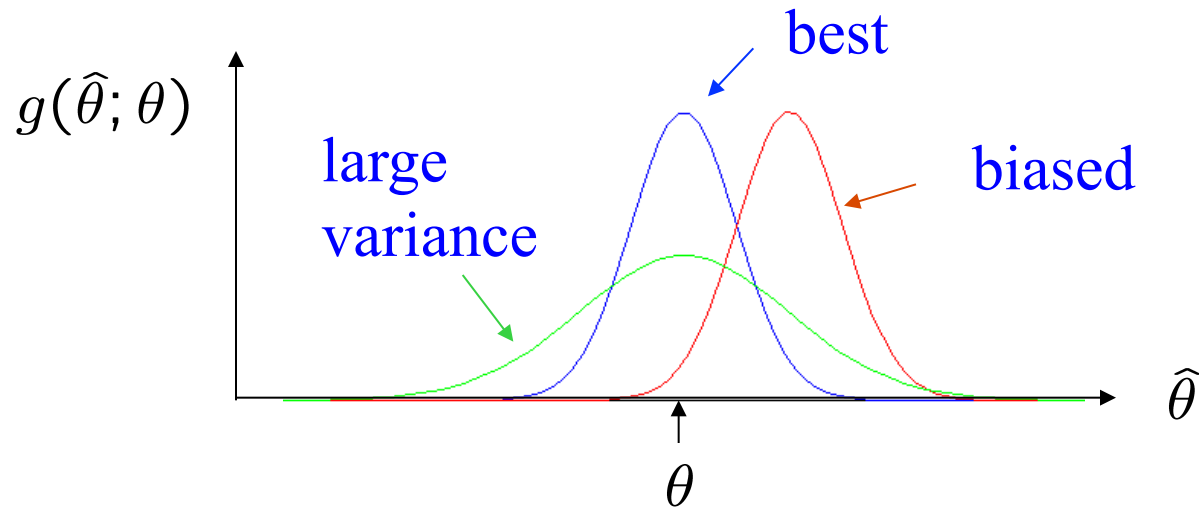


$$\hat{\theta} = \underset{\theta}{\mathrm{argmax}}\, L(x|\theta)$$

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

# Properties of estimators

Estimators are functions of the data and thus characterized by a sampling distribution with a given (co)variance:

$$g(\widehat{\theta}; \theta)$$

best

large variance

biased

$$\widehat{\theta}$$

$$\theta$$

In general they may have a nonzero bias:   $b = E[\hat{\theta}] - \theta$

Under conditions usually satisfied in practice, bias of ML estimators is zero in the large sample limit, and the variance is as small as possible for unbiased estimators.

# ML example: parameter of exponential pdf

Consider exponential pdf, $\quad f(t; \tau) = \dfrac{1}{\tau} e^{-t/\tau}$

and suppose we have i.i.d. data, $\quad t_1, \ldots, t_n$

The likelihood function is $\quad L(\tau) = \displaystyle\prod_{i=1}^{n} \frac{1}{\tau} e^{-t_i/\tau}$

The value of $\tau$ for which $L(\tau)$ is maximum also gives the maximum value of its logarithm (the log-likelihood function):

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

# ML example: parameter of exponential pdf (2)

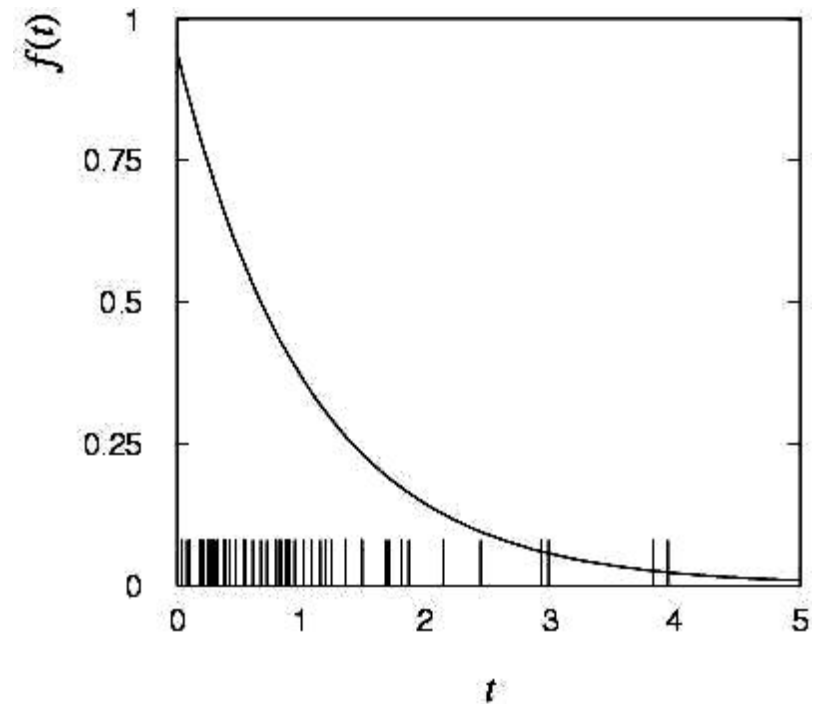Find its maximum by setting $\dfrac{\partial \ln L(\tau)}{\partial \tau} = 0$ ,

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

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

We find the ML estimate:

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

# ML example:  parameter of exponential pdf (3)

For the exponential distribution one has for mean, variance:

$$E[t] = \int_0^\infty t \, \frac{1}{\tau} e^{-t/\tau} \, dt = \tau$$

$$V[t] = \int_0^\infty (t-\tau)^2 \, \frac{1}{\tau} e^{-t/\tau} \, dt = \tau^2$$

For the ML estimator $\quad \hat{\tau} = \dfrac{1}{n} \sum_{i=1}^n t_i \quad$ 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 from information inequality

The information inequality (RCF) sets a lower bound on the variance of any estimator (not only ML):

Minimum Variance Bound (MVB)

$$V[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 \bigg/ E\left[-\frac{\partial^2 \ln L}{\partial \theta^2}\right]$$

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

Often the bias $b$ is small, and equality either holds exactly or is a good approximation (e.g. large data sample limit). Then,

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

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

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

# Variance of estimators: graphical method

Expand $\ln L(\theta)$ about its maximum:

$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[\frac{\partial \ln L}{\partial \theta}\right]_{\theta=\hat{\theta}}(\theta - \hat{\theta}) + \frac{1}{2!}\left[\frac{\partial^2 \ln L}{\partial \theta^2}\right]_{\theta=\hat{\theta}}(\theta - \hat{\theta})^2 + \ldots$$

First term is $\ln L_{\mathrm{max}}$, second term is zero, for third term use information inequality (assume equality):

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

i.e., $\qquad \ln L(\hat{\theta} \pm \hat{\sigma}_{\hat{\theta}}) \approx \ln L_{\mathrm{max}} - \frac{1}{2}$

$\rightarrow$ to get $\hat{\sigma}_{\hat{\theta}}$, change $\theta$ away from $\hat{\theta}$ until $\ln L$ decreases by 1/2.
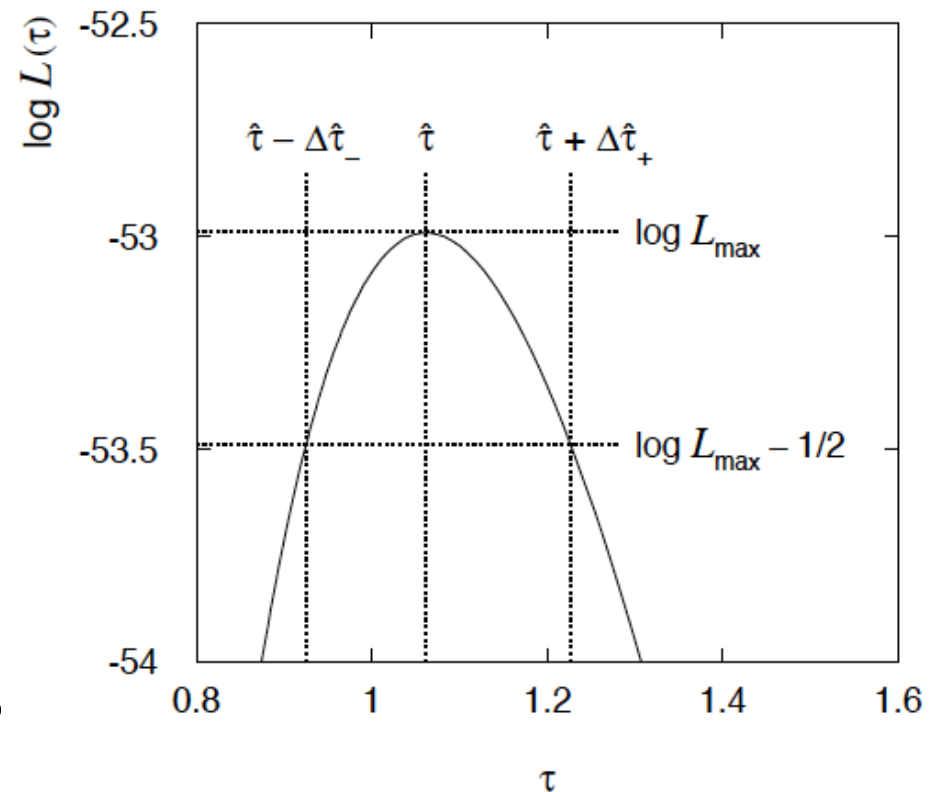
# Example of variance by graphical method

ML example with exponential:



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

$$\Delta\hat{\tau}_- = 0.137$$

$$\Delta\hat{\tau}_+ = 0.165$$

$$\hat{\sigma}_{\hat{\tau}} \approx \Delta\hat{\tau}_- \approx \Delta\hat{\tau}_+ \approx 0.15$$

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

# Information inequality for *N* parameters

Suppose we have estimated $N$ parameters $\vec{\theta} = (\theta_1, \ldots, \theta_N)$ .

The (inverse) minimum variance bound is given by the Fisher information matrix:

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

The information inequality then states that $V - I^{-1}$ is a positive semi-definite matrix, where $V_{ij} = \mathsf{cov}[\hat{\theta}_i, \hat{\theta}_j]$ .   Therefore

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

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

# Frequentist statistical tests
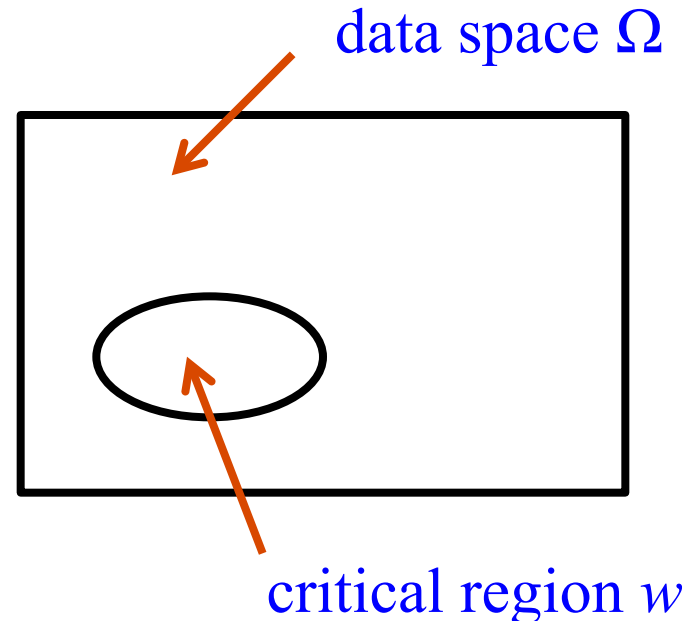
Consider a hypothesis $H_0$ and alternative $H_1$.

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 $\alpha$, assuming $H_0$ is correct, to observe the data there, i.e.,

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

Need inequality if data are discrete.

$\alpha$ is called the size or significance level of the test.

If $x$ is observed in the critical region, reject $H_0$.
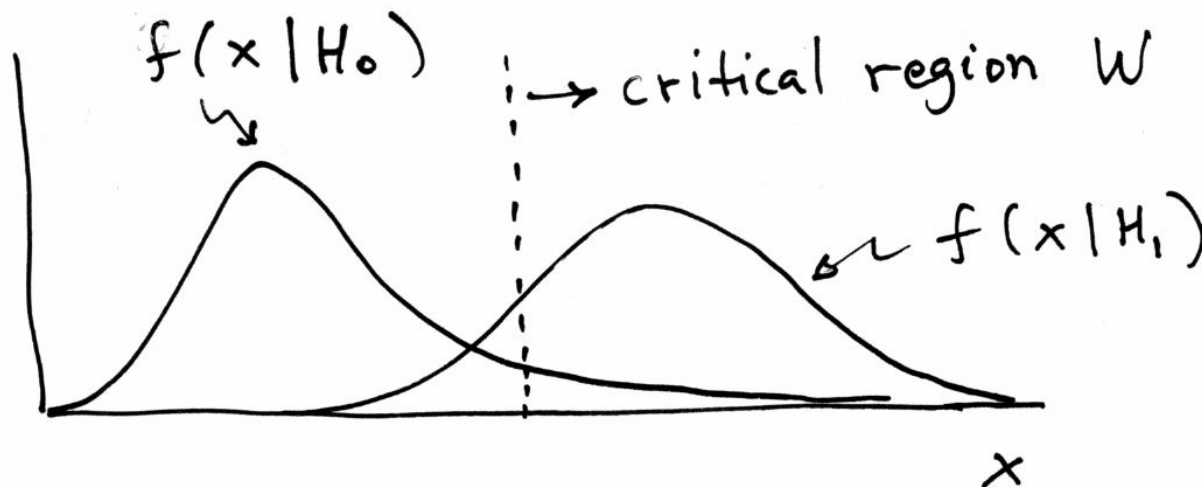
data space $\Omega$

critical region $w$

# Definition of a test (2)

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

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:

# 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:

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

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

$$P(x \in S - W \mid H_1) = \beta$$

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

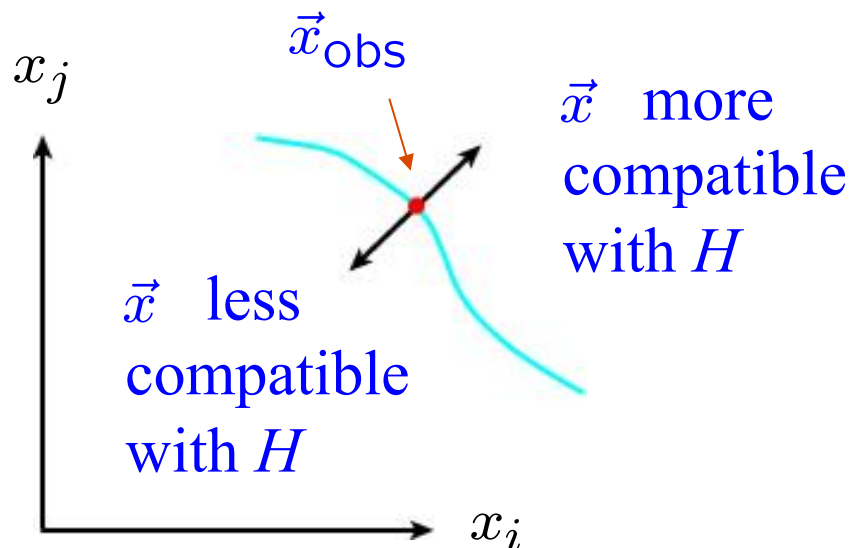$$\text{Power} = 1 - \beta$$

# Testing significance / goodness-of-fit

Suppose hypothesis $H$ predicts pdf $f(\vec{x}|H)$ for a set of observations $\vec{x} = (x_1, \ldots, x_n)$ .

We observe a single point in this space: $\vec{x}_{\mathrm{obs}}$

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 $\vec{x}_{\mathrm{obs}}$ .

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

$x_j$

$\vec{x}_{\mathrm{obs}}$

$\vec{x}$ more compatible with $H$

$\vec{x}$ less compatible with $H$

$x_i$

# *p*-values

Express 'goodness-of-fit' by giving the *p*-value for *H*:

$p$ = probability, under assumption of *H*, to observe data with equal or lesser compatibility with *H* relative to the data we got.

⚠️  This is not the probability that *H* is true!

In frequentist statistics we don't talk about $P(H)$ (unless *H* represents a repeatable observation). In Bayesian statistics we do; use Bayes' theorem to obtain
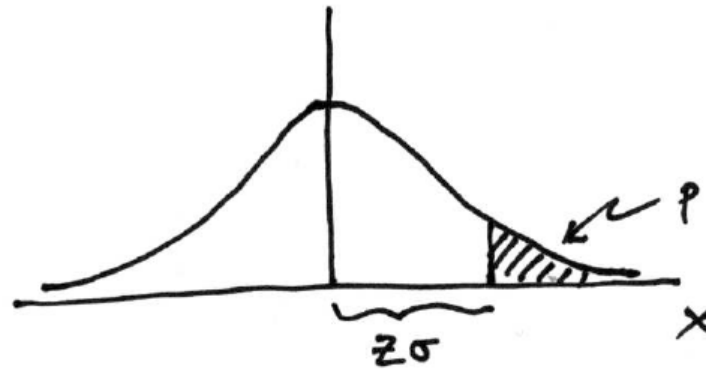
$$P(H|\vec{x}) = \frac{P(\vec{x}|H)\pi(H)}{\int P(\vec{x}|H)\pi(H)\,dH}$$

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

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

# Significance from *p*-value

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



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

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

# Test statistics and *p*-values

Consider a parameter $\mu$ proportional to rate of signal process.

Often define a function of the data (test statistic) $q_\mu$ that reflects level of agreement between the data and the hypothesized value $\mu$.

Usually define $q_\mu$ so that higher values increasingly incompatibility with the data (more compatible with a relevant alternative).

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

$$p_\mu = \int_{q_{\mu,\text{obs}}}^{\infty} f(q_\mu | \mu)\, dq_\mu$$

observed value of $q_\mu$                    pdf of $q_\mu$ assuming $\mu$

Equivalent formulation of test:  reject $\mu$ if $p_\mu < \alpha$.

# Confidence interval from inversion of a test

Carry out a test of size $\alpha$ for all values of $\mu$.

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

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

The interval will cover the true value of $\mu$ with probability $\geq 1 - \alpha$.

Equivalently, the parameter values in the confidence interval have *p*-values of at least $\alpha$.

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

# 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)

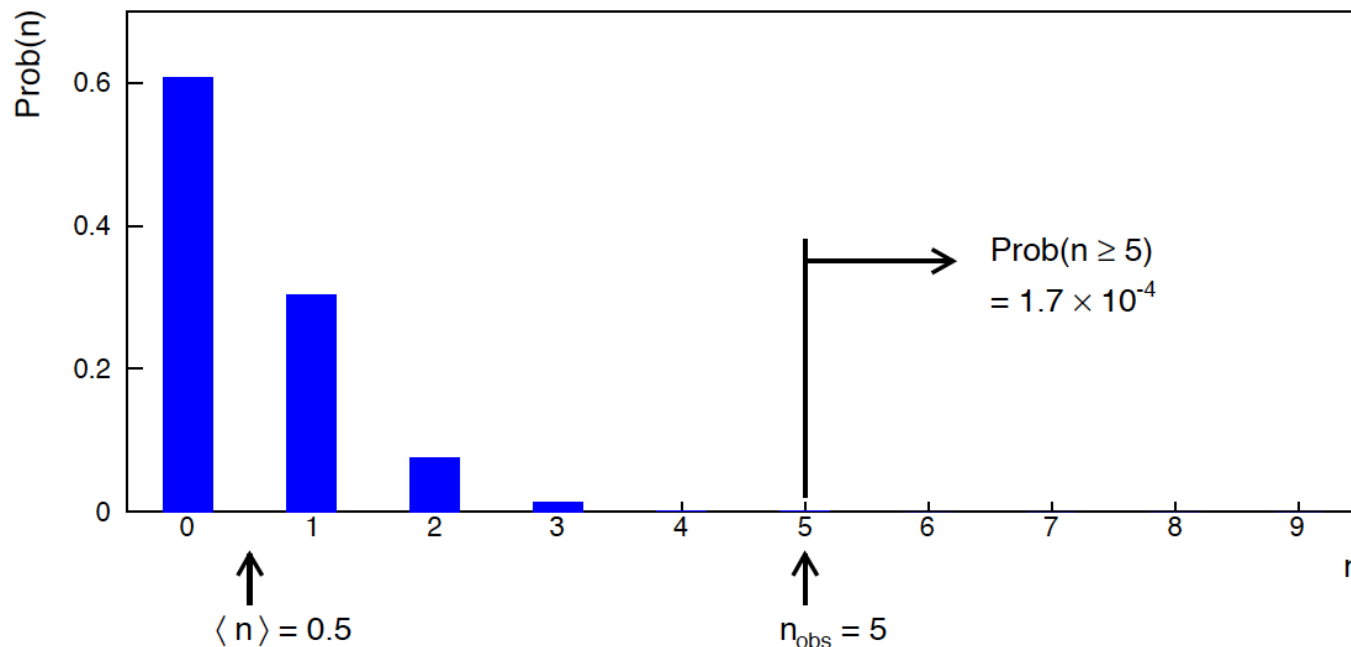In both cases need to ask what is relevant alternative hypothesis.

# Poisson counting experiment: discovery $p$-value

Suppose $b = 0.5$ (known), and we observe $n_{\text{obs}} = 5$.

Should we claim evidence for a new discovery?

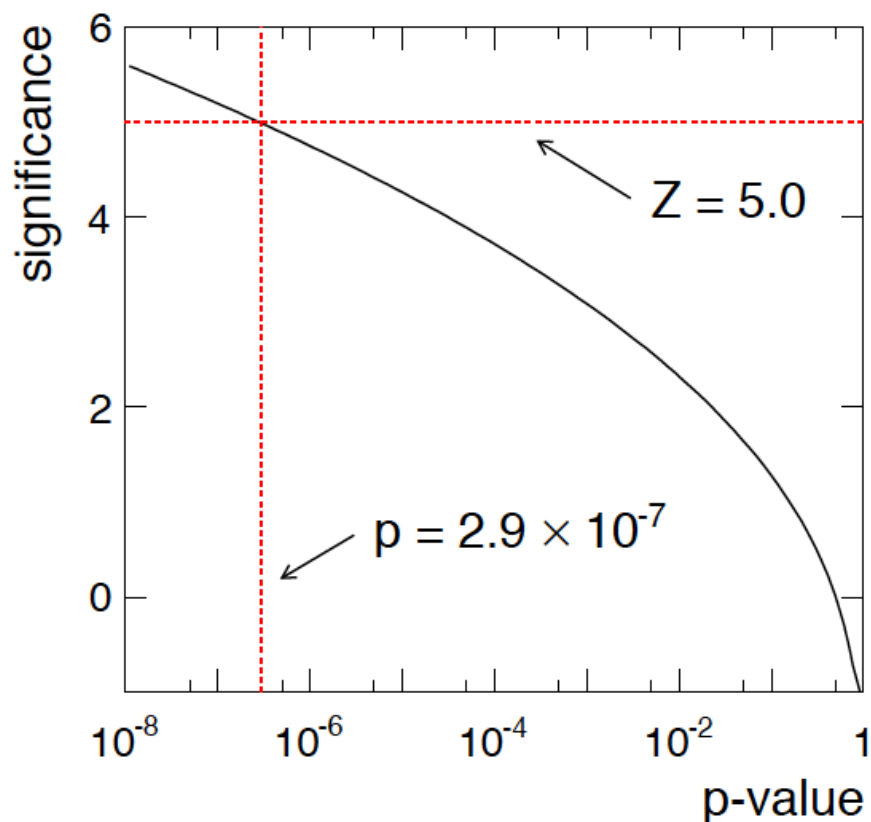Take $n$ itself as the test statistic, $p$-value for hypothesis $s = 0$ is

$$p\text{-value} = P(n \geq 5; b = 0.5, s = 0)$$

$$= 1.7 \times 10^{-4} \neq P(s = 0)!$$

# Poisson counting experiment: discovery significance

Equivalent significance for $p = 1.7 \times 10^{-4}$: $\quad Z = \Phi^{-1}(1 - p) = 3.6$

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 signal-like 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.

# 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.

Relevant alternative is $s = 0$ (critical region at low $n$)

$p$-value of hypothesized $s$ is $P(n \leq n_{obs}; s, b)$

Upper limit $s_{up}$ at CL $= 1 - \alpha$ found by solving $p_s = \alpha$ for $s$:

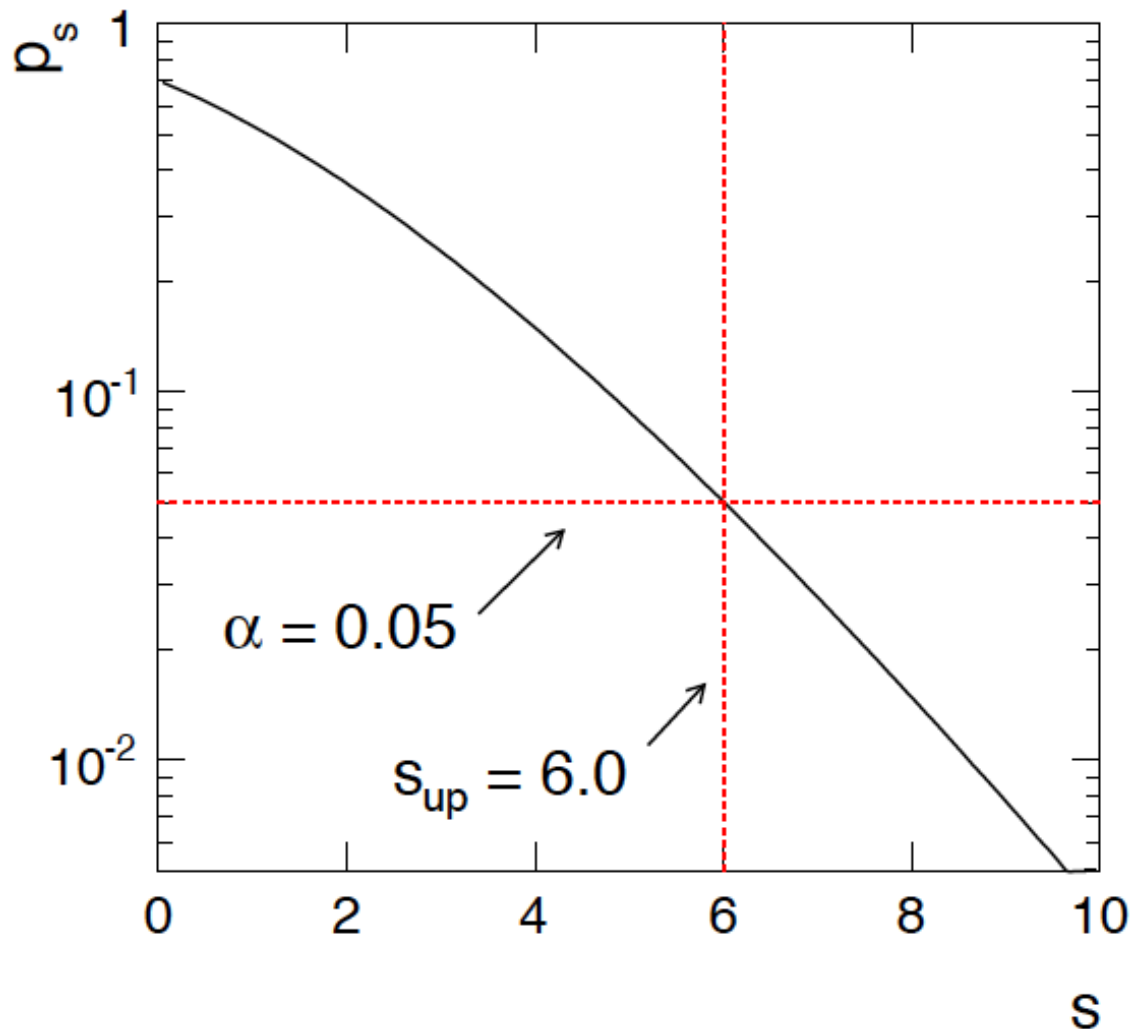$$\alpha = P(n \leq n_{obs}; s_{up}, b) = \sum_{n=0}^{n_{obs}} \frac{(s_{up} + b)^n}{n!} e^{-(s_{up}+b)}$$

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

$$= \frac{1}{2} F_{\chi^2}^{-1}(0.95; 2(5 + 1)) - 4.5 = 6.0$$
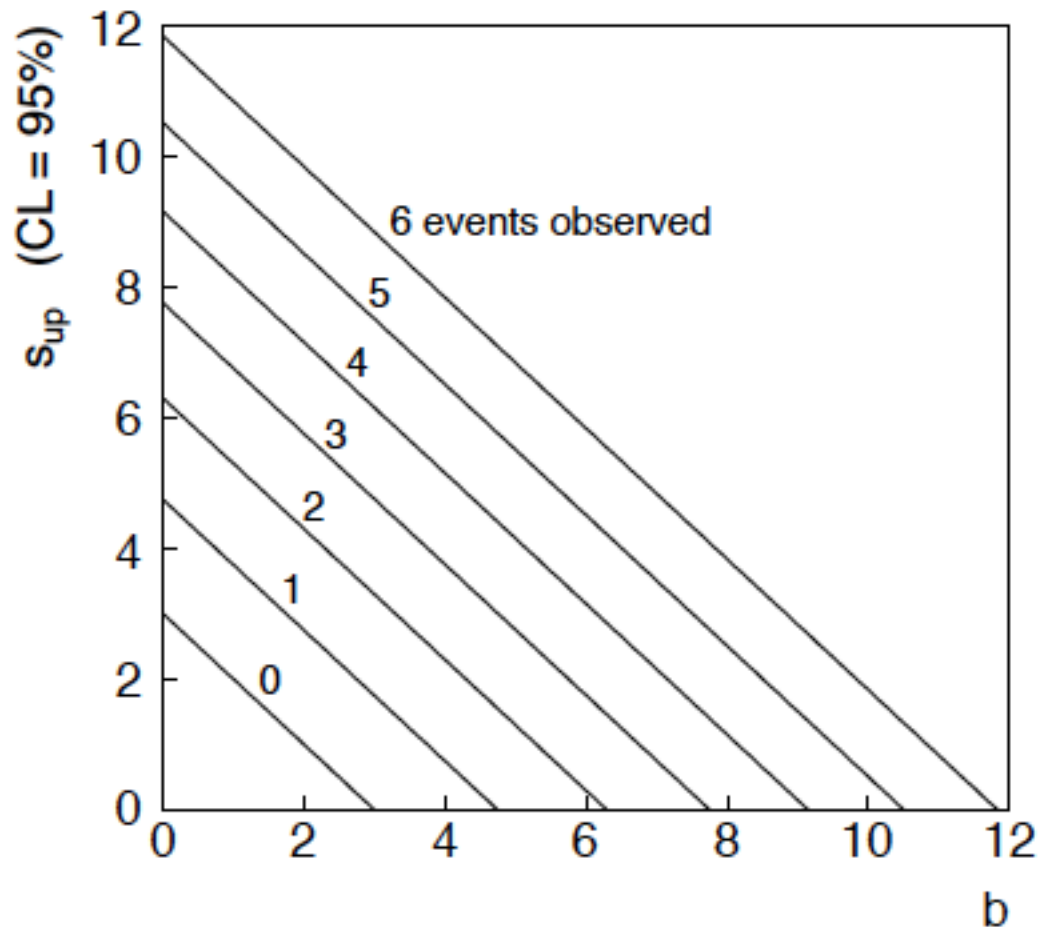
# Frequentist upper limit on Poisson parameter

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



$n_{obs} = 5,$

$b = 4.5$

# $n \sim$ 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.

# 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 $s_{up}$

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

Physicist:

We already knew $s \geq 0$ 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 $s$.
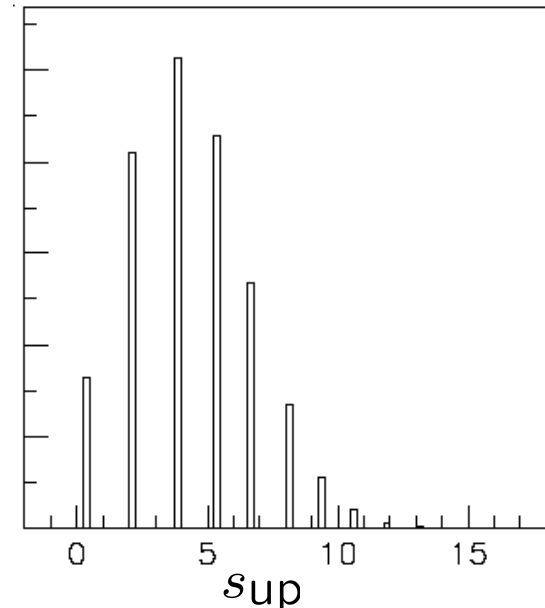
# Expected limit for $s = 0$

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?

Look at the mean limit for the no-signal hypothesis ($s = 0$) (sensitivity).

Distribution of 95% CL limits with $b = 2.5$, $s = 0$.
Mean upper limit = 4.44

# The Bayesian approach to limits

In Bayesian statistics need to start with 'prior pdf' $\pi(\theta)$, this reflects degree of belief about $\theta$ 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. $n \sim$ Poisson($s+b$), 95% CL upper limit on $s$ from

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

# Bayesian prior for Poisson parameter

Include knowledge that $s \geq 0$ by setting prior $\pi(s) = 0$ for $s < 0$.

Could try to reflect 'prior ignorance' with e.g.

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

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 $s_{\text{up}}$:

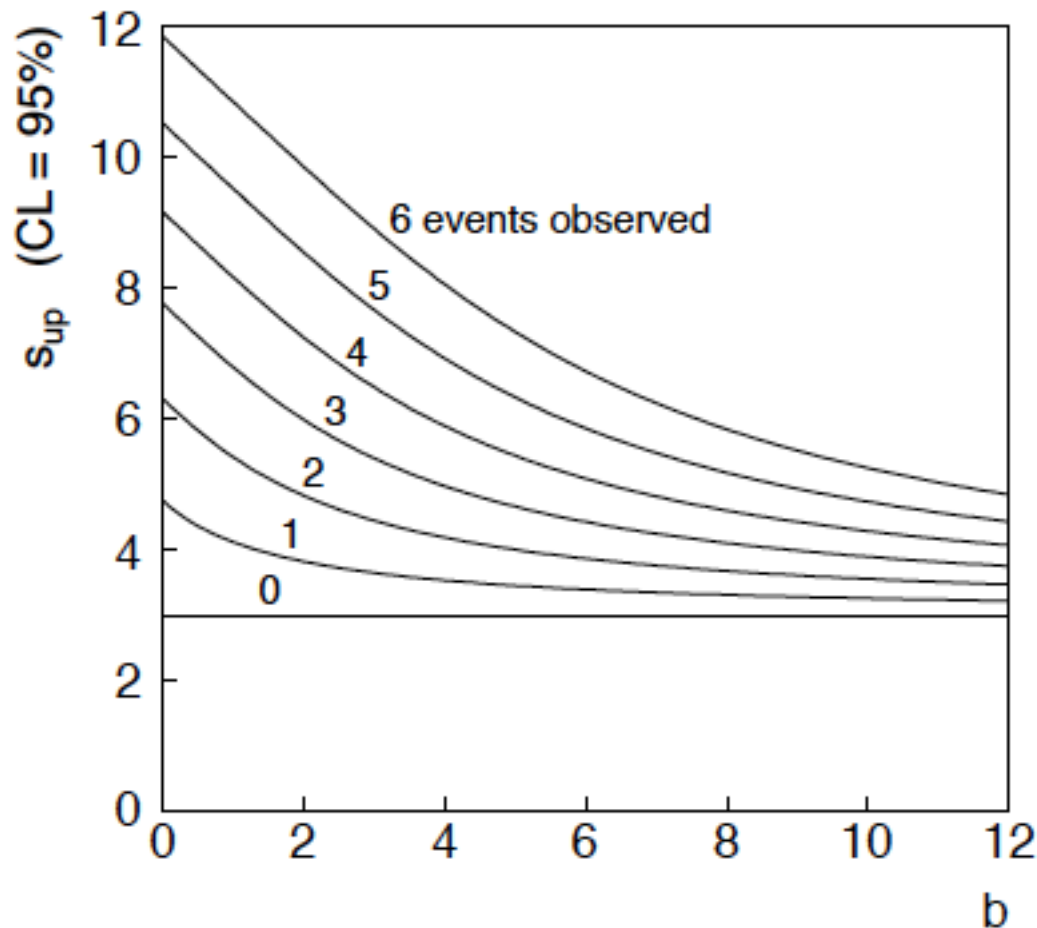$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1} [p, 2(n+1)] - b$$

where

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

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 $n = 0$.

# 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.

# 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

# 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 ($\boldsymbol{\theta}_\mathrm{s}$, $\boldsymbol{\theta}_\mathrm{b}$, $b_\mathrm{tot}$)

Likelihood function is

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

# The profile likelihood ratio

Base significance test on the profile likelihood ratio:

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

maximizes $L$ for specified $\mu$

maximize $L$

Define critical region of test of $\mu$ 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 $\mu \geq 0$.

So take critical region for test of $\mu = 0$ corresponding to high $q_0$ and $\hat{\mu} > 0$ (data characteristic for $\mu \geq 0$).

That is, to test background-only hypothesis define statistic

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

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

Note that even though here physically $\mu \geq 0$, 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.

# 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:

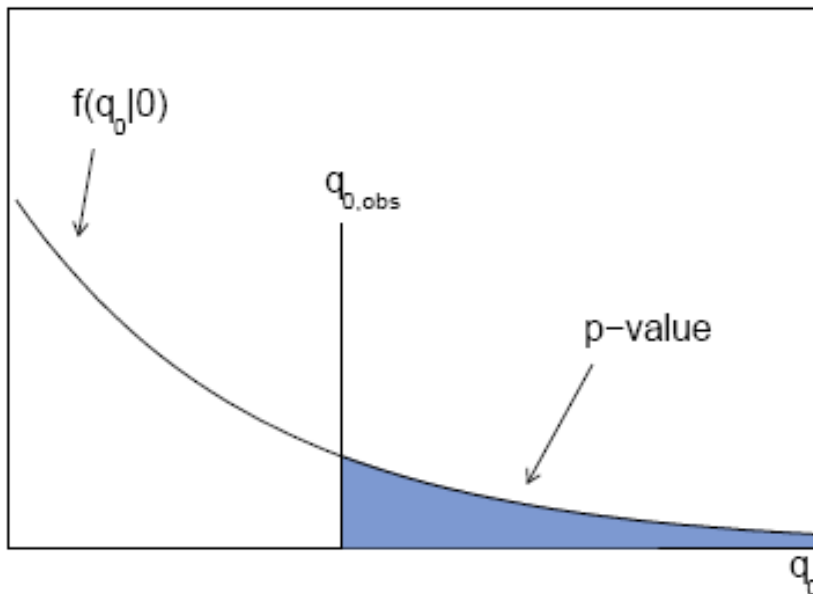$$f(q_0|0) = \frac{1}{2}\delta(q_0) + \frac{1}{2}\frac{1}{\sqrt{2\pi}}\frac{1}{\sqrt{q_0}}e^{-q_0/2}$$

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

# *p*-value for discovery

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

$$p_0 = \int_{q_{0,\text{obs}}}^{\infty} f(q_0|0) \, dq_0$$

use e.g. asymptotic formula

From *p*-value get equivalent significance,

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

# 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

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

# Monte Carlo test of asymptotic formula

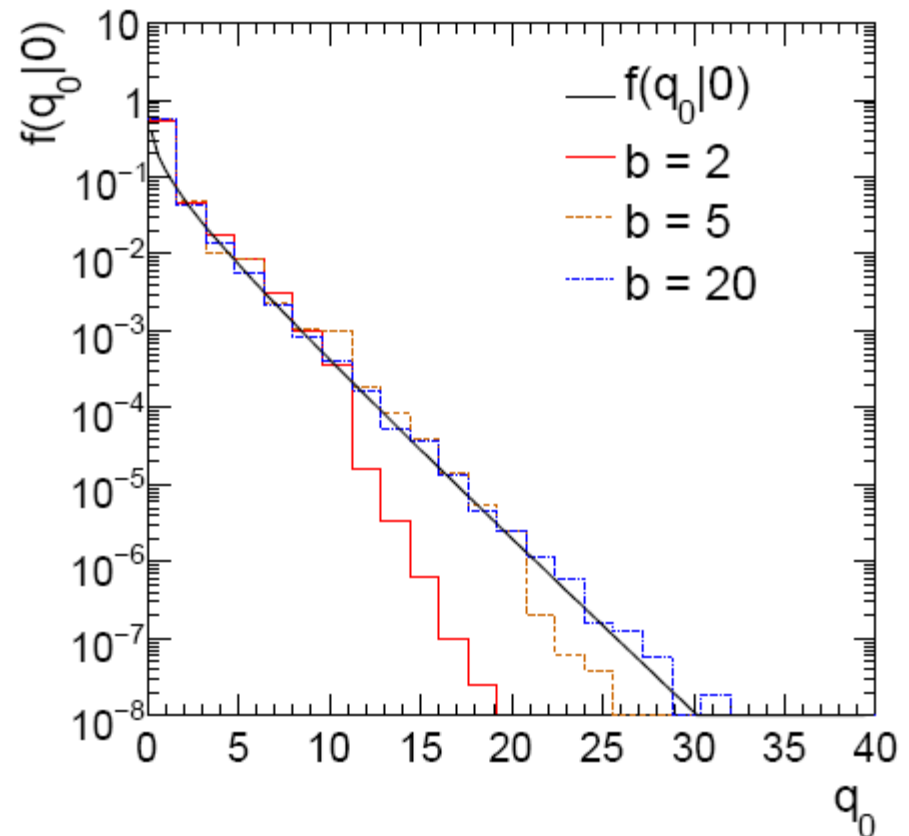$$n \sim \text{Poisson}(\mu s + b)$$

$$m \sim \text{Poisson}(\tau b)$$

$\mu$ = param. of interest
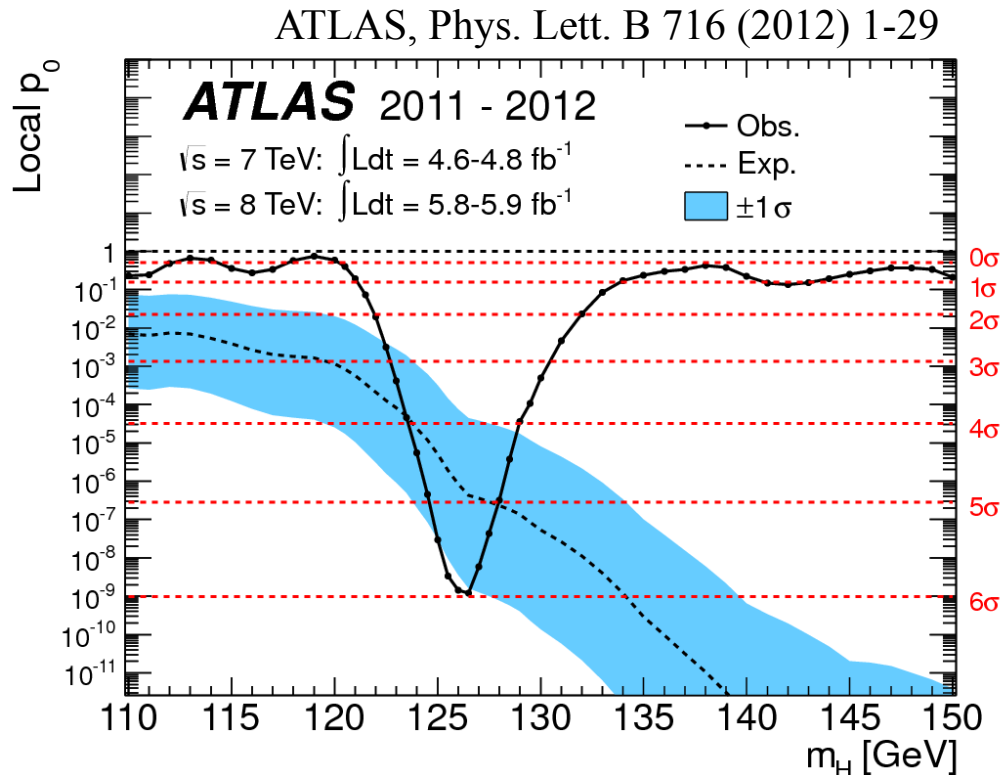
$b$ = nuisance parameter

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

Asymptotic formula is good approximation to $5\sigma$ 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 $m_H$, 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 $m_H$.

ATLAS, Phys. Lett. B 716 (2012) 1-29



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

# Test statistic for upper limits

For purposes of setting an upper limit on $\mu$ use

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

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

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

Large sample approximation: $\qquad \boxed{p_\mu = 1 - \Phi\left(\sqrt{q_\mu}\right)}$

95% CL upper limit on $\mu$ is highest value for which $p$-value is not less than 0.05.
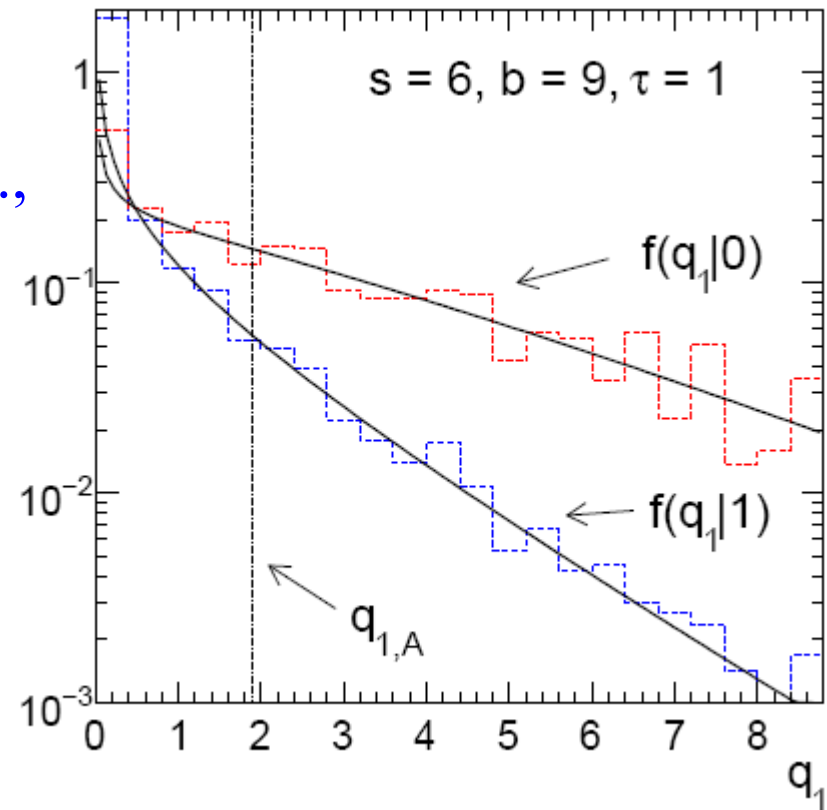
# Monte Carlo test of asymptotic formulae

Consider again $n \sim$ Poisson $(\mu s + b)$, $m \sim$ Poisson$(\tau b)$
Use $q_\mu$ to find $p$-value of hypothesized $\mu$ 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".
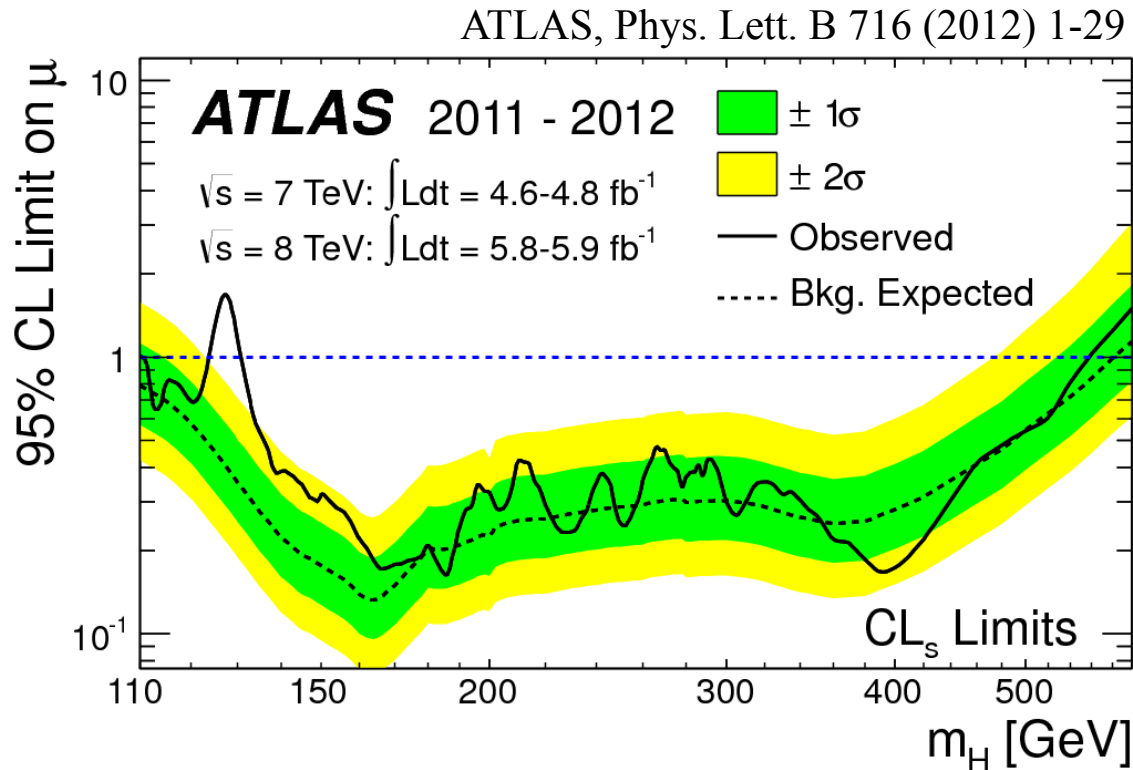
Here asymptotic formulae good for $s = 6$, $b = 9$.

# How to read the green and yellow limit plots

For every value of $m_{\mathrm{H}}$, find the upper limit on $\mu$.

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

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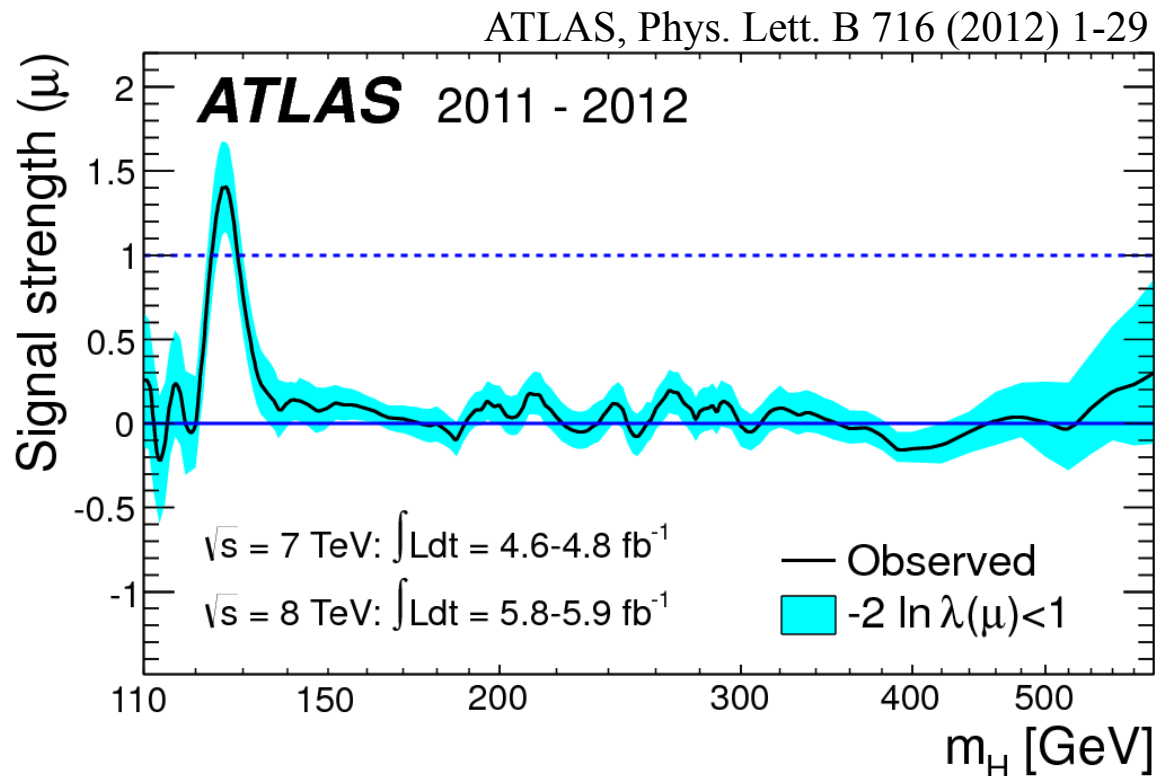


ATLAS, Phys. Lett. B 716 (2012) 1-29

# How to read the "blue band"

On the plot of $\hat{\mu}$ versus $m_H$, the blue band is defined by

$$-2\ln\lambda(\mu) = -2\ln(L(\mu)/L(\hat{\mu})) < 1 \text{ i.e., } \ln L(\mu) > \ln L(\hat{\mu}) - \frac{1}{2}$$

i.e., it approximates the 1-sigma error band (68.3% CL conf. int.)

ATLAS, Phys. Lett. B 716 (2012) 1-29

# Finishing Lecture 1

So far we have introduced the basic ideas of:

Probability (frequentist, subjective)

Parameter estimation (maximum likelihood)

Statistical tests (reject $H$ if data found in critical region)

Confidence intervals (region of parameter space not rejected by a test of each parameter value)

We saw tests based on the profile likelihood ratio statistic

Sampling distribution independent of nuisance parameters in large sample limit; simple formulae for p-value.

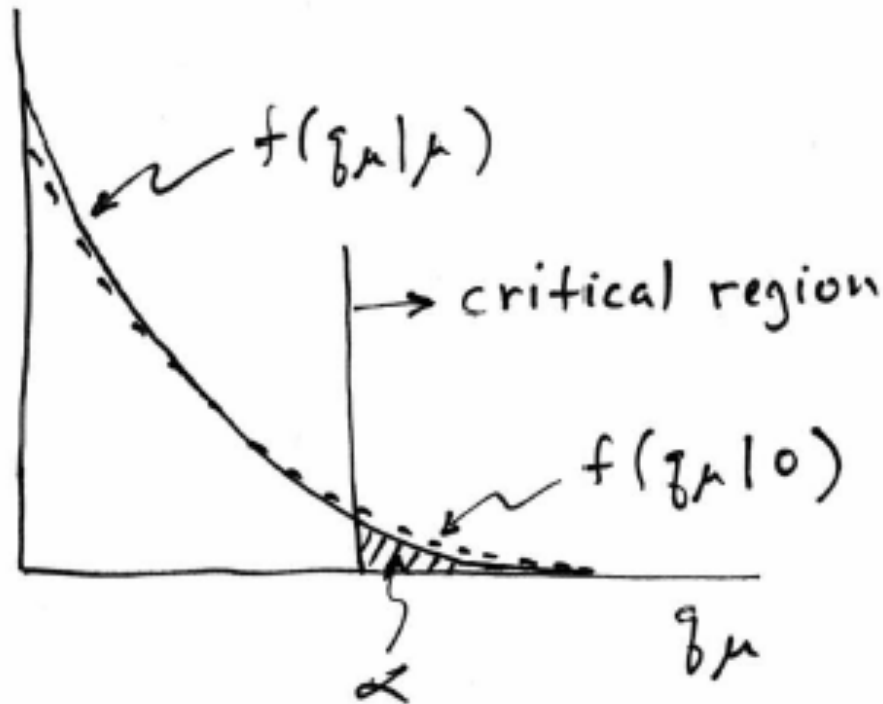Formula for upper limit can give empty confidence interval if e.g. data fluctuate low relative to expected background. Can be avoided/mitigated using Bayesian, CLs, unified,...

# Extra slides

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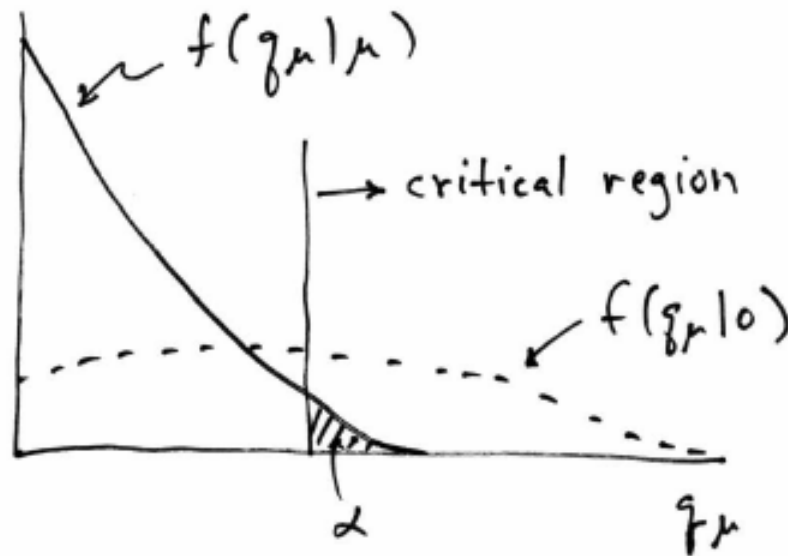
It can be that the effect of a given hypothesized $\mu$ is very small relative to the background-only ($\mu = 0$) prediction.

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

# Having sufficient sensitivity

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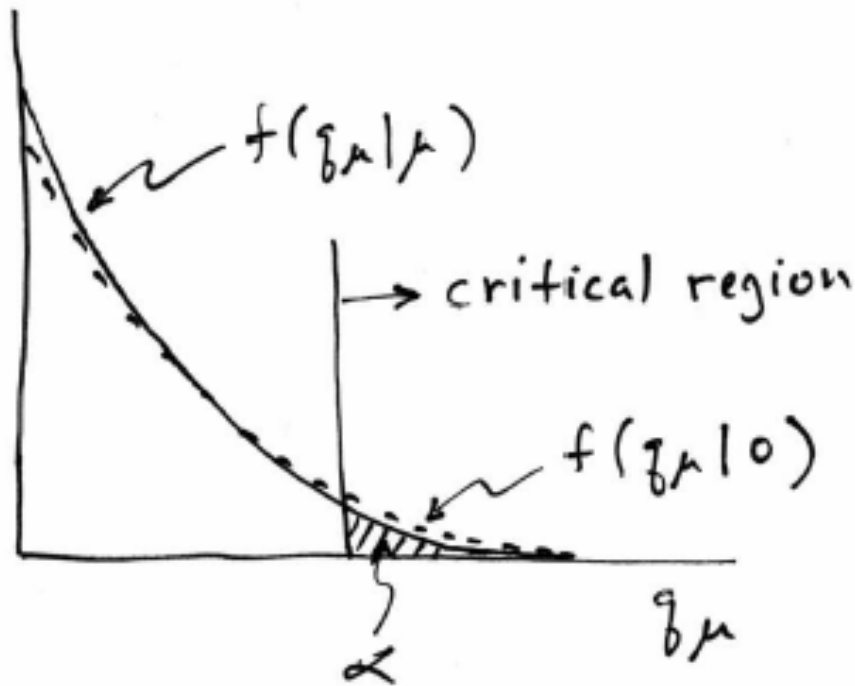


That is, the power (probability to reject $\mu$ if $\mu = 0$) is substantially higher than $\alpha$. Use this power as a measure of the sensitivity.

# Spurious exclusion

Consider again the case of low sensitivity. By construction the probability to reject $\mu$ if $\mu$ is true is $\alpha$ (e.g., 5%).

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



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

"Spurious exclusion"

# Ways of addressing spurious exclusion

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

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

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

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

and led to the "CL$_s$" procedure for upper limits.

Unified intervals also effectively reduce spurious exclusion by the particular choice of critical region.
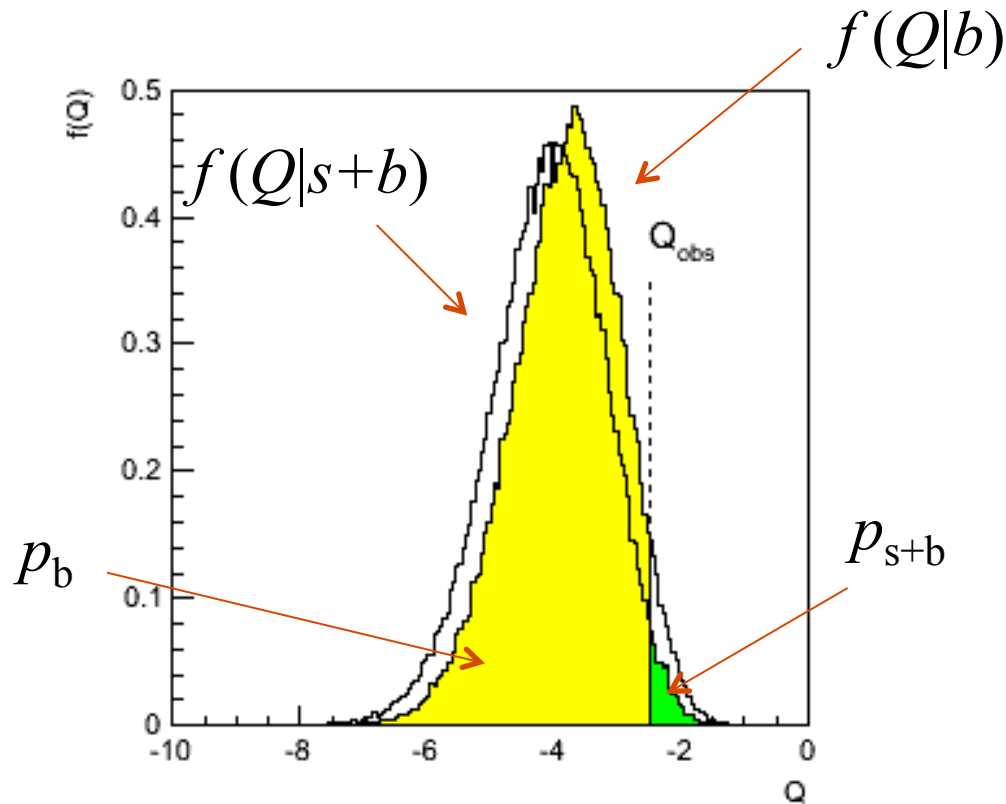
# The CL$_s$ procedure

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

# The CL$_s$ procedure (2)

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

$f(Q|b)$

$f(Q|s+b)$

Q$_{obs}$

$p_b$

$p_{s+b}$

# 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.,

Define:

$$\text{CL}_s = \frac{\text{CL}_{s+b}}{\text{CL}_b}$$

$$= \frac{p_{s+b}}{1 - p_b}$$

Reject s+b hypothesis if:

$$\text{CL}_s \leq \alpha$$



$f(Q|s+b)$    $f(Q|b)$

$1-\text{CL}_b = p_b$

$\text{CL}_{s+b} = p_{s+b}$

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

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

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

In, e.g., a Higgs search, this is done for each value of $m_H$.

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



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

Vertical lines from asymptotic formulae.

# Approximate confidence intervals/regions from the likelihood function

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

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

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

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

so higher $t_{\boldsymbol{\theta}}$ means worse agreement between $\boldsymbol{\theta}$ and the data.

$p$-value of $\boldsymbol{\theta}$ therefore $\qquad p_{\boldsymbol{\theta}} = \int_{t_{\boldsymbol{\theta},\mathrm{obs}}}^{\infty} f(t_{\boldsymbol{\theta}}|\boldsymbol{\theta}) \, dt_{\boldsymbol{\theta}}$

need pdf

# Confidence region from Wilks' theorem

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

$$f(t_\theta | \theta) \sim \chi_n^2$$

chi-square dist. with # d.o.f. = # of components in $\boldsymbol{\theta} = (\theta_1, ..., \theta_n)$.

Assuming this holds, the *p*-value is

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

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

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

# Confidence region from Wilks' theorem (cont.)

i.e., boundary of confidence region in $\boldsymbol{\theta}$ space is where

$$\ln L(\boldsymbol{\theta}) = \ln L(\hat{\boldsymbol{\theta}}) - \tfrac{1}{2} F^{-1}_{\chi^2_n}(1-\alpha)$$

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

$$F^{-1}_{\chi^2_1}(0.683) = 1$$

and so the 68.3% confidence level interval is determined by

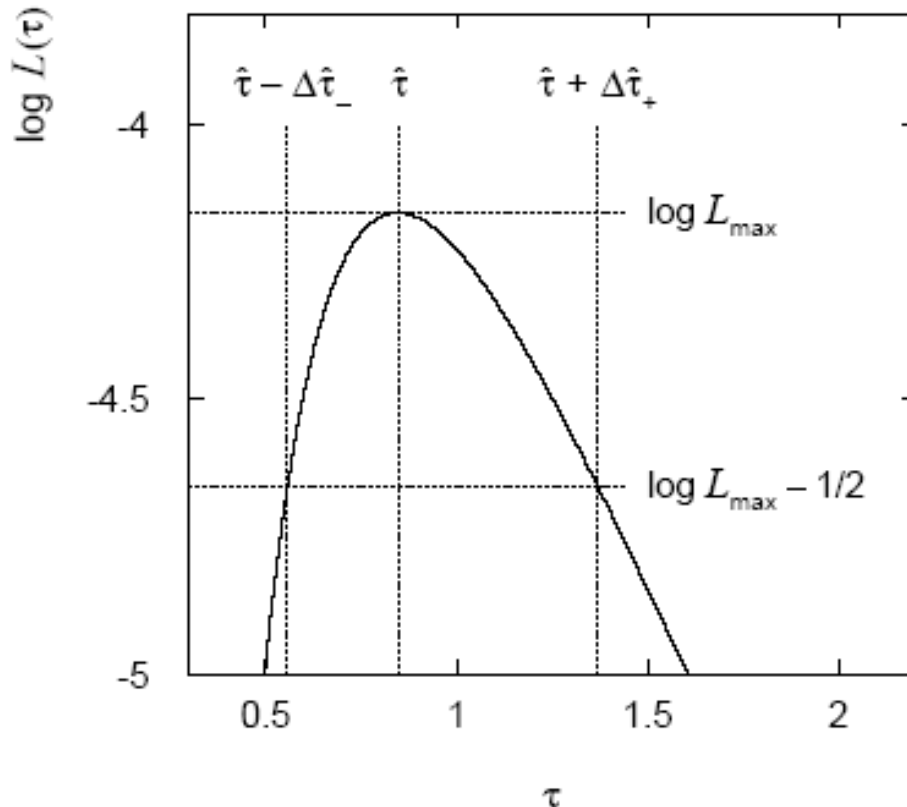$$\ln L(\theta) = \ln L(\hat{\theta}) - \frac{1}{2}$$

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 $L$

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

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

| $Q_\alpha$ | $1 - \alpha$ | | | | |
|---|---|---|---|---|---|
| | $n = 1$ | $n = 2$ | $n = 3$ | $n = 4$ | $n = 5$ |
| 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, $Q_\alpha$ increases with *n* for a given CL $= 1 - \alpha$.

| $1 - \alpha$ | $Q_\alpha$ | | | | |
|---|---|---|---|---|---|
| | $n = 1$ | $n = 2$ | $n = 3$ | $n = 4$ | $n = 5$ |
| 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 |

# Large sample distribution of the profile likelihood ratio (Wilks' theorem, cont.)

Suppose problem has likelihood $L(\boldsymbol{\theta}, \boldsymbol{v})$, with

$$\boldsymbol{\theta} = (\theta_1, \ldots, \theta_N) \qquad \leftarrow \text{ parameters of interest}$$

$$\boldsymbol{\nu} = (\nu_1, \ldots, \nu_M) \qquad \leftarrow \text{ nuisance parameters}$$

Want to test point in $\boldsymbol{\theta}$-space. Define profile likelihood ratio:

$$\lambda(\boldsymbol{\theta}) = \frac{L(\boldsymbol{\theta}, \hat{\hat{\nu}}(\boldsymbol{\theta}))}{L(\hat{\boldsymbol{\theta}}, \hat{\nu})}, \quad \text{where} \quad \hat{\hat{\nu}}(\boldsymbol{\theta}) = \underset{\nu}{\arg\max} L(\boldsymbol{\theta}, \nu)$$

"profiled" values of $\boldsymbol{v}$

and define $q_\theta = -2 \ln \lambda(\boldsymbol{\theta})$.

Wilks' theorem says that distribution $f(q_\theta | \boldsymbol{\theta}, \boldsymbol{v})$ approaches the chi-square pdf for $N$ degrees of freedom for large sample (and regularity conditions), independent of the nuisance parameters $\boldsymbol{v}$.

# *p*-values in cases with nuisance parameters

Suppose we have a statistic $q_\theta$ that we use to test a hypothesized value of a parameter $\theta$, such that the *p*-value of $\theta$ is

$$p_\theta = \int_{q_{\theta,\text{obs}}}^{\infty} f(q_\theta | \theta, \nu)\, dq_\theta$$

Fundamentally we want to reject $\theta$ only if $p_\theta < \alpha$ for all $\nu$.

$\rightarrow$ "exact" confidence interval

Recall that for statistics based on the profile likelihood ratio, the distribution $f(q_\theta | \theta, \nu)$ becomes independent of the nuisance parameters in the large-sample limit.

But in general for finite data samples this is not true; one may be unable to reject some $\theta$ values if all values of $\nu$ must be considered, even those strongly disfavoured by the data (resulting interval for $\theta$ "overcovers").

HCPSS 2016 / Statistics Lecture 1

# Profile construction ("hybrid resampling")

Approximate procedure is to reject $\theta$ if $p_\theta \leq \alpha$ where the $p$-value is computed assuming the profiled values of the nuisance parameters:

$$\hat{\hat{\nu}}(\theta)$$

"double hat" notation means value of parameter that maximizes likelihood for the given $\theta$.

The resulting confidence interval will have the correct coverage for the points $(\theta, \hat{\hat{\nu}}(\theta))$ .

Elsewhere it may under- or overcover, but this is usually as good as we can do (check with MC if crucial or small sample problem).